

ウェブクローラ向けの効率的な重複URL 検出手法

Efficient Duplicated URL Detection for Web Crawlers

久保田 展行^{*} 上田 高德^{*}
山名 早人^{*}

Nobuyuki KUBOTA Takanori UEDA
Hayato YAMANA

ウェブクローラにおいて、冗長なクロールを避けるための重複 URL 検出は必須の処理である。重複 URL 検出処理には、URL のハッシュ値を求め、その集合を元に重複検出を行う手法が一般的に用いられている。しかし、インターネット上には数億の URL が存在するため、最終的には100GB 以上のハッシュ値を管理する必要があり、メモリのみで処理することは困難となる。本稿では、新たな乱択データ構造である Stable Bloom Filter Light を提案し、ストレージ上でのハッシュ値集合操作回数を削減するための近似的なキャッシング手法を実現する。実験の結果、既存手法と比較して、同等のキャッシュヒット率のときにメモリ使用量を約 63%削減した。

It is necessary for Web crawlers to avoid redundant crawling by detecting duplicated URLs. Web crawlers usually use the set of hash values generated from URLs to detect duplicates. However, the size of hash values will be over 100GB because the Web has billions of URLs. This large amount of hash values disables Web crawlers to detect duplicated URLs without storage accesses. In this paper, we present a novel randomized data structure called Stable Bloom Filter Light and develop a method using Stable Bloom Filter Light to reduce storage accesses. Experimental results show that our method cuts 63% of the memory size compared to existing methods for the same cache hit rate.

1. はじめに

ウェブクローラにとって、クロール済みのウェブページを再クロールするような、冗長なクロールは致命的である。同一ウェブページに複数回アクセスすることは、クロール対象のウェブサーバに不必要な負荷をかけるだけでなく、ウェブクローラ自体の性能も低下させる。同一ウェブページに対して短時間で複数回クロールしないようにするために、ウェブクローラはクロール対象のウェブページがクロール済みであるかどうかを判断する必要がある。ウェブクローラでは、クロール済みのウェブページのURLを集合として

持ち、クロール対象のウェブページにアクセスする前に、URL集合に対してアクセスしようとしているURLの包含判定を行うことで、そのウェブページへアクセスすべきか判断する。本稿では、URL集合への包含判定処理を、重複URL検出と呼ぶ。

通常、重複URL検出をするにあたり、まず、URLは数バイトのハッシュ値へと変換される。我々の調査によるとURLの平均長は約70バイトであるため、小さなハッシュ値へ変換することで、空間使用量を削減することが可能となる。次に、収集が完了したURLのハッシュ値は重複検出用の集合に入れられる。この集合に、アクセスする予定であるURLのハッシュ値を含んでいるかどうかを問い合わせることで、重複URL検出が行われる。

URLのハッシュ値集合を操作するにあたり問題となるのは、URLの数が膨大なことである。ウェブ上には、数億のURLが存在することがわかっている[10]。そのため、64ビットのハッシュ値を用いた場合でも、最終的なハッシュ値集合の容量は100GBを超えることとなる。100GBを超えるようなハッシュ値集合をメモリ上で扱うには、ウェブクローラを複数台のマシンで並列実行し、マシン毎にクロールするドメインを限定することで、マシン1台あたりのメモリ使用量を抑えるという手段が考えられる。しかし、クロール中に使用可能なマシンの台数が限られている場合は、メモリ上のみでハッシュ値操作を行うのは困難である。また、マシンの台数にゆとりがある場合も、重複URL検出を少ないメモリ使用量で処理することができれば、100GB以上のメモリ空間の節約となるため、得られる恩恵は大きい。そのため、本稿では使用可能なメモリが限られている場合でも効率的に動作する重複URL検出処理を取り扱う。

少ないメモリ使用量で重複URL検出を行うためには、ストレージ上でハッシュ値集合を管理する必要がある。ただし、ストレージ上でのハッシュ値集合操作は、一般的にメモリ上での操作よりも処理コストが大きい。よって、ストレージ上でのハッシュ値集合操作、すなわちストレージアクセスを減らすことが重要となる。

URLのハッシュ値集合操作によって発生するストレージアクセス回数を削減するための解決策として、URLをキャッシュする方法がある[2]。クロール中に出現するURLをキャッシュしておくことで、高い確率でキャッシュヒットすることが知られている。www.yahoo.comなどの頻出URLがキャッシュに載りやすいことや、相互リンクなどの、共起関係にあるURLの組み合わせが多いことが、高いキャッシュヒット率の理由とされている。例えば、LRUなどの一般的なキャッシング手法を使用することで、約80%のキャッシュヒット率を実現可能である。そのため、キャッシュをメモリ上に確保することで、ストレージ上にあるURLのハッシュ値集合に対する問い合わせ回数と、ハッシュ値集合の更新回数を大幅に減らすことが可能となる。

