

Soluciones para acelerar servicios IP y optimizar el acceso a Internet

Jornadas Técnicas RedIRIS 2001 y XII Grupos de Trabajo
Pamplona, 22-26 Octubre 2001

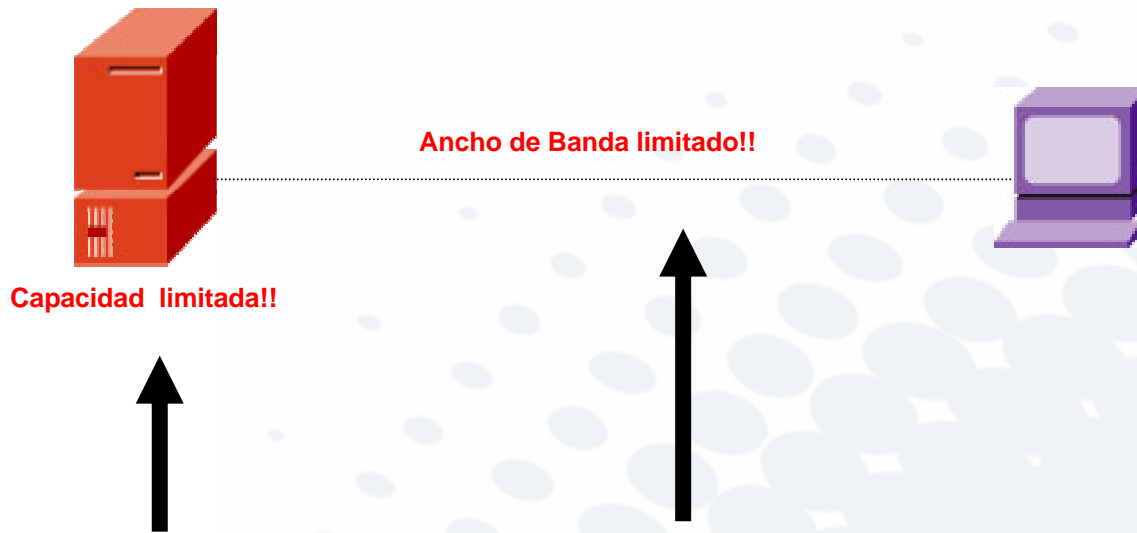
Unitronics Comunicaciones

Adolfo García Yagüe
agy@unitronics.es



v.3. Octubre 2001
v.2. Noviembre 2000
v.1. Abril 2000

Puntos clave de la aceleración y optimización



Capacidad limitada!!

Ancho de Banda limitado!!

Solución:

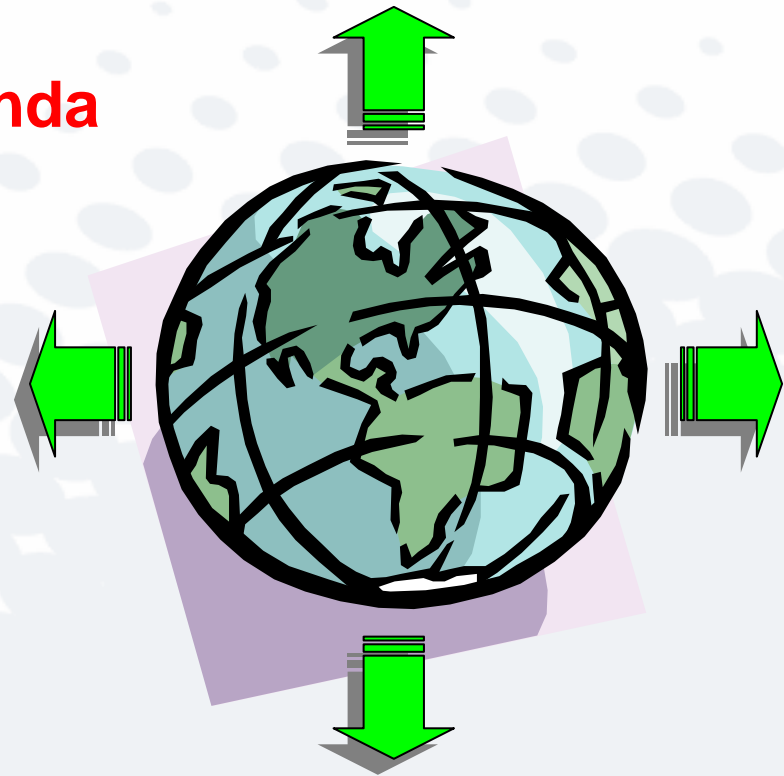
- Balanceo de Carga
- Aceleración de SSL
- Proxy Inverso

Solución:

- Priorización de aplicaciones y usuarios
- Compresión del tráfico
- Optimización de TCP slow start
- Caching de contenidos
- HTTP 1.1 Pipelining

Agenda

- **Gestores de Ancho de Banda**
- Caché de Contenido
- Balanceadores de Carga
- Aceleradores SSL



Gestores de Ancho de Banda

- Aportan a IP un mecanismo de Clase de Servicio, ofreciendo un tratamiento diferenciado a usuarios y aplicaciones
- Fases y parámetros involucrados en la Clase de Servicio:
 - Definición de políticas
 - Clasificación de tráfico
 - Etiquetado
 - Gestión de *buffers* (colas)
 - Cálculos de latencia, control de flujos y ventana (Control TCP)
 - Monitorización
 - Reporting
- Aproximaciones tecnológicas:
 - Modelo de priorización en el tratamiento de unos datos frente a otros a través de la gestión de colas
 - Control de flujo TCP y negociación del tamaño de ventana

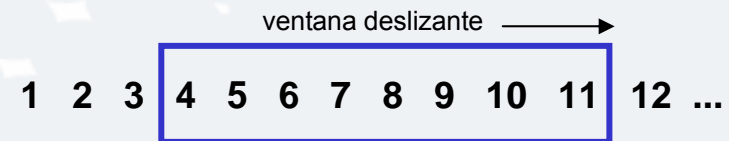
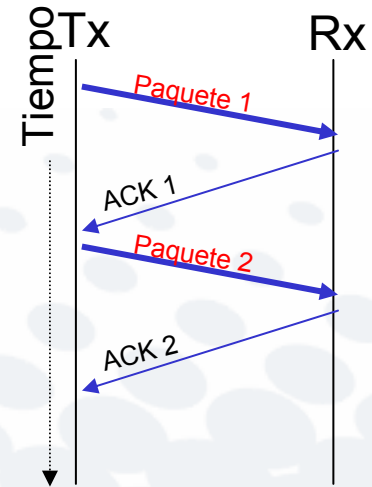
¿Como administra TCP el ancho de banda?

- **Entrega confiable**

Tras la recepción del paquete, el receptor envía un mensaje de confirmación al emisor (ACK)

- **Control de Flujo**

Mecanismo extremo a extremo que permite al receptor restringir la transmisión hasta que disponga de recursos para procesar el tráfico



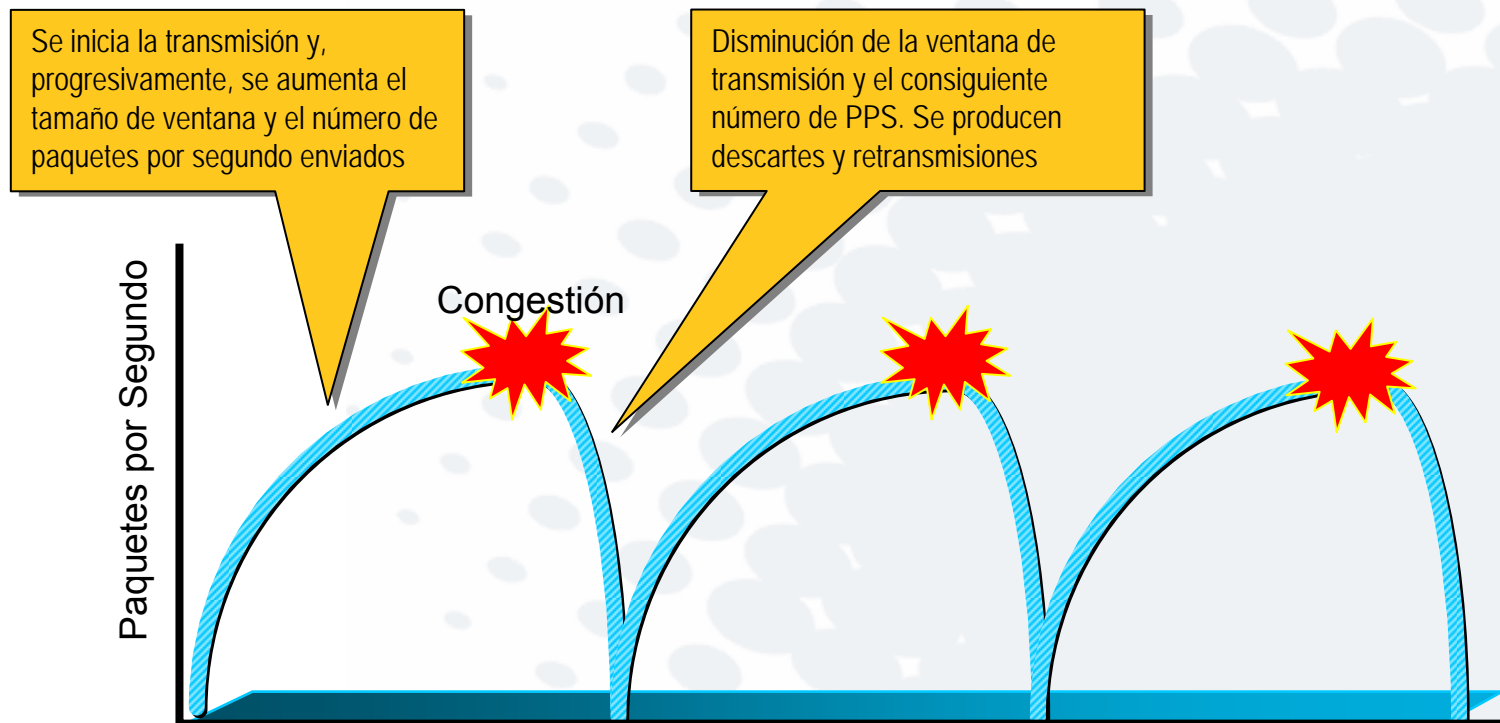
El emisor establece una ventana dentro de la cual se encuentran todos los bytes que ha enviado pero sobre los que aún no ha recibido acuse de recibo. Según va recibiendo acuses de recibo va deslizando la ventana hacia la derecha.

El receptor establece el tamaño de la ventana. TCP permite que el tamaño de ventana cambie. Con cada acuse de recibo se incluye un campo de aviso de ventana donde el receptor puede indicar al emisor el nuevo tamaño.

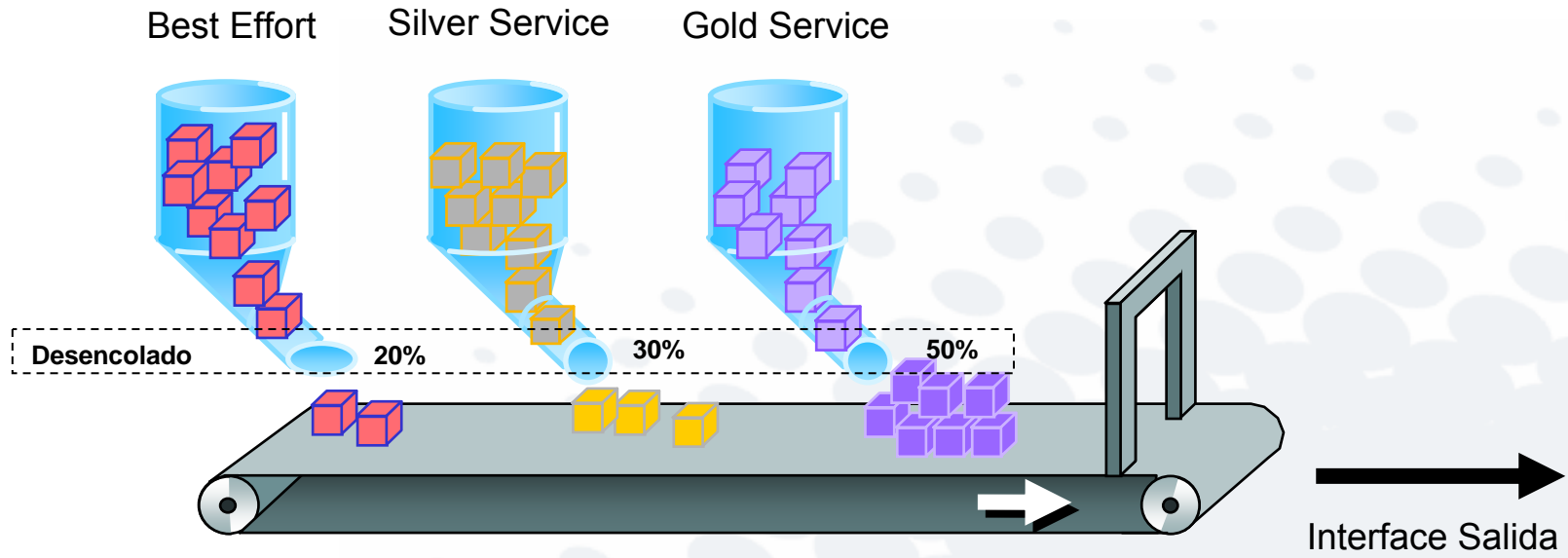
¿Como administra TCP el ancho de banda?

- **Control de Congestión**

Mediante la sincronización TCP o *slow-start* el emisor comienza la transmisión de sólo un paquete, tras la recepción del ACK va incrementado el caudal. En caso de existir una congestión en la ruta no recibirá ACKs, en cuyo caso disminuye el envío de paquetes.



Gestor de ancho de banda basado en colas



El empleo de buffers o colas en el interface de salida, junto con información de nivel de servicio y el adecuado algoritmo de desencolado, permiten priorizar la salida de un determinado tráfico frente a otros

Gestor de ancho de banda dotado de control de flujo y gestión de tamaño de ventana TCP

Estos gestores de ancho de banda desarrollan las capacidades nativas de TCP orientadas al control de la conexión:

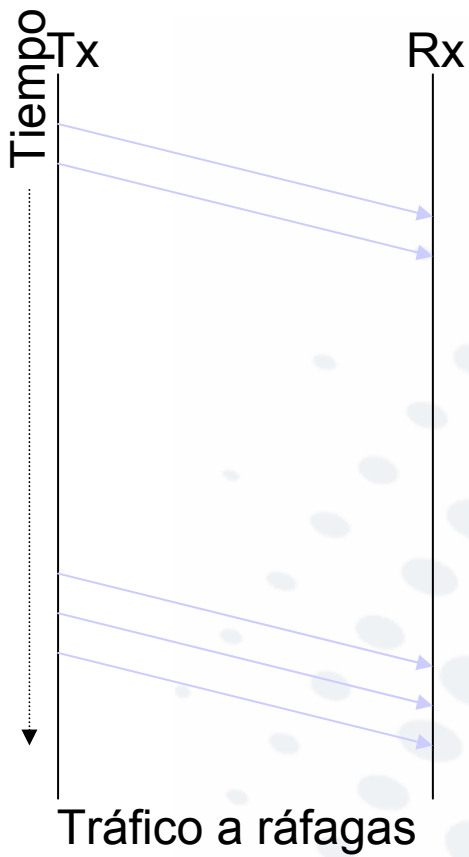
- Acuse de recibo (ACK)
- Control de flujo
- Sincronización TCP

Aportando además:

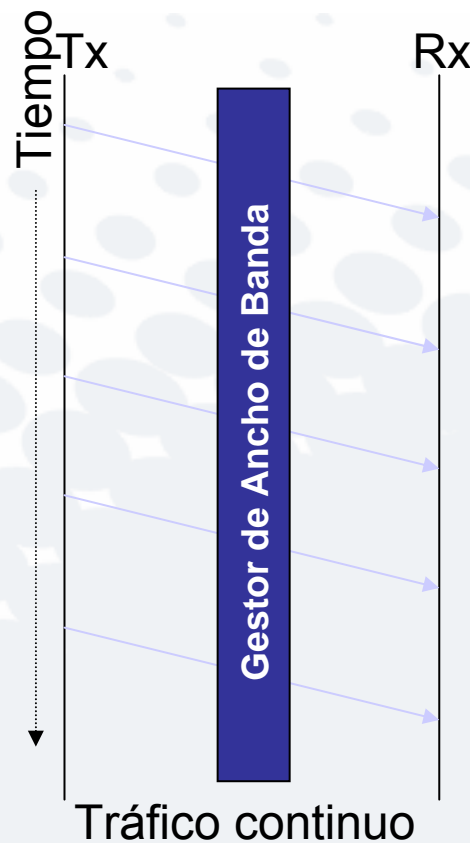
- Medición del tiempo de latencia extremo a extremo
- Cálculo de la cantidad de paquetes a transmitir para reunir las condiciones de latencia y flujo garantizado
- Negociación del tamaño de ventana
- Controla la transmisión de nuevas tramas mediante la interpretación de los ACKs

Gestor de ancho de banda dotado de control de flujo y gestión de tamaño de ventana TCP

Sin gestor de ancho de banda



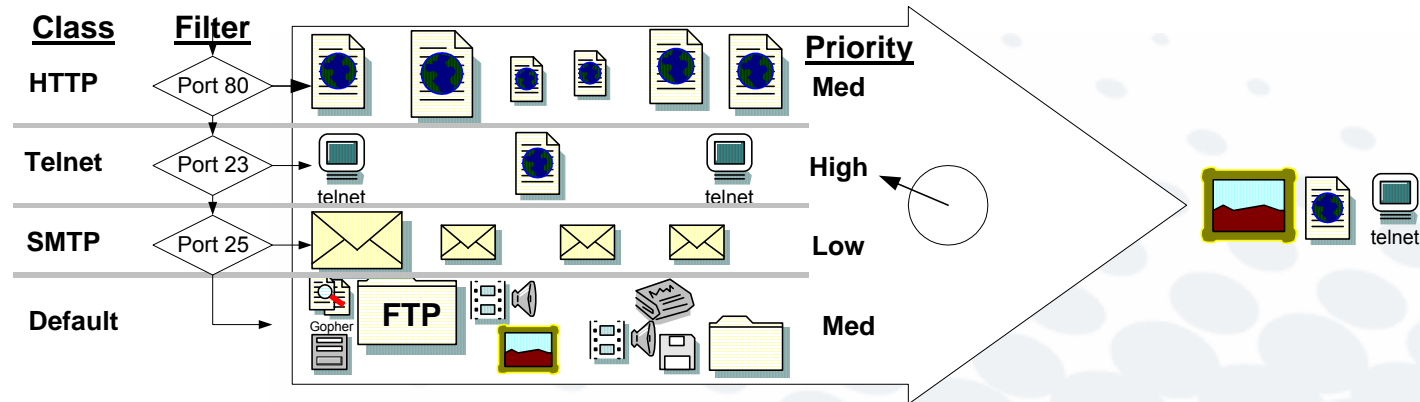
Con gestor de ancho de banda



Gestión de colas frente al control TCP

- Tecnológicamente no hay una solución superior a otra
- Evaluar la aplicación: ¿Priorizar o gestión extremo a extremo?
- Para priorizar ambas soluciones son apropiadas
- En redes como Internet (múltiples equipos intermedios y best effort), si se desea llevar un control preciso extremo a extremo, son más eficientes las soluciones basadas en el control de TCP
- El control TCP se desarrolla para cada conexión TCP, para conseguir la misma capacidad en la gestión de colas es preciso disponer de una cola por conexión
- Sólo el modelo de colas permite la gestión del tráfico UDP y no-IP
- Las funcionalidades Traffic Shaping y Policing son más precisas en el modelo de colas con algoritmos de WFQ

