

## ОГЛАВЛЕНИЕ

|   |           |
|---|-----------|
| <b>ВВЕДЕНИЕ</b>   | <b>5</b>  |
| <b>1 АНАЛИЗ ДИСЦИПЛИН ОБСЛУЖИВАНИЯ</b>  | <b>7</b>  |
| 1.1 Дисциплины обслуживания   | 7         |
| 1.2 Приоритетные очереди  | 7         |
| 1.3 Алгоритм управления очередями на основе классов                                 | 9         |
| 1.4 Алгоритм иерархического маркерного ведра  | 11        |
| 1.5 Алгоритм иерархических честных кривых обслуживания                              | 12        |
| 1.6 Взвешенный алгоритм честного обслуживания очередей                              | 14        |
| 1.7 Взвешенный алгоритм честного обслуживания очередей на основе классов            | 20        |
| 1.8 Выводы  | 22        |
| <b>2 РЕАЛИЗАЦИЯ CLASS-BASED WFQ В ЯДРЕ LINUX</b>                                    | <b>24</b> |
| 2.1 Описание устройства подсистемы планировки в ядре Linux                          | 24        |
| 2.2 Интерфейс управления трафиком   | 25        |
| 2.3 Описание интерфейса   | 27        |
| 2.4 Алгоритм CBWFQ  | 29        |
| 2.4.1 Структуры хранения данных Class-Based WFQ                                     | 29        |
| 2.4.2 Добавление пакета в очередь   | 31        |
| 2.4.3 Удаление пакета из очереди  | 33        |
| <b>3 ТЕСТИРОВАНИЕ РАЗРАБОТАННОГО МОДУЛЯ ДИСЦИПЛИНЫ ОБСЛУЖИВАНИЯ CLASS-BASED WFQ</b> | <b>35</b> |
| 3.1 Описание тестовой среды   | 35        |
| 3.2 Анализ точности выделения канала при конкурирующем трафике                      | 36        |
| 3.3 Анализ точности выделения канала при независимом трафике                        | 38        |
| <b>ЗАКЛЮЧЕНИЕ</b>   | <b>40</b> |
| <b>СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ</b>   | <b>41</b> |
| <b>СПИСОК СОКРАЩЕНИЙ И УСЛОВНЫХ ОБОЗНАЧЕНИЙ</b>                                     | <b>44</b> |

|                               |           |
|-------------------------------|-----------|
| <b>ПРИЛОЖЕНИЕ А . . . . .</b> | <b>45</b> |
| <b>ПРИЛОЖЕНИЕ Б . . . . .</b> | <b>46</b> |
| <b>ПРИЛОЖЕНИЕ В . . . . .</b> | <b>50</b> |

## ВВЕДЕНИЕ

В наши дни на фоне высокого спроса на высокоскоростную и надёжную передачу данных остро встаёт проблема обработки трафика и обеспечения должного уровня обслуживания в узлах компьютерных сетей. Крупные производители сетевого оборудования предоставляют эффективные решения обозначенных проблем, однако оборудование и программное обеспечение стоят немалых денег, что делает их недоступными для массового пользователя. Большое распространение в качестве сетевого ПО получили дистрибутивы Linux, которые из-за многолетнего использования имеют богатый набор возможностей в этой области. Таким образом, встаёт вопрос о интеграции возможностей, предоставляемых гигантами индустрии, с широко распространённым открытым вариантом.

В связи с этим было решено изучить дисциплину обслуживания, разработанную компанией Cisco, известным вендором на рынке сетевых решений, и внедрить её в ядро последней стабильной версии Linux (4.15.3). В качестве дисциплины обслуживания был выбран взвешенный алгоритм честного обслуживания основанного на классах (Class-Based Weighted Fair Queuing, CBWFQ).

CBWFQ является расширением функциональности широко известного взвешенного алгоритма честного обслуживания очереди (Weighted Fair Queuing). Он поддерживает определение пользовательских классов трафика на основе ряда критериев соответствия (протокол, входящий интерфейс и т.д.) и назначение их характеристик, которые отвечают за выделяемые классу ресурсы (вес, пропускная способность, максимальное количество пакетов в очереди, задержка). Такой подход предоставляет гибкую настройку распределения пропускной способности канала между классами трафика и оказывается весьма эффективным в передаче данных в сравнении с рядом других популярных дисциплин.

**Цель работы** – реализовать алгоритм Class-Based WFQ в виде модуля ядра Linux, основываясь на имеющихся описаниях в соответствующей литературе. Для выполнения цели работы необходимо выполнить следующие задачи:

1. Проанализировать и сравнить дисциплины обслуживания PQ, CBQ, HTB, HFSC, FWFQ, CBWFQ.
2. Восстановить алгоритмы Class-Based WFQ.
3. Настроить среду для реализации и тестирования.

4. Реализовать модуль ядра CBWFQ в ядре Linux.
5. Реализовать интерфейс утилиты tc для управления модулем.
6. Провести тестирование.

Сравнительный анализ дисциплин обслуживания в работе основываются на научных статьях, документации и исходных кодах; архитектура подсистемы Linux по управлению качеством обслуживания затрагивается в книге [1], однако в силу отсутствия полной документации все выводы о структуре делались на основе исходного кода ядра Linux. Алгоритм CBWFQ воссоздавался на основе документации Cisco, книг [1] и [2] и исходного кода существующих дисциплин обслуживания.

# 1 АНАЛИЗ ДИСЦИПЛИН ОБСЛУЖИВАНИЯ

## 1.1 Дисциплины обслуживания

В теории массового обслуживания дисциплиной обслуживания очередей (ДО) называют правило выбора заявок из очереди для обслуживания[3], определяющее способ отсылки данных. На практике дисциплина обслуживания определяется не только порядком обслуживания, но и способом организации очередей. Очередь – это базовый элемент управления трафиком, являющийся неотъемлемым элементом системы планирования. Обычно, под очередью подразумевают буфер, где пакеты ожидают передачи устройством. [4]

Дисциплины обслуживания делятся на классовые и бесклассовые дисциплины.

1. Бесклассовые дисциплины обслуживания могут принимать трафик, перепланировать его, задерживать или отбрасывать. Этот тип дисциплин обычно используется по умолчанию или для ограничения трафика через узел. Такой тип дисциплин обслуживания малофункционален и имеет множество недостатков; обычно их используют для обслуживания классов в классовых дисциплинах.
2. Классовые дисциплины обслуживания разделяют трафик на классы с помощью процесса классификации, основанного в основном на фильтрах; такой подход позволяет дифференцировать трафик, отдавая канал более приоритетизированному.[5]

Рассмотрим и проанализируем наиболее известные дисциплины обслуживания (PQ, CBQ, HTB, HFSC) и дисциплины семейства WFQ (FWFQ и CBWFQ).

## 1.2 Приоритетные очереди

Приоритетные очереди (Priority Queueing, PQ) – это техника обслуживания, при которой используется множество очередей с разными приоритетами. Очереди обслуживаются в циклическом порядке (алгоритмом Round-robin) от самого высокого до самого низкого приоритета; обслуживание следующей по порядку очереди происходит, если более приоритетные очереди пусты. Каждая очередь внутри обслуживается в порядке FIFO (First-In, First-Out). В случае пе-

реполнения отбрасываются пакеты из очереди с более низким приоритетом.[6]  
 Схема дисциплины представлена на рисунке 1.

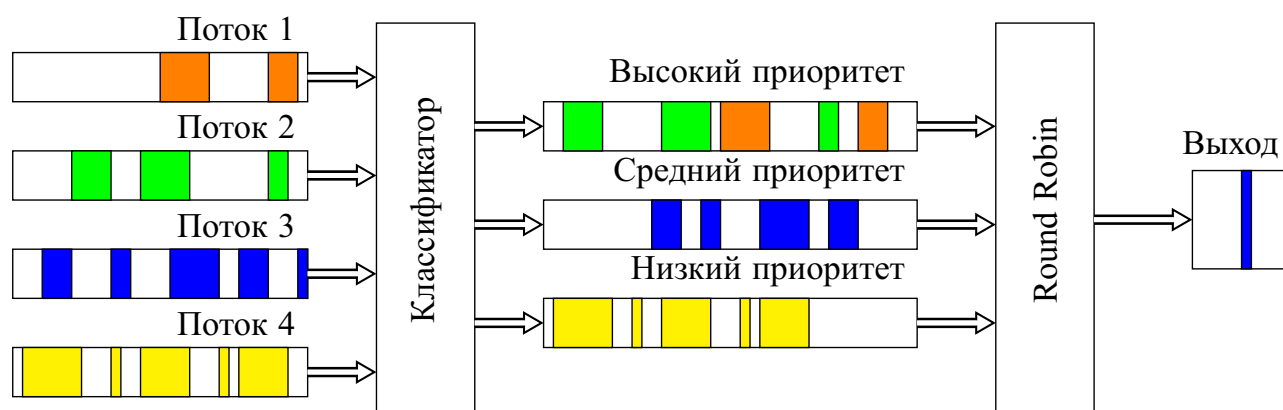


Рисунок 1 — Схема обслуживания алгоритмом приоритетных очередей

Дисциплина используется, чтобы снизить время отклика, когда нет нужды замедлять трафик[7].

В Linux алгоритм реализован в виде дисциплины `prio`, которая создаёт фиксированное значение очередей обслуживания (групп или `bands`), планируемые дисциплиной `pfifo_fast`, и управляет очередями в соответствии с картой приоритетов.[7]

Преимущества алгоритма состоят в следующем:

- возможность понижения времени отклика, когда нет необходимости замедлять трафик [7];
- наиболее простая в реализации классовая дисциплина обслуживания;
- для `software-based` маршрутизаторов PQ предоставляет относительно небольшую вычислительную нагрузку на систему в сравнении с более сложными ДО;
- PQ позволяет маршрутизаторам организовывать буферизацию пакетов и обслуживать один класс трафика отдельно от других. [8]

Однако приоритетные очереди обладают рядом существенных недостатков:

- возникает проблема простоя канала (отсутствие обслуживания в течение продолжительного времени) для низкоприоритетного трафика при избытке высокоприоритетного[6];

- избыточный высокоприоритетный трафик может значительно увеличивать задержку и джиттер для менее приоритетного трафика;
  - не решается проблема с TCP и UDP, когда TCP-трафику назначается высокий приоритет и он пытается поглотить всю пропускную способность.
- [8]

### 1.3 Алгоритм управления очередями на основе классов

Алгоритм управления очередями на основе классов (Class Based Queueing, CBQ) – это классовая дисциплина обслуживания, которая реализует иерархическое разделение канала между классами и позволяет шейпинг трафика.[9]

Главная цель CBQ – это планировка пакетов в очередях, гарантия определённой скорости передачи и разделение канала. Если в очереди нет пакетов, её пропускная способность становится доступной для других очередей. Сила этого метода состоит в том, что он позволяет справляться со значительно различными требованиями к пропускной способности канала среди потоков. Это реализовано путём назначения определённого процента доступной ширины канала каждой очереди. CBQ избегает проблему простоя канала, которой страдает алгоритм PQ, так как как минимум один пакет обслуживается от каждой очереди в течение цикла обслуживания.[6]

Алгоритм CBQ представляет канал в виде иерархической структуры [10] пример которой представлен на рисунке 2. Голубым цветом обозначен узел, представляющий собой основной канал; он разделяется между двумя классами трафика: интерактивным (левый узел) и остальным (правый узел), – которым назначается процент пропускной способности от остального канала. Весь трафик, относящийся к классу, будет получать выделенную пропускную способность для этого класса. Эти классы трафика могут разделяться на подклассы и так далее. Если класс не использует пропускную способность, она будет выделяться классу-соседу. Этот механизм называется механизмом разделения канала [10].

Алгоритм CBQ состоит в следующем. Сначала пакеты классифицируются в классы обслуживания в соответствии с определёнными критериями и сохраняются в соответствующей очереди. Очереди обслуживаются циклически. Различное количество пропускной способности может быть назначено для каждой очереди двумя различными способами: с помощью позволения очереди отправлять

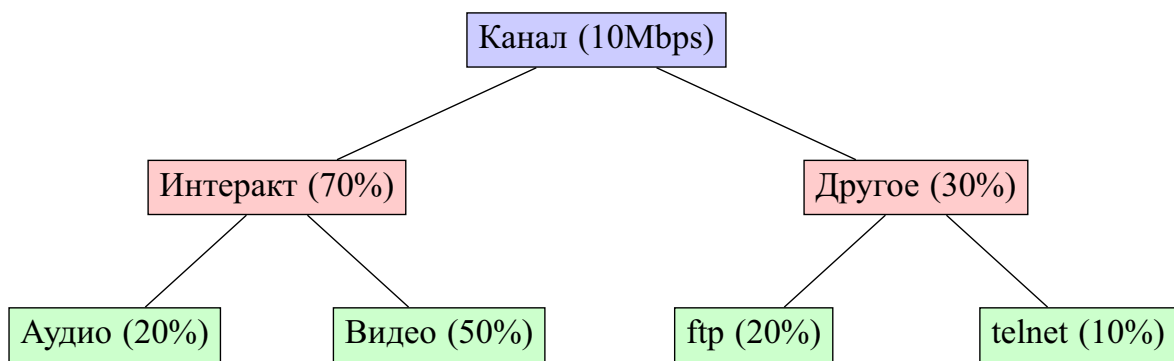


Рисунок 2 — Пример работы механизма разделения канала.

более чем один пакет на каждый цикл обслуживания или с помощью позволения очереди отправлять только один пакет за цикл, но при этом очередь может быть обслужена несколько раз за цикл.[6]

Вычисления в CBQ основываются на вычислении времени в микросекундах между запросами, на основе которого рассчитывается средняя загрузка канала; в этом и состоит главная проблема неточности CBQ в Linux.[5]

Преимущества алгоритма состоят в следующем:

- позволяет контролировать количество пропускной способности для каждого класса обслуживания;
- каждый класс получает обслуживание, вне зависимости от других классов. Это помогает избежать проблемы PQ, когда при избытке высокоприоритетизированного трафика низкоприоритетизированный не обслуживался вообще.[6]

Недостатки же в большей следуют из особенностей реализации алгоритма в системе Linux:

- честное выделение пропускной способности происходит, только если пакеты из всех очередей имеют сравнительно одинаковый размер. Если один класс обслуживания содержит пакет, который длиннее остальных, этот класс обслуживания получит большую пропускную способность, чем сконфигурированное значение [6];
- высокая сложность реализации. В ядре Linux реализация CBQ приближённая и в некоторых случаях может давать неверные результаты.[5]



## 1.4 Алгоритм иерархического маркерного ведра

Алгоритм иерархического маркерного ведра (Hierarchical Token Bucket, НТВ) – дисциплина обслуживания с иерархическим разделением канала между классами.

НТВ, подобно СВQ, использует механизм разделения канала. НТВ обеспечивает, что количество обслуживания, предоставляемое каждому классу, является, минимальным значением из запрошенного количества и назначенного классу. Когда класс запрашивает меньше, чем ему выделено, оставшаяся пропускная способность распределяется между другими классами, которые требуют обслуживание.[11]

Отличительная особенность НТВ от СВQ состоит в том, что в НТВ принцип работы основывается на определении объема трафика[5], что даёт более точные результаты.

