

# N.M-gram: ハッシュ値付き N-gram 法による転置インデックスの実現

平 林 幹 雄<sup>†</sup> 江 渡 浩 一 郎<sup>††</sup>

全文検索システムの転置インデックスを実現するにあたり、テキストデータから N-gram 法によって切り出したトークンを検索キーにする手法が広く用いられている。この手法には、言語中立性や再現率の完全性という利点がある反面、検索対象の文書群から抽出するトークンの数が膨大になるために、転置インデックスのサイズが肥大化して空間効率が悪化するという欠点がある。検索の際にクエリから切り出した各トークンが対象文書のテキスト内でも接続しているかどうかを判断するためには、転置インデックス内にトークンの文書内での出現位置を記録しておくことが必要となるが、この位置情報が転置インデックスの肥大化の一因となっている。本稿では、N-gram 法の欠点である転置インデックスの空間効率を改善する手法として、N.M-gram 法を提案する。N.M-gram 法では、各トークンの文書内での位置情報のかわりに後続のトークンのハッシュ値を用いることによって、N-gram 法の利点である言語中立性や再現率の完全性を保持したまま、空間効率を改善することができる。

## N.M-gram: Implementation of Inverted Index Using N-gram with Hash Values

MIKIO HIRABAYASHI<sup>†</sup> and KOUICHIROU ETO<sup>††</sup>

When constructing inverted index for full-text search system, using N-gram is very popular for tokenizing text data of target documents. Although the method has many advantages like language neutrality and perfect recall ratio, it has also shortage that the inverted index becomes large. The tokens extracted from documents tend to be enormous. The system needs to record each offset of tokens into inverted index because the offset is used for checking adjacency of tokens. The inverted index tends to be large because of the offset. In this paper, we describe *N.M-gram* method, which improves space efficiency of N-gram. The method uses hash values of succeeding tokens instead of offset in each document. The method can improve space efficiency without losing advantages of N-gram.

### 1. はじめに

近年、電子文書の増大に伴い全文検索システムの需要は増大している。大規模な文書群を高速に検索するためには、検索対象の文書群のテキストからトークンを抽出して転置することによって生成される転置インデックスを用いるのが一般的である。テキストからトークンを抽出する主要な手法のひとつとして、N-gram 法が挙げられる<sup>1)2)</sup>。N-gram 法は、文字数のみに基づいてトークンを切り出すため、対象文書の言語に対して中立であり、分かち書き（形態素解析）に見られるような解析ミスが起きず、検索の再現率が完全になるという利点がある<sup>3)</sup>。一方、N-gram 法には、抽出するトークンの数が膨大になるために、転置インデックスのサイズが肥大化して空間効率が悪化するという欠点がある。N-gram 法のインデックスに対して検索を行う際には、クエリから切り出した各トークンが対象文書のテキス

ト内でも接続しているかどうかを判定する必要がある。そのため、転置インデックス内にトークンの文書内での出現位置を記録しておくことが必要となるが、この位置情報が転置インデックスの肥大化の一因となっている。

本稿では、N-gram 法による転置インデックスの肥大化を軽減して空間効率を改善する手法として、N.M-gram 法を提案する。これは、N-gram 法によって切り出したトークンの位置情報の代用としてハッシュ値を用いることで、転置インデックスに記録する情報量の抑制を図る手法である。N.M-gram 法を用いることによって、N-gram 法の利点である言語中立性や再現率の完全性を保持したまま、空間効率を改善することができる。転置インデックスの空間効率は全文検索システムのスケーラビリティと直結するので、空間効率の改善は、大規模な全文検索システムへの要請が高まる中で意義を持つ。また、本稿では、全文検索システム Hyper Estraier における N.M-gram 法の実装について述べるとともに、N.M-gram 法と N-gram 法の空間効率を比較した評価についても述べる。

<sup>†</sup> 株式会社ミクシイ  
Mixi, Inc.

<sup>††</sup> 独立行政法人 産業技術総合研究所

National Institute of Advanced Industrial Science and Technology (AIST)

図 1 N-gram 法の転置インデックスの構造  
Fig. 1 Structure of Inverted Index of N-gram



## 2. N.M-gram 法のアルゴリズム

本節では、まず最初に N.M-gram 法の基盤となる N-gram 法のモデルについて説明し、その上で N.M-gram 法のアルゴリズムおよびその特徴について述べる。

### 2.1 N-gram 法におけるデータ構造

N-gram 法は、全文検索システムにおいて対象文書のテキストからトークンを抽出して転置インデックスを構築するための主要な手法のひとつである。N-gram 法においては、対象文書のテキストから、連続した一定の文字数の文字列をトークンとして切り出す<sup>4)</sup>。テキストを 1 文字ずつ分割する方法を 1-gram (uni-gram) 法、2 文字ずつ分割する方法を 2-gram (bi-gram) 法、3 文字ずつ分割する方法を 3-gram (tri-gram) 法と言い、N 文字ずつ分割する方法を総称して N-gram 法と言う。トークンを構成する文字が何文字であっても、トークンは 1 文字ずつずらしながら重複して切り出される。例えば「東洲齋写楽」という文字列に 2-gram 法を適用すると、「東洲」「洲齋」「齋写」「写楽」というトークンが得られる。