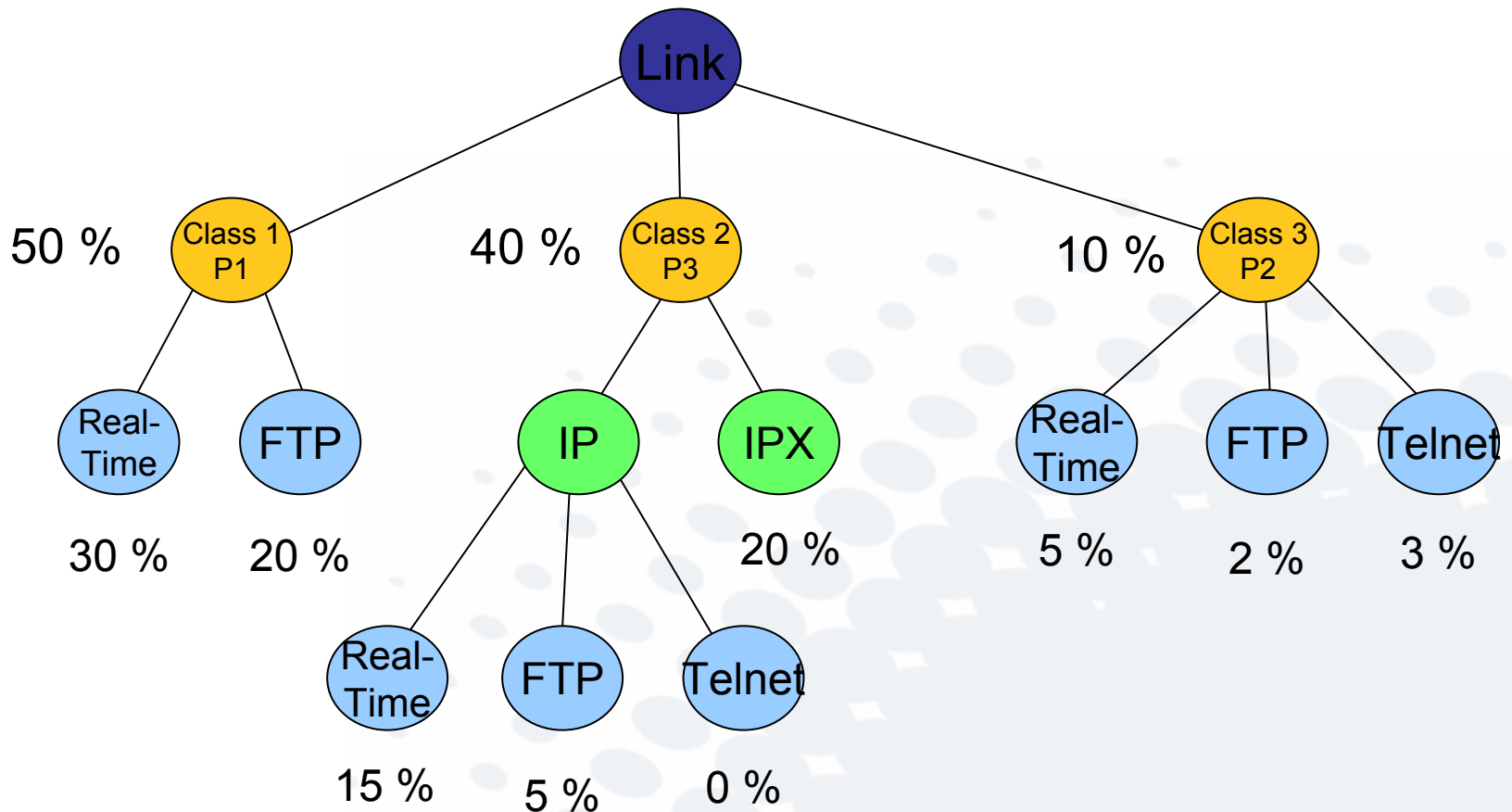
Clasificación en el gestor de ancho de banda



Antes de poder tratar de manera diferenciada un determinado flujo IP es necesario identificarlo y aislarlo. Un Gestor de Ancho de Banda es capaz de clasificar flujos a partir de alguna de las siguientes informaciones o combinación de ellas:

- URL y sub-URL
- Número de puerto TCP/UDP
- Dirección IP
- Contenido de campo TOS: CISCO *IP Precedence* y IETF *DiffServ*
- Campo 802.1p
- MAC address

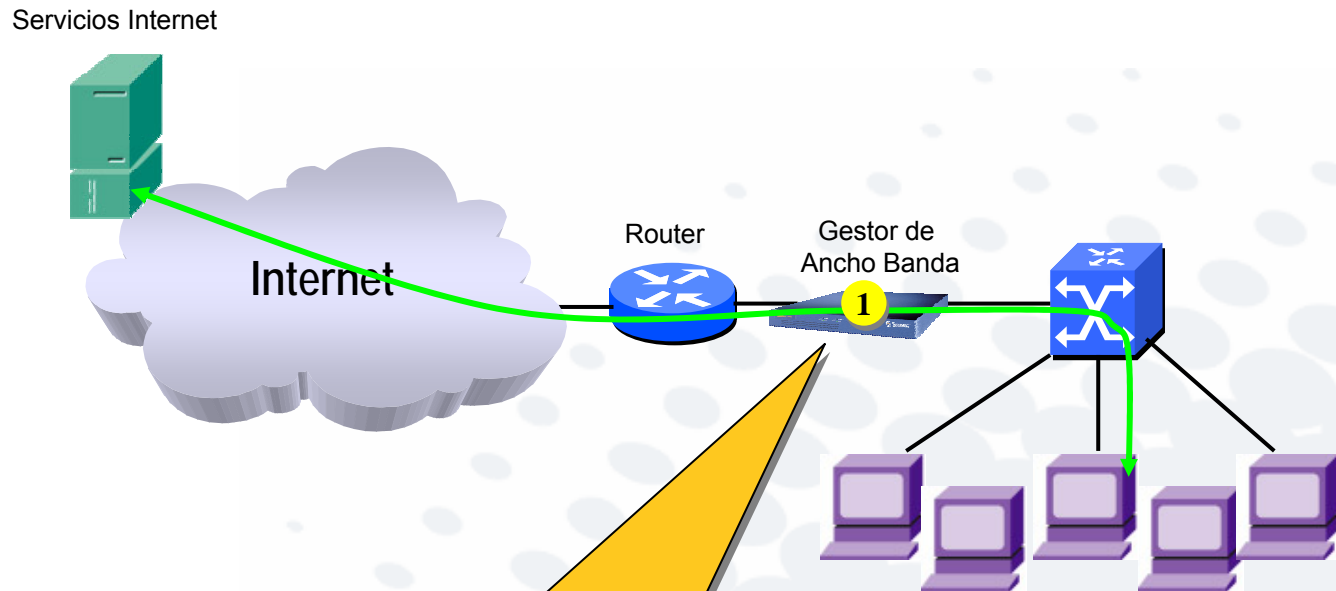
Árboles de clasificación jerárquica



- Establecimiento de una o más clases de servicio
- Asignación de ancho de banda para cada clase de servicio
- Creación perfiles que definen a un tráfico dentro de cada clase
- Asignación de ancho de banda a cada tráfico

Escenario acceso a Internet

Gestión basada en control sesión TCP



Políticas de priorización

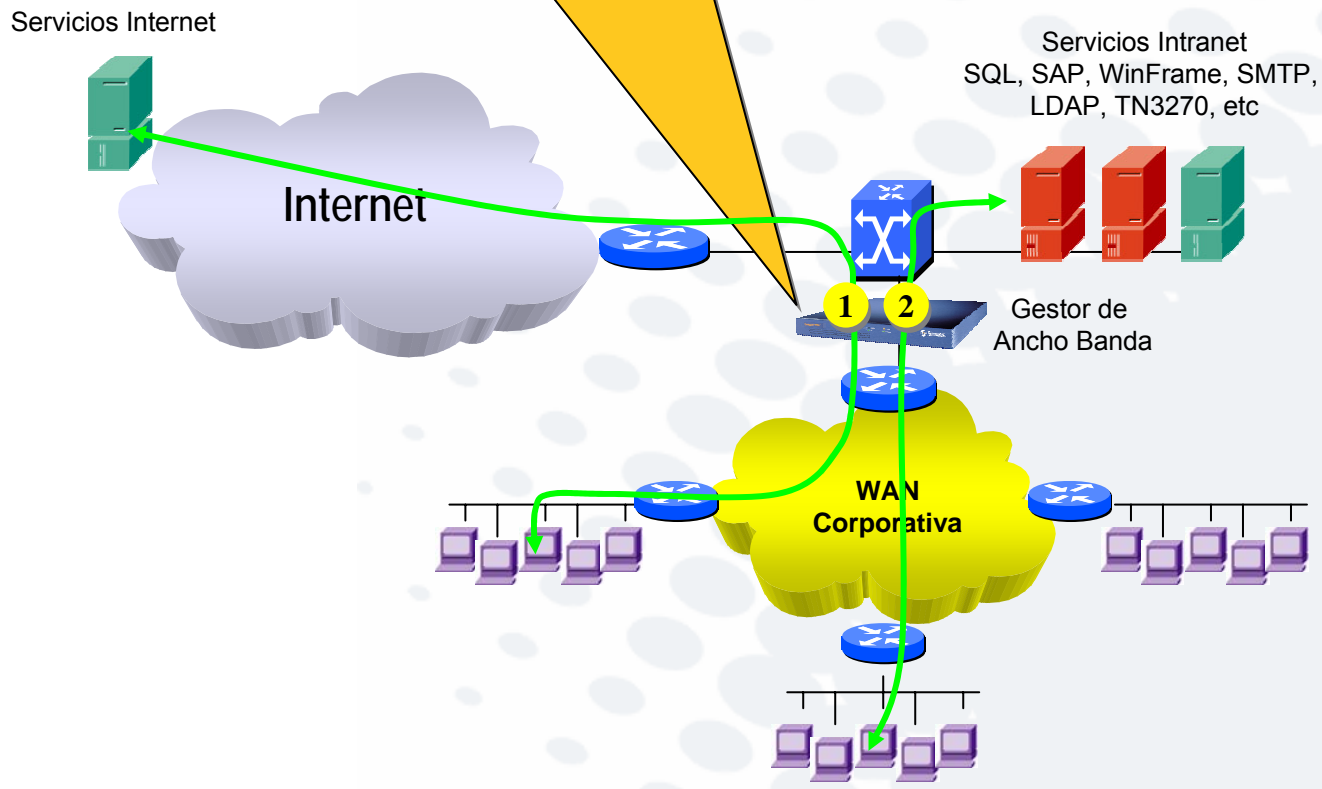
- Dirección IP individual o grupo de direcciones IP
- Servicios HTTP, SMTP, IRC, FTP, H.323, P2P
- Franja Horaria

1 – Gestión de la sesión TCP entre usuario y servidor Internet

Escenario Internet e Intranet

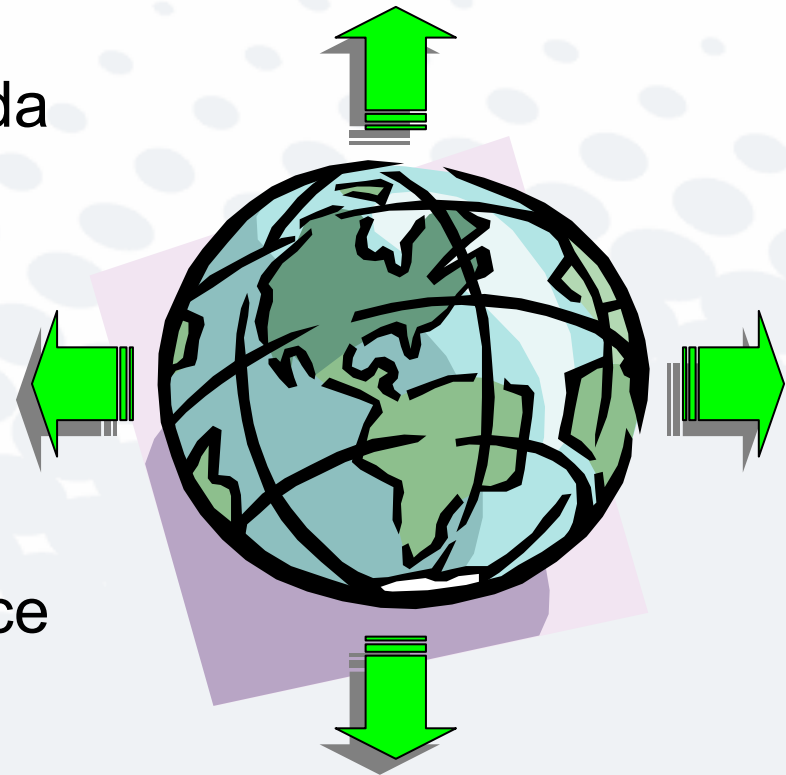
Gestión basada en control sesión TCP

- Políticas de priorización**
- Dirección IP individual o grupo de direcciones IP
 - Servicios Corporativos de la Intranet
 - Servicios Internet HTTP, SMTP, IRC, FTP, H.323
 - Franja Horaria y Fechas
- 1 – Control sesión TCP entre usuario y servidor Internet
2 – Control sesión TCP entre usuario y servidor Intranet

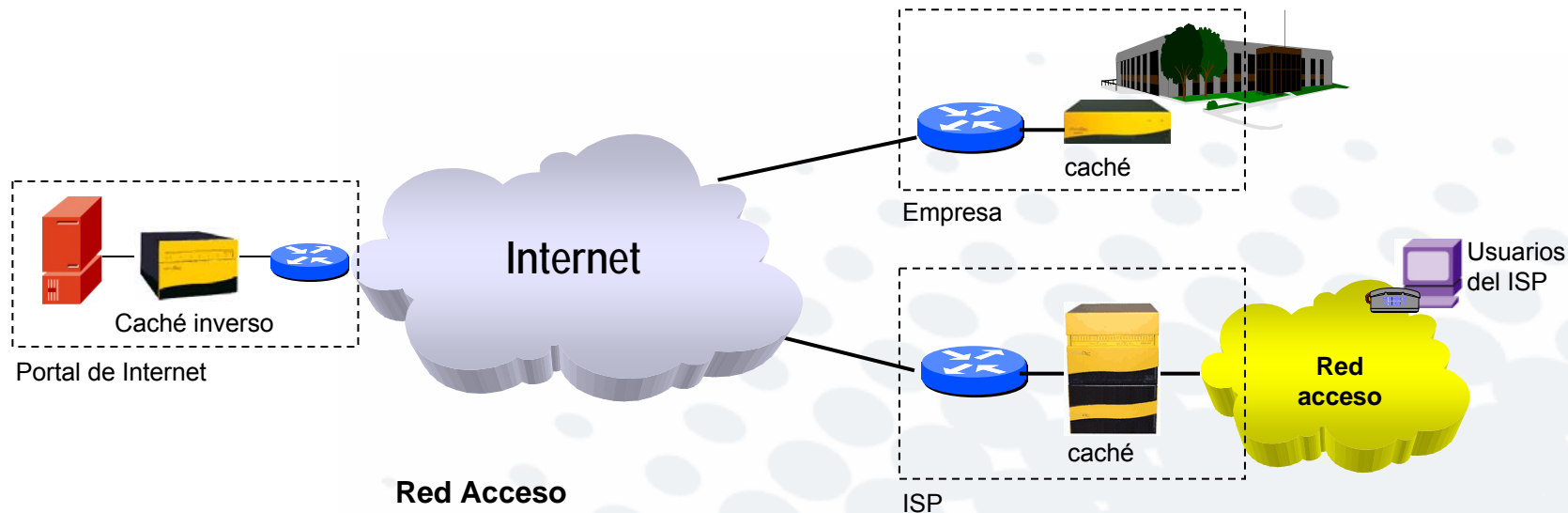


Agenda

- Gestores de Ancho de Banda
- **Caché de contenidos**
- Balanceadores de Carga
- Aceleradores de e-commerce



Caché de contenidos



- Un Proxy-caché almacena los objetos y páginas Web más visitadas por los usuarios de Internet
- Se sitúa entre usuarios e Internet, pudiéndose aplicar en escenarios de ISP y Empresa
- También es posible activar capacidades de caché entre servicios Web y usuarios (proxy inverso)

Características a valorar de una caché

- Caché basado en Software o *Appliance*
- Sistema de almacenamiento
- Control de vigencia y actualización del contenido almacenado
- Soporte HTTP 1.1
- Rendimiento: ancho de banda, objetos/segundo y número de sesiones TCP/IP simultaneas
- Modos de funcionamiento:
 - Transparente
 - Proxy
 - Integración con Routers y Switches multinivel
 - CDN
- Interoperabilidad con inspectores de contenido

Sobre el sistema de almacenamiento

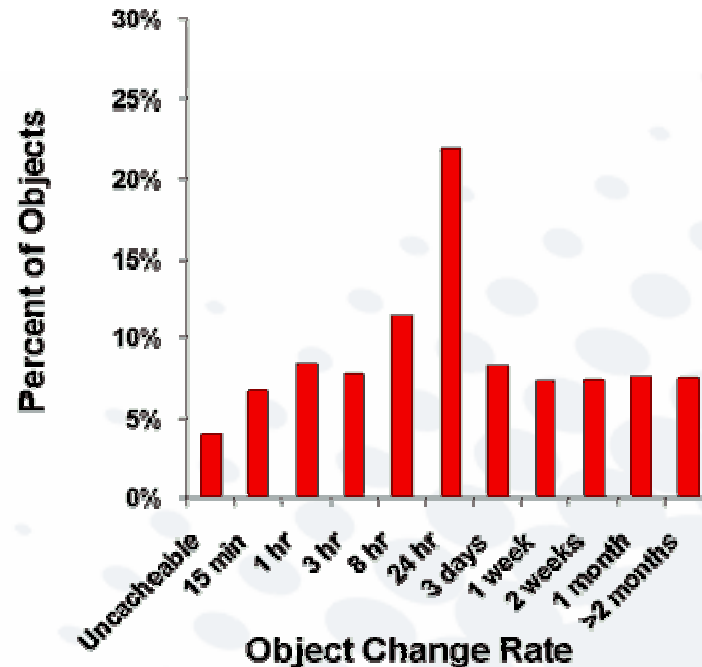
- **Prestaciones hardware**
 - Discos hot-swap
 - Tecnología SCSI Ultra-3
 - Detección pro-activa de fallos
 - Conexión con cabinas externas
- **Sistema de almacenamiento orientado a objetos HTML**
 - Diseñado para operar con miles de ficheros pequeños en constante cambio y/o actualización
 - Sobre cada objeto que se almacena en disco es necesario incluir nuevos atributos, como su fecha de caducidad y renovación
 - No existe una FAT convencional. Los objetos son accesibles a través de una tabla de vectores que se almacena en RAM
 - El sistema de “bloques y clusters” clásico impone un tamaño mínimo de unidad de almacenamiento que, en muchas ocasiones, es mayor que un fichero html u objeto Web (desperdicio de capacidad)
 - Para conseguir mejor tasa de acierto las caches tienden a llenar sus discos con todos los URLs visitados. Esta circunstancia obliga a que el sistema de almacenamiento gestione eficazmente la situación
 - Acceso balanceado: en cada disco se almacenan objetos que pertenecen a un misma página. Cuando un disco falla sólo los objetos afectados son actualizados, en lugar de toda la página

Control vigencia del contenido

Sobre cualquier solución de caching es necesario valorar los mecanismos disponibles para asegurar que los objetos almacenados en la caché se encuentran actualizados. Hay diferentes técnicas:

- **Inspección de modificación de contenidos.** Ante una consulta, la caché siempre comprueba la vigencia del contenido contra la Web original. *HTTP 1.0 Conditional Request "if-modified-since"*
- **Meta Tags incluidos en el código HTML:** *"pragma: no cache"*
- **Cabeceras HTTP 1.1 enviadas por el servidor web:** *max-age, private, no-store, etc..*
- **Reglas Heurísticas de cada fabricante y producto.** Generalmente establecen la pauta de actualización/refresco a partir de la popularidad del objeto.

