

**Centro de Investigación Científica y de
Educación Superior de Ensenada**



**MANEJO DE RECURSOS MEDIANTE LA ASIGNACION
DINAMICA DE ANCHO DE BANDA EN REDES
DE COMUNICACION DE DATOS ORIENTADO
A LA SATISFACCION DE GARANTIAS DE
CALIDAD DE SERVICIO**

**TESIS
MAESTRIA EN CIENCIAS**

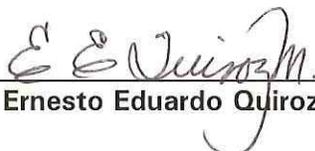
MARLENNE ANGULO BERNAL

Ensenada, Baja Cfa., Mexico.

Diciembre de 1999.



TESIS DEFENDIDA POR
Marlene Angulo Bernal
Y APROBADA POR EL SIGUIENTE COMITÉ



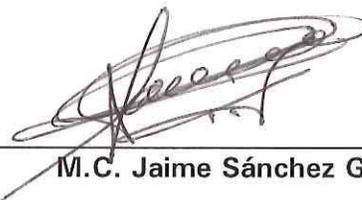
M.C. Ernesto Eduardo Quiroz Morones

Director del Comité



Dr. Jesús Favela Vara

Miembro del Comité



M.C. Jaime Sánchez García

Miembro del Comité



Dr. José Luis Medina Monroy

*Jefe del Departamento de Electrónica
y Telecomunicaciones*

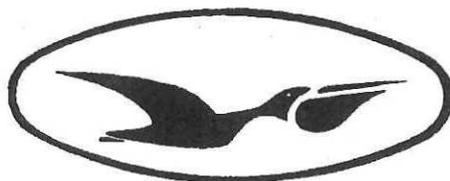


Dr. Federico Graef Ziehl

Director de Estudios de Posgrado

13 de diciembre de 1999

CENTRO DE INVESTIGACION CIENTIFICA Y DE EDUCACION
SUPERIOR DE ENSENADA
DIVISION DE FISICA APLICADA



C I C E S E

DEPARTAMENTO DE
ELECTRONICA Y TELECOMUNICACIONES

MANEJO DE RECURSOS MEDIANTE LA ASIGNACION DINAMICA DE
ANCHO DE BANDA EN REDES DE COMUNICACIÓN DE DATOS
ORIENTADO A LA SATISFACCION DE GARANTIAS DE CALIDAD DE
SERVICIO

TESIS

Que para cubrir parcialmente los requisitos necesarios para
obtener el grado de MAESTRO EN CIENCIAS presenta:

MARLENNE ANGULO BERNAL

Ensenada Baja California, México. Diciembre de 1999

RESUMEN de la Tesis de **Marlenne Angulo Bernal**, presentado como requisito parcial para la obtención del grado de **MAESTRO en CIENCIAS en ELECTRONICA Y TELECOMUNICACIONES**. Ensenada, Baja California, México, Diciembre de 1999.

MANEJO DE RECURSOS MEDIANTE LA ASIGNACION DINAMICA DE ANCHO DE BANDA EN REDES DE COMUNICACION DE DATOS ORIENTADO A LA SATISFACCION DE GARANTIAS DE CALIDAD DE SERVICIO

Resumen aprobado por:



M.C. Ernesto E. Quiroz Moronez
Director de Tesis

El rápido desarrollo de las redes de comunicación de datos y la necesidad de servir diferentes aplicaciones tales como video, voz y datos, etc.; demanda a las redes mas populares como la de Modo de Transferencia Asíncrona (ATM, por sus siglas en inglés), y la red Internet a mejorar su servicio para cubrir los diferentes requerimientos de las aplicaciones. La necesidad de proveer garantías de calidad de servicio, requiere mecanismos de control y vigilancia muy sofisticados, tales como mecanismos de ventana, algoritmo de cubeta con goteo, etc.

En el presente trabajo de tesis, se presenta una modificación al algoritmo de Cubeta con Goteo, esta mejora se denomina algoritmo de Cubeta con Goteo Adaptivo-Ayudado por predicción, el cual permite la predicción de procesos estocásticos de largo rango de dependencia con distribución alfa-estable. Basado en este predictor, el esquema adaptivo de Calidad de Servicio (QoS, por sus siglas en inglés) se diseña para mejorar el desempeño de los sistemas de comunicación que transportan tráfico auto-similar con largo rango de dependencia.

Se implementa el algoritmo de predicción de tráfico alfa-estable así como el algoritmo dinámico de asignación de ancho de banda basado en esta predicción, aplicado a la satisfacción de garantías de Calidad de Servicio. Esto es implementado en redes orientadas a conexión como la tecnología ATM y redes Internet que ofrecen servicios diferenciados (DiffServ), esta implementación se realizó sobre el programa de computo Herramientas de ingeniería para optimización de redes (OPNET, por sus siglas en inglés).

El desempeño del esquema de QoS propuesto se evalúa, obteniendose resultados muy satisfactorios en términos de utilización de recursos y razón de pérdida de información.

Palabras Clave: QoS, ATM, DiffServ, métodos de predicción de tráfico, distribuciones alfa-estables, tráfico auto-similar.

ABSTRACT of the Thesis of **Marlene Angulo Bernal**, presented as partial requirement to obtain the **MASTER OF SCIENCES** degree in **ELECTRONIC AND TELECOMMUNICATIONS**. Ensenada, Baja California, México, December 1999.

DYNAMIC BANDWIDTH ALLOCATION AS A RESOURCE MANAGEMENT STRATEGY AIMED AT SATISFYING QUALITY OF SERVICE GUARANTEES IN DATA COMMUNICATION NETWORKS

The quick development of the communication networks and the necessity to serve different applications such as video, voice, data, etc., requires from the most popular networks in the world (Asynchronous Transfer Mode and Internet), to improve their service, in order to provide different application requirements.

The need to provide QoS guarantees demands very sophisticated control and policing mechanisms for the traffic streams such as window mechanisms, leaky bucket algorithm, etc.

This thesis work presents a modification in the Leaky Bucket policing algorithm which is called Adaptive Leaky Bucket-Prediction helped. This algorithm allows the network to predict the long range dependence stochastic processes with alfa-stable distribution. Based on this predictor the adaptive Quality of service scheme is designed to improve network performance.

It implements an algorithm to predict an alfa-stable traffic and the dynamic bandwidth allocation scheme based on this prediction, applied to satisfy the QoS guarantees, this mechanism is implemented in connection oriented network such as ATM technology and Internet providing Differentiated Services DiffServ, this is implemented by Optimized Network Tools (OPnet).

The performance of the proposed QoS scheme is evaluated, for which very satisfactory results are obtained in terms of resource utilization and buffer overflows.

Keywords: QoS, ATM DiffServ, Traffic prediction methods, alfa-stable distributions, self-similar traffic.

DEDICATORIA

Dedico la presente tesis a dos personas excepcionales:
El Sr. Gregorio Angulo Bernal y la Sra. Hermelinda Bernal de Angulo, de quienes tengo el honor de ser hija y quienes con su amor y responsabilidad son el mejor ejemplo a seguir.

Dedico tambien el presente trabajo de tesis a mis hermanos Gregorio, Lizbeth y Guillermo.
Especialmente a Marco Aurelio Turrubiartes Reynaga.

AGRADECIMIENTOS

A mi director de tesis el M.C. Ernesto Quiroz Morones por su apoyo y esenciales colaboraciones a este trabajo de tesis.

A los miembros de mi comite M.C. Jaime Sanchez y Dr. Jesús Favela por sus muy acertadas correcciones al presente trabajo.

Al Dr. Dimitris Makrakis por todo el apoyo recibido a través del laboratorio *Advanced Communication Engineering Center* en UWO. Por todas sus discusiones, retroalimentación y colaboración a mi trabajo de tesis muchas gracias.

Al M.C. Jose Rosario Gallardo, por su paciencia e invaluable aportaciones a este trabajo de tesis.

A Marco A. Turrubiarres Reynaga, Erica Ruiz, Milka Acosta, Edith García, Juan Pablo Santiago, Leonardo Acho, Maged Zaki. Por supuesto a la *mexican gang*: Alex y Julio.

A mis amigos de siempre Mónica Cebrian, Fernando Vazquez, Irasema Gallo y Juan Pablo Anguas.

Gracias tambien de quienes recibí apoyo moral invaluable durante la elaboración de la tesis: la familia Gallardo, familia López Zatarain, Sra Cristina Alvarez y la familia Fluter.

Al Centro de Investigación Científica y de Educación Superior de Ensenada

Al Consejo Nacional de Ciencia y Tecnología por el apoyo economico recibido

CONTENIDO

	Página
I. INTRODUCCION.....	1
I.1 Antecedentes.....	1
I.2 Motivación y Justificación.....	2
I.3 Objetivo.....	2
I.4 Organización de la tesis.....	3
II. MODELO DE TRAFICO.....	5
II.1 Introducción.....	5
II.2 Modelado de tráfico de voz paquetizada utilizando cadenas de Markov	5
II.2.1 Introducción.....	5
II.2.2 Modelo markoviano de dos estados para fuentes Individuales de voz.....	6
II.3 Modelado de tráfico utilizando procesos auto-similares con Distribución alfa-estable.....	7
II.3.1 Introducción.....	7
II.3.2 Procesos auto-similares con distribución alfa-estable.....	8
II.3.3 Modelado de tráfico para redes de banda amplia.....	12
II.3.4 Predicción de Tráfico.....	13
II.3.5 Análisis de portabilidad del esquema de predicción.....	14
III CALIDAD DE SERVICIO.....	16
III.1 Introducción.....	16
III.2 Definición de Calidad de Servicio.....	16
III.3 Calidad de servicio en redes ATM.....	18
III.3.1 Introducción a redes de Modo de Transferencia Asíncrona ATM.....	19
III.3.2 Categorías de Servicio ATM.....	22
III.3.3 Atributos de Tráfico.....	24
III.4 Calidad de Servicio en redes IP.....	25
III.4.1 Introducción.....	25
III.4.2 Servicios Integrados: El primer paso.....	25
III.4.3 Servicios Diferenciados <i>DiffServ</i>	26
IV MECANISMOS PARA PROVEER GARANTIAS DE CALIDAD DE SERVICIO.....	32
IV.1 Introducción.....	32
IV.2 Mecanismos de Asignamiento para proveer QoS.....	32
IV.2.1 Primer arribo primer servicio (FIFO por sus siglas en inglés).....	32
IV.2.2 Asignamiento de colas por prioridad.....	33

CONTENIDO (continuación)

	Página
IV.2.3 Asignamientos Round-Robin / Colas con servicio Ponderado.....	34
IV.3 Mecanismos de vigilancia y control de congestión.....	35
IV.3.1 Mecanismo de salto de ventana.....	35
IV.3.2 Mecanismo de salto de ventana por disparo.....	35
IV.3.3 Mecanismo de movimiento promedio exponencialmente Ponderado.....	36
IV.3.4 Mecanismo de ventana deslizante.....	36
IV.3.5 Mecanismo de Vigilancia: Cubeta de goteo.....	37
IV.3.6 Mecanismo de Cubeta de Goteo Virtual Modificado.....	38
IV.3.7 Mecanismo Cubeta de goteo Adaptivo.....	39
IV.3.8 Mecanismo de Cubeta de goteo adaptivo ayudado por Predicción.....	41
V DISEÑO DE LOS MODELOS DE LOS NODOS ATM E IP UTILIZANDO ASIGNACIÓN DINÁMICA DE ANCHO DE BANDA.	45
V.1 Tráfico limitado en tiempo.....	45
V.1.1 Tráfico de voz.....	47
V.1.2 Tráfico de video.....	50
V.2 Tráfico no limitado en Tiempo.....	51
V.2.1 Tráfico Ethernet y WWW.....	52
V.3 Parámetros de fuentes de tráfico.....	52
V.4 Diseño del conmutador ATM.....	52
V.5 Diseño del enrutador IP.....	55
VI. IMPLEMENTACIÓN Y RESULTADOS DEL MODELO DE OPNET.....	59
VI.1 Implementación del conmutador ATM.....	59
VI.1.1 Análisis Nodal del conmutador ATM.....	60
VI.1.2 Resultados.....	63
VI.2 Implementación del enrutador IP.....	69
VI.2.1 Análisis Nodal.....	70
VI.2.2 Resultados.....	70
VI.3 Análisis de Resultados.....	77
VII. CONCLUSIONES.....	79
VII.1 Conclusiones.....	79
VII.2 Aportaciones.....	80
VII.3 Trabajo a Futuro.....	81
LITERATURA CITADA.....	82

LISTA DE FIGURAS

Figura		Página
1	Fuente de voz paquetizada	6
2	Esquema simplificado de una red de servicios diferenciados	28
3	Esquema funcional del Nodo limitante	28
4	Esquema de asignación FIFO	33
5	Esquema de asignación de colas por prioridad	33
6	Esquema de asignación de cola con servicio ponderado.....	35
7	Mecanismo de Cubeta con Goteo.....	37
8	Mecanismo de cubeta con goteo virtual.....	39
9	Escenario para esquema de QoS cubeta con goteo ayudado por predicción.....	42
10	Formato de paquete de voz en ATM.....	48
11	Formato de paquete de voz en redes IP.....	49
12	Implementación del modelo del conmutador ATM.....	54
13	Diagrama de flujo del mecanismo de marcación.....	56
14	Modelo del conmutador ATM implementado en OPnet.....	60
15	Porcentaje de utilización del enlace de salida del conmutador ATM.	64
16	Estadística acumulativa de celdas descartadas en la cola de salida....	65
17	Estadística acumulativa de celdas descartadas en el conmutador.....	65
18	Retardo de voz en el conmutador ATM.....	66
19	Retardo de video en el conmutador ATM.....	67

LISTA DE FIGURAS (continuación)

Figura		Página
20	Retardo promedio en el tráfico de voz.....	67
21	Retardo promedio en tráfico de video.....	67
22	Jitter en celdas de voz.....	68
23	Jitter en celdas de video.....	68
24	Modelo del enrutador IP implementado en OPnet.....	69
25	Porcentaje de utilización del enlace de salida del enrutador IP.....	72
26	Bytes descartados en el enrutador.....	76
27	Bytes descartados en la cola de salida.....	76
28	Retardo de video bajo el algoritmo CLB.....	74
29	Retardo de video bajo el algoritmo ALB.....	74
30	Retardo de video bajo el algoritmo ALB-PH.....	75
31	Retardo en voz bajo el algoritmo CLB.....	75
32	Retardo en voz bajo el algoritmo ALB.....	76
33	Retardo en voz bajo el algoritmo ALB-PH.....	76

LISTA DE TABLAS

Tabla		Página
I	Parámetros de modelado para tráfico alfa-estable.....	52
II	Parámetros de trafico.....	52
III	Parámetros de contrato de tráfico.....	52
IV	Tiempo de interarribo en generación de estafeta – ATM.....	54
V	Distribución de buffer – ATM.....	55
VI	Tiempo de interarribo en generación de estafeta - IP.....	58
VII	Distribución de buffer - IP.....	58
VIII	Parámetros de entrada, para fuentes de tráfico modeladas como procesos auto-similares con distribución alfa-estable.....	61
IX	Parámetros de entrada para las fuentes de voz unitarias modeladas por cadenas de markov.....	61

MANEJO DE RECURSOS MEDIANTE LA ASIGNACION DINÁMICA DE ANCHO DE BANDA EN REDES DE COMUNICACIÓN DE DATOS, ORIENTADO A LA SATISFACCION DE GARANTIAS DE CALIDAD DE SERVICIO

I INTRODUCCIÓN

El requerimiento de los usuarios de aplicaciones mas complejas que la simple transferencia de archivos, convierte a la Calidad de Servicio (QoS) en un área en la cual la investigación, el desarrollo e implementación de mecanismos que permitan proveerla se convierten en una *necesidad* en las redes actuales.

Debido a la globalización de las comunicaciones es muy importante que todos los mecanismos que busquen satisfacer una calidad de servicio, presenten escalabilidad en implementación, es decir que no se requiera de cambios drásticos en las redes actuales para proveer la calidad de servicio deseada.

Se observan dos tendencias fundamentales en las redes de datos que proveen QoS: Las redes privadas y la Internet pública; los cuales se adoptan en la presente tesis, donde los escenarios sobre los cuales se implementa el mecanismo de manejo de recursos son: el conmutador de Modo de Transferencia Asíncrona (ATM) y el enrutador de la red Internet que provee servicios diferenciados.

I.1 Antecedentes

Históricamente la calidad de servicio no había sido incorporada en el campo de las redes de telecomunicaciones, debido a que las herramientas que se tenían eran muy primitivas y la implementación de estas en enlaces de alta velocidad, tradicionalmente tenían un impacto negativo sobre el desempeño de la red [Ferguson y Huston, 1998].

Además, no habían sido elaborados mecanismos de medición de desempeño que permitieran a los proveedores demostrar la diferenciación en los niveles de QoS.

El concepto QoS dentro de las redes de comunicaciones de datos, fue introducido por las redes ATM, ya que estas fueron diseñadas para ofrecer servicios para diferentes tipos de tráfico.

1.2 Motivación o justificación

El surgimiento de nuevas necesidades de información por parte de los usuarios, no solo de transmisión de datos, sino también aplicaciones de teléfono, video interactivo y aplicaciones en demanda, requiere que las redes proporcionen servicios no sólo de mejor esfuerzo, sino de diferentes clases.

En el caso del requerimiento de ancho de banda de las aplicaciones, una posible solución que permite ofrecer cierta clase de prioridades, es brindar un margen de ancho de banda mayor al requerido, para *proteger* la aplicación, sin embargo el desperdicio de ancho de banda es un lujo que no se pueden permitir las redes actuales.

De aquí se desprende la principal motivación del presente trabajo de tesis, que consiste en el desarrollo de un mecanismo dinámico de manejo de congestión de redes, el cual permita maximizar u optimizar la utilización de recursos, orientándose a satisfacer los requerimientos de QoS de las aplicaciones de usuario.

1.3 Objetivo

Proponer y evaluar el desempeño de un esquema (algoritmo) de asignación dinámica de recursos enfocado a satisfacer garantías de QoS en redes de comunicación de datos. Este

sistema debe proporcionar escalabilidad, diferenciación de clases de servicio, asignación de colas y manejo de congestión, así como ser capaz de ofrecer servicio a aplicaciones limitadas y no limitadas en tiempo.

1.4 Organización de la tesis.

El resto del documento de tesis se organiza como sigue: En el capítulo 2 se introducen los modelos de tráfico utilizados, simulando las fuentes individuales de voz mediante un modelo markoviano de dos estados [Angulo *et al.*, 2000] y [Beran *et al.*, 1995], también se describen brevemente los procesos autosimilares con distribución alfa-estable que permiten modelar fielmente el tráfico agregado de aplicaciones como: datos en redes Ethernet de área local, datos de red mundial (WWW por sus siglas en inglés), y video [Gallardo *et al.*, 1998b]. Posteriormente se describe un algoritmo rápido de generación de tráfico que permite reducir los tiempos de generación de muestras artificiales de tráfico [Gallardo *et al.*, 1998a]. Por último en este capítulo se describe el algoritmo que estima una muestra futura, basándose en el largo rango de dependencia estadística que presenta el tráfico. El capítulo 3 presenta una introducción a los conceptos de QoS, posteriormente describe brevemente el funcionamiento de la redes: ATM y de Protocolo Internet IP, así como estas implementan el concepto de QoS. El capítulo 4 examina varios mecanismos para proveer garantías de QoS a las redes de comunicación de datos tales como mecanismos de asignamiento y vigilancia. Este capítulo introduce un nuevo algoritmo de vigilancia basado en predicción de tráfico, que permite realizar asignación dinámica de ancho de banda. El capítulo 5 describe las premisas del diseño de un conmutador ATM, y de un nodo limitante IP que provee servicios diferenciados, tales

como requerimientos de ancho de banda y latencia de diferentes tipos de tráfico. En este capítulo se presenta también el diseño estructural de los nodos.

El capítulo 6 describe como los modelos propuestos en el capítulo cinco son simulados mediante computadora utilizando el programa OPNET, y los resultados de la evaluación de desempeño de los modelos antes mencionados. Finalmente el capítulo 7 presenta el análisis de los resultados presentados en el capítulo anterior, así como principales conclusiones, aportaciones y trabajo a futuro como resultado del presente trabajo de tesis.

