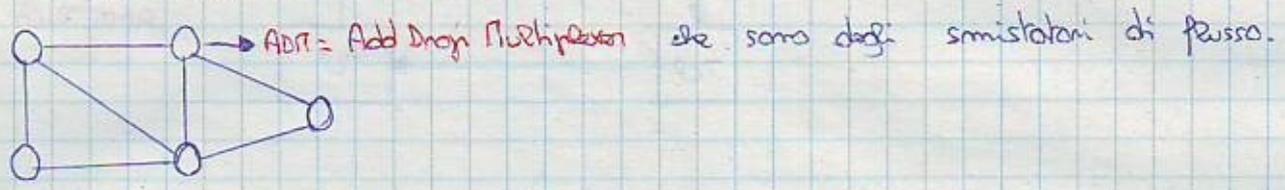


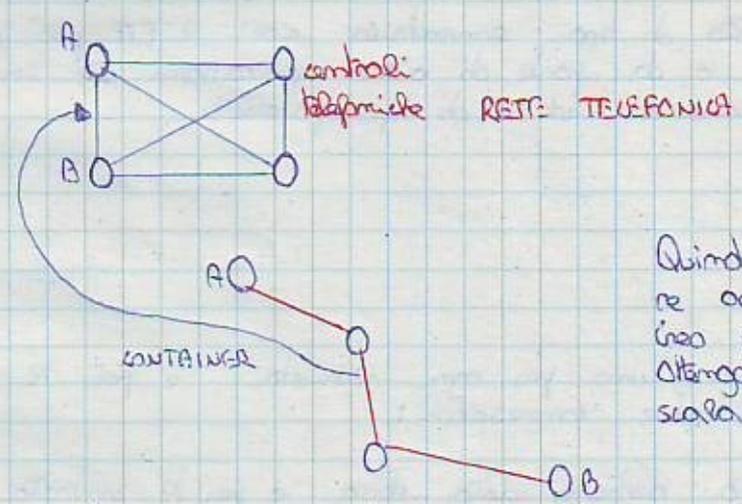
Succede che abbiamo un accodamento intenso e quindi i pacchetti vocali subiscono un ritardo notevole. Siccome:

diventa inutile la lotta che tali pacchetti arrivano lo stesso a destinazione.  
DELAY  $\leq$  10-20 ms

Questa cosa non accade sulle reti telefoniche dorsali cioè quelle che funzionano su un circuito cioè quelle che hanno un canale dedicato. Questo perché il ritardo è deterministico e piccolo. Un problema che rappresenta una sfida per il futuro è il problema legato alla capacità di garantire la qualità del servizio. Tutti sanno che gli ISP coincidono con i vecchi operatori di telefonia. Questi ultimi però hanno una rete. Questi hanno una rete fisica che è costituita da un insieme di nodi e di link:



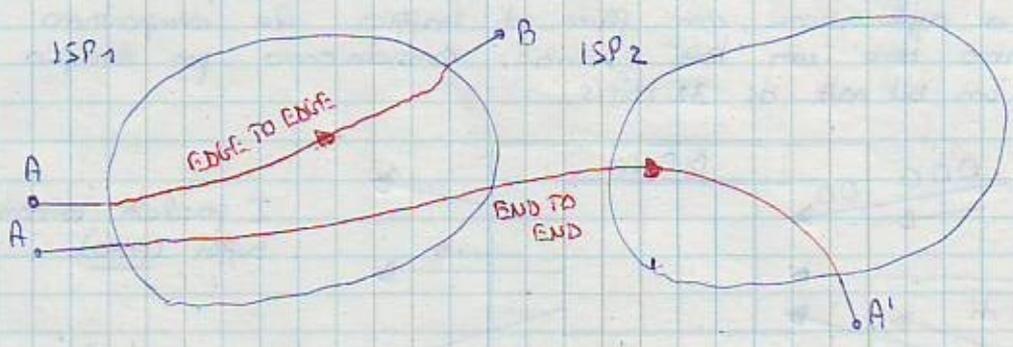
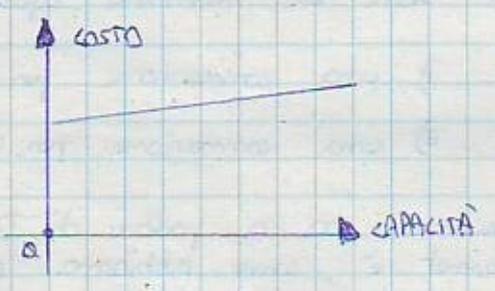
ed hanno una rete telefonica che è sovrapposta a quella fisica. Quindi:



Quindi con una rete fisica posso servire molte più reti fisiche sovrapposte. Crea un container che unisce tali reti. Ottengo così una grande economia di scala.