切り出された各々のトークンには、その位置情報が付与される。位置情報は、そのトークンを含む文書の識別子と、そのトークンが文章中の何番目に出現したかというオフセットからなる。転置インデックスにトークンの情報を記録する際には、文書内における同一パターンのトークンの情報をまとめてレコードを構成する。レコードはトークンの文字列を検索キーとして転置され、ハッシュ表や B+木などのインデックスを伴ったデータベースに記録される (図 1)。

記録した位置情報は、検索時に、トークンの存在と連続性を調べるために使われる。例えば「東洲齋」という検索要求に対しては、「東洲」および「洲齋」が同一の文書に含まれるかどうかをまず調べ、次に「東洲」を N 番目としたときに「洲齋」が N+1 番目にあるものを該当とみなすことになる。

文書内オフセットにおいて、変域を制限せずに効率的に値を表現するためには、可変長のデータ型を用いるのが一般的である<sup>5)6)</sup>。可変長のデータ型においては、小さい値は短いビット長で表現できるが、大きい値を表現するには長いビット長が必要となる。したがって、個々の文書のテキストが長くなるほど文書内オフセットのビット長は長くなる。N-gram 法によるトークンの数は文書内の文字数とほぼ等しくなるので、文書識別子の情報量よりも文書内オフセットの情報量の方が支配的になる。したがって、N-gram 法の転置インデックスにおいては、対象文書のテキストが長くなるほど、転置インデックスのサイズが肥大化し、すなわち空間効率が悪化する。なお、文書内オフセットのリストは単調増加列であるため、より値の低い差分列に可逆変換することが可能であり、それによって情報量を抑制することもできる。文書内オフセットを固定長のデータ型で表現した場合は、転置インデックスの空間効率は対象文書のテキストの長さに依存せずに済む。しかし、固定長データ型において十分な変域のビット長を設定すると、可変長のデータ型を用いるよりも空間効率が悪くなる。

N-gram 法の転置インデックスを使った検索においては、検索語の文字数が N 文字でない場合に時間効率が悪化するという問題がある。N 文字未満の検索語で検索を行う場合には、その検索語に前方一致する全てのトークンによる検索結果の和集合をとることになる。N 文字を超える検索語で検索を行う場合には、検索語を N 文字からなる複数のトークンに分割し、その各々の検索結果の積集合をとり、さらに位置情報を参照してトークンの連続性を確認する必要がある。

### 2.2 N.M-gram 法におけるデータ構造

本稿で提案する N.M-gram 法とは、言語中立性や再現率の完全性という N-gram 法の特徴を保持しつつ、空間効率を改善する方式である。N.M-gram 法においても、N-gram 法と同様に N 文字のトークンを 1 文字ずつずらしながら切り出し、各トークンに対して文書の識別子と、文書内オフセットの代わりに後続トークンのハッシュ値を付与して転置インデックスに記録する。ただし、文書内オフセットの代用として、より短いビット長のハッシュ値を用いることにより、記憶領域の節約を図る。ここで、後続 1 個のトークンのハッシュ値を付与するものを N.1-gram 法、後続 2 個のトークンのハッシュ値を付与するものを N.2-gram 法などと呼び、後続 M 個のトークンのハッシュ値を付与する手法を総称して N.M-gram 法と呼ぶことにする。例えば「東洲齋写楽」に 2.2-gram 法を適用すると、「東洲」には「洲齋」のハッシュ値と「齋写」のハッシュ値が付与され、「洲

図 2 N.M-gram 法の転置インデックスの構造  
Fig.2 Structure of Inverted Index of N.M-gram



齋」には「齋写」のハッシュ値と「写楽」のハッシュ値が付与され、「齋写」には「写楽」のハッシュ値が付与される。

N-gram 法による文書内オフセットと同様に、後続のトークンのハッシュ値を用いてもトークンの連続性を調べることができる。例えば「東洲齋」という検索要求に対しては、「東洲」および「洲齋」が同一の文書に含まれ、「東洲」の直後のトークンのハッシュ値が「洲齋」のハッシュ値と一致するものを該当とみなせばよい。

転置インデックスにトークンの情報を記録する際には、N-gram 法と同様に、文書内における同一パターンのトークンの情報をまとめてレコードを構成する(図 2)。例えば 2.2-gram 法において同一文書内に「アメリカ」および「アメヨコ」という文字列が現れた場合、「アメ」を検索キーとして、文書の識別子および「メリ」「リカ」「メヨ」「ヨコ」の各々のハッシュ値を連結したレコードが記録される。

