

# Requisitos de Banda de Rede para Transmissão Otimizada de Mídia Contínua para Usuários Interativos

Marcelo Maia, Marcus Rocha, Ítalo Cunha, Jussara Almeida, Sérgio Campos

Departamento de Ciência da Computação  
Universidade Federal de Minas Gerais  
Av. Antônio Carlos, 6627, Pampulha  
Belo Horizonte, MG, Brasil, 31270-010

{mmaia, mvrocha, cunha, jussara, scampos}@dcc.ufmg.br

## ABSTRACT

In order to minimize bandwidth requirements and improve on demand streaming media distribution scalability, several distribution protocols based on stream sharing by multiple users have been proposed. Despite the great scalability of these protocols proven for workloads where users access the entire media with no interruptions, several studies show that protocol scalability is severely degraded under interactive scenarios, where users access media segments. These scenarios are commonly observed in real streaming media server workloads. In reply to these results, a number of extensions to the original protocols, optimized for interactive access, have been proposed and evaluated, indicating significant server bandwidth savings even under highly interactive scenarios. However, one of the most promising class of optimizations may send extra streams through the network during inactive client periods. Therefore, network bandwidth requirements for these optimizations become a critical issue that has not been evaluated yet. This work presents an evaluation of the network bandwidth requirements for the interactive access optimized protocols, considering realistic workloads of different interactivity levels and canonical and real network topologies. Results indicate that the optimizations also lead to a significant reduction in the average network bandwidth requirements (peak of 60%).

## RESUMO

Visando minimizar os requisitos de banda e aumentar a escalabilidade da distribuição sob demanda de mídia contínua, vários protocolos de distribuição baseados em compartilhamento de fluxos por múltiplos usuários foram propostos. Embora esteja comprovada a grande escalabilidade desses protocolos para cargas onde os usuários acessam toda a mídia sem interrupções, vários estudos mostraram que a mesma é severamente degradada em cenários interativos, onde os usuários acessam segmentos da mídia. Cenários comumente observados em cargas de acesso a servidores reais de mídia contínua. Em resposta a esse resultado, várias extensões dos protocolos originais, otimizadas para acesso interativo, foram propostas e avaliadas, indicando uma grande economia de banda de servidor mesmo em cenários com alta interatividade. Entretanto, uma das classes de otimizações mais promissoras pode enviar fluxos extras pela rede durante períodos inativos do cliente. Assim, os requisitos de banda de rede dessas otimizações se torna um ponto crítico que ainda não foi avaliado. Este artigo avalia o consumo de banda de rede pelos protocolos otimizados para acessos interativos, considerando cargas de trabalho realistas com diferentes níveis de interatividade e topologias de redes canônicas e reais. Os resultados indicam que as otimizações também levam a uma significativa redução no consumo médio de banda de rede (pico de até 60%).

## 1. INTRODUÇÃO

Uma das aplicações que mais crescem na Internet atualmente é a de conteúdo multimídia. Aulas virtuais ajudam tanto no treinamento pessoal quanto na inclusão digital. Noticiários, seriados e rádios virtuais são outros exemplos notórios. O grande desafio está em entregar todo o conteúdo, cujos requisitos de banda de servidor e rede são altos, com a maior qualidade possível e uma quantidade vasta de usuários exigentes. Nesse contexto, diversas questões estão envolvidas, tais como arquiteturas eficientes de cache ([13] e referências contidas), mecanismos de codificação [22] e protocolos eficientes de distribuição de conteúdo [2, 3, 5, 7, 8, 9, 10, 11, 14].

Considerando acesso seqüencial por parte dos usuários, ou seja, o acesso completo do início ao fim da mídia sem interrupções, dois protocolos de distribuição sob demanda de mídia contínua com compartilhamento de fluxos se destacam. Além de permitirem serviço imediato, *Patching* [10, 11] e *Bandwidth Skimming* [8, 9] são capazes de reduzir significativamente os requisitos de banda de servidor, se comparados com a transmissão *unicast* que impõe um crescimento linear no consumo de banda com a taxa de chegada de requisições. Essa redução de banda é obtida através de fusões de fluxos que transmitem conteúdos próximos no tempo. Dizemos que dois fluxos foram fundidos em um só se os clientes que recebiam o conteúdo de ambos passam a receber de apenas um, sendo um dos fluxos interrompido. Dizemos também que o fluxo restante passa a ser compartilhado entre os clientes.

Diante do mesmo cenário seqüencial, em [23] são derivados os requisitos mínimos de banda de rede necessários para a transmissão do conteúdo multimídia em árvores de distribuição mais simples. Para topologias sintéticas mais complexas, também é mostrado que o *Bandwidth Skimming* é eficiente não só para reduzir os requisitos de banda média de servidor, mas também os de banda de rede.

No entanto, trabalhos de caracterização de cargas de servidores reais como [1, 4, 6, 16] evidenciam a forte presença de interatividade, ou seja, acessos a segmentos da mídia. Por exemplo, pausas e saltos são muito freqüentes em cargas educacionais, onde os usuários tendem a rever o conteúdo com a finalidade de tirar as dúvidas das aulas. Essas freqüentes interrupções atrapalham as fusões dos fluxos, comprometendo o compartilhamento no servidor e, conseqüentemente, levando a um aumento dos requisitos de banda [12, 21, 17].

Em [17, 18, 19] são propostas estratégias de otimização para o protocolo *Bandwidth Skimming* especializadas para acesso interativo. As estratégias foram extensamente avaliadas quanto aos requisitos de banda de servidor. Foi mostrado que o armazenamento de

conteúdo em *buffer* para de evitar retransmissão pode significativamente reduzir os requisitos de banda. Economias expressivas de até 66% sobre o protocolo original puderam ser alcançadas com uma área de armazenamento ilimitada. Economia de até 42% ainda pode ser obtida restringindo a área de armazenamento a apenas 30% da duração da mídia. Entretanto, uma das classes de otimizações mais promissoras pode enviar fluxos extras pela rede durante períodos de inatividade do cliente. Assim, os requisitos de banda de rede dessas otimizações se torna um ponto crítico, que ainda não foi avaliado. Uma questão chave que surge é: *quais os requisitos de banda de rede necessários para que essas estratégias de otimização consigam a expressiva economia de banda obtida no servidor?*

Em [18, 19] é feita uma primeira abordagem do impacto das estratégias de otimização sobre a banda de rede considerando apenas uma topologia canônica bem simples. Os resultados preliminares indicam que para as cargas interativas também é possível obter simultaneamente grandes economias de banda de servidor e de rede em relação ao protocolo *Bandwidth Skimming* original.

O objetivo deste trabalho é fazer uma avaliação, via simulação, bem mais extensa do desempenho das otimizações para o protocolo *Bandwidth Skimming* propostas em [17, 18, 19] quanto aos requisitos de banda de rede, bem como apresentar os principais compromissos envolvidos entre os requisitos de banda de servidor e de rede. Foram usados 36 perfis de acesso a servidores reais de mídia contínua de diferentes tipos de aplicações (entretenimento, educacional, etc) caracterizados em [6]. Cada um deles foi avaliado em um amplo intervalo de taxa de chegada de requisições de usuários e em 4 topologias diferentes de rede, sendo 2 reais e 2 canônicas.

As principais contribuições deste artigo são:

- Avaliação de estratégias de otimização para o protocolo de distribuição sob demanda de mídia contínua com compartilhamento de fluxos *Bandwidth Skimming* quanto aos requisitos de banda média de rede para cargas interativas em topologias canônicas e reais.
- Identificação da estratégia HYBRID AGGRESSIVE (HA) como a que tem maior economia de banda média de rede (até 60%), porém limitada a cargas com taxa de chegada de requisições mais alta, e das estratégias HYBRID SYNERGIC (HS) e HYBRID CONSERVATIVE (HC) como as que apresentam as melhores relações entre economia simultânea de banda média de servidor, até 54% (HS) e 52% (HC), e rede, até 43% (HS) e 48% (HC) por faixa mais ampla de taxa de chegada.