II MODELO DE TRAFICO

II.1 Introducción

El impacto del tráfico en el desempeño de las colas de espera es una característica dominante de los problemas de ingeniería de tráfico, por lo cual, los modelos de tráfico son de marcada importancia. El utilizar modelos de tráfico real en la simulación de un sistema de comunicación de datos descarta la posibilidad de que el buen desempeño de la red es ocasionado por condiciones ideales de tráfico.

Es por ello que los modelos de tráfico utilizados para la simulación de la propuesta de manejo de recursos mediante asignación dinámica de ancho de banda, son parte importante de la propuesta misma ya que se realiza la evaluación del desempeño de la red bajo condiciones reales de tráfico. Es importante mencionar que los modelos de tráfico utilizados en el presente trabajo de tesis muestran una mayor demanda de los recursos de la red en comparación con los modelos de tráfico anteriores tal como el modelo gaussiano que propone teorema del límite central.

II.2 Modelado de tráfico de voz paquetizada utilizando cadenas de Markov.

II.2.1 Introducción

El amplio desarrollo de las redes basadas en conmutación de paquetes ha estimulado el interés de transmitir voz paquetizada a través de la red, siendo que las fuentes de voz activas generan ráfagas periódicas de paquetes (como se muestra en la figura 1), las

propiedades estadísticas del arribo de paquetes de voz difieren de las propiedades del tráfico de datos que pasa a través de la red de conmutación de paquetes.

Para diseñar una red que permita ofrecer retardos inferiores a los mínimos establecidos para una aceptable reconstrucción de la voz, se requiere la utilización de modelos que modelen las propiedades estadísticas de la voz paquetizada.

II.2.2 Modelo markoviano de dos estados para fuentes individuales de voz

El comportamiento típico de las fuentes de voz paquetizada es como sigue [Daigle y Langford, 1996]:

Una fuente se dice *activa* cuando el interlocutor se encuentra literalmente hablando, durante este periodo la fuente de voz genera paquetes de longitud fija a intervalos regulares, seguido a este se encuentra el periodo de silencio, en el cual la fuente cambia su estado a *inactiva* y no se producen paquetes. Al localizarse estos dos *estados* del proceso, es posible modelar la voz mediante cadenas de Markov, ya que estas nos permiten caracterizar a una variable discreta que cambia en el tiempo, mediante sus probabilidades de transición de estados.

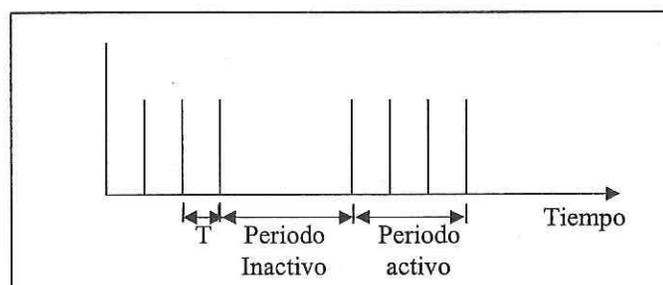


Figura 1. Fuente de voz paquetizada

Una ráfaga de paquetes de una fuente de voz es modelada por arribos a intervalos fijos con una duración de T mseg, y un periodo de silencio en el cual no se tiene arribo (Figura 1). De acuerdo a Heffes [Heffes y Lucantoni, 1986] una fuente de voz es modelada por la siguiente ecuación:

$$F(t) = [(1 - \alpha T) + \alpha T (1 - e^{-\beta(t-T)})] U(t-T) \quad (1)$$

Donde $U(t)$ es la función escalón unitario, α^{-1} es el valor promedio del periodo activo, β^{-1} es el valor promedio del periodo de silencio; ambos periodos son descritos por funciones de distribución exponencial, la tasa de interarribo es de $1/\alpha T$ paquetes y corresponde a una distribución geométrica.

II.3 Modelado de tráfico utilizando procesos auto-similares con distribución alfa-estable.

II.3.1 Introducción

Los modelos tradicionales en fuentes de tráfico eran definidos principalmente por distribuciones gaussianas, sin embargo estos modelos no logran captar la variabilidad que presenta la mayor parte del tráfico real (debido a la limitante de varianza finita), por lo que surge la necesidad de modelos auto-similares no gaussianos.

Los procesos estocásticos alfa-estables y de dependencia estadística de largo plazo (DELP), tales como los procesos llamados ruido estable fraccional lineal (REFL), o ruido estable fraccional logarítmico (REF-log), están cobrando importancia en el escenario de las telecomunicaciones, debido a que ambos introducen el concepto de manejo de largas escalas de tiempo, y modela de manera explícita el tráfico agregado de las redes.

Clasificandose este tráfico en base al comportamiento a ráfagas en: tráfico Ethernet, WWW y video [Gallardo et. al., 1998a]..

El tráfico que exhibe dependencia de largo plazo, es en general mas demandante en términos de recursos de red, tales como ancho de banda y buffer para colas de espera. Sin embargo esta característica del tráfico implica la existencia de información estadística acerca del futuro, contenida en mediciones del pasado y del presente, esta información puede ser explotada usando técnicas de estimación y predicción.

II.3.2 Procesos auto-similares con distribución alfa-estable

Procesos auto-similares

Análisis matemáticos muestran que el origen de la DELP es la varianza infinita de las distribuciones del modelo "ON-OFF" para las fuentes individuales. El término auto-similar, fue definido por Mandelbrot para caracterizar los procesos que son escalables en tiempo, y no pierden sus propiedades estadísticas.

Un proceso auto-similar $X(t)$ satisface la siguiente ecuación:

$$\{ X(at), t \in T \} \stackrel{d}{=} a^H X(t), t \in T \quad (2)$$

La ecuación anterior denota que las variables aleatorias $X(at)$ y $a^H X(t)$ tienen distribuciones idénticas [Taqqu y Samorodnisky, 1994]. El parámetro de auto-similitud H , también conocido como parámetro Hurst es positivo, y su límite superior depende de la

distribución marginal del proceso. En particular $H \in (0, 1/\alpha]$ cuando el proceso tiene distribución alfa-estable.

Variables aleatorias alfa-estables.

El Teorema Generalizado del Límite Central propone las distribuciones alfa-estables para modelar la contribución agregada de muchas variables aleatorias, sin restringir a estas variables a tener una varianza finita (como lo hace el Teorema del Límite Central).

Todas las distribuciones alfa-estables no Gaussianas poseen varianza infinita, y por lo tanto, poseen mayor variabilidad que la distribución Gaussiana.

Las distribuciones alfa-estables carecen de función de densidad de probabilidad en forma cerrada. Se definen mediante su función característica, dada por [Taqqu y Samordnisky, 1994]:

$$\Phi_x(\theta) = \begin{cases} \exp[j\mu\theta - |\sigma\theta|^\alpha \cdot (1 + j\beta \cdot \text{sign}(\theta) \cdot \tan(\pi\alpha/2))]; & \alpha \neq 1 \\ \exp[j\mu\theta - |\sigma\theta| \cdot (1 + j\beta \cdot \frac{2}{\pi} \cdot \text{sign}(\theta) \cdot \ln|\theta|)]; & \alpha = 1 \end{cases} \quad (3)$$

El parámetro α es el índice de estabilidad ($0 \leq \alpha \leq 2$), σ es conocido como el parámetro de escala ($\sigma \geq 0$), μ es el parámetro de corrimiento, y β es identificado como el parámetro de asimetría ($-1 \leq \beta \leq 1$).

Se han definido dos procesos auto-similares con función de distribución alfa-estable con incrementos estacionarios: Movimiento estable fraccional lineal y Movimiento estable

fraccional logarítmico. Los correspondientes incrementos de estos procesos, son referidos como ruido estable fraccional lineal (REFL), o ruido estable fraccional logarítmico (REF-log).

La expresión exacta del proceso Y_j de ruido estable fraccional (REF) es el siguiente [Beran *et al.*, 1995]:

$$Y_j = \int_{-\infty}^{\infty} g(j, x) M(dx) \quad (4)$$

Donde $M(dx)$ es una variable aleatoria alfa-estable y:

$$g(j, x) = \begin{cases} \ln|j+1-x| - \ln|j-x| & ; REF - Log \\ |j+1-x|^{H-1/\alpha} - |j-x|^{H-1/\alpha} & ; REF balanceado \\ [j+1-x]^{(H-1/\alpha)} - [j-x]^{(H-1/\alpha)} & ; REF anti-balanceado \end{cases} \quad (5)$$

En la ecuación anterior, se utiliza la siguiente notación:

$$z^{(a)} \triangleq |z|^a \cdot \text{sign}(z) \quad (6)$$

La cual es válida para cualquier número real z , y para $a \geq 0$.

Nótese que ambos procesos REFL y REF-log son procesos estacionarios, ambos pueden tener picos positivos y negativos con amplitudes similares, por otro lado los trazos de tráfico real tienen los picos negativos menos pronunciados [Gallardo *et al.*, 1998b], por lo cual se tuvieron que realizar modificaciones al modelo de los procesos alfa-estables para algunos de los casos analizados, obteniendo con ello un comportamiento más cercano al real. El modelo modificado consiste en el valor absoluto de los procesos originales

(ecuación 5), donde bajo esta modificación no se pierden las propiedades de auto-similitud y dependencia de largo plazo [Gallardo *et al.*, 1998b]. Una alternativa es truncar todos los valores negativos, pero en este caso el proceso perdería su propiedad de auto-similitud.

Generación rápida de muestras artificiales de tráfico.

En el presente trabajo se utilizó un algoritmo rápido generador de muestras artificiales de tráfico propuesto en [Gallardo *et al.*, 1998a], en este artículo se propone utilizar un modelo auto-regresivo como una aproximación al proceso original (representado por la ecuación 4), basado en el principio de mínima dispersión, el cual provee una generación de muestras de tráfico apegada a la realidad, y además mejora la eficiencia respecto al tiempo de cómputo, dado que las muestras artificiales no se generan a través de una integral.

En este algoritmo, un proceso REF puede ser aproximado por la siguiente expresión:

$$Y_j = \sum_{i=1}^N a_i Y_{j-i} + (\gamma_\varepsilon)^{1/\alpha} u_j \quad (7)$$

Donde N es un número positivo entero que denota el orden del proceso autoregresivo, u_j 's son variables aleatorias idénticamente distribuidas $S_\alpha(1,0,0)$ de acuerdo con la notación utilizada en [Beran *et al.*, 1995]. Finalmente los coeficientes a_i 's y el parámetro de dispersión γ_ε son calculados utilizando el criterio de mínima dispersión. La ecuación 7 es equivalente a decir que para un proceso REF depende del parámetro de

dispersión, y del valor promedio condicional $\hat{Y}_j = \sum a_i y_{j-i}$ que se obtiene de las muestras pasadas $\{y_{j-1}, y_{j-2}, \dots, y_{j-N}\}$. La dependencia de estos dos parámetros (dispersión y promedio condicional) es representada en la siguiente ecuación:

$$Y_j = S_{\alpha}^d \left((\gamma_{\varepsilon})^{1/\alpha}, 0, \hat{Y}_j \right) \quad (8)$$

II.3.3 Modelado de tráfico para redes de banda amplia

Recientes estudios de mediciones de tráfico de banda amplia muestran la presencia de propiedades auto-similares en redes de área local y en redes de área amplia, esto es cierto para diferentes tipos de tráfico incluyendo tráfico de video, WWW, y Ethernet LAN/WAN [Beran et. al., 1995], [Leland et al., 1994], [Willinger et al., 1997]. El tráfico generado por los usuarios de redes de telecomunicaciones puede ser modelado en una forma mas precisa utilizando procesos aleatorios auto-similares con distribuciones alfa-estables para un comportamiento marginal.

El modelo de tráfico propuesto en [Gallardo et al., 1998b] para la generación de datos (Ethernet, LAN/WAN, WWW) es definido por la siguiente ecuación:

$$A_j = m + a \left(|Y_j| - \mu_Y \right) \quad (9)$$

Donde A_j representa el número de arribos de tráfico generados durante el intervalo de tiempo j -th, m es el valor promedio del número de arribos por unidad de tiempo, a es el

factor de escala, Y_j es un proceso REF y μ_Y es el valor promedio de $|Y_j|$. Nótese que μ_Y no depende de j debido a que el proceso Y es estacionario.

El modelo propuesto en [Gallardo *et al.*, 1998b] para tráfico de video se expresa en una forma mas simplificada por:

$$A_j = m + a \cdot Y_j \quad (10)$$

donde m , a , y Y_j tienen el mismo significado que en la ecuación 9.

II.3.4 Predicción de Tráfico

El tráfico que exhibe dependencia estadística de largo plazo es en general mas demandante en términos de recursos de la red, sin embargo esta característica presenta algunas ventajas como la información estadística de las futuras muestras.

En [Gallardo *et al.*, 1999] se propone un algoritmo que estima la muestra de un proceso estocástico de dependencia estadística de largo plazo en el tiempo $(j+L)$ -th, basado en N muestras del pasado como sigue:

$$\hat{Y}_{j+L} = \sum_{i=1}^N a_i \cdot Y_{j-i} \quad (11)$$

Donde a_i 's son constantes y L es un número entero.

El error de predicción puede ser definido por:

$$\varepsilon_{j+L} = Y_{j+L} - \hat{Y}_{j+L} \quad (12)$$

Donde ε_{j+L} es una variable aleatoria alfa-estable de promedio cero y Y_j es un proceso alfa-estable, la ecuación 12 es equivalente a decir que el promedio condicional y el parámetro de dispersión de Y_{j+L} , donde los valores pasados $\{Y_{j-1}, Y_{j-2}, \dots, Y_{j-N}\}$, son dados por \hat{Y}_{j+L} y $\gamma_{\varepsilon_{j+L}}$, respectivamente. Para aplicar el criterio de *Mínima dispersión* los valores de las constantes $a_i, i = 1, 2, \dots, N$, son calculados de forma que el parámetro de dispersión ε_{j+L} se minimice [Gallardo *et al.*, 1999].

II.3.5 *Análisis de portabilidad del esquema de predicción.*

El modelo de predicción para tráfico auto-similar con distribución alfa-estable es un importante punto de apoyo para el manejo de los recursos a la entrada de la red, sin embargo se presenta un conflicto entre la dependencia estadística y el comportamiento de los mecanismos de control de la red, y con ello la disyuntiva de que el algoritmo de predicción presente el mismo grado de utilidad dentro de la red en comparación con los nodos a la entrada de la red [Beran *et al.*, 1995]. Los mecanismos de control presentan el comportamiento de filtro pasa-bajas (formadores, algoritmos de vigilancia), sin embargo esto no es suficiente para eliminar la variabilidad en muchas escalas de tiempo [Jena *et al.*, 1996], de acuerdo con Pruthi *et al.* [Pruthi y Popescu, 1997] el tráfico que pasa a través de la red, no presenta cambios significativos en el parámetro H.

Cuando un proceso fractal pasa a través de un filtro pasa bajas, el proceso resultante tiene dos efectos: el efecto no lineal y el lineal. La parte no lineal se manifiesta como

transiente, un proceso oscilatorio con larga desviación, la característica de este proceso depende del espectro de potencia del proceso incidente, y del filtro paso-bajas. Mientras las estructuras de correlación de dependencia de corto plazo pueden ser significativamente alteradas por los mecanismos de control, la dependencia estadística de largo plazo no es alterada [Jena *et al.*, 1996].

De acuerdo a Erramilli [Erramilli *et al.*, 1996] la función de autocorrelación de un tráfico de paquetes de un proceso X decae en forma hiperbólica.

$$r_x(k) \sim |k|^{-\beta} \quad (13)$$

Esta ecuación se satisface cuando $|k| \rightarrow \infty$ donde $0 < \beta < 1$.

El proceso estocástico que satisface la ecuación anterior se dice que exhibe dependencia estadística de largo plazo.

En el dominio de la frecuencia, la DELP se manifiesta a si misma con una densidad espectral:

$$s_x(\omega) \sim \sum_K r_k(k) e^{ik\omega} \quad (14)$$

se aproxima a $s_x(\omega) \sim |\omega|^{-\gamma}$, cuando $\omega \rightarrow \infty$ donde $0 < \gamma < 1$, por otro lado los procesos de dependencia a corto plazo son caracterizados porque su densidad espectral permanece finita aun cuando $\omega \rightarrow 0$.

Esta distinción estadística, nos permite distinguir entre el tráfico de red medido, y los modelos tradicionales de tráfico.

III CALIDAD DE SERVICIO

III.1 Introducción

Las primeras redes de datos tales como ARPANET o NSFNET tomaban la información y mensajes de capas superiores, fragmentándolas en unidades lógicas denominadas paquetes. Cada paquete era transmitido sin importar que tipo de información transportara; a esta primera clase de servicio se le denomina “*mejor esfuerzo*”; denominado así porque la red realizará el mejor esfuerzo de entregar los paquetes a su destino, sin embargo no garantiza a ningún paquete recursos disponibles (tales como ancho de banda) [Metz Chris, 1999].

Con el desarrollo de las redes de paquetes se empezó a dar acomodo a otras aplicaciones, tales como voz, videoconferencia, multimedia, etc., las cuales tienen requerimientos que no puede satisfacer el esquema de mejor esfuerzo.

Es aquí donde surge la necesidad de brindar servicios con diversas calidades de servicio de acuerdo al tipo de aplicación. ATM fue la primer arquitectura que ofreció QoS para cubrir los diferentes requerimientos de las aplicaciones, dando la pauta a que posteriormente la comunidad Internet propusiera arquitecturas en las cuales se incluye QoS. En el presente capítulo se presenta una breve introducción a QoS, y como incluye en redes orientadas y no orientadas a conexión.

III.2 Definición de Calidad de Servicio

Existen diferentes definiciones de QoS, tantas como tan variados sean los enfoques de los proveedores o de los usuarios de servicios, así como la forma cualitativa o

cuantitativa de medir o percibir la calidad. En [Campbell et. al., 1994] QoS es definido como: El desempeño cualitativo del proveedor de servicios y su capacidad de proveer diferentes niveles de servicio. En el caso particular de redes de comunicación Giroux [Giroux y Sudhakar, 1999], brinda un concepto mas completo de QoS, definiéndolo como la especificación de los requerimientos de una aplicación, los cuales debe satisfacer el sistema que realiza el servicio de transporte de datos para lograr la calidad deseada de la aplicación.

La especificación cualitativa se refiere a como el usuario percibe el desempeño del sistema y es tan subjetivo como la opinión de: bueno, malo, regular. Sin embargo esta especificación puede definirse en forma cuantitativa para el requerimiento de determinados parámetros tales como retardo, jitter, probabilidad de pérdida de celda o paquete.

Una red generalizada que ofrece QoS debe permitir a la QoS ser configurable, predecible y mantenerse sobre todas las capas, para ofrecer un determinado nivel de QoS de extremo a extremo, donde las aplicaciones deben aislarse de la complejidad de las especificaciones y el manejo.

El procedimiento para proveer QoS es el siguiente [Campbell *et al.*, 1994]:

Mediante especificación de QoS se capturan los requerimientos de la capa de aplicación, tales como el desempeño de flujo, el nivel de servicio, la vigilancia del manejo de QoS, y el costo del servicio, los cuales se proveen mediante mecanismos que son seleccionados y configurados de acuerdo a la especificación proporcionada por QoS, la disponibilidad de recursos y el mecanismo de vigilancia de tráfico. Los mecanismos QoS son categorizados en: *Estáticos* y *Dinámicos*, donde los primeros actúan durante la fase de

establecimiento y renegociación, mecanismo definido como provisión. Los mecanismos Dinámicos (control y manejo de QoS) son implementados en la fase de transferencia de datos.

Mecanismo de Provisión.

Realiza la traslación entre la representación de QoS a diferentes niveles de sistemas (operación del sistema, capa de transporte, red), en términos de especificaciones de alto nivel [Aurrecochea *et al.*, 1995], realiza la prueba de Admisión que compara los requerimientos de QoS contra los recursos disponibles del sistema. En caso de que lo especifique el protocolo, realiza la reservación de recursos.

Mecanismo de Control

Estos mecanismos operan bajo diferentes escalas de tiempo [Campbell *et al.*, 1994]. Pueden proveer control de tráfico en tiempo real basados en los niveles requeridos de QoS establecidos durante la fase de provisión de QoS, sus funciones son regular, controlar y sincronizar el flujo, manejando el crecimiento de los flujos del sistema final en forma integrada a la red [Giroux y Sudhakar, 1999].