URLをキャッシュするときに問題となるのは、キャッシュの更新と、要素の包含判定の操作の計算量である。例えばLRUキャッシュの場合は $O(\log N)$ の計算量が必要である。そのため、キャッシュサイズを大きくすると、ヒット率は向上するが、キャッシュの操作にかかる計算時間が増加する。その結果、一定のキャッシュサイズを超えると、ストレージアクセス以上の計算時間がかかり、キャッシュを利用する利点がなくなる[6]。しかし、キャッシュの操作コストを $O(1)$ にすることができれば、計算時間をほとんど変えることなくキャッシュサイズを増加させることが可能となり、キャッシュヒット率を向上させることができる。 $O(1)$ の計算量でキャッシュ操作を実現するためには、Bloom Filter[3]を利用する近似的な手法が考えられる。Bloom Filterとは、集合の包含判定を行う際に、含まれていない要素が含まれていると誤判定するfalse positiveの発生を許容する代わりに空間効率を改善

^{*} 非会員 株式会社 Preferred Infrastructure
nobu@preferred.jp

^{*} 学生会員 早稲田大学大学院 基幹理工学研究所, 早稲田大学 メディアネットワークセンター
ueda@yama.info.waseda.ac.jp

^{*} 正会員 早稲田大学理工学術院, 国立情報学研究所
yamana@yama.info.waseda.ac.jp

した乱択(ランダムイズド)データ構造である。

本稿では、新しい乱択データ構造であるStable Bloom Filter Lightを提案し、従来から提案されているURLをキャッシュする手法に加えてStable Bloom Filter Lightを用いることで、近似的な重複URL検出を行う手法を実現する。提案手法を用いることにより、重複URL検出処理にLRUキャッシュのみを用いた場合と比較して、ストレージ上の重複URL検出用ハッシュ値へのアクセス回数を削減することができる。

以下、2節では用語定義を行う。3節では、キャッシュ、Bloom Filterを利用した重複URL検出の関連研究について述べる。4節では、Stable Bloom Filter Lightを用いた重複URL検出手法を提案する。5節では、ウェブグラフを用いたクロールシミュレーションにより、提案手法を評価する。

2. 用語定義

本稿では、false positive と false negative を以下のように定義する。

- false positive (FP)
未クロールのウェブページをクロール済みと誤って判定した場合
- false negative (FN)
クロール済みのウェブページを未クロールと誤って判定した場合

3. 関連研究

3.1 キャッシュを利用した重複 URL 検出

2003年にBroderらによって行われた実験により、LRUやCLOCKをウェブクローラにおける重複URL検出に利用することで、高いキャッシュヒット率を実現できることが示された[2]。実験では、サイズが 2^{18} 程度のキャッシュを用いることで、80%程度のキャッシュヒット率を実現した。

ストレージを利用した重複URL検出を行う前に、キャッシュを用いて重複URL検出を行うことで、ストレージアクセスを減らすことが可能となる。ただし、キャッシュの操作にかかる計算量は一般的にキャッシュサイズに依存するため、キャッシュサイズが大きすぎると、ストレージのみを利用した重複URL検出よりも遅くなるという実験結果が得られている[6]。

3.2 Bloom Filter を利用した重複 URL 検出

Bloom Filterとは、1970年にBloomによって提案された、FPを許容する代わりに空間効率を改善した、集合に対する要素の包含判定を近似的に実現する乱択(ランダムイズド)データ構造である[3]。Bloom Filterは長さ m のビット列である。最初に、Bloom Filterのビットはすべて0で初期化される。Bloom Filterに要素を挿入するには、 k 個のハッシュ関数を用いる。挿入する要素のハッシュ値(mod m)を k 個求め、ハッシュ値に対応するBloom Filterのビットをすべて1にすることで、挿入処理が完了する。ある要素がBloom Filterに挿入されているかどうかを判定するためには、まず、挿入時と同様に判定する要素の k 個のハッシュ値を計算する。次に、そのハッシュ値に対応するビットを調べ、対応するすべてのビットが1であったときのみ、要素が挿入済みであると判断される。図1は、要素 x 、 y が挿入されたBloom Filterにおける要素 z の包含判定である。図では、要素 z のハッシュ値に対応するビットの1つが0であるため、要素 z はBloom Filterには含まれないと判断される。