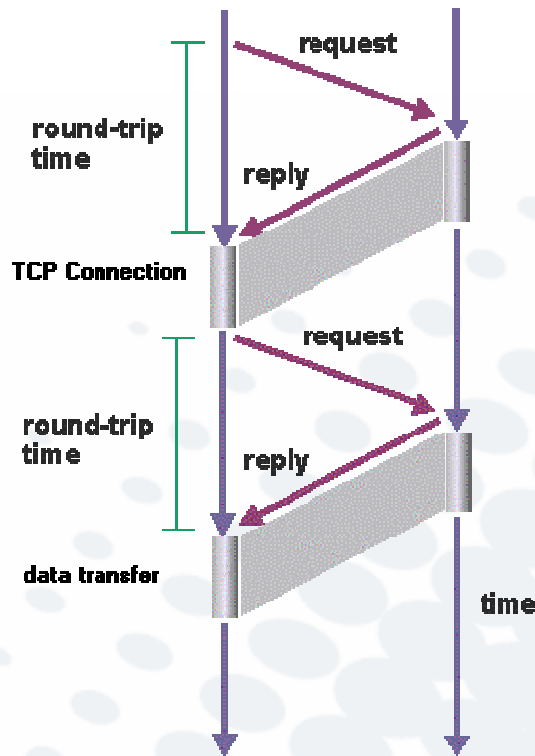
Ejemplo de un algoritmo de refresco basado en reglas heurísticas



Sobre cada página almacenada en la caché, ésta establece el patrón de refresco apropiado a partir de las siguientes variables:

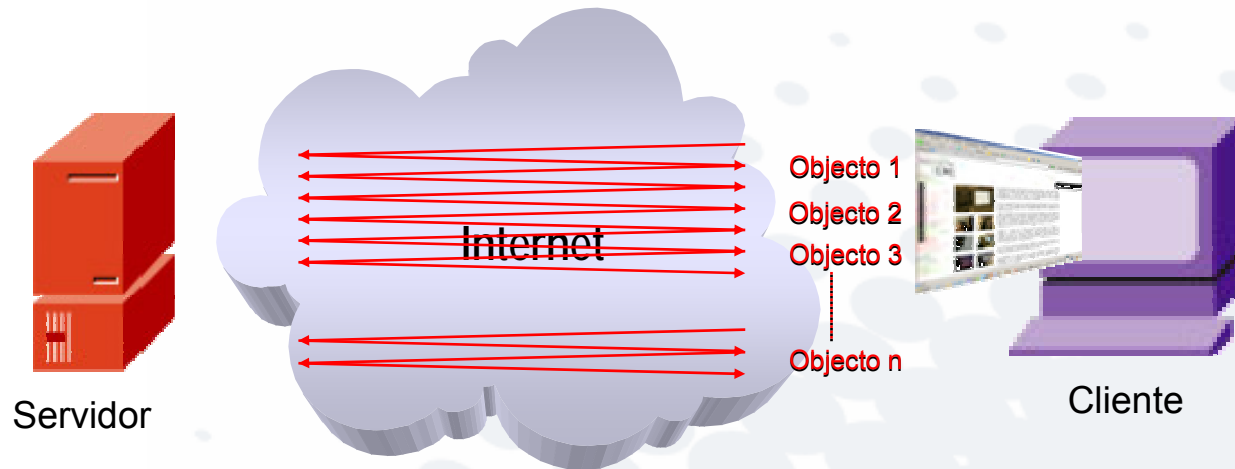
- **Pauta de cambio:** ¿Cuándo suele cambiar la página original?
- **Modelo de uso:** ¿A qué horas suelen visitar esta web los usuarios?

Funcionamiento de HTTP v.1.0



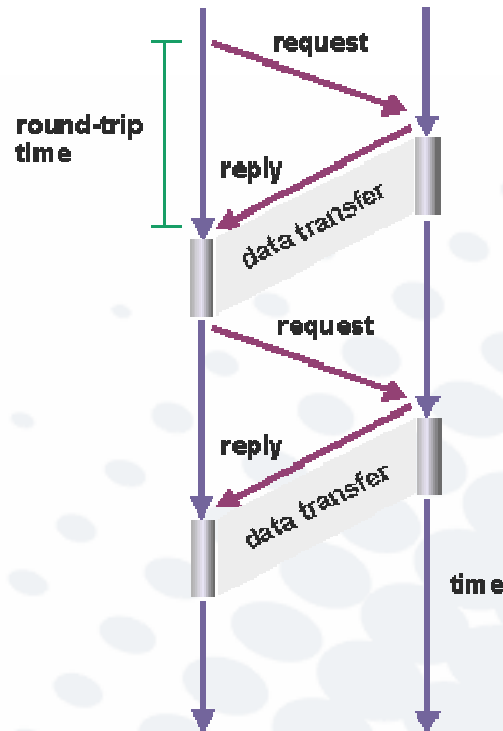
Una sesión TCP por objeto. En HTTP 1.0 los browser sólo son capaces de solicitar cuatro objetos simultáneamente (limitación autoimpuesta)

HTTP v.1.0 Object Round Trip



- La representación de cada objeto de una página supone un trayecto de ida y otro de retorno desde el servidor.
- El proceso *TCP Slow Start* asociado a cada sesión TCP impide alcanza la “velocidad máxima”, ya que cada objeto enviado desde el servidor no es lo suficientemente grande

Funcionamiento de HTTP v.1.1 Persistent



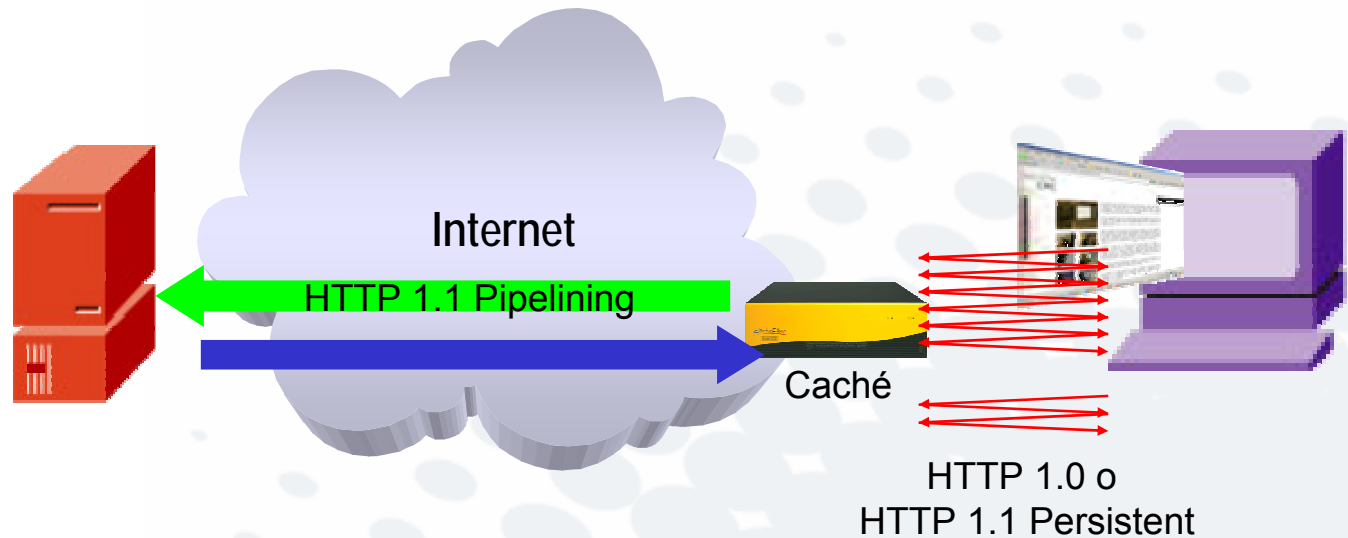
- Sobre una única sesión TCP se descargan todos los objetos de manera secuencial, es decir, cuando se recibe un objeto se solicita el siguiente.
- RFC 2068 limita al navegador o browser establecer más de dos conexiones HTTP contra un mismo servidor.
- El funcionamiento habitual de los browser es HTTP 1.1 *Persistent*

Funcionamiento de HTTP v.1.1 Pipelining



- Sobre una única sesión TCP se descargan todos los objetos de manera simultánea sin esperar la llegada del objeto previo.
- Netscape y IE no soportan HTTP 1.1 Pipelining
- Cualquier servidor compatible HTTP 1.1 debe soportar Pipelining

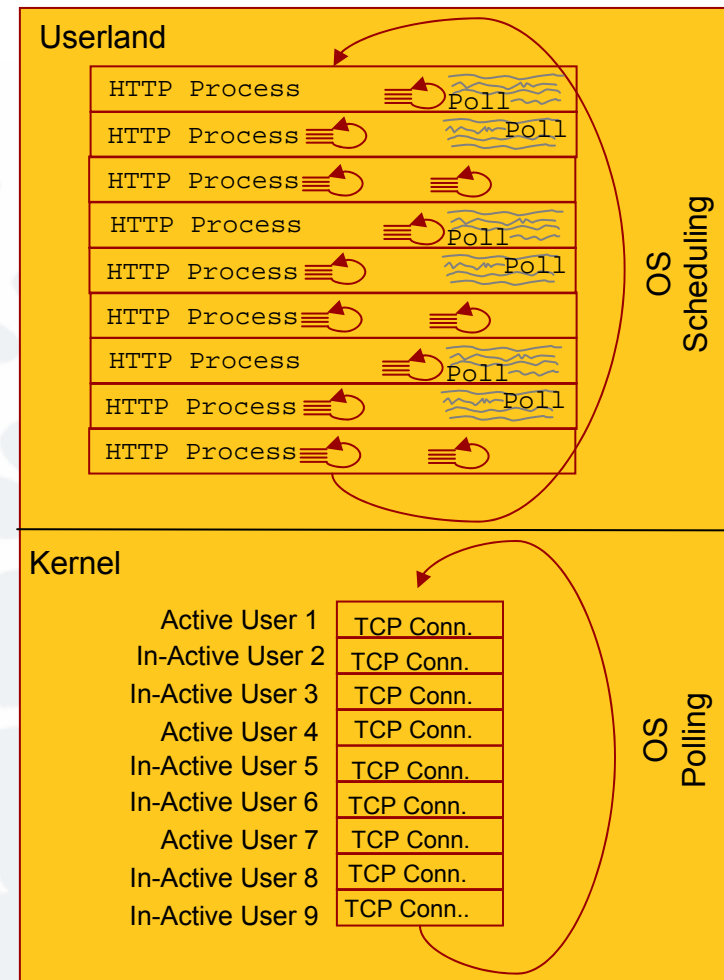
HTTP 1.1 Pipelining en la caché



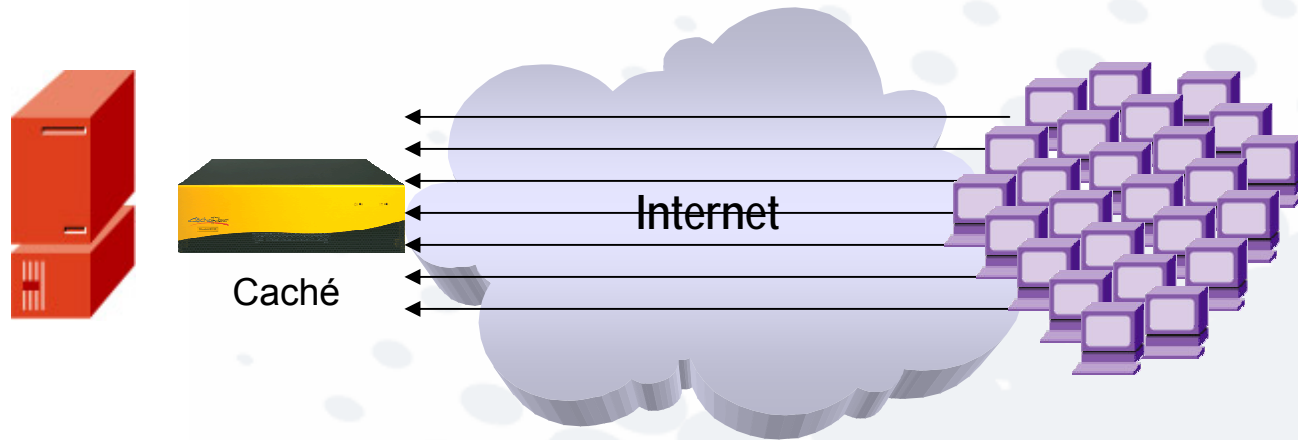
Acelera y optimiza la descarga de páginas mediante la solicitud en paralelo de todos los objetos contenidos en una página (*HTTP 1.1 Pipelining*)

HTTP v.1.1 en el lado del servidor

- HTTP 1.1 es más eficaz para el backbone y para el cliente:
 - Menos conexiones que abrir y cerrar: Una única sesión TCP sobre la que se piden todos los objetos
- HTTP 1.1 sobrecarga los servidores debido a la gestión de todas las conexiones TCP: deficiencias de polling y scheduling
 - Los clientes lentos usan la mayoría de los recursos del servidor
 - HTTP 1.1 consume los recursos de los servidores

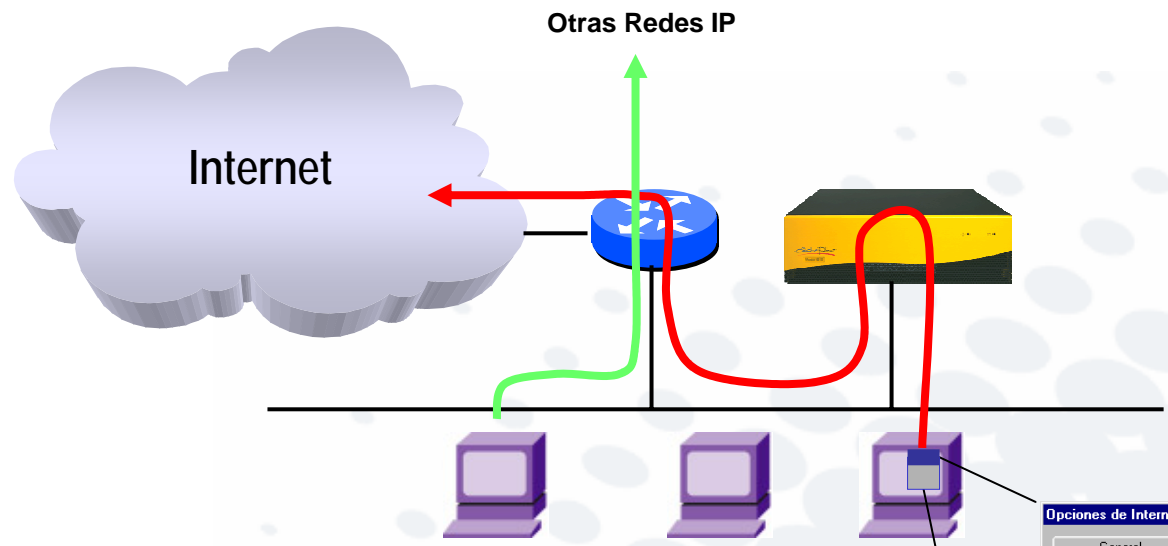


Caché inverso



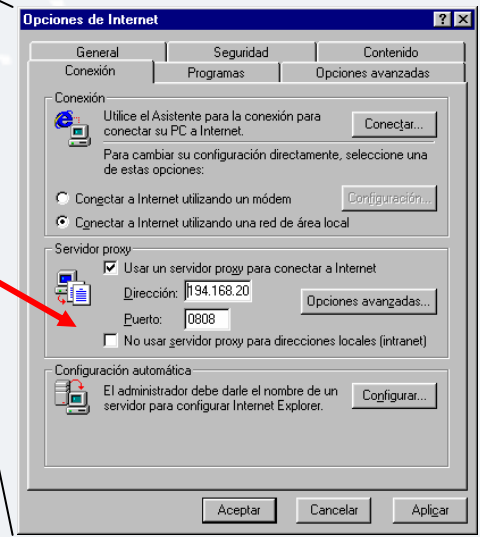
Finaliza y gestiona miles de sesiones HTTP de los usuarios

Funcionamiento como Proxy HTTP

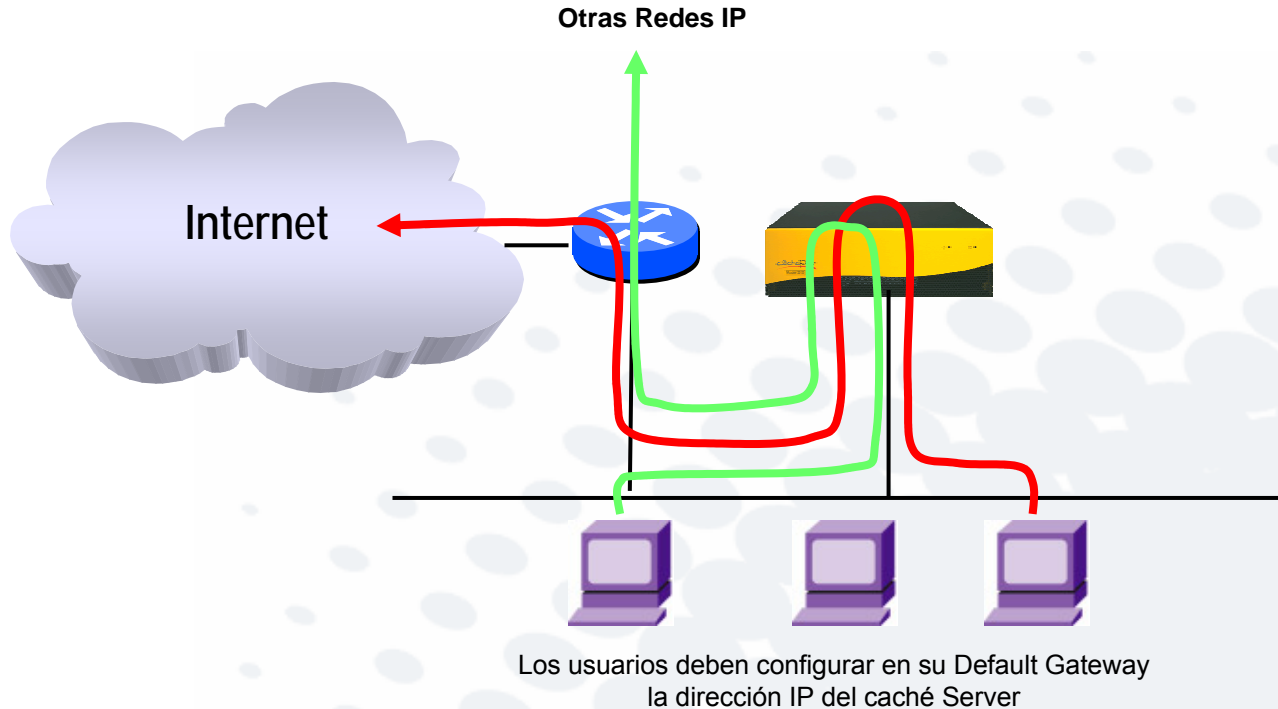


Los usuarios deben configurar sus Browsers para salir a Internet a través de la caché