Mecanismo de Manejo

Este mecanismo compara la QoS monitoreada contra el desempeño esperado y hace un reajuste de recursos (si se requiere), determina si los niveles bajos han fallado para

mantener una QoS en el flujo de tráfico y permite a la aplicación especificar el intervalo sobre el cual pueden ser monitoreados uno o más parámetros.

III.3 Calidad de servicio en redes ATM

III.3.1 Introducción a redes de Modo de Transferencia Asíncrona ATM

ATM es una red de paquetes orientada a conexión en la cual la información es transmitida utilizando celdas de tamaño fijo y pequeño, donde el tamaño fijo reduce la varianza en el retardo de la celda [Pandya y Sen, 1999]. Dos de las principales características de la tecnología ATM son la habilidad para proveer garantías de QoS y el uso de multicanalización estadística de celdas [Kasiolas, 1999], lo cual mejora la utilización de los recursos de la red. Sin embargo, la necesidad de proveer garantías de QoS mientras se aplica la multicanalización estadística en el tráfico, demanda que la red realice una efectiva vigilancia y control en el tráfico.

Las celdas generadas por diferentes emisores son multicanalizadas en enlaces de velocidad general pero no necesariamente elevada. La multicanalización se efectúa por medio de colas de espera. Cuando no se tiene información para enviar se insertan celdas vacías, por lo que ATM es síncrono a nivel físico, pero asíncrono a nivel informacional puesto que el contenido de las cargas útiles no tiene ninguna relación con la posición temporal de las celdas.

Conmutación.

La conmutación ATM es temporal asíncrona y opera en celdas (conexiones virtuales). Las celdas son conmutadas entre un acceso de entrada y uno o varios de salida en función de su identificador y tras la validación del encabezado [Kyas, 1995].

Congestión.

Dado que ATM permite el transporte de tráfico heterogéneo, las fuentes pueden tener desde comunicación interactiva de datos, hasta imágenes de video comprimido, con lo cual este tipo de tráfico puede causar una congestión severa.

El control de congestión juega un papel importante en el manejo efectivo de tráfico en las redes ATM. La congestión puede ser definida como un estado en el cual la red no tiene la seguridad de que la QoS negociada para las conexiones existentes continúe, y además no puede ofrecer una QoS para nuevas peticiones de conexión.

Existen dos mecanismos complementarios para reducir los efectos de la congestión [Pandya y Sen, 1999]:

- 1) El control de admisión de llamada (CAC)
- 2) El control de parámetros de uso (UPC)

El CAC decide durante la fase de conexión si se acepta o no la nueva conexión, dependiendo de los recursos disponibles y de los requerimientos de QoS de la nueva conexión. Si la QoS requerida permite garantizar la QoS para las conexiones existentes y para la nueva conexión, entonces CAC le asigna una porción del ancho de banda que tiene disponible.

El UPC se requiere para asegurar que cada fuente esta enviando el tráfico de acuerdo a los parámetros negociados. La función UPC puede ser definida como un

conjunto de acciones que son tomadas por la red durante la fase de llamada para monitorear y controlar el tráfico ofrecido.

El propósito principal en el manejo de tráfico es que la red sea capaz de detectar violaciones a los parámetros negociados y sea autorizada para tomar la acción apropiada, la cual puede ser marcar o descartar la celda.

Además de UPC y CAC existen otras funciones para el manejo de tráfico [Raj, 1996]:

- *Formación de tráfico.*

Mecanismo que permite delimitar las ráfagas de tráfico a una longitud máxima mediante buffers.

- *Manejo de recursos de la red.*

Este mecanismo se realiza mediante disciplina de asignación y servicio a colas, mecanismos de asignación tales como FIFO, asignación de colas por prioridad, servicio Round Robin y asignamiento de colas con servicio ponderado, estos mecanismos se describen ampliamente en el siguiente capítulo.

- *Control de Prioridad.*

El control de prioridad se realiza mediante el bit de prioridad de pérdida de celda (CLP, del inglés Cell Loss Priority) que se encuentra en el encabezado de la celda ATM. Cuando una celda arriba a un nodo en periodo de congestión y el campo CLP tiene valor uno, esta celda será descartada antes de descartar una celda con CLP cero.

- *Mecanismos de control de retroalimentación*

Solicita a las fuentes incrementar o decrementar sus cargas, dentro de los cuales se pueden mencionar:

- a) Indicación explícita de congestión en forma adelantada.
- b) Velocidad explícita.
- c) Notificación explícita de congestión en retroceso.

III.3.2 Categorías de Servicio ATM

El comité técnico *Manejo de Tráfico* del Foro ATM, se encarga de definir (y redefinir) las categorías de servicio ATM, que hasta marzo de 1999 [Foro ATM, 1999b] eran las siguientes: Tasa de bit constante (CBR), Tasa de bit variable Tiempo no-real (nrt-VBR), tasa de bit variable en tiempo real (rt-VBR) este tipo de ,Tasa de bit disponible (ABR) y Tasa de bit no especificado (UBR).

Servicio CBR

Esta categoría es utilizada para emular conmutación de circuitos, donde la tasa de arribo de celdas es constante. La razón de pérdida de la celda es especificada para $CLP=0$ y puede no ser especificada para $CLP=1$. Ejemplos de aplicaciones que utilizan CBR son el teléfono, video conferencia y televisión [Raj, 1996].

Servicio rt-VBR (Aplicaciones con restricción de tiempo)

Esta categoría permite a los usuarios enviar información a una tasa variable, la multicanalización estadística es utilizada permitiendo definir una probabilidad de pérdida pequeña diferente de cero. Los parámetros especificados para este servicio son: el retardo

máximo, y el pico de la variación de retardo de celda. Dentro de las aplicaciones que utilizan este servicio se encuentran el video comprimido y video conferencia.

Servicio nrt-VBR (Aplicaciones sin restricción de tiempo)

Esta categoría (al igual que la anterior) permite a los usuarios enviar información a una tasa variable. El único parámetro que se especifica es el retardo promedio de la celda. Un ejemplo de aplicación que utilice este servicio es el correo multimedia.

Servicio ABR

Esta categoría es diseñada para un tráfico tradicional de datos, como transferencia de archivos y correo electrónico. La especificación no garantiza el retardo de transferencia de celda, ni tampoco la razón de pérdida de celda [Bambarelli, 1996].

Dependiendo del estado de la red, se requiere a la fuente controlar su velocidad de generación de celdas. Al usuario se le permite declarar una velocidad mínima de celda, la cual es garantizada.

Tasa de bit no especificado

Esta categoría es definida para aquellas aplicaciones de datos que quieren utilizar los recursos remanentes del conmutador, y que no son sensitivos a pérdida de celdas y retardo. Este servicio no realiza control de admisión, ni tampoco monitorea el comportamiento del tráfico. Durante periodos de congestión las celdas son descartadas, pero no se espera que la fuente reduzca su velocidad de generación de celda. Aplicaciones de usuario que pueden ser transportadas son: transferencia de archivos, correo electrónico, etc.

Parámetros de Calidad de Servicio

Los parámetros de QoS que deben ser negociados entre los sistemas terminales y la red de acuerdo con el Foro ATM [Foro ATM, 1999a] son: (a) máximo retardo de transferencia de celda, (b) variación de retardo y (c) razón de pérdida de celda.

(a) El *retardo de transferencia de celda* es el tiempo que transcurre desde que las celdas son generadas hasta que arriban a su destino, e incluye tiempos de propagación en el medio, y permanencia en colas de espera. Las aplicaciones limitadas en tiempo tienen especificado el máximo retardo de transferencia [Raj, 1996].

(b) La *variación de retardo* es la medida de varianza del retardo. Una alta variación implica largos tiempos de espera en colas para tráfico sensitivo a retardo tales como voz y video. Existen varias formas de calcular la variación del retardo, una de ellas es el retardo pico a pico obteniéndose mediante la diferencia entre el retardo máximo y el mínimo.

(c) La *razón de pérdida de celda* (CLR por sus siglas en inglés) es el porcentaje de celdas que se pierden en la red debido a corrupción de información en la capa física y a congestión en la red. Este valor de desempeño de la red es medido en el tiempo que dura la conexión.

$$CLR = \frac{\text{Celdas perdidas}}{\text{Celdas Transmitidas}} \quad (20)$$

Cada celda ATM tiene un bit de prioridad de descarte de celda en el encabezado. Durante congestión la red primero descarta celdas que tienen este bit colocado en uno. El

CLR puede ser especificado en forma separada para celdas con prioridad de descarte igual a uno ó cero.

III.3.3 Atributos de Tráfico

Los usuarios, al establecer el contrato de tráfico deben especificar algunos de los siguientes atributos de tráfico, de acuerdo al tipo de servicio requerido.

- a) Tasa pico de celda. Es la tasa máxima a la cual un usuario puede transmitir.
- b) Tasa sostenida de celda. Es la tasa de celda promedio medida sobre un intervalo de tiempo largo.
- c) Tasa mínima de celda: Es la tasa mínima deseada por un usuario.
- d) Tamaño máximo de ráfaga. Es el número máximo de celdas continuas que son enviadas a velocidad pico de celda, sin violar la velocidad sostenida de celda.

III.4 Calidad de Servicio en redes IP

III.4.1 Introducción.

La globalización en las redes de comunicación esta empujando a la red Internet a convertirse en una red capaz de manejar todo tipo de tráfico, ofreciendo los servicios requeridos para cada usuario o aplicación. Para satisfacer estos requerimientos de QoS, se han propuesto distintos mecanismos de congestión como filtros tipo Cubeta de goteo [Kim *et al.*, 1999] y ventana deslizante [Parekh y Gallager, 1993], los cuales son algoritmos que proveen una noción de soporte de QoS en esquemas de congestión. No es sino hasta 1994, cuando la comunidad Internet a través del Grupo de Trabajo de Ingeniería en Internet (IETF por sus siglas es inglés) propuso una arquitectura completa de Servicios Integrados

(IntServ) [Clark *et al.*, 1994]. Posteriormente fue propuesta una estructura de Servicios Diferenciados (DiffServ) [Blake *et al.*, 1998] que proporciona escalabilidad y modelos de servicio flexibles, lo cual no ofrecen los servicios integrados.

III.4.2 Servicios Integrados: El primer paso

La arquitectura de servicios integrados “IntServ”, propone que una ráfaga de paquetes con dirección fuente, dirección destino, y número de puertos comunes, reciba un nivel deseado de servicio en términos cuantitativos de ancho de banda o retardo, para lo cual es necesario establecer y mantener un estado de flujo específico en la red. Con esto, se garantiza la QoS, y permite a la red Internet soportar flujos de tiempo real, así como ofrecer servicios de mejor esfuerzo. El principal problema que presenta esta arquitectura es que hoy en día solo un pequeño número de enrutadores son capaces de manejar señalización para realizar la reservación de recursos.

Las dificultades asociadas con la reservación de recursos por flujo de tráfico son las siguientes:

- Escalabilidad. La reservación de recursos por flujo implica la necesidad de un enrutador para procesar las reservaciones y mantener la QoS de cada flujo de tráfico que pasa a través del enrutador.
- Modelos de servicio flexibles. Los servicios integrados definen servicios específicos, o clases de servicio, lo cual los hace menos flexibles.

III.4.3 Servicios Diferenciados (*DiffServ*).

Servicios Diferenciados es una arquitectura que propone escalabilidad mediante estados de clasificación de tráfico agregado, implementándose en la capa de red, donde la categoría de un paquete es marcada en uno o varios campos del encabezado de la capa IP.

Los servicios diferenciados pueden ser implementados bajo dos modos: Envío expedito y envío asegurado. El modo de *envío expedito* fue diseñado para soportar conexiones con bajas pérdidas, retardos y jitter bajos. Este modo muestra un comportamiento parecido al servicio de líneas dedicadas virtuales entre dos puntos finales con un ancho de banda pico. El *envío asegurado* define cuatro clases relativas de servicio, donde cada una soporta tres niveles de prioridad de descarte. Cuando una congestión ocurre en el enrutador, los paquetes con precedencia de descarte alta son descartados primero, seguidos por los paquetes con una precedencia menor. Las cuatro clases de envío asegurado no tienen definidos los valores de ancho de banda o retardo.

La arquitectura propuesta a través del RFC 2475 [Blake *et al.*, 1998], dispone que los servicios diferenciados clasifican micro-flujos individuales a la entrada de la red en una de cuatro clases de servicios. El ingreso a la red se basa en el análisis de uno o mas campos en el paquete. El paquete es entonces marcado (cambiando algunos bits o campos del encabezado) para especificar que pertenece a una determinada clase de servicio, y luego se introduce en la red.

El nodo a la entrada de la red (figura 2) es definido como nodo *limitante*, los nodos interiores son llamados núcleos. Los enrutadores núcleo envían los paquetes examinando el encabezado del paquete para determinar como debe ser tratado basándose en un

comportamiento por salto, el comportamiento por salto es definido como el servicio que recibe el paquete a cada salto y como es enviado a través de la red. Los nodos interiores utilizan típicamente el manejo de colas y una disciplina de asignación para proveer el comportamiento por salto. El *Nodo limitante* es responsable de la clasificación del paquete, la medición, la marcación de este, así como mecanismos de formación y descarte. Los administradores de redes son responsables de la configuración del clasificador, el cual define los campos que serán examinados en cada paquete, así como otras acciones necesarias para enviar el paquete al siguiente nodo.

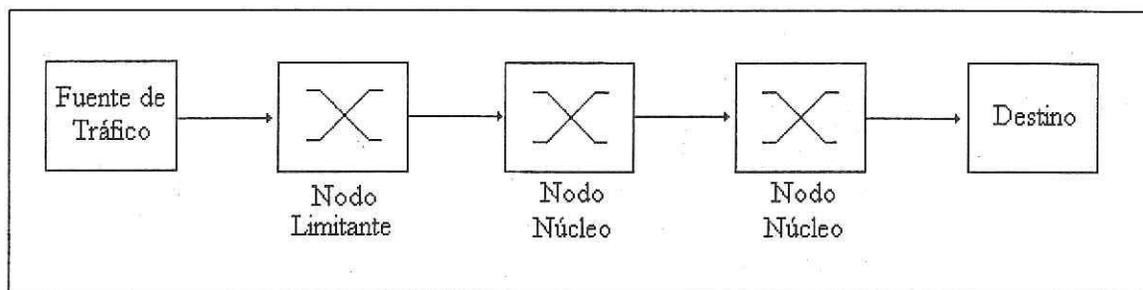


Figura 2. Esquema simplificado de una red de servicios diferenciados

El presente trabajo se encuentra enfocado a los nodos de entrada de la red, por lo que a continuación se describen a detalle las funciones que deben ser implementadas en este, de acuerdo con la arquitectura descrita en el RFC 2475.

Los nodos de entrada (o nodos limitantes), cumplen con las funciones de clasificación, medición, marcación, formación y descarte. La figura siguiente nos muestra un esquema funcional del nodo limitante.

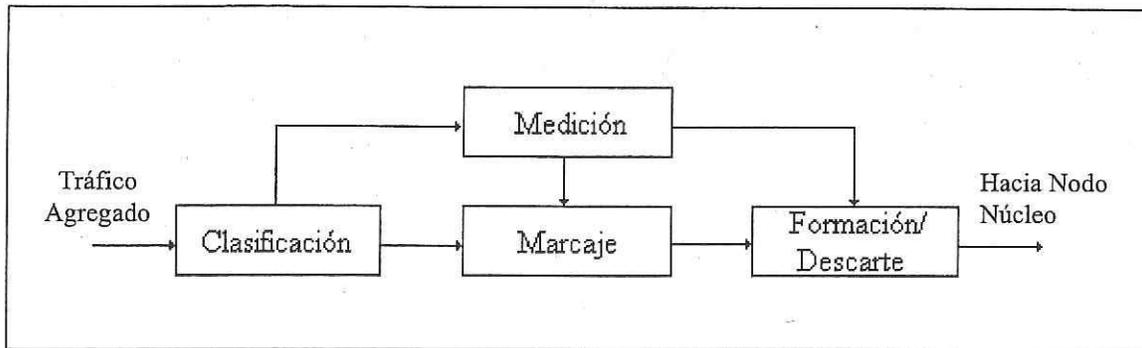


Figura 3. Esquema funcional del Nodo limitante.

Clasificación

La arquitectura de servicios diferenciados no propone reglas para tal clasificación, por lo cual esto puede ser implementado en el enrutador de entrada por el administrador de la red, mediante una tabla con las direcciones de fuentes de tráfico, que necesitan ser marcadas para dar una determinada preferencia.

En la arquitectura de servicios diferenciados el paquete es marcado en el campo llamado de servicios diferenciados (DS) el cual se encuentra en el encabezado de paquete de IPv4 o IPv6. Este campo se compone de un subcampo de seis bits llamado código de servicios diferenciados DSCP y otro sub-campo de dos bits que actualmente se encuentra sin uso.

Los clasificadores seleccionan los paquetes basados en los valores de uno o mas campos del encabezado IP como dirección fuente, dirección destino, DSCP, etc.

Medición.

Esta función compara el flujo de paquetes a la entrada del enrutador con el contrato de tráfico negociado y determina si el paquete se encuentra cumpliendo el contrato. Dentro

de los mecanismos para comparar conformidad con el contrato de tráfico se pueden utilizar mecanismos de vigilancia como: Estimador promedio de velocidad (basado en ventana) y Cubeta de Goteo [Karlsson, 1996], los cuales serán presentados en el siguiente capítulo de forma detallada. Para el caso particular del medidor, cuando el estimador promedio de velocidad concluye que el flujo de tráfico excede el contrato, este mecanismo marcará el paquete como fuera de contrato, con una probabilidad de descarte que se incrementa linealmente. El algoritmo de Cubeta de Goteo, muestra mejor desempeño ya que permite la transmisión de una ráfaga determinística.

Marcaje.

La arquitectura propuesta en [Blake *et al.*, 1998] no define un mecanismo de marcación, permitiendo cualquier algoritmo siempre y cuando cumpla con los requerimientos de la arquitectura propuesta. La comunidad Internet ha publicado trabajos en progreso para marcación de dos y tres prioridades de descarte, descritos en [Kim, 1999] y [Heinaneen, 1999] respectivamente.

En el algoritmo de *dos prioridades de descarte*, se tienen los siguientes parámetros de contrato: velocidad de bit y tamaño de ráfaga. Las unidades del primero son bytes/segundo de paquetes IP, que incluye el encabezado, el tamaño de la ráfaga es medida en bytes y debe ser configurada para ser mayor que el tamaño del paquete IP mas grande de la ráfaga.

El mecanismo sigue la siguiente regla [Kim, 1999]: si el número de bytes del paquete es mayor que el tamaño de la ráfaga, el paquete es marcado con baja prioridad de descarte.

El mecanismo que permite desarrollar *tres prioridades de descarte* (verde, amarillo, rojo), se basa en tres parámetros de tráfico: Velocidad de información comprometida (VIC), Tamaño de ráfaga comprometida (TRC), y Exceso de tamaño de ráfaga (ETR). El paquete es marcado como verde si no excede el VIC, amarillo si excede el VIC pero no el TRC, de otra manera es marcado como rojo.

Formación

Los formadores o espaciadores retardan algunos o todos los paquetes de una ráfaga de tráfico para forzar a este a cumplir con el contrato de tráfico. Un formador usualmente es un buffer de tamaño finito, en el cual los paquetes son descartados si no se encuentra espacio suficiente para mantenerlos o retrasarlos.

Descarte.

Este mecanismo descarta alguno o todos los paquetes de una ráfaga de tráfico que no cumpla con el contrato de tráfico, en esta función la IETF, deja abierta la opción de utilizar cualquier tipo de mecanismo de vigilancia. En el siguiente capítulo se describen en forma extensa los algoritmos de vigilancia tales como Ventana Fija, Cubeta de goteo y diferentes modificaciones al algoritmo original de Cubeta de goteo.

IV MECANISMOS PARA PROVEER GARANTIAS DE CALIDAD DE SERVICIO

IV.1 Introducción

En el capítulo anterior se describieron servicios importantes que permiten a las redes ATM e IP proveer QoS a aplicaciones diversas. En el presente capítulo examinaremos varios mecanismos que permiten ofrecer estas garantías, tales como mecanismos de asignación y de vigilancia.