後続のトークンのハッシュ値は固定長のデータ型で表現するため、対象文書のテキストの長さにかかわらず一定の長さになる。また、変域を狭めて短いビット長にすることができるため、N.M-gram 法におけるハッシュ値の情報量は N-gram 法における文書内オフセットの情報量に比べて小さくなる。さらに、後続のハッシュ値の同一の組み合わせが同一文書内で複数回現れた場合はそれをひとつに代表させることができるため、組み合わせが重複したトークン群の情報量が節約できる。例えば「アメリカ」という文字列が 2 回現われた場合、「アメ」という検索キーには「メリ」「リカ」「メリ」「リカ」の各々のハッシュ値を連結したレコードではなく、「メリ」「リカ」の各々のハッシュ値を連結したレコードが記録される。以上の性質により、N.M-gram 法の転置インデックスの空間効率性は、N-gram 法の転置インデックスの空間効率性よりも良くなる。

転置インデックスの空間効率が良いということは、大規模な文書群を対象にした転置インデックスを作成した場合

にも、そのサイズを小さく抑えられるということである。転置インデックスのサイズが小さいことはハードディスク等の補助記憶装置の節約になるだけでなく、主記憶装置上にキャッシュできる情報量を増やすことができるため、検索速度や転置インデックスの更新速度の高速化にも寄与する。

N.M-gram 法では、N 文字から N+M 文字までのトークンによる検索を効率的に行うことができる。例えば 2.2-gram 法では、2-gram 法でトークンを分割しながらも、4-gram 法と同等の絞り込みを行う能力を持つ。4-gram 法だと、3 文字以下の検索語では検索できないが、検索処理に多大な時間を要することになるが、2-gram 法で効率が悪くなるのは検索語が 1 文字の場合のみである。すなわち、N.M-gram 法では、N-gram の N を小さくしても効率的に検索できる。N-gram 法においても、N.M-gram 法においても、転置インデックスの構築にかかる時間計算量および領域計算量は、対象文書の規模に対して  $O(n)$  であるが、N を減らしてトークンの文字数を少なくするとトークンの種類(異なり語数)が少なくなるため、転置インデックスにおける検索キーの情報量が少なくて済む。検索キーが少ないほど、各トークンの出現情報を参照および更新する際の時間効率は向上する。

### 2.3 N.M-gram 法の特長

後続のトークンについてはハッシュ値で表現するので、異なるパターンのトークンのハッシュ値が衝突する可能性があり、その場合は意図せぬ文書が検索結果に現れてしまう。ただし、その確率は、ハッシュ値を 1 バイト(オクテット)で表現するならば  $1/256^M$  でしかないため、実用的には無視することもできると考えられる。意図せぬ結果が許容できない場合は、該当した文書の実データを参照して、検索語が実際に含まれているかどうかを調べればよい。この検査は結果の表示件数の分だけ行えばよいので、時間計算量は転置インデックスの規模に対して  $O(1)$  で済む。

不確定ではあるが高い確率で検索対象文書の絞り込みができるという点においては、N.M-gram 法をシグネチャファイルの延長ととらえることもできる。シグネチャファイルにおいては、文書内に出現した各トークンに対してハッシュ値を算出し、unary 符号で表現したハッシュ値のビット論理和を文書のシグネチャとして記録する<sup>7)</sup>。検索時には、検索クエリのシグネチャと各対象文書のシグネチャを比較し、検索クエリにおいて立っているビットの全てが対象文書において立っていれば、その文書を該当とみなす。ただし、unary 符号のハッシュ値は binary 符号のハッシュ値に比べて記憶領域あたりのエントロピーが低く、ハッシュ値の衝突率が高くなるため、十分な精度を出すためにはシグネチャのサイズを大きくする必要がある。N.M-gram 法は、転置イン

デックスのレコードに binary 符号のシグネチャを分散させて記録する手法だとも言える。対象文書の全てのシグネチャを比較する必要はなく、検索クエリの先頭のトークンで絞り込んだ候補のシグネチャのみを比較するので、N.M-gram 法はシグネチャファイルに比べて検索精度が高く、時間効率も優れている。

利便性のために検索結果において該当文書のスニペット（紹介文）を表示することが考えられる。スニペットは該当文書のテキストの中から検索語の周辺部分を抽出することによって生成されるが、位置情報を保持しない N.M-gram 法の実装ではその処理を効率的に行うことはできない。しかし、トークンの連続性の検査と同様に、スニペット生成の計算量は転置インデックスの規模に対して  $O(1)$  で済む。

検索結果において該当文書のスコアリング（順位付け）を行うことも考えられる。N-gram 法ではトークンの位置情報を用いてスコアの重み付けを行うことができるが、位置情報を保持しない N.M-gram 法の実装ではそれは不可能である。この問題に対しては、スコア情報をトークンに別途付与することで対処できる。スコア値の情報量は 8 ビット程度でも十分な精度が出るため、文書内オフセットを保持するよりはオーバーヘッドは小さいと考えられる。

### 3. N.M-gram 法の実装

本節では、全文検索システム **Hyper Estraier**<sup>8)</sup> における N.M-gram 法の実装について詳述する。

#### 3.1 転置インデックスの構築

Hyper Estraier は 2.2-gram 法による転置インデックスを実装している。転置インデックスを構築する際には、対象文書から取り出したテキストから、2-gram 法に基づいてトークンを切り出す。すなわち、連続する 2 文字のパターンを切り出す処理を、先頭から 1 文字ずつずらしながら末尾まで行う。ただし、末尾のみは 1 文字のトークンとする。