O restante do artigo está organizado como a seguir. A seção 2 discute os trabalhos relacionados. A seção 3 apresenta as estratégias já propostas e resultados obtidos. A metodologia empregada, incluindo a descrição das cargas realistas e topologias de rede usadas, é mostrada na seção 4. A avaliação das estratégias de otimização nas topologias utilizadas e os principais compromissos é feita na seção 5. A seção 6 conclui o trabalho e aponta futuras direções.

## 2. TRABALHOS RELACIONADOS

A fim de contornar a escalabilidade limitada da transmissão *unicast* de mídia contínua, vários protocolos baseados no compartilhamento de fluxos foram propostos [2, 3, 5, 7, 8, 9, 10, 11, 14]. Dois despertam interesse por permitirem serviço imediato. *Patching* [10, 11] foi o primeiro a explorar a transmissão dois fluxos a um mesmo usuário, um para serviço imediato e outro, o alvo, para

adiantar conteúdo. *Bandwidth Skimming* [8, 9], diferentemente de *Patching*, permite compartilhamento hierárquico de fluxos.

Para uma carga seqüencial e processo de chegadas Poisson, a banda média requerida pelo servidor utilizando os protocolos *Patching* e *Bandwidth Skimming* para a distribuição da mídia pode ser calculada analiticamente [9]. Ela depende apenas da taxa de chegada de requisições normalizada pela duração da mídia  $N = \lambda T$ , onde  $\lambda$  é a taxa de chegada de requisições e  $T$  é o tamanho da mídia. Diante do contexto seqüencial, o protocolo *Bandwidth Skimming*, com a política de escolha de alvo *Closest Target* [9], tem escalabilidade superior em termos de banda média requerida pelo servidor [18].

Em [23] são derivados os requisitos mínimos de banda de rede necessários para a distribuição do conteúdo multimídia em árvores de transmissão multicast canônicas para acesso seqüencial. Através de topologias canônicas simples e sintéticas mais complexas é mostrado que é possível alcançar simultaneamente grande economia de banda de servidor e de rede.

No entanto, caracterizações de servidores reais como [6, 1, 4, 16] mostram que a interatividade é freqüente. Os usuários realizam pausas e saltos e não necessariamente assistem toda a mídia. Estudos analíticos, com modelos de interatividade simplificados [21], e experimentos realistas [1, 17] mostram que a interatividade causa forte queda no desempenho do *Bandwidth Skimming*, fazendo com que seu consumo de banda média seja superior ao logarítmico obtido com o acesso seqüencial [9]. Nesse contexto, com a finalidade de minimizar esse efeito negativo da interatividade, as estratégias de otimização para os protocolos se fazem necessárias.

Em [15] é proposto o *Patching Interativo* e em [20] os autores o otimizam propondo o *Patching Interativo Eficiente* e o *Patching Interativo Completo*. Esses últimos avaliam limiares de tempo na decisão de fusão e abertura de novos fluxos. Em [14] é proposto um novo protocolo, BEP (*Best Effort Patching*), onde a hierarquia de fusões, apesar de ainda limitada, é estendida para um terceiro nível, em contraste com os dois níveis originais do *Patching*.

Em [18] é mostrado que o protocolo *Bandwidth Skimming* apresenta melhor escalabilidade também para cargas interativas. Assim esse protocolo foi escolhido como base para a proposta das otimizações em [17, 18, 19]. Elas foram amplamente avaliadas quanto aos requisitos de banda média de servidor. As avaliações foram conduzidas utilizando cargas realistas de diversos perfis de interatividade e variados tamanhos de *buffer* e taxa de chegada de requisições. Comparada com o protocolo original, a melhor otimização, HYBRID AGGRESSIVE, economiza a banda média de servidor em até 66%, utilizando uma área de armazenamento ilimitada. Economia significativa de até 42% ainda é obtida restringindo essa área a apenas 30% da duração da mídia.

Em [17, 18, 19] foi mostrado que a utilização de *buffer* na tentativa de evitar retransmissão é capaz de economizar significativamente a banda média de servidor. Entretanto, a fragmentação das requisições decorrente do seu uso apresenta um efeito degenerativo, pois prejudica o processo de fusão dos fluxos. Ele força a divisão das requisições enviadas pelos usuários em vários segmentos de durações reduzidas. Como consequência, a duração média das requisições percebidas pelo servidor é menor, dificultando o compartilhamento. Existe, portanto, um compromisso no uso do *buffer*. Ele evita retransmissão de conteúdo, mas fragmenta a carga percebida pelo servidor dificultando o compartilhamento de fluxos.

As estratégias de otimização para *Bandwidth Skimming* que obtêm os melhores resultados fazem uso de períodos em que não haveria consumo de banda por parte dos usuários para adiantar conteúdo e podem, portanto, gerar fluxos extra na rede. Em [18, 19], utilizando uma topologia canônica bem simples, é feita uma primeira abordagem desse impacto sobre os requisitos de banda de rede. Os resultados preliminares com as cargas interativas indicam que as otimizações, que conseguem expressiva economia de banda no servidor, simultaneamente também são capazes de economizar a banda de rede da ordem de 30% a 40% em relação ao protocolo *Bandwidth Skimming* original. O presente trabalho realiza uma extensa avaliação dos requisitos de banda de rede das estratégias de otimização diante de acessos interativos por parte dos usuários e através de topologias reais e canônicas.

### 3. PROTOCOLOS OTIMIZADOS PARA INTERATIVIDADE

O protocolo de distribuição sob demanda de mídia contínua *Bandwidth Skimming* [8, 9] baseia-se na fusão hierárquica de fluxos para diminuir os requisitos de banda de servidor. Os dados são recebidos de dois fluxos simultaneamente. O cliente exibe os dados recebidos de um fluxo criado no instante de sua chegada para serviço imediato e armazena os dados recebidos por outro já ativo. O dado armazenado permite que o cliente alcance o do fluxo já ativo. Nesse ponto os fluxos são fundidos e os clientes compartilham o fluxo restante até o fim da transmissão.

As seções seguintes descrevem as estratégias de otimização para o protocolo *Bandwidth Skimming* propostas em [17, 18, 19] e os principais resultados obtidos. Todas as estratégias têm em comum o uso de *buffer* com a finalidade de evitar retransmissão de conteúdo já recebido. Quando um cliente envia uma requisição ao servidor, ele primeiro verifica quais segmentos já possui. O conteúdo armazenado no *buffer* varia com cada estratégia e existe apenas dentro de uma sessão do usuário, sendo descartado no final dela.

#### 3.1 Estratégias para período ativo

Períodos ativos são aqueles em que o cliente está realmente requisitando conteúdo ao servidor e, conseqüentemente, consumindo banda. De forma inversa, períodos inativos são aqueles em que o usuário faz uma pausa ou o cliente está exibindo conteúdo para o usuário diretamente do *buffer* e, portanto, não está consumindo banda. A seguir é dada a descrição das duas estratégias propostas que exploram os períodos ativos.

**LOCALITY (LOC):** Armazena no *buffer* toda a mídia requisitada pelo usuário e recebida para visualização. Esta estratégia explora a alta localidade de acesso presente em cargas interativas [6].

**KEEP MERGE BUFFER (KMB):** Mantém no *buffer* os segmentos da mídia que seriam descartados devido a uma tentativa de fusão mal sucedida.

#### 3.2 Estratégias para período inativo

Longos períodos de inatividade, observados na caracterização das cargas em [6], motivaram o desenvolvimento dessas estratégias. Durante esses períodos o cliente possui banda livre para adiantar conteúdo que está sendo transmitido para outros clientes ativos. As estratégias podem, portanto, incorrer em um possível consumo extra de banda de rede. A diferença principal entre as estratégias está no método de escolha dos fluxos ativos. Preferência é dada para aqueles cujo conteúdo transmitido esteja mais próximo do último

ponto visualizado pelo usuário, explorando assim, a alta localidade de acesso observada em [6]. As estratégias escutam no máximo 2 fluxos, como na configuração original do *Bandwidth Skimming*.