Bloom FilterではFPが発生する。 m ビットのBloom Filterに k 個のハッシュ関数を用いて n 個の要素が既に挿入されているとき、新たな要素の挿入によるFP発生率は

$$\left(1 - \left(1 - \frac{1}{m}\right)^{kn}\right)^k \dots (1)$$

によって求まる。また、 m ビットのBloom Filterに n 個の要

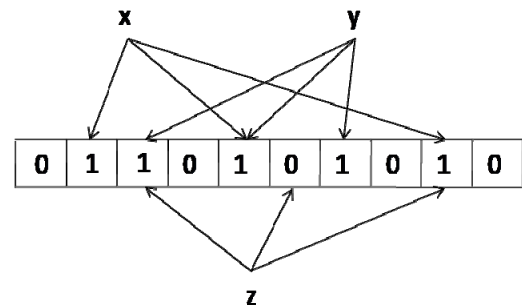


図1 要素 x, y の挿入と z の包含判定
Fig.1 Insertion of x and y and a membership test of z

素を挿入することがわかっている場合、最適なハッシュ関数の数 k は

$$k = \frac{m}{n} \ln 2 \dots (2)$$

で求めることが可能である[3]。

3.2.1 Counting Bloom Filter (CBF)

Bloom Filterの特徴として、挿入された要素の削除ができないという点が挙げられる。2000年にFanらによって提案されたCounting Bloom Filter (CBF)は、Bloom Filterが使用している1ビットのフラグの代わりに、数ビットのカウンタを用いることで、削除機能を実現したデータ構造である[7]。CBFでは、要素を追加する際に、 k 個のハッシュ値に対応するカウンタをインクリメントする。要素の包含判定は、 k 個のハッシュ値に対応するカウンタがすべて0でないかどうかを調べることで行われる。そして、要素の削除は、 k 個のハッシュ値に対応するカウンタをデクリメントすることで実現可能である。CBFでは要素の削除が可能であるが、カウンタに使用するビット数を d とすると、通常のBloom Filterの d 倍の空間を使用することになる。

3.2.2 ストリームにおける重複検出

2003年にMetwallyらによって、ランドマークウィンドウモデル(landmark window model)とスライドウィンドウモデル(sliding window model)、ジャンピングウィンドウモデル(jumping window model)におけるウィンドウに含まれる要素の包含判定にBloom Filterを利用する研究が行われた[1]。

スライドウィンドウモデルとジャンピングウィンドウモデルでは、古い情報をウィンドウから削除するためにCBFを利用している。スライドウィンドウモデルでは要素を削除するための情報を持つ追加領域が必要となる。ジャンピングウィンドウモデルでは、複数個のCBFを合成して利用することで、スライドウィンドウモデルと比較して少ないメモリ使用量でウィンドウサイズを拡大することを可能にした。

ウェブクローラにおいても、クロール中の過程で得られるクロール対象のURLはストリームであると考え、Metwallyらの手法を用いて重複URL検出を行うこともできる。

3.2.3 Stable Bloom Filter (SBF)

Bloom FilterのFP発生率を低くする別の手法として、2006年にDengらがStable Bloom Filter (SBF)を提案した[5]。SBFはCBFをベースにした手法である。SBFでは、FP発生率を減らすため、古い要素を削除するのではなく、要素の挿入時にランダムにいくつかのカウンタをデクリメントする。この処理により、FP、FN発生率がある一定の値で安定することが数学的に示されている。

SBFにおいて、要素挿入時にデクリメントするカウンタの最適な個数 p は

$$p = \left\{ \left(1 - FPR^{\frac{1}{k}} \right)^{\frac{1}{Max}} - 1 \right\}^{-1} \left(\frac{1}{k} - \frac{1}{m} \right)^{-1} \dots (3)$$

として求めることができる。mはSBFで使用するカウンタの個数、kは使用するハッシュ関数の数、FPRは許容するFP発生率、Maxはカウンタの最大値である。式(3)で求めるpを、デクリメントするカウンタの個数として使用することで、FN発生率を最小化することが可能となる。

SBFでは、Metwallyらの手法と違い、SBFに含まれている要素を管理するための追加領域や、複数のCBFを使用する必要がない。また、SBFに必要なカウンタの数mは、挿入する予定の要素数に依存しない。そのため、SBFの空間使用量は最悪でもCBFと同等となる。その代償として、要素挿入時のデクリメントに乱数生成が必要なため、計算量はO(1)だが、Metwallyらの手法と比較して計算時間が増加する。

4. 提案手法

本節では、まず、既存手法の問題点とその改善案を示す。次に、SBFの近似的な実装であるStable Bloom Filter Light (SBFL)を提案する。最後に、SBFLを用いたストレージアクセスの少ない重複URL検出手法を提案する。

4.1 既存手法の問題点と改善案

以下では、URLのキャッシング手法、Bloom Filterの既存手法の問題点と、それに対する改善案を示す。

4.1.1 キャッシュを利用した重複URL検出

URLをキャッシュすることで、重複URL検出処理を行うときに80%のキャッシュヒット率を実現できる[2]。しかし、これは十分大きなキャッシュを用いた場合であり、一般的にヒット率の高いキャッシュでは、要素の挿入・包含判定の処理をO(1)の計算量で実装することは困難である。従って、キャッシュサイズが増加すると計算時間が増加する。また、キャッシュサイズを小さくすれば計算時間を削減することは可能であるが、キャッシュヒット率が低下する。既存手法では、キャッシュヒット率の向上と計算時間の削減を同時に実現することは不可能である。