トークンは全て UTF-8 のバイト列として扱う。抽出した各々のトークンに対し、後続の 2 つのトークンに別々のハッシュ関数を適用し、各 1 バイトずつのハッシュ値を算出する。ただし、変域は 0x01 から 0xFC までとする。また、後続がない場合は 0xFF とする。

次に、トークンを整理して、同じ文字列のトークンに連なるハッシュ値を単一のレコードにまとめる。レコードの先頭には、対象文書の識別子を 4 バイトの値として置く。その後には、同じトークンの出現数を数えた上で、その値の 1 バイトをレコードを置く。この値は検索時のスコアリングに用いられる。スコア値の後には出現したハッシュ値のリストを連結する。ただし、連結するハッシュ値が既に

表 1 トークン毎に整理されたレコード群

Table 1 Records Organized by Tokens

キー	値
ファ	0x00, 0x00, 0x00, 0x01, 0x02, 0x99, 0xF6, 0x00
アイ	0x00, 0x00, 0x00, 0x01, 0x02, 0x8F, 0x02, 0x8F, 0x6D, 0x00
イル	0x00, 0x00, 0x00, 0x01, 0x02, 0x88, 0x66, 0xB1, 0x1D, 0x00
ルト	0x00, 0x00, 0x00, 0x01, 0x01, 0x42, 0x77, 0x00
とフ	0x00, 0x00, 0x00, 0x01, 0x01, 0xF2, 0x48, 0x00
ル保	0x00, 0x00, 0x00, 0x01, 0x01, 0x20, 0x84, 0x00
保存	0x00, 0x00, 0x00, 0x01, 0x01, 0x1D, 0xFF, 0x00
存	0x00, 0x00, 0x00, 0x01, 0x01, 0xFF, 0xFF, 0x00

出現したものと全く同じ場合は連結しない。さらに、各レコードの末尾に番兵として 0x00 を加え、複数の文書のデータを番兵の後に次々と連結できるようにする。

スコア値は文書内のトークンの数の平方で割ることによって調整される。その理由は、テキストが長くなるほどそれに含まれるトークンの数は多くなり、検索語に含まれるトークンが出現する確率も高まるが、検索語とは関係ない内容の文章の量も多くなる傾向があるため、その文書を検索結果の上位に提示すると検索精度が下がると考えられるからである。

以上の処理を「ファイルとファイルの保存」という文書に上述の処理を行ったとすると、表 1 に示すレコード群が得られる。

トークンをレコードのキーとし、出現数とハッシュ値と番兵を連結したデータをレコードの値として、B+木構造のデータベースに保存する。このデータベースが転置インデックスとなる。なお、B+木の実装にはリーフノードを圧縮してファイルに書き込む圧縮 B+木を用いる。圧縮アルゴリズムには deflate (LZ77 およびハフマン符号) を用いる。複数の対象文書を転置インデックスに登録する際には、番兵によって区切りがわかるので、追加するレコードの値を既存のレコードの値の末尾に連結していけばよい。

#### 3.2 検索

検索時には、検索語に対しても転置インデックス作成時と同様の N-gram 法の解析をした上で、各トークンのハッシュ値の算出を行い、各トークンに対応するレコードを転置インデックスから取り出し、ハッシュ値に検索クエリの条件と一致する部分があるか検査する（図 3）。

ハッシュ値が一致する全てのトークンに含まれる文書識別子があれば、その文書には検索語が存在することになる。該当の文書群はスコア値を用いてソートされて、検索結果として提示される。なお、検索語に含まれるトークンで実際に検査を行うものは、N-gram 法では N 個毎でよく、N+M-gram 法では N+M 個毎でよい。例えば「ファイル」という

図 3 検索時の処理  
Fig. 3 Processing of Search

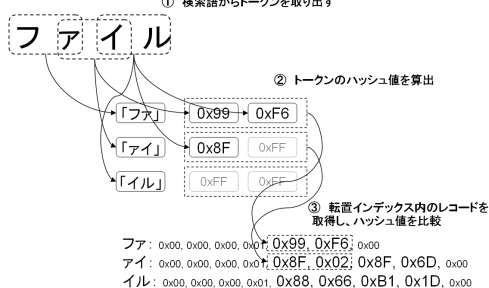


表 2 テストプログラムの実行環境

Table 2 Running Environment of the Test Program

プロセッサ	1.70GHz (Pentium M)
主記憶装置	1G バイト (PC-2700)
ハードディスク	30GB (ATA-100, 4200rpm/9.5mm)
OS	Linux 2.4.31 (ext2 ファイルシステム)
ライブラリ	GLIBC 2.3.3, QDBM 1.8.60, ZLIB 1.2.3