**SILENT PREFETCH (SP):** A partir do último ponto visualizado, o cliente seleciona o fluxo ativo que esteja transmitindo o conteúdo mais próximo. O seu alvo (no sentido de *Bandwidth Skimming*) também é selecionado, caso ele exista.

**GREEDY PREFETCH (GP):** Tem como premissa adiantar e armazenar a maior quantidade possível de dados, a fim de evitar ao máximo retransmissão de conteúdo. Os fluxos são selecionados independentemente. São selecionados os fluxos que transmitem o dado mais próximo do último ponto visualizado pelo usuário (para frente ou para trás) e cujo conteúdo transmitido seja inédito. GP não previne contra a fragmentação das requisições.

**COOPERATIVE PREFETCH (CP):** Semelhante a GP, no entanto CP preocupa-se com a fragmentação das requisições. O primeiro fluxo é escolhido segundo os mesmos critérios que GP, mas o segundo é selecionado apenas se ele não introduzir mais fragmentação. Isso somente é alcançado caso o conteúdo transmitido pelo fluxo recaia imediatamente após uma região preenchida do *buffer* ou se o fluxo escolhido começar sobrescrevendo algum conteúdo até que o novos dados sejam recebidos.

### 3.3 Estratégias híbridas

As estratégias híbridas combinam as otimizações com a finalidade de explorar as suas diferentes características individuais.

**HYBRID AGGRESSIVE (HA):** Combina as estratégias KMB e LOC do período ativo e GP do período inativo visando adiantar e armazenar a maior quantidade possível de dados em *buffer*, evitando assim ao máximo a retransmissão.

**HYBRID SYNERGIC (HS):** Combina as estratégias do período ativo KMB e LOC com a do período inativo que é menos suscetível ao efeito da fragmentação das requisições, SP.

**HYBRID CONSERVATIVE (HC):** Combina as estratégias de ambos os períodos menos suscetíveis ao efeito da fragmentação das requisições, LOC do ativo e SP do inativo.

### 3.4 Resultados prévios

A banda de servidor, que independe da topologia de rede, foi calculada medindo a quantidade de fluxos necessários para atender todas as requisições dos usuários. Os principais resultados obtidos em [17, 18, 19] são resumidos a seguir.

O desempenho de cada estratégia de otimização depende principalmente da relação entre a redução de retransmissão e a fragmentação das requisições. Quanto mais conteúdo armazenado no *buffer*, maiores as oportunidades de utilizá-lo e evitar retransmissão. Quanto mais conteúdo do *buffer*, menor a fragmentação imposta na carga e maiores as chances de economia de banda.

Otimização como GP, KMB e as híbridas que as incorporam, povam mais agressivamente o *buffer* e, portanto, evitam mais retransmissão. Elas apresentam as maiores economias de banda média se comparadas com o protocolo *Bandwidth Skimming* original, GP com até 56% entre as individuais e HA com até 66% entre as híbridas. Entretanto elas sofrem mais o efeito da fragmentação e perdem desempenho em taxas de chegada mais elevadas.

Estratégias de otimização que têm como objetivo fragmentar menos a carga percebida pelo servidor apresentam picos de economia tipicamente menores, mas economizam a banda média de servidor em uma faixa mais ampla de taxas de chegada. Economia de banda média ainda expressiva de até 51% pode ser obtida com HS.

## 4. METODOLOGIA

O mesmo simulador do *Bandwidth Skimming*, desenvolvido e validado em [18], foi estendido para avaliar os requisitos de banda de rede de cada estratégia de otimização. Para cada perfil de acesso interativo, foram geradas 5 cargas sintéticas com diferentes sementes aleatórias. Os resultados apresentados são uma média das cargas, sendo o desvio padrão inferior a 2% das médias em todos os casos.

### 4.1 Cargas interativas

A avaliação baseou-se em cargas com vários níveis de interatividade de 5 servidores reais caracterizadas detalhadamente em [6]. MANIC<sup>1</sup> e eTeach<sup>2</sup> são sistemas educacionais. Possuem vídeos de tamanhos variados, desde anúncios de 5 minutos até aulas de 1 hora. Outras 3 cargas de 2 grandes provedores de serviço contêm áudio e vídeo de entretenimento com tamanhos tipicamente menores que 10 minutos. Um deles é o UOL<sup>3</sup>, o outro autorizou o uso dos dados, mas não sua identificação.

As cargas reais apresentam variação limitada da taxa de chegada de requisições. Dessa forma, um gerador de cargas, validado em [18] com erros inferiores a 27%, foi utilizado para construir um rico conjunto de cargas sintéticas realistas. Essas cargas sintéticas foram então usadas na avaliação dos requisitos de banda das estratégias. Em [18] elas foram divididas em três níveis de interatividade: Alto (IA), Médio (IM) e Baixo (IB). Essa divisão foi feita com a finalidade de facilitar as análises, uma vez que os resultados entre cada grupo são qualitativamente similares.

### 4.2 Topologias de rede

A banda de rede é calculada como o somatório do número de fluxos ao longo do tempo que atravessam todos os saltos da rede. Ela é afetada por dois fatores: o número de fluxos criados pelo servidor e a quantidade adicional de fluxos de rede necessários para compartilhar os fluxos de servidor por múltiplos *sites*. Considera-se aqui um *site* como sendo a possível localização de um ou mais clientes. O custo de transmitir um fluxo em cada salto é igual para todos.

Os fluxos sempre percorrem a menor distância entre o servidor e o cliente. Em [23] é mostrado que a construção da rede de forma dinâmica, maximizando a quantidade de links compartilhados, reduz os requisitos de banda de rede da ordem de 3% a 16%, se comparado com a construção estática do menor caminho entre o servidor e cada cliente. No entanto, como a construção dinâmica da árvore de distribuição tem maior custo de manutenção, optou-se pela construção utilizando o menor caminho.

Foram utilizadas 4 topologias de rede, sendo 2 canônicas e 2 reais. Uma vez que as topologias canônicas podem ser combinadas para formar topologias mais gerais, os seus resultados foram usados para facilitar a avaliação das topologias reais. Foram usadas as duas topologias canônicas, discutidas a seguir: *Fan-out-K* e *Daisy-Chain-K*, onde  $K$  é o número de *sites* atingidos. Assim como em [23], apenas um servidor é utilizado na avaliação.

<sup>1</sup>RIPPLES/MANIC - <http://manic.cs.umass.edu> (2006)

<sup>2</sup>Learning on Demand - <http://eteach.cs.wisc.edu> (2006)

<sup>3</sup>Universo Online - <http://www.uol.com.br> (2006)

Em ambas as topologias canônicas, os clientes são distribuídos uniformemente entre os  $K$  *sites*, de forma que a taxa de chegada de requisições  $N_i$  vindas de cada *site*  $i$  seja também uniformemente distribuída, ou seja,

$$N_1 = N_2 = \dots = N_K \quad \text{e} \quad N = \sum_{i=1}^K N_i.$$

A distribuição uniforme dos clientes pelos múltiplos *sites* representa o pior caso de distribuição. Quanto mais tendenciosa ela for, maior será a aglomeração de clientes em um único *site* e menor será a quantidade de fluxos gerados pelo compartilhamento entre clientes de *sites* diferentes. Quanto mais distribuídos os clientes estiverem, mais fluxos de rede serão criados para compartilhar os fluxos criados pelo servidor.

A figura 1 mostra o desenho esquemático da topologia *Fan-out-K*. Nesta topologia o servidor conecta-se diretamente a todos os  $K$  *sites* através de um único salto. Um fluxo que é compartilhado no servidor somente o será na rede se os seus clientes estiverem localizados no mesmo *site*, caso contrário, um novo fluxo de rede será criado para permitir o compartilhamento no servidor.