La forma en la que los paquetes que se encuentran en la cola de espera son seleccionados para ser transmitidos en el enlace es conocido como disciplina de asignación de enlace.

La función de vigilancia debe estar disponible para cada conexión o flujo durante todo el tiempo de duración, y debe operar en tiempo real. Estos requerimientos implican que el mecanismo utilizado debe ser rápido, simple y con un costo que permita la implementación física.

Idealmente, los mecanismos de vigilancia deben de ser transparentes a las conexiones o flujo y respetar el contrato de tráfico. Otro requerimiento es que el tiempo de reacción dinámica de los mecanismos de vigilancia debe de ser corto para contrarrestar la avalancha de paquetes debido a las colas de espera relativamente pequeñas.

IV.2 Mecanismos de Asignación para proveer QoS

IV.2.1 Primer arribo primer servicio (FIFO por sus siglas en inglés)

La disciplina de asignación FIFO selecciona los paquetes del enlace de transmisión en el mismo orden en el que arriban a la cola (ver figura 4).

Los paquetes que arriban al enlace cuando el servidor está ocupado ofreciendo servicio a otro paquete, tendrán que ser almacenados en la cola para su posterior transmisión. En caso de no existir suficiente espacio en la cola para almacenar el paquete, este será descartado.

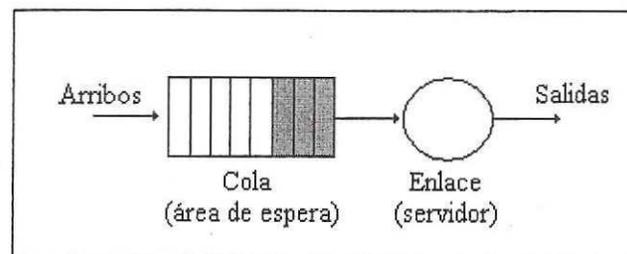


Figura 4. Esquema de asignación FIFO

IV.2.2 Asignación de colas por prioridad.

Bajo la disciplina de Asignación de colas por prioridad, los paquetes que arriban al enlace de salida son clasificados en una de dos o mas clases de prioridad como se muestra en la figura 5 [Demers *et al.*, 1990].

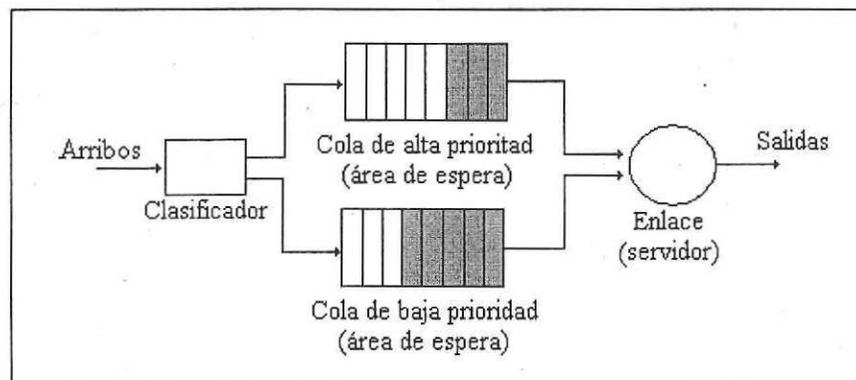


Figura 5. Esquema de asignación de colas por prioridad

La clasificación de prioridad en paquetes dependerá del marcaje explícito que tenga el encabezado del paquete. Por ejemplo en el caso de los paquetes IPv4 el valor en el campo "Tipo de Servicio".

Al transmitir, la disciplina de asignación tomará los paquetes de la clase de prioridad más alta hasta que esta cola se encuentre vacía, y continuará ofreciendo servicio al paquete que tenga la prioridad próxima mas baja.

IV.2.3 Asignación Round-Robin / Colas con servicio ponderado

Bajo la disciplina de asignación Round Robin los paquetes son clasificados y enviados hacia una de las colas de espera basados en su prioridad, sin embargo no se tiene una prioridad estricta de servicio entre las colas, este asignador alterna el servicio entre las clases, transmitiendo un paquete de clase 1, uno de la clase 2, hasta la última clase, repitiendo este patrón indefinidamente. La asignación de colas con servicio ponderado (figura 6) realiza la misma forma de clasificación que Round-Robin, sin embargo difiere en como recibe la cantidad *diferencial* de servicio en un intervalo de tiempo durante el cual hay paquetes de clase i para enviar. A la clase i se le garantizará la asignación de una fracción de ancho de banda igual a w_i/S_{wi} , donde w_i es la ponderación de la clase y S_{wi} es la suma de los pesos de todas las clases. En el peor caso, cuando todas las colas tengan paquetes, la clase i seguirá teniendo una fracción w_i/S_{wi} del ancho de banda [Demers *et al.*, 1990], [Kyas, 1995].

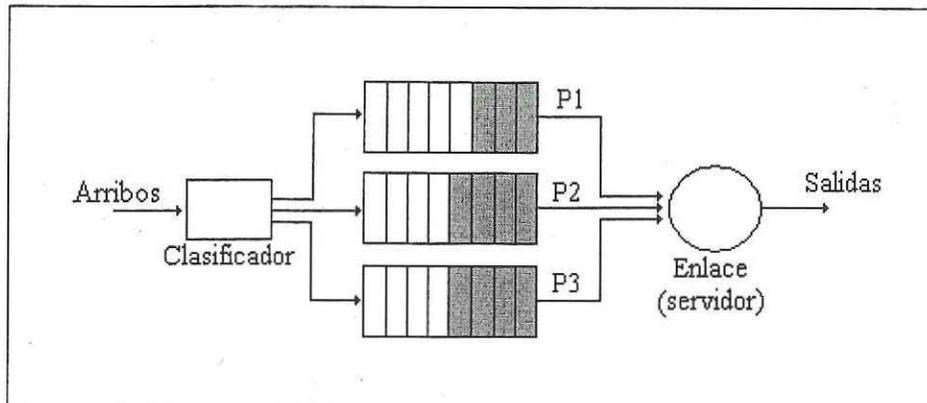


Figura 6. Esquema de asignación de cola con servicio ponderado.

IV.3 Mecanismos de vigilancia y control de congestión

IV.3.1 Mecanismo de salto de ventana

El mecanismo de salto de ventana limita el número máximo de paquetes que son aceptados de una fuente en un intervalo de tiempo (ventana) a un número máximo N [Parekh y Gallager, 1993]. Un intervalo nuevo inicia inmediatamente al final del intervalo precedente, donde el contador asociado inicializa su valor a cero. La implementación del presente mecanismo necesita la medición del intervalo T , contar el número de arribos, y requiere dos variables, una para el manejo del contador y la otra para la medición del intervalo T .

IV.3.2 Mecanismo de salto de ventana por disparo.

En este mecanismo el tiempo de ventana no es sincronizado con la actividad de la fuente, como lo realiza el mecanismo de salto de ventana, aquí las ventanas de tiempo no son consecutivas.

IV.3.3 Mecanismo de movimiento promedio exponencialmente ponderado

Este mecanismo utiliza ventanas consecutivas con tiempo fijo, tal como lo hace el mecanismo de salto de ventana, la diferencia es que el máximo número de paquetes aceptados en la i -ésima ventana es N_i y es función del número promedio de paquetes (N) permitido por intervalo y la suma ponderada exponencialmente del número de paquetes aceptados en los intervalos precedentes (X_{i-1}) de acuerdo con la siguiente ecuación [Parekh y Gallager, 1993]:

$$N_i = \frac{N - \gamma S_{i-1}}{1 - \gamma} \quad 0 \leq \gamma < 1 \quad (21)$$

$$\text{donde } S_{i-1} = (1 - \gamma)X_{i-1} + \gamma S_{i-2} \quad (22)$$

El factor γ controla la flexibilidad del algoritmo con respecto al comportamiento a ráfagas del tráfico. Si $\gamma = 0$ y N_i es constante, el algoritmo es idéntico al mecanismo de salto de ventana.

IV.3.4 Mecanismo de ventana deslizante

Similar al mecanismo de salto de ventana, este algoritmo limita el número de arribos de paquetes en un intervalo T . El tiempo de arribo de cada paquete es almacenado, y el contador es incrementado en uno por cada arribo. Exactamente T unidades de tiempo después del arribo de un paquete aceptado, el contador es decrementado en uno. Este mecanismo puede ser interpretado como una ventana, la cual se mueve a través del eje del tiempo, requiere que los tiempos de arribo hasta N paquetes sea almacenado en la duración

de una ventana, por lo que la complejidad del algoritmo es considerablemente mas alta que otros mecanismos de ventana.

IV.3.5 Mecanismo de Vigilancia: Cubeta con goteo.

En el algoritmo original de Cubeta con goteo que se muestra en la figura 6, los paquetes que llegan al nodo son aceptados solo si pueden tomar una estafeta del depósito de estafetas. Si el depósito de estafetas se encuentra vacío, entonces el paquete es descartado. Las estafetas son generadas a una velocidad de transmisión promedio de la red R , y almacenadas en el depósito de estafetas. El depósito tiene un tamaño finito B , después de que se completa el depósito de estafetas, las estafetas adicionales son descartadas. El tamaño del depósito puede ser visto como la máxima longitud de ráfaga permisible que puede ser transmitida a la vez. El depósito de estafetas puede ser implementado utilizando un contador que se incrementa cuando las estafetas son generadas y se decrementa cuando las estafetas son utilizadas [Kim *et al.*, 1999].

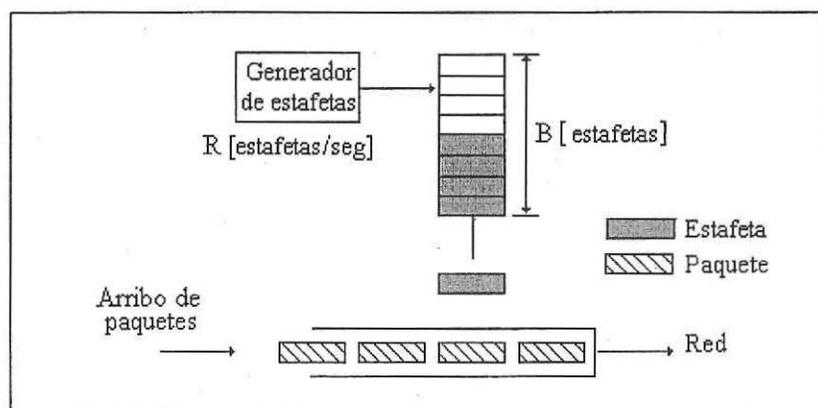


Figura 7. Mecanismo de Cubeta con Goteo

El mecanismo Cubeta con goteo convencional puede combinar el monitoreo con el control de flujo, ya que en este sistema los paquetes son descartados cuando el depósito de estafetas esta vacío.

La probabilidad de bloqueo de paquetes es la probabilidad de que una celda arribe cuando el depósito de estafetas se encuentre vacío, esto depende del tamaño del buffer de paquetes y el tamaño del depósito de estafetas.

Es deseable que la red pueda manejar ráfagas de tráfico grandes y al mismo tiempo que el retardo de los paquetes pueda ser reducido, ya que el costo de implementación de un gran depósito de estafetas es menor que el costo de un buffer de paquetes grande.

IV.3.6 Mecanismo de Cubeta con Goteo Virtual Modificado

Este mecanismo presenta una modificación del algoritmo original de cubeta con goteo, el cual es un algoritmo de vigilancia. En la Cubeta con goteo virtual, se permite hacer uso de los recursos del nodo si estos se encuentran disponibles aunque el paquete no cumpla con el contrato establecido, sin embargo, se prevé no afectar a otros paquetes (conexiones o flujo) que si cumplen con el contrato.

Existe una técnica muy utilizada, en la cual se diferencian los paquetes de acuerdo a su cumplimiento con el contrato de tráfico, basado en determinados parámetros del mismo. La diferenciación puede darse en dos o tres colores, para el caso de dos colores los paquetes que violan el ancho de banda disponible para su conexión o flujo son marcados como rojo, mientras que la velocidad de bit por algún tiempo puede ser excedida hasta que

el depósito de estafetas se vacíe. Cuando ambos depósitos de estafeta se encuentren vacíos el paquete se descartará inmediatamente.

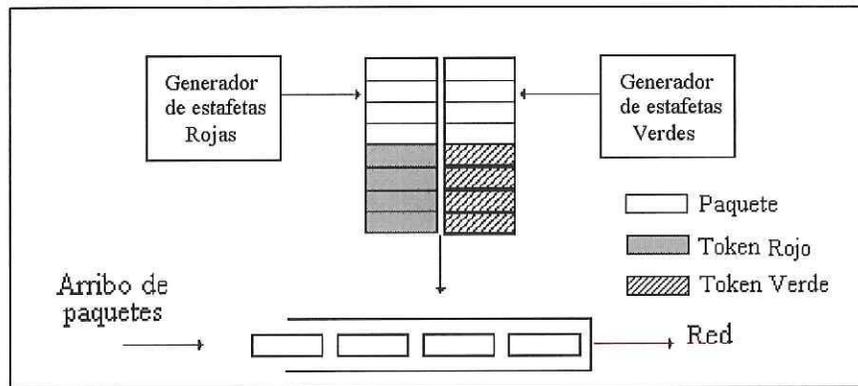


Figura 8. Mecanismo de cubeta con goteo virtual

Cuando los paquetes marcados como rojo arriben a un enlace congestionado serán los primeros en ser descartados, entonces el caudal eficaz de los paquetes marcados como verdes no es afectado significativamente.

Se puede pensar que marcar degrada a los paquetes que cumplen con el contrato, sin embargo los esquemas de manejo de buffer que existen, permiten un óptimo desempeño de los paquetes no marcados [Karlsson, 1996].

IV.3.7 Mecanismo Cubeta con goteo Adaptivo

Es una mejora de la técnica tradicional Cubeta con goteo, en el cual la generación de estafetas es función de la carga de la red, o de la congestión. Este algoritmo se basa en

observar que tan ocupado se encuentra el buffer. El buffer es monitoreado regularmente, y se toma una acción apropiada después de cada observación.

Se pueden utilizar dos niveles de decisión (el alto y el bajo), cuando el nivel de ocupación de buffer se encuentra dentro de estos dos niveles, la velocidad de generación de estafetas permanece sin cambio. Cuando el nivel de ocupación cae por debajo del nivel de decisión, la velocidad de generación de estafeta debe ser incrementada. En el caso opuesto, la velocidad se decrementa. Los niveles pueden ser seleccionados si ciertas medidas de desempeño son optimizadas.

La congestión en la red se detecta cuando la velocidad excede al nivel de decisión superior. En este caso, el mecanismo cubeta con goteo limita la transmisión de tráfico decrementando la velocidad de generación de estafeta, en el caso opuesto la generación de estafeta es incrementada [Ibañez y Nichols, 1998].

El algoritmo que describe el comportamiento de Cubeta con goteo Adaptivo es el siguiente:

Entrada:

occupancy	/*El valor de ocupación del buffer */
y_k	/*La velocidad de generación de estafeta en la ranura k*/
a_1, a_2	/*Constantes donde $a_1 > 1$ y $a_2 < 1$ */
thr_low, thr_up	/*Umbrales que indican el valor bajo (thr_low) y alto (thr_up) de ocupación del buffer */

Salida:

y_{k+1}	/*La velocidad de generación de estafeta en la ranura k+1 */
-----------	--

Algoritmo:

si (occupancy < thr_low)	$y_{k+1} = a_1 y_k$
de otra forma si (occupancy > thr_up)	$y_{k+1} = a_2 y_k$
de otra forma si (thr_low < occupancy < thr_up)	$y_{k+1} = y_k$

IV.3.8 Mecanismo de Cubeta con goteo adaptivo ayudado por predicción

La función principal de Cubeta con goteo es complementada mediante la predicción de tráfico (véase capítulo 2) para realizar un mejor manejo de recursos dentro de un nodo de acuerdo al contrato de QoS, bajo el principio de imparcialidad [Angulo *et al.*, 1999], este principio describe que la red debe proporcionar el servicio contratado sin inclinarse por los servicios con restricción de tiempo, logrando una convivencia pacífica entre todos los tipos de tráfico. Este algoritmo se puede aplicar sobre fuentes de tráfico tales como video, WWW y Ethernet. Fuentes que pueden ser modeladas por procesos auto-similares con distribución alfa-estable, este algoritmo permite realizar una decisión basada en el estado futuro del nodo (conmutador o enrutador).

En el presente mecanismo el incremento de la velocidad de generación de estafeta se realiza en cada canal virtual (ATM) o flujo (Servicios diferenciados) de acuerdo al estado futuro de las colas (en base a la predicción de tráfico) y al estado actual del buffer de salida.

La decisión de considerar los paquetes que arriban como excesivos o no, obliga a realizar una evaluación lógica de la utilización de umbrales fijos en los algoritmos de control, que nos provee una acción dinámica ineficiente, tomando en cuenta que la congestión no es una función estática.

Bajo la restricción de recursos fijos por parte de la red, si el tráfico se encuentra violando el contrato, y los otros tráficos (VC o flujos) están utilizando sus recursos

correspondientes, los paquetes serán descartados para no perjudicar los otros flujos que si se encuentran cumpliendo su contrato, todo esto bajo el principio de imparcialidad.

El escenario que se asume para un esquema de QoS basado en predicción de tráfico es el que se muestra en la figura 9 [Angulo *et al.*, 1999], donde se asume que el nodo tiene un buffer de entrada por cada tipo de tráfico y todos los tráficos comparten un mismo buffer de salida. En esta propuesta, los buffers en forma individual son servidos por la disciplina FIFO, al ser observados en conjunto, son servidos por asignación de colas ponderadas.

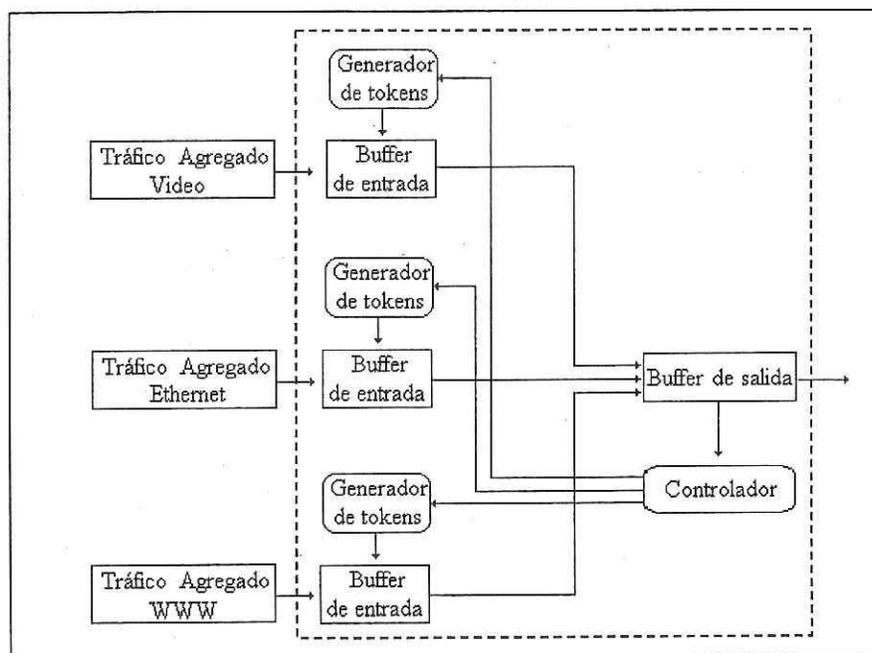


Figura 9. Escenario para esquema de QoS cubeta con goteo ayudado por predicción

La inteligencia del mecanismo de Cubeta con goteo adaptivo ayudado por predicción es implementado en el controlador, el cual calcula la tasa de generación de

estafetas para cada generador. La velocidad de generación de estafetas es incrementada o decrementada, dependiendo del futuro estado del nodo, el cual es estimado basado en el estado actual y en la predicción del tráfico [Angulo *et al.*, 1999].

En el algoritmo de Cubeta con goteo adaptivo, la decisión de modificar la tasa de generación de estafetas se toma después de recibir muchos paquetes o muy pocos, es decir hasta que uno de los umbrales sea cruzado, lo cual hace que las acciones se realicen muy tarde para resolver la situación. Una posible solución sería tener los umbrales muy cerca uno del otro, sin embargo esto situaría al sistema en el otro extremo acusando sobre-acción del mismo.

