

ПРОТОКОЛЫ И СЕРВИСЫ

QoS

Осознание необходимости QoS

Изменение способа предоставления интеллектуальных услуг

Необходимость



Высокая доступность

Безопасность

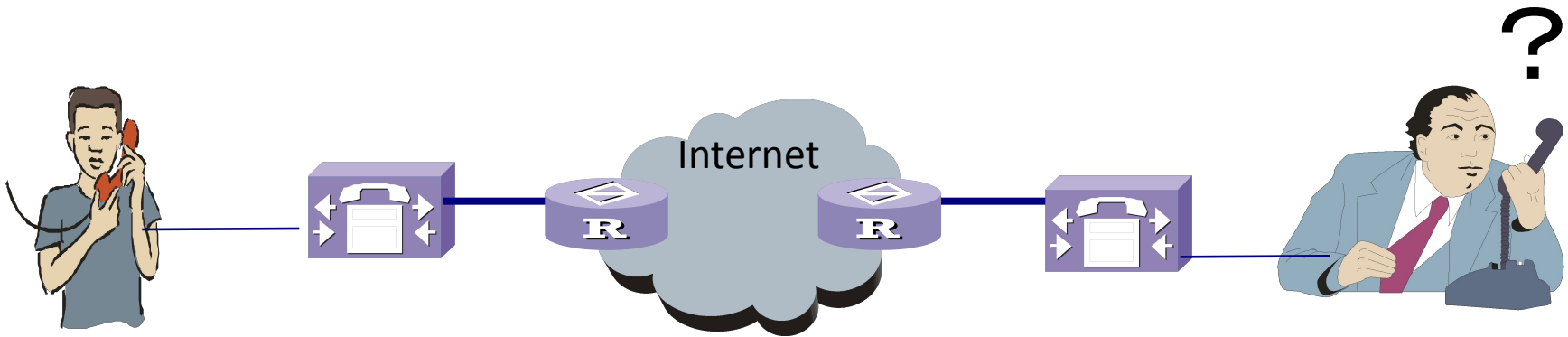


Роскошь

Качество обслуживания



Потеря пакетов...



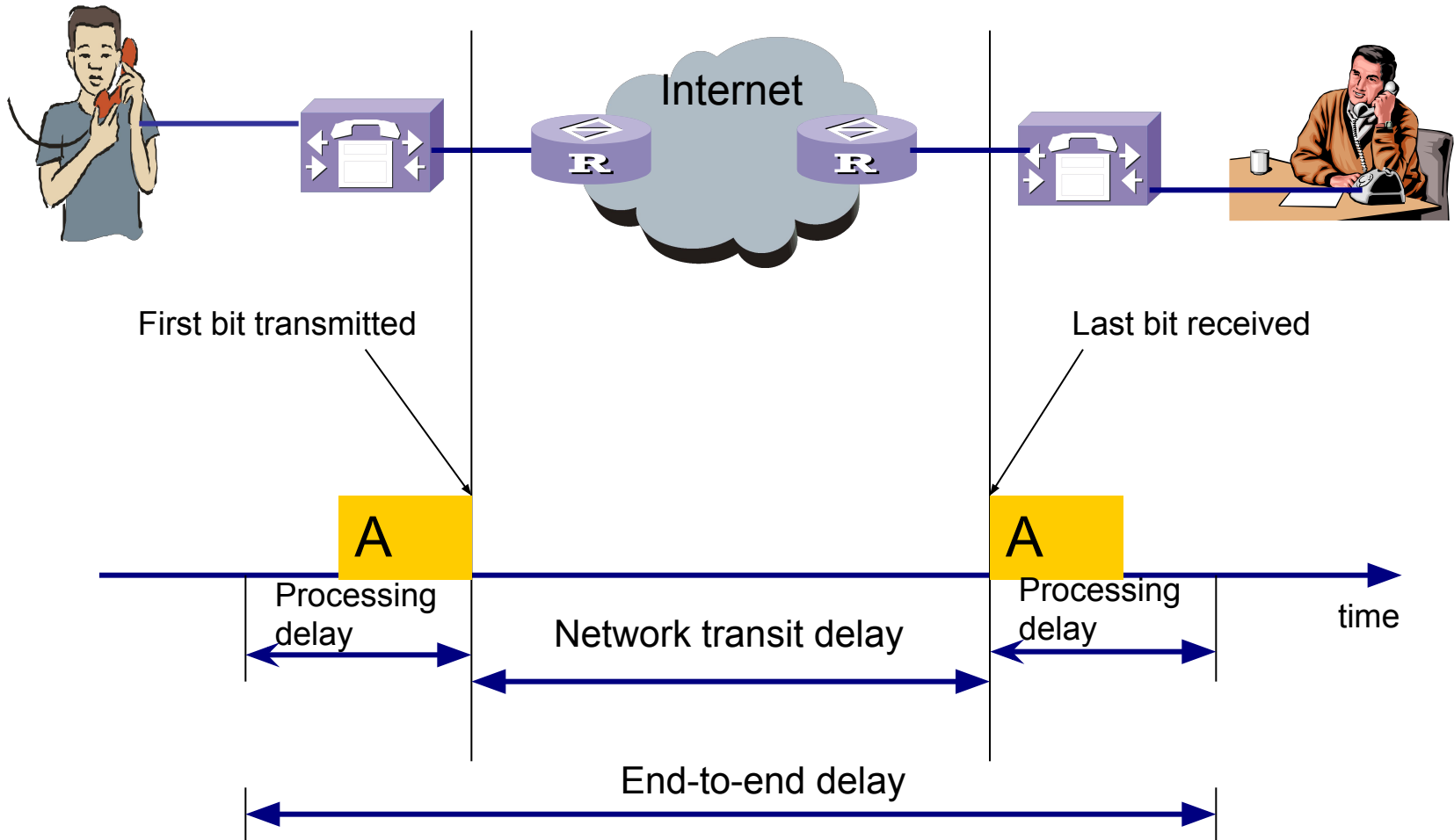
This Is John Smith Speaking

One party said,

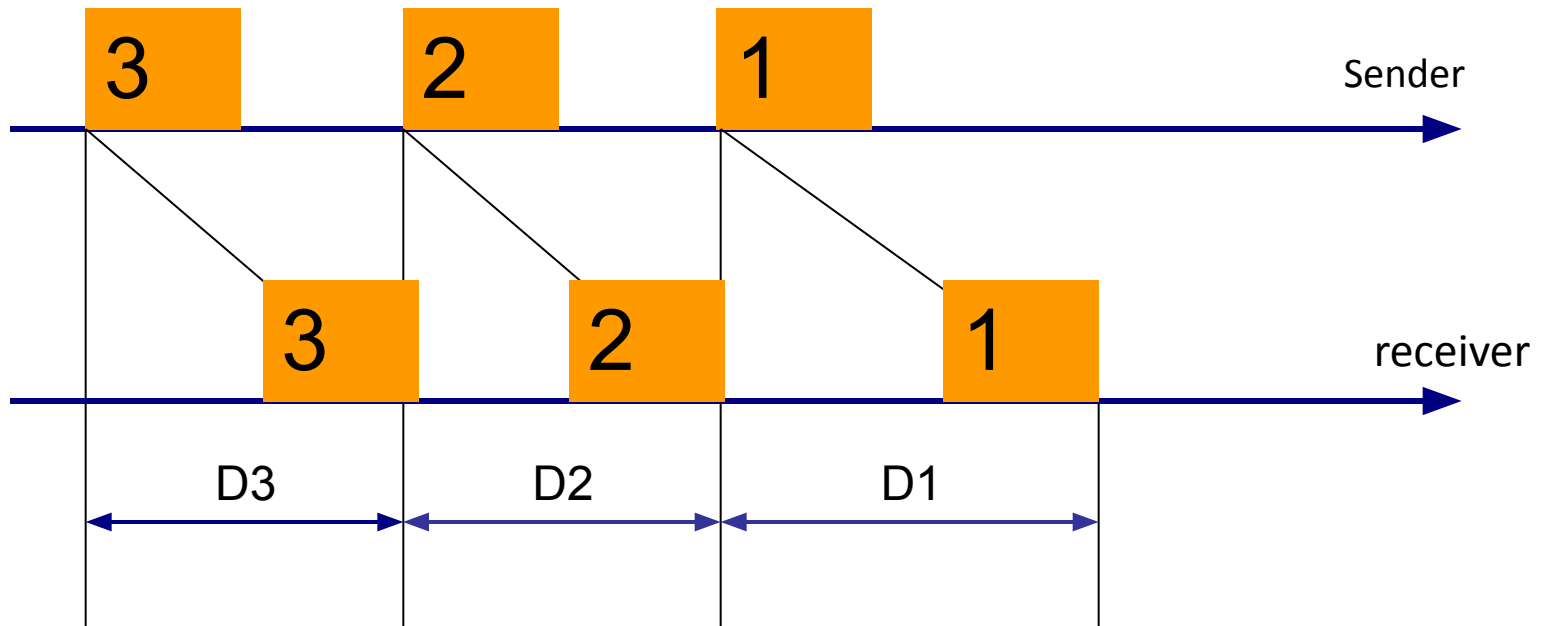
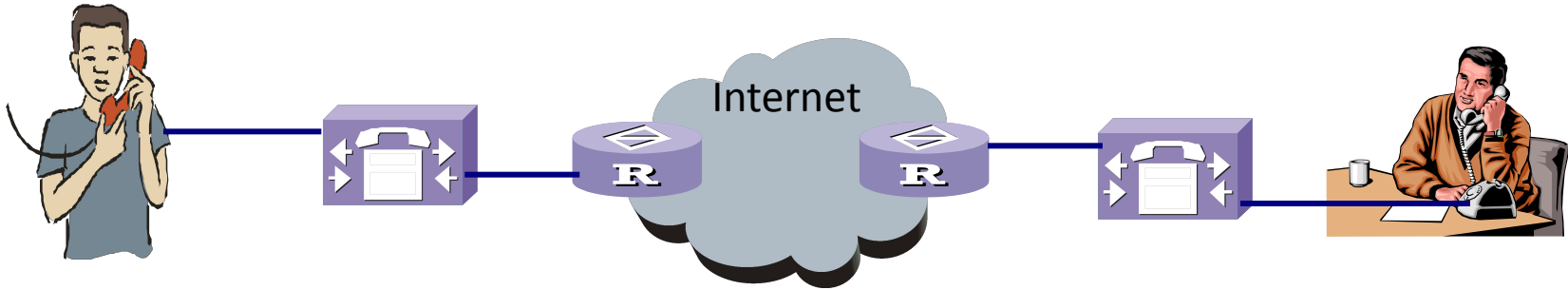
This Is Smith Speaking....

The opposite party heard.....

Delay...