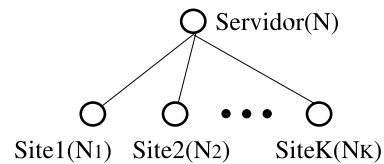


Figura 1: Topologia *Fan-out-K*

A figura 2 mostra o esquema da topologia *Daisy-Chain-K*. Nesse caso, formando uma corrente, os saltos na rede unem os *sites* uns aos outros, até que o último conecta-se ao servidor. Essa topologia, apesar de ter um caminho médio entre o servidor e os clientes mais longo, apresenta maior possibilidade de compartilhamento de fluxos de rede. Um novo cliente, que compartilha um fluxo no servidor com um outro cliente qualquer, conseguirá compartilhar o mesmo fluxo na rede se o seu *site* estiver mais próximo do servidor do que o *site* do outro cliente cujo fluxo foi compartilhado. Em outras palavras, o novo cliente não criará fluxos adicionais na rede se o seu *site* estiver localizado em algum ponto do caminho na rede do fluxo compartilhado, pois assim o fluxo já estaria sendo enviado pelo *site* do novo cliente.

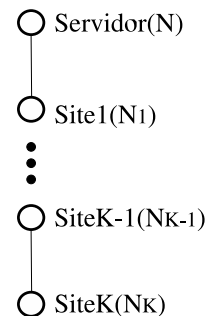
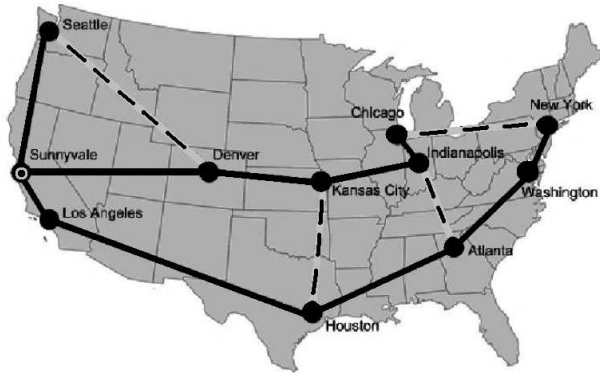
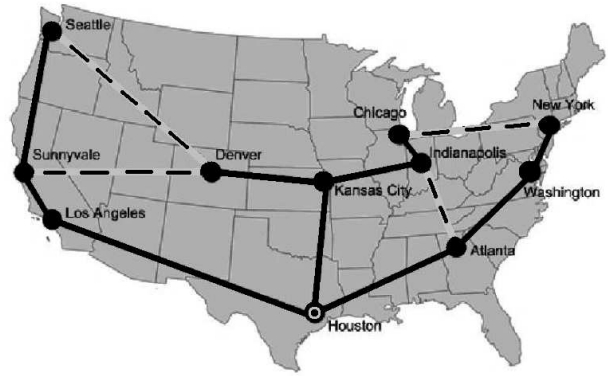


Figura 2: Topologia *Daisy-Chain-K*





(a) Árvore com o servidor em Sunnyvale



(b) Árvore com o servidor em Houston

Figura 3: Topologia do backbone da Abilene Internet2

As topologias reais usadas são derivadas do backbone da Abilene<sup>4</sup> Internet2. Essa topologia possui 11 nodos espalhados pelo território dos Estados Unidos da América. Com a finalidade de gerar duas árvores de distribuição diferentes utilizando a mesma topologia real de rede, dois nodos foram selecionados para sediar o servidor. Os nodos escolhidos foram *Sunnyvale*, figura 3(a), e *Houston*, figura 3(b). Essas figuras mostram com um traço preto mais forte as árvores de distribuição criadas assumindo o caminho mais curto. As linhas tracejadas indicam ligações da topologia que não são utilizadas para a distribuição da mídia. A avaliação de cada topologia foi feita distribuindo também uniformemente os clientes nos 10 nodos restantes. A escolha do nodo servidor privilegiou os nodos com maior vazão e que gerassem árvores de distribuição o mais diferente possível.

## 5. AVALIAÇÃO

A seção 5.1 mostra o comportamento do protocolo *Bandwidth Skimming* original em cada uma das duas topologias canônicas para carga interativa. Esses resultados são usados na avaliação das estratégias de otimização. As seções 5.2, 5.3 e 5.4 avaliam as estratégias de período ativo, período inativo e híbridas, respectivamente, comparando o seu desempenho em cada uma das 4 topologias de rede com o do protocolo *Bandwidth Skimming* original.

<sup>4</sup><http://abilene.internet2.edu> (2006)

Elas mostram que a banda de rede é afetada principalmente por dois fatores: o número de fluxos criados pelo servidor e a quantidade adicional de fluxos de rede necessários para compartilhar os fluxos de servidor por múltiplos *sites*. A seção 5.5 faz um sumário dos principais compromissos e resultados obtidos.

### 5.1 Comparando as topologias canônicas

A figura 4 mostra, em função da taxa de chegada, o impacto da variação de  $K$  na banda média de rede requerida pelo protocolo *Bandwidth Skimming* original medida em número de fluxos na topologia *Fan-out*. Quanto maior o valor de  $K$ , maior é a probabilidade de que dois clientes, compartilhando um mesmo fluxo de servidor, sejam de *sites* diferentes. Assim, para compartilhar o mesmo fluxo do servidor entre os dois *sites*, um fluxo adicional na rede necessariamente deve ser criado.

A figura 5 mostra, em função da taxa de chegada, a banda média de rede requerida pelo protocolo *Bandwidth Skimming* original medida em número de fluxos na topologia *Daisy-Chain*. Um comportamento semelhante ao da topologia *Fan-out* pode ser observado, pois quanto maior o valor de  $K$ , menor a probabilidade de que dois clientes, que compartilham um mesmo fluxo no servidor, sejam do mesmo *site*. O caminho médio do servidor até os clientes nesta topologia é maior se comparado à topologia *Fan-out*, portanto um número maior de fluxos de rede é necessário.

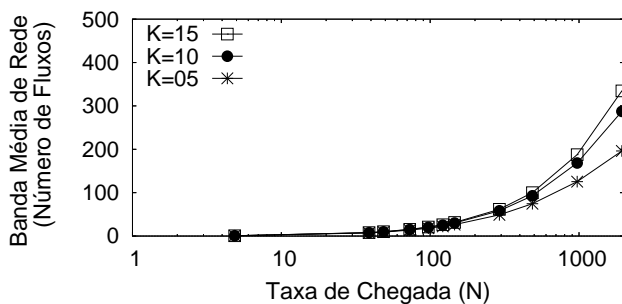


Figura 4: Banda média de rede requerida para *Bandwidth Skimming* original na topologia *Fan-out* (carga IA típica).

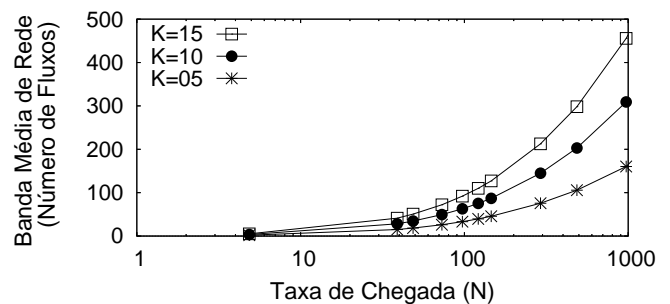


Figura 5: Banda média de rede requerida para *Bandwidth Skimming* original na topologia *Daisy-Chain* (carga IA típica).

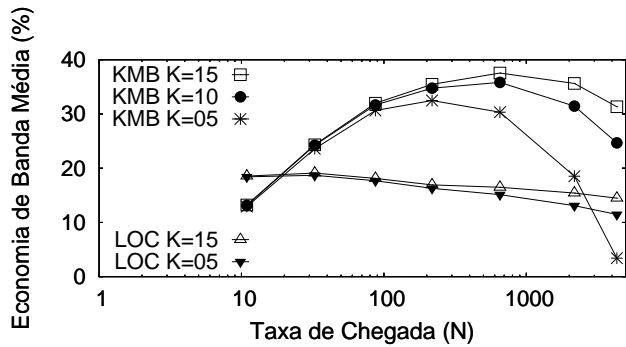


Figura 6: Economia de banda média de rede das estratégias de período ativo para a topologia *Fan-out* (carga IA).