キャッシュヒット率の向上と計算時間の削減を同時に実現するためには、要素の挿入と包含判定をO(1)の計算量で行う近似的なキャッシュを用いれば良い。これにより、キャッシュサイズを増やしても計算時間が増加しなくなるため、既存手法と比較して少ない計算時間で大きなサイズのキャッシュを操作することが可能となる。同時に、キャッシュサイズを増やすことにより、キャッシュヒット率が向上する。

4.1.2 Bloom Filterを利用した重複URL検出

Bloom Filterには、挿入された要素数が増加するほどFP発生率が大きくなるという問題がある[3]。この問題は、あらかじめBloom Filterに挿入される予定の要素数がわかっているならば、挿入予定の要素数に対して十分長いビット列を用いることで回避できる。しかし、ウェブクロウリングにおいて重複検出処理の対象となるURL数は未知である。

挿入される要素数が未知である場合でも一定のFP発生率を実現するために、CBFを利用する手法がMetwallyらによって提案されている[1]。Metwallyらの手法ではCBFに挿入される要素数に最大値を設けることで、入力されるURLの数に依存せずFP発生率を固定することができる。Metwallyらの手法では、一定のルールに従って、CBFへ挿入された要素を削除していく。しかし、削除の操作によりFNが発生するという問題が起こる。

ウェブクロウラにおいては、FNが発生する問題に対する改善策として、Metwallyらの手法の前段にLRUキャッシュを置く手法が考えられる。LRUキャッシュによって頻出URLの重

複検出を使うことにより、頻出URLに対するFNの発生率を抑えることが可能となる。

4.1.3 Stable Bloom Filter (SBF)

Metwallyらの手法では、要素の挿入と包含判定はO(1)で行うことが可能であるが、挿入されている要素を管理するための追加領域が必要となる。SBFは、追加領域を使用せずに、O(1)の計算量でMetwallyらの手法と同等の処理を行うことができる[5]。

SBFでは要素の挿入と同時に、ランダムにカウンタをデクリメントする。しかし、デクリメントするカウンタをランダムに選択するために、多くの乱数を生成する必要がある。乱数の生成には計算時間がかかるため、同じ計算量のMetwallyらの手法と比較し、計算時間が増加する。例えば、SBFで用いるカウンタの数を2³⁰、ハッシュ関数の数を8、許容するFP発生率を0.001%、カウンタの最大値を3としたときには、要素を挿入するたびに約84個の乱数を生成する必要がある。

要素挿入時の乱数生成によって計算時間が増加する問題を解決するためには、乱数を生成しないという手段が考えられる。具体的には、SBFへの要素挿入時にデクリメントされるカウンタを規則的に選択することで、乱数生成の必要性をなくし、高速化することが可能となる。

SBFも、FNが発生する問題に関しては、Metwallyらの手法と同じである。

4.2 Stable Bloom Filter Light (SBFL)

SBF[5]の近似的な実装である、Stable Bloom Filter Light (SBFL)を提案する。SBFLによって、SBFを利用することにより発生する、乱数生成による速度低下の問題を解決することが可能となる。

以下では、まず、SBFLのアルゴリズムについて説明し、次にSBFLがSBFとほぼ同等の特性を持つことを証明する。

4.2.1 アルゴリズム

SBFLでは、SBFにおいて、要素挿入時に発生する乱数生成をなくすことで高速化を図っている。SBFとSBFLの主な違いは、要素挿入時に必要なカウンタのデクリメントに乱数を用いるかどうかである。

SBFでは、要素挿入時にランダムにp個のカウンタを選択しデクリメントすることで、FP発生率を抑えている。一方、SBFLでは、要素挿入時に、シーケンシャルかつ循環的にp個のカウンタをデクリメントすることで、FP発生率を抑える(図2)。SBFLでは、デクリメントするカウンタは規則的に選択されるため、乱数生成が不要である。そのため、乱数生成にかかる時間の分だけ、SBFよりも高速であると言える。

4.2.2 同等性の証明

SBFLがSBFとほぼ同等の特性を有することを証明する。この事実を証明するためには、挿入された要素のハッシュ値に対応するカウンタがデクリメントされる確率が、SBFとSBFLでほぼ等しいことを示せば良い。