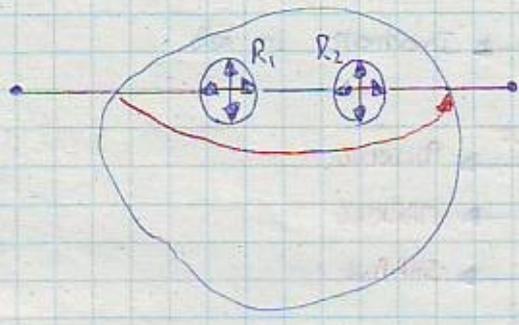
Quindi se raddoppio la capacità, le costi non raddoppia.

La gestione del traffico è la chiave per garantire la qualità del servizio ma è piuttosto complicata. Consideriamo:



Un conto è garantire la qualità del servizio END TO END ed un conto è garantire la

qualità del servizio EDGFA EDGFA. È più facile garantire la qualità di servizio per EDGFA EDGFA. Consideriamo:

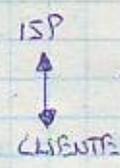


Ogni router avrà l'occasione di non avere la qualità del servizio a causa delle proprie politiche di buffering. Deve quindi garantire la qualità del servizio a livello locale.

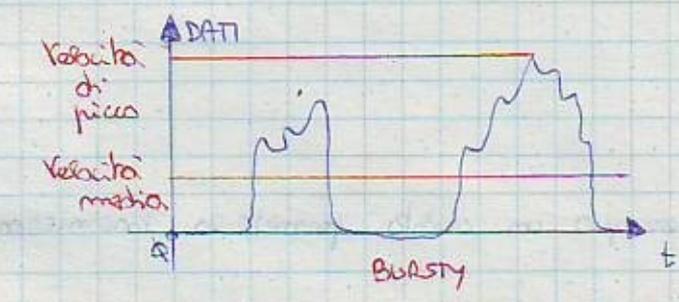
Per prima cosa bisogna specificare quale qualità di servizio si vuole garantire. Quindi:

Questa SLA dipende dalla banda minima richiesta se si pensa di dati, o dipende dalla soglia di ritardo se si pensa di telesema.

Per quanto riguarda per esempio le VIDEO/AUDIO STREAMING bisogna garantire sia la banda sia la soglia. Lo streaming purtroppo non può essere fatto in real-time. Quindi il contratto tra cliente e ISP viene specificato attraverso SLA. Bisogna poi anche specificare il TRAFFIC PROFILE che ci dice come è fatto il traffico inserito dall'utente. Si ricorda che la rete deve passare dalla struttura in grado di reggere un determinato TRAFFIC PROFILE se abbiamo una sorgente di dati si ha:



SLA: Service Level Agreement mi definisce la qualità del servizio che l'utente si aspetta.

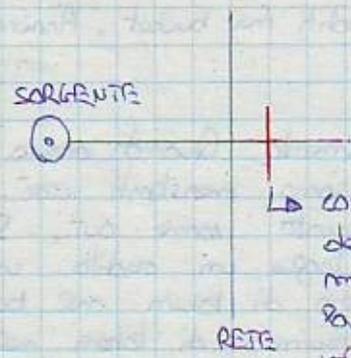


Per esempio: traffico vocale.

La velocità che incide sulla banda è la velocità media, e non quella di picco. Definiamo:

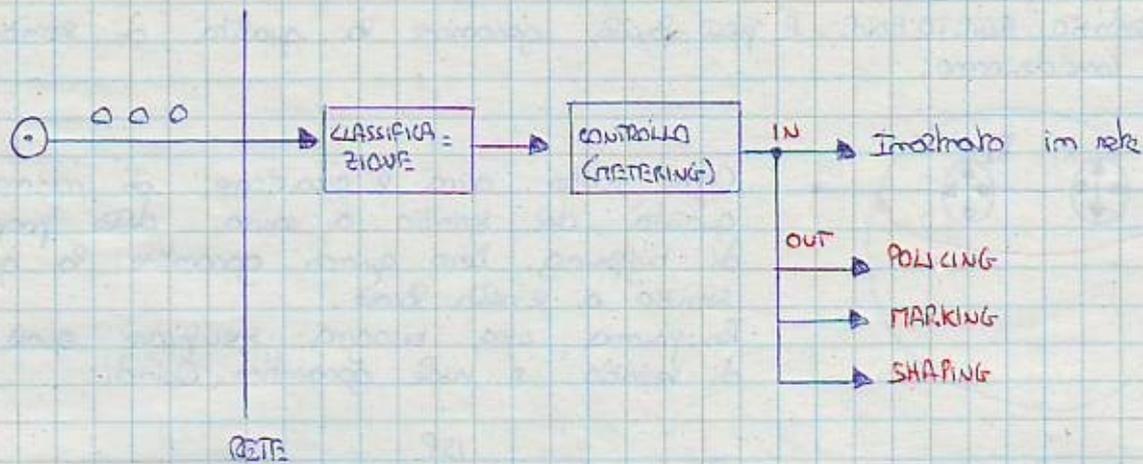
$$\text{BURSTINESS} = \frac{V_{\text{picco}}}{V_{\text{media}}}$$

Per abbiamo sorgenti domestiche più "adattabili". Per gestire una singola sorgente di traffico bisogna controllare il traffico.

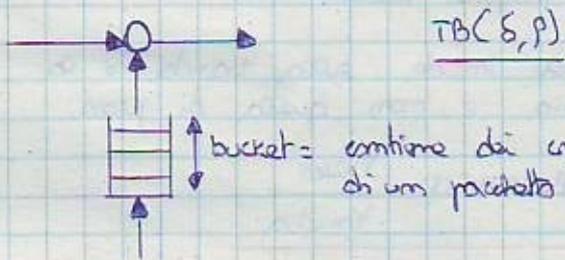


CONTROLLO DEL TRAFFICO. per fare in modo che il traffico sia dentro o delle specifiche ben definite. Tale controllo mi controlla la velocità media, la velocità di picco, la lunghezza massima dei burst dei pacchetti ossia velocità di picco. fissa a loro su cui che è stato dichiarato nel traffic profile.

Inoltre il controllo contiene e identifica dall'utente che viene o sia, rete controllata. Vediamo brevemente il controllo del traffico immesso in rete:



La classificazione serve per vedere a quale controllo appartiene il pacchetto. Il controllo serve per capire se il singolo pacchetto è aderente al profilo di traffico dichiarato. Se un pacchetto è fuori dal profilo, posso ammorzarlo attraverso il POLICING, oppure posso marcarlo come pacchetto fuori profilo mettendo il campo TOS nel header IP. I pacchetti marcati possono subire se la rete è libera altrimenti vengono buttati. Alternativa, mentre un pacchetto può essere ritardato in un buffer attraverso lo shaping. Quest'ultima operazione marca avanti il pacchetto solo quando è possibile farlo rimanendo aderenti al profilo del traffico. Vediamo ora un meccanismo che realizza il marking e che viene chiamato **token bucket**. Graficamente si ha:

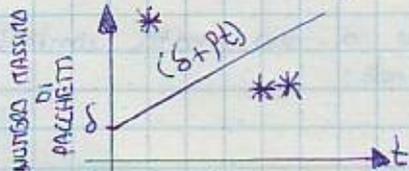


bucket = contiene dei crediti dove per esempio un credito permette la trasmissione di un pacchetto

Il bucket ha due parametri:

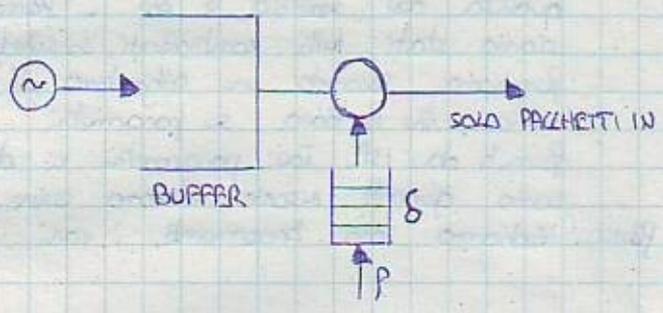
- 1)  $S$ : che mi dà il numero massimo di crediti accumulabili
- 2)  $P$ : mi fornisce la velocità di arrivo dei crediti nel bucket. Arriva un credito ogni  $1/P$  secondi.