НТВ состоит из произвольного числа иерархически организованных фильтров маркерного ведра (Token Bucket Filter, TBF)[6], однако реализация не использует готовый модуль tbf: алгоритм маркерного ядра встроен в код реализации НТВ, что повышает его эффективность. Внутренние классы содержат фильтры, которые распределяют пакеты по очередям и метаданные для разделения канала. Листовые классы содержат очереди, которые содержат очереди, которые управляются сконфигурированными дисциплинами обслуживания (по умолчанию pfifo\_fast). Пример сконфигурированного дерева НТВ представлен на рисунке 3.

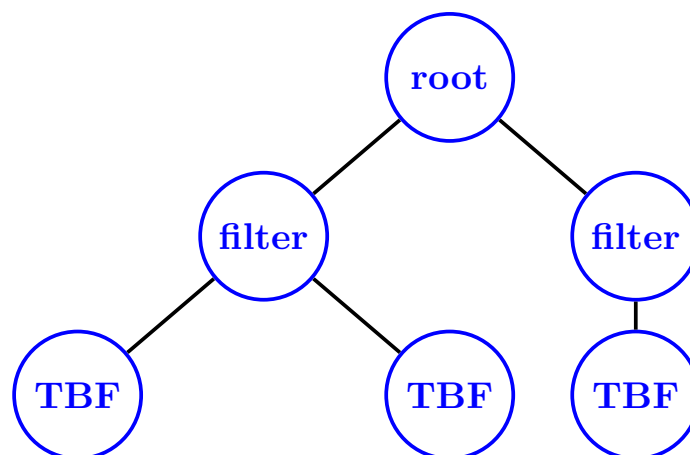


Рисунок 3 — Пример иерархии классов при использовании дисциплины НТВ

При добавлении пакета в очередь НТВ начинает обход дерева от корня для определения подходящей очереди: в каждом узле происходит поиск инструкций, и затем происходит переход в узел, на который ссылается инструкция. Обход заканчивается, когда алгоритм доходит до листа, в очередь которого помещается пакет.[12] В реализации алгоритма существует прямая очередь, которая используется не только в качестве очереди с наивысшим приоритетом, но и как очередь, в которую попадают пакеты, не определённые в другую очередь. Это мера не самая удачная, но используется для избежания ошибок.

Преимущества алгоритма НТВ приведены ниже.

- Наиболее используемая дисциплина обслуживания в Linux, так как НТВ эффективно справляется с обработкой пакетов, а конфигурация НТВ легко масштабируется.[5]
- Иерархическая структура предоставляет гибкую возможность конфигурировать трафик.
- Не зависит от характеристик интерфейса и не нуждается в знании о лежащей в основе пропускной способности выходного интерфейса из-за свойств TBF. [12]
- Вычислительно проще, чем алгоритм CBQ.[11]

Недостатки.

- Медленнее CBQ в  $N$  раз, где  $N$  – глубина дерева разделения, что, однако, компенсируется простотой вычислений.[11]
- Нужно что-то ещё весомое.

### **1.5 Алгоритм иерархических честных кривых обслуживания**

HFSC (Hierarchical Fair-Service Curve) – иерархический алгоритм планирования пакетов, основанный на математической модели честных кривых обслуживания (Fair Service Curve), где под термином "кривая обслуживания" подразумевается зависящая от времени неубывающая функция, которая служит нижней границей количества обслуживания, предоставляемого системой.[13]

HFSC ставит перед собой цели:

- гарантировать точное выделение пропускной способности и задержки для всех листовых классов (критерий реального времени);
- честно выделять избыточную пропускную способность так, как указано классовой иерархией (критерий разделения канала);
- минимизировать несоответствие кривой обслуживания идеальной модели и действительного количества обслуживания.[14]

Алгоритм планировки основан на двух критериях: критерий реального времени (real-time) и критерий разделения канала (link-sharing). Критерии реального времени используются для выбора пакета в условиях, когда есть потенциальная опасность, что гарантия обслуживания для листового класса нарушается. В ином случае используется критерий разделения канала.[14]

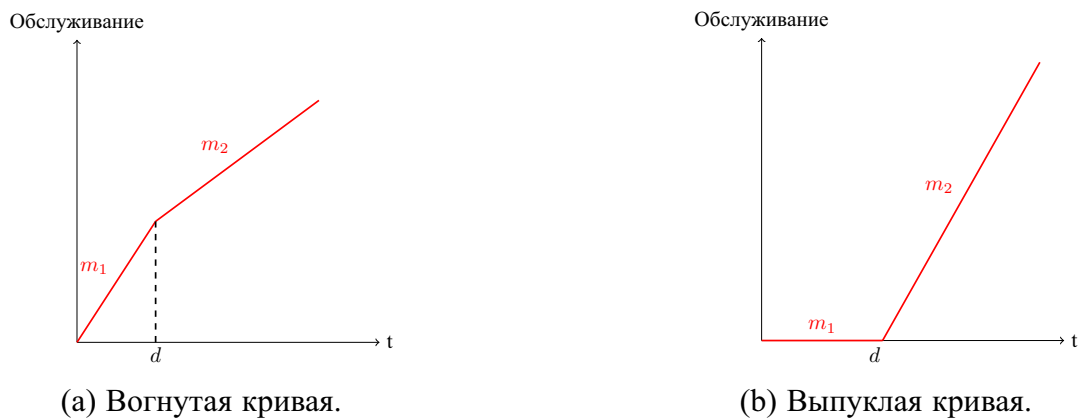


Рисунок 4 — Примеры кривых обслуживания.  $m_1$  – скорость в стационарном состоянии,  $m_2$  – скорость в режиме burst,  $d$  – время, за которое происходит передача в режиме burst.[15]

На рисунке 4 изображены примеры кривых обслуживания, используемых в дисциплине HFSC; параметры  $m_1$ ,  $m_2$  и  $d$ , отображённые на графиках, задаются при конфигурации дисциплины.[15]

HFSC использует три типа временных параметров: время крайнего срока (deadline time), "подходящее" время (eligible time) и виртуальное время (virtual time). Время крайнего срока назначается таким образом, чтобы, если крайние сроки всех пакетов сессии выполнены, его кривая была гарантирована. "Подходящее" время используется для выбора критерия планировки для следующего пакета. Виртуальное время показывает нормализованное количество обслуживания, которое получил класс. Виртуальное время присуще всем вершинам дерева

классов, так как является важным параметром при критерии разделение канала, при котором должно минимизироваться несоответствие между виртуальным временем класса и временами его соседей (так как в идеальной модели виртуальное время соседей одинаково); при выборе критерия разделения канала алгоритм рекурсивно, начиная с корня, обходит всё дерево, переходя в вершины с наименьшим виртуальным временем. Время крайнего срока и «подходящее» время используются дополнительно в листовых классах, так как в этих вершинах непосредственно содержатся очереди.[13]

Основное преимущество алгоритма состоит в том, что он основан на формальной модели с доказанными нижними границами. Он даёт гарантированные результаты и вычисляет более точно, чем дисциплины CBQ и HTB, которые служат схожим целям.

Главные же недостатки HFSC заключены в его достоинстве. Алгоритм основан на формальной модели и имеет множество параметров, требующих дополнительных расчётов и времени на подготовку к конфигурации. Также он довольно сложен в реализации и поддержке.

## **1.6 Взвешенный алгоритм честного обслуживания очередей**

WFQ (Weighted Fair Queueing) – динамический метод планировки пакетов, который предоставляет честное разделение пропускной способности всем потокам трафика. WFQ применяет вес, чтобы идентифицировать и классифицировать трафик в поток и определить, как много выделить пропускной способности каждому потоку относительно других потоков. WFQ на месте планирует интерактивный трафик в начало очереди, уменьшая тем самым время ответа, и честно делит оставшуюся пропускную способность между остальными потоками. [16]

WFQ представляет собой аппроксимацию обобщённой схемы разделения процессорного времени (General Processor Sharing, GPS). GPS – это схема, обеспечивающая честное обслуживание по типу взвешенной максиминной схемы равномерного распределения ресурсов (схема, при которой каждому пользователю назначается определённый вес и выделяется равномерная доля ресурсов, пропорциональная этому весу). В соответствии со схемой GPS каждый поток трафика помещается в собственную логическую очередь, после чего бесконечно малый объём данных из каждой непустой очереди обслуживается по круговому принципу. Необходимость обработки бесконечно малого объёма данных на

каждом круге обусловлена требованием обслуживания всех непустых очередей на любом конечном временном интервале. Из-за чего схема GPS является справедливой в любой момент времени. Однако технически невозможно реализовать данную схему на практике; WFQ же предлагает брать за рабочую единицу пакет.[2]

Алгоритм WFQ использует концепцию виртуального времени, чтобы отслеживать расчёты схемы GPS. В рамках этой концепции используется понятие события  $j$ , которое обозначает прибытие или отправление пакета во время  $t_j$ . Виртуальное время идёт не с той же скоростью, с какой идёт реальное; его скорость зависит от активных потоков в рассматриваемый реальный промежуток времени. Разница представлена на рисунке 5.

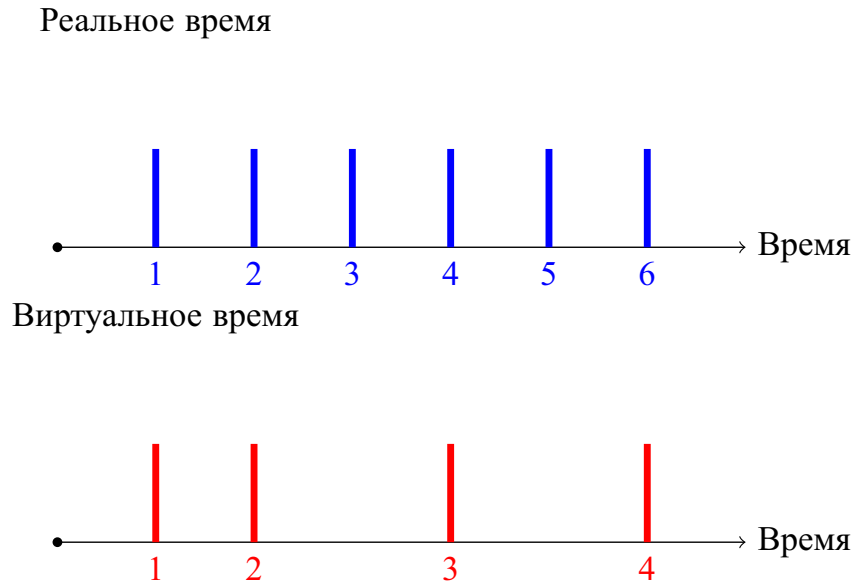


Рисунок 5 — Сравнение движений реального и виртуальных времён.

По Формуле (1) рассчитывается движение виртуального времени в GPS. Скорость виртуального времени зависит от количества активных сессий.

$$V_{t_{j-1}+\tau} = V_{t_{j-1}} + \frac{\tau}{\sum_{i \in B_j} w^i}, \tau \leq t_j - t_{j-1} \quad (1)$$

$$V(0) = 0 \quad (2)$$

где

$w^i$  — вес потока  $i$ .

$B_j$  — обозначает, является ли поток на данный момент активным.

Виртуальное время завершения обслуживания обозначает виртуальное время, когда должно завершиться обслуживание в соответствии со схемой GPS. Пакет с наименьшим виртуальным временем первым покинет систему обслуживания. Для вычисления виртуального времени завершения обслуживания по Формуле (4) необходимо вычислить виртуальное время начала обслуживания по Формуле (3).

$$S_i^k = \max\{F_i^{k-1}, V_{a_i^k}\}, \quad (3)$$

где

$S_i^k$  — виртуальное время начала обработки пакета под номером  $k$  из потока  $i$ .

$F_i^k$  — виртуальное время конца обработки пакета под номером  $k$  из потока  $i$ .

$a_i^k$  — время прибытия пакета под номером  $k$  из потока  $i$ .

$$F_i^k = S_i^k + \frac{L_i^k}{w_i}, \quad (4)$$

$$F_i^0 = 0,$$

где

$L_i^k$  — длина пакета под номером  $k$  из потока  $i$ .

$w_i$  — вес потока  $i$ .

В Таблице 1 приведён пример вычисления виртуального времени заверше-

ния обслуживания. По столбцу  $F$  видно, что систему раньше покинут пакеты с большим весом.

Таблица 1 — Пример вычисления времени завершения обслуживания.

|              | Поток 1             |       |       | Поток 2             |       |       |
|--------------|---------------------|-------|-------|---------------------|-------|-------|
| Вес          | $w_1 = \frac{1}{3}$ |       |       | $w_2 = \frac{2}{3}$ |       |       |
| Номер пакета | $a_1$               | $L_1$ | $F_1$ | $a_2$               | $L_2$ | $F_2$ |
| 1            | 1                   | 1     | 4     | 0                   | 3     | 4     |
| 2            | 2                   | 1     | 5     | 3                   | 2     | 8     |
| 3            | 4                   | 2     | 9     | 5                   | 2     | 14    |
| 4            | 7                   | 2     | 13    | -                   | -     | -     |

Для отправки каждый цикл обслуживания выбирается пакет с наименьшим виртуальным временем конца обработки. Итоговая пропускная способность  $r_i$  рассчитывается на основе заданных весов по Формуле 5[17]

$$r_i = \frac{w_i}{\sum_{i=0}^N w_i} \cdot R \quad (5)$$

где

$w_i$  — вес, назначенный потоку  $i$ .

$N$  — количество потоков.

$R$  — полная пропускная способность канала.

Рассмотрим алгоритм WFQ на основе потока (Flow-based WFQ, FWFQ). Название алгоритма отсылает к методу классификации пакетов, при котором очередь выделяется для пакетов одного потока.[2] Схема планировщика представлена на рисунке 6.

В реализации FWFQ от Cisco используется алгоритм WFQ на основе порядкового номера пакета. Алгоритм поддерживает два счётчика: счётчик цикла, который определяет количество пройденных циклов побайтового планировщика (и равняется порядковому номеру последнего обслуженного пакета), и значение наибольшего порядкового номера пакета, поставленного в очередь потока. На основе этих значений высчитывается поряд-

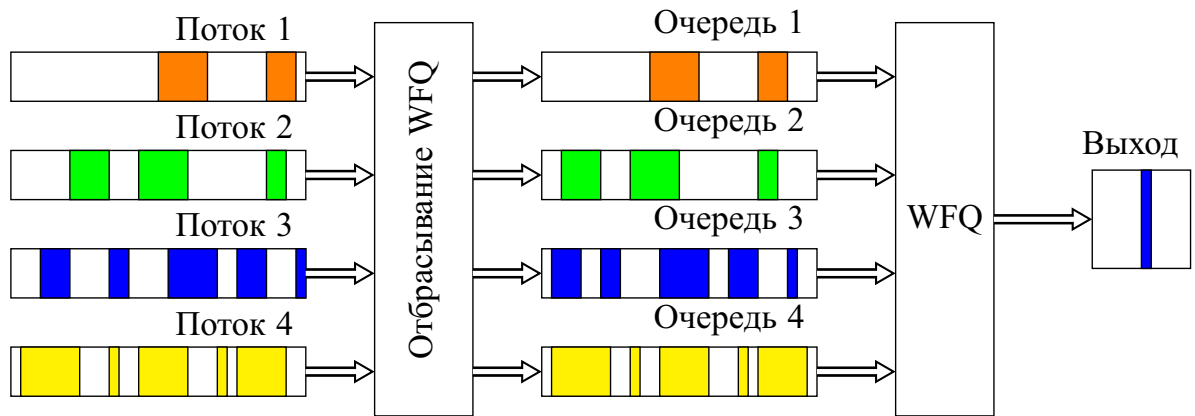


Рисунок 6 — Схема WFQ системы на основе потоков.