検索語を 2-gram 法で考えると、「ファ」が 1 番目に含まれ、「イル」が 3 番目に含まれるならば、「アイ」が 2 番目に含まれることは検査しなくてもわかる。同じく「ファイル」を 2.2-gram 法で考えると、「ファ」のレコードを取り出せば「アイ」と「イル」が後続するかどうかわかるので、その時点で検索処理を完了できる。

#### 4. 評価

本節では、N.M-gram 法の典型例として 2.2-gram 法を取り上げ、その性能を評価する。まず最初に実験環境について述べる。次に、単一の文書での空間効率について述べ、さらに、多数の文書群の転置インデックスを生成した際の空間効率および時間効率について述べる。

##### 4.1 実験環境

本実験で用いたテストプログラムの実行環境は、表 2 の通りである。転置インデックスはデータベースライブラリ QDBM<sup>9)</sup> の圧縮 B+木の API を用いて実装した。圧縮アルゴリズムには deflate 圧縮を用い、B+木の論理ページは異なり語 99 個毎に割り当てるようにした。

N.M-gram 法の転置インデックスでは、レコードの先頭には、対象文書の識別子を 4 バイトの値として置き、その後には出現したハッシュ値のリストを連結する。ただし、連結するハッシュ値が前にあるものと全く同じ場合は連結しない。ハッシュ値は 1 バイトで表現する。さらに、各レコードに番兵として 0x00 を末尾に加える。

N-gram 法の転置インデックスでは、レコードの先頭には、対象文書の識別子を 4 バイトの値として置き、その後

表 3 単一文書における 2.2-gram のトークンの出現傾向

Table 3 Appearance Tendency of 2.2-gram Tokens in a Document

文書のサイズ	6028.01 バイト
トークンの数	2387.83
異なり語数	1136.10
レコードのキーの総サイズ	6164.79 バイト
レコードの値の総サイズ	9291.34 バイト

には、トークンの文書内での出現位置のリストを差分列に変換した上で、byte aligned 符号による 128 進数の可変長数値として表現したものを連結する。さらに、各レコードに番兵として 0x00 を末尾に加える。

検索対象文書群のコーパスとして、Wikipedia 日本語版<sup>10)</sup> から無作為に抽出した 10000 件の記事のデータを用意した。UTF-8 としてのデータサイズの合計は 57.48M バイトであり、ファイルシステム上では 75.70M バイトの領域を占める。

##### 4.2 単一の文書における評価

単一の文書における 2.2-gram 法のトークンの出現傾向を調べたところ、表 3 の結果が得られた。なお、10000 件の各々の値の相加平均を代表値とした。

この結果により、一般的な日本語の文書では、2-gram 法によるトークンの数の 40%程度が異なり語数になることが予想される。その場合、レコードのキーの総サイズは対象文書のサイズとほぼ同等になり、値のサイズは対象文書のサイズの 165%程度になる。なお、同一文書内で同じパターンのトークンが複数回出現する確率は、テキストが長くなるほど大きくなる。したがって、N.M-gram 法の転置インデックスの空間効率は、対象文書群を扱ったテキストが長くなるほど向上する。また、N-gram 法と N.M-gram 法の双方において、多数の文書を対象として転置インデックスを作成した場合は、レコードのキーのサイズが全体に占める割合は小さくなり、レコードの値のサイズが支配的となる。

##### 4.3 文書群における評価

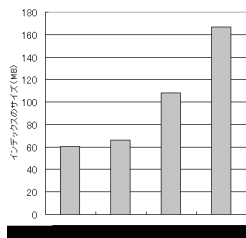
多数の文書群を対象として、2.2-gram 法と N-gram 法のインデックスの空間効率を比較する実験を行った。本実験の結果を表 4 および図 4 に示す。

異なり語数が同じである 2.2-gram 法と 2-gram 法を比較して考える。この結果から、レコードの値のサイズの合計は、双方でほぼ同じになることがわかる。一方、圧縮 B+木のサイズを見ると、2.2-gram 法の転置インデックスの方が小さく、90%ほどになっていることがわかる。その理由としては、N-gram 法のレコードは差分列に byte aligned 符号化をして既に圧縮されているのに対して、N.M-gram 法のレコード内のハッシュ値はトークンの出現傾向の偏りを保持しているために、エントロピーが低いことが考えられ

表 4 2.2-gram 法と N-gram の空間効率の比較  
Table 4 Space Efficiency of 2.2-gram and N-gram

方式	2.2-gram	2-gram	3-gram	4-gram
異なり語数	560,149	560,149	3,596,486	8,370,344
レコードサイズ合計	91.72MB	91.92MB	140.83MB	212.95MB
キーサイズ小計	3.06MB	3.06MB	28.82MB	87.42MB
値サイズ小計	88.65MB	88.86MB	112.01MB	125.53MB
圧縮 B+木のサイズ	60.44MB	65.85MB	108.08MB	166.91MB
圧縮 B+木の圧縮率	0.6590	0.7163	0.7674	0.7837
構築時間	158 秒	165 秒	1,551 秒	7,510 秒