Si noti che  $S$  si misura in crediti. Vediamone il funzionamento. Quando arriva un pacchetto si va a controllare i crediti nel bucket. Se tali crediti sono inesistenti cioè se il numero di crediti nel bucket è zero, allora il pacchetto viene marcato come out. Se nel bucket c'è almeno un credito il pacchetto passa come in e si toglie un credito dal bucket. Se mi arriva un token (credito) e se il numero di token nel bucket è minore di  $S$ , aumento di un credito il bucket altrimenti se il numero di token nel bucket è uguale a  $S$ , non posso ricevere nuovi crediti e il bucket rimane così. Ci chiediamo ora, quale è il numero massimo di pacchetti in permesso in funzione del tempo.



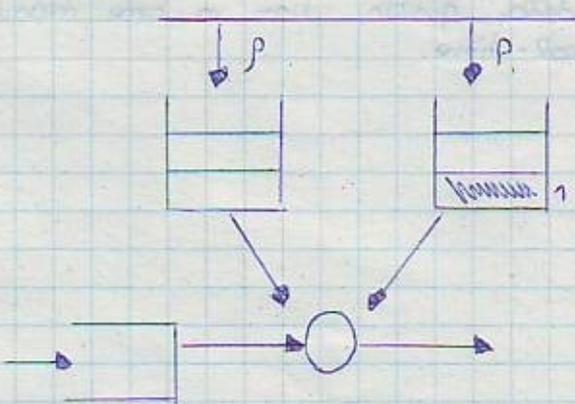
\* = OUT OF PROFILE  
 \*\* = IN PROFILE

Il processo di **marking** quindi ci consente di verificare le burst dei  $\delta$  pacchetti in profilo e la velocità media dei pacchetti in pari a  $P$  pacchetti al secondo. Il profilo di traffico è quindi un token bucket con un certo  $P$  e un certo  $\delta$ . Analizziamo ora cosa succede nel caso delle operazioni di **shaping**. Arriviamo di un nuovo pacchetto se il buffer ha già dei pacchetti quest'ultimo viene accodato. Se il buffer è vuoto resta a vedere il contenuto del bucket. Se il bucket ha almeno un credito allora passa il pacchetto come in  $1$  e si toglie un credito, altrimenti il pacchetto resta in coda.



Quando arriva un token se il bucket è vuoto resta a vedere il buffer di ingresso. Se il buffer ha dei pacchetti si serve un pacchetto e non si incrementa il bucket. Se il buffer è vuoto allora incrementa il bucket di un credito. Se il bucket ha meno di  $\delta$  crediti incrementa il bucket e se il bucket ha  $\delta$  crediti si comoda il nuovo credito.

Tutte queste operazioni rendono prevedibile il traffico che entra in rete. Questo permette anche la allocazione delle risorse. Ma come si limita la velocità di picco, visto che il token bucket non lo fa? Si affiancano due token bucket.

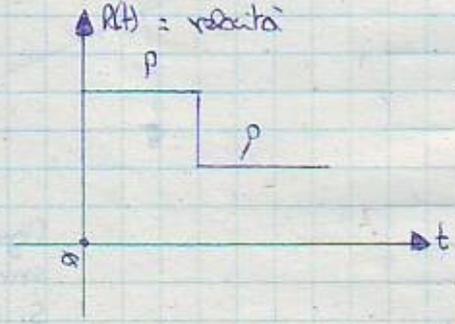
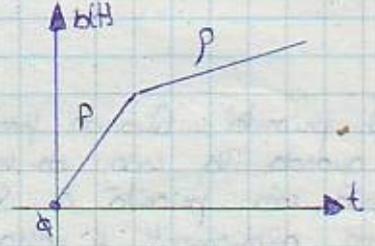


$P = p$  con  $P =$  velocità di picco  
 $p =$  velocità media

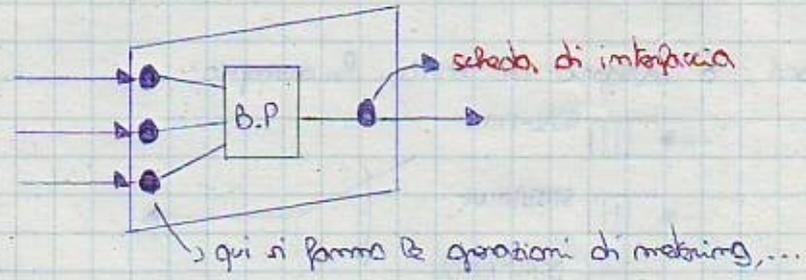
Entrambe queste velocità vengono controllate.

