

モバイル環境におけるキーワード抽出システムの検討

荒金 陽助[†] 金井 敦[†]

[†] NTT 情報流通プラットフォーム研究所

あらまし モバイル技術の発達に伴い、携帯電話をはじめとした様々な機器によってモバイル環境からの情報アクセスが可能となっている。これら情報機器の入出力インタフェースについても種々の検討や実装が行われているが、そのひとつに音声認識による入力インタフェースがある。音声認識はユーザへの負荷が非常に軽いというインタフェースとして優れた特徴を持つが、モバイル環境のような高雑音環境下での利用に際しては、ユーザが発声（入力）しようとするキーワード（認識語彙）をあらかじめ適切に設定する必要があったり、認識率とトレードオフの関係となる認識語彙数のある程度の範囲に押さえる必要がある、という課題が存在し、音声認識は情報検索のためのインタフェースとしては適格であるとは言い難かった。そこで本論文では、ニュース記事検索を対象タスクとして、適切な認識語彙を抽出する手法を提案し、被験者評価による有効性検証を行う。

キーワード 認識語彙抽出, 複合名詞, モバイル環境, 認識語彙数削減, 一致率

A Study for Keywords Extraction System for Mobile Conditions

Yosuke Aragane[†] Atsushi Kanai[†]

[†] NTT Information Sharing Platform Laboratories

Abstract According to development of mobile technology, we are able to access information in mobile environments by mobile equipments such as cellular phone or PDA. There are many studies for the input/output interfaces of these equipments. The voice recognition input interface is a one of the most major technology. It has a good characteristics of low load for users. However, in high noisy environments such as mobile environments, it needs pre-defined vocabularies for recognition. Since, there is a trade-off between a number of recognition vocabularies and the ration of recognition succeedings, it is needed that decrease of the number of recognition vocabularies. In this paper, we choose a news report searching as a application for information search using voice recognition on mobile environments. We developed a news report searching prototype system using our proposed methods of a keywords extraction method and decreasing vocabularies method. Finally, we evaluated our methods by subjective experiments.

Keywords Extract keywords, Compound nouns, Mobile environments, Decrease of recognition vocabulary, Hit ratio

1 はじめに

オフィスやホームと比較して過酷な入出力環境となるモバイル環境での入力インタフェースとして様々な検討が行われている [1, 2]. そのような入力インタフェースの中で, 音声認識インタフェースは手操作を要求しないという特徴を持つことから, カーナビゲーション装置やボイスポータルシステムなどにおいて多く利用されている [3, 4].

この音声認識インタフェースを利用したシステムでは, 入力目的の違いなどにより, 自由文認識が可能なディクテーション方式と, 高認識率を確保しやすい単語認識方式が存在する. モバイル環境という高雑音環境下での音声認識では, 高認識率を確保するために単語認識方式が採用されることが多い. この単語認識方式では, メニュー選択などユーザとシステムの間で決められた単語をインタフェースとして利用するアプリケーションにおいて, 利便性の高いサービスを提供可能である. 単語認識方式ではあらかじめ設定した認識語彙(キーワード)に合致した単語のみを認識可能であるという機能的特徴がある. そのため, 高雑音化や不明瞭な発音などの悪条件であっても, 近い認識語彙を認識結果として返すことが可能なため, 結果として物わがりの良い(高認識率を持つ)システムとなる. 一方, 利用者が音声認識を利用する直前までに認識語彙を定義しておく必要があるため, インタネットブラウジングなど動的なコンテンツを対象として検索を行う場合などは, 音声認識操作が行われる前までに認識用の単語を抽出することが必要となる.

一般のテキストベースのインタネット検索エンジンなどでは膨大な個数の単語を抽出し, それを対象として部分一致などの文字検索を行うことで検索結果を該当させ, 検索率の向上を図っている. しかし, 音声認識では, 音声認識率と認識対象語彙数がトレードオフの関係にあり, 実用的なサービスを考えるならば, 認識語彙はあまり増やせない. さらに, 認識対象語彙において部分一致で認識することが困難であるため, 認識対象語彙を正確に設定する必要がある.