## 5.2 KMB x LOC

Ao contrário do *Bandwidth Skimming* original (figura 4), as estratégias KMB e LOC obtêm maior economia de banda de rede com o aumento do valor de  $K$ . Essa relação é ilustrada na figura 6. Esse comportamento deve-se ao fato desses protocolos diminuírem o compartilhamento de fluxos ao utilizar o conteúdo armazenado em *buffer* e é mais evidente na estratégia KMB do que em LOC, pois KMB fragmenta mais a carga percebida pelo servidor.

A figura 7 mostra, em função da taxa de chegada  $N$ , a economia de banda média de rede das otimizações sobre o protocolo original na topologia *Daisy-Chain* com  $K = 5$ . Para todas as estratégias avaliadas nesse trabalho a maior diferença na economia de banda média foi inferior a 5%. A variação de  $K$  já não afeta significativamente o desempenho das otimizações de período ativo. O desempenho delas quanto a economia de banda média de rede segue o mesmo da banda de servidor. O compartilhamento de fluxos na rede tem pouco impacto porque um fluxo tende a visitar todos os *sites*.

A figura 8 mostra a economia de banda obtida sobre o protocolo *Bandwidth Skimming* original nas duas topologias reais de rede. Para uma carga de interatividade alta típica, economia de banda de rede de até 31% e 18% pode ser obtida com o servidor em Houston utilizando KMB e LOC, respectivamente. A diferença entre as duas topologias é mais evidente para taxas de chegada mais altas, onde a fragmentação das requisições é mais elevada. Isso ocorre porque os caminhos de rede saindo de Sunnyvale até chegar nos clientes são mais longos e existem apenas 2 conexões principais com o servidor.

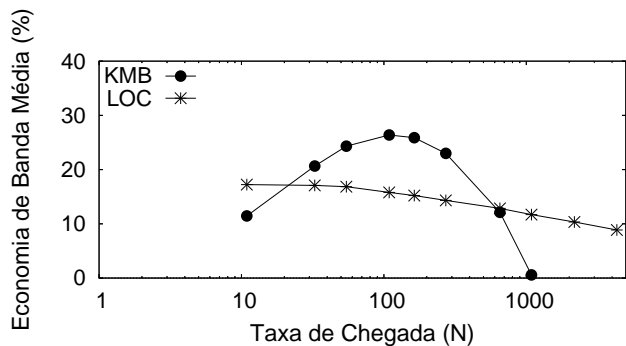


Figura 7: Economia de banda média de rede das estratégias de período ativo para a topologia *Daisy-Chain* (carga IA).

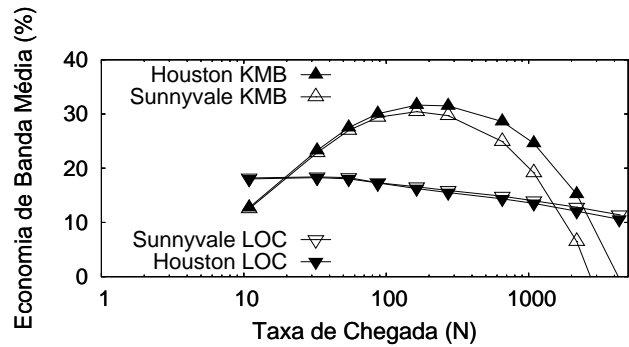


Figura 8: Economia de banda média de rede das estratégias de período ativo para as topologias reais (carga IA).

Esse caminho se assemelha mais a uma topologia *Daisy-Chain*. Já os caminhos de rede saindo de Houston são mais curtos e existem 3 principais conexões com o servidor. Este caminho se assemelha mais à topologia *Fan-out*.

## 5.3 SP x GP x CP

A figura 9 mostra, em função da taxa de chegada, a parcela do período inativo que cada estratégia de otimização gasta adiantando conteúdo. As estratégias GP e CP têm um maior consumo de banda de rede no período inativo do que SP. A partir de um certo valor de  $N$ , GP e CP escutam fluxos ininterruptamente, ou seja, utilizam 100% do período inativo adiantando conteúdo, enquanto SP ocupa a totalidade apenas para taxas de chegada mais elevadas.

As figuras 10(a-c) mostram, em função da taxa de chegada  $N$ , a economia de banda média de rede obtida pelas estratégias de período inativo para a topologia *Fan-out* e diferentes valores de  $K$ . Com o aumento do valor de  $K$ , a quantidade adicional de fluxos de rede necessários para compartilhar os fluxos de servidor aumenta. Comparando com GP e CP, SP apresenta melhor desempenho porque além de manter economia expressiva de banda média de servidor, utiliza menos o período inativo, conforme visto na figura 9.

A figura 11 mostra, em função da taxa de chegada  $N$ , a economia de banda média de rede sobre o protocolo *Bandwidth Skimming* original para a topologia *Daisy-Chain* com  $K = 5$ . Na topologia *Daisy-Chain-K*, a variação de  $K$  também não afeta significativamente o desempenho das otimizações.

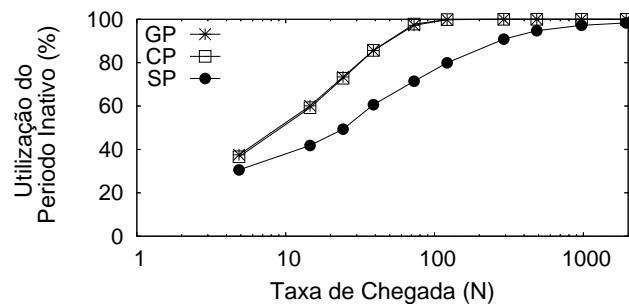


Figura 9: Utilização do período inativo para diferentes taxas de chegada (carga IA típica).

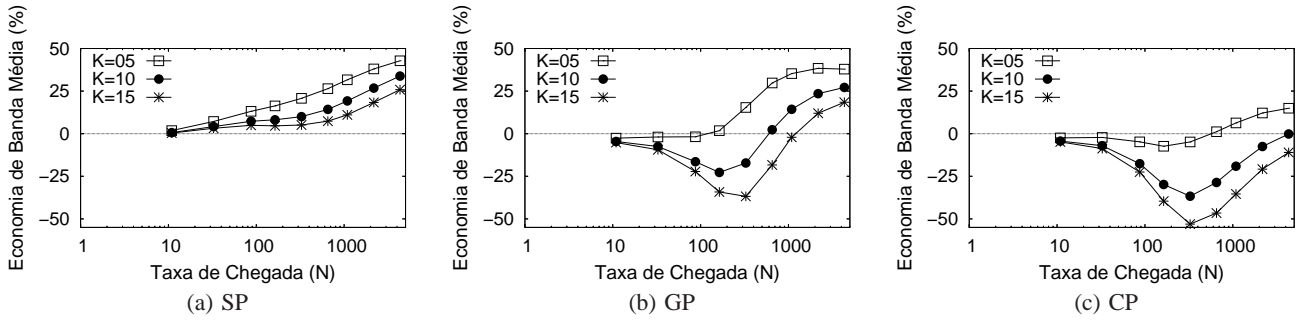


Figura 10: Economia de banda média de rede das estratégias de período inativo para a topologia *Fan-out* (carga IA).

As figuras 12(a) e 12(b) e as figuras 13(a) e 13(b) mostram, em função da taxa de chegada  $N$ , a economia de banda média de rede obtida sobre o protocolo *Bandwidth Skimming* original nas duas topologias reais de rede e para duas cargas educacionais de interatividade alta. A diferença entre essas duas cargas é que a do MANIC possui a razão entre os períodos ativo e inativo menor do que a carga do eTeach. Dessa forma, a carga do MANIC apresenta maiores oportunidades de adiantar conteúdo durante o período inativo e as estratégias GP e CP têm pior desempenho (figura 9). As figuras mostram a diferença na economia de banda média obtida variando tanto a posição do servidor (comparando 12(a) com 12(b) e 13(a) com 13(b)) quanto proporção entre os períodos ativo e inativo (comparando 12(a) com 13(a) e 12(b) com 13(b)).