ковый номер пакета, пакет с минимальным значением покидает систему. Вычисление порядкового номера пакета вычисляется в зависимости от того, активный ли поток на момент прибытия нового пакета. Если поток был неактивным, то *порядковый номер пакета = размер пакета в байтах · вес + значение счётчика цикла на момент поступления пакета*; если поток активный, то *порядковый номер пакета = размер пакета в байтах · вес + значение наибольшего порядкового номера пакета в потоке*. Вес пакета строко зависит от его приоритета (определяется по соответствующему полю заголовка IP) и не может быть изменён.[2] Условная схема работы алгоритма представлена на рисунке 7.

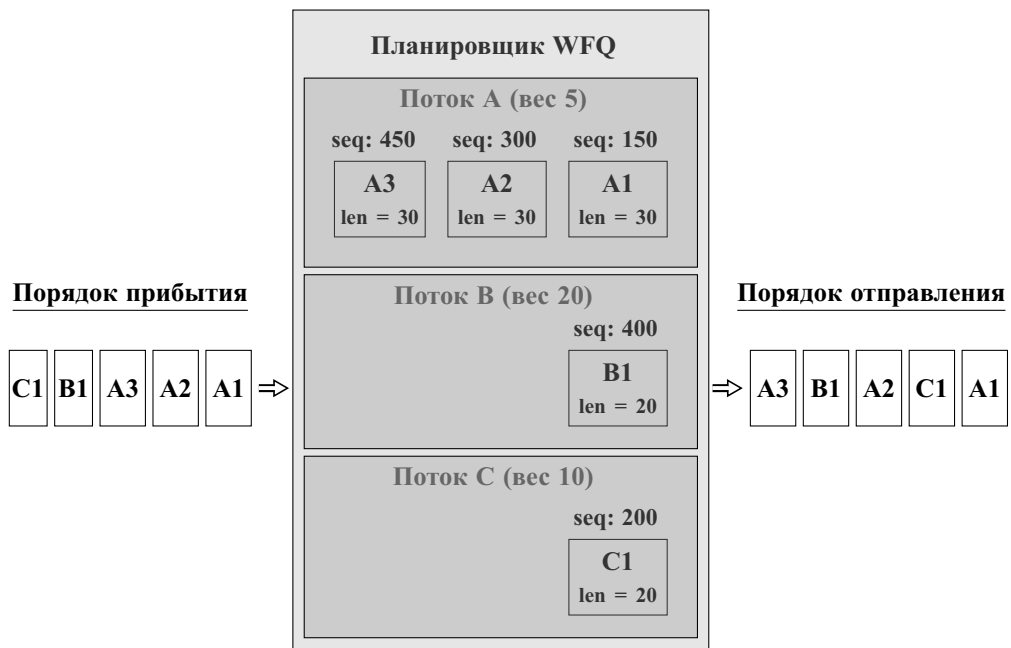


Рисунок 7 — Условная схема работы алгоритма WFQ на вычисления порядкового номера.



Планировщик не нарушает порядка обработки пакетов, принадлежащих одному потоку, даже в том случае, если они имеют различный приоритет. В целях планировки в WFQ длина очереди измеряется не в пакетах, а во времени, которое заняла бы передача всех пакетов в очереди.[2]

WFQ использует два метода отбрасывания пакетов: ранее (Early Dropping) и агрессивное (Aggressive Dropping) отбрасывания. Раннее отбрасывание срабатывает тогда, когда достигается congestive discard threshold (CDT); CDT – это количество пакетов, которые могут находиться в системе WFQ перед тем, как начнётся отбрасывание новых пакетов из самой длинной очереди; используется, чтобы начать отбрасывание пакетов из наиболее агрессивного потока, даже перед тем, как он достигнет предел hold queue out (HCO). HCO – это максимальное количество пакетов, которое может быть во всех выходящих очередях в интерфейсе в любое время; при достижении HCO срабатывает агрессивный режим отбрасывания. Алгоритм представлен на рисунке 8. [18]

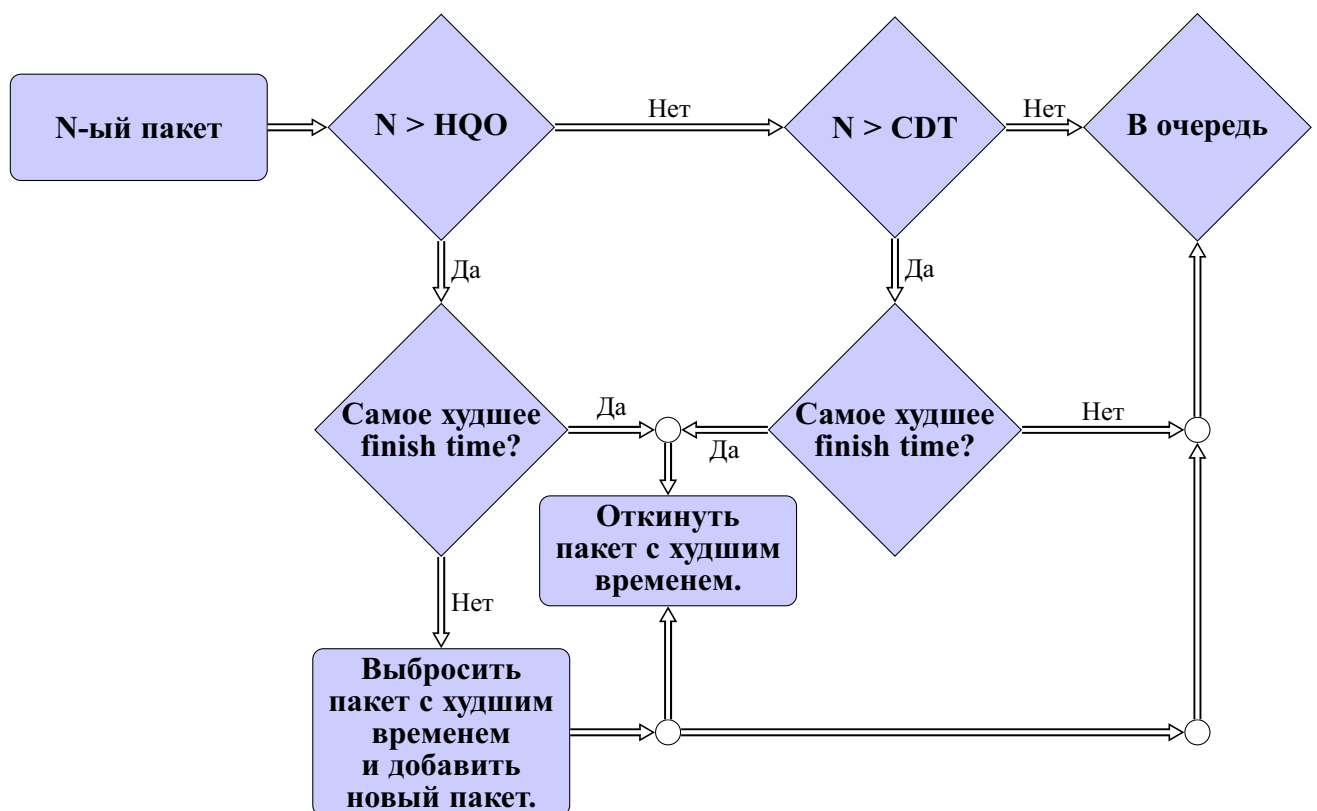


Рисунок 8 — Схема отбрасывания пакетов WFQ.

Преимущества WFQ.

- Простая конфигурация.

- Отбрасывание пакетов из более агрессивных потоков, что предотвращает перегрузки.
- Из-за честного обслуживания отсутствует проблема голодания потоков.

WFQ страдает от нескольких недостатков.

- Трафик не может регулироваться на основе достигнутых определённых классов обслуживания.
- Не поддерживает задание определённой пропускной способности для типа трафика.
- В Cisco системах WFQ поддерживается только на медленных каналах.[19]

Эти ограничения были исправлены CBWFQ.

### **1.7 Взвешенный алгоритм честного обслуживания очередей на основе классов**

CBWFQ (Class-based weighted fair queueing) – основанный на классах взвешенный алгоритм равномерного обслуживания очередей[2]; является расширением функциональности дисциплины обслуживания WFQ, основанной на потоках, для предоставления определяемых пользователями классов трафика.

Class-Based WFQ – это механизм, использующийся для гарантировании пропускной способности для класса. Для CBWFQ класс трафика определяется на основе заданных критериев соответствия: список контроля доступа (ACL), протокол, входящий интерфейс и т.п. Пакеты, удовлетворяющие критериям класса, составляют трафика для этого класса. Дисциплина позволяет задавать до 64-х пользовательских классов.

Схема дисциплины обслуживания представлена на рисунке 9.

После определения класса, ему назначаются характеристики, которые определяют политику очереди: пропускная способность, выделенная классу, максимальная длина очереди и так далее. Алгоритм CBWFQ позволяет явно указать требуемую минимальную полосу пропускания для каждого класса трафика. Полоса пропускания используется в качестве веса класса. Вес можно задать в абсолютной (опция `bandwidth`), в процентной (опция `bandwidth percent`) и в доле от оставшейся полосы пропускания (опция `bandwidth remaining percent`) величинах. Кроме пользовательских классов CBWFQ

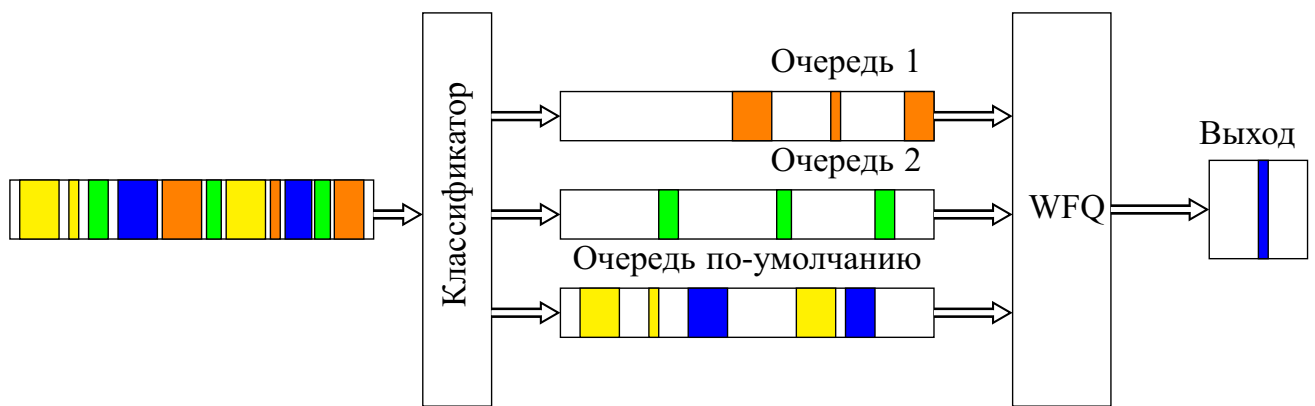


Рисунок 9 — Схема дисциплины обслуживания CBWFQ.

предоставляет стандартный класс (default class), в который попадает весь трафик, который не был классифицирован. В стандартном классе управление очередью может осуществляться с помощью алгоритмов FIFO и FQ (Fair Queueing). [16]

В случае переполнения очередей начинает работать алгоритм отбрасывания пакетов. В качестве политики отбрасывания пакетов по умолчанию используется отбрасывание конца очереди (Tail Drop), однако допускается сконфигурировать работу алгоритм взвешенного произвольного раннего обнаружения (Weighted Random Early Detection, WRED) для каждого класса.[16]

Преимущества:

- позволяет явно задать полосу пропускания для класса;
- позволяет создавать классы трафика и настраивать их в соответствии с требованиями;
- простая конфигурация вследствие небольшого числа параметров.[2][16]

Недостатки:

- нет поддержки работы с интерактивным трафиком (что исправляется в дисциплине обслуживания Low Latency Queueing (LLQ), которая является развитием CBWFQ);
- в Cisco реализации наблюдается ограничение на количество пользовательских классов (до 64-х классов);[2]
- отсутствие открытой реализации, что усложняет реализацию алгоритма в других системах и требует его полного воссоздания на основе имеющихся источников.

## 1.8 Выводы

Таблица 2 — Сравнительная таблица дисциплин обслуживания.

| Свойство                   | PQ   | CBQ | HTB  | HFSC  | FWFQ  | CBWFQ   |
|----------------------------|------|-----|------|-------|-------|---------|
| Метод планирования         | RR   | WRR | RR   | RT/LS | WFQ   | WFQ     |
| Отбрасывание               | TD   | TD  | TD   | TD    | ED/AD | TD/WRED |
| Честность (справедливость) | -    | -   | -    | -     | +     | +       |
| Разделение канала          | -    | +   | +    | +     | -     | -       |
| Решение проблемы голодания | -    | +   | +    | +     | +     | +       |
| Сложность реализации       | Низк | Выс | Сред | Выс   | Сред  | Сред    |
| Сложность конфигурации     | Низк | Выс | Сред | Выс   | Низк  | Низк    |
| Конфигурация классов       | -    | +   | +    | +     | -     | +       |
| Реализация в Linux         | +    | +   | +    | +     | -     | -       |

Каждая из рассмотренных ДО обладает своими достоинствами и недостатками. В Таблице 2 приведено сравнение основных элементов проанализированных дисциплин обслуживания.

Особое внимание следует уделить механизмам честного обслуживания и разделения канала. Оба метода служат решением проблемы голодания, однако используют разные по сложности подходы. Механизм разделения канала использует сложные вычисления и требует не только тщательной реализации и поддержки (что явно видно при исследовании исходного кода приведённых выше дисциплин CBQ, HTB и HFSC), но и кропотливой конфигурации. В то время, когда честное обслуживание не требует реализации сложных механизмов. У него есть свои недостатки: оно не позволяет построения сложных иерархических связей и разделения пропускной способности между классами одного уровня. Однако это требуется не во всех ситуациях; для простой конфигурации в нетривиальных, но не требующих тщательного контроля, случаях лучше всего

использовать несложные в конфигурации механизмы, дающие схожий результат. Поэтому реализация CBWFQ в ядре Linux целесообразна.