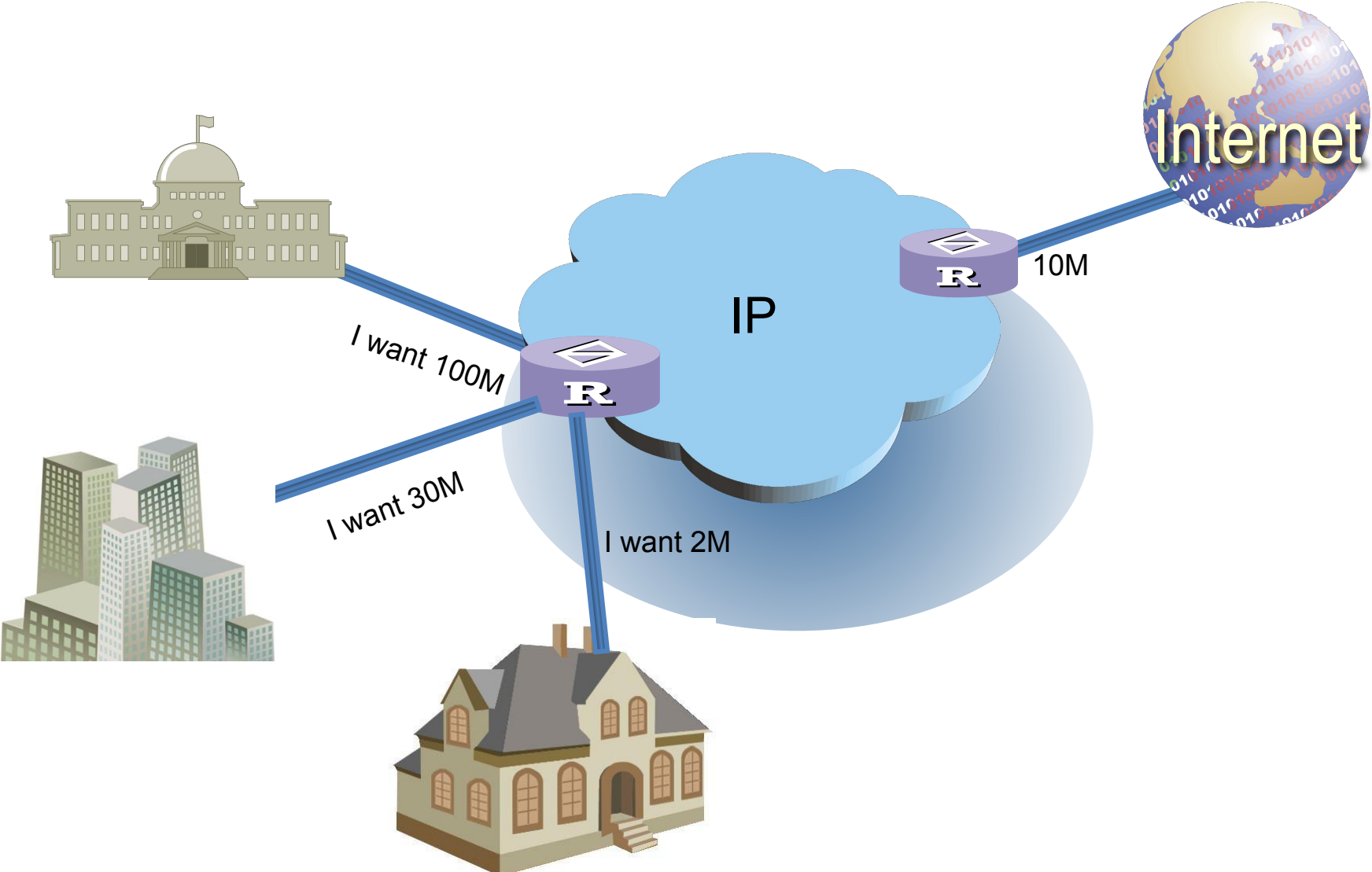


Jitter...



$$D3 = D2 = D1$$

Bandwidth Limit...



Что такое QoS?

- **Качество обслуживания:** Термин, обозначающий набор параметров производительности, характеризующих поток данных на заданном соединении
- **Качество обслуживания:** Доступный для измерения набор параметров, определяющий уровень предоставляемых услуг, который обязуется поддерживать провайдер
- **Качество обслуживания:** Мера производительности системы передачи данных, отражающая качество передачи и доступность сервиса

КАЧЕСТВО ОБСЛУЖИВАНИЯ

QoS (Quality of Servers) рассматривается как «суммарный эффект *рабочих характеристик обслуживания*, который определяет степень удовлетворенности пользователя этой службой» (E.800)

Задача: обеспечить заданное качество обслуживания в сквозном соединении (end-to-end) для различных видов трафика.

Условие: заданное качество обслуживания должны поддерживать все сетевые устройства на всем сквозном соединении.

Службы QoS

- **Best effort** – обработка информации как можно быстрее, но без дополнительных усилий (FIFO, drop tail)
- **Мягкий QoS** – сервис с предпочтениями. Приоритетное обслуживание, значения параметров QoS зависят от характеристик трафика.
- **Жесткий QoS** – гарантированный сервис. Основан на предварительном резервировании ресурсов для каждого потока.

Логические плоскости механизмов QoS

Контрольная плоскость

- Управление допустимостью соединения
- QoS-маршрутизация
- Резервирование ресурсов

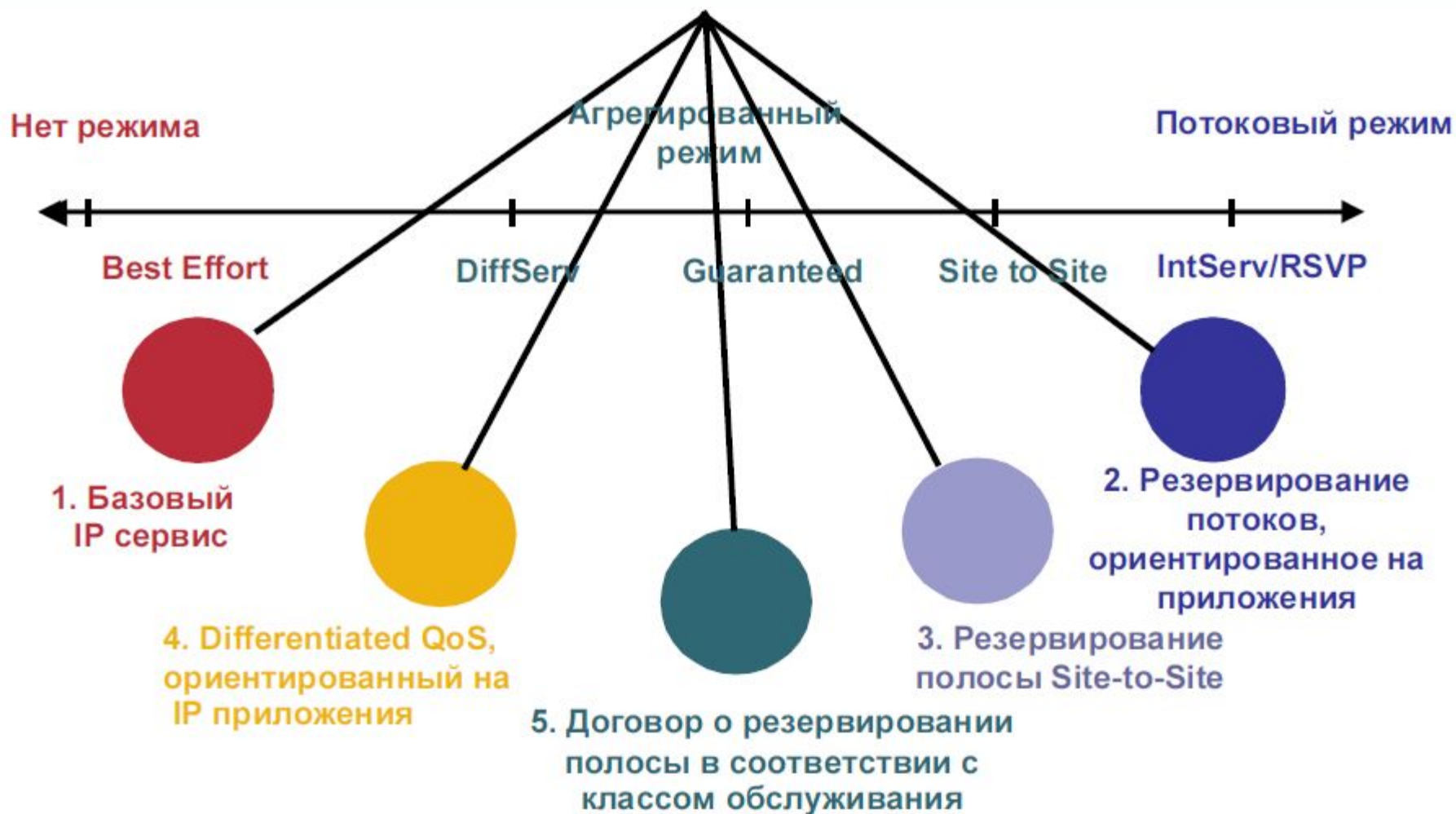
Плоскость данных

- Предотвращение перегрузок
- Управление буфером
- Классификация трафика
- Маркировка пакетов
- Управление характеристиками трафика
- Организация и планирование очередей

Плоскость менеджмента

- Измерения
- Восстановление трафика
- Соглашение об уровне обслуживания

Маятник QoS



Факторы QoS

Атрибуты, требующие исключительного уровня обслуживания

Задержка
Delay
(Latency)

Переменная
задержка
**Delay-
Variation**
(Jitter)

Потеря
пакетов
**Packet
Loss**

Работа с качеством обслуживания

Как работает функциональность QoS?

Классификация и
маркировка

Заполнение очередей и
(выборочный) сброс

shaping/сжатие/
фрагментация/interleaving

Какие возможности есть у QoS?

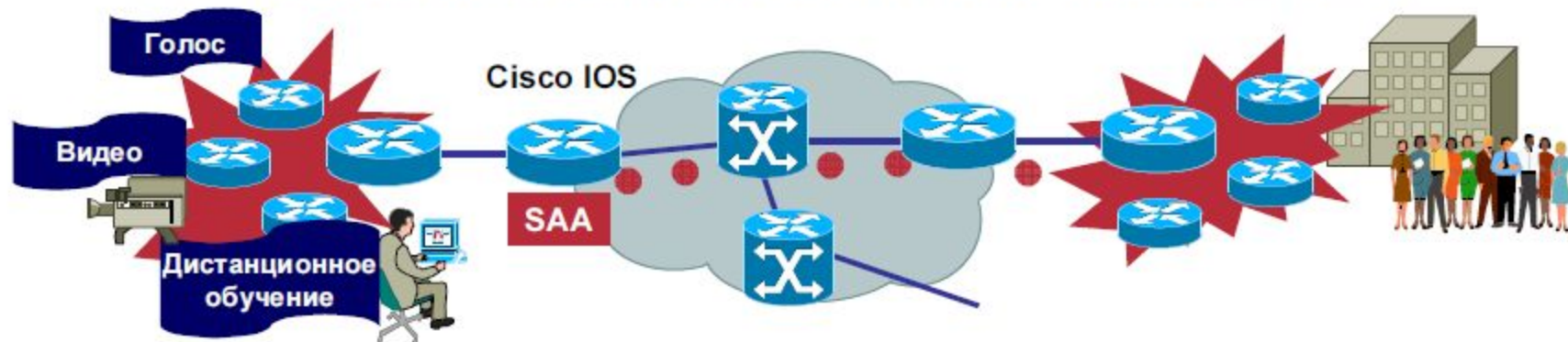
Modular QoS CLI (MQC):

- **Классификация и маркировка:** способность дифференцировать пакеты путём установки определённых значений в заголовок второго или третьего уровня, например IP precedence или L2 class of service или drop priority
- **Policing:** используется для сброса или перемаркировки в более низкий приоритет IP precedence или MPLS EXP битов в потоке данных, который превышает контрактный объем
- **Размещение в очередях:** управление заторами путём выдачи корректных приоритетов классам трафика, с целью улучшения передачи данных приложений, чувствительных к времени, без ухудшения передачи трафика с более низким приоритетом (CBWFQ, MDRR, и т.д..)
- **Shaping:** при превышении скорости передачи, вместо policing могут быть использованы буферизация или помещение в очереди

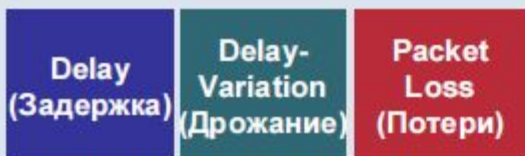
Service Assurance Agent (SAA)

Активное наблюдение за сетевой инфраструктурой

- Допустима ли потеря пакетов?
- Какова задержка в сети?
- Хорошо ли функционируют сетевые приложения?
- Возможен ли контроль за выполнением SLA?



