

WWW キャッシュサーバのヒット率予想について

1F-2

鍋島 公章

NTTソフトウェア研究所

1. はじめに

WWW のトラフィック抑制のために WWW キャッシュサーバが広く使われはじめている。また、キャッシュの性質上、より多くのリクエストを扱うほど、ヒット率が高くなる。そのため、ISP のバックボーンのようなトラフィックの多い所でキャッシュサーバを運用することが望まれる。しかし、現状では、トラフィックが多い所では、ほとんどキャッシュサーバは運用されていない。この一つの理由は、処理すべきトラフィックが多くなるほど、キャッシュサーバへの投資が大きくなるが、キャッシュサーバを導入して得られるトラフィック抑制率(回線費の節約額)を予想する方法がないことである。

本稿では、ヒット率(トラフィック抑制率)予想方法の一つとして、パケットモニタリングを使用する方法を提案し、その評価実験と結果について報告する。

2. ヒット率予想

今回提案するヒット率予想では、パケットモニタリングを利用する。つまり、FDDI やイーサネットのようなマルチアクセス媒体上を流れているパケットから HTTP リクエスト部分を抜き出す。そして、その抜き出したリクエストを仮想的なキャッシュサーバに投入した時のヒット率をシミュレーションする方法をとる。

3. 実験

97年9月1日からの5日間(平日)に、IMnet のバックボーンの一つである大手町 NOC の FDDI リング上でパケットのモニタリングを行なった。今回の実験では、まず、リクエストをディスクに保存し、後日、ヒット率をシミュレーションする方法をとった。リクエストのモニタリング処理には、pcap ライブラリを使った専用のプログラムを使用した。これは、パケットの中でディスティネーションポートが 80 番(HTTP のデフォルトポート)で

あり TCP ペイロード部分が GET という文字列で始まるものをディスクに保存する。

全モニタパケット数は約 5.3 億個、モニタリング時のパケットドロップ率は平均 0.36%であった。この中の HTTP リクエスト数は 719 万リクエスト(1 日平均約 144 万リクエスト、ピークでは約 30 リクエスト/秒)であった。また、リクエスト中に出現した URL(CGI ページは除く)は 203 万個であった。CGI ページ(通常はキャッシュされるべきではない)へのリクエストは全体の 7.3%である52万リクエストであった¹。

4. 実験結果

今回の 5 日間にわたるモニタリングにおける、同一 URL の出現回数の頻度分布をグラフ 1 に示す。また、仮想キャッシュにおけるヒット率の変移をグラフ 2 に示す。ここでの仮想的なキャッシュは、サイズ無限、キャッシュの中身は古くならない、という理想的な環境を仮定している(これについては 5 節において考察する)。

1 日目で、60%以上のヒット率となっており、それ以降は、あまりヒット率は上がっていない。そのため、このバックボーン上のキャッシュサーバでは、1 日程度のトラフィックの保存で十分であると予想できる。また、この時のキャッシュ中のオブジェクト数は約 50 万個であり、これに必要なディスク容量は、7GB 程度となると予想できる²。

早朝のヒット率低下は、オートパイロットプログラム(通常ネットワークが空いているこの時間帯に実行されることが多い)が影響しているものと思われる。

An Estimation Method for WWW Cache Hit ratio
Masaaki NABESHIMA (nabe@slab.ntt.co.jp)
NTT Software Laboratories

本研究は、科学技術庁の平成 9 年度科学技術振興調整費による「生活工学アプリケーション研究」の一環として行われている。

¹ クライアント側では、リクエストしたページが CGI により生成されたかどうかは、厳密には判定できない。しかし、CGI で生成されるページの多くは、“CGI”もしくは“?”という文字列を URL に含むため、これらの文字列を含む URL を持つページを CGI で生成されたものとした。

² キャッシュサーバにおける平均オブジェクトサイズは 17KB 程度(cache.imnet.ad.jp におけるデータ)である。ただし、これには FTP オブジェクトも含まれているため、HTTP オブジェクトのみだと、これよりも小さくなる。

5. 考察

今回のヒット率予想方法は、いくつかの面で不十分である。これらの実際のヒット率との違いを生む要因について考察する。

5.1 ヒット率を大きく予想する要因

5.1.1 古くなったデータ

実際のキャッシュサーバでは、キャッシュ中のデータは、時間が経つに連れて古くなる。そして、IMS などの整合性照会リクエストを受けた時には MISS 扱い(新しい情報を WWW サーバから取り出す)となる。