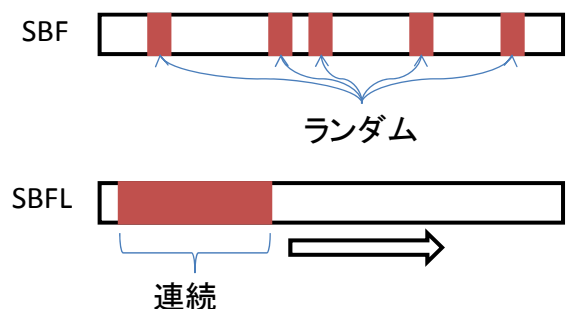


図2 SBFとSBFLにおけるデクリメント位置
Fig.2 Decrement positions of SBF and SBFL

まず、デクリメントされるカウンタを重複ありでランダムに選択した場合と、重複なしでランダムに選択した場合において、個々のカウンタがデクリメントされる確率がほぼ等しいことを示す。

SBF では、要素挿入時にランダムにいくつかのカウンタを選択しデクリメントすることで、FP 発生率を抑えている。デクリメントされるカウンタは重複して選択される可能性がある。このとき、SBF で使用しているカウンタの個数を m 、1 度の試行でデクリメントするカウンタの数を p とすると、1 度の試行で個々のカウンタが少なくとも 1 回選択される確率は

$$1 - \left(\frac{m-1}{m}\right)^p \dots (4)$$

となる。

SBFL では、デクリメントされるカウンタは重複して選択されない。よって、デクリメントされるカウンタがランダムに選択されたとすると、1 度の試行で個々のカウンタがデクリメントされる確率は

$$1 - \prod_{k=0}^{p-1} \frac{m-k-1}{m-k} \dots (5)$$

となる。しかし、 $p \ll m$ のとき、式(5)は式(4)とほぼ等しくなる。よって、 m が p と比較して十分に大きいときは、デクリメントするカウンタが重複ありで選択される確率と、重複なしで選択される確率はほぼ等しいと言える。

次に、デクリメントされる重複のない p 個のカウンタをどのように選んでも、特定の要素のハッシュ値に対応するカウンタがデクリメントされる確率が変わらないことを示す。このとき、ハッシュ関数によって生成されるハッシュ値は一樣に分布していると仮定する。

SBF では、要素の挿入は k 個のハッシュ値によって行われる。そのため、1 つの要素に対応するカウンタは、SBF 上に k 個存在する。ハッシュ値は一樣に分布しているため、 k 個のカウンタがデクリメントされる確率はそれぞれ等しい。よって簡単のため、1 つのカウンタがデクリメントされる確率を考える。このとき、ある 1 つのハッシュ値に対応するカウンタがデクリメントされる確率は p/m である。ハッシュ値が一樣に分布しているため、この確率は、 p 個のカウンタをランダムに選択しても規則的に選択しても不変である。

以上より、ハッシュ関数によって生成されるハッシュ値は一樣に分布しているという仮定の下では、SBFL は SBF とほぼ同等の特性を有することが示された。

4.3 重複 URL 検出手法

次に、LRU キャッシュと SBFL、ストレージ上の 64 ビットハッシュ値集合を用いた、ストレージアクセスの少ない近似的な重複 URL 検出手法を提案する。提案手法により、キャッシュサイズを増加させることによってキャッシュヒット率を向上させること、また、SBFL における FN 発生率を低下させることが可能となる。結果として、冗長なストレージアクセスを削減することができる。