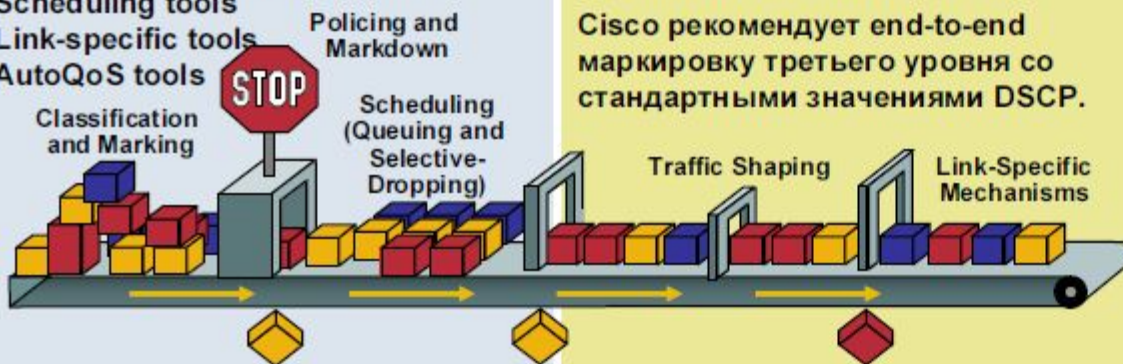
QoS – мера качества передачи и доступности сервиса в сети. Качество передачи в сети определяется следующими факторами: задержка, дрожание и потери.



Технологии QoS – это набор приемов и возможностей управления сетевыми ресурсами, эти технологии являются ключевыми технологиями для прозрачной одновременной передачи по сети Видео, Голоса и Данных. Дополнительно, возможности QoS могут играть стратегическую роль в снижении эффективности или даже предотвращении DoS атак и распространения червей

Cisco's QoS Toolset состоит из :

- Classification and marking tools
- Policing and markdowm tools
- Scheduling tools
- Link-specific tools
- AutoQoS tools



Возможности QoS

Классификация выполняется на 2-7

уровнях:

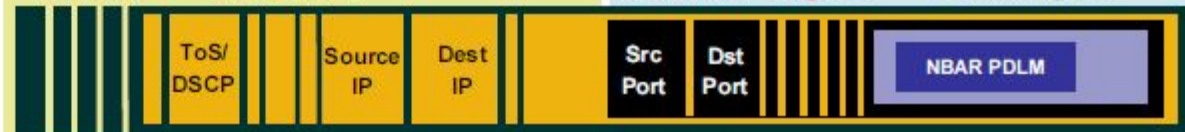
L2 Frame

L3 IP Packet

Policing tools дополняют marking tools измеряя маркированные потоки и пометая не соответствующий контракту трафик.

L4 TCP/UDP Segment

L7 Data Payload

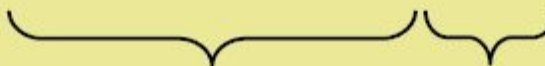
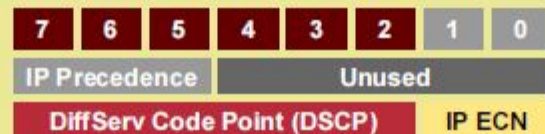


Маркировка выполняется на 2 или 3 уровне:

Layer 2: 802.1Q/p CoS, MPLS EXP

Layer 3: IP Precedence, DSCP и/или IP ECN

Layer 3 (IP ToS Byte) Marking Options:



RFC 2474
DiffServ Extensions

RFC 3168
IP ECN Bits

Cisco рекомендует end-to-end маркировку третьего уровня со стандартными значениями DSCP.

Policer'ы относят трафик к одной из трех категорий:

- Conform: Поток данных соответствует заданной скорости передачи (зеленый свет)
- Exceed: Допустимые всплески средней величины (жёлтый свет)
- Violate: Не допускается передача данных сверх этого лимита (красный свет)



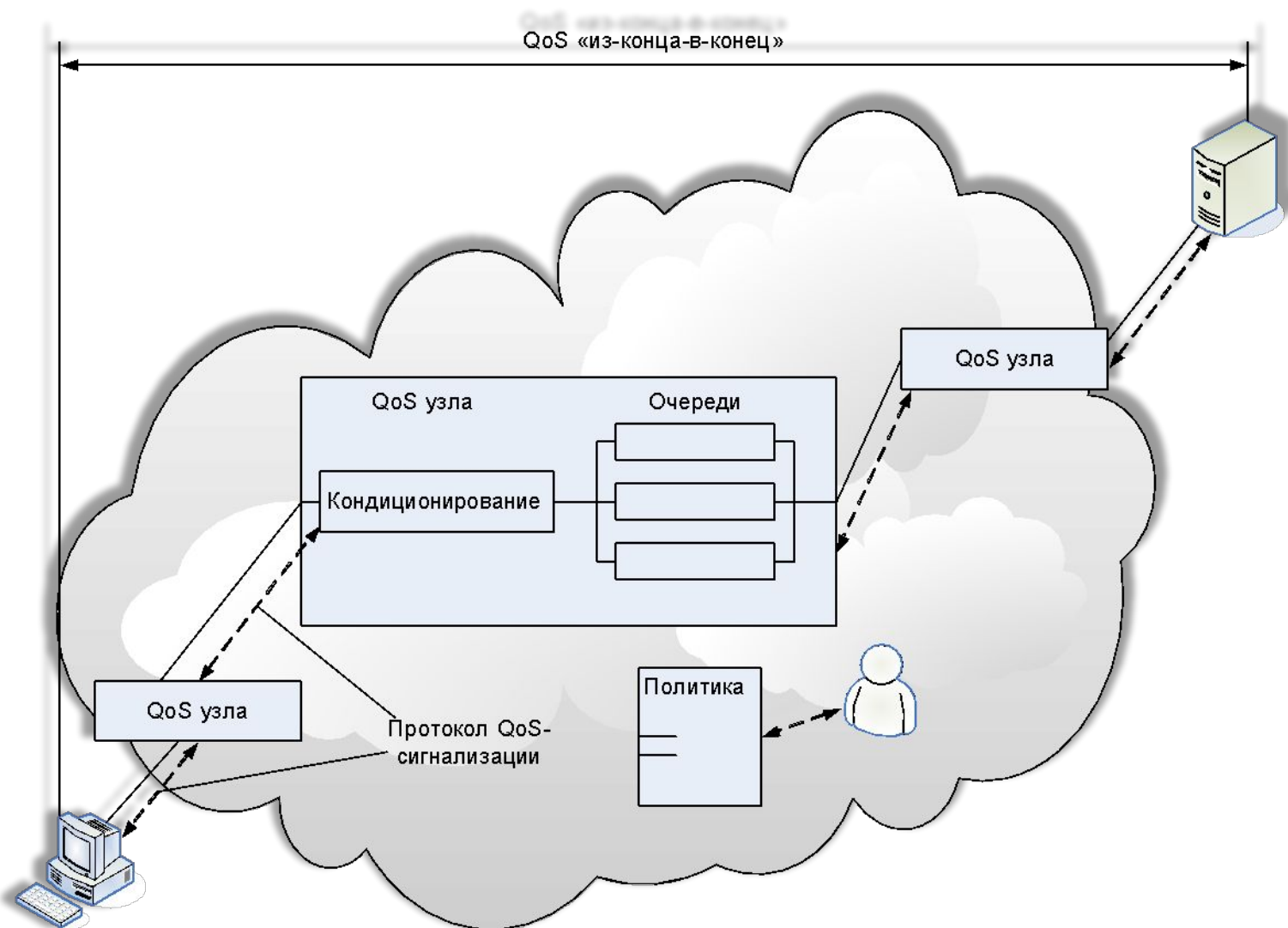
Scheduling tools меняют место пакета в очереди и селективно отбрасывают пакеты в случае возникновения затора



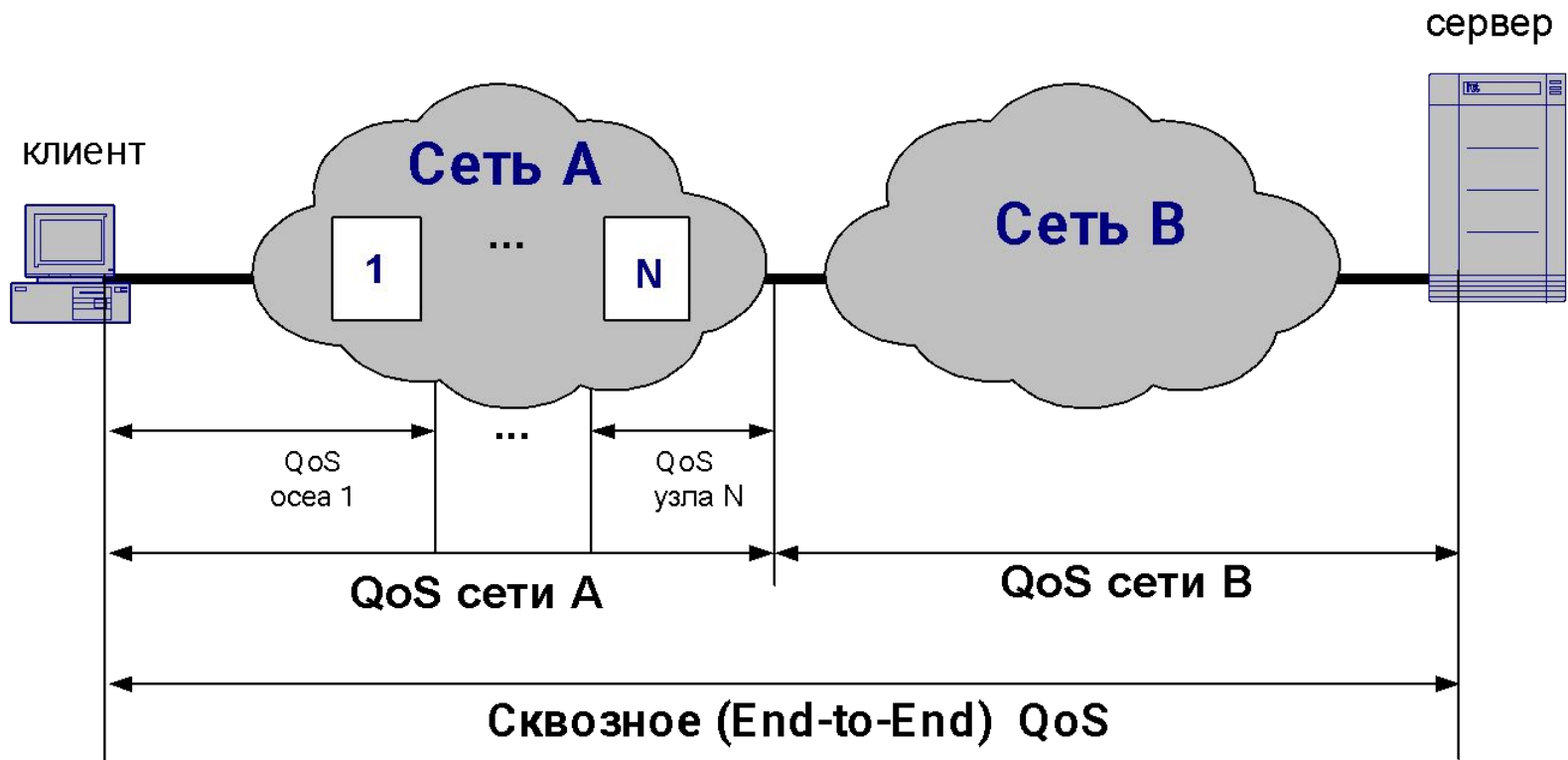
Link-Specific tools полезны на низкоскоростных WAN/VPN соединениях, включает shaping, компрессию, фрагментацию и interleaving.