図 4 転置インデックスのサイズのグラフ  
Fig. 4 Graph of Size of Inverted Indexes



る．転置インデックスの構築にかかる時間は，2.2-gram 法と 2-gram 法では大差ないことがわかった．

検索においてトークンを評価する回数と同じである 2.2-gram 法と 4-gram 法を比較すると，2.2-gram 法の空間効率および時間効率が非常に優れていることがわかる．4-gram 法に対して 2.2-gram 法は，転置インデックスのサイズが 3 分の 1 ほどになる．また，転置インデックスの構築にかかる時間は 50 分の 1 程度になる．

転置インデックスのサイズはレコードのサイズの合計と論理ページの圧縮率によって決まる．レコードのサイズはキーのサイズと値のサイズに分けて考えられる．レコードのキーの総サイズは異なり語数に比例し，異なり語数は N-gram 法の N の増加にともなって等比級数的に増加する．2-gram 法でのキーのサイズを基準とすると，3-gram 法では 6.5 倍程度，4-gram 法では 15 倍程度になっていることからその事が確認できる．2.2-gram 法と 2-gram 法の異なり語数は等しいので，検索においてトークンを評価する回数が同じである 4-gram 法に比べて 2.2-gram 法は有利である．なお，転置インデックスを trie 構造で表現して検索キーの領域を節約することも考えられるが，trie の探索はランダムアクセスを伴うので，時間効率が悪化してしまう．

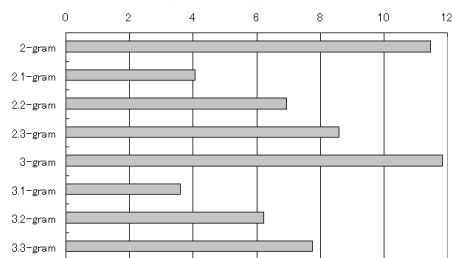
## 5. 考察

本節では，前節における実験および性能評価の過程で得られた各種の事象について考察し，N.M-gram 法の転置インデックスを実装する際に留意すべきことについて述べる．

表 5 N-gram 法および N.M-gram 法の各トークンのエントロピー  
Table 5 Entropy of Each Token of N-gram and N.M-gram

方式	エントロピー	レコード総計	圧縮 B+木	圧縮率
2-gram	11.4605 ビット	996.90MB	723.56MB	0.7258
2.1-gram	4.0604 ビット	799.13MB	475.89MB	0.5955
2.2-gram	6.9408 ビット	1008.45MB	655.45MB	0.6499
2.3-gram	8.5933 ビット	1237.40MB	853.84MB	0.6900
3-gram	11.8391 ビット	1336.31MB	1141.09MB	0.8539
3.1-gram	3.6007 ビット	1132.43MB	883.42MB	0.7801
3.2-gram	6.2250 ビット	1347.44MB	1090.23MB	0.8091
3.3-gram	7.7471 ビット	1575.97MB	1316.71MB	0.8354

図 5 各トークンのエントロピーのグラフ  
Fig. 5 Graph of Entropy of Each Token



### 5.1 エントロピーと圧縮率

各トークンに付与されるデータの性質および圧縮アルゴリズムの性質によって，レコードを圧縮してデータベースに記録する際の圧縮率は変化する．圧縮率の理論的な限界はエントロピーに基づいて求めることができる．例えば，Wikipedia 日本語版の全文書（214,358 件）をコーパスとして 2-gram 法の転置インデックスを作成すると，各文書のトークン数の平均は 1173.48 となり，文書内オフセットのエントロピーは  $\log_2 1173.48 = 10.1965$  ビットとなる．なお，可変長データ型を配列の要素にする場合，各要素のサイズを符号中に表現する必要があるため，最終的なエントロピーは 10.1965 よりも大きくなる．一方で N.M-gram 法のハッシュ値は固定長データ型を要素とした配列で表現されるため，各要素のサイズを記録するためのオーバーヘッドはかからない．

N.M-gram 法におけるエントロピーは，コーパスの各文書のトークン数だけでなく，語彙の偏りにも依存する．また，N-gram 法における文書オフセットも，差分列への変換を行った際のエントロピーはコーパスの語彙の偏りに依存することになる．そこで，同じく Wikipedia 日本語版の全文書にして各トークンに付与される文書内オフセットもしくはハッシュ値のエントロピーを集計してみたところ，表 5 および図 5 に示す結果が得られた．N-gram 法の文書内オフセットのエントロピーは，byte aligned 符号化した文書内オフセットのバイト単位のエントロピーの平均に符号の長

さの平均を掛けたものとした．N.M-gram 法のハッシュ値のエントロピーは，後続トークンの M 個のハッシュ値を M バイト単位の数値として扱って算出した．圧縮率は，レコード全体を deflate 圧縮した際の圧縮率である．