提案手法は、(1)重複 URL 検出、(2)LRU キャッシュと SBFL の更新 2 つの過程にわけて実行される。以下では、それぞれの過程について説明する。

4.3.1 重複 URL 検出

重複 URL 検出では、まず URL が入力される。入力された URL は

1. LRU キャッシュ
2. Stable Bloom Filter Light (SBFL)
3. ストレージ上の 64 ビットハッシュ値集合

の順に重複 URL 検出処理にかけられる(図 3)。1 から 3 のい

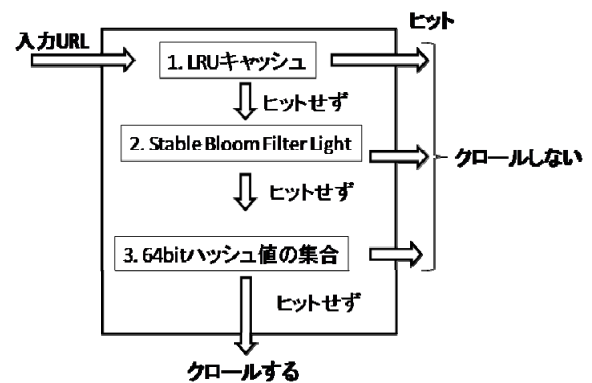


図 3 提案手法の構成
Fig.3 Framework of out method

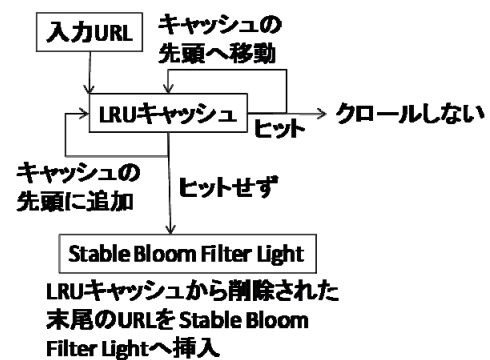


図 4 キャッシュの更新
Fig.4 Cache replacement

ずれかで重複 URL であると判定された場合は、その URL に対してクローリングを行わない。

重複 URL ではないと判定された場合は、その URL の 64 ビットハッシュ値が、ストレージ上のハッシュ値集合へ挿入される。

4.3.2 LRU キャッシュと SBFL の更新

提案手法では、重複 URL 検出を行う際に、同時に LRU キャッシュ内の要素の並べ替えと、SBFL への要素の挿入を行う(図 4)。

LRU キャッシュで URL がヒットした場合、ヒットした URL を LRU キャッシュの先頭へ移動する。

LRU キャッシュで URL がヒットしなかった場合は、URL をキャッシュの先頭へ追加する。その際に、LRU キャッシュに含まれる要素数が設定値を超えた場合には、LRU キャッシュから末尾の要素が削除され SBFL へ挿入される。

5. 実験と評価

5.1 評価指標

本稿では、キャッシュヒット率と、メモリ使用量によって、提案手法を評価する。

5.2 データセット

538, 137, 220 ホストから構成される、ホスト単位のウェブグラフを用いた。ホスト単位のウェブグラフとは、ウェブページ間のリンク関係を、ホスト単位でまとめたものである。

実験に用いたウェブグラフは e-Society プロジェクト[4]によって収集されたウェブデータを元に作成したものである。このウェブグラフに含まれる、アウトリンクを持つユニークなホストの数は 34, 848, 026 であり、アウトリンク (辺)

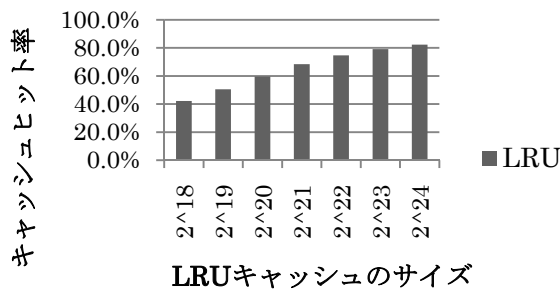


図5 LRUのキャッシュヒット率
Fig.5 Cache hit ratio of LRU

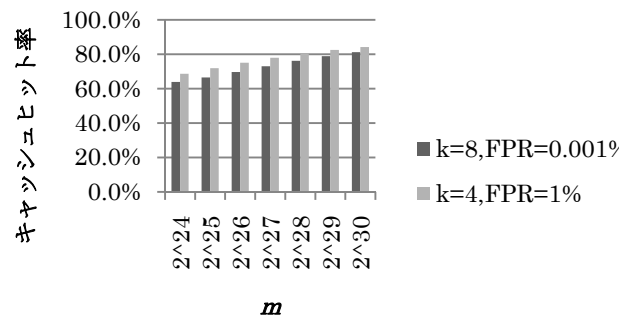


図8 提案手法のキャッシュヒット率
Fig.8 Cache hit ratio of our method

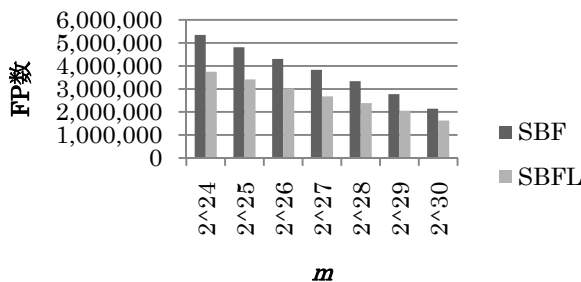


図6 SBFとSBFLのFP発生回数 (k=4, FPR=1%)
Fig.6 Number of FPs (k=4, FPR=1%)

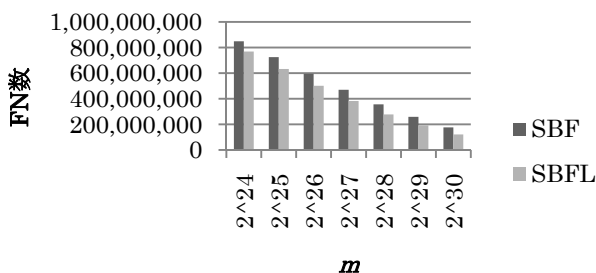


図7 SBFとSBFLのFN発生回数 (k=4, FPR=1%)
Fig.7 Number of FNs (k=4, FPR=1%)

表1 SBFとSBFLのデクリメント時間の比較
Table 1 Comparison of the time spent to decrement

手法	時間(秒)
SBF	320.7
SBFL	9.435
SBFL (32 並列)	0.9672

の総数は 4,155,528,923 である。アウトリンクを持たないホストは、リンクはされているが存在しないホスト、もしくは、未クロールのホストである。

5.3 実験方法

5.2 のウェブグラフを用いて、www.soumu.go.jp を起点として、幅優先型のウェブクロウリングシミュレーションを実行し、FP、FNの発生回数、キャッシュアクセス回数、キャッシュヒット回数を計測した。