Функции AutoQoS автоматически настраивает рекомендуемые Cisco значения QoS на коммутаторах Catalyst® и маршрутизаторах под управлением Cisco IOS за одну-две команды.

Модель службы QoS



Эталонная модель СКВОЗНОГО QoS



End-to-End QoS

- End-to-end QoS включает в себя QoS в сети организации и в сети оператора связи
- Подобно цепи, QoS крепок настолько же, насколько крепко его самое слабое звено (канал)



- Для организации эффективного QoS необходима координация работы оператора связи и предприятия

Характеристики QoS (Y.1540)

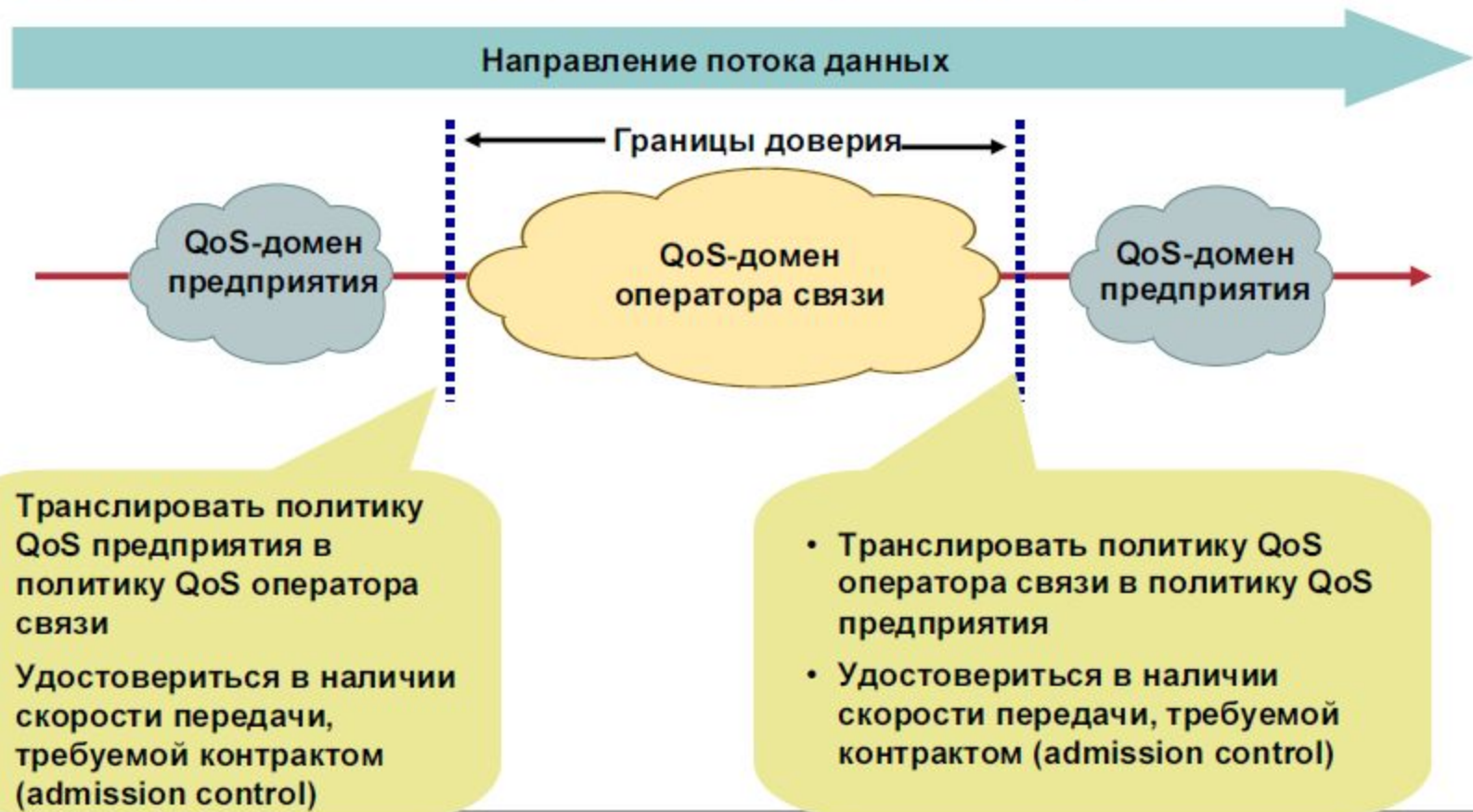
- Задержки и джиттер* задержки
- Величина потерь
- Производительность сети
- Надежность сетевых элементов

G.1000 – определяет структуру связей между рабочими характеристиками QoS.

* джиттер задержки – отклонение значений задержки от заданной величины

Граница доверия оператора связи: Действия, выполняемые в рамках QoS

Действия, выполняемые в рамках QoS, на входящей и выходящей сторонах границы доверия



Соглашение об уровне сервиса (SLA)

Соглашение об уровне сервиса (SLA)

Формализация “качества обслуживания” в контракте между заказчиком и оператором связи

Обычно SLA определяет:

- Классы трафика и критерии идентификации (DSCP, IPP, VLAN, ATM VC)
- Выделение полосы пропускания Per-class/aggregate
- Гарантированную доступность, MTTR
- Процедуру эскалирования проблем
- Гарантированные максимальные потери, задержка, вариацию задержки per-class
- Per-class/aggregate критерии доступности и действия в случае превышения трафика
- Методы измерения SLA и отчётность — точки замеров, методологию, варианты отчётности (web, e-mail) периоды составления отчётов, содержимое отчётов, критерии неисправностей и штрафные неустойки

Атрибуты соглашения об уровне сервиса

Имеющие отношение к QoS	Остальные атрибуты SLA
Потери пакетов	Доступность
Задержка (Delay)	Среднее время восстановления (MTTR)
Вариации задержки (Jitter)	Среднее время наработки на отказ (MTBF)
Критерий Admission Control	
Оговоренная полоса пропускания	
Сохранение потоковой последовательности пакетов	
Пропускная способность	

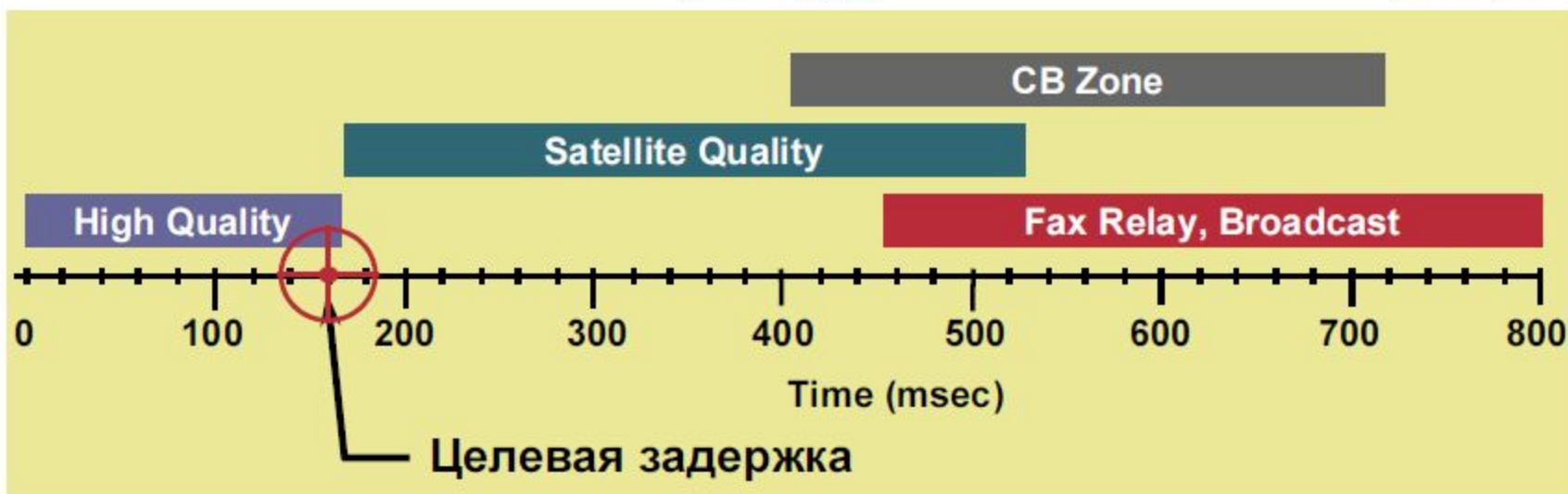
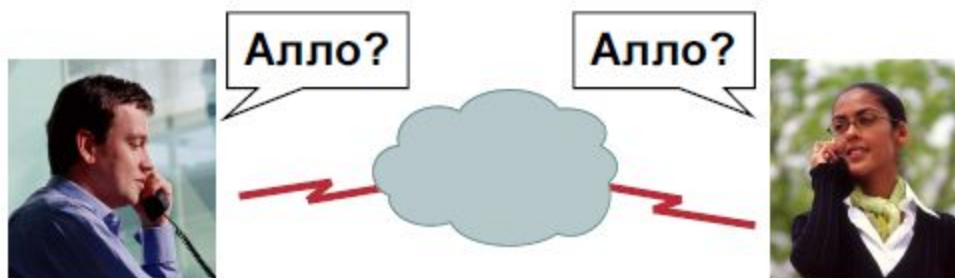
Потери пакетов:

1. Потеря пакетов: обычно рассчитывается как максимальный процент потерь, допустимый для данного класса трафика (как правило менее 0.1%) — исключая трафик вне контракта
2. Потеря пакетов более вероятна на краю сети, чем в ядре (обычно)

Атрибуты соглашения об уровне сервиса

Требования по End-to-End задержке при передаче голоса

Предотвратить
“Human Ethernet”