5.1.2 隠れ CGI

今回は URL 中に“CGI”もしくは“?”を含むページを CGI ページとしている。しかし、実際にはこれらの文字列を URL 中に持たない CGI ページが存在する。

5.1.3 リクエストの方向

ISP から出て行くリクエストと、入ってくるリクエストを比較すると、一般に、出て行くリクエストのヒット率が低いと考えられる。今回の見積りの対象は ISP から出て行くリクエストであるが、これらを区別していない。

5.2 ヒット率を低く予想する要因

5.2.1 ラウンドロビンサーバ

リクエスト数の多い人気サーバの多くは、複数のミラーサーバを用意し、DNS ラウンドロビンを使ってリクエストを分散させている。つまり、一つのホスト名に対して複数の IP アドレスが割振られている。一方、パケットモニタリングによる HTTP リクエストの取得では、URL 中のホスト部分が IP アドレス形式となる。このため、IP アドレス形式では異なる URL へのリクエストが、ホスト形式では同一 URL へのリクエストである可能性がある。

5.2.2 キャッシュサーバからのリクエスト

今回の計測には、同じ FDDI リング上にある 1 台のキャッシュサーバからの 30 万リクエスト/日が含まれている。キャッシュサーバからのリクエストは、既に多くのリクエストがキャッシュサーバでヒットしている。そのため、ここから出るリクエストはヒットしにくい。

6. おわりに

現状の見積り方法は精度の荒いものである。そのため、実験により得られたヒット率は、実際のものよりもかなり高いと思われる。今後は、5 節の考察項目について詳細を詰め、他のヒット率予想方法[1]を参考にしながら、見積り方法の精度を上げる予定である。また、HTTP リクエストだけでなく、サーバからのコンテンツを含むパケットをモニタリングし、オブジェクトの平均更新

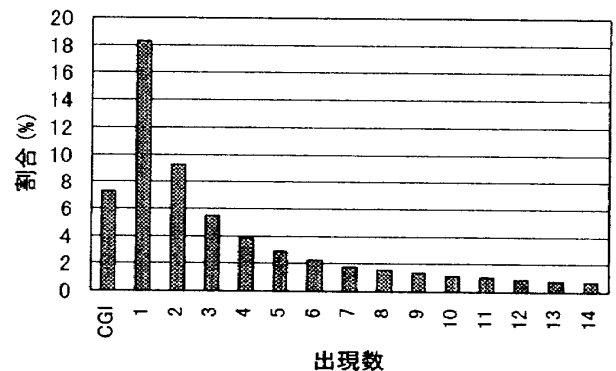
間隔や、平均オブジェクトサイズのデータを集め、それを見積り方法に反映させる予定である。そして、もっと流量の多いネットワークや、リクエストの偏差が少ない海外リンクにおいて見積り方法の適用実験を行なう予定である。

謝辞

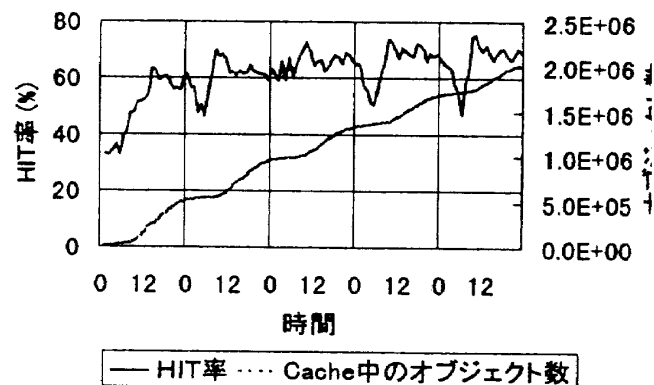
有用なコメントをいただいた後藤教授をはじめとする早稲田大学後藤研究室の方々に感謝します。

参考文献

- [1] Bradley M. Duska, David Marwood and Michael J. Feeley, “The measured access characteristics of World-Wide-Web client proxy caches”, USENIX Symposium on Internet Technologies and Systems, December 1997



グラフ 1: URL 出現回数の分布



グラフ 2: ヒット率の変移