## 2 РЕАЛИЗАЦИЯ CLASS-BASED WFQ В ЯДРЕ LINUX

### 2.1 Описание устройства подсистемы планировки в ядре Linux

В операционной системе Linux дисциплина обслуживания, обозначаемая термином `qdisc`, используется для выбора пакетов из выходящей очереди для отправки на выходной интерфейс. Схема движения пакета приведена на Рисунке 10. Выходная очередь обозначена термином `egress`; именно на этом этапе следования пакета и работает механизм `qdisc`. [5]

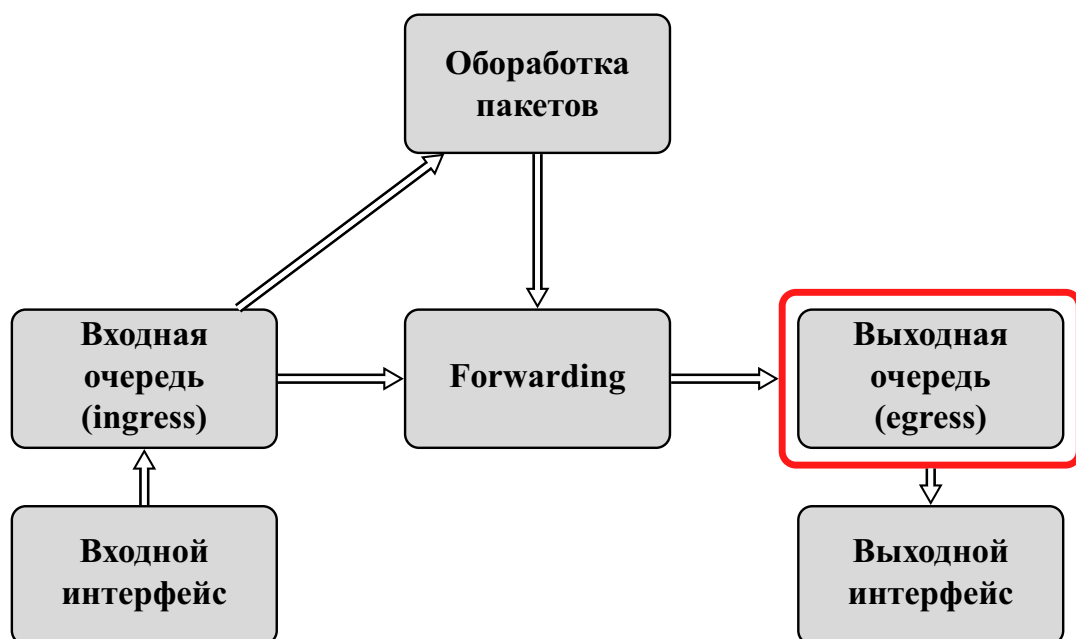


Рисунок 10 — Схема движения пакета в системе Linux[1]. Красным отмечена стадия, в которой работает механизм контроля качества обслуживания.

В общем случае, дисциплина обслуживания — это чёрный ящик, который может принимать поток пакетов и выпускать пакеты, когда устройство готово к отправке, в порядке и во время, определёнными спрятанным в ящике алгоритмом. В ядре Linux дисциплины обслуживания представляются в качестве модулей ядра, которые реализуют предоставляемый ядром интерфейс.

Linux поддерживает классовые и бесклассовые дисциплины обслуживания. Примером бесклассовой дисциплины служит `pfifo_fast`, классовой — `htb`. [5]

Классы представляют собой отдельные сущности в иерархии основной дисциплины. Если структура представляет собой дерево, то в классах-узлах мо-

гут содержаться фильтры, которые определяют пакет в нужный класс-потомок. В классах-листьях непосредственно располагаются очереди, которые управляются внутренней дисциплиной обслуживания. По умолчанию это `pfifo_fast`, но можно назначить другие.

Каждый интерфейс имеет корневую дисциплину, которой назначается идентификатор (`handle`), который используется для обращения к дисциплине. Этот идентификатор состоит из двух частей: мажорной (MAJ) и минорной (MIN); мажорная часть определяет родителя, минорная — непосредственно класс. На Рисунке 11 представлен пример иерархии.

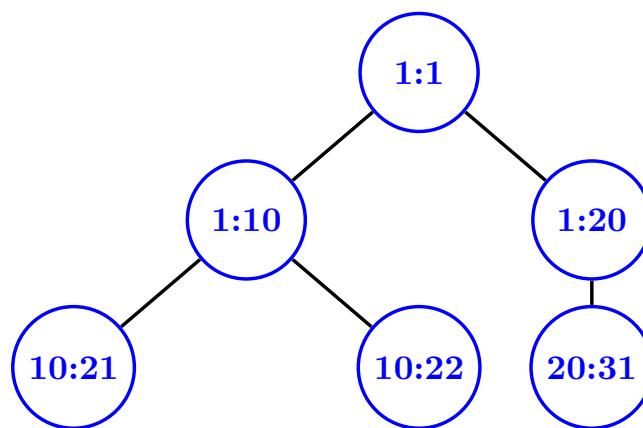


Рисунок 11 — Схема классовой иерархии с использованием идентификаторов MAJ:MIN

Идентификатор класса называется `classid` (к примеру, `1:10`), а идентификатор его родителя — `parenid` (`1:1` для классов `1:10` и `1:20`). По этим идентификаторам происходит поиск нужного класса внутри дисциплины.

Такая иерархия позволяет организовать гибкую систему классификации с набором классов и их подклассов, пакеты в которые назначаются фильтрами, которые предоставляются ядром.

## 2.2 Интерфейс управления трафиком

В Linux управление трафиком осуществляется с помощью подсистемы Traffic Control, которая предоставляет пользовательский интерфейс с помощью утилиты `tc`. `tc` — это пользовательская программа, которая позволяет настраивать дисциплины обслуживания в Linux. Она использует Netlink в качестве коммуникационного канала для взаимодействия между пользовательским пространством

и пространством ядра. tc добавляет новые дисциплины обслуживания, классы трафика, фильтры и предоставляет команды для управления всеми обозначенными объектами.[1]

tc предоставляет интерфейс для дисциплины обслуживания, представленный структурой `struct qdisc_util`, которая описывает функции для отправления команд и соответствующих параметров ядру и вывода сообщений о настройке дисциплины, списках классов и их настройке, а также статистику от ядра. Сообщение, помимо общей информации для подсистемы, содержит специфичную для дисциплины структуру с опциями, описываемую в заголовке ядра `pkt_sched.h`:

- структура `struct tc_cbwfq_glob` определяет глобальные настройки дисциплины обслуживания; структура передаётся по Netlink к модулю дисциплины при инициализации и изменении настроек;
- структура `struct tc_cbwfq_copt` определяет настройку класса и используется при добавлении или изменении класса.

Патч для заголовка в Приложении А.

Для назначения новой дисциплины обслуживания на интерфейс используется команда “tc qdisc add” системными параметрами (к примеру, название интерфейса), названием дисциплины и её локальными параметрами, которые определяются и обрабатываются в модуле дисциплины для утилиты tc. Для внесения изменений и удаления используются соответственно “tc change” и “tc delete”.

Опции для настройки дисциплины обслуживания:

- “bandwidth”— пропускная способность В канала;
- “default”— ключевое слово, определяющее, что далее пойдёт настройка класса по умолчанию; опции те же самые, что и при настройке класса.

Для классовых дисциплин используется команда “tc class” с под командами “add”, “change” и так далее. Классы обычно имеют параметры, отличные от параметров всей дисциплины обслуживания, поэтому нуждаются в отдельной структуре данных и функции обработчике.

Опции для настройки класса:

- “rate”— минимальная пропускная способность для класса; задаётся в единицах скорости (Mbps, Kbps, bps) или в процентах от размера канала (при использовании ключевого слова “percent”);



- “limit”— максимальное число пакетов в очереди.

Модуль для утилиты tc и ядра Linux, обеспечивающие взаимодействие между пользовательским пространством и дисциплиной представлен в Приложении Б.

## 2.3 Описание интерфейса

API ядра для подсистемы qdisc предоставляет две функции: `register_qdisc ( struct Qdisc_ops *ops)` и обратную — `unregister_qdisc ( struct Qdisc_ops *ops)`, которые регистрируют и разрегистрируют дисциплину обслуживания на интерфейсе. Важно отметить, что обе эти функции принимают в качестве аргумента структуру `struct Qdisc_ops`, которая явным образом идентифицирует дисциплину обслуживания в ядре.

Структура `struct Qdisc_ops` помимо метаинформации (в виде наименования дисциплины) содержит указатели на функции, которые должен реализовывать модуль дисциплины обслуживания для работы в ядре. Если не реализовать некоторые функции, то ядро в некоторых случаях попытается использовать функции по умолчанию, однако для особенно важных (к примеру, изменение конфигурации дисциплины или класса) сообщит пользователю, что операция не реализована.

Поля структуры `Qdisc_ops` представляют собой указатели на функции представленными ниже сигнатурами.

- `enqueue`

```
int enqueue(struct sk_buff *skb, struct Qdisc *sch, struct sk_buff **to_free);
```

Функция добавляет пакет в очередь. Если пакет был отброшен, функция возвращает код ошибки, говорящий о том, был отброшен пришедший пакет или иной, чье место занял новый.

- `dequeue`

```
struct sk_buff *dequeue(struct Qdisc *sch);
```

Функция, возвращающая пакет из очереди на отправку. Дисциплина может не передавать пакет при вызове этой функции по решению алгоритма, в таком случае вернув нулевой указатель; однако то же значение алгоритм

возвращает в случае, если очередь пуста, поэтому в таком случае дополнительно проверяется длина очереди.

– peek

```
struct sk_buff *peek( struct Qdisc *sch);
```

Функция возвращает пакет из очереди на отправку, не удаляя его из реальной очереди, как это делает функция dequeue.

– init

```
int init( struct Qdisc *sch, struct nlattr *arg);
```

Функция инициализирует вновь созданный экземпляр дисциплины обслуживания sch. Вторым аргументом функции является конфигурация дисциплины обслуживания, передаваемая в ядро с помощью подсистемы Netlink.

– change

```
int change( struct Qdisc *sch, struct nlattr *arg);
```

Функция изменяет текущие настройки дисциплины обслуживания.

– dump

```
int dump( struct Qdisc *sch, struct sk_buff *skb);
```

Функция отправляет по Netlink статистику дисциплины обслуживания.

Также структура содержит указатель на struct Qdisc\_class\_ops, которая описывает указатели функции исключительно для классовых дисциплин. Ниже приведены наиболее важные сигнатуры и их описания.

– find

```
unsigned long find( struct Qdisc *sch, u32 classid );
```

Функция возвращает приведённый к unsigned long адресс класса по его идентификатору ( classid ).

– change

```
int change( struct Qdisc *sch, u32 classid , u32 parentid , struct nlattr * attr , unsigned long *arg);
```

Функция используется для изменения и добавления новых классов в иерархию классов.

– tcf\_block, bind\_tcf, unbind\_tcf

В данном случае, описание сигнатур не даст какой-либо значимой инфор-

мации; практически для всех дисциплин обслуживания они идентичны. Эти функции предназначены для работы системы фильтрации.

– `dump_class`

```
int dump_class(struct Qdisc *sch, unsigned long cl, struct sk_buff *skb,
               struct tcmsg *tcm);
```

Функция предназначена для передачи по Netlink информации о классе и дополнительной статистики, собранной во время функционирования класса.

Для классовых дисциплин, помимо описанного, реализуют классификацию пакетов, которая определяет класс, куда попадает пакет. Классификация обычно выражается в функции `classify`, которая вызывается при добавлении пакета в очередь (функция `enqueue`) определяет, какому классу принадлежит пакет, и возвращает указатель на этот класс. Экземпляр структур для дисциплины обслуживания CBWFQ приведён в патче, представленном в Приложении В.

## 2.4 Алгоритм CBWFQ

Реализация CBWFQ требует:

- вычисление порядкового номера для каждого пакета в очередях и содержание глобального счётчика циклов;
- поддержку классов и классовых операций;
- фильтрацию для классификации трафика по классам.

### 2.4.1 Структуры хранения данных Class-Based WFQ

Обычно классовые дисциплины обслуживания содержат две основные структуры: для описания непосредственно дисциплины и для описания класса. Структура дисциплины содержит в себе данные, которые описывают всю дисциплину: это могут быть структура данных с классами (в виде списка или дерева), ограничения на очереди, статистика по всей дисциплине и так далее. Структура класса, соответственно, содержит непосредственно очередь и описывающие класс параметры.

Описание полей структуры дисциплины обслуживания `struct cbwfq_sched_data`.

- `struct Qdisc_class_hash clhash`  
Хэш-таблица для хранения классов.
  - `struct tcf_proto * filter_list` и `struct tcf_block *block`  
Структуры для хранения и обработки фильтров. Используются при выборе класса, в который поместить пакет.
  - `struct cbwfq_class *default_queue`  
Ссылка на очередь по умолчанию для быстрого доступа.
  - `enum cbwfq_rate_type rtype`  
Определяет тип, в котором указаны пропускная способность канала и пропускная способность для класса:
    - `TCA_CBWFQ_RT_BYTE` — ПС задана в байтах;
    - `TCA_CBWFQ_RT_PERCENT` — ПС задана в процентах.
- Используется для поддержания консистентности конфигурации.
- `u32 ifrate`  
Пропускная способность канала.
  - `u32 active_rate`  
Суммарная пропускная способность всех классов. Используется для определения состояния системы, при котором ни один поток не активный.
  - `u64 sch_sn`  
Счётчик циклов; используется для определения порядкового номера пакета, принадлежащий неактивному классу.
- Описание структуры полей класса.
- `struct Qdisc_class_common common`  
Структура, используемая для управления в хэш-таблице `clhash`. Содержит в себе идентификатор класса (`classid`) и метainформацию для таблицы.
  - `struct Qdisc *queue`  
Внутренняя дисциплина обслуживания, непосредственно содержащая пакеты. Настраивается на дисциплину `pfast_fifo`.
  - `u32 limit`  
Максимальное количество пакетов в очереди класса.

– u32 rate

Минимальна пропускная способность, выделенная классу.

– u64 cl\_sn

Значение последнего порядкового номера пакета в очереди класса; используется для определения порядкового номера в случае, когда класс активен.

– bool is\_active

Флаг активности класса.

Определение структур и реализация функций представлены в Приложении В.

### 2.4.2 Добавление пакета в очередь

Алгоритм добавления пакета обычно состоит из схожих действий: классификация и добавление в очередь, если есть место в очереди для пакета. Ниже приведен алгоритм на псевдо-языке.

```

1: function ENQUEUE(Q, pkt)
2:   c ← CLASSIFY(Q, pkt)
3:   if c.queue_len < c.limit then
4:     DROP(Q, pkt)
5:   else if Q.ENQUEUE(pkt) then
6:     
$$v1 \leftarrow \text{pkt.len} \cdot \frac{Q.\text{ifrate}}{c.\text{rate}}$$

7:     if cl не активен then
8:       c.sn ← Q.cycle + v1
9:     else
10:      c.sn ← c.sn + v1
11:    end if
12:    pkt.sn ← c.sn
13:  else
14:    DROP(Q, pkt)
15:  end if
16: end function

```

Сначала нужно классифицировать пакет в очередь. В строке 2 функция классификации определяет очередь, которой соответствует пакет, с помощью

заданных фильтров и возвращает указатель на класс. В строке 3 проверяем, достигла ли очередь предела, отбрасываем его, если достигла (в строке 4), иначе пытаемся добавить его в очередь (в строке 5). Если добавить пакет в очередь удалось, вычисляем его порядковый номер. Для начала вычисляется виртуальная длина пакета в строке 6. Далее в зависимости от того, была ли очередь активна, вычитываем текущий порядковый номер пакета (строки 9 и 11). Если же добавить в очередь не получилось, отбрасываем пакет (строка 14).

Блок-схема алгоритма приведена на рисунке 12.