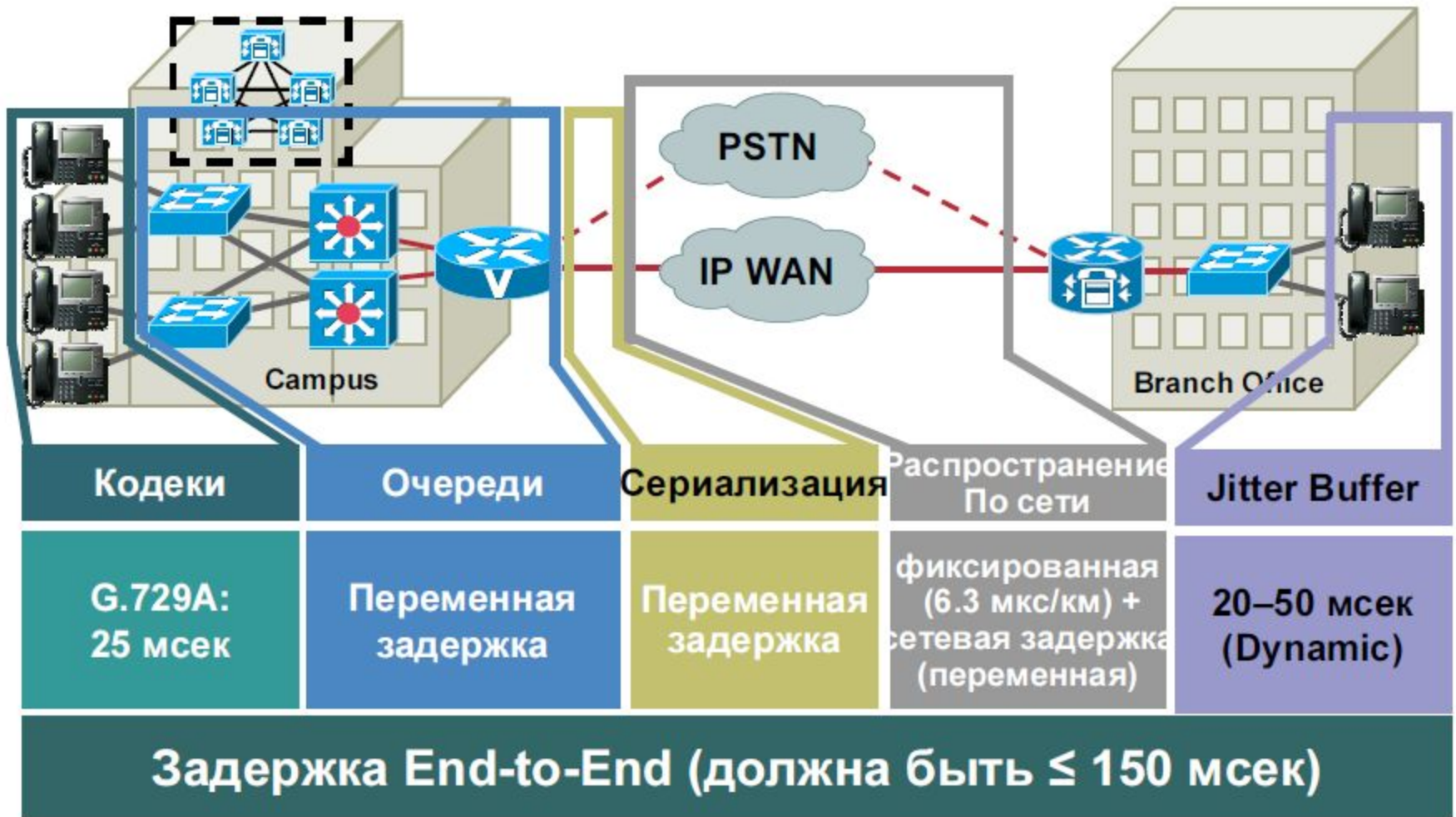
Рекомендации ITU's G.114 : односторонняя задержка ≤ 150 мсек
для качественной передачи голоса

Классификация трафика мультисервисной IP-сети по приложениям

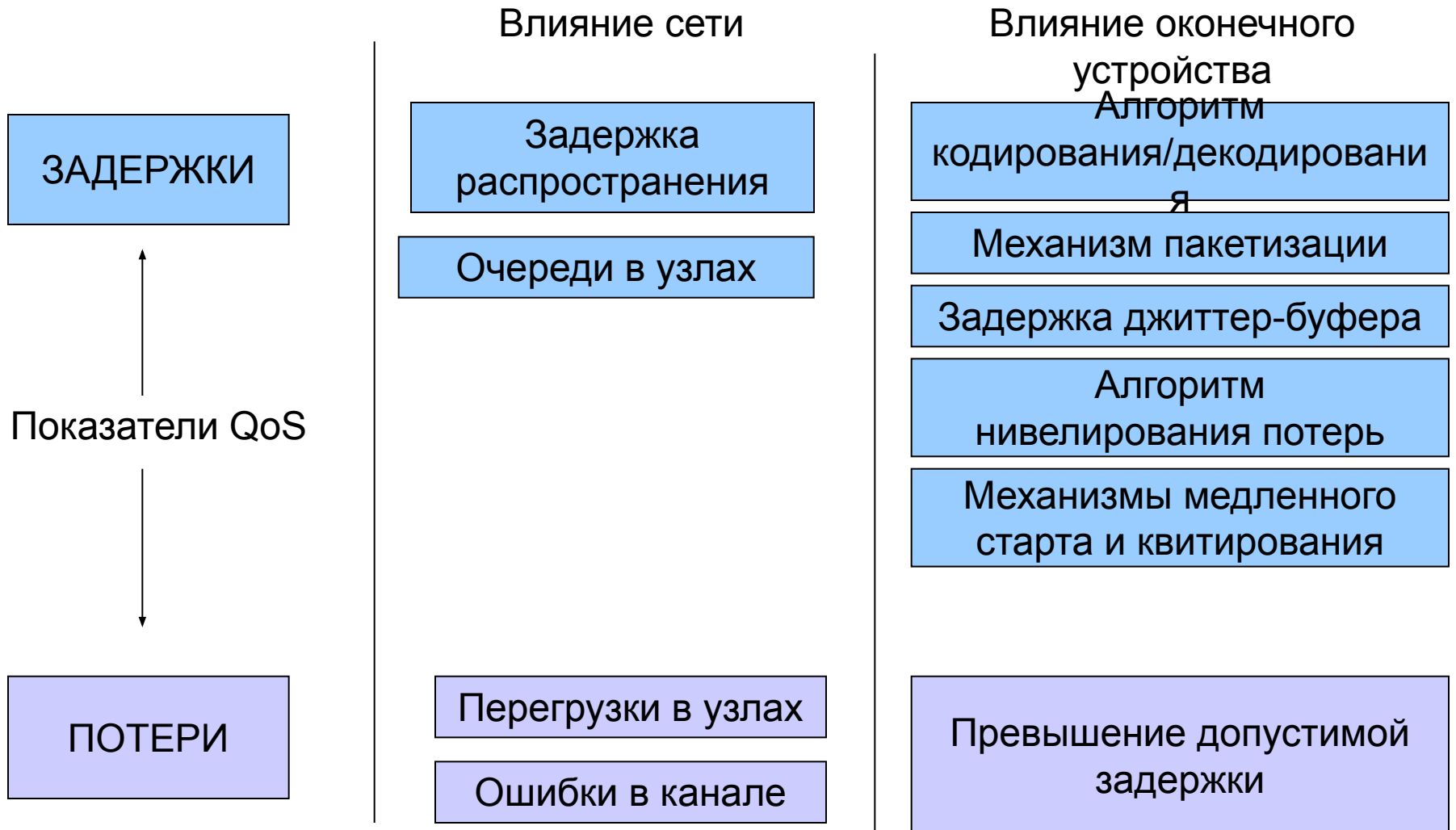
Тип трафика	Приложения	Требования	Протоколы транспортного уровня
Реального времени	IP-телефония и видеоконференцсвязь	<ul style="list-style-type: none"> -Чувствительность к задержке -Чувствительность к джиттеру задержки -Малая чувствительность к потерям 	RSVP, RTP, RTCP, UDP
	Процессы управления, игры on-line	<ul style="list-style-type: none"> -Чувствительность к задержке -Чувствительность к джиттеру задержки -Чувствительность к потерям 	UDP, TCP
Потоковый	Аудио по требованию Видео по требованию Интернет-вещание	<ul style="list-style-type: none"> -Малая чувствительность к задержке -Чувствительность к джиттеру задержки -Чувствительность к потерям 	RSVP, SCTP, UDP, TCP
Эластичный	Конференция документов	<ul style="list-style-type: none"> -Малая чувствительность к задержке -Малая чувствительность к джиттеру задержки -Высокая чувствительность к потерям 	TCP
	Анимация Передача файлов E-mail	<ul style="list-style-type: none"> -Очень малая чувствительность к задержке -Малая чувствительность к джиттеру задержки -Высокая чувствительность к потерям 	

Атрибуты соглашения об уровне сервиса

Факторы влияющие на задержку при передаче голоса



Показатели качества обслуживания, учитываемые при передаче мультимедийного трафика, и механизмы их формирования



Бюджет по задержке для класса realtime: Пример

- Бюджет задержки голосового класса рассчитывается первым
- В соответствии с G.114, максимальная задержка ≤ 150 мсек для качественной передачи голосовых данных
(Может быть > 150 мсек, для “менее, чем бизнес-качественных” звонков, или для межконтинентальных звонков, когда люди ожидают большие задержки при разговоре)
- Бюджет по задержке в соответствии со инженерным SLA более узкий, чем бюджет в соответствии со SLA



Классы QoS и соответствующие им приложения (Y.1541)

- **Класс 0:** Приложения реального времени, чувствительные к джиттеру, характеризующиеся высоким уровнем интерактивности (VoIP, видеоконференции)
- **Класс 1:** Приложения реального времени, чувствительные к джиттеру, интерактивные (VoIP, видеоконференции)
- **Класс 2:** Транзакции данных, характеризующиеся высоким уровнем интерактивности (например, сигнализация)
- **Класс 3:** Транзакции данных, интерактивные приложения
- **Класс 4:** Приложения, допускающие низкий уровень потерь (короткие транзакции, массивы данных, потоковое видео)
- **Класс 5:** Традиционные применения сетей IP

Нормы на параметры доставки пакетов IP с разделением по классам обслуживания, модель ITU-T

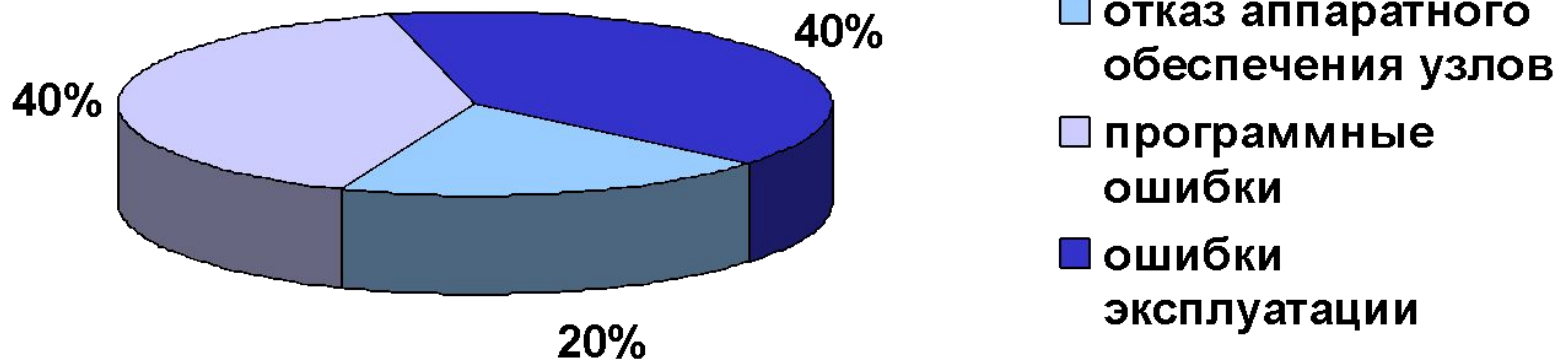
Сетевые характеристики	Классы QoS					
	0	1	2	3	4	5
Задержка доставки пакета IP, IPTD	100 мс	400 мс	100 мс	400 мс	1 с	Н
Вариация задержки пакета IP, IPDV	50 мс	50 мс	Н	Н	Н	Н
Коэффициент потери пакетов IP, IPLR	1×10^{-3}	1×10^{-3}	1×10^{-3}	1×10^{-3}	1×10^{-3}	Н
Коэффициент ошибок пакетов IP, IPER	1×10^{-4}	1×10^{-4}	1×10^{-4}	1×10^{-4}	1×10^{-4}	Н

Примечание. Н - не нормировано. Значения параметров представляют собой верхние границы для средних задержек, джиттера, потерь и ошибок пакетов.