Un pacchetto passa se è presente almeno un token in entrambi i bucket.

Sia  $bit$  il numero di bit trasmessi.

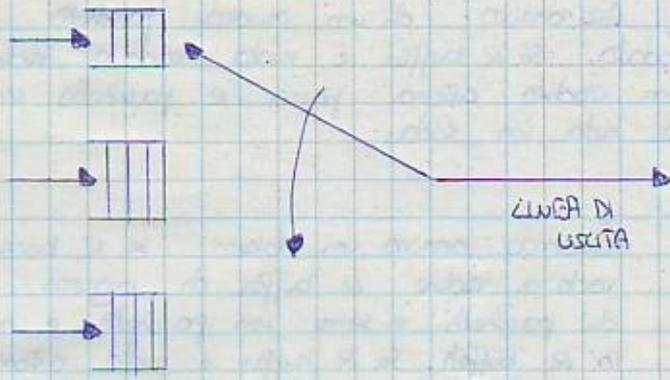


Definiamo **edge router** un router che controlla se il utente si comporta bene rispetto al controllo. Schematicamente si ha:



Nella scheda di interfaccia c'è lo scheduler.

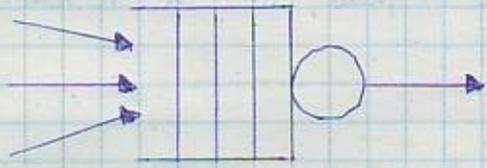
Lo scheduler sostanzialmente è quel dispositivo che decide quale coda servire e lo fa seguendo delle specifiche politiche.



È qui che si determina la qualità dei servizi ricevuti dai flussi. Una condizione necessaria e sufficiente per la garanzia della qualità del servizio è che i flussi siano stati tutti controllati. Lo scheduler funziona secondo un algoritmo di scelta che si basa su parametri forniti da ISP. Tali parametri si dicono come quante risorse devono essere

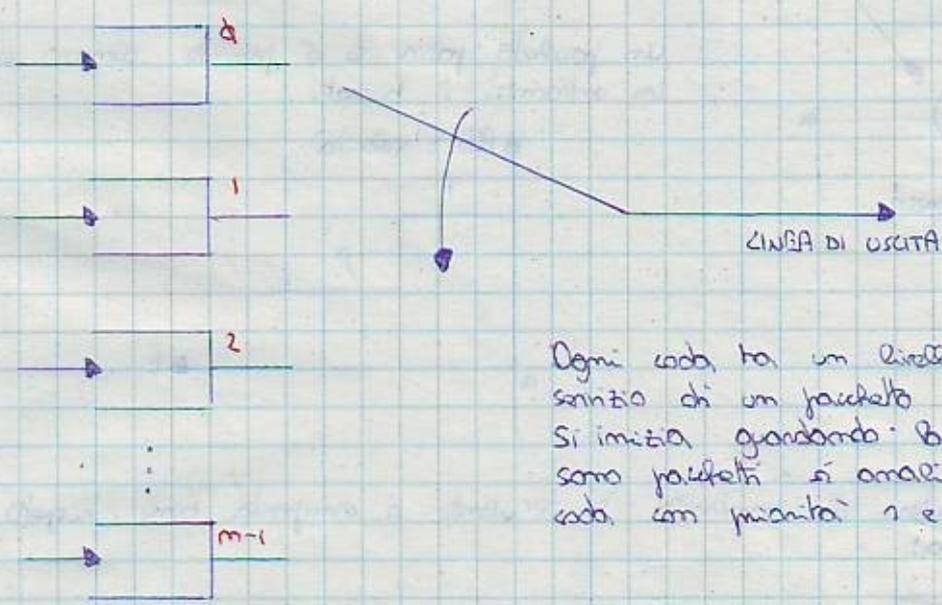
assegnate alla singola coda e quindi al singolo flusso. Vediamo ora brevemente i vari tipi di scheduler:

1) FIFO



Questo scheduler non serve a niente. Va bene per il servizio best effort ma non va bene per il controllo e servizio della qualità. Non va bene neanche per il traffico real-time.