SP consegue índices de economia de banda média de rede mais altos porque, além de utilizar totalmente o período inativo apenas em taxas de chegada mais altas, se beneficia do fato de selecionar dois fluxos consecutivos, ou seja, ao escolher um fluxo, também escuta o seu alvo. Essa seleção permite que o conteúdo adiantado seja mais contíguo, fragmentando menos as futuras requisições. Outro fator é que SP preserva o padrão de acesso do protocolo original, minimizando a sobrecarga de rede caso o fluxo principal já esteja sendo enviado ao *site* do cliente. Para as topologias reais, SP apresenta economia de banda média de servidor e de rede em relação protocolo original de até 46% e 42%, respectivamente.

Apesar de GP ser mais agressivo ao povoar o *buffer* e evitar mais retransmissão, ela impõe uma carga maior na rede do que SP. Assim, a economia de banda média de rede de GP torna-se mais restrita. A economia de banda de rede com CP é reduzida porque ela falha em economizar banda de servidor.

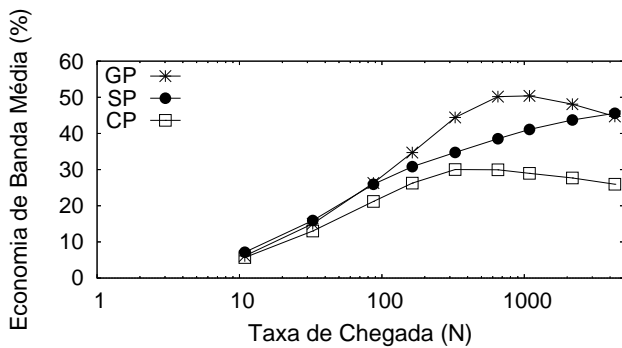


Figura 11: Economia de banda média de rede das estratégias de período inativo para a topologia *Daisy-Chain* (carga IA).

Um rápido crescimento na economia de banda pode ser observado nas figuras 13(a) e 13(b) próximo de  $N = 100$ . Quando a taxa de chegadas é baixa, GP e CP não evitam retransmissão suficiente para compensar os fluxos adicionais de rede. A partir do ponto em que as estratégias utilizam todo o período inativo, a quantidade de fluxos extras na rede atinge o seu máximo (figura 9). Como a economia obtida com a redução de retransmissão continua crescendo, o resultado é o crescimento observado. Para as topologias reais, GP apresenta economia de banda média de servidor e de rede em relação protocolo original de até 55% e 40%, respectivamente.

#### 5.4 HA x HS x HC

Conforme [17, 18, 19], ao combinar as estratégias de otimização busca-se economia de banda média de servidor significativa por uma faixa mais ampla de configurações que qualquer estratégia individualmente. Nesta seção é feita a avaliação nas diferentes topologias dos requisitos de banda de rede destas estratégias híbridas, que conseguem economia de banda média de até 66% no servidor.

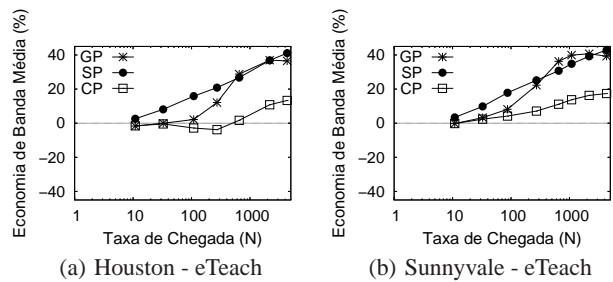


Figura 12: Economia de banda média de rede das estratégias de período inativo para as topologias reais (carga eTeach).

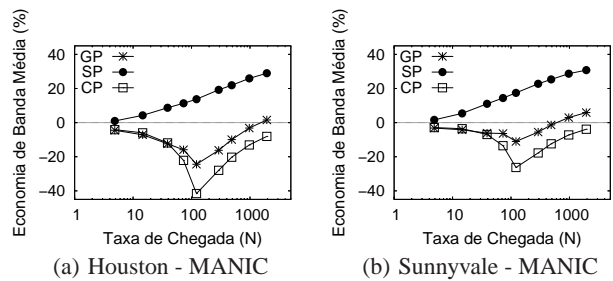


Figura 13: Economia de banda média de rede das estratégias de período inativo para as topologias reais (carga MANIC).

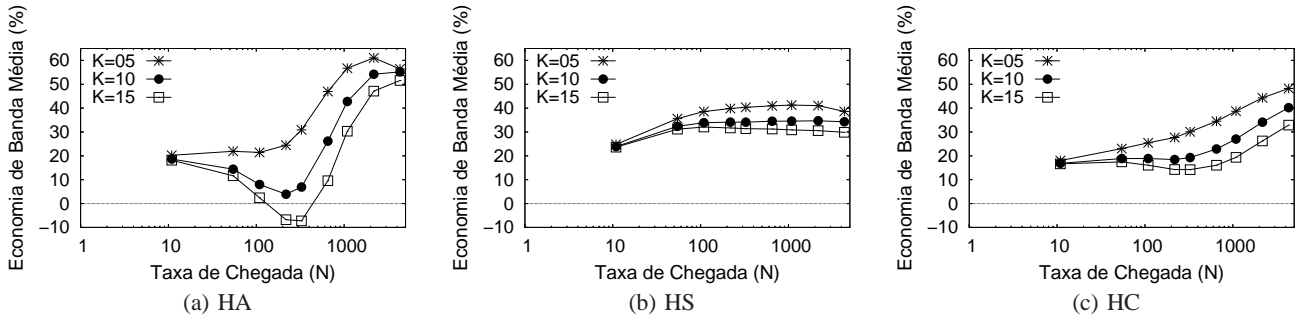


Figura 14: Economia de banda média de rede das estratégias híbridas para a topologia *Fan-out* (carga IA).

As figuras 14(a-c) mostram a economia de banda média de rede das estratégias híbridas para a topologia *Fan-out*, em função da taxa de chegada e variando o valor de  $K$ . A figura mostra os resultados para a carga do eTeach, sendo que para a carga do MANIC, bem como as demais cargas de interatividade alta, os resultados são qualitativamente semelhantes. Assim como na figura 10, a quantidade adicional de fluxos de rede necessários para compartilhar os fluxos de servidor aumenta com  $K$ .

O crescimento rápido na economia de banda média de rede também pode ser observado por HA na figura 14(a) em função da ação conjunta da expressiva economia de banda no servidor e do fato de GP ter alcançado o aproveitamento total do período inativo escutando fluxos (figura 13).

Apesar de HS apresentar um pico de economia menor do que HA, HS obtém economia de banda média de rede significativa por uma faixa maior de taxas de chegada. Essa economia é obtida porque SP utiliza menos o período inativo do que GP. Já HC obtém menos economia de banda média de rede para baixas taxas de chegada, se comparada a HS, porque HC não incorpora KMB, que é eficiente para adiantar conteúdo e evitar retransmissão quando  $N$  é baixo.

As figuras 15(a) e 15(b) mostram, em função da taxa de chegada  $N$ , a economia de banda média de rede das estratégias híbridas sobre o protocolo *Bandwidth Skimming* original para a topologia *Daisy-Chain* com  $K = 5$ . Nesta topologia, a variação de  $K$  novamente não afeta significativamente o desempenho das otimizações.

As figuras 16(a) e 16(b) e as figuras 17(a) e 17(b) mostram a economia de banda média de rede obtida pelas estratégias híbridas sobre o protocolo *Bandwidth Skimming* original nas duas topologias reais de rede para as cargas do eTeach e do MANIC, respectivamente. O

mesmo crescimento abrupto na economia de banda média de rede de HA pode ser observado.

HA apresenta economias de banda média de rede mais expressivas apenas para cargas onde o número de requisição por sessão é alto e o período inativo é menor do que o ativo. Assim, o usuário tem mais oportunidade de utilizar o conteúdo armazenado e gera menos fluxos adicionais na rede.