- Sólo el tráfico HTTP pasa por el Proxy-caché
- Posibilidad de aplicar políticas de control
- Es necesario configurar todos los clientes
- Si la caché falla se pierde la conexión a Internet

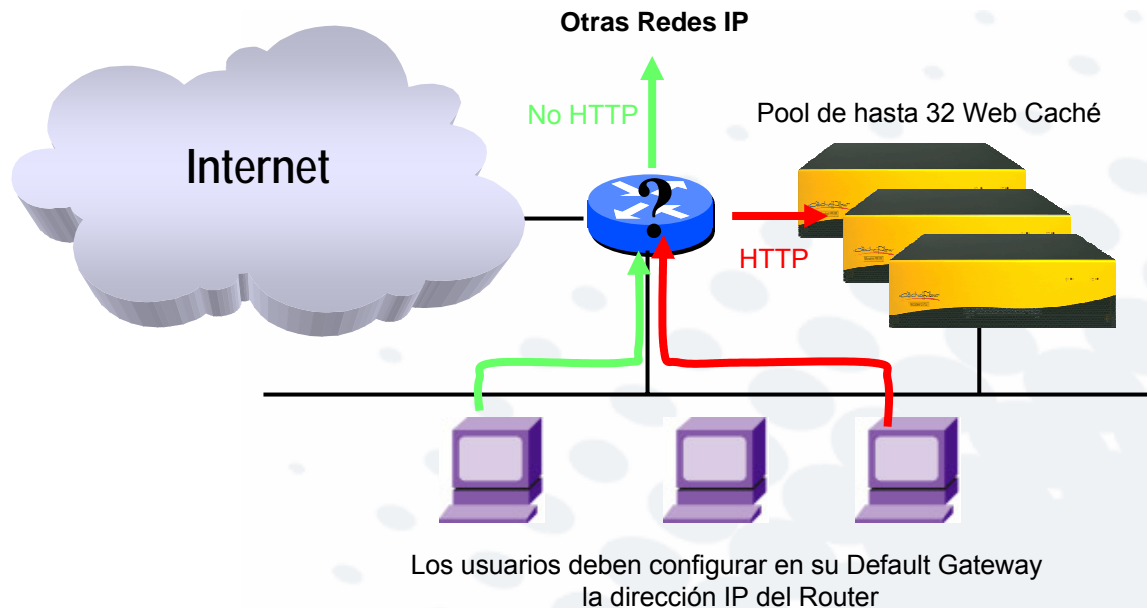


Funcionamiento Transparente



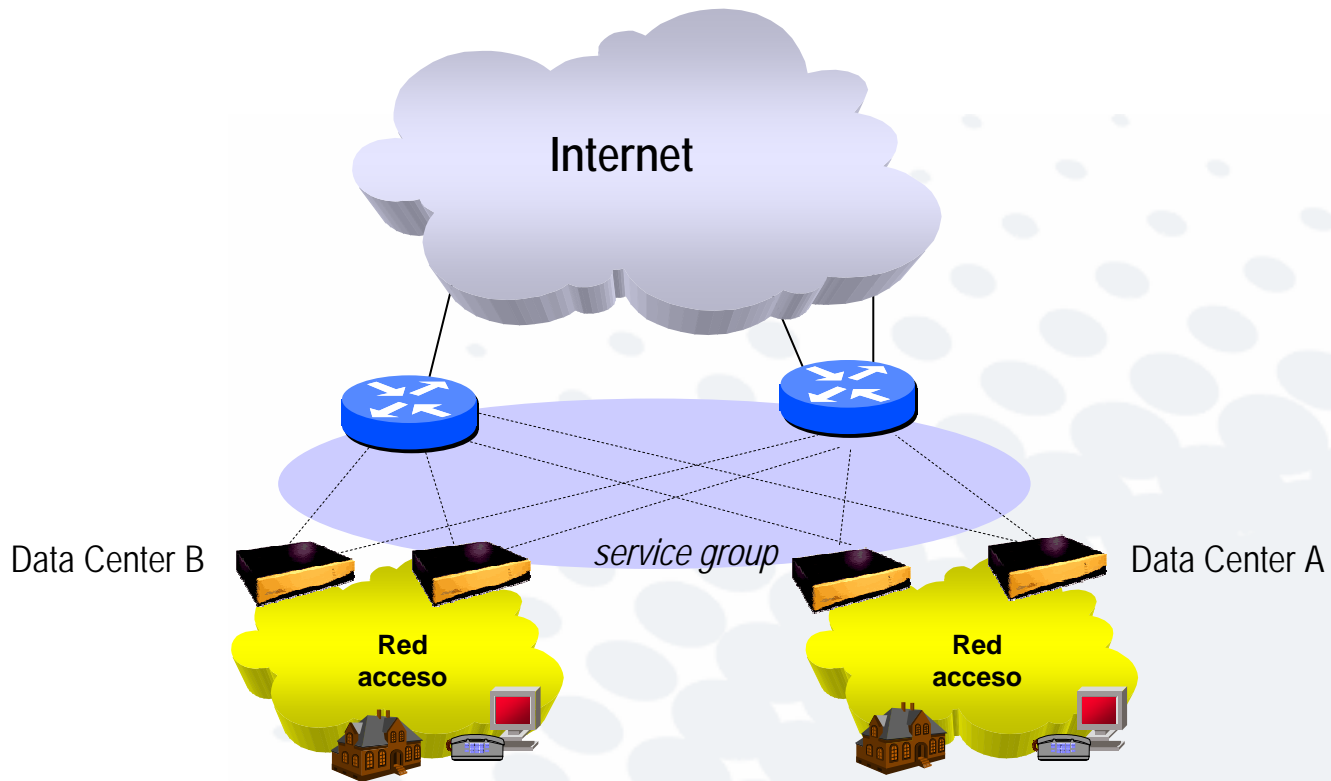
- Resulta transparente a los usuarios
- Todo el tráfico pasa por el Proxy
- El proxy se comporta como un Router
- Representa un punto singular de fallo y un cuello de botella

Funcionamiento transparente con WCCP v.1



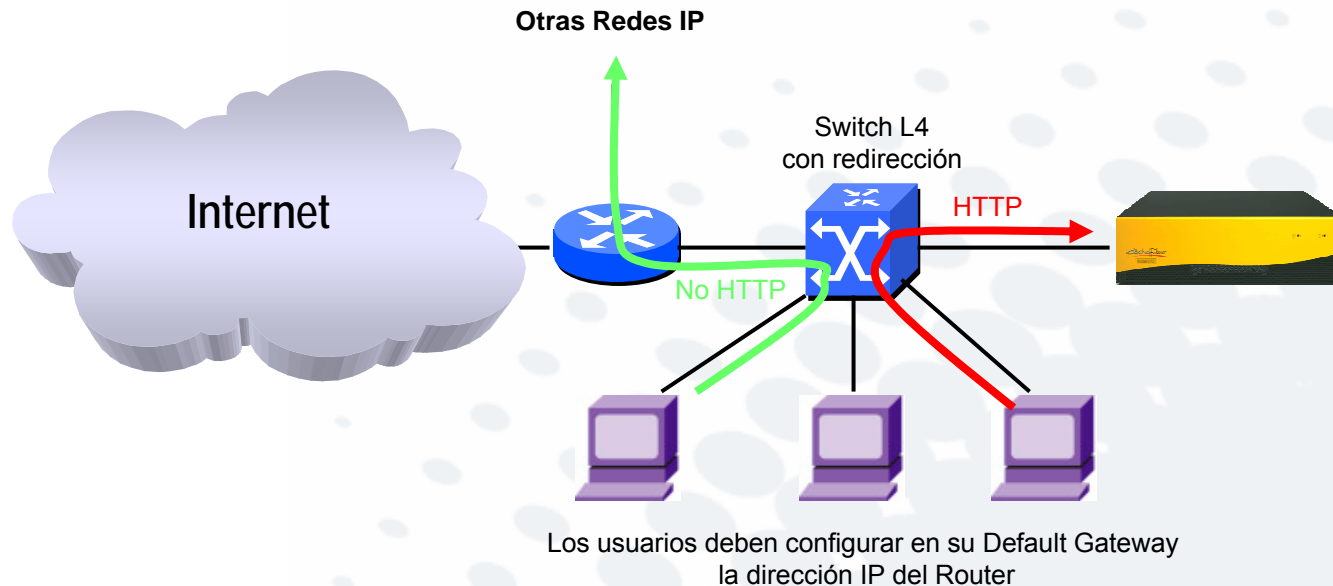
- Empleo de WCCP (Web Cache Control Protocol)
- Resulta transparente a los usuarios
- El tráfico HTTP es redireccionado por el Router hacia el Web Caché
- Gracias a WCCP es posible establecer un pool de hasta 32 caches
- Inserción y eliminación transparente de caches en el pool

Clustering de caches con WCCP v.2



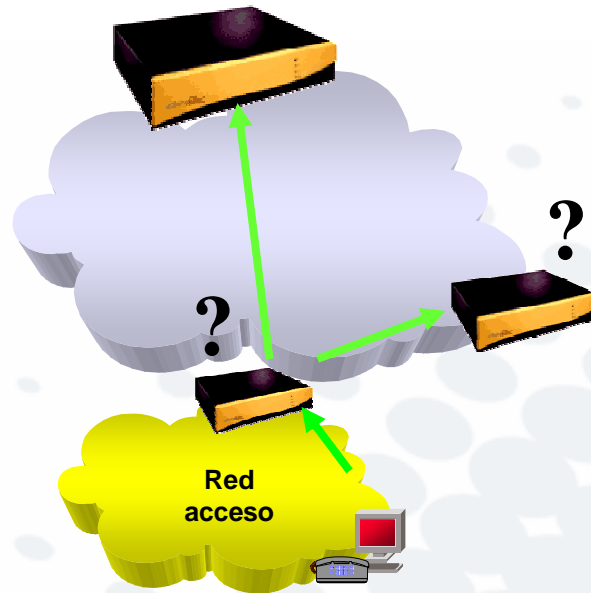
- Cluster de 32 routers que pueden redireccionar tráfico hacia 32 cachés
- Seguridad MD5 para autenticación de los diversos componentes de un “service group”

Funcionamiento transparente con Switch L4



- Transparente a los usuarios
- El tráfico HTTP es redireccionado por el switch hacia la caché
- Importante: Valorar si existe comprobación del estado de la caché desde switch

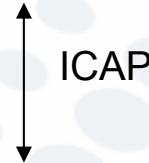
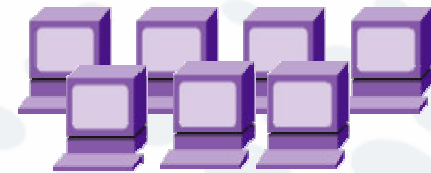
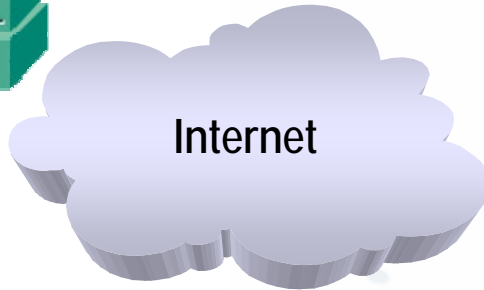
ICP (Internet Caching Protocol)



- Establecimiento de un grupo de caches
- Si una caché dentro del grupo no dispone de un objeto, puede solicitarlo a otra caché previa determinación de que vecino dispone del objeto. También puede conocer quien ofrece la mejor respuesta
- Tras una petición ICP el caché solicitante debe esperar 2 segundos para recibir las respuestas de los restantes caches

ICAP (Internet Content Adaptation Protocol)

Servidor origen

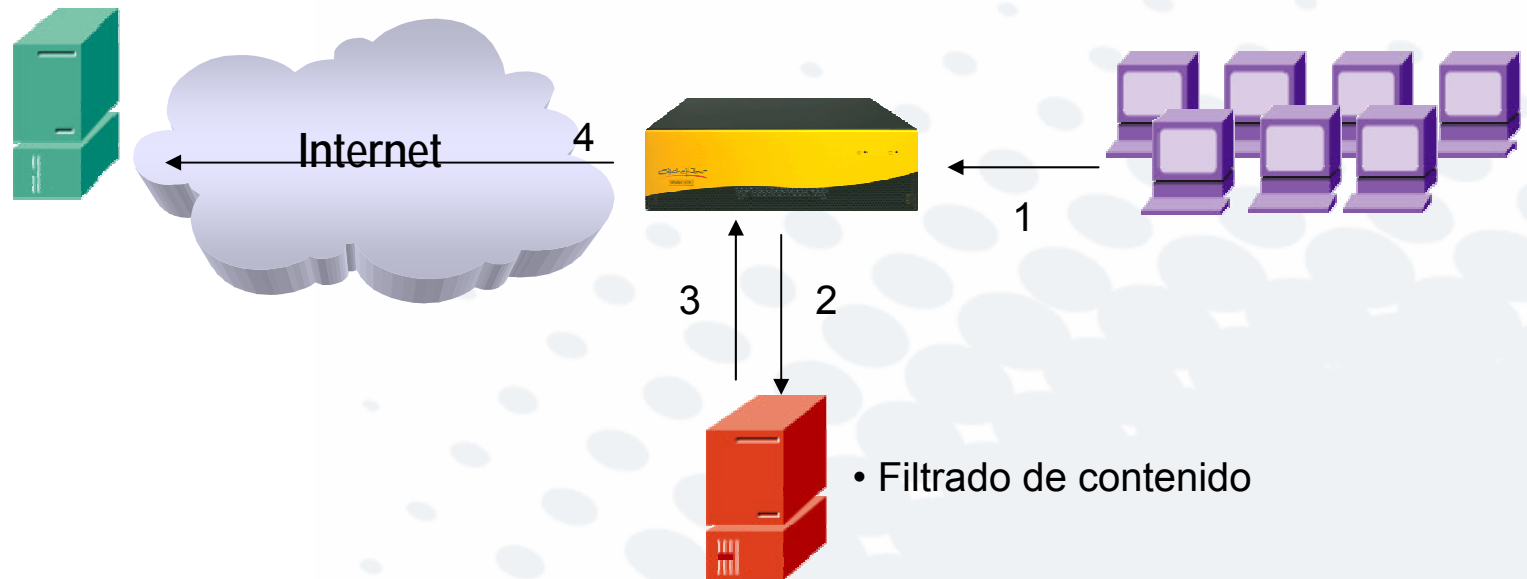


Servidores con soporte ICAP

Protocolo estandarizado a través del cual la caché traslada una petición de un usuario, o la respuesta del servidor origen, a un servidor para tratar de forma especializada el contenido

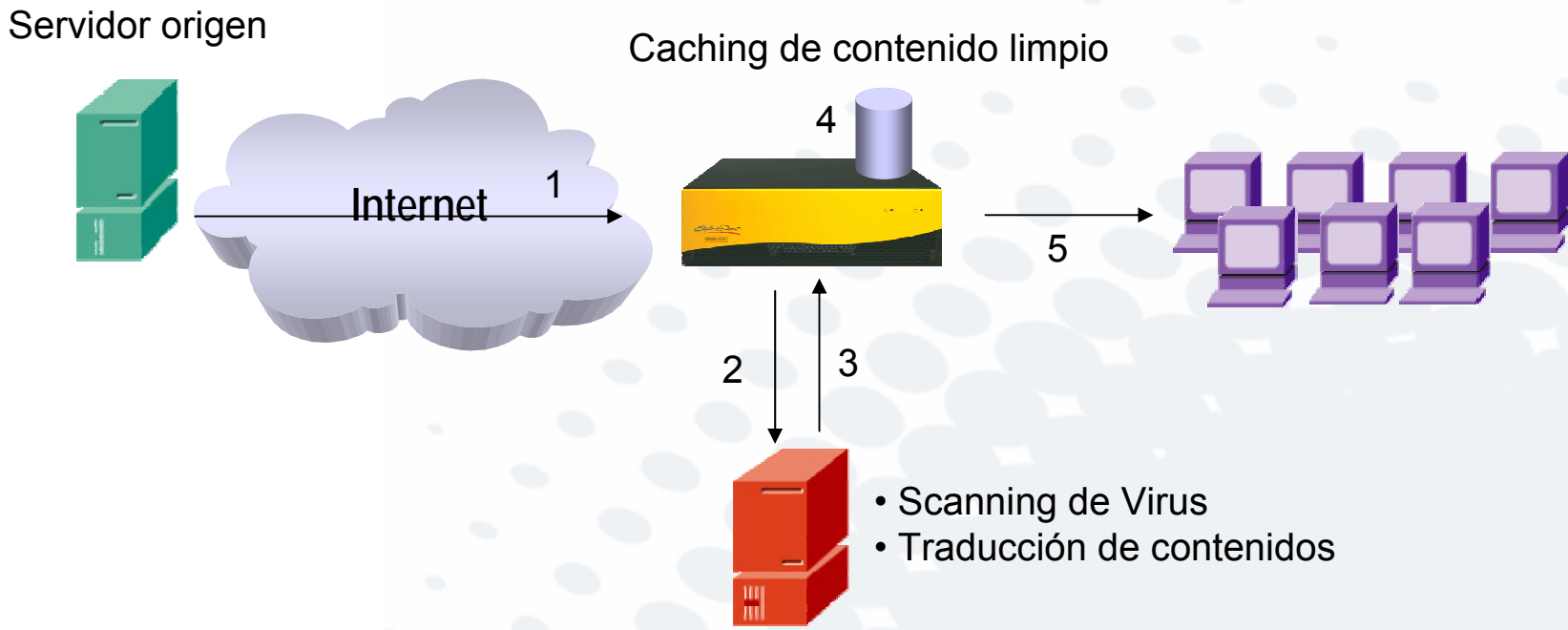
ICAP Request Modification

Servidor origen



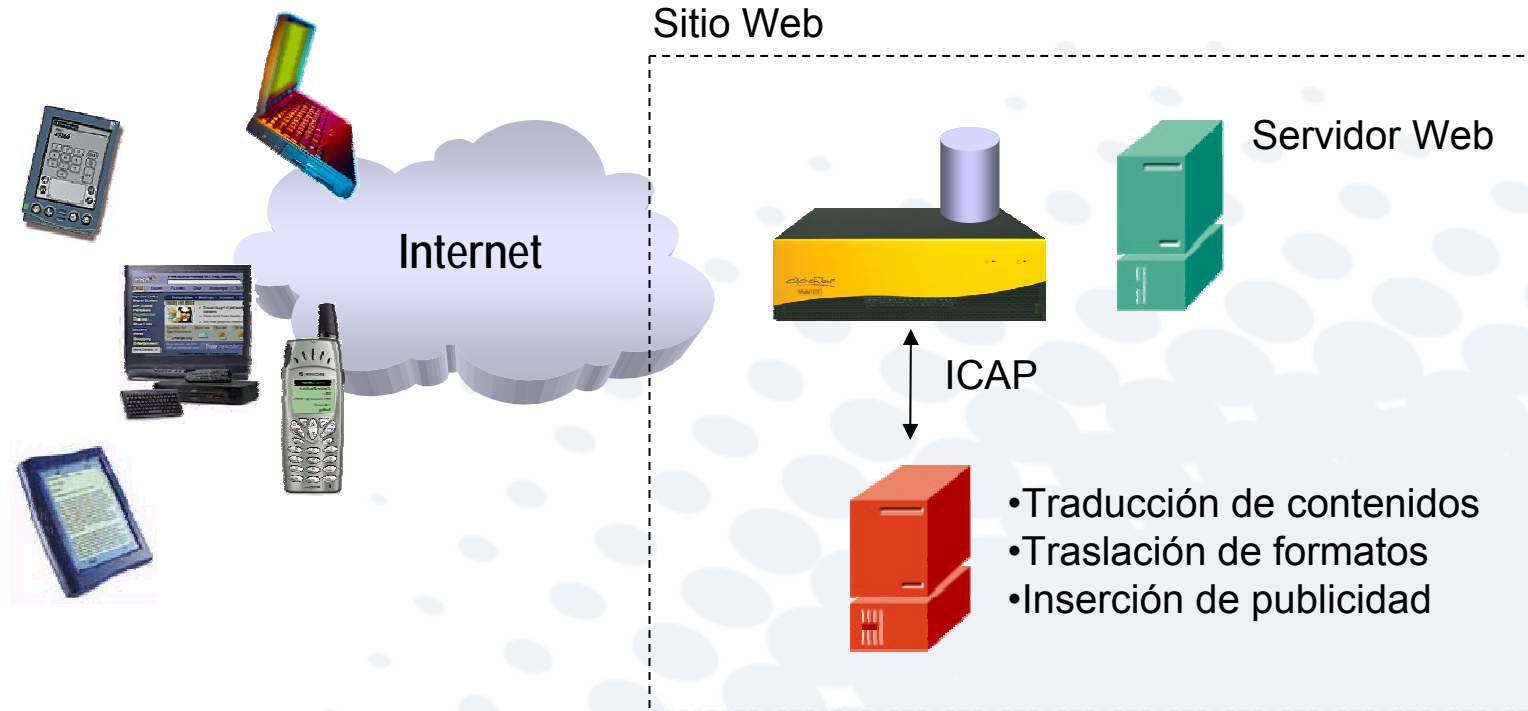
La caché traslada la petición de un usuario a un servidor de filtrado de contenidos donde se comprueba si está permitida la visita al contenido almacenado en el servidor/sitio de Internet