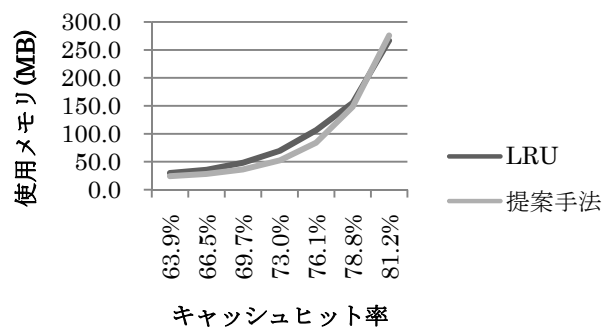


図9 使用メモリの比較(k=8, FPR=0.001%)
Fig.9 Comparison of memory usage (k=8, FPR=0.001%)

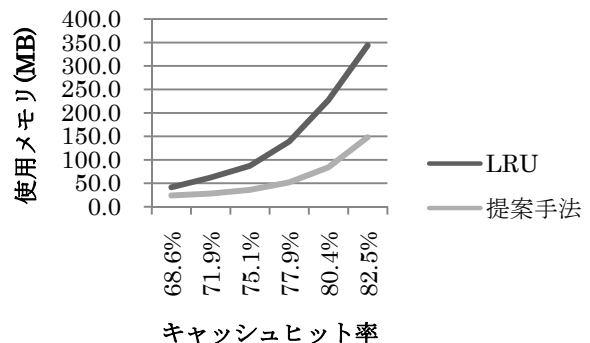


図10 使用メモリの比較(k=4, FPR=1%)
Fig.10 Comparison of memory usage (k=4, FPR=1%)

提案手法のシミュレーションは、LRU キャッシュのサイズを 2²⁰ に固定して行った。SBFL のカウンタの個数 m は 2ⁿ (n=24, 25, ..., 30) で変化させた。SBFL で使用するハッシュ関数の数は、4 個、もしくは 8 個とした。LRU キャッシュが必要とするメモリ使用量は実装により異なるが、本稿では、LRU キャッシュが必要とするメモリ使用量は、1 要素あたり 20 バイトとする。また、SBFL のカウンタは 1 個あたり 2 ビット使用する。

比較対象として、LRU キャッシュを単体で用いて実験を行った。LRU のサイズは 2ⁿ (n=18, 19, ..., 24) で変化させた。

5.4 LRU キャッシュの効果

LRU キャッシュ単体を用いてシミュレーションを行った際

のキャッシュヒット率を図5に示す。LRU キャッシュを単体で用いた際には、FPは発生しなかった。

5.5 SBFとSBFL(提案手法)の比較

まず、SBFとSBFLの実行速度を比較する。SBFとSBFLが実装上異なる部分は、要素挿入時のカウンタデクリメントである。その他の部分は全く同じ実装をしているため、カウンタをデクリメントする速度を計測し、比較する。SBFでは、乱数生成にMersenne Twister[8]を利用した。SBFLでは、2ビットのカウンタを個別にデクリメントする方法と、2ビットのカウンタ32個を並列に分岐命令無し飽和デクリメントする方法[9]を使用した。また、実行速度の計測は、Intel Xeon 2.66GHz, 16GBのメモリを搭載したマシンで行った。2³⁰個のカウンタ(メモリ使用量256MB)を20億回デクリメントしたときの計測結果を表1に示す。

計測結果より、SBFのカウンタデクリメントに要する時間は、SBFLの約34倍であることがわかる。また、SBFLにおいて、32並列のデクリメントを行った場合、SBFの約332倍の速度で動作することがわかる。

次に、 $k=4$, $FPR=1\%$ としたときのSBFとSBFLにおける、FP, FN発生回数の比較結果を図6と図7に示す。今回のシミュレーションではSBFLの方がSBFよりもFP, FNの発生率を低く抑えることができた。これは、デクリメントするカウンタの選択方法の関係上、SBFLがSBFと比較して、LRUやCLOCKに近い特性を備えているためだと考えられる。

5.6 LRU+SBLFの効果

SBFLの前にLRUを置いた場合(提案手法)において、 $k=8$, $FPR=0.001\%$ としたときと、 $k=4$, $FPR=1\%$ としたときのキャッシュヒット率を図8に示す。FPの発生回数は、 $m=2^{30}$ のとき、 $k=8$ では約10,000、 $k=4$ では図6に示すように約150万であった。