O desempenho das estratégias HS e HC se mostrou qualitativamente similar quando a carga é variada. A presença de KMB em HS faz com que HS consiga economizar significativamente a banda de rede para baixas taxas de chegada, mas diminui a eficiência de HS quando  $N$  é alto (figura 8). HC, que é menos suscetível à fragmentação das requisições, obtém economia média de banda de rede até taxas de  $N$  mais altas. Para as topologias reais, HA, HS e HC economizam a banda média de rede em até 60%, 43% e 48%, respectivamente.

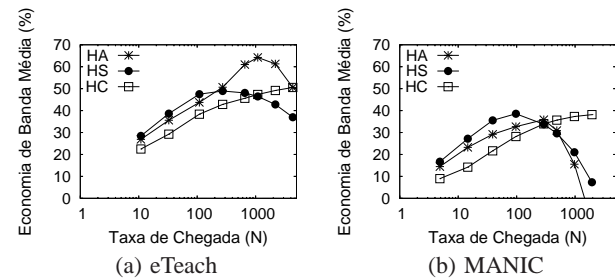


Figura 15: Economia de banda média de rede das estratégias híbridas para a topologia *Daisy-Chain* (duas cargas IA).

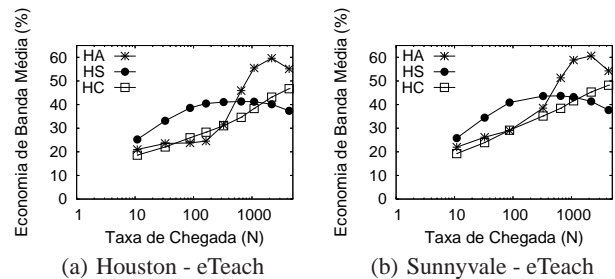


Figura 16: Economia de banda média de rede das estratégias híbridas para as topologias reais (duas cargas IA).

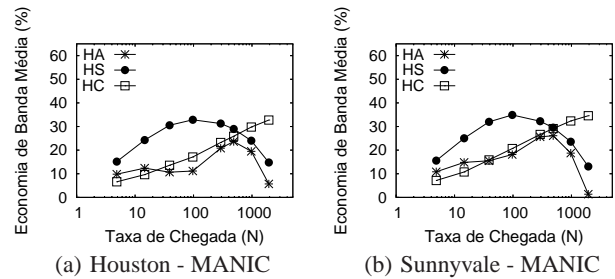


Figura 17: Economia de banda média de rede das estratégias híbridas para as topologias reais (duas cargas IA).



## 5.5 Sumário

As seções anteriores apresentaram o desempenho de cada estratégia de otimização quanto aos requisitos de banda média de rede para cada uma das 4 diferentes topologias de rede utilizadas. Foi mostrado que o desempenho de cada estratégia de otimização é afetado principalmente por dois fatores: o número de fluxos criado pelo servidor e a quantidade adicional de fluxos de rede necessária para compartilhar os fluxos de servidor por múltiplos *sites*. A banda de servidor, por sua vez, depende principalmente da redução de retransmissão e da fragmentação das requisições. Quanto mais conteúdo é armazenado no *buffer*, maiores as oportunidades de utilizá-lo e evitar retransmissão. Quanto mais contíguo esse conteúdo estiver, menor a fragmentação imposta na carga.

As estratégias KMB e LOC têm melhor desempenho na topologia *Fan-out* quando  $K$  aumenta porque elas diminuem o compartilhamento de fluxos ao utilizar o conteúdo armazenado em *buffer*. Assim a rede pode se beneficiar das durações reduzidas dos fluxos resultantes (figura 6). Para as demais estratégias o aumento de  $K$  prejudica o compartilhamento de fluxos e compromete seus desempenhos.

Nenhuma estratégia de otimização apresenta variação significativa (inferior a 5%) de desempenho na topologia *Daisy-Chain* quando o valor de  $K$  é variado. O desempenho das estratégias quanto a economia de banda média de rede segue o mesmo da banda de servidor. O compartilhamento de fluxos na rede tem pouco impacto porque um fluxo tende a visitar todos os *sites*.

A figura 18 mostra, em função da taxa de chegada, a economia de banda média de rede da estratégia *Hybrid Aggressive* sobre o protocolo *Bandwidth Skimming* original nas duas topologias de rede reais do backbone da Abilene Internet2 utilizadas. Os resultados dos requisitos de banda média para as duas topologias são qualitativamente semelhantes. Entretanto, na topologia com a localização do servidor em Sunnyvale, as estratégias conseguem um desempenho geral melhor do que com a localização do servidor em Houston. O desempenho das estratégias com o servidor em Sunnyvale chega a ser superior em até 76%. Isso ocorre porque os caminhos de rede saindo de Sunnyvale até chegar nos clientes se assemelham mais a uma topologia *Daisy-Chain*. Já os caminhos de rede saindo de Houston se assemelham mais à topologia *Fan-out*. Essa diferença mostra como a escolha do local para sediar o servidor deve ser cuidadosa, pois árvores de distribuição diferentes podem incorrer em requisitos de banda de rede significativamente diferentes.

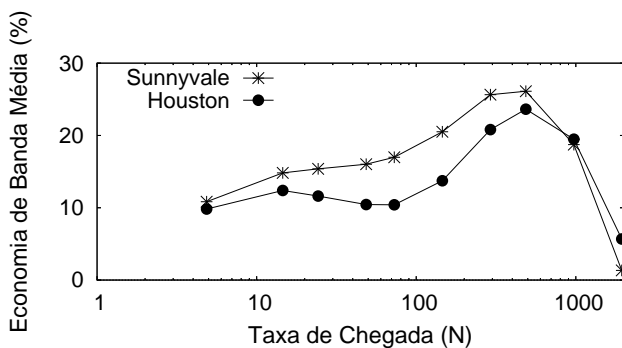


Figura 18: Diferença na economia de banda média de rede para as duas topologias reais (estratégia HA e carga IA típica).

Foi mostrado o impacto da proporção entre os períodos ativo e inativo na redução de retransmissão das estratégias que exploram períodos inativos (figuras 12 e 13). Quanto maior o período ativo em relação ao inativo, maiores as oportunidades das estratégias para escutar fluxos ativos, armazenar conteúdo e evitar retransmissão. Entretanto, quanto maior a utilização do período inativo, maior a carga extra imposta na rede. Portanto, existe um compromisso. Para economizar banda de rede o cliente deve ter sido capaz de povoar o *buffer* o suficiente para que os fluxos adicionais na rede sejam compensados evitando retransmissões no servidor.

Entre as estratégias que exploram os períodos inativos, SP apresenta a melhor relação entre a economia de banda média de servidor e a de rede (figuras 12 e 13) porque SP tem um desempenho na rede melhor do que GP e a economia de banda média de servidor de SP ainda é bastante expressiva (até 46%). Entre as estratégias híbridas, apesar de HA ter melhor desempenho no servidor [19], sua aplicação na rede é limitada a cargas com alta taxa de chegada (com economia de até 60%). HS, que incorpora KMB, obtém expressiva economia de banda média de rede para baixas taxas de chegada, apesar de sofrer com o efeito da fragmentação quando  $N$  aumenta. HC, menos suscetível à fragmentação, consegue sustentar significativa economia de banda média de rede até taxas de chegada mais altas. HS e HC apresentam as melhores relações entre a economia de banda média de servidor e rede simultaneamente por faixa mais ampla de taxa de chegada. São obtidas com HS e HC economias de até 54% e 52% no servidor [19] e de até 43% e 48% na rede (figuras 16 e 17), respectivamente.

## 6. CONCLUSÕES E TRABALHOS FUTUROS

Este trabalho investigou, através de cargas interativas realistas e topologias de rede reais e canônicas, os requisitos de banda de rede para estratégias de otimização do protocolo de distribuição de mídia contínua sob demanda *Bandwidth Skimming*.