ICAP Response Modification



La caché traslada la respuesta de un servidor/sitio de Internet a un servidor ICAP especializado

ICAP Response Modification en Sitios Web



La caché traslada la respuesta de un servidor/sitio de Internet a un servidor ICAP especializado en el tratamiento del contenido

Caching y entrega de contenidos multimedia

- Real Server Proxy
- Microsoft Media Proxy
- Proxy RTSP nativo

- Formatos multimedia soportados:
 - Audio y Vídeo Real Networks
 - Audio y Vídeo Microsoft Windows Media
 - Audio y Vídeo Apple QuickTime
 - Vídeo MPEG-2
 - Audio MP3

- Entrega de Vídeo en Stream Splitting

- Entrega de contenidos pregrabados

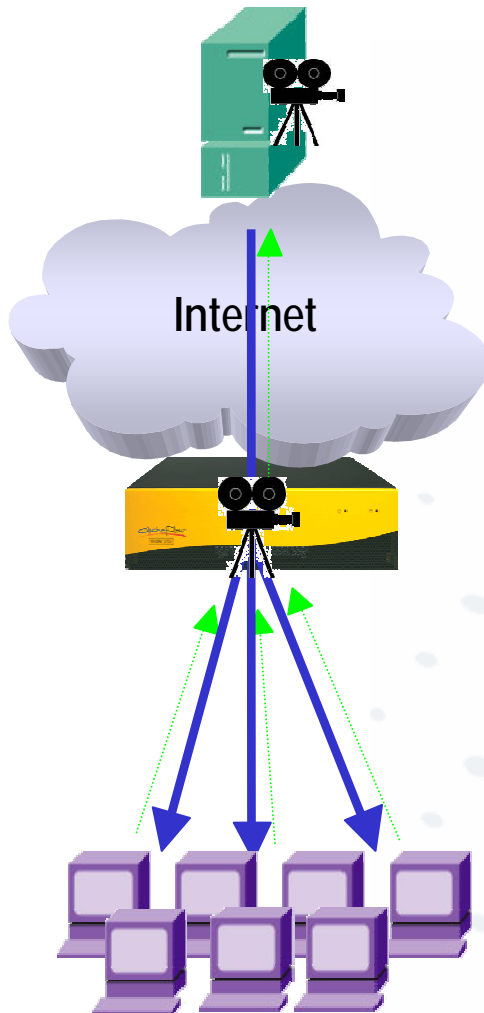
- Traslación Multicast-Unicast

- SureStream e Intelligent Stream

- Gestión de ancho de banda

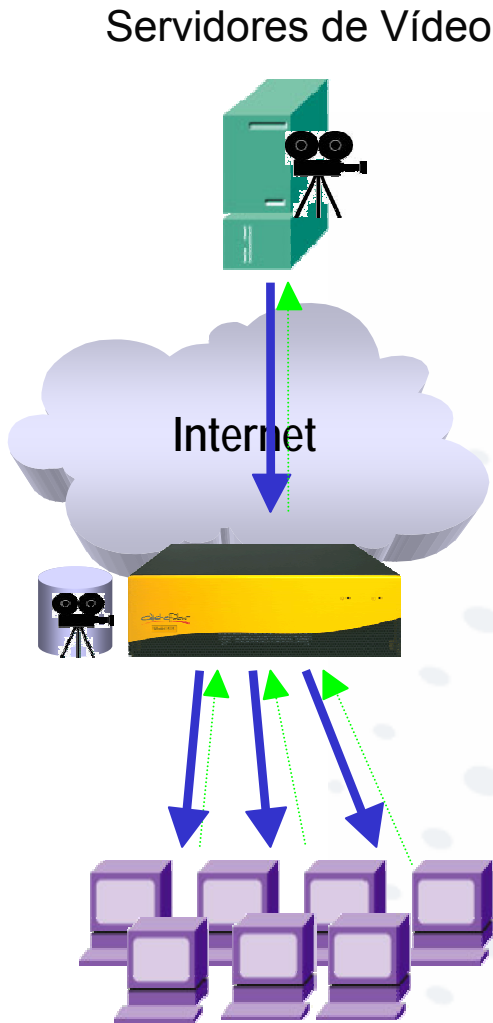
Stream Splitting

Servidores de Vídeo



- Un único flujo de vídeo recibido en tiempo real es reenviado a “n” clientes
- Traslación origen unicast hacia usuarios multicast
- Traslación origen multicast hacia usuarios unicast
- Ahorro de ancho de banda
- Gestión de ancho de banda dedicado a bajar el contenido desde la fuente y el ancho de banda destinado a cada usuario
- Control centralizado de los contenidos accesibles por los usuarios

Contenido pregrabado



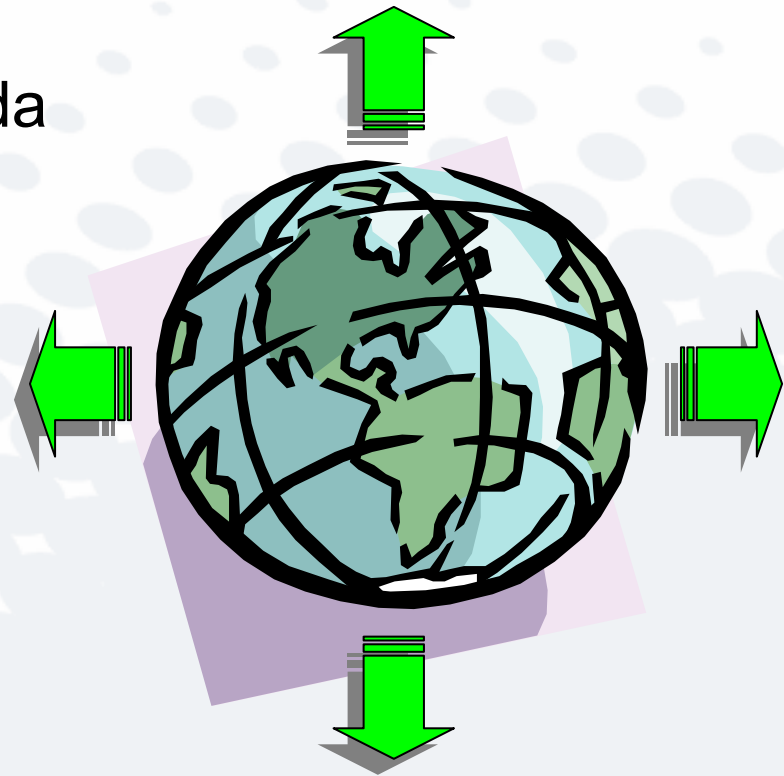
- Ahorro de ancho de banda
- Entrega rápida y óptima del contenido a los usuarios según sus características de acceso. SureStream e Intelligent Stream
- Si el servidor origen lo requiere se mantienen tantas sesiones de control como usuarios visualicen el contenido
- Control centralizado de los contenidos accesibles por los usuarios

Conclusiones sobre el Caching

- Ahorro de ancho de banda WAN en un 30-40%
- El ratio de aciertos de un caché de contenido se estima en un 50% aproximadamente
- Alrededor de un 40% de los contenidos Web no son cacheables. p.j. páginas generadas dinámicamente CGI, ASP, PHP, JSP, etc
- Existen páginas marcadas como no-cacheables
- Beneficios Empresa: Reducir costes de WAN, acelerar el acceso a Internet (contenido cerca de usuarios) y seguridad
- Beneficios ISP: Reducir costes WAN y mejorar el acceso usuarios
- Beneficios e-commerce: aumentar el rendimiento/respuesta de servicios del portal

Agenda

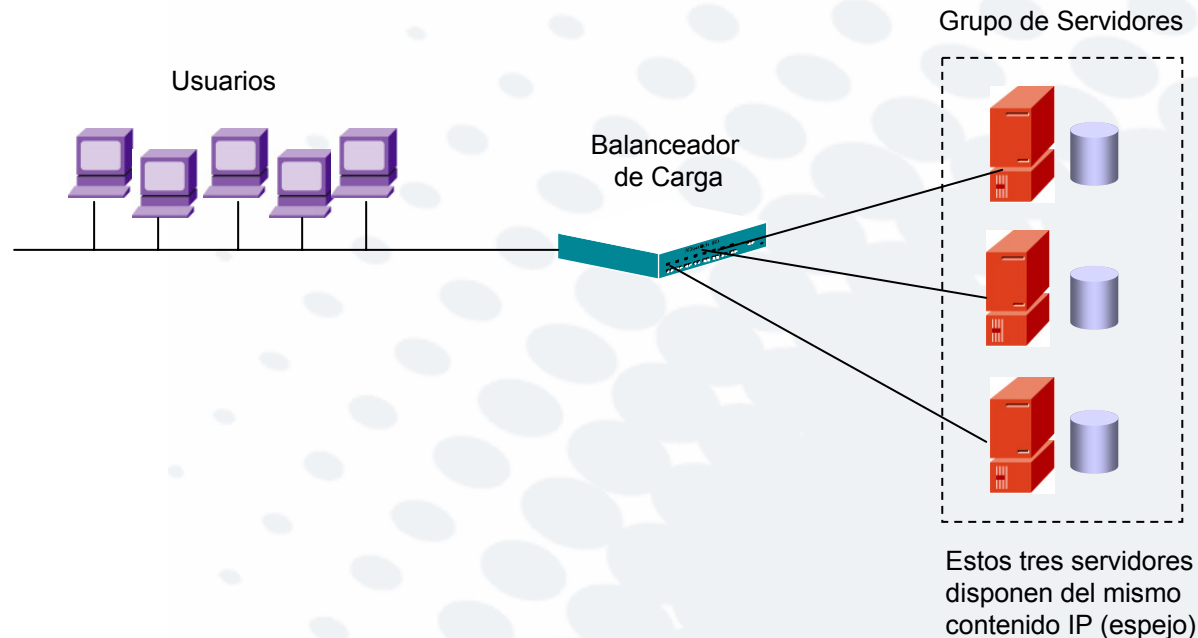
- Gestores de Ancho de Banda
- Caché de contenidos
- **Balanceadores de Carga**
- Aceleradores SSL



¿Qué es un balanceador de carga?

Dispositivo que se sitúa entre los servicios IP y los usuarios, distribuyendo las sesiones de estos hacia un grupo de servidores que ejecutan una aplicación común

La capacidad del grupo crece con el número de servidores, no es necesario que el hardware y sistema operativo del grupo sea homogéneo



El balanceador de carga identifica el comienzo de una sesión y le asigna a uno de los servidores del grupo hasta que finalice la citada sesión

Dónde, cuándo y cómo

¿Dónde?

- Cualquier empresa donde la disponibilidad y rendimiento del servicio IP sea vital: e-commerce, redes CDN, portales y negocios que se apoyen en Internet

¿Cuándo?

- Asegurar la fiabilidad de un servicio IP
- Incrementar rendimiento, acelerar y personalizar la entrega de contenidos
- Facilitar la escalabilidad, flexibilidad y el mantenimiento del servicio

¿Cómo?

- Balanceo local hacia un grupo de servidores. *Content Switching*
- Balanceo global hacia múltiples servidores distribuidos en diferentes sedes. *Content Routing*

Métodos de *Load Balancing*

- **Round Robin DNS**

- Se vincula un nombre a varias direcciones IP. El DNS resuelve cíclicamente la dirección IP, asignando una sesión a cada servidor del grupo
- ↓ El DNS no conoce el estado del servidor
- ↓ Limitaciones en el número de direcciones IP asignadas a un nombre
- ↓ Tiempo de permanencia del vínculo nombre-dirección IP en las caches DNS

- **Asignación dinámica con monitorización pasiva**

- Se distribuyen las sesiones en función del estado de cada servidor: nº de sesiones abiertas en cada servidor, servicios levantados, medición del tiempo de respuesta y el perfil hardware de cada máquina entre otros parámetros

- **Asignación dinámica con monitorización activa**

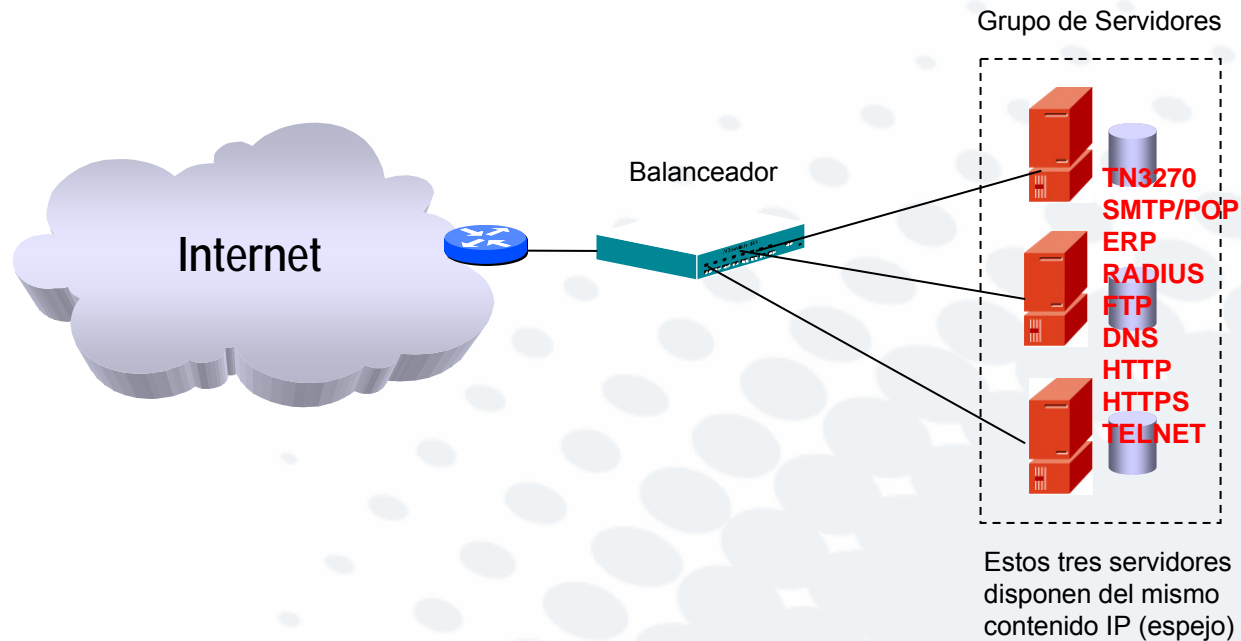
- Igual que la monitorización pasiva pero añadiendo agentes software en los servidores para informar al balanceador del estado real del procesador, memoria, servicios, etc.