Коэффициенты готовности и значения времени простоя оборудования

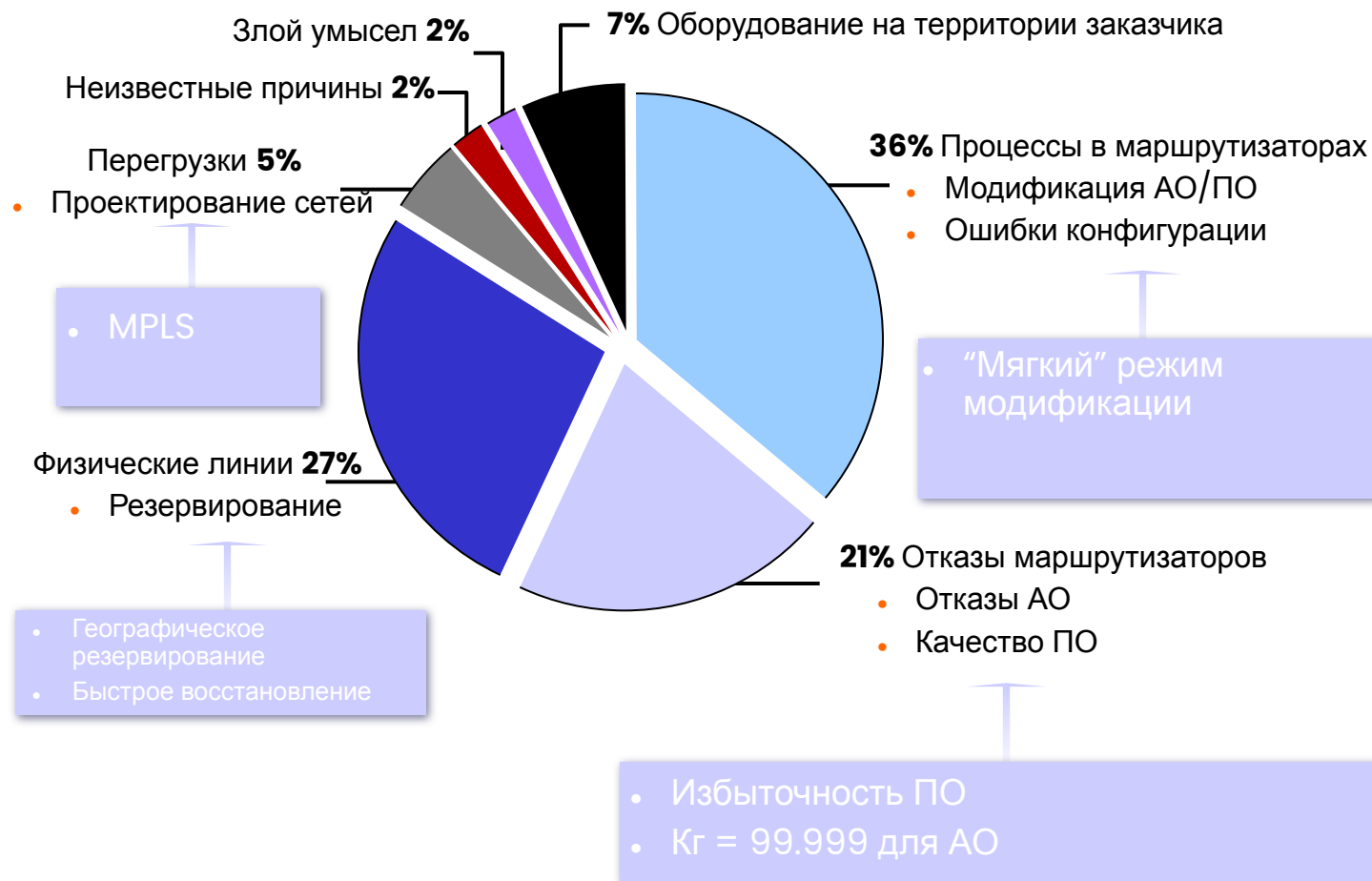
Коэффициент готовности		Время простоя
0,99	“две девятки”	3,7 дней в год
0,999	“три девятки”	9 часов в год
0,9999	“четыре девятки”	53 минуты в год
0,99999	“пять девяток”	5,5 минут в год
0,999999	“шесть девяток”	30 секунд в год

Причины системной ненадежности

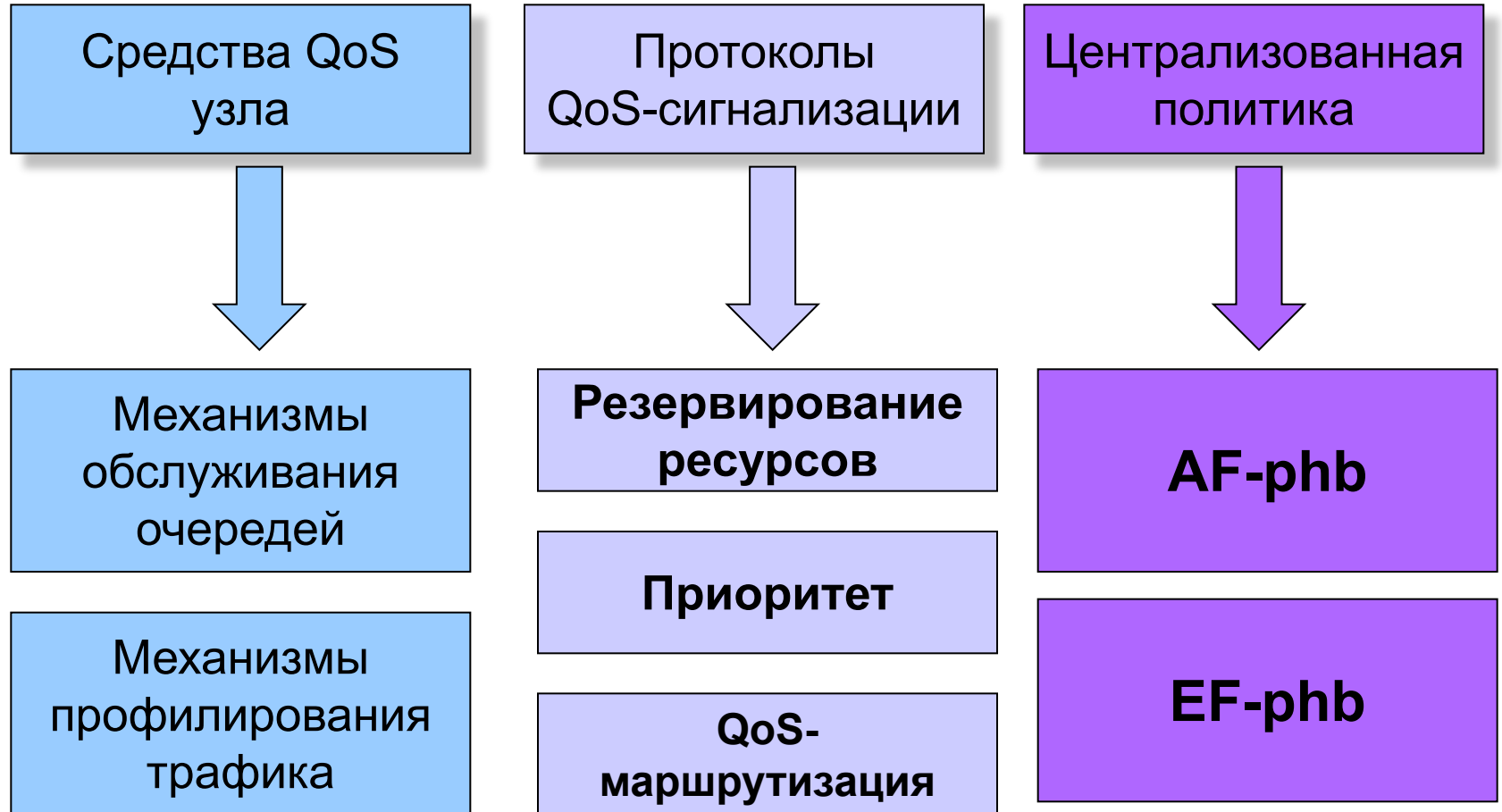


Источник: Gartner Group

Причины отказов в IP-сетях



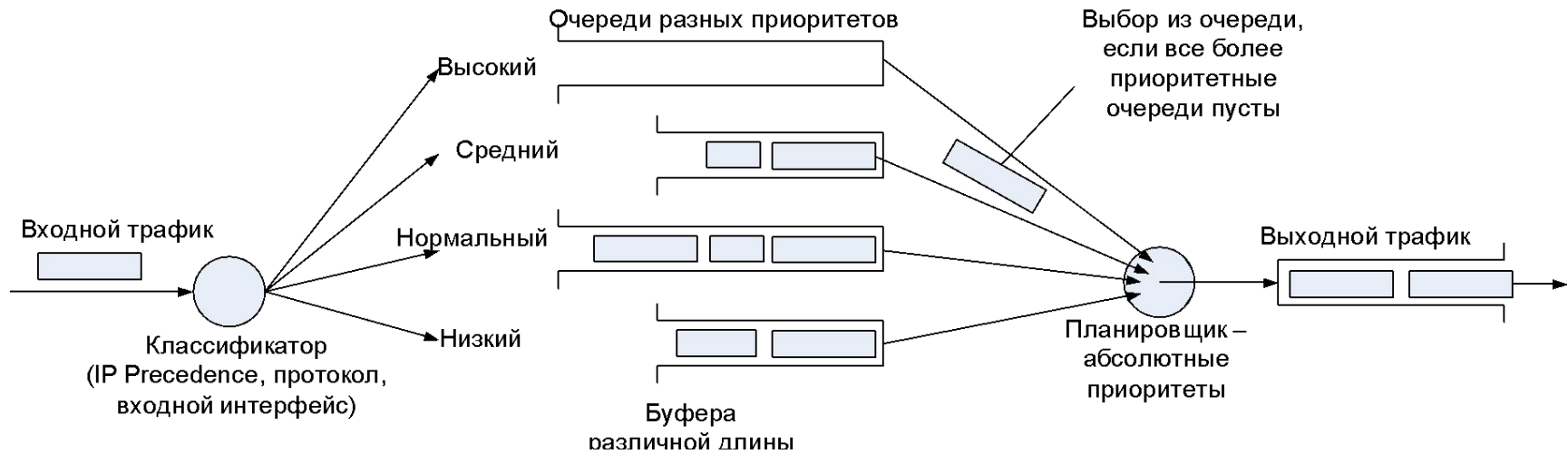
Базовая архитектура службы QoS



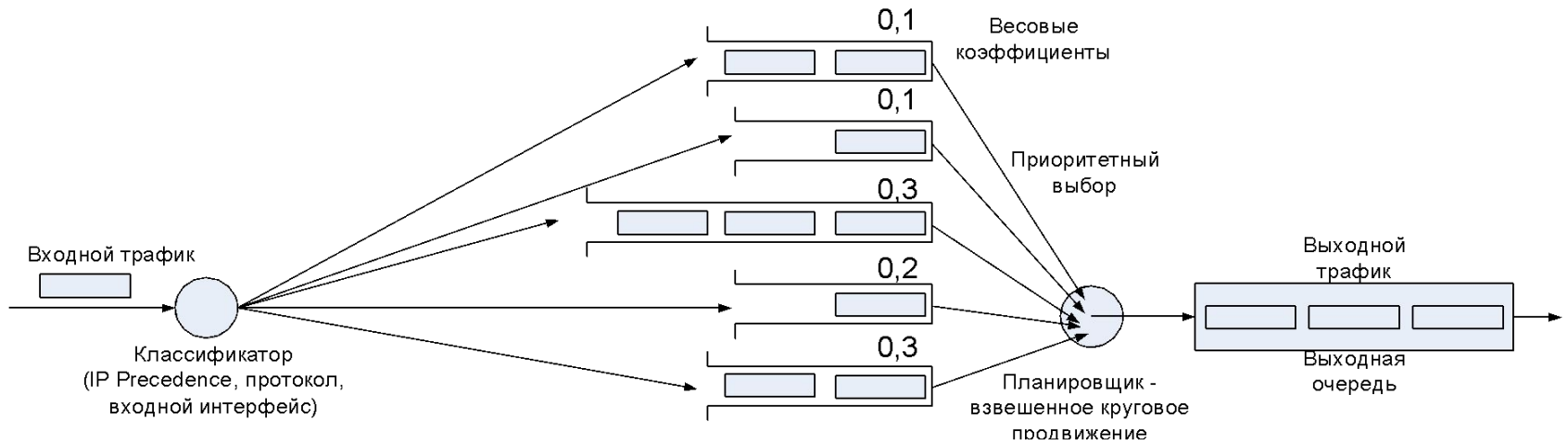
Механизмы обслуживания очередей

- **FIFO** (First In First Out) – без использования дополнительных возможностей, используется в best effort
- **PQ** (Priority Queuing) – приоритетные очереди, вводится приоритет трафика (1-8)
- **CQ** (Custom Queuing) – настраиваемые очереди, используется при резервировании ресурсов
- **WFQ** (Weighting Fair Queuing) – взвешенное справедливое обслуживание, позволяет динамически управлять ресурсами