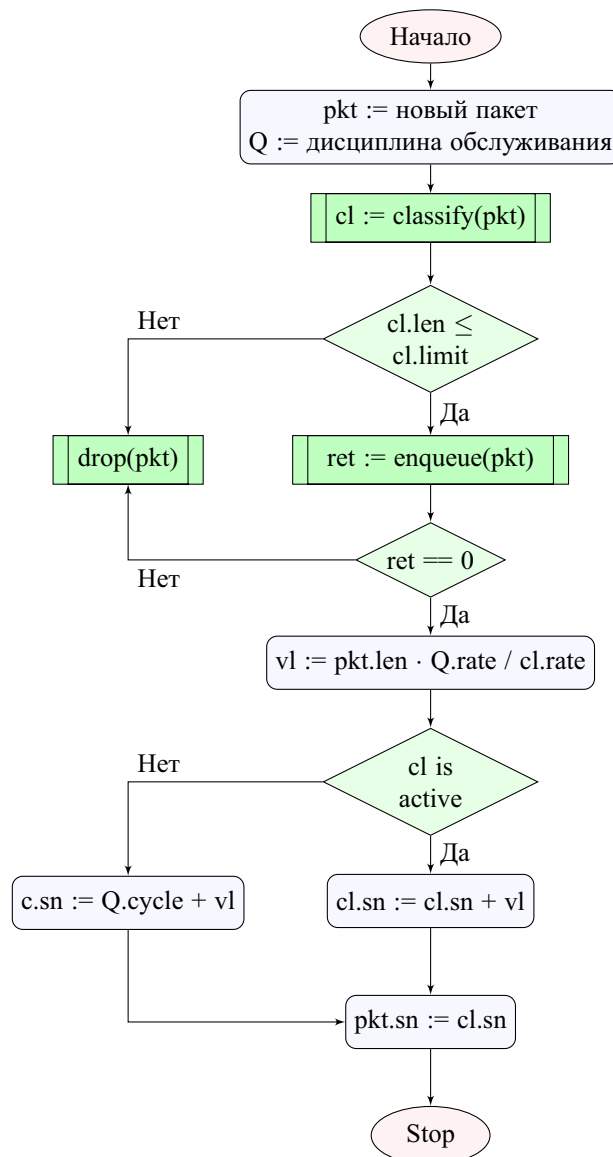


Рисунок 12 — Блок-схема алгоритма добавления пакета в очередь.

### 2.4.3 Удаление пакета из очереди

Функция удаления пакета из очереди непосредственно реализует планировщик WFQ.

```

1: function DEQUEUE(Q)
2:   C ← FIND_MIN(Q)
3:   if C is null then
4:     return null
5:   end if
6:   pkt ← C.QUEUE.DEQUEUE(C.queue)
7:   if c.queue.len == 0 then
8:     cl.sn ← 0
9:   end if
10:  if все классы не активны then
11:    Q.sn ← 0
12:  else
13:    Q.sn ← pkt.sn
14:  end if
15:  return pkt
16: end function

```

В первую очередь в строке 2 находим класс с наименьшим порядковым номером пакета в голове очереди. Если класс не был найден, значит все очереди пусты (строка 3), и алгоритм не может вернуть пакет. Иначе достаём пакет из очереди (строка 6). Проверяем, пуста ли очередь, в строке 8. Если пуста, то класс становится неактивным и счётчик сбрасывается (строка 11). В ином случае увеличиваем счётчик циклов всей дисциплины в строке 13. И возвращаем пакет. Блок-схема алгоритма приведена на рисунке 13.

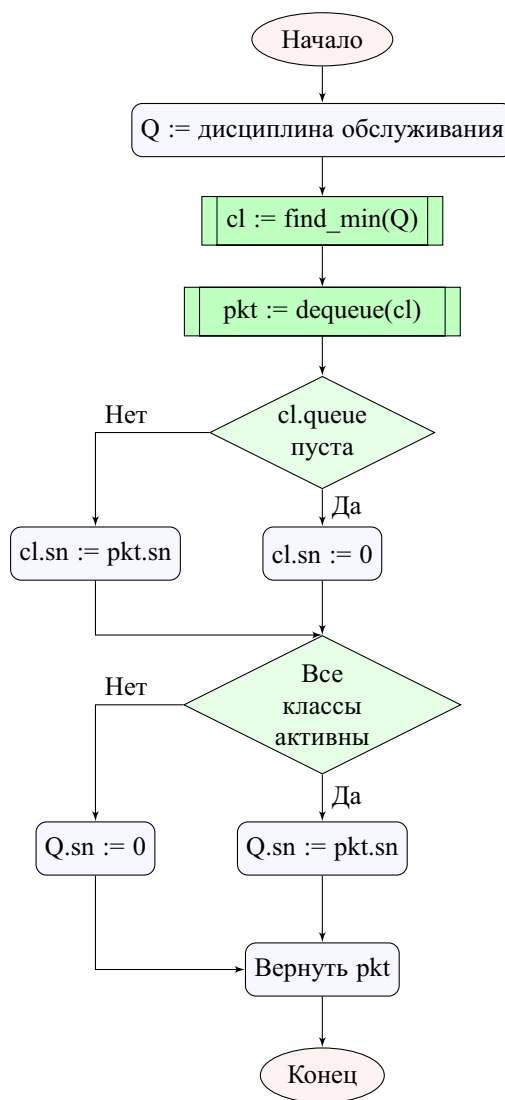


Рисунок 13 — Блок-схема алгоритма удаления пакета из очереди.



### 3 ТЕСТИРОВАНИЕ РАЗРАБОТАННОГО МОДУЛЯ ДИСЦИПЛИНЫ ОБСЛУЖИВАНИЯ CLASS-BASED WFQ

#### 3.1 Описание тестовой среды

Для тестирования модуля ядра была создана система виртуальных машин на основе системы эмуляции программного обеспечения QEMU. Схема тестовой среды представлена на Рисунке 14. Источником служит узел, от которого исходит трафик; таблицы маршрутизации настроены таким образом, чтобы весь трафик, который должен попасть на узел-цель шёл через промежуточный узел, на котором настроена тестируемая дисциплина обслуживания CBWFQ.

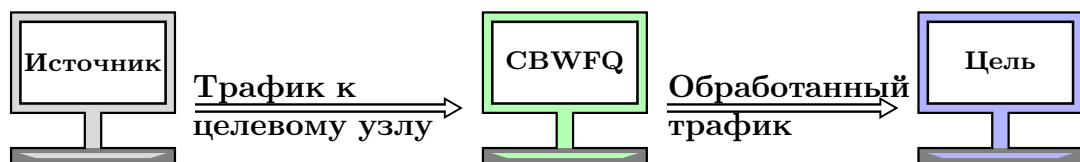


Рисунок 14 — Схема тестовой среды.

В Листинге 1 приведена система настройки дисциплины на промежуточном узле. Переменная окружения TESTPORT содержит в себе номер порта, с

---

```

1
2 tc qdisc add dev $IFACE root handle 1: cbwfq bandwidth
   100Mbps\
3     default rate 5Mbps
4 tc class add dev $IFACE parent 1: classid 1:2 cbwfq rate
   25Mbps
5 tc class add dev $IFACE parent 1: classid 1:3 cbwfq rate
   70Mbps
6 tc filter add dev ens4 parent 1:1 protocol ip u32 match \
7     ip dport $TESTPORT1 0xffff flowid 1:2
8 tc filter add dev ens4 parent 1:0 protocol ip u32 match \
9     ip dport $TESTPORT2 0xffff flowid 1:3
  
```

---

Листинг 1 — Список команд для конфигурации дисциплины обслуживания CBWFQ.

которого будет отправлен трафик на сервер. При добавлении дисциплины на интерфейс необходимо указать пропускную способность канала; размер очереди по умолчанию назначается опционально. При конфигурации класса во второй строке происходит назначения полосы пропускания для класса в процентах (возможно также назначение в bps). Так как CBWFQ всегда имеет класс по умолчанию, то при добавлении нового класса классу по умолчанию будет доставаться оставшаяся пропускная способность. В третьей строке происходит назначение фильтра, который будет направлять трафик с исходным портом TESTPORT в очередь класса 1:2. На один класс можно назначит множество фильтров (Linux предоставляет гибкие возможности по настройке фильтрации); это не контролируется непосредственно дисциплиной и находится в компетенции пользователя.

### 3.2 Анализ точности выделения канала при конкурирующем трафике

Первый эксперимент заключается в исследовании разделения канала между двумя потоками трафика, приходящими со скоростью, больше скорости канала на промежуточном узле.

С использованием описанной конфигурации и указанных в Листингах 2 и 3 команд произвелось десять испытаний в рамках эксперимента (схема эксперимента представлена на Рисунке 15).

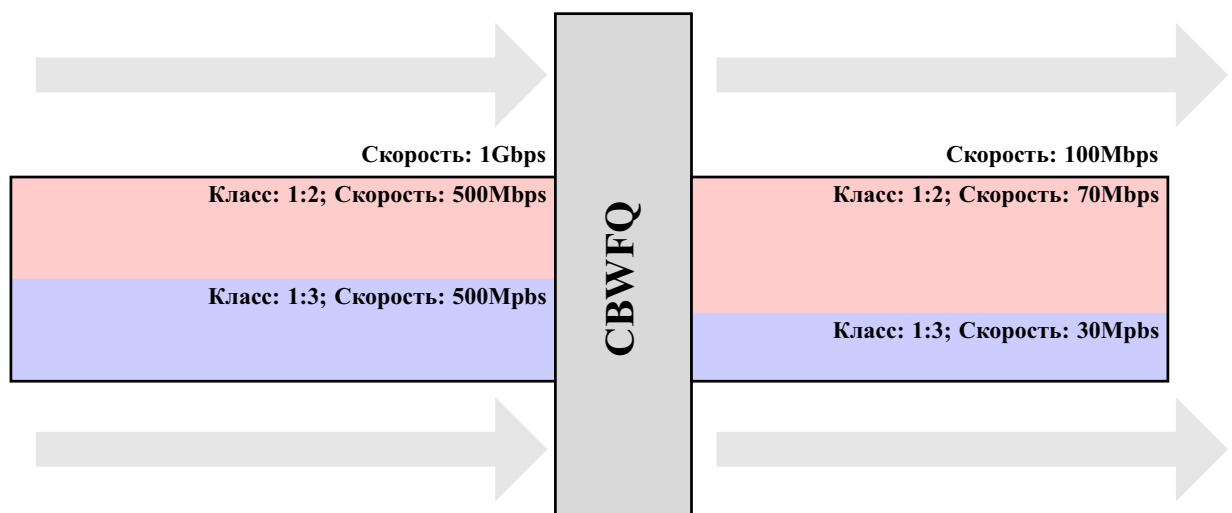


Рисунок 15 — Условная схема эксперимента 1.

В течение экспериментов собирались отчёты от утилиты `iperf` со стороны сервера в течение 90-ти секунд. Отчёты представляют собой таблицу с полями:

---

```
iperf3 -c $SERVERIP -p $TESRPORT1 -b 500M -u -t 90
iperf3 -c $SERVERIP -p $TESRPORT2 -b 500M -u -t 90
```

---

Листинг 2 — Команда iperf на узле-источнике (клиентская сторона).

---

```
iperf3 -s -p $TESTPOR1 --logfile class2 "$EXPNUM".log
iperf3 -s -p $TESTPOR2 --logfile class3 "$EXPNUM".log
```

---

Листинг 3 — Команда iperf на узле-цели (серверная сторона).

временной интервал (в секундах), количество переданных данных и пропускная способность. На каждый временной интервал таблица содержит описанную информацию для двух потоков.

Результаты эксперимента можно наблюдать на Рисунке 16. Доля пропускной способности для первого и второго класса соответствует минимальной сконфигурированной доле, выделенной этим классам (или их весу  $\frac{70\text{Mbps}}{100\text{Mbps}} = 0.7$  и  $\frac{25\text{Mbps}}{100\text{Mbps}} = 0.25$  соответственно).

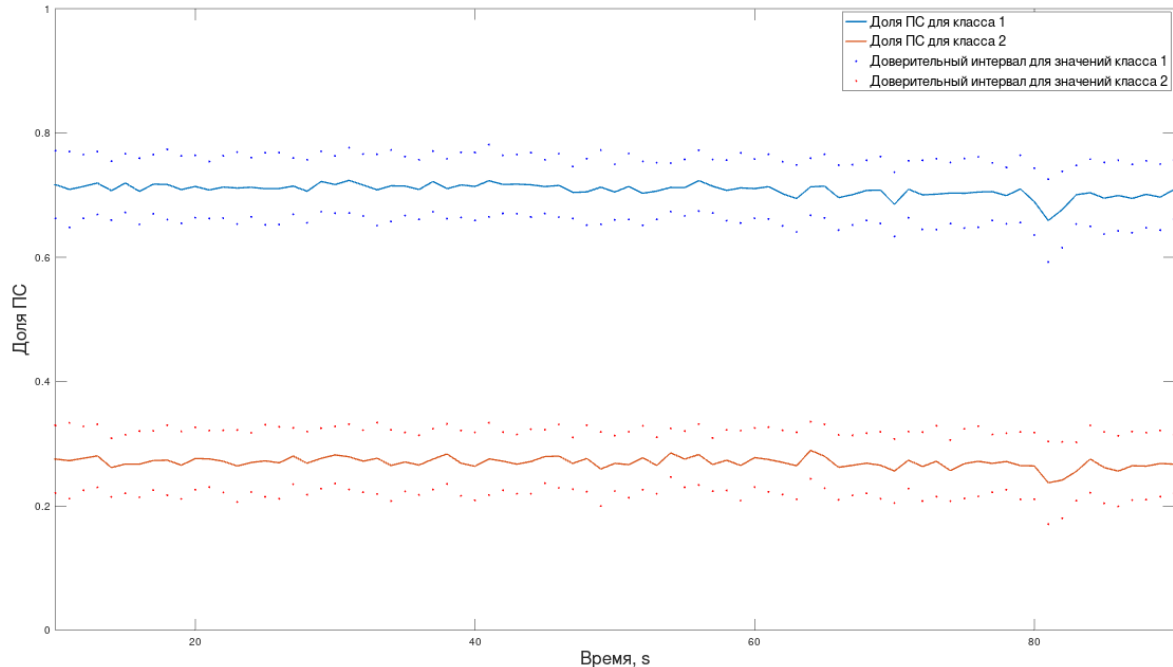


Рисунок 16 — График распределения доли пропускной способности (ПС) по типам трафика в течение времени. Среднее значение процента ПС для класса 1:  $71 \pm 3\%$  ( $P = 0.95$ ). Среднее значение процента ПС для класса 2:  $27 \pm 2\%$  ( $P = 0.95$ ).

### 3.3 Анализ точности выделения канала при независимом трафике

Второй эксперимент заключается в исследовании разделения канала между двумя потоками трафика, приходящими со скоростью, суммарно меньше скорости канала и соответствующие конфигурационным параметрам.

С использованием описанной конфигурации и указанных в Листингах 4 и 5 команд произвелось десять испытаний в рамках эксперимента (схема эксперимента представлена на Рисунке 17).

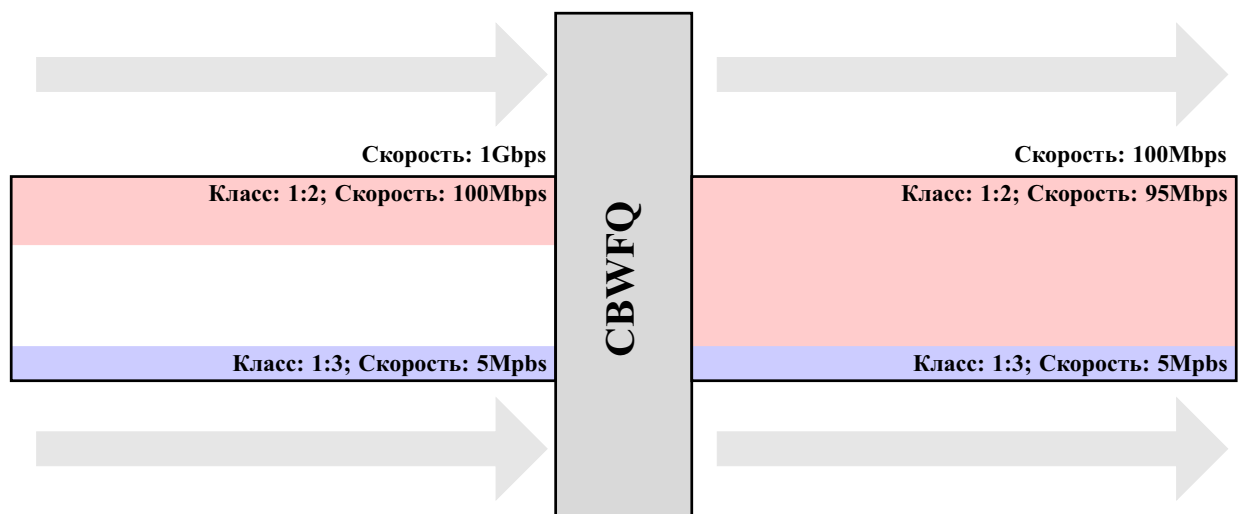


Рисунок 17 — Условная схема эксперимента 2.