Técnicas de monitorización

Los algoritmos de monitorización permiten al balanceador conocer el estado de los servidores que constituyen el grupo. Además de comprobar esta disponibilidad, la información obtenida se emplea para tomar decisiones en asignación de sesiones:

- Control de respuesta mediante pings
- Apertura y cierre de servicios/puertos TCP y UDP
- Registro del número de sesiones abiertas en cada servidor
- Medición del tiempo de respuesta de algunos servicios IP mediante peticiones preconfiguradas (HTTP, Radius, IMAP, SSL, etc)
- Monitorización del estado de la memoria y procesador de cada servidor a través de agentes instalados en los propios servidores

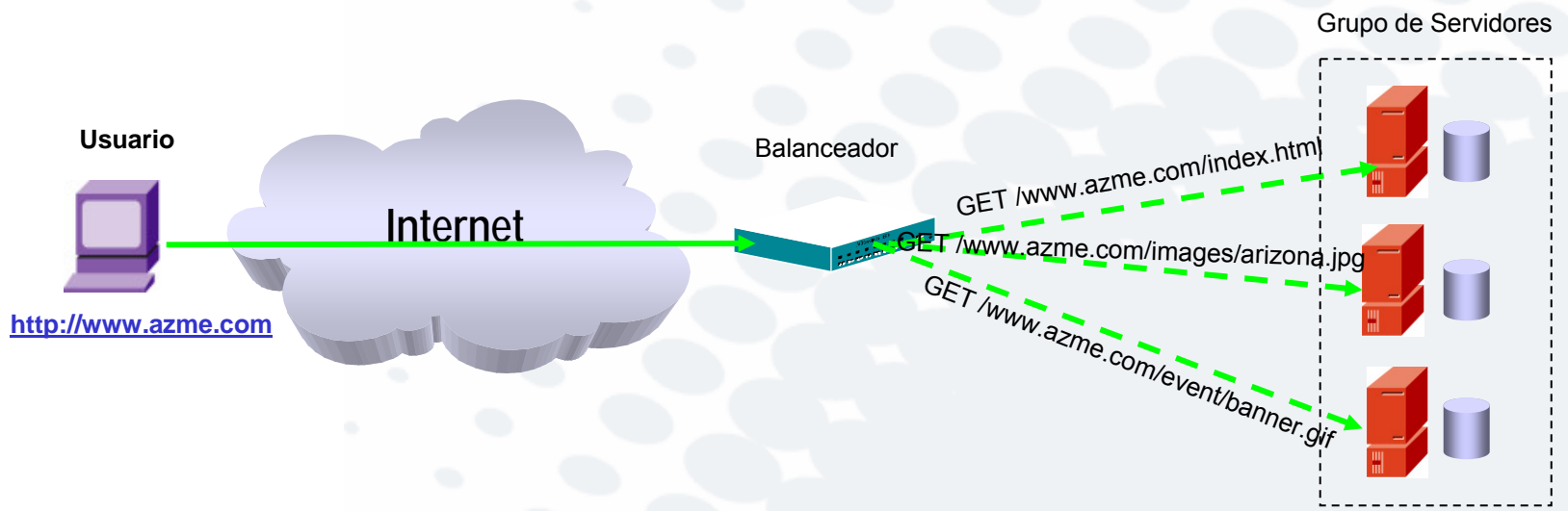
Servicios “balanceables”



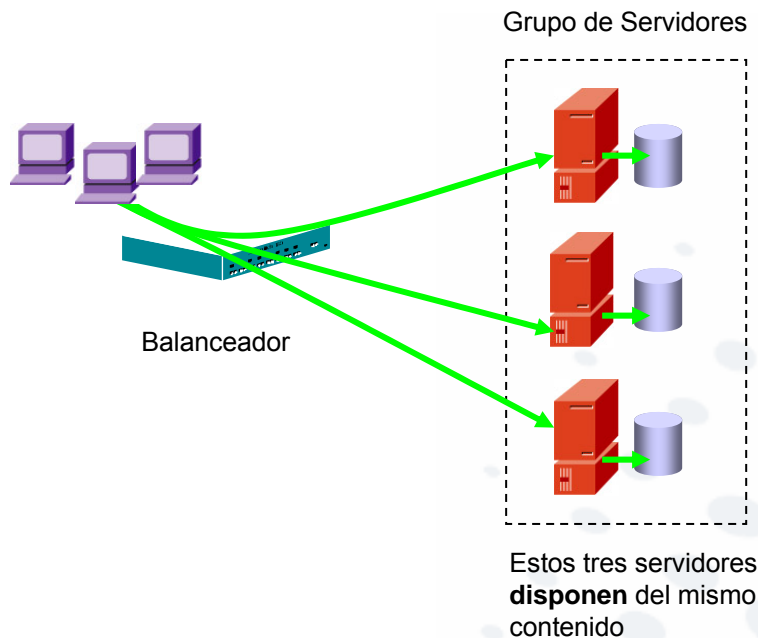
Mientras que en los servidores exista homogeneidad en el contenido es posible hacer *balancing* de cualquier servicio TCP y UDP, aunque en la práctica el uso mayoritario se orienta a servidores web, ssl, caches y firewalls

URL Balancing y Content Switching

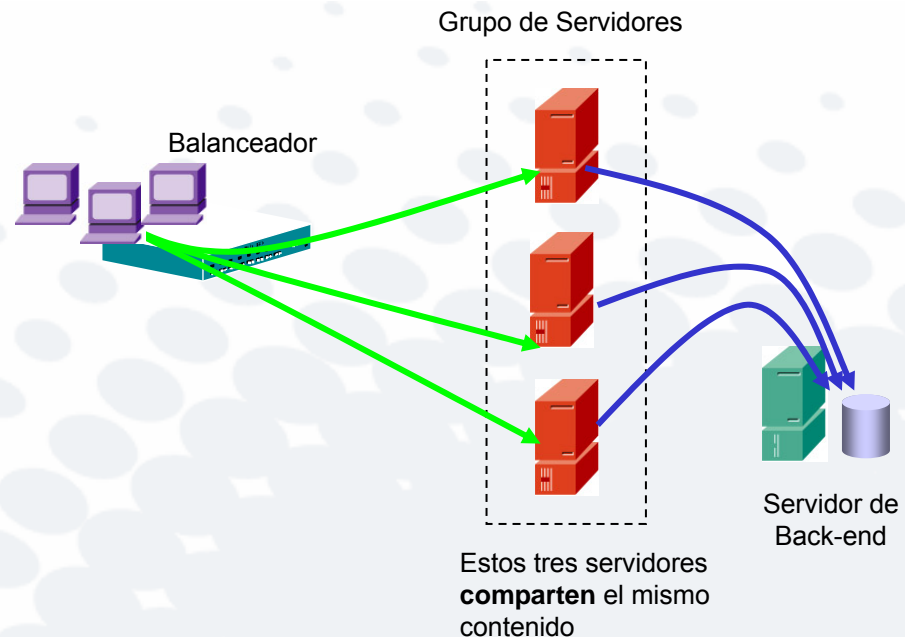
Es posible extender las capacidades de balanceo a nivel URL. De esta forma es posible distribuir la petición de los objetos que componen una página Web hacia varios servidores



El contenido en el grupo de servidores

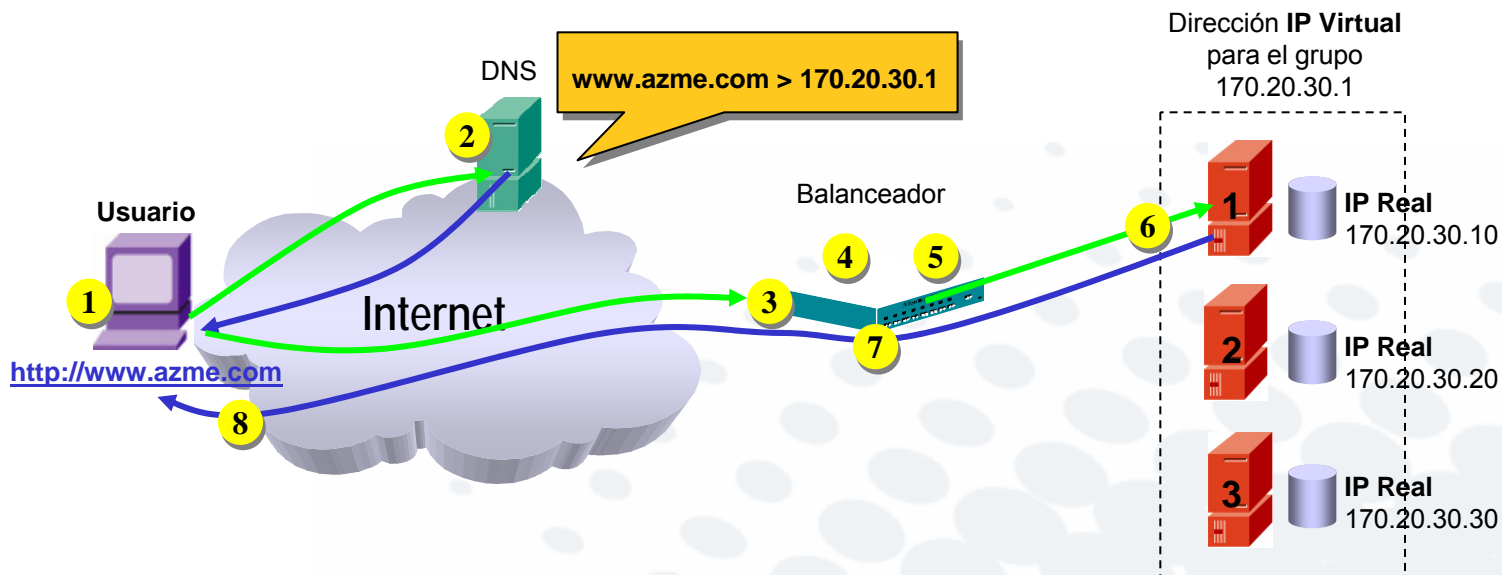


El pool de servidores dispone del mismo Contenido. Estos permanecen sincronizados o simplemente el Web master actualiza todas las máquinas simultáneamente. Este modelo es recomendable solo cuando los servicios IP son solo de lectura



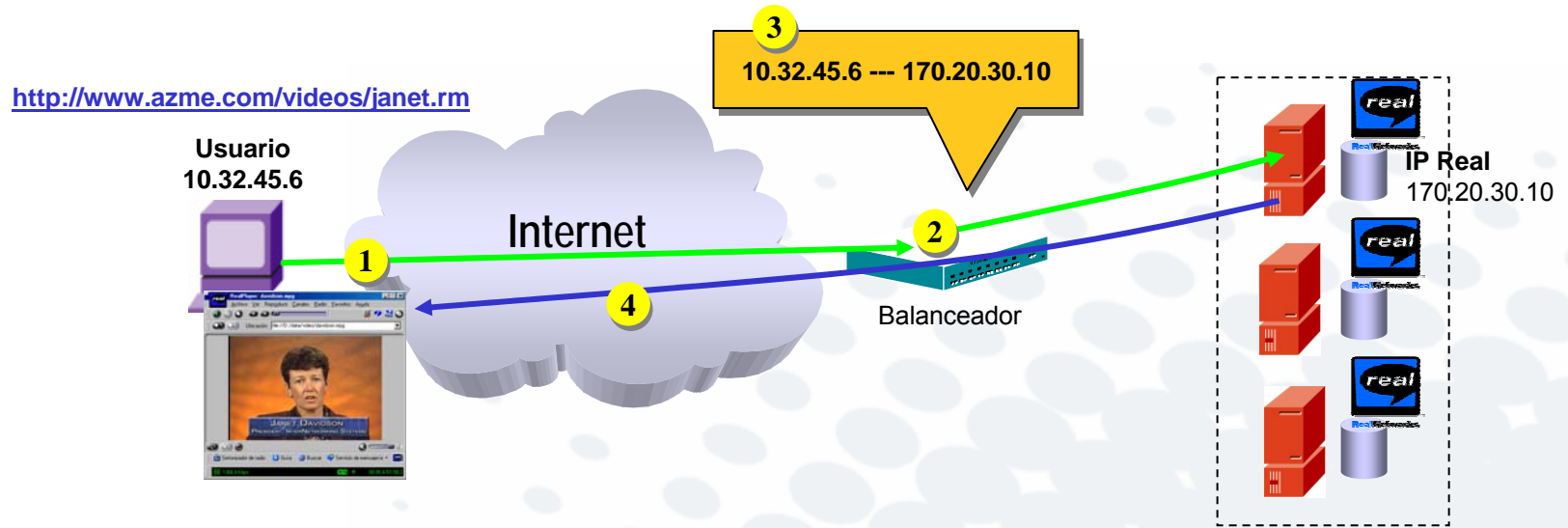
En este escenario todos los servidores acceden y comparten una base de datos común. En los servidores únicamente permanecen las páginas y objetos estáticos. Recomendado en aplicaciones que interactúan con bases de datos en lectura y escritura

Detalle de funcionamiento



1. Un usuario inicia una sesión contra www.azme.com
2. El DNS resuelve en nombre www.azme.com con la dirección IP 170.20.30.1
3. El usuario hace un HTTP contra 170.20.30.1 (IP virtual del grupo)
4. El balanceador identifica la sesión del usuario mediante el flag de estado TCP (SYN/FIN), dirección IP origen/destino y puerto TCP
5. El balanceador monitoriza el estado de cada servidor, y en función de ello, asigna la sesión al servidor 1. Además, el balanceador sustituye la IP destino 170.20.30.1 por la 170.20.30.10
6. El servidor 1 procesa la sesión del usuario hasta que esta finalice
7. El tráfico de retorno hacia el usuario sufre la translación de dirección IP origen 170.20.30.10 por 70.20.30.1
8. El usuario recibe el contenido de la página www.azme.com resultando transparente el proceso

Tratamiento de tráfico UDP



- UDP es un protocolo no orientado a la conexión, por lo tanto no existe mecanismo alguno que señale el comienzo o fin de una sesión (1).
- Para hacer un seguimiento de estos tráficos el balanceador detecta el comienzo de la sesión en el momento que se cursa tráfico UDP (2).
- A continuación, establece una entrada en una tabla donde vincula dirección IP del usuario con la del servidor que atenderá el servicio (3).
- El vínculo y la comunicación se mantendrá mientras que el balanceador detecte tráfico UDP (4).

Persistencia de Sesiones

La persistencia de sesiones permite al balanceador asignar un usuario a un servidor específico:

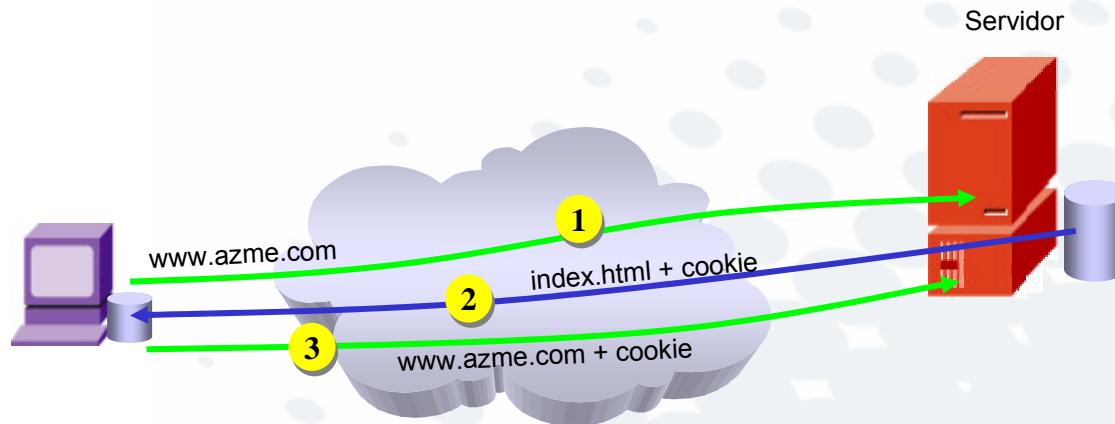
- Cuando se necesita restablecer una sesión de SSL entre un usuario y el servidor con el que inicialmente negoció la sesión
- Si el servidor dispone de ciertos datos del usuario que no han sido o no serán compartidos con los restantes servidores del grupo
- En un servidor concreto, trato diferenciado a un usuario a partir de su identidad

Tres técnicas permiten identificar a un usuario para lograr la persistencia de la sesión:

- Vínculo de IP usuario e IP del servidor (incompatible con Proxys IP y NAT)
- *Cookies* HTTP
- *SSL Sesión ID* en HTTPS

Persistencia de sesiones HTTP (*cookies*)

La entrega de *cookies* a un usuario por parte del servidor, permite a este último identificar a los usuarios en posteriores visitas.



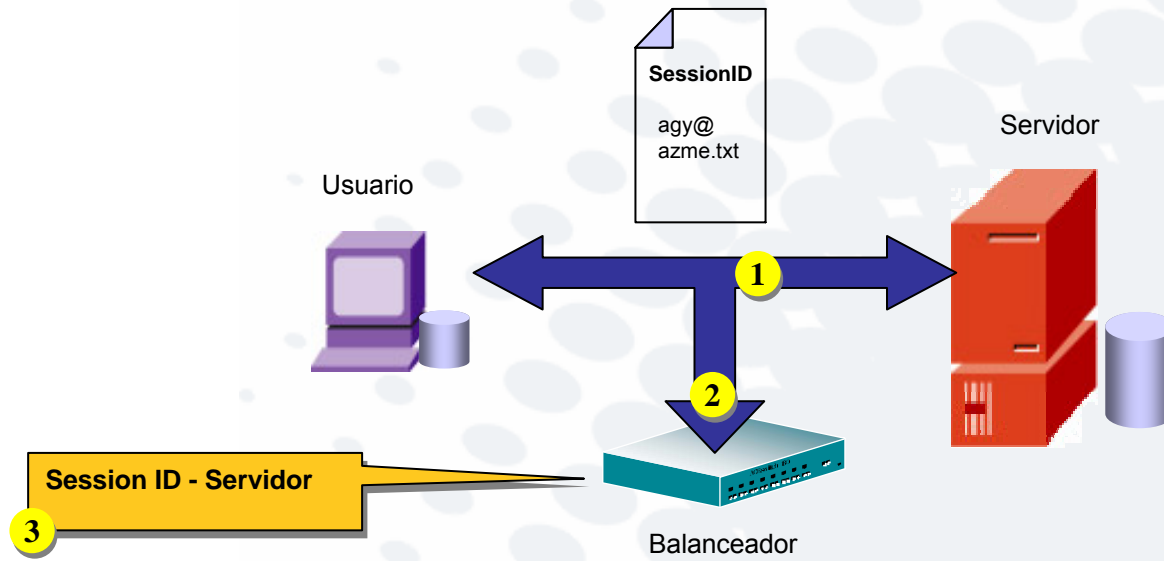
Mediante la interpretación de una *cookie*, el balanceador puede asignar una sesión de usuario a un servidor específico. Hay dos métodos de operación:

- Pasivo
- Activo

Persistencia con *cookies* en modo pasivo

En el modo pasivo, el balanceador observa el cruce de cookies (1-2) y confecciona una tabla de vínculos donde anota el *SessionID* de cada *cookie* con el servidor que la generó y el usuario que la posee (3).

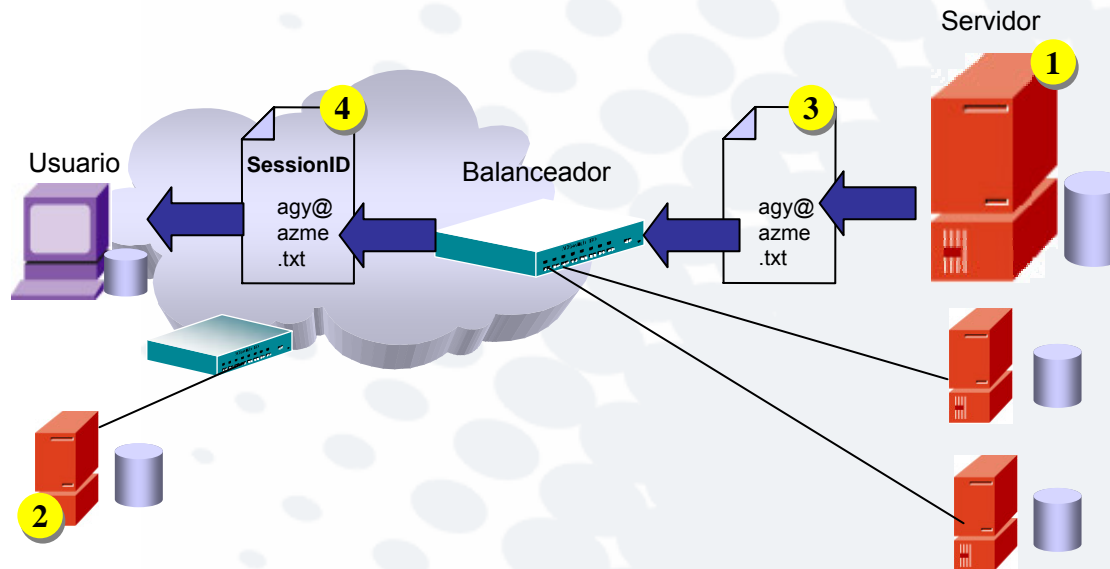
En sucesivas visitas de los usuarios, mediante esta tabla de vínculos, el balanceador conoce hacia que servidor encaminará la sesión



Persistencia con *cookies* en modo activo

El balanceador busca quien es el servidor más apropiado para tratar una sesión, el servidor puede ser local (1) o global (2). En este modo el servidor pone la *cookie* (3) y el balanceador es quien genera los *SessionID* (4) para cada *cookie* a favor del servidor.