El mecanismo de Cubeta con goteo adaptivo ayudado por predicción tiene que asegurarse de calcular las velocidades de paquete a conexiones individuales y tiene que limitar la velocidad de paquete de todas las conexiones para no exceder la capacidad de ancho de banda del enlace de salida. El algoritmo es descrito a continuación:

Entrada:

occupancy	/* El valor de ocupación del buffer */
i-th_qoc	/* La predicción del estado de la i-ésima cola en la ranura k */
i-th_yk	/* La velocidad de generación del i-ésimo estafeta en la ranura k*/
i-th_thr_low, i-th_thr_up	/* Umbrales que indican el valor bajo (thr_low) y alto (thr_up) de ocupación de la i-ésima cola */
a ₁ , a ₂	/* Constantes donde a ₁ >1 y a ₂ <1 */
thr_low, thr_up	/* Umbrales que indican el bajo (thr_low) y el alto (thr_up) valor de ocupación del buffer */

Salida:

i-th y _{k+1}	/* La tasa de generación de estafetas en la ranura k+1 */
-----------------------	---

Algoritmo:

```

si (occupancy < thr_low)
{
    si ( i-th_qoc > i-th_thr_up)            $y_{k+1} = a_1 \cdot a_1 \cdot y_k ;$ 
    de otro modo:  $y_{k+1} = a_1 y_k ;$ 
}

de otro modo si(occupancy > thr_up)
{
    si (i-th_qoc < i-th_thr_low)            $y_{k+1} = a_2 \cdot a_2 \cdot y_k ;$ 
    de otro modo  $y_{k+1} = a_2 y_k ;$ 
}

de otro modo si(thr_low < occupancy < thr_up)
{
    si ( i-th_qoc > i-th_thr_up)            $y_{k+1} = a_1 \cdot y_k ;$ 
    de otro modo si (i-th_qoc < i-th_thr_low)  $y_{k+1} = a_2 \cdot y_k ;$ 
    de otro modo  $y_{k+1} = y_k ;$ 
}

```

V DISEÑO DE LOS MODELOS DE LOS NODOS ATM E IP UTILIZANDO ASIGNACIÓN DINÁMICA DE ANCHO DE BANDA.

V.1 Tráfico limitado en tiempo

El tráfico limitado en tiempo, presenta algunos requerimientos límite en los siguientes parámetros de desempeño de la red: Pérdida de paquetes, Retardo extremo a extremo y Variación de retardo.

Pérdida de paquetes

Existen dos posibles situaciones en las cuales un paquete se pierde: que exista corrupción en algunos de sus bits, o que uno de los buffers en la ruta del paquete se encuentre lleno y no admita el paquete.

Para el caso de ATM el factor de corrupción de paquete no es un factor de peso, dado que la capa física sobre la cual trabajan la mayoría de las redes ATM está basada en líneas de transmisión confiables. Sin embargo, respecto al descarte de celdas por colas llenas debe existir un mecanismo de manejo de congestión para evitar el descarte y garantizar el caudal eficaz contratado por la aplicación. Para eliminar el efecto de pérdida de paquetes en redes IP se pueden enviar los paquetes sobre una capa de transporte TCP en lugar de UDP, sin embargo los mecanismos de retransmisión generalmente no son aceptados por aplicaciones limitadas en tiempo debido a que se incrementa el retardo punto final a punto final [Bolot y Vega-Garcia, 1996].

Retardo extremo a extremo

Es la acumulación de los retardos de procesamiento en las colas de los enrutadores (nodos IP) o conmutadores (nodos ATM), mas el retardo de propagación en las líneas de transmisión y retardos de procesamiento del sistema final.

Los tres componentes del retardo en la información son [Foro ATM, 1999a]: retardo de codificación, retardo de paquetización y retardo en las colas de espera.

Retardo de codificación es el retardo que implica la codificación de una señal analógica en una señal digital. Puede tener distintos componentes como el retardo de procesamiento, el retardo del tamaño de la trama y el retardo de correlación entre una trama y la siguiente. El *Retardo de paquetización* es el retardo necesario para llenar un paquete con tramas antes de ser transmitido. Este retardo es claramente definido por la longitud mínima que el paquete puede enviar y la velocidad de bit de salida del codificador. En el caso de ATM donde las celdas tienen tamaño fijo, si la velocidad de salida del codificador es muy pequeña se provoca un retardo de paquetización muy grande donde una posible solución es enviar las celdas parcialmente llenas. *Retardo debido a colas de espera* es el tiempo que transcurre desde que un paquete arriba a la cola de espera, hasta que es transmitido.

Variación en retardo

Uno de los componentes del retardo de extremo a extremo es el retardo aleatorio de las colas de espera en los nodos de interconexión, debido a esa aleatoriedad el tiempo desde que el paquete es generado hasta que es recibido en la fuente puede fluctuar de paquete a paquete, este fenómeno es llamado jitter. En redes IP, la principal afectación de una

cantidad de jitter mayor al tiempo en el que se crea un encapsulado, es que un paquete arribe al destino antes que el paquete que fue creado en un tiempo anterior en la fuente, esto debido a que los paquetes toman diversas rutas y por lo tanto los tiempos de recorrido varían. Si el receptor ignora la presencia del jitter y reproduce las muestras tan pronto como estas arriban, el resultado es una calidad no deseada por el receptor.

Remover el jitter en el receptor.

En aplicaciones de Internet tales como teléfono, audio en demanda, video interactivo o video en demanda, el receptor puede proveer sincronización en la reproducción de los encapsulados de voz, en presencia de jitter aleatorio de la red. Esto es típicamente realizado combinando los siguientes mecanismos: Secuencia de números, estampas de tiempo y retardo de reproducción. La secuencia de números y las estampas de tiempo ocupan campos en el encabezado del encapsulado de la muestra [Ramjee *et al.*, 1997].

V.1.1 Tráfico de voz

Restricciones del tráfico de voz

Para aplicaciones de audio altamente interactivas, retardos menores que 150 mseg no son percibidos por el oído humano, entre 150 y 400 mseg puede ser aceptable pero no ideal, retardos mayores 400 mseg resulta en conversaciones de voz ininteligible [Foro ATM, 1999a], [Kostas, 1998] y [Bolot y Vega-García, 1996].

El porcentaje máximo de pérdida de paquetes admitida por aplicaciones de voz es de 1%. Como se comentó anteriormente el jitter máximo permitido es menor que el periodo de encapsulamiento de muestras de voz.

Modelado de Tráfico de voz en redes ATM

Para transmisión de voz sobre ATM, se simularon las ráfagas de voz mediante modelo markoviano descrito en el capítulo 2, donde el ancho de banda promedio que utiliza cada usuario es 57.5% de la velocidad de bit proporcionada por el codificador, dado que la mayoría de los mecanismos de transmisión de voz presentan reconocimiento de periodos de silencio. En la actualidad un codificador de 64kbps de salida (PCM) se considera que desperdicia el ancho de banda, por lo que en su lugar se utiliza otro codificador que permita utilizar menor ancho de banda: *modulación por codificación de pulsos delta adaptivo (ADPCM)* estandarizado por la Unión Internacional de Telecomunicaciones ITU como G.726 [ITU, 1996b]. Con el cual se obtiene una reducción en la velocidad de datos de 64kbps a 32kbps.

No se adoptó en este trabajo una tasa menor a 32kbps, ya que al ser las celdas de tamaño fijo, los 48 bytes de carga útil se envían de todas formas y al utilizar 8kbps el retardo de paquetización se incrementaría a 47msegundos, utilizándose 4 veces menos ancho de banda, pero incrementando el retardo en 6×10^5 veces (comparación realizada con TP11368 National Semiconductor G.726 y G.729A SX800LCS Lucent). El utilizar ADPCM en lugar de PCM, nos brinda un retardo de codificación del 0.010% del retardo total permitido para voz (150mSec) y reduce a la mitad el ancho de banda utilizado.

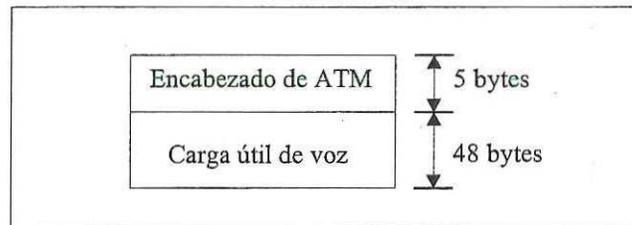


Figura 10. Formato de paquete de voz en ATM

Tráfico de voz en redes IP

En el caso de IP no se tiene un tamaño fijo de carga útil, por lo que esta puede ser tan pequeña como sea requerido, el inconveniente que presenta el enviar paquetes con poca carga útil es que el caudal eficaz se reduce considerablemente al utilizar una carga muy pequeña, dado que el encabezado formaría gran parte del tamaño del paquete y no representa información.

Al utilizar un codificador del tipo G.729A propuesto por ITU [ITU, 1996b], el formato del paquete utilizado para efectos de simulación se presenta en la figura 11, como se comentó anteriormente al enviarse la voz por paquetes IPv6 se utiliza la capa de transporte del tipo datagrama UDP que representa 4 bytes de encabezado:

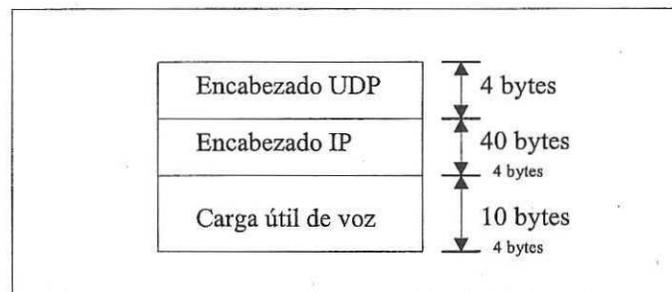


Figura 11. Formato de paquete de voz en redes IP

Por que no enviar mas de una trama de voz en el mismo paquete de IP?

Dado que el retardo que aporta el codificador G.729A es del orden de 25 mseg, es inadecuado añadir a ello un retardo por paquetización. Se utiliza este codificador (G.729A) porque representa una reducción del ancho de banda con respecto a PCM del orden de 8 veces (en la carga útil).

V.1.2 Tráfico de video

El video es producido por el despliegue de tramas a una velocidad fija conocida como velocidad de reproducción, esta velocidad varia de un formato de video a otro. En algunos formatos de video estandarizados se incluye NTSC (30 tramas/seg) y PAL (25 tramas/seg).

Requerimientos de Caudal Eficaz.

El caudal eficaz mínimo es dado por la velocidad de bit promedio, la cual es conocida para muchas técnicas de compresión (1.5 Mbps para MPEG-1 y 5Mbps para MPEG-2). Sin embargo debido a las fluctuaciones de velocidad de bit, los requerimientos de caudal eficaz son típicamente mas altos que la velocidad promedio [Krunz, 1999].

Requerimientos de retardo.

Las comunicaciones interactivas de video tienen fuertes requerimientos de retardo en términos de retardo máximo y de variación del retardo (jitter). Los valores aceptables de retardo máximo se encuentran en el rango de 150-400 mseg. Las limitaciones del jitter dependen de los requerimientos de los suavizadores de imagen en la fuente y el receptor.

Asumiendo una velocidad constante de reproducción de video de f tramas/segundo y suponiendo no suavizadores, el receptor requiere $1/f$ segundos de tiempo de almacenamiento para mitigar el efecto del jitter.

En suma, la sincronización entre video y audio requiere que el tiempo de interarribo de los paquetes de video y los paquetes de audio sea menor que 80 msec para un sentido de la comunicación.

Requerimientos de pérdidas.

En general los requerimientos de pérdidas de paquete son del orden de 10^{-2} y 10^{-6} .

Modelado de Tráfico de video en redes ATM

Para la generación de trazos de video, se utilizo el modelado de un sistema auto-similar con distribución alfa estable de Gallardo [Gallardo *et al.*, 1999] para tramas de video MPEG-1, donde los parámetros alfa y H son descritos en la tabla 1. La generación de las muestras artificiales de tráfico se realiza con el algoritmo propuesto en [Gallardo *et al.*, 1998a], para reducir tiempos de simulación.

Tráfico de video en redes IP

La utilización de redes de servicios diferenciados aplicado al manejo de tráfico agregado, nos permite situarnos en el núcleo de redes de área amplia, por lo que para simplificar la simulación se utilizará un formato de paquete con tamaño fijo de 576 bytes. Para reproducir el comportamiento alfa-estable se varía el tiempo de interarribo de los paquetes.

V.2 Tráfico no limitado en Tiempo

Dado el trato preferencial del tráfico limitado en tiempo, es muy importante observar que la red no afecte la calidad del servicio al tráfico no limitado en tiempo, garantizando el caudal eficaz requerido por estos. Un dato interesante, como se expone en Apéndice A, es que entre los diferentes tipos de tráfico el comportamiento a ráfagas difiere y por lo tanto sus difieren sus demandas a la red.

V.2.1 Tráfico Ethernet y WWW

Las fuentes de Ethernet LAN/WAN modeladas por Gallardo [Gallardo *et al.*, 1998b] basado en las muestras obtenidas por investigadores de Bellcore, fueron simuladas bajo los parámetros descritos en la siguiente tabla:

Tabla I. Parámetros de modelado para tráfico alfa-estable [Gallardo *et al.*, 1998b].

TRAFICO	Parámetro H	Parámetro Alfa
Video	0.903	1.90
WWW	0.833	1.28
Ethernet	0.834	1.95

V.3 Parámetros de fuentes de tráfico.

Tabla II. Parámetros de tráfico.

Parámetros de Fuentes de Tráfico	Tipo de Tráfico			
	Video	Ethernet	WWW	Voz
Velocidad sostenida (bps)	6,870,528	51,542,400	4,608	184,000
Velocidad pico (bps)	10,199,380	173,903,615	393,265	320,000

Tabla III. Parámetros de contrato de tráfico.

Parámetros de Contrato de Tráfico	Tipo de Tráfico			
	Video	Ethernet	WWW	Voz
Velocidad sostenida (bps)	6,500,000	51,500,000	4,500	184,000
Velocidad pico (bps)	8,000,000	80,000,000	10,000	320,000

V.4 Diseño del conmutador ATM

El conmutador ATM se compone de los siguientes elementos:

Un buffer de entrada para cada tipo de tráfico, un buffer de salida, un multicanalizador y un controlador. Todos los buffers son servidos por disciplina FIFO, para los buffers de entrada la velocidad de servicio es dada por sus correspondientes generadores de estafetas. La velocidad de servicio del buffer de salida es la suma de las velocidades sostenidas de cada canal virtual entrando al conmutador.

Se asume que los enlaces particulares de tráfico son cuatro canales virtuales independientes, donde las celdas que arriban al conmutador son almacenadas en el buffer apropiado, asociado con su específico VC/VP

Las fuentes de tráfico de video, Ethernet LAN y WWW son implementadas en base al modelo de tráfico agregado alfa-estable. El tráfico agregado de voz es obtenido mediante la suma de fuentes individuales.

La velocidad inicial de generación de estafeta, es calculada en base a la velocidad sostenida de celda, dado que cada estafeta corresponde a una celda. Los valores alto y bajo de generación de estafeta son ajustados por los requerimientos de tráfico [Angulo *et al.*, 1999]. El valor mas alto de generación de estafeta es calculado de acuerdo con la máxima velocidad de celda, el valor mas bajo de generación de estafeta es ajustado de acuerdo con la mínima velocidad de celda y el máximo retardo de espera en cola que puede permitir el paquete.

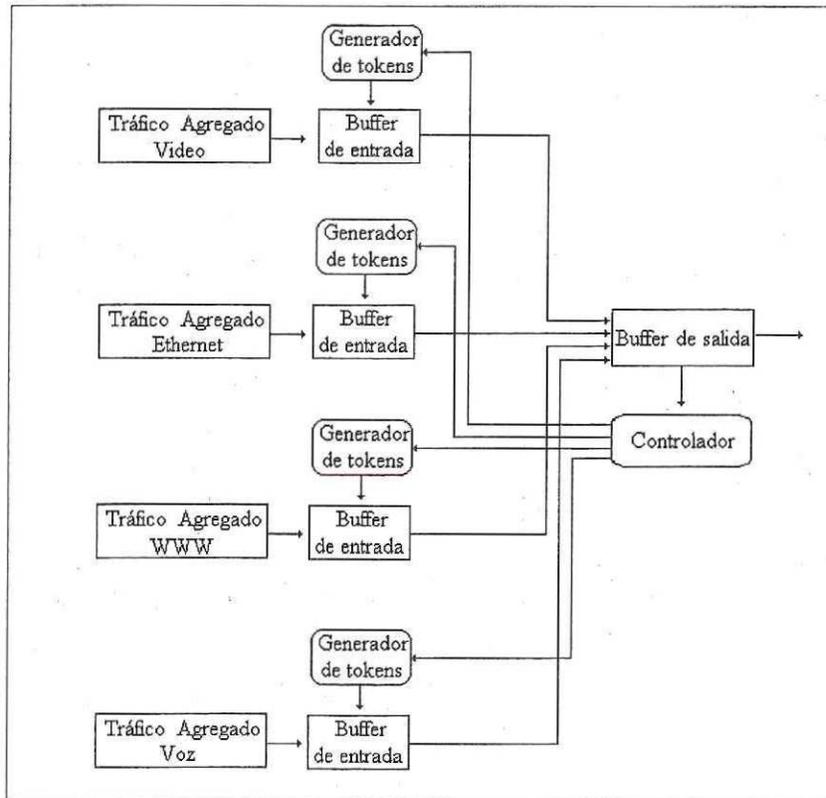


Figura 12. Implementación del modelo del conmutador ATM

El máximo retardo del paquete sufrido en el conmutador es calculado por el tamaño del buffer de la red y su respectiva velocidad de servicio, así como el valor límite superior aplicado al buffer de salida.

Tabla IV. Tiempo de interarribo en generación de estafeta

Parámetros de generación de estafeta	Tipo de Tráfico			
	Video	Ethernet	WWW	Voz
Valor inicial (seg)	5.58×10^{-5}	7.9×10^{-6}	2×10^{-1}	0.0026
Límite superior (seg)	6×10^{-5}	1.2×10^{-5}	0.4	0.0052
Límite inferior (seg)	4×10^{-5}	3×10^{-6}	1.3×10^{-2}	0.0013

Tabla V. Distribución de Buffer

Parámetros de generación de estafetas	Tipo de Tráfico				
	Video	Voz	Ethernet	WWW	Output Buffer
Tamaño de buffer distribuido(celdas)	40	10	340	10	210
Tamaño del depósito de estafetas (estafetas)	30	10	115	5	***

V.5 Diseño del enrutador IP

En las redes IP que proveen servicios diferenciados la clasificación de los paquetes se basa en el encabezado de los paquetes en la capa de red de acuerdo con RFC 2475 [Blake *et al.*, 1998]. Existen dos clasificadores: El de comportamiento agregado (basado en el código DS) y el de multi-campo. En nuestra propuesta se utiliza el modo de comportamiento agregado, dado que uno de los fundamentos de la arquitectura de servicios diferenciados es la simplicidad del modelo. En [Bernet *et al.*, 1999] Bernet *et al.*, proponen que la función de medición de tráfico puede utilizar mecanismos de ventana o de cubeta de goteo. Sin embargo el modelo de asignamiento dinámico de ancho de banda [Angulo *et al.*, 2000] utiliza sólo mecanismos de cubeta de goteo, dado que estos proporcionan la velocidad de tráfico instantánea y la velocidad de tráfico a largo plazo, a diferencia de los mecanismos de ventana, que no brindan la medición instantánea del tráfico, por lo que las decisiones del enrutador se realizan demasiado tarde.

Para la función de marcación, la comunidad Internet ha desarrollado algunas propuestas a través de trabajos en progreso (Draft) como Marcador de tres colores y una sola velocidad [Heinanen, 1999], o Marcador de dos colores [Kim, 1999]. Para efectos del

modelo general y para no perder el enfoque de este trabajo que es el aplicar los diferentes algoritmos de cubeta de goteo para lograr un mejoramiento en el manejo de recursos.

Para el marcador, se utiliza el mecanismo de cubeta de goteo en cascada [Heinanen, 1999], marcando de acuerdo con los parámetros de contrato de tráfico tales como velocidad sostenida y exceso de ráfaga. El algoritmo se presenta en la figura 13, notese que el depósito de estafetas 1 tiene un tamaño unitario, mientras el tamaño del depósito 2 es el del exceso de ráfaga.