В течение экспериментов собирались отчёты от утилиты `iperf` со стороны сервера в течение 90-ти секунд. Отчёты представляют собой таблицу с полями: временной интервал (в секундах), количество переданных данных и пропускная способность. На каждый временной интервал таблица содержит описанную информацию для двух потоков.

---

```
iperf3 -c $SERVERIP -p $TESRPORT1 -b 100M -u -t 90
iperf3 -c $SERVERIP -p $TESRPORT2 -b 5M -u -t 90
```

---

Листинг 4 — Команда `iperf` на узле-источнике (клиентская сторона).

Результаты эксперимента можно наблюдать на Рисунке 18. В данном случае второй класс получает всю требуемую ему пропускную способность; в этом факте и состоит отсутствие конкуренции: второй поток не пытается занять боль-

---

```
iperf3 -s -p $TESTPOR1 -- logfile class2 "$EXPNUM".log
iperf3 -s -p $TESTPOR2 -- logfile class3 "$EXPNUM".log
```

---

Листинг 5 — Команда iperf на узле-цели (серверная сторона).

ше, чем ему назначенно. Поток первого класса занимает оставшуюся свободную часть канала.

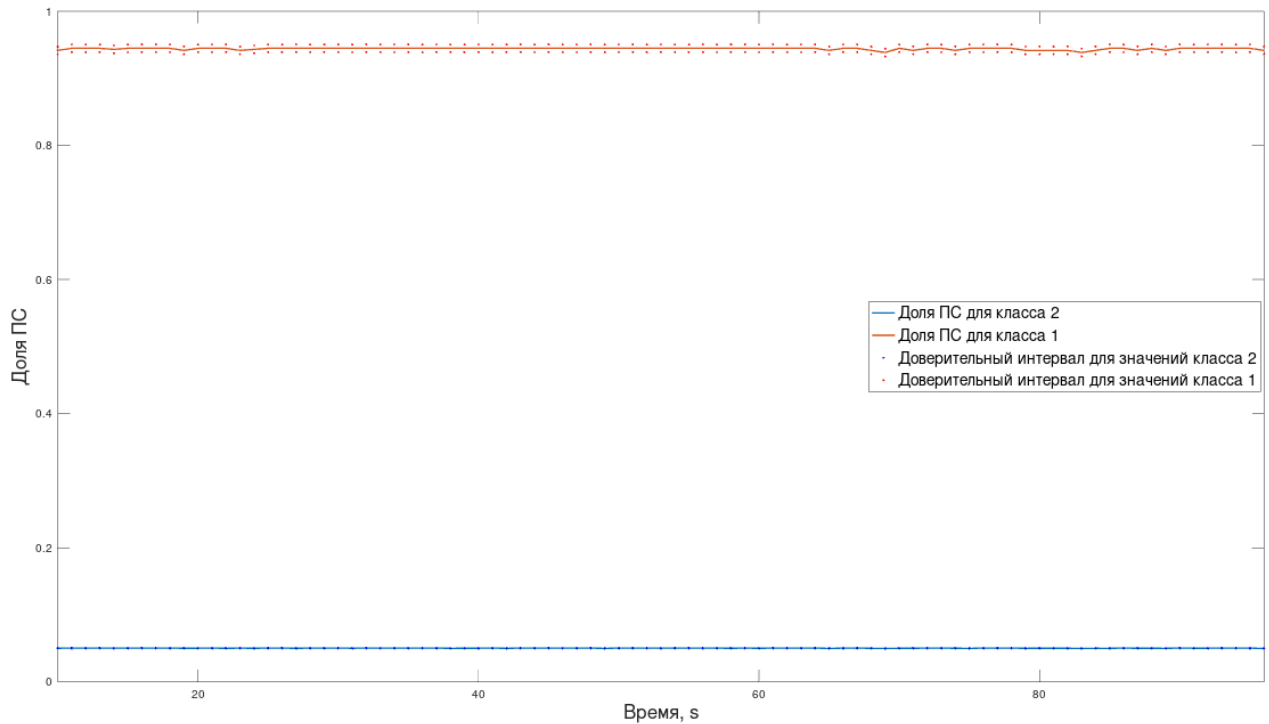


Рисунок 18 — График распределения доли пропускной способности (ПС) по типам трафика в течение времени. Среднее значение процента ПС для класса 1:  $94 \pm 0.5\%$  ( $P = 0.95$ ). Среднее значение процента ПС для класса 2:  $5 \pm 0.07\%$  ( $P = 0.95$ )

## ЗАКЛЮЧЕНИЕ

В результате выполнения работы были достигнуты следующие цели:

- проведён сравнительный анализ дисциплин обслуживания PQ, CBQ, HTB, HFSC, FWFQ, CBWFQ; дано краткое описание алгоритмов и приведены их слабые и сильные стороны;
- проведено исследование архитектуры подсистемы контроля качества обслуживания ядра Linux, описаны механизмы работы подсистемы;
- реализован и протестирован модуль дисциплины обслуживания Class-Based WFQ для ядра Linux;
- реализован модуль для утилиты tc для взаимодействия с модулем дисциплины обслуживания.

В дальнейшем работу можно развить в следующих направлениях.

1. Реализация взвешенного алгоритма раннего обнаружения (WRED) для возможности конфигурации политики отбрасывания для модуля CBWFQ;
2. Доработка работы до дисциплины Low-Latency Queuing (LLQ), которая совершенствует дисциплину CBWFQ с помощью добавления приоритетных очередей, использующиеся для чувствительного к задержкам трафика.

## СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

- [1] Sameer Seth. *TCP/IP Architecture, Design, and Implementation in Linux*. / Sameer Seth, M. Ajaykumar Venkatesulu – Wiley-IEEE Computer Society Press, 2008. – 772 p.
- [2] Вегешна Ш. *Качество обслуживания в сетях IP*. М.: Издательский дом Вильямс, 2003. - 368 с.
- [3] Алиев Т.И. *Основы моделирования дискретных систем*. – СПб.: СПбГУ ИТМО, 2009.
- [4] Иван Песин. Повесть о linux и управлении трафиком, 2003. [Электронный ресурс] – Режим доступа: <http://computerlib.narod.ru/html/traffic.htm> (дата обращения 10.04.2018).
- [5] Bert Hubert. *Linux Advanced Routing and Traffic Control, 2012*. [Электронный ресурс] – Режим доступа: [lartc.org/lartc.pdf](http://lartc.org/lartc.pdf) (дата обращения 19.03.2018).
- [6] Davide Astuti. Packet handling. // Seminar on Transport of Multimedia Streams in Wireless Internet, 2003.
- [7] Alexey N. Kuznetsov. tc-prio (8), . [Электронный ресурс] // Linux Man Pages – Режим доступа: <https://www.systutorials.com/docs/linux/man/8-tc-prio/> (дата обращения 10.04.2018).
- [8] Chuck Semeria. Supporting differentiated service classes, 2001. [Электронный ресурс] // Режим доступа: <https://pdfs.semanticscholar.org/dd1f/27c4b1e0b5395d67520e65737c9835a7bced.pdf> (дата обращения 30.03.2018).
- [9] Alexey N. Kuznetsov. tc-cbq (8), . [Электронный ресурс] // Linux Man Pages – Режим доступа: <https://www.systutorials.com/docs/linux/man/8-tc-cbq/> (дата обращения 10.04.2018).
- [10] Van Jacobson Sally Floyd. Link-sharing and resource management models for packet networks. // IEEE/ACM Transactions on Networking, Vol. 3 No. 4, 1995.

- [11] Martin Devera. Htb linux queuing discipline manual - user guide, 2002, . [Электронный ресурс] Режим доступа: <http://luxik.cdi.cz/~devik/qos/htb/manual/userg.htm> (дата обращения 21.04.2018).
- [12] Martin Devera. tc-htb (8), . [Электронный ресурс] // Linux Man Pages – Режим доступа: <https://www.systutorials.com/docs/linux/man/8-tc-htb/> (дата обращения 11.04.2018).
- [13] T.S. Eugene Ng Ion Stoica, Hui Zhang. A hierarchical fair service curve algorithm for link-sharing, real-time and priority service, 1997. [Электронный ресурс] – Режим доступа: <https://www.trash.net/~kaber/hfsc/SIGCOM97.pdf> (дата обращения 13.03.2018).
- [14] Michal Soltys. tc-hfsc (7). [Электронный ресурс] // Linux Man Pages – Режим доступа: <https://www.systutorials.com/docs/linux/man/7-tc-hfsc/> (дата обращения 11.04.2018).
- [15] Linux-tc-notes. notes on the linux traffic control engine, 2014. [Электронный ресурс] – Режим доступа: [http://linux-tc-notes.sourceforge.net/tc/doc/sch\\_hfsc.txt](http://linux-tc-notes.sourceforge.net/tc/doc/sch_hfsc.txt) (дата обращения 23.03.2018).
- [16] Cisco IOS Quality of Service Solutions Configuration Guide, Release 12.2 [Электронный ресурс] – 2009. – Режим доступа: [https://www.cisco.com/c/en/us/td/docs/ios/qos/configuration/guide/12\\_2sr/qos\\_12\\_2sr\\_book.pdf](https://www.cisco.com/c/en/us/td/docs/ios/qos/configuration/guide/12_2sr/qos_12_2sr_book.pdf) (дата обращения 18.04.2018).
- [17] Abhay K. Parekh. A generalized processor sharing approach to flow control in integrated services networks: The single-node case. / Abhay K. Parekh, Robert G. Gallager // IEEE/ACM Transactions on Networking, Vol. 1, No. 3, 1993.
- [18] Quality of service, part 10 – weighted fair queuing, 2010. [Электронный ресурс] – Режим доступа: <http://blog.globalknowledge.com/2010/02/12/quality-of-service-part-10-weighted-fair-queuing/> (дата обращения 04.03.2018).
- [19] Aaron Blachunas. Qos and queueing, 2010. [Электронный ресурс]



- Режим доступа: [http://www.routeralley.com/guides/qos\\_queueing.pdf](http://www.routeralley.com/guides/qos_queueing.pdf) (дата обращения 09.04.2018).

## **СПИСОК СОКРАЩЕНИЙ И УСЛОВНЫХ ОБОЗНАЧЕНИЙ**

- ДО – Дисциплина обслуживания.
- ПС – пропускная способность.
- CBWFQ – Class Based Weighted Fair Queueing.
- WFQ – Weighted Fair Queueing.
- PQ – Priority Queueing.
- CBQ – Class Based Queueing.
- HTB – Hierarchical Token Bucket.
- HFSC – Hierarchical Fair-Service Curve

## ПРИЛОЖЕНИЕ А

```

1  937,987d936
2  < enum {
3      TCA_CBWFQ_UNSPEC,
4      TCA_CBWFQ_PARAMS,
5      TCA_CBWFQ_INIT,
6      __TCA_CBWFQ_MAX
7  < };
8  < #define TCA_CBWFQ_MAX (__TCA_CBWFQ_MAX - 1)
9  <
10 < enum cbwfq_rate_type {
11     TCA_CBWFQ_RT_BYTE,
12     TCA_CBWFQ_RT_PERCENT
13 < };
14 <
15 < struct tc_cbwfq_glob {
16     __u32 cbwfq_gl_default_limit;
17     __u32 cbwfq_gl_default_rate;
18     __u32 cbwfq_gl_total_rate;
19     enum cbwfq_rate_type cbwfq_gl_rate_type;
20 < };
21 <
22 < struct tc_cbwfq_copt {
23     __u32 cbwfq_cl_rate;
24     __u32 cbwfq_cl_limit;
25     enum cbwfq_rate_type cbwfq_cl_rate_type;
26 < };
27 <

```

Листинг 6 — Патч для заголовочного файла pkt\_sched.h.

## ПРИЛОЖЕНИЕ Б

```

1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <unistd.h>
4  #include <syslog.h>
5  #include <fcntl.h>
6  #include <sys/socket.h>
7  #include <netinet/in.h>
8  #include <arpa/inet.h>
9  #include <string.h>
10
11 #include "utils.h"
12 #include "tc_util.h"
13
14 static void explain(void)
15 {
16     fprintf(stderr,
17 "Usage: ... qdisc add .. cbwfq bandwidth B default rate R [limit L] \n"
18 "\tbandwidth          bandwidth of the link (Mbps, Kbps, bps)\n"
19 "\tdefault            configuration for default class (see description below)\n"
20 "    \n"
21 "... class add ... cbwfq rate R [limit L]\n"
22 "\trate R [percent]   rate of the class in Kbit; use 'percent' to declare in\n"
23 "    percent\n"
24 "\tlimit              max queue length (in packets)\n"
25 ");
26 }
27
28 static void explain1(char *arg)
29 {
30     fprintf(stderr, "Illegal \"%s\"\n", arg);
31     explain();
32 }
33
34 static int cbwfq_parse_opt(struct qdisc_util *qu, int argc, char **argv,
35                          struct nlmsg_hdr *n)
36 {
37     struct tc_cbwfq_glob opt;
38     struct rtattr *tail;
39
40     memset(&opt, 0, sizeof(opt));
41     while (argc > 0) {
42         if (matches(*argv, "default") == 0) {
43             NEXT_ARG();
44             if (matches(*argv, "rate") == 0) {
45                 NEXT_ARG();
46                 if (get_rate(&opt.cbwfq_gl_default_rate, *argv)) {
47                     explain1("rate");
48                     return -1;
49                 }
50             }
51             argv++; argc--;
52             if (argc <= 0) {
53                 break;
54             }
55             if (matches(*argv, "percent") == 0) {
56                 opt.cbwfq_gl_rate_type = TCA_CBWFQ_RT_PERCENT;
57             }
58         }
59     }
60 }