ここで、LRUキャッシュを単体で用いた場合と、提案手法を用いた場合において、同等のキャッシュヒット率のときに、それぞれの手法が必要とするメモリ使用量を比較する。

まず、提案手法において、 $k=8$, $FPR=0.001\%$ のときと、LRUキャッシュを単体で用いたときのキャッシュヒット率を比較した結果を図9に示す。なお、特定のキャッシュヒット率における、LRUキャッシュのサイズは、キャッシュヒット率を実際に計測したキャッシュサイズ(2^n)に基づいて線形補間をすることで、およそのサイズを求めた。提案手法で発生したFPは10,000前後であった。実験結果より、キャッシュヒット率が低いときには、提案手法を用いる事で使用メモリを削減できることがわかった。ただし、キャッシュヒット率が高くなるとLRUを単体で用いた方がメモリ使用量を抑えることができた。

同様に、 $k=4$, $FPR=1\%$ としたときの結果を図10に示す。提案手法においてFP発生率を高めに設定した場合は、LRUキャッシュを単体で用いたときと比較して、メモリ使用量を最大約63%削減することができた。

6. おわりに

本稿では、Stable Bloom Filterの近似的な実装であるStable Bloom Filter Lightを提案し、LRUキャッシュとStable Bloom Filter Lightを用いた、ウェブクローラ向けの効率的な重複URL検出手法を提案した。本手法を用いることで、LRUキャッシュを単体で用いたときと比較して、同等のキャッシュヒット率のときに、FPを1%許容することにより、メモリ使用量を約63%削減できることを示した。

提案手法はいくつかの課題を残している。1つは、提案手法の実行時間に関する評価である。提案手法では、LRUを単

体で用いたときと比較して、使用メモリを削減できることはわかったが、クローリングに用いたときの実行時間に関しては未調査である。

もう1つの課題は、FPの扱いである。提案手法では重複URL検出時にFPが発生する。このとき、重要なウェブページがFPと判定されることは問題である。そのため、重要なウェブページのFP発生率を下げるか、クローラ後に得られるウェブグラフから、未クローラの重要なウェブページを検出する手法が必要となる。

【謝辞】

本研究の一部は、科学研究費補助金「情報爆発に対応する高度にスケーラブルなモニタリングアーキテクチャ」によるものである。

【文献】

- [1] Ahmed Metwally, Divyakant Agrawal and Amr El Abbadi: "Duplicate Detection in Click Streams", Proc. of WWW2005, pp.12-21 (2005).
- [2] Andrei Z. Broder, Marc Najork and Janet L. Wiener: "Efficient URL Caching for World Wide Web Crawling", Proc. of WWW2003, pp.679-689 (2003).
- [3] Burton H. Bloom: "Space/Time Trade-offs in Hash Coding with Allowable Errors", CACM, Volume 13, pp.422-426 (1970).
- [4] "e-Society プロジェクト", <http://www.yama.info.waseda.ac.jp/e-society/>.
- [5] Fan Deng and Davood Rafiei: "Approximately Detecting Duplicates for Streaming Data Using Stable Bloom Filters", Proc. of ACM SIGMOD, pp.25-36 (2006).
- [6] Hsin-Tsang Lee, Derek Leonard, Xiaoming Wang and Dmitri Loguinov: "IRLbot: scaling to 6 billion pages and beyond", Proc. of WWW2008, pp.427-436 (2008).
- [7] Li Fan, Pei Cao, Jussara Almeida and Andrei Z. Broder: "Summary Cache: A Scalable Wide-Area Web Cache Sharing Protocol", IEEE/ACM Trans. on Networking, Volume 8, No. 3, pp.281-293 (2000).
- [8] Makoto Matsumoto and Takuji Nishimura: "Mersenne Twister: A 623-Dimensionally Equidistributed Uniform Pseudo-Random Number Generator", ACM Trans. on Modeling and Computer Simulation (TOMACS), pp.3-30 (1998).
- [9] Shinichi Tsuruta: "複数のビットフィールドを持つ数値の並列演算", http://www.emit.jp/prog/prog_b.html, (2000).
- [10] Hayato Yamana: "検索エンジンの信頼性", 人工知能学会誌 Volume 23 No.6, pp.752-759 (2008).

久保田 展行 Nobuyuki KUBOTA

2009 早稲田大学理工学部コンピュータネットワーク工学科卒業。2009 株式会社 Preferred Infrastructure エンジニア。大規模全文検索エンジン Sedue の開発に従事。

上田 高德 Takanori UEDA

2009 早大・基幹理工学研究科修士課程了。同研究科博士後期課程在学中。同大・メディアネットワークセンター助手。ACM, IEEE, IEICE, IPSJ, DBSJ 各学生会員。

山名 早人 Hayato YAMANA

1993 早大・理工学研究科博士課程了。博士(工学)。1993-2000 電総研。2000 早大・理工助教授。2005 同大理工学術院教授。IEEE, ACM, IEICE, IPSJ, DBSJ 各会員。