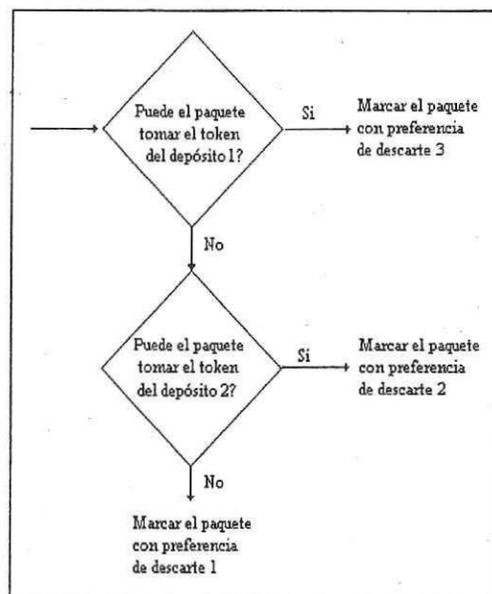


Figura 13. Diagrama de flujo del mecanismo de marcación

La función de Formación es desempeñada por los buffers, donde Bernet *et al.*, en [Bernet *et al.*, 1999] recomienda el uso de diferentes colas para cada tipo de tráfico y para cada preferencia de descarte, para aislar las ráfagas de tráfico. La función de descarte no restringe específicamente ningún mecanismo para desempeñarla; en esta función es utilizado el esquema de asignación dinámica de ancho de banda basada en predicción (algoritmo descrito en el capítulo 4).

Requerimientos de Tráfico limitado en tiempo.

La emigración de las redes Internet actuales a redes que proveen servicios diferenciados propone un incremento en el retardo de procesamiento en los nodos de la red. Sin embargo no todos los enrutadores provocan el mismo retardo. Se estima que los nodos núcleo realicen casi las mismas funciones que los enrutadores actuales, donde los nodos limitantes efectuaran las funciones complejas de procesamiento del paquete, lo cual no agrega mucho retardo dado que el tráfico agregado requiere pasar a través del nodo limitante sólo una vez.

El retardo del tráfico limitado en tiempo se debe a la paquetización, codificación, y retardo por almacenamiento en las colas [Foro ATM, 1999a]. Para estas tres situaciones, la única que el manejo de red puede optimizar es la última tratando de que el retardo sea menor que el máximo permitido para el tipo de aplicación.

Un parámetro que cobra importancia en las redes no orientadas a conexión es el Jitter. Como el caso de las aplicaciones de voz en demanda, es bien conocido que nuestro sistema puede soportar largos tiempos para reproducir el mensaje original, pero no puede esperar entre una ejecución de paquete y otra [Raj, 1996].

Como se mencionó anteriormente utilizamos diferentes buffers de entrada para cada preferencia de descarte en cada clase (sumando 12 buffers en total), la disciplina de asignamiento implementada es *Asignamiento ponderado de colas*.

El tamaño de los buffers es el mismo (en bits) que en el conmutador ATM, donde la velocidad de generación de estafetas (tabla VI) es calculada de la misma forma que para el

conmutador ATM, en este caso generando un estafeta por cada paquete en lugar de un estafeta por cada celda.

Tabla VI. Tiempo de interarribo de generación de estafeta

Parámetros de generación de estafetas	Tipo de Tráfico			
	Video	Ethernet	WWW	Voz
Valor inicial (seg)	60.65×10^{-5}	85.87×10^{-6}	2.174	0.02826
Límite superior (seg)	65.22×10^{-5}	13.04×10^{-5}	4.348	0.05652
Límite inferior (seg)	43.48×10^{-5}	32.61×10^{-6}	14.13×10^{-2}	0.01413

Tabla VII. Distribución de Buffer

Parámetros de generación de estafetas	Tipo de Tráfico				
	Video	Voz	Ethernet	WWW	Output Buffer
Tamaño de buffer distribuido (bits)	5760	500	19584	11520	12096
Tamaño de la cubeta de estafetas (estafetas)	30	10	115	5	***

VI. IMPLEMENTACION Y RESULTADOS DEL MODELO DE OPnet

VI.1 Implementación del conmutador ATM

Dos de las principales características de ATM son el uso de multicanalización estadística de las ráfagas de tráfico y la habilidad para proveer garantías de calidad de servicio. Estas dos características tienen el potencial de optimizar los recursos de la red, pero al mismo tiempo demandan un control sofisticado y mecanismos de vigilancia para el tráfico.

En la presente tesis se evalúa el desempeño de dos de los tipos de servicio que se presentan en redes ATM con contrato de tasa de bit variable para aplicaciones limitadas en tiempo y aplicaciones no limitadas en tiempo (VBR-rt y VBR-nrt) así como la convivencia entre ellas, buscando respetar el principio de imparcialidad y cumplimiento de contrato de tráfico de la red y de los usuarios.

Evaluación de desempeño del conmutador ATM

Es muy importante hacer notar que estamos forzando al sistema a trabajar bajo condición de congestión. El conmutador implementado (figura 14) es compuesto por un buffer de entrada para cada tipo de tráfico, donde el tráfico comparte el buffer de salida. Todos los buffers son servidos por disciplina FIFO. Para los buffers de entrada, la velocidad de servicio es dada por sus correspondientes generadores de estafetas. La velocidad de servicio del buffer de salida es la suma de las velocidades sostenidas de cada canal virtual entrando al conmutador.

Esta tesis presenta un análisis del desempeño de tres diferentes tipos de mecanismos de vigilancia basados en Cubeta con goteo: Cubeta con goteo constante, Cubeta con goteo adaptivo y Cubeta con goteo adaptiva ayudada por predicción. Estos mecanismos de vigilancia son implementados para minimizar el número de violaciones del contrato de tráfico por medio del mecanismo de descarte, y también optimizar la utilización de los recursos que satisface la QoS acordada en las diferentes conexiones.

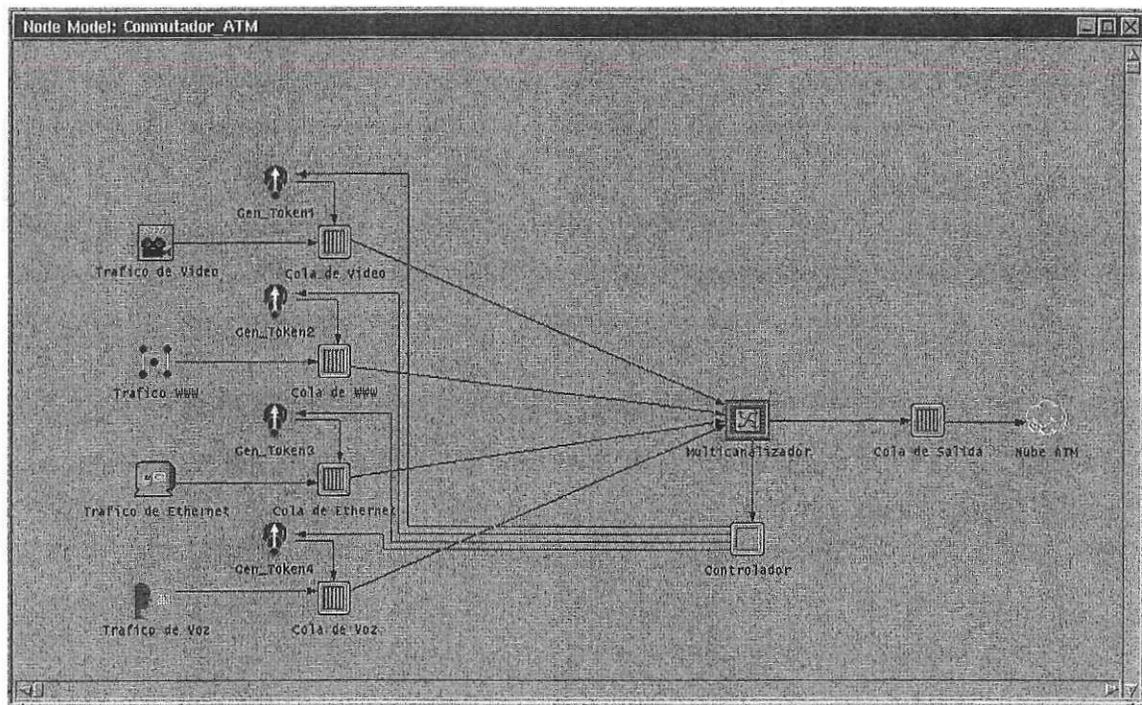


Figura 14. Modelo del conmutador ATM implementado en OPnet

VI.2 Análisis Nodal del conmutador ATM

Fuentes De Tráfico

Las fuentes de tráfico generan muestras artificiales mediante el algoritmo de generación rápida propuesto por Gallardo [Gallardo *et al.*, 1999] para emular tráfico de

Video, WWW y Ethernet, en cuyo caso se genera tráfico con distribución alfa-estable, basado en los parámetros determinados por Gallardo [Gallardo *et al.*, 1998a]

Los parámetros de entrada de los nodos generadores de tráfico alfa-estable se muestran en la tabla VIII.

El tráfico agregado de voz como se menciona en el capítulo 2, es emulado por fuentes individuales de voz que se suman para ingresar al conmutador. Los parámetros requeridos por cada fuente individual se presentan en la Tabla IX.

Tabla VIII. Parametros de entrada, para las fuentes de tráfico modeladas como procesos auto-similares con distribución alfa-estable.

	Tráfico agregado Video	Tráfico agregado WWW	Tráfico agregado Ethernet
Intensidad promedio de tráfico	913918.625	6497.85	268,452.8
Intensidad pico de tráfico	1351418	542,788	905,708
Modo [BAL/No Bal]	Balanceado	Anti-balanceado	Balanceado
Unidad de tiempo original	0.96	10	2
Ciclo de reloj	0.005	0.075	0.0005
H	0.903	0.833	0.834
Alfa	1.9	1.28	1.95

Tabla IX. Parámetros de entrada para las fuentes unitarias de voz modeladas por cadenas de markov

Parámetro de entrada	Valor
Intervalo de interrupción	[8kHz] ⁻¹
Duración promedio del periodo de silencio	1.35
Duración promedio del periodo de habla	1.00

Colas de entrada

Se compone de buffer para cada tipo de tráfico (con tamaño finito), que son servidos por la disciplina FIFO, además este nodo realiza el proceso del algoritmo de cubeta con goteo. Cuando una celda arriba a la cola de espera (proveniente de la fuente de tráfico), es insertado inmediatamente en la cola, si esta no se encuentra llena. Si arriba una estafeta procedente del generador de estafetas correspondiente, se incrementa el número de estafetas en el depósito. En caso de que exceda el tamaño correspondiente al depósito, esta estafeta sera eliminado. Si se encuentra una celda en la cola, y existe una estafeta disponible en el depósito, esta será *servida* y enviada al multicanalizador. Los parámetros de tamaño de colas de entrada en el conmutador ATM, se encuentran en la tabla V.

Generadores de estafetas

El nodo generador de estafetas recibe instrucciones del controlador como aumentar, disminuir o mantener una cierta tasa de generación de estafetas. En caso de que la instrucción del generador exceda los límites de generación marcados por la tabla IV del capítulo anterior, la estafeta se generará con la tasa límite indicada en el parámetro de simulación.

Los parámetros de simulación que se requieren en este nodo son: el valor inicial de tiempo de interarribo de generación de estafetas, límites máximo y mínimo en los cuales puede variar la tasa de generación; Los valores particulares para cada generador de estafetas son descritos por la tabla IV.

Multicanalizador

Realiza la función principal de un conmutador ATM, que es establecer una conexión virtual y conmutar celdas de un puerto de entrada hacia un puerto de salida, permitiendo así que estas lleguen a su destino. Para efectos de simulación el tráfico total converge al mismo puerto de salida. El parámetro requerido por este nodo es la velocidad de servicio requerida en bps.

Cola de Salida

Esta cola tiene una capacidad de almacenamiento de 210 celdas, y es servida bajo una disciplina FIFO.

Controlador

En este nodo reside la inteligencia del conmutador, aquí se realiza el manejo de recursos donde la entrada de este es el estado del conmutador, El controlador calcula los valores óptimos de asignación de ancho de banda (de acuerdo al algoritmo utilizado). Los algoritmos que este controlador puede utilizar son: Cubeta con Goteo, Cubeta con goteo adaptivo y Cubeta con goteo adaptivo ayudado por predicción. Los cuales fueron descritos en la sección IV.3.

VI.1.2 Resultados

A continuación se presentan resultados de la medición de desempeño del conmutador ATM.

Utilización del enlace de salida del conmutador ATM

Debido a que el objetivo de la tesis es el asignamiento dinámico de ancho de banda del enlace de salida, el parámetro principal de observación es el porcentaje de utilización

del enlace de salida; el cual nos permite determinar en forma cuantitativa la eficacia de los algoritmos simulados. La figura 15 muestra el porcentaje de utilización del enlace de salida del conmutador.

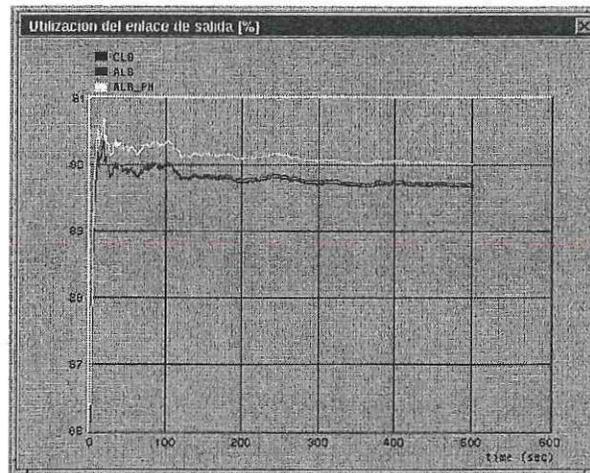


Figura 15. Porcentaje de utilización del enlace de salida del conmutador ATM

Desbordamiento de la cola de salida

Las pérdidas en el buffer de salida se deben a periodos de congestión dentro del conmutador, los cuales son provocados principalmente por la naturaleza de ráfagas del tráfico. El desbordamiento de la cola de salida nos brinda un índice de que tan bien trabaja el sistema, dado que estas pérdidas de celdas no son descartes por algún mecanismo de vigilancia, y al no tener control sobre las pérdidas estas pueden afectar la *imparcialidad* del sistema. La figura 16 muestra la estadística de celdas descartadas en un periodo de 500 segundos.

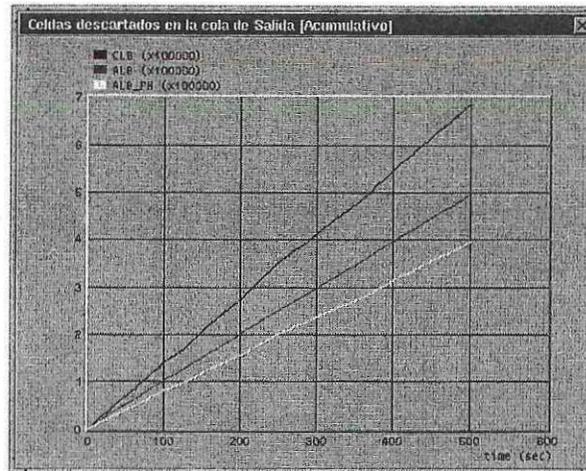


Figura 16. Estadística acumulativa de celdas descartadas en la cola de salida

Celdas descartadas en el conmutador

La figura 17 muestra la estadística del número celdas descartadas en el conmutador, en una ventana de 500 segundos. Esta estadística incluye las celdas descartadas en la cola de salida y en la cola de entrada.

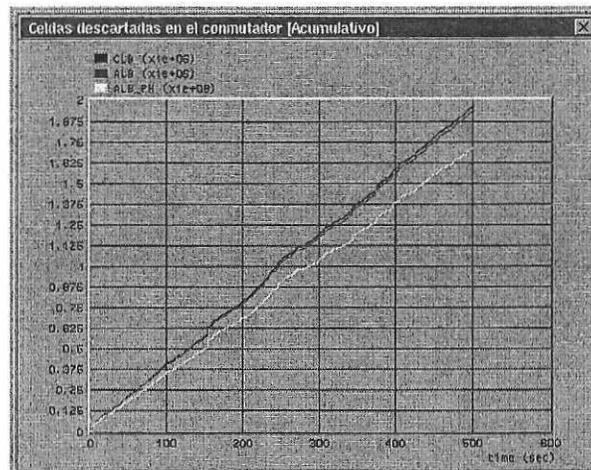


Figura 17. Estadística acumulativa de celdas descartadas en el conmutador

Retardo en tráfico sensitivo a tiempo.

Se recolectaron las estadísticas de retardo del tráfico sensitivo a tiempo como audio y video. Este retardo es calculado por la diferencia de tiempos desde el ingreso de la celda a la cola de entrada, hasta que esta celda es enviada por el enlace de salida. Donde cada punto en la gráfica representa el retardo sufrido por cada una de las celdas en el conmutador. Las figuras 18 y 19 muestran las estadísticas de retardo de las aplicaciones de voz y video respectivamente, donde fácilmente pueden ser observados el retardo máximo y mínimo sufrido por las celdas. Las figuras 20 y 21 muestran estadísticas del retardo promedio de celdas de voz y video.

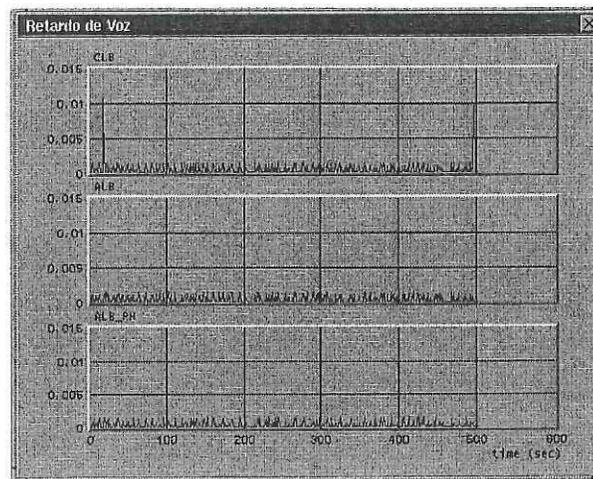


Figura 18. Retardo de voz en el conmutador ATM

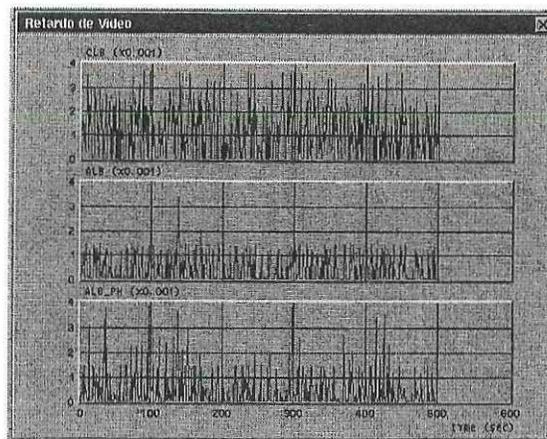


Figura 19. Retardo de video en el conmutador ATM

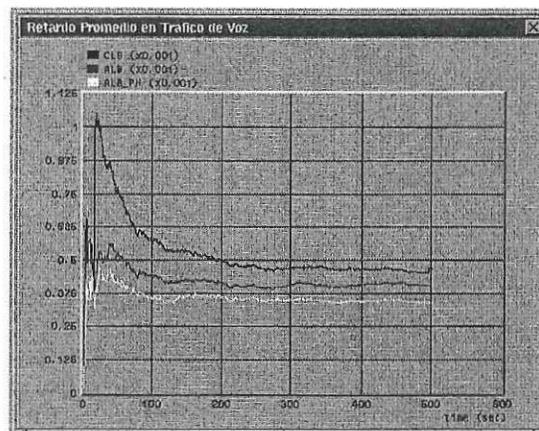


Figura 20. Retardo promedio en el tráfico de voz

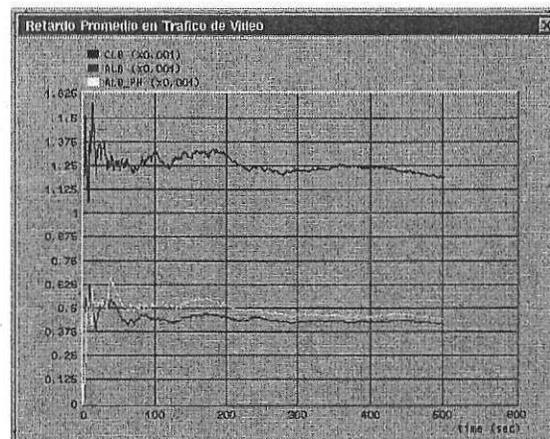


Figura 21. Retardo promedio en tráfico de video

Jitter en tráfico sensitivo a tiempo

Las figuras 22 y 23 muestran las variaciones de retardo en celdas de voz y video respectivamente. El jitter es un parámetro que permite evaluar el desempeño de aplicaciones que son sensitivas a las variaciones de retardo.