```

```

54         argv -- ;
55         if (get_u32(&opt.cbwfq_gl_default_rate , *argv , 10)) {
56             explain1("percent");
57             return -1;
58         }
59         argv++;
60     } else {
61         opt.cbwfq_gl_rate_type = TCA_CBWFQ_RT_BYTE;
62     }
63 }
64 if (matches(*argv , "limit") == 0) {
65     NEXT_ARG();
66     if (get_u32(&opt.cbwfq_gl_default_limit , *argv , 10)) {
67         explain1("limit");
68         return -1;
69     }
70 } else {
71     fprintf(stderr , "Unknown default parameter: \"%s\".\n" , *argv
72 );
73     explain();
74     return -1;
75 }
76 } else if (matches(*argv , "bandwidth") == 0) {
77     NEXT_ARG();
78     if (get_rate(&opt.cbwfq_gl_total_rate , *argv)) {
79         explain1("bandwidth");
80         return -1;
81     }
82 } else {
83     fprintf(stderr , "What is \"%s\"?\n" , *argv);
84     explain();
85     return -1;
86 }
87 }
88
89 if (opt.cbwfq_gl_total_rate <= 0) {
90     fprintf(stderr , "Bandwidth must be set!\n");
91     return -1;
92 }
93
94 if (opt.cbwfq_gl_default_rate <= 0) {
95     fprintf(stderr , "Default rate must be set!\n");
96     return -1;
97 }
98
99 tail = NLMSG_TAIL(n);
100 addattr_l(n, 1024, TCA_OPTIONS, NULL, 0);
101 addattr_l(n, 2024, TCA_CBWFQ_INIT, &opt, NLMSG_ALIGN(sizeof(opt)));
102 tail->rta_len = (void *) NLMSG_TAIL(n) - (void *) tail;
103 return 0;
104 }
105
106 static int cbwfq_parse_class_opt(struct qdisc_util *qu, int argc, char **argv
107 ,
108                                struct nlmsg_hdr *n)
109 {
110     struct tc_cbwfq_copt opt;
111     struct rtattr *tail;

```

```

112     memset(&opt, 0, sizeof(opt));
113     while (argc > 0) {
114         if (matches(*argv, "rate") == 0) {
115             NEXT_ARG();
116             opt.cbwfq_cl_rate_type = TCA_CBWFQ_RT_BYTE;
117             if (get_rate(&opt.cbwfq_cl_rate, *argv)) {
118                 explain1("rate");
119                 return -1;
120             }
121
122             argv++; argc--;
123             if (argc <= 0) {
124                 break;
125             }
126
127             if (matches(*argv, "percent") == 0) {
128                 opt.cbwfq_cl_rate_type = TCA_CBWFQ_RT_PERCENT;
129                 argv--;
130                 if (get_u32(&opt.cbwfq_cl_rate, *argv, 10)) {
131                     explain1("percent");
132                     return -1;
133                 }
134                 argv++;
135             } else {
136                 fprintf(stderr, "What is \"%s\"?\n", *argv);
137                 explain();
138                 return -1;
139             }
140         } else if (matches(*argv, "limit") == 0) {
141             NEXT_ARG();
142             if (get_u32(&opt.cbwfq_cl_limit, *argv, 10)) {
143                 explain1("limit");
144                 return -1;
145             }
146         } else {
147             fprintf(stderr, "What is \"%s\"?\n", *argv);
148             explain();
149             return -1;
150         }
151         argc--; argv++;
152     }
153
154     if (opt.cbwfq_cl_rate <= 0) {
155         fprintf(stderr, "Rate must be set!\n");
156         explain();
157         return -1;
158     }
159
160     tail = NLMSG_TAIL(n);
161     addattr_l(n, 1024, TCA_OPTIONS, NULL, 0);
162     addattr_l(n, 1024, TCA_CBWFQ_PARAMS, &opt, sizeof(opt));
163     tail->rta_len = (void *) NLMSG_TAIL(n) - (void *) tail;
164     return 0;
165 }
166
167 int cbwfq_print_opt(struct qdisc_util *qu, FILE *f, struct rtattr *opt)
168 {
169     struct tc_cbwfq_glob *qopt = NULL;
170     struct rtattr *tb[TCA_CBWFQ_MAX+1];

```

```

172
173     if (opt == NULL) {
174         return 0;
175     }
176
177     if (parse_rtattr_nested(tb, TCA_CBWFQ_MAX, opt)) {
178         return -1;
179     }
180
181     if (tb[TCA_CBWFQ_INIT] == NULL) {
182         return -1;
183     }
184
185     if (RTA_PAYLOAD(tb[TCA_CBWFQ_INIT]) < sizeof(*opt)) {
186         fprintf(stderr, "qdisc opt is too short\n");
187     } else {
188         qopt = RTA_DATA(tb[TCA_CBWFQ_INIT]);
189     }
190
191     if (qopt != NULL) {
192         fprintf(f, "total rate %d ", qopt->cbwfq_gl_total_rate);
193     }
194     return 0;
195 }
196
197 static int cbwfq_print_copt(struct qdisc_util *qu, FILE *f, struct rtattr *
198     opt)
199 {
200     struct rtattr *tb[TCA_CBWFQ_MAX+1];
201     struct tc_cbwfq_copt *copt = NULL;
202
203     if (opt == NULL)
204         return 0;
205
206     if (parse_rtattr_nested(tb, TCA_CBWFQ_MAX, opt))
207         return -1;
208
209     if (tb[TCA_CBWFQ_PARAMS] != NULL) {
210         if (RTA_PAYLOAD(tb[TCA_CBWFQ_PARAMS]) < sizeof(*opt))
211             fprintf(stderr, "CBWFQ: class opt is too short\n");
212         else
213             copt = RTA_DATA(tb[TCA_CBWFQ_PARAMS]);
214     }
215
216     if (copt) {
217         fprintf(f, "limit %d rate %d ", copt->cbwfq_cl_limit,
218             copt->cbwfq_cl_rate);
219     }
220
221     return 0;
222 }
223
224 struct qdisc_util cbwfq_qdisc_util = {
225     .id = "cbwfq",
226     .parse_copt = cbwfq_parse_class_opt,
227     .parse_qopt = cbwfq_parse_opt,
228     .print_qopt = cbwfq_print_opt,
229     .print_copt = cbwfq_print_copt,
230 };

```

Листинг 7 — Модуль CBWFQ для утилиты tc.

## ПРИЛОЖЕНИЕ В

```

1  #include <linux/module.h>
2  #include <linux/slab.h>
3  #include <linux/types.h>
4  #include <linux/kernel.h>
5  #include <linux/string.h>
6  #include <linux/errno.h>
7  #include <linux/skbuff.h>
8  #include <net/netlink.h>
9  #include <net/pkt_sched.h>
10 #include <net/pkt_cls.h>
11
12 #include <linux/list.h>
13
14 #define MAX(a, b) ((a) > (b) ? (a) : (b))
15
16 #define DEFAULT_CL_ID      65537
17
18 /**
19  * cbwfq_class -- class description
20  * @common      Common qdisc data. Used in hash-table.
21  * @queue       Class queue.
22  *
23  * @limit       Max amount of packets.
24  * @rate       Assigned rate.
25  *
26  * @cl_sn      Sequence number of the last enqueued packet.
27  *
28  * @is_active   Set if class is in active state (transmit packets).
29  */
30 struct cbwfq_class {
31     struct Qdisc_class_common common;
32     struct Qdisc *queue;
33
34     u64 limit;
35     u64 rate;
36
37     u64 cl_sn;
38
39     bool is_active;
40 };
41
42 /**
43  * cbwfq_sched_data -- scheduler data
44  *
45  * @clhash      Hash table of classes.
46  *
47  * @filter_list List of attached filters.
48  * @block      Field used for filters to work.
49  *
50  * @default_queue Default class with the default queue.
51  *
52  * @ifrate      Total rate of the link.
53  * @active_rate Rate of all active classes. Used to determine idle.
54  *
55  * @sch_sn      Sequence number of the last dequeued packet; cycle number.

```



```

56  */
57  struct cbwfq_sched_data {
58      struct Qdisc_class_hash clhash;
59
60      struct tcf_proto __rcu *filter_list;
61      struct tcf_block *block;
62
63      struct cbwfq_class *default_queue;
64
65      enum cbwfq_rate_type rtype;
66      u64 ifrate;
67
68      u32 active_rate;
69
70      u64 sch_sn;
71  };
72
73
74  /* For parsing netlink messages. */
75  static const struct nla_policy cbwfq_policy[TCA_CBWFQ_MAX + 1] = {
76      [TCA_CBWFQ_PARAMS] = { .len = sizeof(struct tc_cbwfq_copt) },
77      [TCA_CBWFQ_INIT]   = { .len = sizeof(struct tc_cbwfq_glob) },
78  };
79
80
81  /**
82   * Add class to the hash table.
83   *
84   * @comment Class is allocated outside.
85   */
86  static void
87  cbwfq_add_class(struct Qdisc *sch, struct cbwfq_class *cl)
88  {
89      struct cbwfq_sched_data *q = qdisc_priv(sch);
90      cl->queue = qdisc_create_dflt(sch->dev_queue,
91                                  &pfifo_qdisc_ops, cl->common.classid);
92      qdisc_class_hash_insert(&q->clhash, &cl->common);
93  }
94
95  /**
96   * Destroy class.
97   *
98   * @se Free memory.
99   */
100 static void
101 cbwfq_destroy_class(struct Qdisc *sch, struct cbwfq_class *cl)
102 {
103     struct cbwfq_sched_data *q = qdisc_priv(sch);
104
105     sch_tree_lock(sch);
106
107     qdisc_tree_reduce_backlog(cl->queue, cl->queue->q.len,
108                              cl->queue->qstats.backlog);
109     qdisc_class_hash_remove(&q->clhash, &cl->common);
110
111     sch_tree_unlock(sch);
112
113     if (cl->queue) {
114         qdisc_destroy(cl->queue);
115     }

```

```

116     kfree(cl);
117 }
118
119 /**
120  * Find class with given classid.
121  */
122 static inline struct cbwfq_class *
123 cbwfq_class_lookup(struct cbwfq_sched_data *q, u32 classid)
124 {
125     struct Qdisc_class_common *clc;
126
127     clc = qdisc_class_find(&q->clhash, classid);
128     if (clc == NULL)
129         return NULL;
130     return container_of(clc, struct cbwfq_class, common);
131 }
132
133 /**
134  * Find class with given id and return its address.
135  */
136 static unsigned long
137 cbwfq_find(struct Qdisc *sch, u32 classid)
138 {
139     struct cbwfq_sched_data *q = qdisc_priv(sch);
140     return (unsigned long)cbwfq_class_lookup(q, classid);
141 }
142
143 /**
144  * Modify class with given options.
145  */
146 static int
147 cbwfq_modify_class(struct Qdisc *sch, struct cbwfq_class *cl,
148                   struct tc_cbwfq_copt *copt)
149 {
150     struct cbwfq_sched_data *q = qdisc_priv(sch);
151
152     if (copt->cbwfq_cl_limit > 0) {
153         cl->limit = copt->cbwfq_cl_limit;
154     }
155
156     if (copt->cbwfq_cl_rate_type != q->rtype) {
157         PRINT_INFO_ARGS("different rate types.");
158         return -EINVAL;
159     }
160     cl->rate = copt->cbwfq_cl_rate;
161
162     return 0;
163 }
164
165 /**
166  * Create new class.
167  */
168 static int
169 cbwfq_class_create(struct Qdisc *sch, struct tc_cbwfq_copt *copt,
170                  unsigned long classid)
171 {
172     struct cbwfq_class *cl;
173     struct cbwfq_sched_data *q = qdisc_priv(sch);
174
175     cl = kmalloc(sizeof(struct cbwfq_class), GFP_KERNEL);

```

```

176     if (cl == NULL)
177         return -ENOMEM;
178
179     cl->common.classid = classid;
180     cl->limit          = 1000;
181     cl->rate           = 0;
182     cl->cl_sn          = 0;
183     cl->is_active      = false;
184
185     if (cbwfq_modify_class(sch, cl, copt) != 0) {
186         kfree(cl);
187         return -EINVAL;
188     }
189
190     sch_tree_lock(sch);
191     cbwfq_add_class(sch, cl);
192     sch_tree_unlock(sch);
193
194     return 0;
195 }
196
197 /**
198  * Add or change a class by given id.
199  *
200  * @se Allocated memory.
201  */
202 static int
203 cbwfq_change_class(struct Qdisc *sch, u32 classid, u32 parentid,
204                   struct nlattr **tca, unsigned long *arg)
205 {
206     struct cbwfq_sched_data *q = qdisc_priv(sch);
207     struct cbwfq_class *cl;
208     struct nlattr *opt = tca[TCA_OPTIONS];
209     struct nlattr *tb[TCA_CBWFQ_MAX + 1];
210     struct tc_cbwfq_copt *copt;
211     int err;
212     int p_maj = TC_H_MAJ(parentid) >> 16;
213     int cid_maj = TC_H_MAJ(classid) >> 16;
214
215     /* Both have to be 1, because there's no class heirarchy. */
216     if (p_maj != 1 || cid_maj != 1)
217         return -EINVAL;
218
219     if (opt == NULL) {
220         return -EINVAL;
221     }
222
223     err = nla_parse_nested(tb, TCA_CBWFQ_MAX, opt, cbwfq_policy, NULL);
224     if (err < 0) {
225         return err;
226     }
227
228     if (tb[TCA_CBWFQ_PARAMS] == NULL) {
229         return -EINVAL;
230     }
231     copt = nla_data(tb[TCA_CBWFQ_PARAMS]);
232
233     cl = cbwfq_class_lookup(q, classid);
234     if (cl != NULL) {
235         return cbwfq_modify_class(sch, cl, copt);

```

```

236     }
237     return cbwfq_class_create(sch, copt, classid, extack);
238 }
239
240 /**
241  * Delete class by given id.
242  *
243  * @se Free memory.
244  */
245 static int
246 cbwfq_delete_class(struct Qdisc *sch, unsigned long arg)
247 {
248     struct cbwfq_class *cl = (struct cbwfq_class *)arg;
249
250     if (cl == NULL || cl->common.classid == DEFAULT_CL_ID) {
251         return -EINVAL;
252     }
253
254     cbwfq_destroy_class(sch, cl);
255     return 0;
256 }
257
258 /**
259  * Classify the given packet to a class.
260  */
261 static struct cbwfq_class *
262 cbwfq_classify(struct sk_buff *skb, struct Qdisc *sch, int *qerr)
263 {
264     struct cbwfq_sched_data *q = qdisc_priv(sch);
265     struct cbwfq_class *cl;
266     struct tcf_result res;
267     struct tcf_proto *fl;
268     int err;
269     u32 classid = TC_H_MAKE(1 << 16, 1);
270
271     *qerr = NET_XMIT_SUCCESS | __NET_XMIT_BYPASS;
272     if (TC_H_MAJ(skb->priority) != sch->handle) {
273         fl = rcu_dereference_bh(q->filter_list);
274         err = tcf_classify(skb, fl, &res, false);
275
276         if (!fl || err < 0) {
277             return q->default_queue;
278         }
279
280     #ifdef CONFIG_NET_CLS_ACT
281         switch (err) {
282             case TC_ACT_STOLEN:
283             case TC_ACT_QUEUED:
284             case TC_ACT_TRAP:
285                 *qerr = NET_XMIT_SUCCESS | __NET_XMIT_STOLEN;
286                 /* fall through */
287             case TC_ACT_SHOT:
288                 return NULL;
289         }
290     #endif
291
292     classid = res.classid;
293 }
294
295 return cbwfq_class_lookup(q, classid);