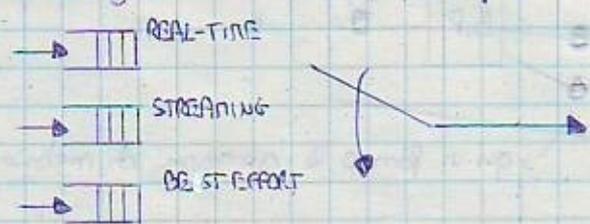
2) STATIC PRIORITY



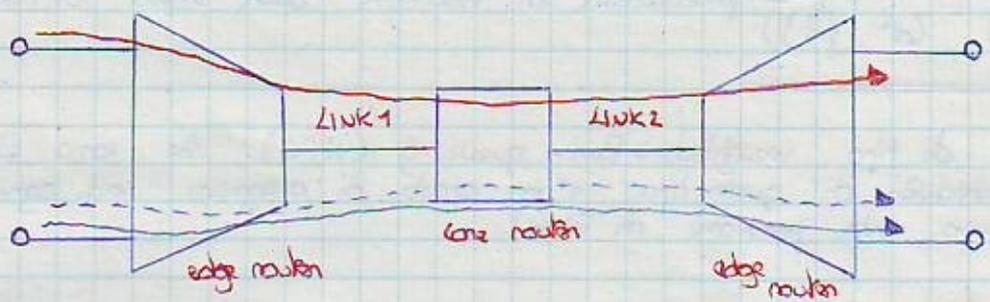
Ogni coda ha un livello di priorità. Quando finisce il servizio di un pacchetto si guarda la coda da servire. Si inizia guardando la coda con priorità 0. Se ci sono pacchetti si analizzano, altrimenti si passa alla coda con priorità 1 e così via.

Così facendo posso anche definire classi o categorie di servizi. Per esempio:

- 1) SERVIZI REALTIME
- 2) SERVIZI DI DATI O STREAMING
- 3) BEST-EFFORT



Alcun ingresso in rete i pacchetti vengono classificati per cui oltre ad individuare le traffic profile viene identificato anche il service profile all'interno del contratto. In merito il campo TOS del header IP viene usato per definire 3 livelli di priorità definendo così tre categorie di servizi. Consideriamo ora:

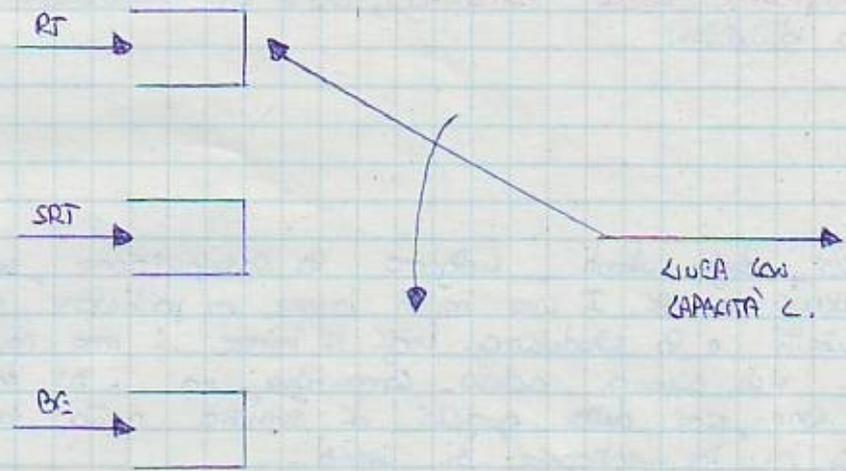


{ ———▶ = Flusso real-time  
 { - - -▶ = Flusso streaming  
 { ———▶ = Flusso best-effort.

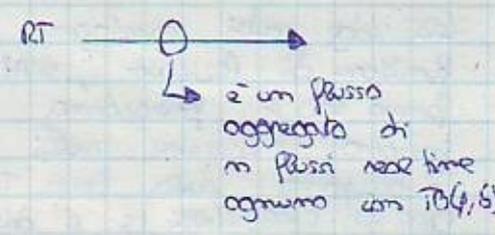
Come dimensioniamo i link e i router in modo che la qualità di servizio sia rispettata? Nel momento di picco cioè nel momento di traffico più intenso si prevede di avere  $N_{RT}$  richieste di servizio real-time,  $N_{ST}$  richieste di servizio di tipo streaming ed  $N_{BE}$  richieste best-effort. Progettare la capacità del link 1 ed del link 2 che sono rispettivamente  $C_1$  e  $C_2$  e scegliere la parametrizzazione degli scheduler in modo da garantire la qualità di servizio nella situazione di massimo traffico. In merito bisogna attuare la procedura CAC = Connection Admission Control che verifica all'ingresso le richieste che eccedono i limiti massimi per i quali la rete è stata progettata. Abbiamo quindi due funzioni di controllo in real-time:

- 1) CAC
- 2) POLICING / SHAPING / MARKING

Così in ogni punto di accesso alla rete ogni flusso è dimensionato. Si noti che le richieste si ha quando arrivano i flussi e non quando li diciamo. Ci chiediamo ora però come posso parametrizzare la linea di uscita e lo scheduler. Supponiamo di avere uno scheduler STATIC PRIORITY:



Dobbiamo garantire i ritardi prestabiliti data questa ordine di servizio. Si noti che:

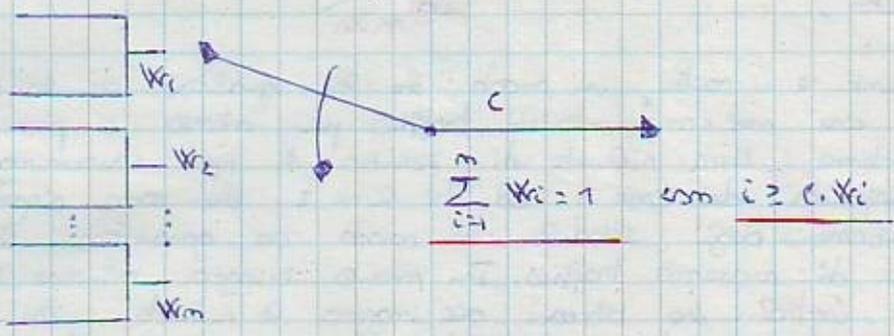


Per ogni flusso definisco l'aggregato massimo. Esiste una formula che, dati il numero di classi di servizio e data una  $TB(p_i, d_i)$  che descrive le proprietà del traffico della classe di servizio  $i$  nel caso di traffico di picco, e dati la qualità di servizio:  $P(D_i > d_i) \leq p_i$ , ci dà la capacità della linea di uscita data lo scheduler di tipo STATIC PRIORITY. Purtroppo lo STATIC PRIORITY funziona con poche classi di servizio, ma si possono trasportare servizi di real-time. Si noti che:

$D^i$  = classe del traffico  
 $d^i$  = soglia  
 $p^i$  = probabilità

Quindi:  $C = P(TB(p, \delta), (d^i, p^i))$    
  $\nearrow$  soglia di ritardo   
  $\nearrow$  probabilità di violazione della soglia di ritardo.

Esistono anche scheduler di tipo **weighted - fair - queuing (WFQ)** che sono scheduler di tipo RATE BASED. Uno scheduler di questo tipo è in grado di assegnare ad ogni coda una certa banda (cioè una certa porzione di link).



Come si può notare è uno scheduler più evoluto. Tutti i router hanno i due tipi di scheduler menzionati. Inoltre questa scheduler consente di gestire il flusso singolo e classificare i flussi al bordo della rete e mantenere e' identità di ogni singolo flusso. Purtroppo esso costa e soffre di scalabilità. Vediamo ora l'architettura per le reti IP con qualità di servizio. Sistema:

- servizi IP integrati (INTSERV) (IS)
- servizi IP differenziati (DIFFSERV) (DS)

Entrambi possono usare MPLS (Multiprotocol Label Switching), che è un meccanismo di routing aggiuntivo. I router vengono divisi in:

- 1) edge router
- 2) core router (router interni)

Gli edge router ricevono il traffico degli utenti, svolgono la classificazione, svolgono funzioni di POLICING, SHAPING, MARKING e CAR. I core router invece, in particolare nei DS fanno solo il forwarding dei pacchetti e lo scheduling. Negli IS invece i core router possono intraprendere tutte le funzioni visto appena adesso. Comunque sia i DS che i IS si occupano della EDGE TO EDGE QoS, cioè della qualità di servizio nella connessione EDGE-TO-EDGE. Negli IS si distinguono in tre categorie di servizi:

- 1) Guaranteed Service (GS)
- 2) Controlled Load Service (CLS)
- 3) Best effort (BE)

Ogni singolo flusso può appartenere a uno di questi tre servizi. C'è una gestione