2.2-gram 法のレコードに含まれるハッシュ値のエントロピーは 2-gram 法のレコードに含まれる文書内オフセットのエントロピーの 60%ほどであり，2.2-gram 法のレコードの法が圧縮しやすいデータだと言える．2.2-gram 法ではエントロピーが 6.9408 のデータを 16 ビットの領域で表現しているので，圧縮率の限界は  $6.9408/16 = 0.4338$  となり，同様に，2.3-gram 法の圧縮率の限界は  $8.5933/24 = 0.3580$  となる．ただし，M 個のハッシュ値を単位として符号を割り当てるといふ知識を用いることによってはじめてこのような低いエントロピーになるのであり，そのような知識を持たない deflate のような汎用の圧縮アルゴリズムで圧縮率の限界に近づくのは困難である．圧縮率を向上させるには，トークン毎に辞書を切り替えるなど，N.M-gram 法の特徴を利用した符号化を行うことが求められる．

## 5.2 時間効率

ハードディスク等の大容量の補助記憶装置は以前より安価に入手できるようになっているため，全文検索システムのスケーラビリティを決めるのは空間効率ではなく時間効率だと言える．時間効率は，検索処理と更新処理に分けて考えることができるが，そのどちらがスケーラビリティを決定するかは全文検索システムの運用の仕方によって異なる．対象文書が頻繁に更新あるいは追加される場合，スケーラビリティは更新処理の時間によって制限される．検索要求が頻繁にある場合，スケーラビリティは検索処理の時間によって制限される．

Wikipedia 日本語版の全文書をコーパスとして作成した転置インデックスを用いて検索処理時間を測定した．Wikipedia 日本語版の全文検索機能を提供するデモサイト<sup>11)</sup> に一般のユーザが入力した 38321 件の検索語で 4 文字以上のものから無作為に抽出した語のリストと，38321 件の全てから無作為に抽出した語のリストを検索条件として用いた．検索語の完全一致条件で検索する処理を 10000 回行って時間を測定し，処理時間の相加平均を算出した結果を表 6 および図 6 に示す．

4 文字以上の語の検索時間の平均を見ると，2-gram 法や 2.M-gram 法よりも 3-gram 法や 3.M-gram 法の方が検索の時間効率が良いことがわかる．これは，トークンの文字数が少ないほどインデックスのレコードを特定するための時間は小さくなるが，逆に各レコードのサイズが大きくなるので，N が小さいほどレコードの読み込みに時間がかかる

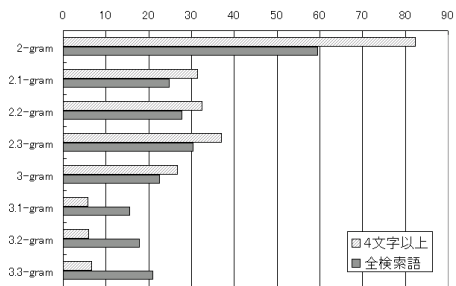
表 6 検索性能の比較

Table 6 Comparison of Search Performance

方式	4 文字以上平均	全検索語平均
2-gram	82.3359 ミリ秒	59.3935 ミリ秒
2.1-gram	31.3859 ミリ秒	24.8091 ミリ秒
2.2-gram	32.3859 ミリ秒	27.7521 ミリ秒
2.3-gram	36.8928 ミリ秒	30.3425 ミリ秒
3-gram	26.7236 ミリ秒	22.5456 ミリ秒
3.1-gram	5.7804 ミリ秒	15.5083 ミリ秒
3.2-gram	6.0102 ミリ秒	17.8223 ミリ秒
3.3-gram	6.5813 ミリ秒	20.8722 ミリ秒

図 6 検索性能のグラフ

Fig. 6 Graph of Search Performance



からである．また，N-gram 法よりも N.M-gram 法の方が検索の時間効率が良いこともわかる．これは，N-gram 法では N 個個のトークンのレコードを取り出さねばならないのに対し，N.M-gram 法では N+M 個個のトークンを取り出せばよいので，N.M-gram 法の方がレコードを読み込む回数が少ないからである．

N.M-gram 法で N が同一の場合，M が小さい方が速度が速いことがわかる．これは，M を大きくすることで減るレコードの読み込み回数の影響よりも，M を大きくすることでレコードのサイズが増大する分のオーバーヘッドの方が大きいことを示している．転置インデックスをファイルシステムに置く場合，レコードサイズの平均とデータベース内におけるページ毎のレコード数がファイルシステムのブロックサイズに等しくなるように N と M を設定すると検索処理の時間効率を最大にできるが，具体的な最適値は文書群の内容や規模によって変わる．

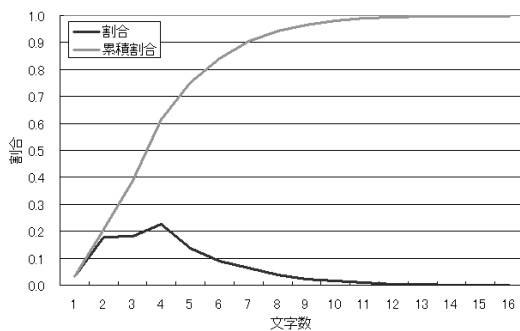
文字数を制限しない全検索語の検索時間の平均を見ると，2-gram 法や 2.M-gram 法の検索性能と 3-gram 法や 3.M 法の検索性能の差は 4 文字以上に制限した場合に比べて小さくなっていることがわかる．N-gram 法や N.M-gram 法の性能は検索語の長さに強く依存するので，実際の検索性能を捉えるには実際に入力される検索語の傾向を考慮する必要がある．そこで，一般のユーザが入力した検索語の各々の文字数を集計したところ，表 7 および図 7 に示す頻度分