上記のような認識語彙数の制限から, あまねくインタネット上の情報を検索するサービスに音声認識を適用することは現実的ではなく, ある一定の範囲のデータを対象として検索を行うサービスに適用することが望ましいと考えられる. そこで本論文で

は, 動的なコンテンツとしてニュース記事を対象として考える. すなわち, インタネット上に存在する新聞社のサイトの過去一定期間のニュース記事データを対象として, そのデータの中から音声認識インタフェースを用いてニュース記事を検索するサービスをアプリケーションとして考える. 本論文では上記のサービスのために, ニュース記事から検索のための適切な認識語彙を抽出する手法を提案する. そして, 被験者による評価を通して音声認識によるニュース記事検索の可能性について検討したので報告する.

2 提案手法

本論文で提案する“適切なキーワードを抽出する手法”および“認識語彙を少数に押さえる手法”について表1の記事を例に説明する.

2.1 キーワード抽出手法

表1の元記事の見出し文および本文を形態素解析にかけ, 一般的にキーワードと考えられる名詞を抽出したものを表2に示す. なお, 本論文では形態素解析として茶釜 [5] を利用した.

表2に示す名詞は断片的であり, 音声認識によるニュース記事検索サービスにおいて, 利用者が発声するとは考えづらい. そこで本論文では, 「複数の名詞からなるある意味を持ったキーワードを発声することが多いだろう」という仮定を置き, 名詞が連続したものを“複合名詞”として認識語彙に採用することを考える.

元記事例から抽出した複合名詞を表3に示す. ただの名詞の場合(表2)と比較して, 「九州新幹線」や「つばめ」などの利用者が発声する可能性の高いキーワードが構成されている.

2.2 認識語彙数低減手法

元記事を対象とした例の場合における抽出された認識語彙数を表4に示す.

複合名詞を採用することで25%程度の認識語彙数削減効果が現れている. しかし, 以下の理由からさらなる認識語彙の削減が必要であると考えられる.

- 「過去一定期間のニュース記事を検索する」と

表 1: 元記事例

見出し	九州新幹線：「つばめ」に乗って来て 南九州へ観光誘致
本文	<p>3月13日に新八代（熊本県） 鹿児島中央（現在は西鹿児島＝鹿児島市）間で部分開業する九州新幹線「つばめ」に期待し、南九州への観光客誘致の動きが活発になっている。JR九州は地元自治体などと協力して観光需要の掘り起こしを目指す。これに対抗して、日本航空システム（JAL）はゴールデンウィーク期間中に福岡 鹿児島間の全便の割引運賃をさらに値下げ。旅行代理店もツアーをそろえ「近くなった九州の旅」が熱を帯びている。「つばめ」は開業区間を約35分で結び、現在の2時間10分から大幅に時間短縮。在来線との乗り継ぎで博多 鹿児島間が最短で約2時間10分程度と近くなる。</p> <p>JR九州は「半分の開業のため投資効果が十分でない部分もあるが、南九州の情報発信はどんどんやる」（石原進社長）と意気込む。韓国・釜山と福岡を結ぶ定期高速船「ビートル」を増便し、4月に開業する韓国新幹線と「つばめ」を乗り比べするツアーも開始。熊本、鹿児島両県は合同でキャンペーンを行い、各地の温泉や名所、特産品をPRし、近畿日本ツーリスト、日本旅行も13日の開業日をはさんだ記念乗車ツアーなどを売り出している。また、鹿児島中央駅構内には約200種類の焼酎をそろえたショットバーも開店する。JALは4月29日～5月5日の福岡 鹿児島線の割引券「特便割引7」を一律7500円と、3500～4000円値下げ。JR九州の「つばめ2枚きっぷ」（2枚つづり、博多鹿児島間は片道7800円）との価格競争も激化しそうだ。</p>

表 2: 元記事から抽出された名詞

見出し	観光, つば, め, 九州, 南, 誘致, 新幹線
本文	<p>短縮, 旅, 価格, 店, 4月, 十分, 大幅, 開店, 分, 日, 便, 八, 記念, 品, 間, 0, 各地, PR, 期間, 鹿児島線, 情報, 片道, 競争, 県, 客, 券, 在来, 特, キャンペーン, 日本旅行, 増便, ゴールデンウィーク, 焼酎, 熱, 名所, 8, 日本航空, 船, ツアー, 割引, 最短, 線, 運賃, 構内, バー, 部分, 自治体, 進, 動き, 4, 石原, 2, JAL, 1, 観光, 旅行, 値下げ, 種類, 開業, 7, 合同, 区間, 効果, 需要, 日, 活発, つば, 5月, きっぷ, ショット, 投資, 韓国, 9, 中央, 代, 対抗, 博多, これ, 開始, 西, め, 程度, 代理, 協力, 鹿児島, 釜山, 激化, 枚, 3月, 熊本, ため, 特産, 半分, 時間, 福岡, 九州, 南, 期待, システム, 掘り起こし, 温泉, そう, 5, 駅, 誘致, 中, 地元, 市, 定期, 一律, 新幹線, 近畿日本ツーリスト, 高速, 全便, 現在, 円, 両県, 社長, JR九州, 3, 乗車, 発信, 乗り継ぎ</p>