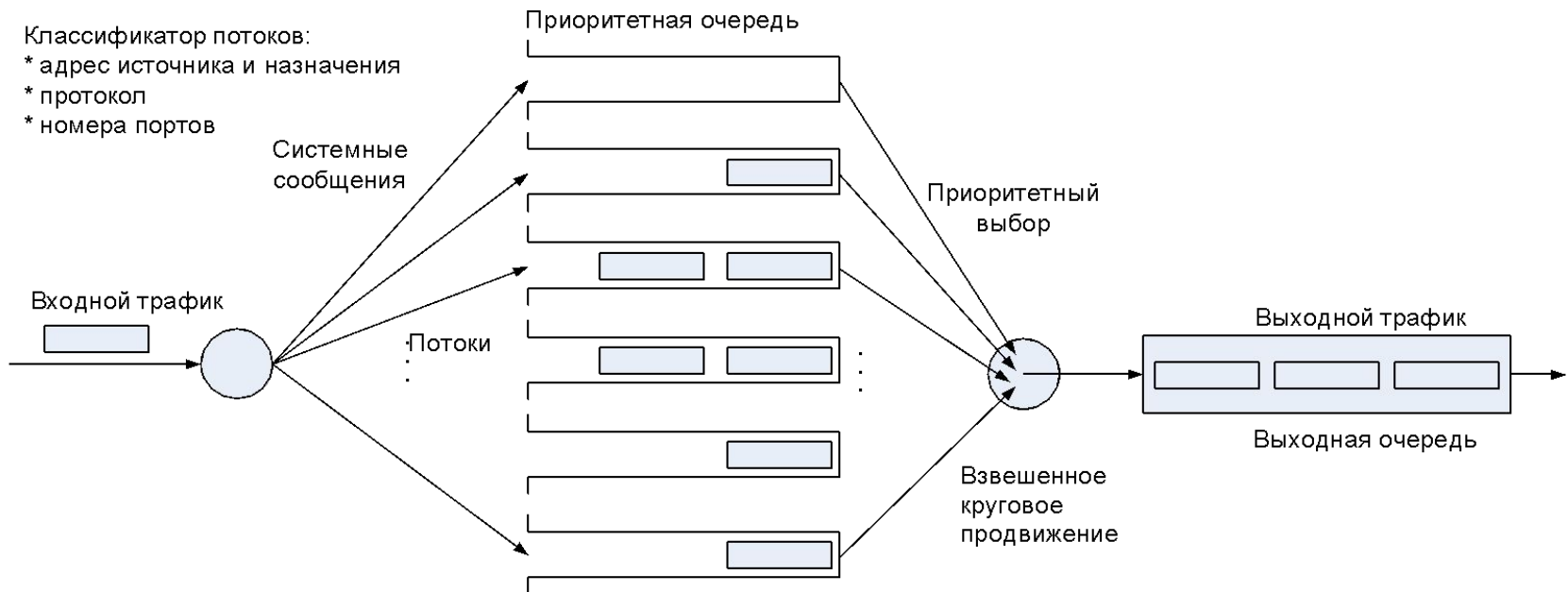
Приоритетное обслуживание



Взвешенные настраиваемые очереди



Взвешенное справедливое обслуживание (WFQ)



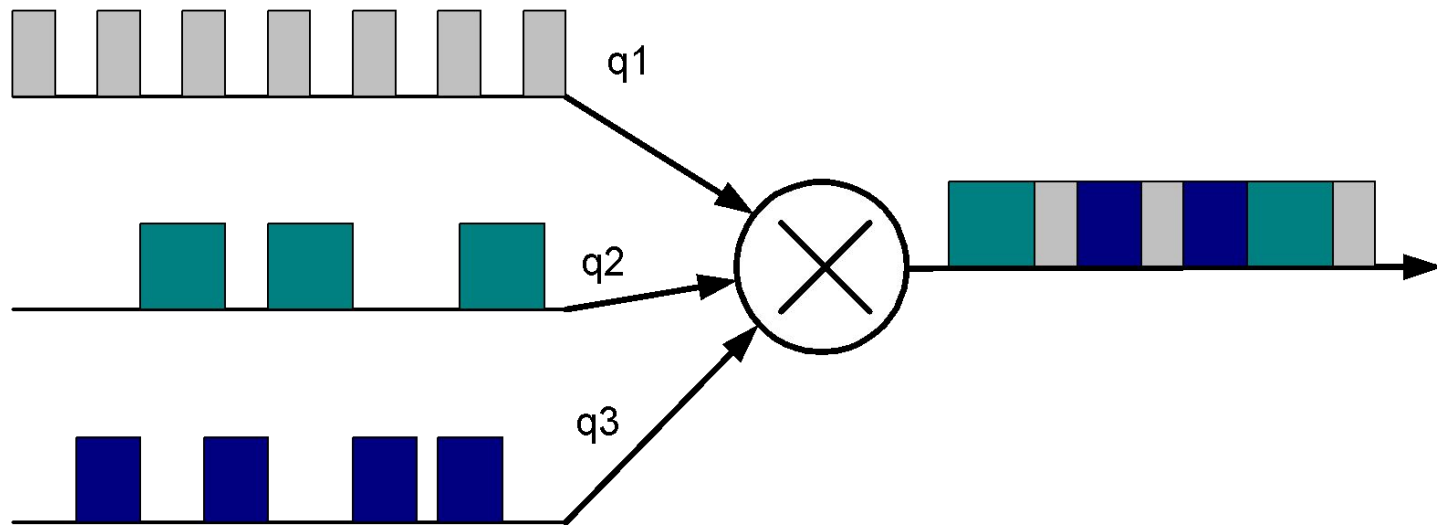
Организация очередей WFQ

Приоритет:

7-8 сигнализация, транзакции

5-6 трафик реального времени

1-4 эластичный трафик



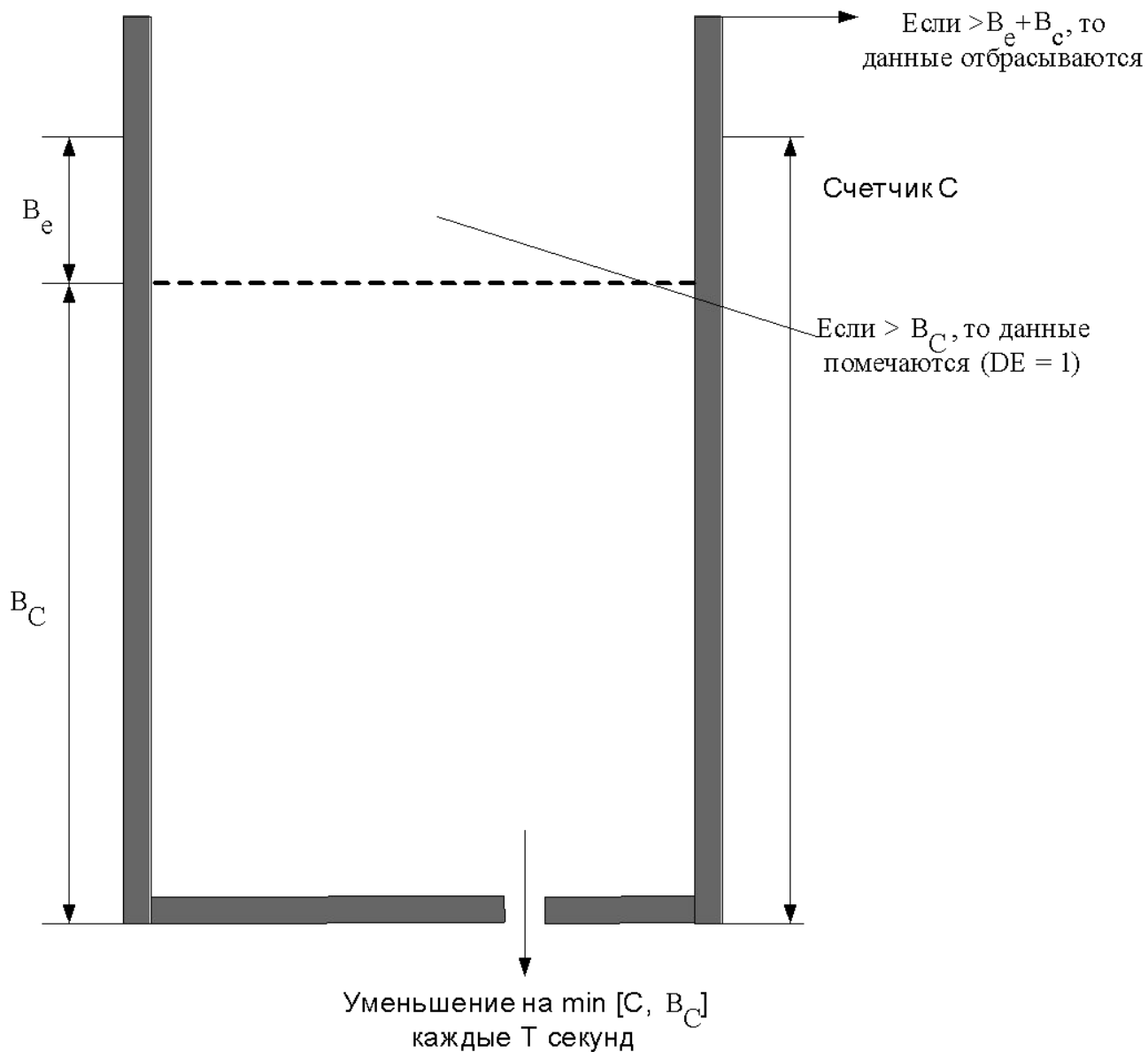
Модификации WFQ

- WFQ на основе вычисления номера пакета
- WFQ на основе потока
- CBWFQ – WFQ на основе класса
- DWFQ – распределенный WFQ
- DWFQ на основе QoS-группы
- CBWFQ с приоритетной очередью (LLQ)
- Заказное обслуживание очередей