表 7 検索語の文字数の頻度分布  
Table 7 Frequency of Length of Search Word

文字数	頻度	割合	累積割合
1 文字	1189	0.031027	0.031027
2 文字	6800	0.177448	0.208476
3 文字	6936	0.180997	0.389473
4 文字	8634	0.225307	0.614780
5 文字	5197	0.135618	0.750398
6 文字	3474	0.090655	0.841053
7 文字	2428	0.063360	0.904413
8 文字	1443	0.037656	0.942068
9 文字	913	0.023825	0.965893
10 文字	591	0.015422	0.981316
11 文字	325	0.008481	0.989797
12 文字	149	0.003888	0.993685
13 文字	72	0.001879	0.995564
14 文字	51	0.001331	0.996895
15 文字	20	0.000522	0.997417
16 文字	33	0.000861	0.998278

図 7 検索語の文字数のグラフ

Fig. 7 Graph of Length of Search Word



布が得られた．平均は 4.4123 文字であった．検索語の長さがトークンの長さを下回る確率は，2-gram 法や 2.M-gram 法の場合は 3.10%ほど，3-gram 法や 3.M-gram 法の場合は 20.84%ほど，4-gram 法や 4.M-gram 法の場合は 38.94%ほどになる．

実際の運用においては，更新性能と検索性能のバランスをとって  $N$  と  $M$  のパラメータを決定することになる．転置インデックスをファイルシステムに置く場合，更新および検索にかかる時間のほとんどはディスクアクセスに費され，ディスクアクセスにかかる時間のほとんどはシーク時間が占める．したがって，時間効率を改善するためにはシークの回数を減らすことが求められる． $N$ -gram 法の転置インデックスの更新処理においては，シークの回数はレコード数すなわち異なり語数に比例するので， $N$  を小さくした方がシーク回数が少なくなる． $N$ -gram 法の検索においては，シークの回数は検索語の長さを  $N$  で割った数となる．すなわち検索処理におけるシーク回数は  $N$  に反比例するの

で， $N$  を大きくした方がシーク回数が少なくなる．一方で， $N$ . $M$ -gram 法においては，更新処理のシーク回数については  $N$  に比例するが，検索処理のシーク回数は  $N+M$  に反比例する．したがって， $N$  を小さくして更新処理を高速化しても， $M$  を大きくすることで検索性能を保つことができる．

## 6. おわりに

本稿では， $N$ -gram 法の空間効率を改善する手法である  $N$ . $M$ -gram 法について述べた． $N$ -gram 法とともに用いられる位置情報の代わりにハッシュ値を用いることにより， $N$ -gram の利点である言語中立性や再現率の完全性を保ったまま，転置インデックスのサイズを縮小できることを示した．また， $N$ -gram 法の転置インデックスと  $N$ . $M$ -gram 法の転置インデックスの空間効率を比較して， $N$ . $M$ -gram 法が優れていることを示した．

## 参考文献

- 1) 原田昌紀，風間一洋，佐藤進也：“Unicode を用いた  $N$ -gram 索引の一実現方式とその評価”，情報処理学会研究会報告，2000-NL-136-17, pp.135-142, 2000.
- 2) 小川泰嗣，松田透，橋本信次：“ $N$ -gram 索引における複合検索条件の効率的な処理方法”，情報処理学会論文誌 Vol. 40 No SIG-5(TOD2), pp.43-52, 1999.
- 3) Min-Soo Kim, Kyu-Young Whang, Jae-Gil Lee, Min-Jae Lee: “ $n$ -Gram/2L: A Space and Time Efficient Two-Level  $n$ -Gram Inverted Index Structure”, Proceedings of the 31st VLDB Conference, pp.325-335, 2005.
- 4) C.E. Shannon: “Predication and Entropy of Printed English”, The Bell System Technical Journal, 1951.
- 5) P. Elias: “Gamma Code, Delta Code: Universal code-words sets and representations of the integers”, IEEE Trans. on Information Theory, IT-21(2), pp.194-203, 1975.
- 6) S.W. Golomb: “Golomb Code: Run-length Encodings”, IEEE Trans. on Information Theory, IT-12(3), pp.399-401, July 1966.
- 7) C. Faloutsos, S. Christodoulakis: “Description and Performance Analysis of Signature File Methods for Office Filing”, ACM Transactions on Office Information Systems, pp.237-257, July 1987.
- 8) “Hyper Estraier”, <http://hyperestraier.sourceforge.net/>
- 9) “Quick Database Manager”, <http://qdbm.sourceforge.net/>
- 10) “Wikipedia 日本語版”, <http://ja.wikipedia.org/>
- 11) “Wikipedia 検索”, <http://athlon64.fsij.org/mikio/wikipedia/>