Este modelo es el recomendado en escenarios de balanceo global



Persistencia de Sesiones HTTPS (SSL)

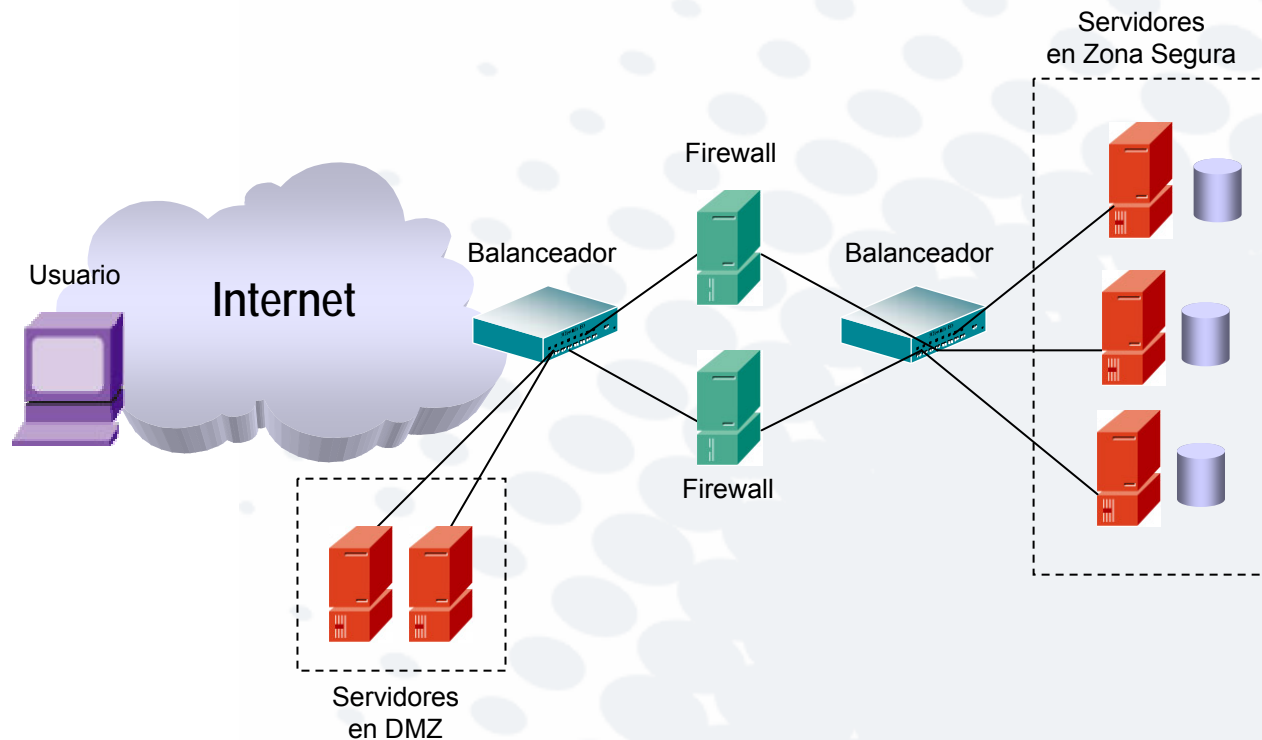
SSL es un mecanismo que permite cifrar una sesión TCP. Este protocolo es actualmente la base para desarrollar transacciones seguras en HTTP (HTTPS).

Durante la negociación SSL entre usuario y servidor se establece un *SessionID* que identifica inequívocamente la sesión segura. Hasta que el *SessionID* expire puede ser empleado para iniciar más de una sesión SSL. Es decir, el *SessionID* puede ser reutilizado entre un usuario y un servidor más de una vez.

El balanceador es capaz de reconocer el *SessionID* y, en función de su valor, es capaz de dirigir a un usuario hacia un servidor concreto. Únicamente entre ambos tendría validez este *SessionID*.

Balancedo de carga de Firewalls

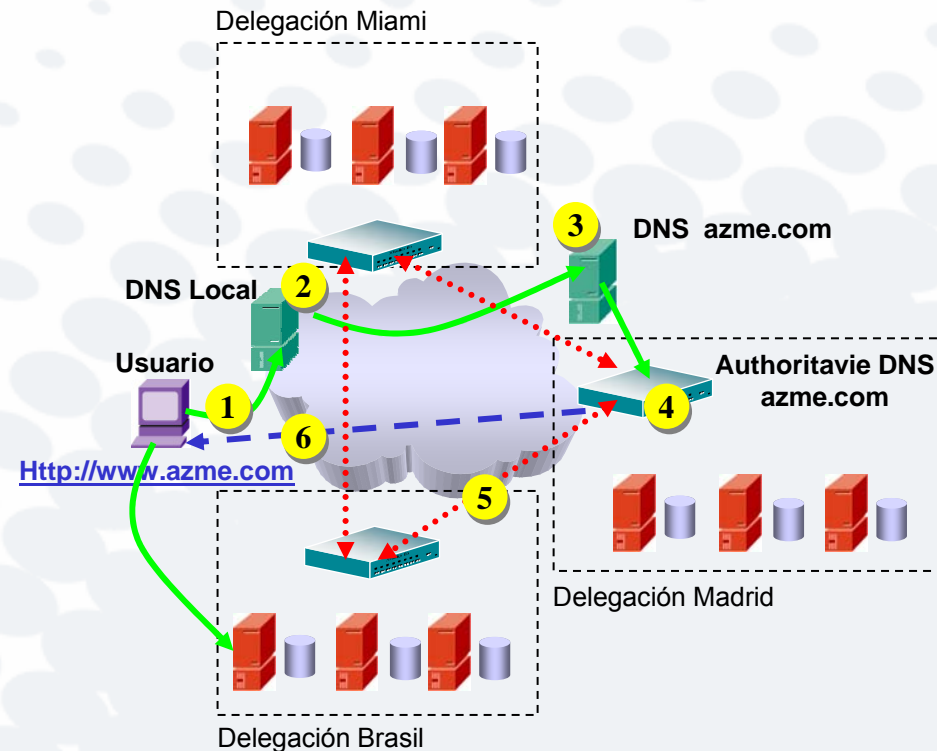
Mediante el empleo de balanceadores es posible disponer de un entorno de firewalling redundado. De esta manera se aumenta el rendimiento y se elimina el punto de fallo que puede representar un único firewall



Balaneo global y Content Routing

El Balanceo Global permite a una organización balancear entre sedes distantes el tráfico de usuarios. En estos escenarios es posible distribuir la carga de forma homogénea, asegurando al usuario el mejor servicio y/o tratamiento idiomático.

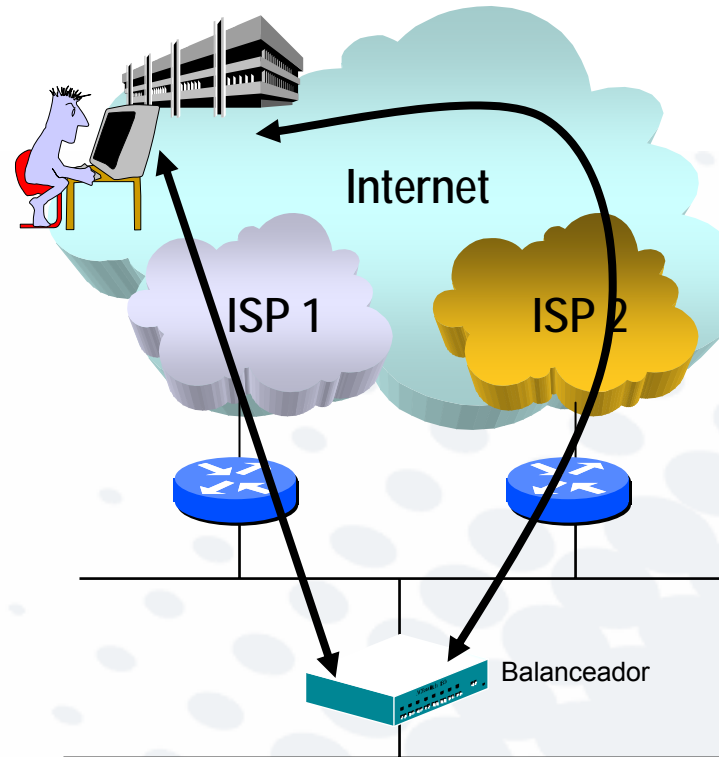
1. El usuario desea iniciar una visita a azme.com, para ello necesita conocer la dirección IP del servidor web de azme
2. El DNS local del usuario escala la resolución hasta llegar al DNS del dominio azme.com
3. El DNS de azme.com traslada la petición al DNS autorizado
4. El balanceador está configurado para desarrollar la funcionalidad de DNS, en este caso es el DNS autorizado de azme.com
5. Entre todos los balanceadores se mantiene una comunicación para conocer el estado de sus respectivos servidores
6. Tras decidir sobre que servidor se desarrollará la sesión al usuario se le comunica la dirección IP válida para www.azme.com



Balancedores acceso a Internet

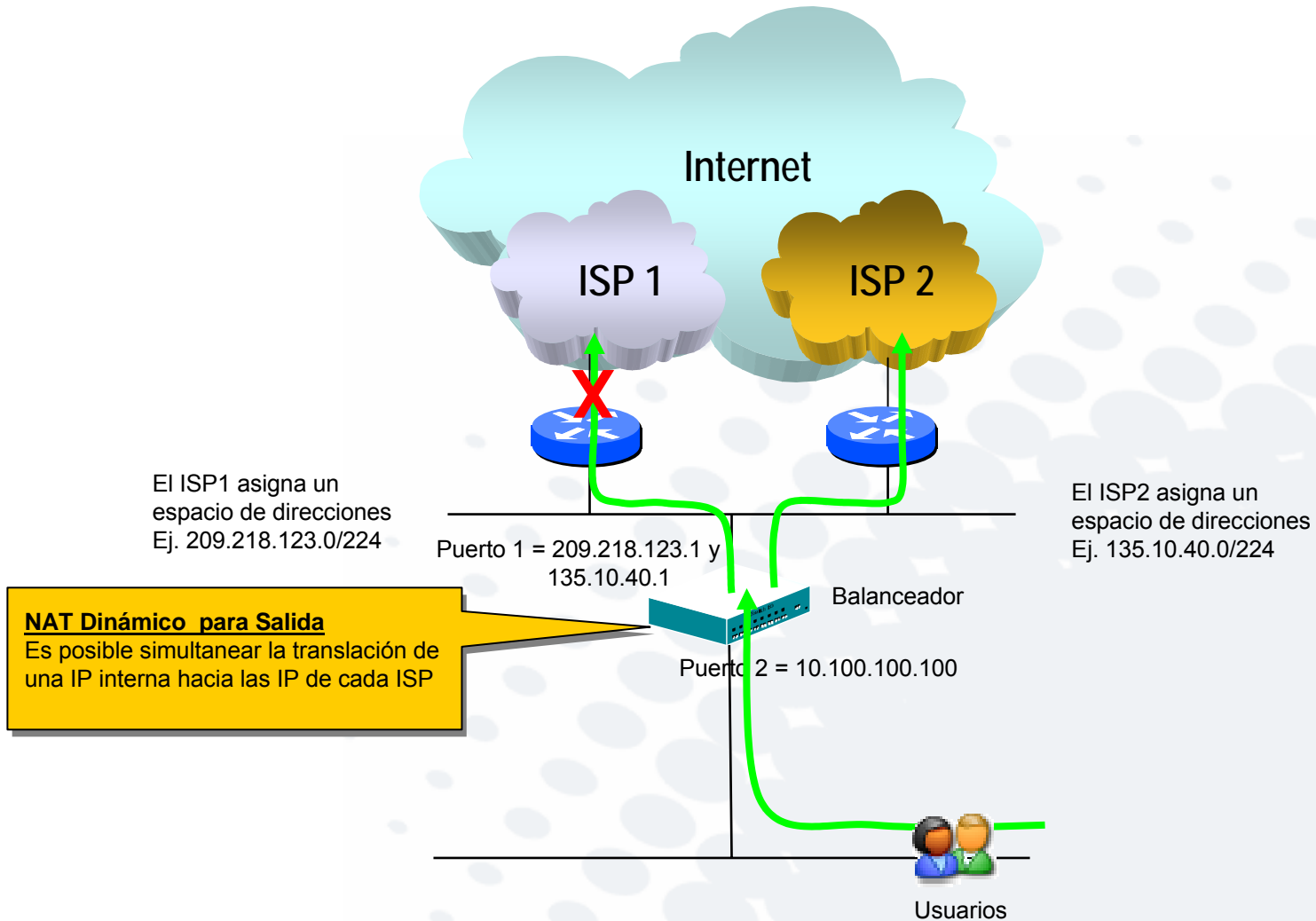
- Permiten balancear el tráfico hacia dos o más conexiones a Internet, a través de diferentes ISP
- Se asegura el uso óptimo y equilibrado de todas las conexiones a Internet
- Antes de encaminar a un usuario hacia un determinado ISP, es posible monitorizar el estado de cada router de salida y la ruta completa
- Criterios de balanceo: N° Sesiones, VLANs, Reglas, Detección de congestión, caída de un enlace, etc.
- Ciertos balanceadores, mediante algoritmos para el cálculo de métricas, son capaces de determinar –en ambos sentidos – la ruta óptima a través de cada ISP

Cálculo de métricas

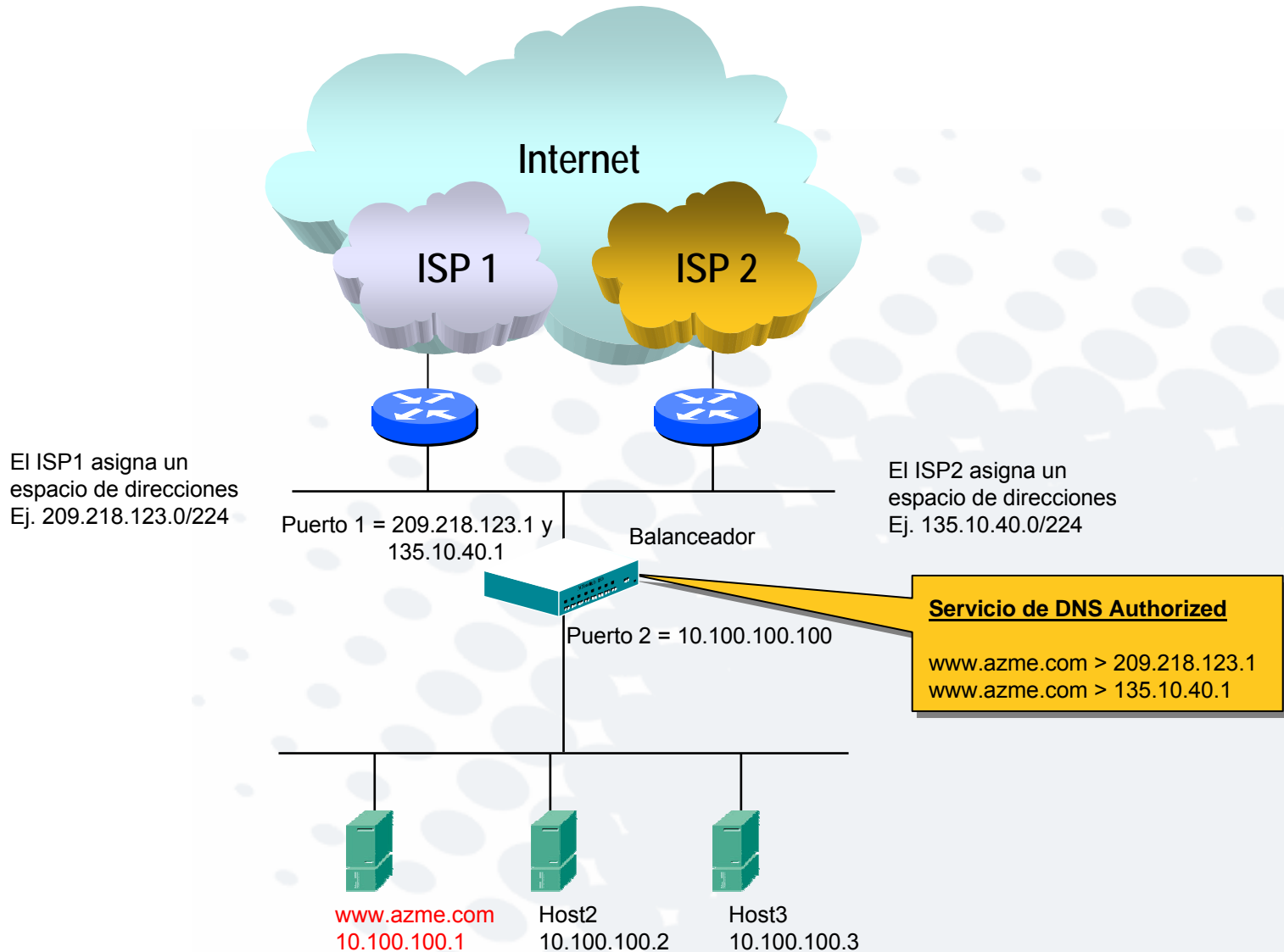


- Para cada usuario que sale a Internet se calcula la ruta óptima, en función de ello se le asigna una IP del ISP que facilite el mejor trayecto
- Para un usuario que visita nuestros servicios, se valorará a través de que ISP alcanza mejor nuestra web. A partir de esta medida, el servicio DNS del balanceador resolverá el nombre de dominio con una IP del rango del ISP más favorable.