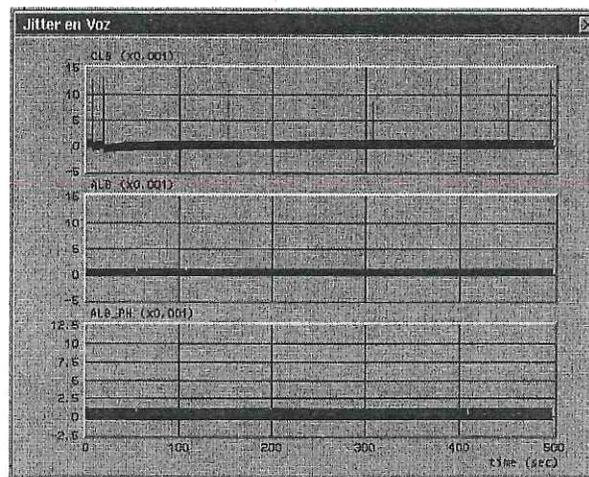


Figura 22. Jitter en celdas de voz

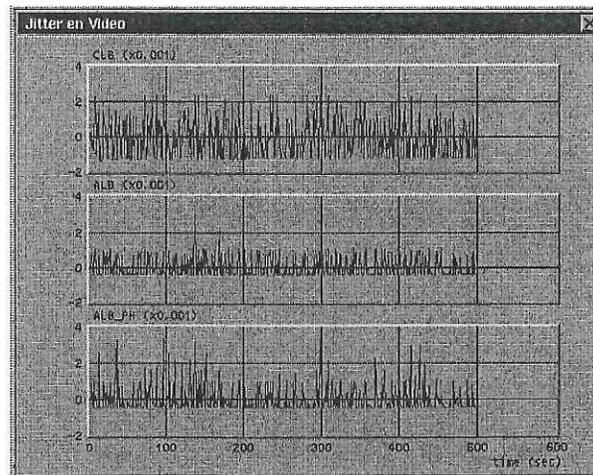


Figura 23. Jitter en celdas de video

VI.2 Implementación del enrutador IP

El enrutador IP fue implementado en OPnet como se muestra en la figura 24, donde se modela un nodo limitador que se localiza a la entrada de la red Internet que ofrece servicios diferenciados. Este enrutador ofrece servicios diferenciados bajo el modo de envío expedito por lo cual se simulan cuatro tipos de tráfico de acuerdo con el RFC 2475, estos cuatro tipos de tráfico al igual que en el conmutador ATM representan al tráfico sensible a tiempo y tráfico no sensible a tiempo.

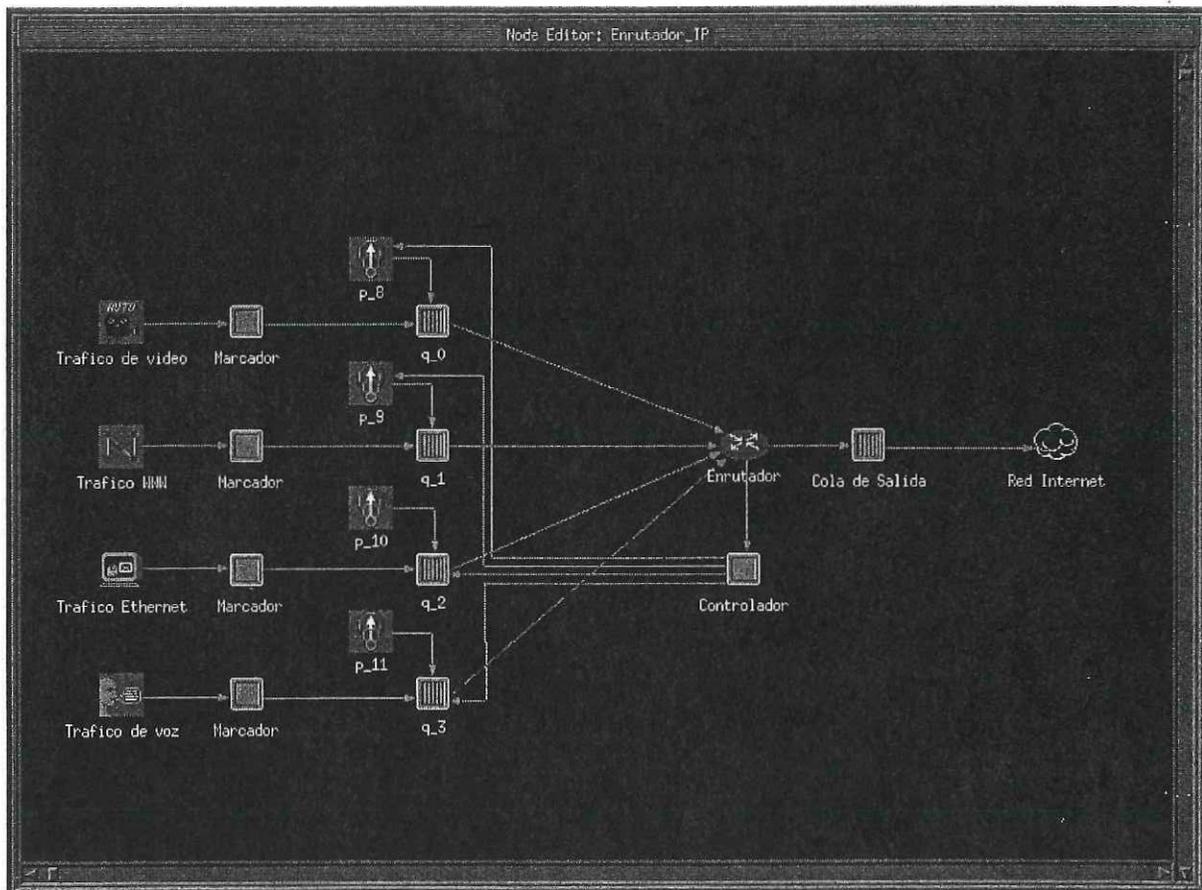


Figura 24. Modelo del enrutador IP implementado en OPnet

Este enrutador presenta algunas similitudes en simulación con el conmutador ATM, como la función de cola de salida y la carga de tráfico generada. Sólo difiere en el formato del paquete de voz, debido al tipo de compresión utilizado como se explica en la sección V.1 pero manteniendo el mismo valor de carga sostenida y de valor pico de tráfico.

VI.2.1 Análisis Nodal

Colas de entrada

Se compone de tres buffers de igual tamaño, donde cada uno almacena paquetes con la misma prioridad respecto a un mismo tipo tráfico, los buffers son servidos por la disciplina de asignación ponderada de colas. Este nodo (colas de entrada) emula el algoritmo de Cubeta con goteo.

Cuando un paquete arriba a la cola de espera es insertado en uno de los buffers dependiendo de su prioridad, si este se encuentra lleno tratará de insertar el paquete en un buffer de menor prioridad (en caso de que exista). Si arriba al nodo una estafeta se incrementa el número de estafetas en el depósito, en caso de que exceda el tamaño correspondiente al depósito, esta estafeta será eliminada.

Si se encuentra un paquete en la cola, y existe una estafeta disponible en el depósito, esta será *servida* y enviada al multicanalizador. Los parámetros de tamaño de colas de entrada en el enrutador IP se encuentran en la tabla VII.

Generadores de estafetas

El nodo generador de estafetas recibe instrucciones del controlador; donde aumenta, disminuye o mantiene una cierta tasa de generación de estafetas. En el caso de que la instrucción del generador exceda los límites de generación marcados por la tabla 6 del capítulo anterior, se generara con una tasa del límite.

Los parámetros de simulación que requieren son: el valor inicial de tiempo de interarribo de generación de estafetas, límites máximo y mínimo en los cuales puede variar la tasa de generación; Los valores particulares para cada generador de estafetas son descritos por la tabla VI.

Enrutador

Realiza la función principal de un enrutador IP, enrutar un paquete recibido en un puerto de entrada hacia un puerto de salida en base a la dirección destino del paquete y a la tabla de enrutamiento. Para efectos de simulación el tráfico total converge al mismo puerto de salida. El parámetro requerido por este nodo es la velocidad de servicio cuyas unidades son bps.

Cola de Salida

Esta cola tiene una capacidad de almacenamiento de 12,180 bytes, servida bajo una disciplina FIFO.

Controlador:

En este nodo reside la inteligencia del enrutador, aquí se realiza el manejo de recursos, donde la entrada de este controlador es el estado del conmutador, y obtiene los valores óptimos (de acuerdo al tipo de algoritmo utilizado) de asignamiento de ancho de

banda. Los algoritmos que este controlador puede utilizar son: Cubeta con Goteo, Cubeta con Goteo adaptivo y Cubeta con Goteo adaptivo ayudado por predicción. Los cuales fueron descritos en la sección IV.3.

VI.2.2 Resultados

Utilización del enlace de salida del enrutador IP

La figura 25 muestra la estadística de porcentaje de utilización del enlace de salida, en la cual se observa el desempeño de los algoritmos CLB, ALB y ALB-PH implementados en el controlador del enrutador IP. Estos resultados serán analizados en la sección VI.3.

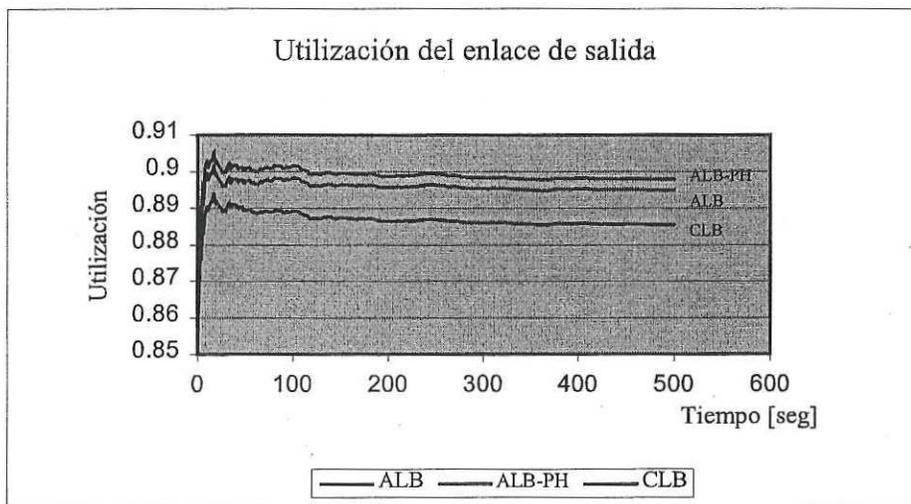


Figura 25. Porcentaje de utilización del enlace de salida del enrutador IP

Bytes descartados en el enrutador

La figura 26 muestra el número de bytes descartados en el enrutador bajo los diferentes algoritmos simulados. Esta estadística es estimada del número de bytes que

fueron descartados por no tener espacio suficiente para ser almacenados en el buffer de entrada o en el buffer de salida.

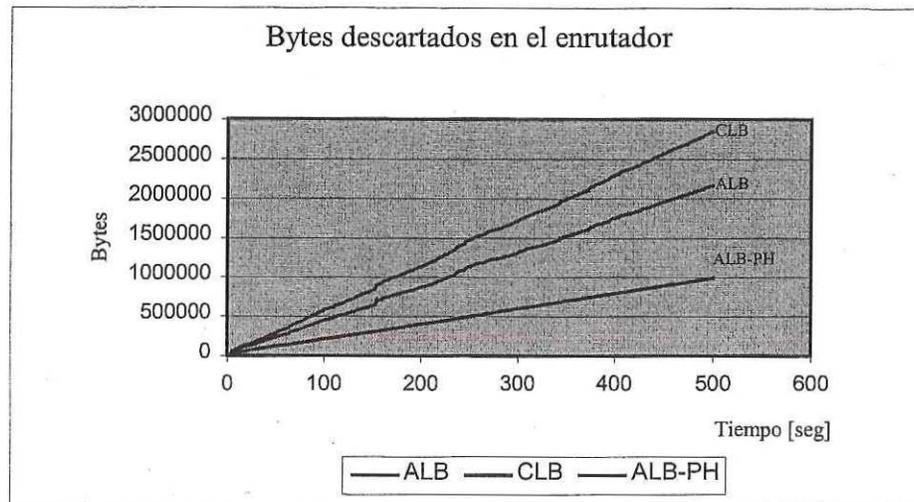


Figura 26. Bytes descartados en el enrutador

La figura 27 muestra el número de bytes descartados en la cola de salida del enrutador. Esta estadística muestra

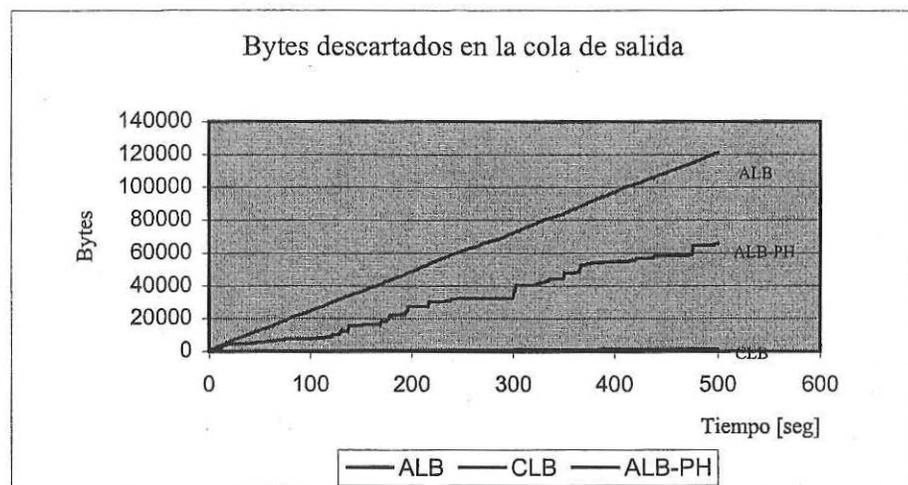


Figura 27. Bytes descartados en la cola de salida

Retardo en tráfico sensitivo a tiempo

Las figuras 28, 29 y 30 muestran el retardo de video sufrido en el enrutador, bajo los algoritmos CLB, ALB y ALB-PH respectivamente.

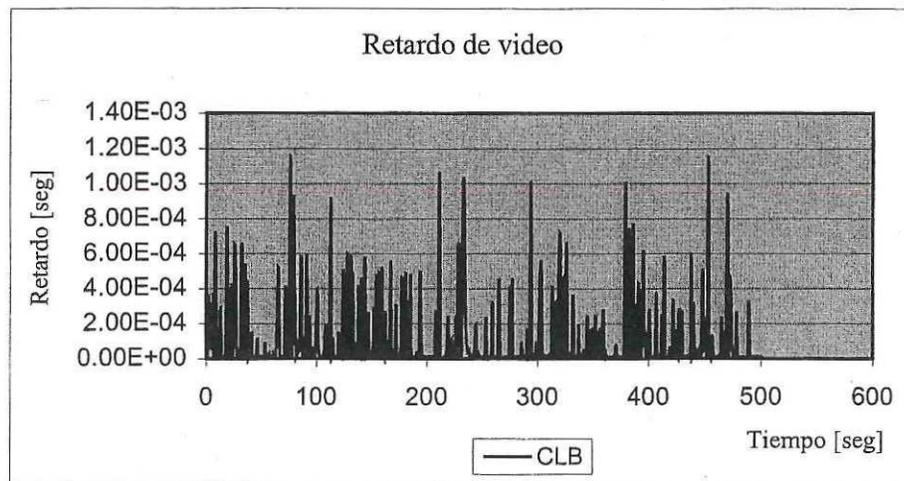


Figura 28. Retardo de video bajo el algoritmo CLB

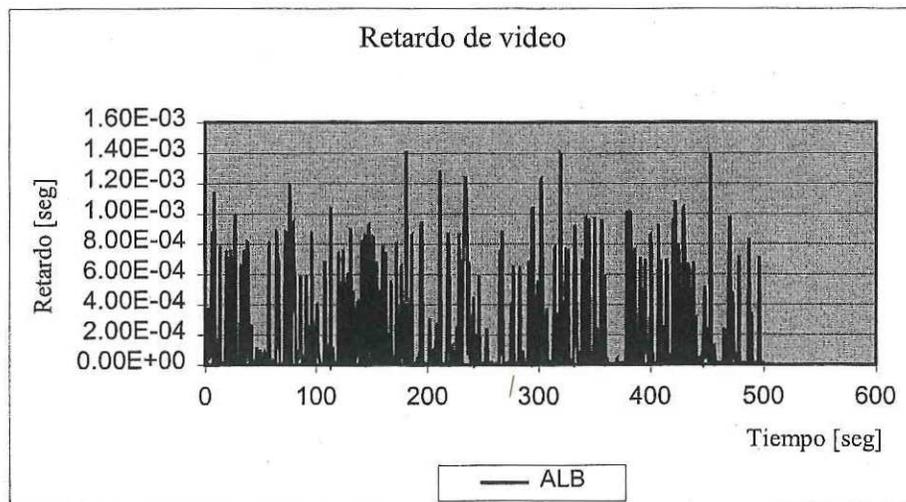


Figura 29. Retardo de video bajo el algoritmo ALB

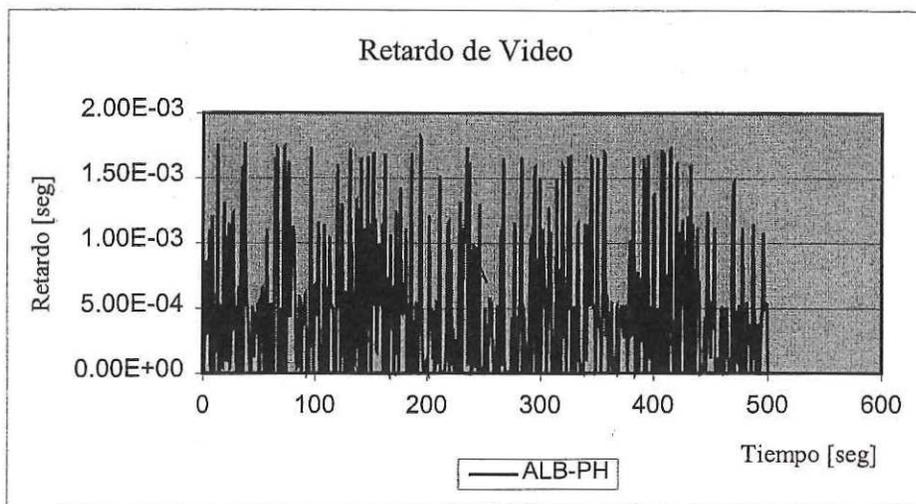


Figura 30. Retardo de video ALB-PH

La estadística del retardo sufrido por los paquetes de voz a través del enrutador fue colectada bajo una ventana de 500 segundos. Este parámetro se evalúa para los algoritmos CLB, ALB y ALB-PH mostrándose los resultados en las figuras 31, 32 y 33 respectivamente.