Trabalhos prévios indicaram que, utilizando o protocolo *Bandwidth Skimming* para acesso seqüencial, é possível obter grandes economias simultâneas de banda de servidor e de rede. No entanto, caracterizações de acesso de usuários de servidores reais evidenciam forte presença de interatividade, comprometendo a economia de banda de servidor obtida com o protocolo *Bandwidth Skimming*. Estratégias de otimização especializadas para acesso interativo foram previamente propostas para minimizar esse impacto. Elas são capazes de economizar em até 66% a banda média de servidor se comparadas com o protocolo original. Entretanto uma das classes de otimizações mais promissoras pode impor uma carga extra na rede, tornando um ponto crítico das estratégias que até então não tinha sido avaliado.

Os resultados obtidos neste trabalho utilizando as estratégias de otimização indicam que é possível obter simultaneamente grandes economias de banda de servidor e de rede. HYBRID AGGRESSIVE é a estratégia que consegue a maior economia de banda média de servidor, até 66% para uma carga altamente interativa. No entanto, seu desempenho na rede é limitado a cargas com alta taxa de chegada (economia de até 60%) porque GP escuta uma quantidade maior de fluxos durante os períodos inativos. Deve-se preferir esta estratégia em cenários onde o recurso com maior restrição seja o servidor e não a rede, como em uma rede local.

HYBRID SYNERGIC obtém expressiva economia de banda média de rede para baixas taxas de chegada. Quando a taxa é alta, a

fragmentação das requisições imposta por KMB degrada a economia de banda média de rede. Essa estratégia deve ser usada em aplicações onde a fragmentação tem menor impacto, como em vídeos educacionais altamente interativos. HYBRID CONSERVATIVE obtém economia de banda média de rede menor do que HS para baixas taxas de chegada, no entanto consegue sustentá-la até taxas mais elevadas. Por ser uma estratégia que fragmenta menos a carga do que HS, ela deve ter uso preferencial em aplicações com um caráter mais sequencial, como áudio, ou em aplicações em que se conheça pouco sobre a carga. HS e HC, respectivamente, conseguem economizar a banda média de servidor em até 54% e 52% e de rede em até 43% e 48% por faixa mais ampla de taxa de chegada.

Trabalhos futuros incluem o desenvolvimento de um protocolo adaptativo que monitora a carga interativa recebida pelo servidor e, baseado apenas nos parâmetros da carga, escolhe dinamicamente qual ou quais estratégias de otimização utilizar, caso haja necessidade de alguma. Dessa forma, pretende-se aproveitar o melhor de cada estratégia para obter economias expressivas de banda média de servidor por uma faixa ainda mais ampla de configurações do que qualquer estratégia sozinha. Uma análise dos requisitos de banda de rede para esse protocolo adaptativo também é deixado como trabalho futuro.

## 7. REFERÊNCIAS

- [1] J. Almeida, J. Krueger, D. Eager, and M. Vernon. Analysis of Educational Media Server Workloads. In *Proc. NOSSDAV*, Port Jefferson, NY, June 2001.
- [2] A. Bar-Noy and R. Ladner. Competitive On-Line Stream Merging Algorithms for Media-on-Demand. *Journal of Algorithms*, 48(1):59–90, Aug. 2003.
- [3] M. Bradshaw, B. Wang, S. Sen, L. Gao, J. Kurose, P. Shenoy, and D. Towsley. Periodic Broadcast and Patching Services-Implementation, Measurement and Analysis in an Internet Streaming Video Testbed. *Multimedia Systems*, 9(1):78–93, July 2003.
- [4] M. Chesire, A. Wolman, G. Voelker, and H. Levy. Measurement and Analysis of a Streaming Media Workload. In *Proc. Symp. on Internet Technologies and Systems*, San Francisco, CA, Mar. 2001.
- [5] E. Coffman, P. Momcilovic, and P. Jelenkovic. The Dyadic Stream Merging Algorithm. *Journal of Algorithms*, 43(1):120–137, Apr. 2002.
- [6] C. Costa, Í. Cunha, A. Borges, C. Ramos, M. Rocha, J. Almeida, and B. Ribeiro-Neto. Analyzing Client Interactivity in Streaming Media. In *Proc. World Wide Web Conference*, New York, NY, May 2004.
- [7] D. Eager and M. Vernon. Dynamic Skyscraper Broadcasts for Video-on-Demand. In *Proc. International Workshop on Advances in Multimedia Information Systems*, Istanbul, Turkey, Sep. 1998.
- [8] D. Eager, M. Vernon, and J. Zahorjan. Bandwidth Skimming: A Technique for Cost-Effective Video on Demand. In *Proc. Multimedia Computing and Networking*, San Jose, CA, Jan. 2000.
- [9] D. Eager, M. Vernon, and J. Zahorjan. Minimizing Bandwidth Requirements for On-Demand Data Delivery. *IEEE Trans. on Knowledge and Data Engineering*, 13(5):742–757, 2001.
- [10] L. Gao and D. Towsley. Supplying Instantaneous Video-on-Demand Services Using Controlled Multicast. In *Proc. IEEE Multimedia Computing Systems*, Florence, Italy, June 1999.
- [11] K. Hua, Y. Cai, and S. Sheu. Patching: A Multicast Technique for True Video-on-Demand Services. In *Proc. ACM MULTIMEDIA*, Bristol, UK, Sep. 1998.
- [12] S. Jin and A. Bestavros. Scalability of Multicast Delivery for Non-sequential Streaming Access. In *Proc. SIGMETRICS*, Marina Del Rey, CA, June 2002.
- [13] J. Liu and J. Xu. Proxy Caching for Media Streaming over the Internet. *IEEE Communications Magazine*, 42(8):88–94, Aug. 2004.
- [14] H. Ma, K. Shin, and W. Wu. Best-effort Patching for Multicast True VoD Service. *Multimedia Tools Applications*, 26(1):101–122, May 2005.
- [15] B. Netto. Patching Interativo: Um Novo Método de Compartilhamento de Recursos para Transmissão de Vídeo com Alta Interatividade. Master's thesis, UFRJ, Rio de Janeiro, RJ, Brasil, 2004.
- [16] J. Padhye and J. Kurose. An Empirical Study of Client Interactions with a Continuous-Media Courseware Server. In *Proc. NOSSDAV*, Cambridge, UK, July 1998.
- [17] M. Rocha, M. Maia, J. Almeida, and S. Campos. Escalabilidade de Protocolos com Compartilhamento de Banda para Cargas de Mídia Contínua Realistas. In *Proc. SBC WPerformance*, São Leopoldo, RS, Brasil, Jul. 2005.
- [18] M. Rocha, M. Maia, Í. Cunha, J. Almeida, and S. Campos. Scalable Media Streaming to Interactive Users. In *Proc. ACM MULTIMEDIA*, Singapore, Singapore, Nov. 2005.
- [19] M. Rocha, M. Maia, Í. Cunha, J. Almeida, and S. Campos. Estratégias Híbridias para Transmissão de Mídia Contínua Interativa com Compartilhamento de Fluxo. In *Proc. SBRC*, Curitiba, PR, Brasil, Jun. 2006.
- [20] C. Rodrigues and R. Leão. Novas Técnicas de Compartilhamento de Banda para Servidores de Vídeo Sob Demanda Com Interatividade. In *Proc. Simpósio Brasileiro de Redes de Computadores*, Fortaleza, CE, Brasil, Mai. 2005.
- [21] H. Tan, D. Eager, and M. Vernon. Delimiting the Range of Effectiveness of Scalable On-Demand Streaming. In *Proc. International Symposium on Computer Performance Modeling and Evaluation*, Rome, Italy, Sep. 2002.
- [22] B. Vickers, C. Albuquerque, and T. Suda. Source-adaptive Multilayered Multicast Algorithms for Real-time Video Distribution. *IEEE/ACM Transactions on Networking*, 8(6):720–733, Dec. 2000.
- [23] Y. Zhao, D. Eager, and M. Vernon. Network Bandwidth Requirements for Scalable On-Demand Streaming. In *Proc. IEEE INFOCOM*, New York, NY, June 2002.