Balancedo tráfico de salida a Internet

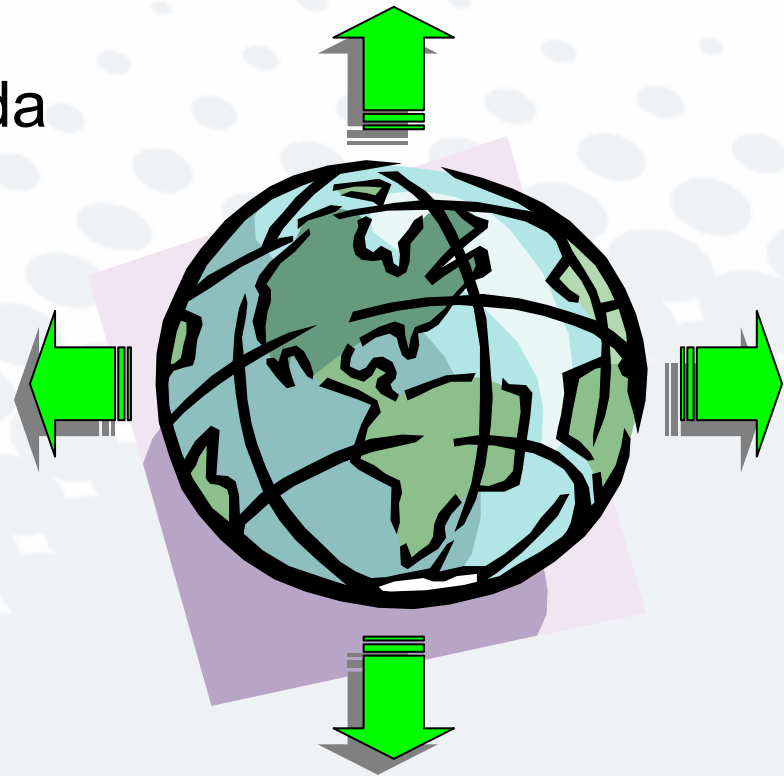


Balaneo tráfico de entrada a Internet



Agenda

- Gestores de Ancho de Banda
- Caché de contenidos
- Balanceadores de Carga
- **Aceleradores SSL**



e-commerce y SSL

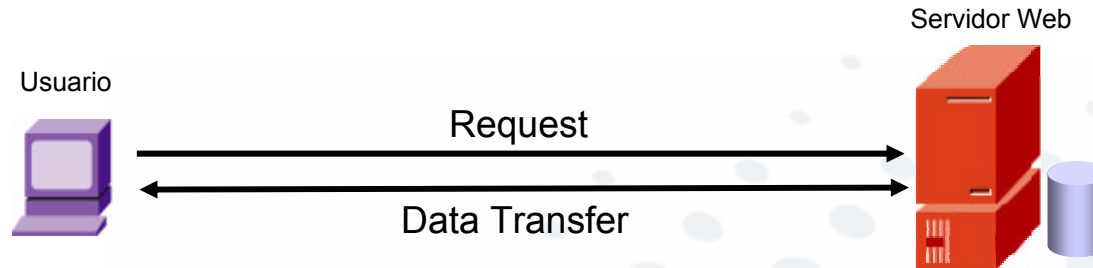
En la actualidad, la base para desarrollar de forma segura transacciones de comercio electrónico sobre Internet es SSL

SSL es un protocolo de propósito general que se utiliza para enviar información cifrada. Las capacidades de encriptación de SSL se desarrollan a nivel 4, siendo posible cifrar sesiones HTTP (https), SMTP (ssmtp) y LDAP (ssl-ldap) entre otras muchas

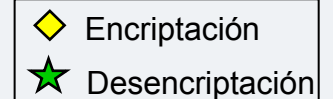
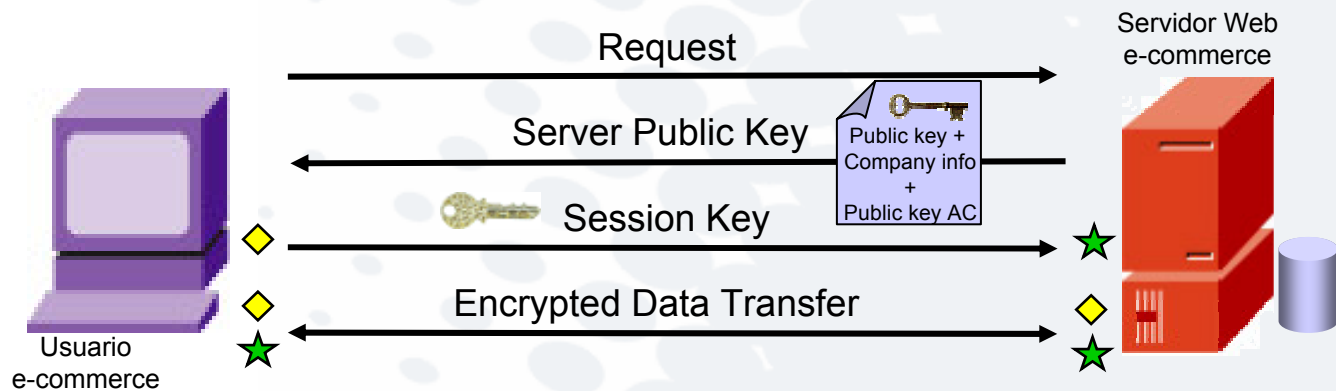
SSL se apoya en el sistema de intercambio de clave pública RSA, además exige el envío de certificados X.509 v3 para autenticar la identidad de la empresa "on-line". Estos certificados X.509 pueden haber sido certificados a su vez por una autoridad de certificación.

Establecimiento de una sesión segura SSL

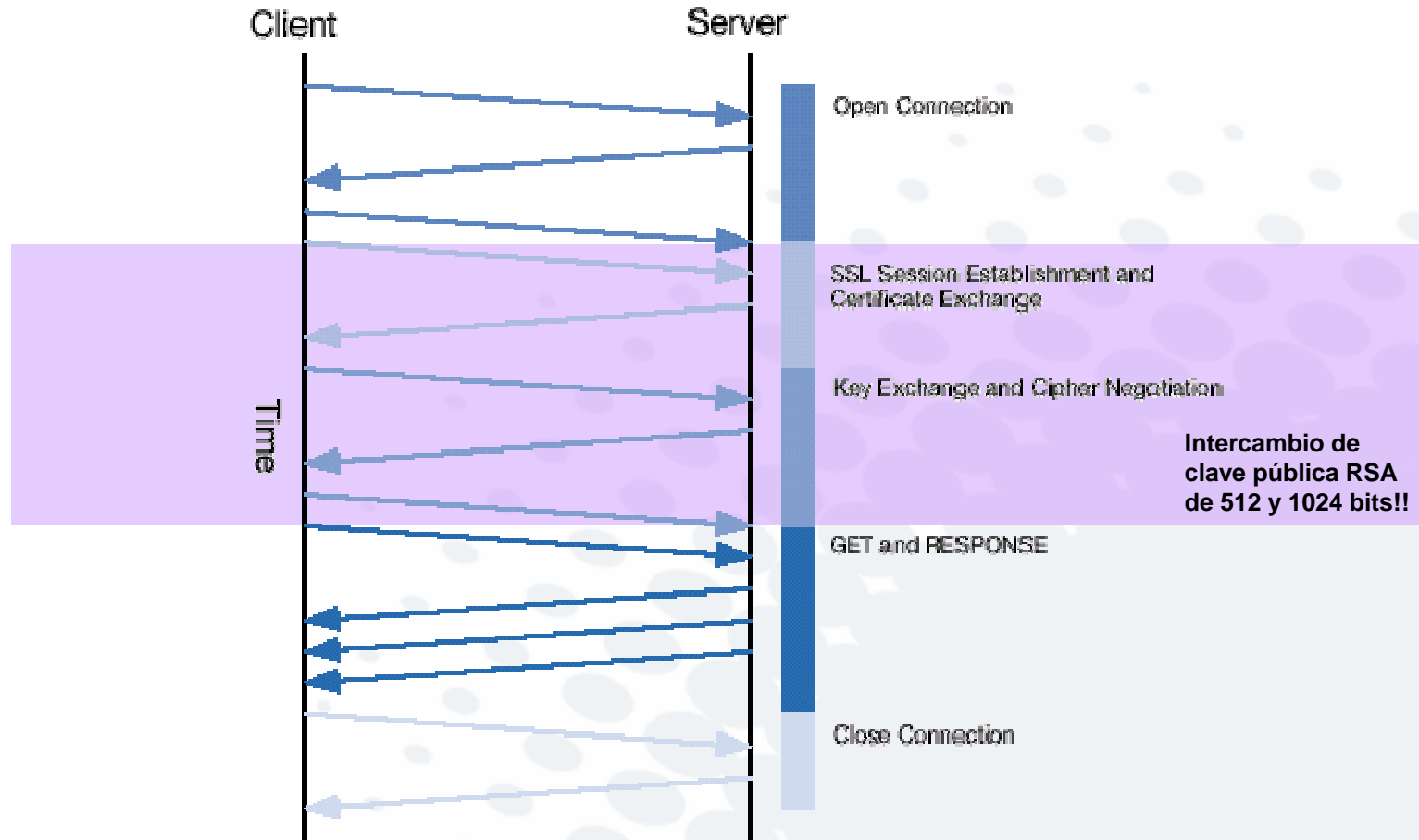
Transacción Normal:



Transacción Segura:



Establecimiento de una sesión segura SSL

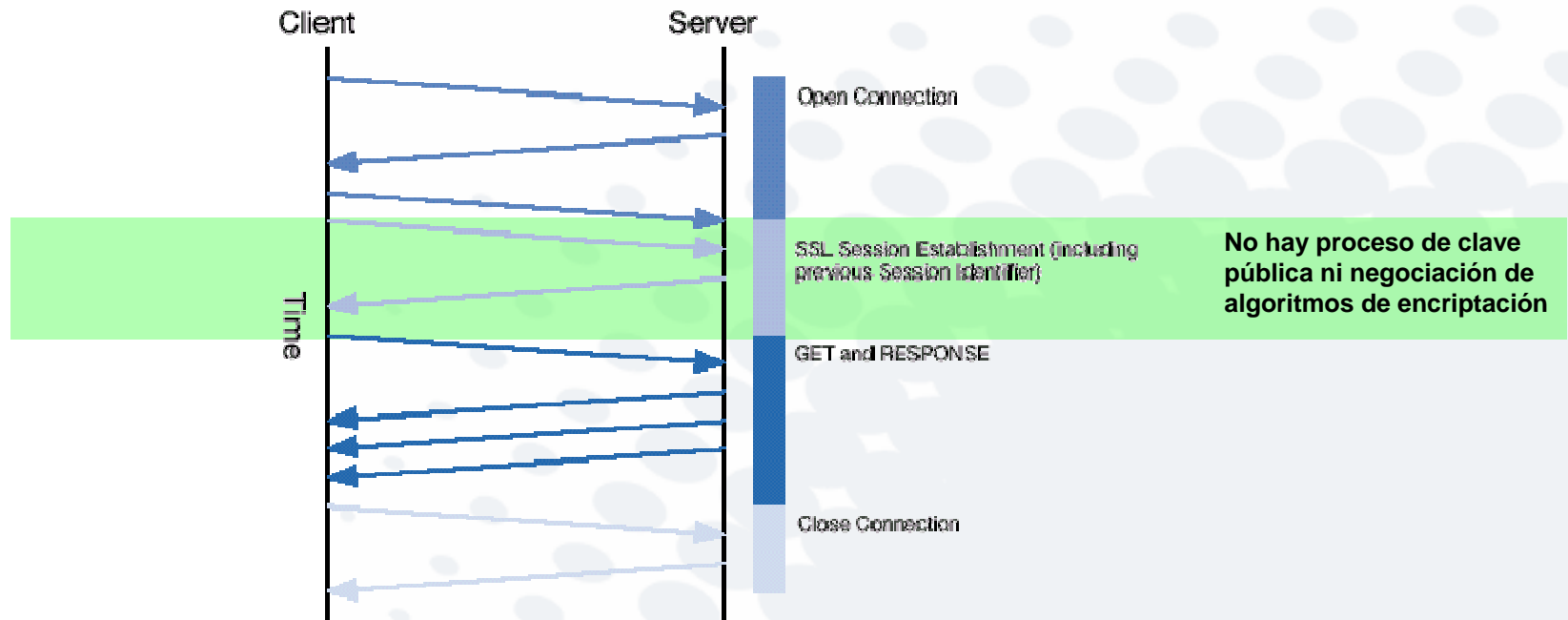


Esta actividad inicial, incluyendo el intercambio de claves y la negociación del protocolo de encriptación, es la que consume más tiempo de proceso. Tras el establecimiento de la sesión es asignado un *Session Identifier (SessionID)*

SessionID

El *SessionID* hace referencia a la clave maestra y al protocolo de encriptación empleado en la comunicación entre cliente y servidor.

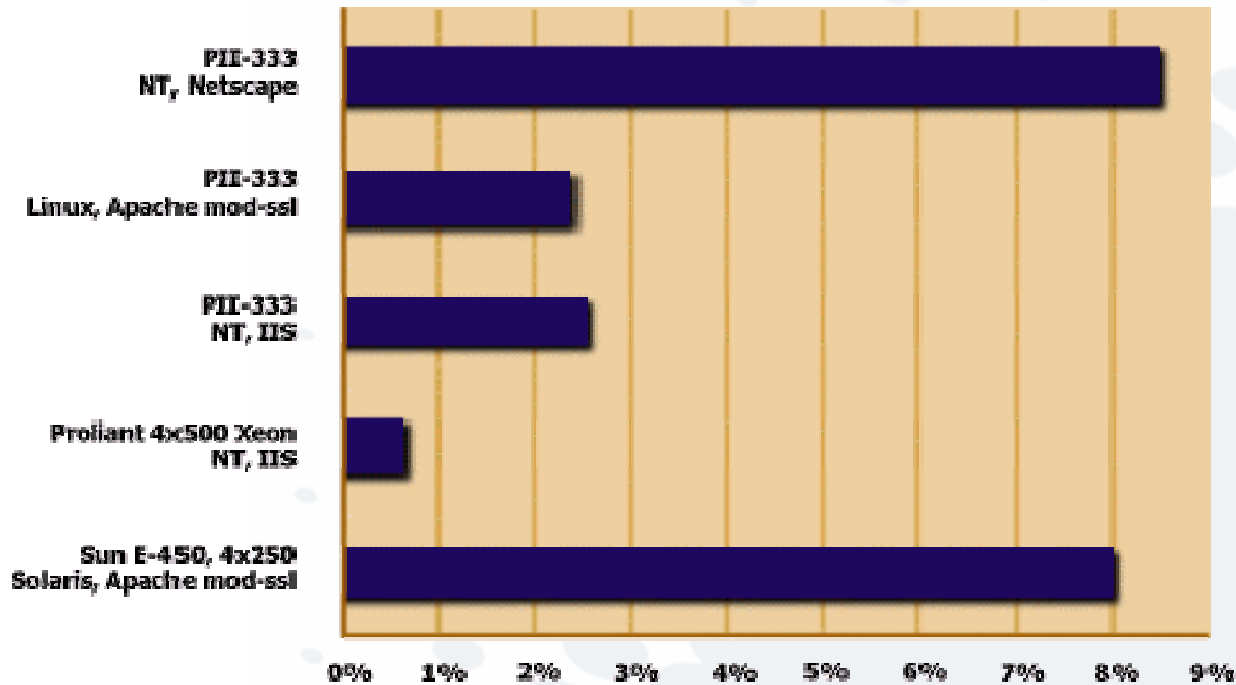
Es conservado por el cliente y el servidor para establecer futuras conexiones SSL, sin tener que volver a iniciar un proceso de negociación con “pesados” sistemas de clave pública. A esta característica de reutilización se la denomina *SessionID re-use*.



Los navegadores intentan reutilizar su *SessionID* tantas veces como les sea posible, de hecho mantienen en caché el *SessionID*. La mayoría de los servidores Web recuerdan los *SessionID* hasta que expire el tiempo establecido en parámetros del tipo *SSLSessionCacheTimeout*

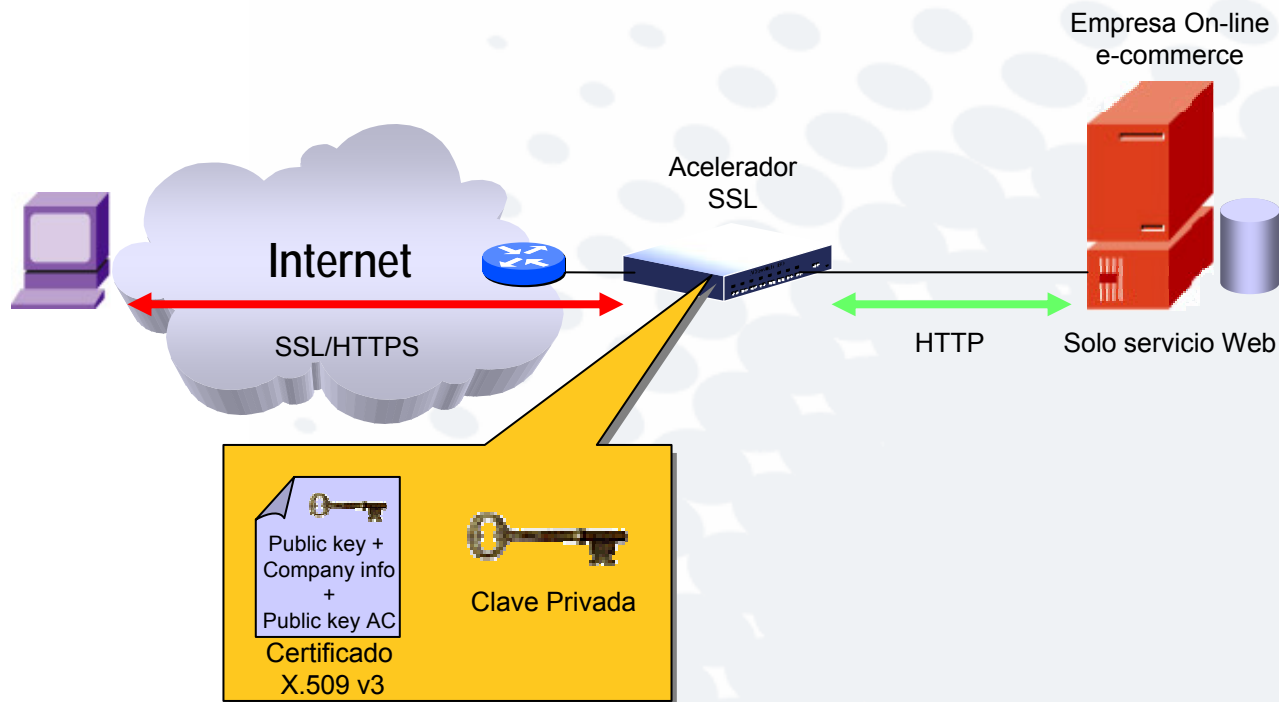
% de CPU consumido por sesión SSL

Aún haciendo uso de la característica *SesiónID re-use*, la negociación SSL penaliza gravemente el rendimiento del servidor, provocando lentitud e inestabilidad en el lado del cliente.



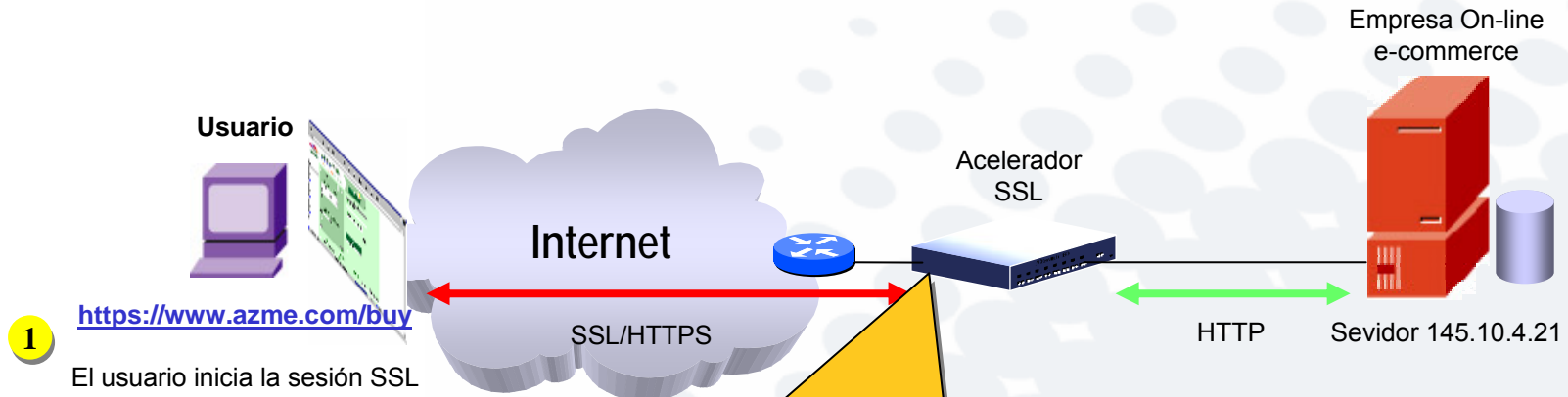
Aceleradores SSL

Las soluciones de aceleración de SSL están orientadas a descargar al servidor Web de todos los procesos de negociación y mantenimiento de *SessionID*



¿Como funciona?

En el acelerador se activa una tabla donde se vincula la dirección IP del servidor de comercio electrónico junto con los servicios donde se desarrollará SSL. Por último, a cada entrada se le anexa el certificado y la llave correspondiente.



1

<https://www.azme.com/buy>
El usuario inicia la sesión SSL

IP Servidor	Tabla de Mapping Servicio	Certificado+Key
145.10.4.21	https(443)	cert1 edTR#45%g
145.10.4.21	http (80)	cert1 edTR#45%g
...		

2

El acelerador intercepta el SSL Request del usuario y atiende la petición, el negocia la sesión SSL

Gracias por su Atención