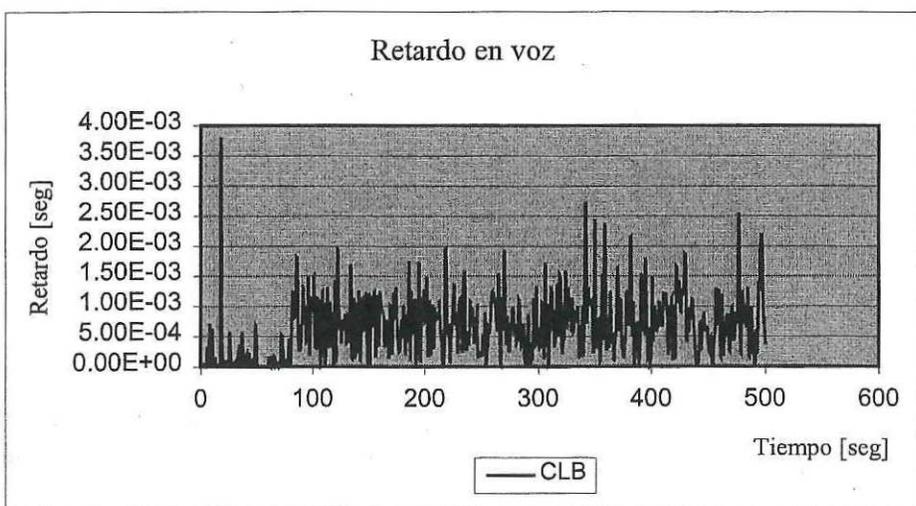


Figura 31. Retardo en voz bajo el algoritmo CLB

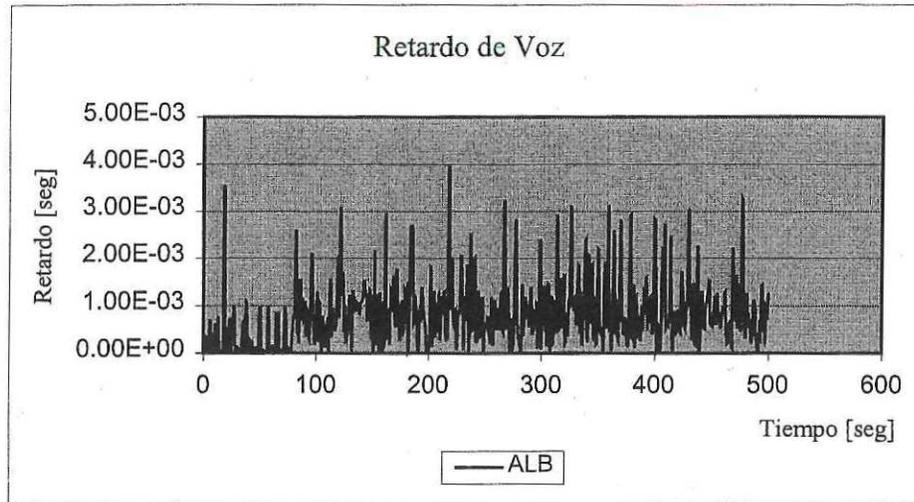


Figura 32. Retardo en voz bajo el algoritmo ALB

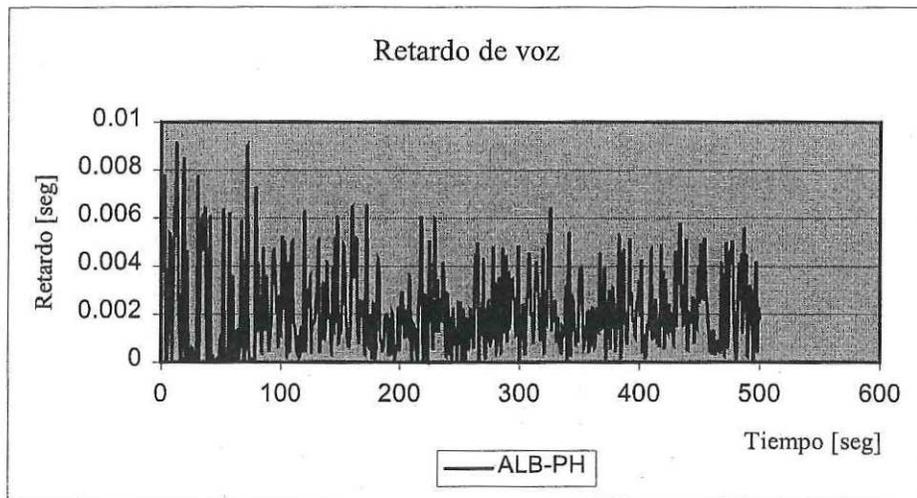


Figura 33. Retardo en voz bajo el algoritmo ALB-PH

VI.3 Análisis de resultados

Desbordamiento de la cola de salida.

Esta estadística es tomada en el buffer de salida representado en las figuras 14 y 24 como cola de salida, mostrando las celdas que son rechazadas dado que el buffer se encuentra lleno. La figura 16 muestra el comportamiento de descarte de celdas bajo los diferentes algoritmos simulados. Esta gráfica representa el número de celdas descartadas en forma acumulativa, lo cual nos permite realizar una comparación cuantitativa entre los diferentes algoritmos de vigilancia de una forma mas sencilla.

Para el algoritmo CLB, a los 500 segs se tiene una pérdida en el buffer de salida en forma acumulativa 680,000 celdas. El algoritmo ALB, realiza un notable decremento de desbordamiento de aproximadamente 500,000 celdas. Por último en el algoritmo de ALB-PH se tiene una perdida de 400,000 celdas aproximadamente el 58% de las pérdidas ofrecidas por CLB.

Estadística del total de celdas perdidas

La figura 17 muestra un comportamiento muy parecido entre los algoritmos CLB y ALB en el conmutador ATM. El número de celdas perdidas en un periodo de 500 seg es de $1.8E+05$ para el caso de CLB, donde se reduce aproximadamente el 16% de las pérdidas para el caso de ALB-PH.

Nótese que este número que suena extremadamente alto es debido a que se fuerza al sistema a trabajar bajo presión.

En el enrutador IP, se observa que el algoritmo ALB muestra una reducción del número de bytes descartados al 40% de lo presentado por el algoritmo CLB.

Utilización del enlace de salida

En el conmutador ATM, la figura 15 muestra un índice de comportamiento de los algoritmos de vigilancia, se observa una mayor utilización de enlace de salida cuando se utiliza el algoritmo de cubeta con goteo ayudado por predicción, seguido por el mecanismo ALB y CLB. En el caso del enrutador IP (figura 25) se observa el mismo comportamiento.

Retardo en tráfico sensitivo a tiempo.

Se recolectaron las estadísticas del retardo sufrido a través del conmutador o enrutador para el tráfico de voz y video, dado que son aplicaciones sensitivas a retardo. Para el caso de voz (figura 20), se considera como un máximo retardo admisible 20 mseg. Sin embargo para el retardo sufrido por el conmutador se tiene un máximo de 2mseg en el caso de video. A lo largo de una ventana de observación de 500 sec. se tiene un promedio de 0.48 mseg. de retardo para CLB, el algoritmo ALB muestra un retardo en el conmutador de 0.42mseg, por último el mecanismo ALB-PH obtiene 0.37 mseg.

Para el caso de video (figura 21), se considera como un máximo retardo admisible de 20 msec. Sin embargo dentro del conmutador, el retardo máximo para tráfico de video es de 2mseg. Para CLB a lo largo de una ventana de 500 sec. se observa el promedio de 1.2mseg. de retardo. Para ALB-PH se tiene un promedio de 0.43mseg., en el caso de ALB el retardo promedio es de 0.4 mseg.

VII. CONCLUSIONES

VII. Conclusiones

El utilizar modelos no gaussianos para generación de tráfico permite confiar en que los buenos resultados obtenidos en la simulación del modelo de asignación de recursos dentro de un conmutador o en un enrutador no son debidos al tráfico *ideal* simulado.

La disciplina de servicio a las colas de espera, así como el esquema de vigilancia cubeta con goteo, no tiene un efecto significativo sobre las propiedades de auto-similitud, lo cual implica que el algoritmo ALB-PH puede ser implementado en cualquier lugar dentro de la red.

En el algoritmo ALB la decisión de cambiar la razón de generación de estafeta es realizada después de que muy pocos o que muchos paquetes son recibidos, lo que crea un bajo desempeño del nodo de la red, especialmente para tráfico con alto contenido de ráfagas.

La adopción de la predicción de tráfico brinda la ventaja de realizar decisiones en base al futuro estado del conmutador o enrutador. El gran rango de dependencia en procesos aleatorios provee un modelo mas cercano a la realidad para varias categorías de tráfico.

El servicio de mejor esfuerzo que ofrece la red Internet es tratar a cada paquete de la misma forma, por lo que no puede asegurar una calidad de servicio, dado que el servicio que provee esta red es altamente variable y depende directamente del estado en que esta se encuentre.

VIII.2 Aportaciones

Las aportaciones del presente trabajo de tesis son sustentadas en la evaluación del desempeño de los modelos propuestos realizada bajo simulación de tráfico real. Las cuales se mencionan a continuación.

- La propuesta de un nuevo enfoque de los diferentes mecanismos de vigilancia, para ser utilizados como mecanismos de manejo de recursos.
- Diseño y evaluación del desempeño de un esquema de asignación dinámica de recursos enfocado a satisfacer garantías de calidad de servicio en redes de comunicación de datos.
- Implementación en OPnet del algoritmo de predicción de tráfico auto-similar con distribución alfa-estable.
- Simulación y evaluación del desempeño de un nodo Internet proveedor de servicios diferenciados
- Implementación de un controlador de conmutador basado en el esquema de asignación dinámica de recursos basado en predicción de tráfico. Por último cabe mencionar la implementación de un esquema que permita la sana convivencia de aplicaciones limitadas en tiempo y no limitadas en tiempo.

VIII.3 Trabajo a futuro

- La propuesta de una red que permita proveer calidad de servicio de extremo a extremo en forma garantizada, mediante la incorporación de servicios diferenciados como *backbone* de una red completa que ofrezca en los extremos de usuario servicios integrados (como lo propone el RFC 2275).
- El análisis de las variaciones en los parámetros de predicción de tráfico alfa-estable, realizando un análisis cuantitativo de su influencia en el comportamiento de los algoritmos de vigilancia.
- La implementación del protocolo de transporte TCP utilizado para servicios no limitados en tiempo en las redes Internet, así como el protocolo RTP/RTCP en el caso de aplicaciones limitadas en tiempo como voz y video.
- Extender la convivencia de tráfico ATM, no sólo a servicios VBR-rt y VBR-nrt, creando un modelo que se base en predicción para realizar la retroalimentación del tráfico ABR.

LITERATURA CITADA

- Angulo M., Gallardo J., Makrakis, 1999. "*Adaptive QoS Scheme Based on Prediction of Alpha-Stable Self-Similar Traffic*", IEEE/IEE International Conference on Telecommunication, Atenas Grecia,
- Angulo M., Gallardo J., Makrakis, 2000. "*Dynamic Bandwidth Allocation Schemes for Internet Node providing Differentiated Services*", por presentar en ICT2000, Acapulco México
- Aurrecochea C., Campbell A. Hauw linda, 1995, "*A survey of Quality of Service Architectures*", Distributed Multimedia Research Group, University of Lancaster, Internal report number MPG-95-18.
- Bates Stephen, Steve Mc Laughlin, 1998, "*The effective Bandwidth of stable Distributions*", ICASSP 98.
- Beran, J. Sherman, R., Taqqu M.S. and Willinger W, 1995, "*Long range Dependence in Variable Bit Rate Video Traffic*", IEEE Transactions on Communications, Vol 43, No. 2-4, pp 1566-1579.
- Bernet Y. Smith A., Blake S., 1999, "*A conceptual Model for DiffServ Routers*", Internet Engineering Task Force, Draft, Trabajo en progreso.
- Bertsekas Dimitri, Gallager Robert, 1992. "*Data Networks*", 2nd Edition, Prentice Hall USA, pp 556.

Blake S., D. Black, M. Carlson, 1998, "*An Architecture for Differentiated Services*", Internet Engineering Task Force, RFC 2475.

<ftp://ftp.isi.edu/in-notes/rfc2475.txt>

Bolot J.C y Adreas Vega-Garcia, 1996, "*Control Mechanisms for Packet Audio in the Internet*", Proceedings of IEEE Infocomm, pp 232-239.

Bolot J.C y T. Turletty, 1994, "*A rate control scheme for packet video in the Internet*", Proceedings of IEEE Infocomm, pp 1216-1223.

Campbell A., Coulson G., Hutchison D. 1994, "*A Quality of Service Architecture*", ACM Computer Communications Review, Volumen 24, Número 2, pp. 6-27, MPG Internal report number MPG-94-08.

Clark D., Braden B., Shenker S., 1994, "*Integrated Services in the Internet Architecture*", Internet Engineering Task Force, RFC 1633, IntServ Working Group.

http://www.iit.nrc.ca/IETF/RSVP_survey/ietf_rsvp-qos_survey_02.txt

Clark David, Shenker Scott, Lixia Zhang, 1992, "*Supporting Real-Time Applications in an Integrated Services Packet Networks: Architecture and Mechanism*" ACM Sigcomm proceedings.

Daigle John N., Langford Joseph d., 1986, "*Models for Analysis of Packet Voice communications systems*". IEEE Journal on selected areas in communications VOL SAC -4, No. 6 Septiembre 1986. 847- 854

Demers A., Keshav S., Shenker S, 1990, "*Analysis and Simulation of a Fair Queueing Algorithm*", internetworking: Research and Experience, Vol 1. No. 1. pp 3-26.

- Demeure I., Farhat J. Gasperoni F., 1996, "*A scheduling Framework for the Automatic Support of Temporal QoS Constraints*", Proceeding of the Fourth International Workshop on Quality of Service, Paris.
- Dharanikota S., Maly Kurt , 1996, "*QUANTA: Quality of Service Architecture for Native TCP/IP over ATM Networks*", HPDC'96 Proceedings; also Old Dominion University Department of Computer Science Technical report # TR-96-01, February 1996.
- Erramilli A., Narayan O. and Willinger W., 1996, "*Experimental Queuing Analysis With Long-Range Dependent Traffic*", IEEE Transactions on Networking, April 1996, p 209-220
- Ferguson Paul, Huston Geoff, 1998, "*Quality of Service-Delivering QoS on the Internet and in Corporate Networks*", John Wiley Computer Publishing, 1998, USA.
- Foro ATM, 1999a, "*A practical guide to carrying voice over ATM*".
- Foro ATM, 1999b, Grupo de trabajo en manejo de tráfico, "*Traffic Management Specification Version 4.1*", Marzo 1999
<http://www.atmforum.com/atmforum/specs/approved.html>
- Gallardo J.R., Makrakis D., Orozco-Barbosa L., 1998a, "*Fast Generation of Artificial Traces of Alpha-stable Long Range Dependent Stochastic Processes*", submitted to IEEE transactions on Signal Processing (Sept. 1998).

- Gallardo José R, Makrakis Dimitris, Orozco-Barbosa Luis, 1999, "*Prediction of Alpha-Stable Long-Range Dependent Stochastic Processes*". IEEE/IEE International Conference on Telecommunications ICT'99. Cheju, Korea, 1999.
- Gallardo José R., Makrakis Dimitris, Orozco-Barbosa Luis, 1998b "*Use of Alpha-Stable Self-Similar Stochastic Processes for Modeling Traffic in Broadband Networks*". SPIE's Symposium on Voice, Video and Data Communications within the Performance and Control of Network Systems program.1998. Boston, Massachusetts, USA.
- Gilligan R, and Callon R, 1995, "*Ipv6 Transition Mechanism Overview*", Connexions, Octubre 1995.
- Giroux Natalie, Sudhakar Ganti, 1999, "*Quality of Service in ATM networks*", Prentice Hall, 1999 USA.
- Heffes Harry, Lucantoni David, 1986 "*A Markov Modulated Characterization of Packetized Voice and Data Traffic and Related Statistical Multiplexer Performance*", IEEE Journal on selected areas in communications VOL SAC -4, No. 6 Septiembre 1986. pp 856- 867
- Heinanen Juha, 1999, "*A single Rate three color marker*", Internet Engineering Task Force, Internet Draft Trabajo en progreso, Marzo 1999. [3 colors]
- Hinden, 1994, "*IP Next Generation Overview*", Internet Engineering Task Force, Internet Draft, Trabajo en progreso, Octubre 1994. [Ipv64]

Hong Duke, Suda Tatsuya, 1991, "*Congestion Control, and Prevention in ATM networks*", July 1991, IEEE Network Magazine, pp10-16.

Karlsson Gunnar, 1996, "*Asynchronous Transfer of Video*", IEEE Communications Magazine, Agosto 1996 pp 118-126

Kruntz Mrwan, 1999, "*Bandwidth Allocation Strategies for Transporting Variable -Bit-Rate Video Traffic*", IEEE Communications Magazine, Enero 1999 pp 40-46

Ibañez, Nichols K, 1998, "*Preliminary Simulation Evaluation of an Assured Service*", Internet Draft, Agosto 1998.

ITU, 1990, Reg. G.726 "*40, 32, 24, 16 kbit/s Adaptive Differential Pulse Code Modulation (ADPCM)*", June 1990.

ITU, 1996b Reg. G.729 "*Coding speech at 8 kbit/s Using Conjugate-Structure Algebraic-Code-Excited Linear-Prediction CS-ACELP*", Marzo 1996

ITU, 1996b, Reg. G.729A "*Reduced Complexity 8kbit/s CS-ACELP Speech codec*", Noviembre 1996.

Jena Ajit K., Popescu Adrian, Parag Pruthi, Ashok Erramilli, "*Traffic Control in ATM Networks: Engineering Impacts of Realistic Traffic Process*", in Proceedings of Nordisk Teletrafik Seminarium NTS - 13, Trondheim, Norway, August 1996.

Kasiolas A., *"Intelligent Control in Multimedia Traffic Policing, Shaping and Congestion Avoidance over Broadband Networks"* M.Eng.Sc. Thesis, The University of Western Ontario, London, Ontario, Canada, April 1999.

Kim Hyogon, *"A fair Marker"*, Internet Engineering Task Force, Internet Draft Internet Engineering Task Force, Trabajo en progreso, April 1999.

[2colors

Kostas T., Borella M., Sidhu I, Shuster G., Grabiec J. "Real Time voice over packet – Switched networks", IEEE Network, January/February 1998.

Kyas Othmar, *"ATM networks"*, International Thomson Publishing, Inglaterra, 1995.

Lambarelli Livio, *"ATM Service Categories: The benefits to the user"*. CSELT, Torino, Italy, Caso de Estudio,

http://www.atmforum/library/service_categories.html

Chris Metz, *"IP QoS: Traveling in First Class on the Internet"*, IEEE Internet Computing, Marzo-Abril 1999.

Parekh a., Gallager R. "A generalized processor sharing approach to flow control in integrated services networks: the single-node case," IEEE/ACM Transactions on Networking, Vol. 1, No. 3, June 1993, pp 344-357.

Pandya Abhijit S., Sen Ercan, *"ATM Technology for Broadband Telecommunications Networks"*, 1999, CRC Press, USA, pp 282

- Perkins C., Hodson O., Hardman V. "*A survey of Packet Loss Recovery Techniques for Streaming Audio*" IEEE Network Magazine, September/October 1998, pp40-47
- Pruthi P., Popescu A., "*Effect of Controls on Self-Similar Traffic*", Technical Report, Department of Telecommunications and Mathematics, University of Karlskrona, Suiza, 1997.
- Raj Jain, "*Congestion Control and Traffic Management in ATM Networks, Recent Advances and Survey*", Department of Computer and Information Science, The Ohio State University, Draft Version, Agosto1996.
- Ramjee, R, Kurose J. Towsley D, Schulzrin, "*Adaptive Playout Mechanisms for Packetized Audio Applications in Wide Area Networks*", Proceedings of IEEE INFOCOMM , 1997, pp 680 688. [Voice 2]
- Rathgeb Erwin P. "*Modeling and Performance Comparison for Policing Mechanism for ATM Networks*", IEEE Journal on Selected Areas in Communications, Vol 9, No3. Abril 1991.
- Rexford Jennifer, Bonomi Flavio, Greenberg Albert, Wong Albert, "*Scalable architectures for Integrated Traffic Shaping and Link Scheduling in High-Speed ATM switches*". IEEE Journal on selected Areas in Communications, Vol 15, 938-950, 1997
- Sriram K., P.K. Varshney and J. G. Shanthikumar, "*Discrete-time analysis of integrated voice-data multiplexers with and without speech activity detection*", IEEE J.

Select. Areas Commun, vol SAC-1 Special Issue on Packet switched voice and Data communications, Dec 1983

W. Stallings, "*Data and Computer Communications*", 5th Edition, Upper Saddle River, NJ. Prentice Hall, 1996. [Ipv61]

Taqqu Murad , Samorodnisky Gennady, "*Stable non-gaussian Random process*", Ed. Chapman & Hall, 1994.