```

```

296 }
297
298 /**
299  * Enqueue packet to the queue.
300  */
301 static int
302 cbwfq_enqueue(struct sk_buff *skb, struct Qdisc *sch, struct sk_buff **
               to_free)
303 {
304     struct cbwfq_sched_data *q = qdisc_priv(sch);
305     struct cbwfq_class *cl;
306     struct Qdisc *qdisc;
307     int ret;
308
309     cl = cbwfq_classify(skb, sch, &ret);
310     if (cl == NULL || cl->queue == NULL) {
311         if (ret & __NET_XMIT_BYPASS)
312             qdisc_qstats_drop(sch);
313         __qdisc_drop(skb, to_free);
314         return ret;
315     }
316
317     if (cl->queue->q.qlen >= cl->limit) {
318         if (net_xmit_drop_count(ret)) {
319             qdisc_qstats_drop(sch);
320         }
321         return qdisc_drop(skb, sch, to_free);
322     }
323
324     qdisc = cl->queue;
325     ret = qdisc_enqueue(skb, qdisc, to_free);
326     if (ret == NET_XMIT_SUCCESS) {
327         u32 virtual_len = qdisc_pkt_len(skb) * (q->ifrate / cl->rate);
328         if (!cl->is_active) {
329             cl->cl_sn = q->sch_sn + virtual_len;
330             cl->is_active = true;
331             q->active_rate += cl->rate;
332         } else {
333             cl->cl_sn += virtual_len;
334         }
335         skb->tstamp = cl->cl_sn;
336
337         sch->q.qlen++;
338         qdisc_qstats_backlog_inc(sch, skb);
339         qdisc_qstats_backlog_inc(cl->queue, skb);
340         return NET_XMIT_SUCCESS;
341     }
342
343     if (net_xmit_drop_count(ret)) {
344         qdisc_qstats_drop(sch);
345         qdisc_qstats_drop(cl->queue);
346     }
347     return ret;
348 }
349
350 /**
351  * Find minimum class with minimum sequence number.
352  */
353 static struct cbwfq_class *
354 cbwfq_find_min(struct cbwfq_sched_data *q)

```

```

355 {
356     struct cbwfq_class *it, *cl = NULL;
357     ktime_t ft = KTIME_MAX;
358     int i;
359
360     for (i = 0; i < q->clhash.hashsize; i++) {
361         hlist_for_each_entry(it, &q->clhash.hash[i], common.hnode) {
362             if (it->is_active) {
363                 struct sk_buff *skb = it->queue->ops->peek(it->queue);
364                 if (ft > skb->tstamp) {
365                     cl = it;
366                     ft = skb->tstamp;
367                 }
368             }
369         }
370     }
371     return cl;
372 }
373
374 /**
375  * Return packet without deleting it from a queue.
376  */
377 static struct sk_buff *
378 cbwfq_peek(struct Qdisc *sch)
379 {
380     struct cbwfq_sched_data *q = qdisc_priv(sch);
381     struct cbwfq_class *cl = NULL;
382
383     cl = cbwfq_find_min(q);
384     if (cl == NULL) {
385         return NULL;
386     }
387     return cl->queue->ops->peek(cl->queue);
388 }
389
390 /**
391  * Dequeue packet.
392  */
393 static struct sk_buff *
394 cbwfq_dequeue(struct Qdisc *sch)
395 {
396     struct cbwfq_sched_data *q = qdisc_priv(sch);
397     struct cbwfq_class *cl = NULL;
398     struct sk_buff *skb;
399
400     cl = cbwfq_find_min(q);
401     if (cl == NULL) {
402         return NULL;
403     }
404
405     skb = cl->queue->ops->dequeue(cl->queue);
406     if (skb == NULL) {
407         return NULL;
408     }
409
410     qdisc_bstats_update(sch, skb);
411     qdisc_qstats_backlog_dec(sch, skb);
412     qdisc_qstats_backlog_dec(cl->queue, skb);
413     sch->q.qlen --;
414

```

```

415     if (cl->queue->q.qlen == 0) {
416         cl->is_active = false;
417         cl->cl_sn      = 0;
418         q->active_rate -= cl->rate;
419     }
420
421     if (q->active_rate == 0) {
422         q->sch_sn = 0;
423     } else {
424         q->sch_sn = skb->tstamp;
425     }
426     return skb;
427 }
428
429 /**
430  * Reset qdisc.
431  */
432 static void
433 cbwfq_reset(struct Qdisc *sch)
434 {
435     int i;
436     struct cbwfq_sched_data *q = qdisc_priv(sch);
437     struct cbwfq_class *it;
438
439     q->sch_sn = 0;
440     for (i = 0; i < q->clhash.hashsize; i++) {
441         hlist_for_each_entry(it, &q->clhash.hash[i], common.hnode) {
442             cl->cl_sn = 0;
443             cl->is_active = false;
444             qdisc_reset(it->queue);
445         }
446     }
447     sch->qstats.backlog = 0;
448     sch->q.qlen = 0;
449 }
450
451 /**
452  * Destroy qdisc.
453  */
454 static void
455 cbwfq_destroy(struct Qdisc *sch)
456 {
457     int i;
458     struct cbwfq_sched_data *q = qdisc_priv(sch);
459     struct cbwfq_class *it;
460     struct hlist_node *next;
461
462     tcf_block_put(q->block);
463
464     for (i = 0; i < q->clhash.hashsize; i++) {
465         hlist_for_each_entry_safe(it, next, &q->clhash.hash[i], common.hnode)
466             {
467                 if (it != NULL) {
468                     cbwfq_destroy_class(sch, it);
469                 }
470             }
471     }
472     qdisc_watchdog_cancel(&q->watchdog);
473     qdisc_class_hash_destroy(&q->clhash);
474 }

```

```

474
475 /**
476  * Change qdisc configuration.
477  */
478 static int
479 cbwfq_change(struct Qdisc *sch, struct nlattr *opt)
480 {
481     struct cbwfq_sched_data *q = qdisc_priv(sch);
482     struct tc_cbwfq_glob *qopt;
483     struct nlattr *tb[TCA_CBWFQ_MAX + 1];
484     int err;
485
486     if (opt == NULL) {
487         return -EINVAL;
488     }
489
490     err = nla_parse_nested(tb, TCA_CBWFQ_MAX, opt, cbwfq_policy, NULL);
491     if (err < 0) {
492         return err;
493     }
494
495     if (tb[TCA_CBWFQ_INIT] == NULL) {
496         return -EINVAL;
497     }
498
499     qopt = nla_data(tb[TCA_CBWFQ_INIT]);
500
501     sch_tree_lock(sch);
502
503     if (qopt->cbwfq_gl_default_limit > 0) {
504         q->default_queue->limit = qopt->cbwfq_gl_default_limit;
505     }
506
507     sch_tree_unlock(sch);
508     return 0;
509 }
510
511 /**
512  * Initilize new instance of qdisc.
513  */
514 static int
515 cbwfq_init(struct Qdisc *sch, struct nlattr *opt)
516 {
517     struct cbwfq_sched_data *q = qdisc_priv(sch);
518     struct cbwfq_class *cl;
519     struct tc_cbwfq_glob *qopt;
520     struct nlattr *tb[TCA_CBWFQ_MAX + 1];
521     int err;
522
523     if (!opt)
524         return -EINVAL;
525
526     /* Init filter system. */
527     err = tcf_block_get(&q->block, &q->filter_list, sch, extack);
528     if (err)
529         return err;
530
531     /* Init hash table for class storing. */
532     err = qdisc_class_hash_init(&q->clhash);
533     if (err < 0)

```



```

534         return err;
535
536     q->active_rate = 0;
537     q->sch_sn      = 0;
538
539     /* Init default queue. */
540     cl = kmalloc( sizeof(struct cbwfq_class), GFP_KERNEL);
541     if (cl == NULL) {
542         return -ENOMEM;
543     }
544     q->default_queue = cl;
545
546     /* Set classid for default class. */
547     cl->common.classid = TC_H_MAKE(1 << 16, 1);
548     cl->limit          = 1024;
549     cl->rate           = 0;
550     cl->is_active      = false;
551     cl->cl_sn          = 0;
552
553     err = nla_parse_nested(tb, TCA_CBWFQ_MAX, opt, cbwfq_policy, NULL);
554     if (err < 0)
555         return err;
556
557     qopt = nla_data(tb[TCA_CBWFQ_INIT]);
558
559     if (qopt->cbwfq_gl_default_limit != 0) {
560         cl->limit = qopt->cbwfq_gl_default_limit;
561     }
562
563     if (qopt->cbwfq_gl_rate_type == TCA_CBWFQ_RT_BYTE) {
564         q->rtype = TCA_CBWFQ_RT_BYTE;
565         q->ifrate = qopt->cbwfq_gl_total_rate;
566         cl->rate = qopt->cbwfq_gl_default_rate;
567     } else {
568         q->rtype = TCA_CBWFQ_RT_PERCENT;
569         q->ifrate = 100;
570         cl->rate = qopt->cbwfq_gl_default_rate;
571     }
572
573     sch_tree_lock(sch);
574     cbwfq_add_class(sch, cl, extack);
575     sch_tree_unlock(sch);
576     qdisc_watchdog_init(&q->watchdog, sch);
577     return 0;
578 }
579
580 /**
581  * Dump qdisc configuration.
582  */
583 static int
584 cbwfq_dump(struct Qdisc *sch, struct sk_buff *skb)
585 {
586     struct cbwfq_sched_data *q = qdisc_priv(sch);
587     unsigned char *b = skb_tail_pointer(skb);
588     struct tc_cbwfq_glob opt;
589     struct nlattr *nest;
590
591     memset(&opt, 0, sizeof(opt));
592     opt.cbwfq_gl_total_rate = q->ifrate;
593

```

```

594     nest = nla_nest_start(skb, TCA_OPTIONS);
595     if (nest == NULL)
596         goto nla_put_failure;
597
598     if (nla_put(skb, TCA_CBWFQ_INIT, sizeof(opt), &opt))
599         goto nla_put_failure;
600
601     return nla_nest_end(skb, nest);
602
603 nla_put_failure:
604     nlmsg_trim(skb, b);
605     return -1;
606 }
607
608 /**
609  * Dump class configuration.
610  */
611 static int
612 cbwfq_dump_class(struct Qdisc *sch, unsigned long cl,
613                 struct sk_buff *skb, struct tcmsg *tcm)
614 {
615     struct cbwfq_class *c = (struct cbwfq_class *)cl;
616     struct nlattr *nest;
617     struct tc_cbwfq_copt opt;
618
619     if (c == NULL) {
620         return -1;
621     }
622
623     tcm->tcm_handle = c->common.classid;
624     tcm->tcm_info = c->queue->handle;
625
626     nest = nla_nest_start(skb, TCA_OPTIONS);
627     if (nest == NULL)
628         goto failure;
629
630     memset(&opt, 0, sizeof(opt));
631     opt.cbwfq_cl_limit = c->limit;
632     opt.cbwfq_cl_rate = c->rate;
633
634     if (nla_put(skb, TCA_CBWFQ_PARAMS, sizeof(opt), &opt))
635         goto failure;
636
637     return nla_nest_end(skb, nest);
638
639 failure:
640     nla_nest_cancel(skb, nest);
641     return -1;
642 }
643
644 /**
645  * Dump class statistics.
646  */
647 static int
648 cbwfq_dump_class_stats(struct Qdisc *sch, unsigned long cl,
649                      struct gnet_dump *d)
650 {
651     struct cbwfq_class *c = (struct cbwfq_class *)cl;
652     int gs_base, gs_queue;
653

```

```

654     if (c == NULL)
655         return -1;
656
657     gs_base = gnet_stats_copy_basic(qdisc_root_sleeping_running(sch),
658                                     d, NULL, &c->queue->bstats);
659     gs_queue = gnet_stats_copy_queue(d, NULL, &c->queue->qstats,
660                                     c->queue->q.qlen);
661
662     return gs_base < 0 || gs_queue < 0 ? -1 : 0;
663 }
664
665 /**
666  * Attach a new qdisc to a class and return the prev attached qdisc.
667  */
668 static int
669 cbwfq_graft(struct Qdisc *sch, unsigned long arg, struct Qdisc *new,
670            struct Qdisc **old)
671 {
672     struct cbwfq_sched_data *q = qdisc_priv(sch);
673     struct cbwfq_class *cl;
674
675     if (new == NULL)
676         new = &noop_qdisc;
677
678     cl = cbwfq_class_lookup(q, arg);
679     if (cl) {
680         *old = qdisc_replace(sch, new, &cl->queue);
681         return 0;
682     }
683     return -1;
684 }
685
686 /**
687  * Returns a pointer to the qdisc of class.
688  */
689 static struct Qdisc *
690 cbwfq_leaf(struct Qdisc *sch, unsigned long arg)
691 {
692     struct cbwfq_sched_data *q = qdisc_priv(sch);
693     struct cbwfq_class *cl = cbwfq_class_lookup(q, arg);
694
695     return cl == NULL? NULL : cl->queue;
696 }
697
698 static unsigned long
699 cbwfq_bind(struct Qdisc *sch, unsigned long parent, u32 classid)
700 {
701     return cbwfq_find(sch, classid);
702 }
703
704
705 static void
706 cbwfq_unbind(struct Qdisc *q, unsigned long cl)
707 { }
708
709 /**
710  * Iterates over all classed of a qdisc.
711  */
712 static void
713 cbwfq_walk(struct Qdisc *sch, struct qdisc_walker *arg)

```

```

714 {
715     struct cbwfq_sched_data *q = qdisc_priv(sch);
716     struct cbwfq_class *cl;
717     int h;
718
719     if (arg->stop)
720         return;
721
722     for (h = 0; h < q->clhash.hashsize; h++) {
723         hlist_for_each_entry(cl, &q->clhash.hash[h], common.hnode) {
724             if (arg->count < arg->skip) {
725                 arg->count++;
726                 continue;
727             }
728             if (arg->fn(sch, (unsigned long)cl, arg) < 0) {
729                 arg->stop = 1;
730                 break;
731             }
732             arg->count++;
733         }
734     }
735 }
736
737 static struct tcf_block *
738 cbwfq_tcf_block(struct Qdisc *sch, unsigned long cl, struct netlink_ext_ack *
739               extack)
740 {
741     return qdisc_priv(sch)->block;
742 }
743
744 static const struct Qdisc_class_ops cbwfq_class_ops = {
745     .graft      = cbwfq_graft,
746     .leaf       = cbwfq_leaf,
747     .find       = cbwfq_find,
748     .walk       = cbwfq_walk,
749     .change     = cbwfq_change_class,
750     .delete     = cbwfq_delete_class,
751     .tcf_block  = cbwfq_tcf_block,
752     .bind_tcf   = cbwfq_bind,
753     .unbind_tcf = cbwfq_unbind,
754     .dump       = cbwfq_dump_class,
755     .dump_stats = cbwfq_dump_class_stats,
756 };
757
758 static struct Qdisc_ops cbwfq_qdisc_ops __read_mostly = {
759     /* Points to next Qdisc_ops. */
760     .next      = NULL,
761     /* Points to structure that provides a set of functions for
762      * a particular class. */
763     .cl_ops    = &cbwfq_class_ops,
764     /* Char array contains identity of the qdisc. */
765     .id        = "cbwfq",
766     .priv_size = sizeof(struct cbwfq_sched_data),
767
768     .enqueue   = cbwfq_enqueue,
769     .dequeue   = cbwfq_dequeue,
770     .peek      = cbwfq_peek,
771     .init      = cbwfq_init,
772     .reset     = cbwfq_reset,
773     .destroy   = cbwfq_destroy,

```

```

773     .change      =   cbwfq_change ,
774     .dump        =   cbwfq_dump ,
775     .owner       =   THIS_MODULE,
776 };
777
778 static int __init
779 cbwfq_module_init(void)
780 {
781     return register_qdisc(&cbwfq_qdisc_ops);
782 }
783
784 static void __exit
785 cbwfq_module_exit(void)
786 {
787     unregister_qdisc(&cbwfq_qdisc_ops);
788 }
789
790 module_init(cbwfq_module_init)
791 module_exit(cbwfq_module_exit)
792
793 MODULE_LICENSE("GPL");

```

Листинг 8 — Модуль CBWFQ для ядра Linux.