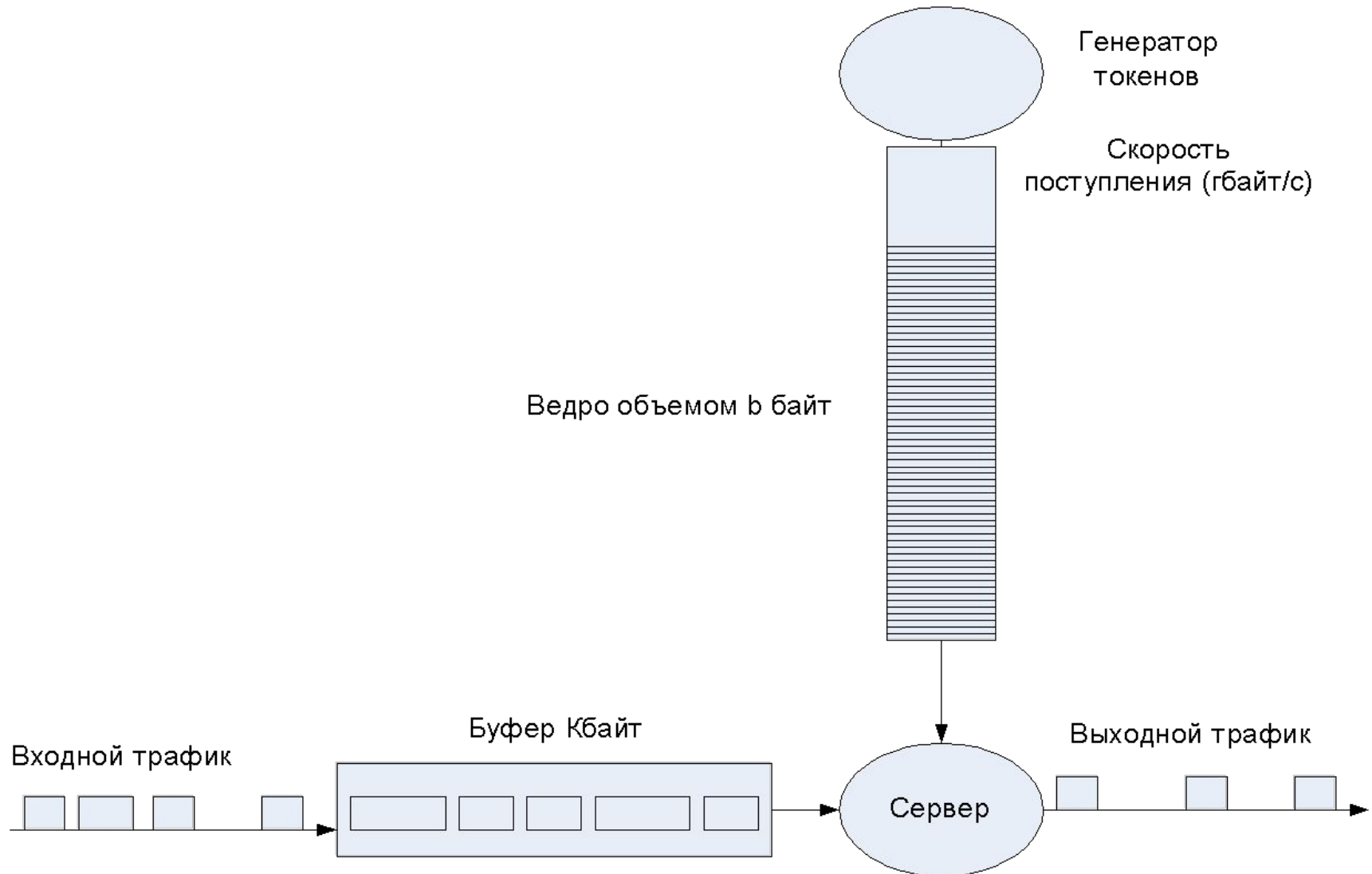
Механизмы профилирования трафика

- **Drop tail** – отбрасывание хвоста: отбрасываются все пакеты, заставшие буфер полным. Используется в best effort.
- **RED** – случайное раннее обнаружение: при угрозе перегрузки пакеты из буфера отбрасываются с ненулевой вероятностью.
- **Дырявое ведро** – отбрасываются пакеты, не обслужившиеся за установленный период.
- **Корзина маркеров (токенов)** – дозирование трафика с целью уменьшения неравномерности продвижения пакетов

Алгоритм "дырявого ведра"



Алгоритм "ведро токенов"



Управление потоками

- **Прерывание передачи:** при перегрузке передача пакетов источниками трафика прерывается на случайный интервал времени, затем возобновляется с той же интенсивностью.
- **Использование динамического окна:** размер окна (количество пакетов, посылаемых источником за период) изменяется в зависимости от загрузки буфера.
- **Медленный старт:** в случае перегрузки источники трафика прекращают передачу, затем посылают пакеты, постепенно увеличивая размер окна.

Модели обеспечения качества обслуживания в сетях IP

- **Модель предоставления интегрированных услуг (IntServ)**
RFC-2205, 1994-1997 г.
- **Модель предоставления дифференцированных услуг (DiffServ)**
RFC 2475, 1998 г.
- **MPLS (Multi-Protocol Label Switching)**

Интегрированные услуги IntServ

Разработана IETF, 1994-1997 г.

RFC 2205, RFC 2210, RFC 2211, RFC 2212

Цель: предоставление приложениям возможности запрашивать сквозные требования у ресурсов.

Недостатки: проблемы масштабирования.

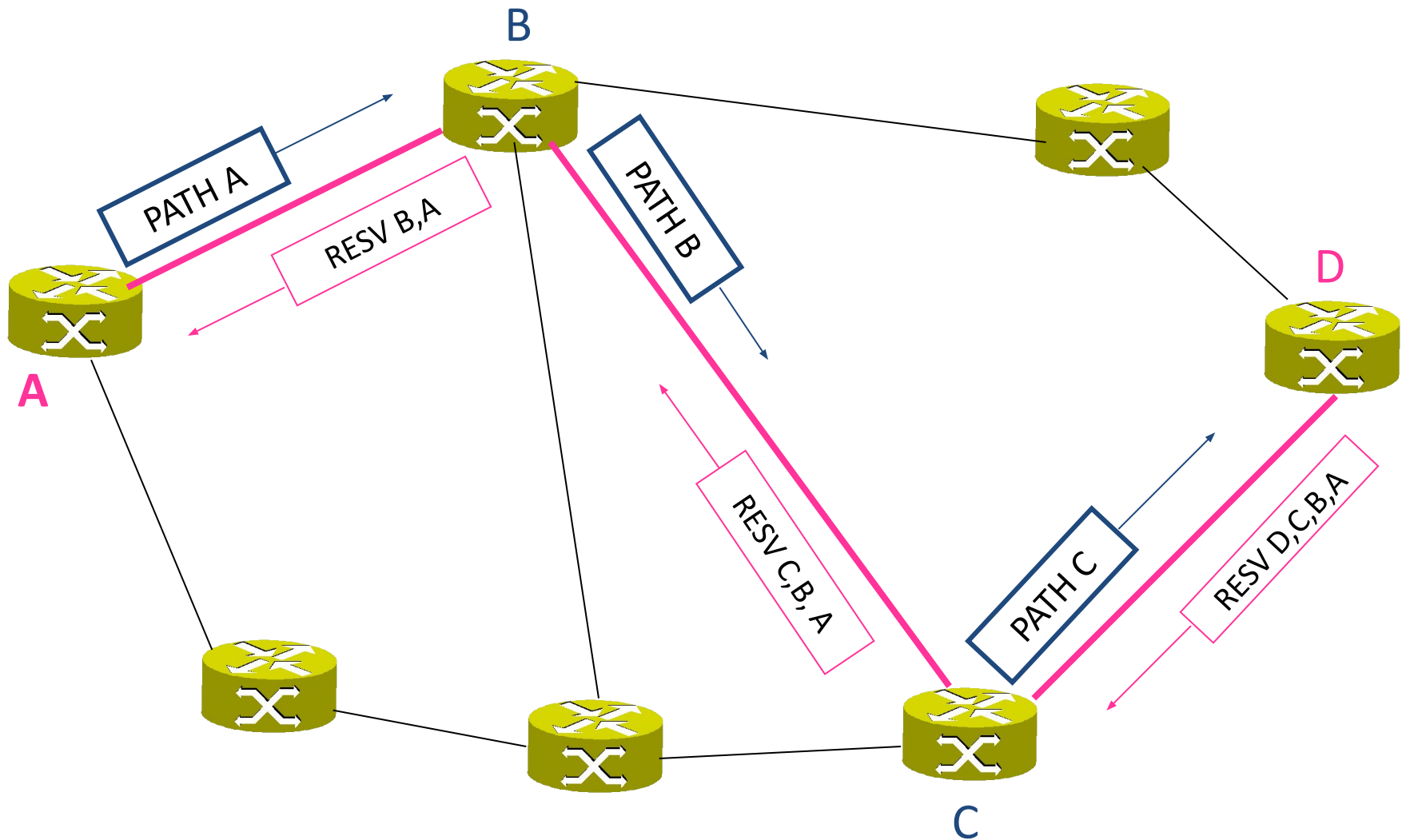
Основной механизм: протокол резервирования ресурсов **RSVP**, в узлах используется WFQ.

RSVP –

Resource Reservation Protocol

- Протокол резервирования ресурсов. Позволяет посылать в сеть информацию о требованиях QoS для каждого потока. Работает совместно с IP.
- Резервирование проводится по адресу получателя. В случае отказа маршрута резервирование происходит заново.
- Работает с двумя видами сообщений:
 - PATH: запрос на резервирование. Содержит:
 - скорость передачи данных;
 - максимально допустимый размер пульсации трафика.
 - RESV: запрос резервирования. Содержит:
 - скорость передачи данных;
 - максимально допустимый размер пульсации трафика.
 - QoS

Организация RSVP-пути



Процесс резервирования пути

- Узел-отправитель посылает запрос PATH как обычный пакет.
- Каждый маршрутизатор прописывает в своей памяти адрес предыдущего и посылает свой адрес в PATH-запросе.
- Получатель в ответ на PATH генерирует RESV и отправляет по прописанному в PATH пути. Т.о. резервирование происходит в обратном порядке, от получателя к отправителю.
- Маршрутизаторы обрабатывают RESV-запросы, пытаясь предоставить требуемые ресурсы. В случае невозможности предоставления ресурсов резервирование начинается сначала.
- Путь считается установленным, когда отправитель получает RESV. После этого начинается сеанс.

Дифференцированные услуги

DiffServ

Разработана IETF, 1998 г.

RFC 1349, RFC 2475, RFC 2597, RFC 2598

Цель: поддержка легко масштабируемых дифференцируемых в Internet

Недостатки: отсутствие гарантированного QoS

Основной механизм: маркировка трафика с использованием бита ToS (Type of Service).

Поддерживает политики поведения сетевого узла: AF-phb и EF-phb (Per-Hop Behavior)

Политики поведения сетевого узла - phb

- **AF-phb** (Assured Forwarding): политика гарантированной доставки – средство, позволяющее обеспечить несколько различных уровней надежности доставки IP-пакетов.

Механизмы: эффективное управление полосой пропускания за счет организации собственной очереди для каждого типа трафика; 3 уровня приоритетов пакетов; RED.

- **EF-phb** (Expedited Forwarding): политика немедленной доставки – обеспечение сквозного QoS для приложений реального времени.

Механизмы: приоритезация трафика; WFQ; распределение ресурсов; RED.

MPLS

(Multi-Protocol Label Switching)

Разрабатывается IETF

RFC 2702, RFC 2283, RFC 2547

Цель: отделение процесса маршрутизации пакета от необходимости анализа IP-адресов в его заголовке, что существенно уменьшает время пребывания пакетов в маршрутизаторе и обеспечивает требуемые показатели QoS для трафика реального времени.

Недостатки: ориентирован на топологию

Основной механизм: коммутация по меткам, туннелирование