表 4: 抽出された認識語彙の数

	名詞	複合名詞
見出し	7	4
本文	123	94

合名詞ではその記事に密接に結びつく特徴的な語彙がその多くを占めることになる。従って、マージする記事数が増加した際にも語彙の重複はそれほど多くないと考えられ、記事数の増加に伴い比較的線形に語彙数も増加してしまうことが想定される。

いったサービスでは、100以上の記事を対象に検索を行うこととなる。従って、認識語彙の重複のために線形でないにしろ、認識語彙の増加はさげられない。

- 表2に示すように、単純な名詞では一般的な語彙が数多く出現する。従って、数多くの記事の認識語彙をマージした際には重複が数多く発生すると想定される。一方表3に示すように、複

そこで本論文では、複合名詞に対して以下の2つの手法を用いることで、認識語彙の増加を抑えることを狙う。

2.2.1 数名詞の削除

音声認識では、認識語彙に対して完全一致となる発音が要求される。モバイル環境では、手元に参照

表 3: 元記事から抽出された複合名詞

見出し	観光誘致, 九州新幹線, 南九州, つばめ
本文	旅, 3月13日, 4月, 十分, 値下げ, 開業, 大幅, 開店, 旅行代理店, 合同, 4月29日, 間, つばめ, 片道7800円, 活発, つばめ2枚きっぷ, 西鹿児島, 4000円値下げ, 各地, ため投資効果, 記念乗車ツアー, 地元自治体, 13日, ショップバー, 韓国, PR, 対抗, 一律7500円, 博多, これ, 鹿児島線, 開始, 2時間10分程度, 特産品, 九州新幹線, ゴールデンウィーク期間中, 韓国新幹線, 協力, 割引券, 釜山, 部分開業, 激化, 5月5日, 開業日, 200種類, 熊本, 時間短縮, 鹿児島中央, 定期高速船, 半分, 観光需要, 福岡, キャンペーン, 日本旅行, 九州, 増便, 2枚, 期待, 在来線, 鹿児島両県, 鹿児島中央駅構内, 観光客誘致, 焼酎, 南九州, 掘り起こし, 熱, 温泉, 名所, そう, 割引運賃, ツアー, 鹿児島間, 2時間10分, 日本航空システム, 石原進社長, 35分, 最短, 特便割引7, 部分, 3500, 価格競争, 鹿児島市, 八代, 動き, 近畿日本ツーリスト, JAL, 全便, 現在, JR九州, 開業区間, 情報発信, 熊本県, 乗り継ぎ

表 5: 削除された数名詞

見出し	0個	
本文	17個	3月13日, 4月, 4月29日, 片道7800円, つばめ2枚きっぷ, 4000円値下げ, 13日, 一律7500円, 2時間10分程度, 5月5日, 開業日, 200種類, 2枚, 2時間10分, 35分, 特便割引7, 3500

表 6: 削除された短名詞

見出し	0個	
本文	6個	旅, 間, 韓国, これ, 熱, そう

以上の2つの手法を併用することで、表3に示す複合名詞による認識語彙は、表7に示す認識語彙となる。認識語彙数にして70個、当初の124個から約45%の認識語彙数の削減となる。

するデータが存在することは少ないと考えられる。また逆に、数字を含むキーワードを知りたいために検索するのであって、数字を含むキーワードで検索することは少ないとも考えられる。そこで、正確な数字を認識語彙として設定しても実行上意味がないと考え、数名詞を始め、数字を含む名詞については削除することとした。表3に示す複合名詞に対して本手法を適用した際に削除されるキーワードを表5に示す。

2.2.2 短単語の削除

「意味の乏しい短いキーワードを発声する可能性は低いだろう」という、複合名詞と対となる仮定に基づくのが本手法である。本論文では2モーラ以下のキーワードについては、発声可能性が低いとして認識語彙から削除することとした。表3に示す複合名詞に対して本手法を適用した際に削除されるキーワードを表6に示す。

3 実験評価

被験者の発声キーワードによる提案手法の有効性評価方法およびその評価結果について説明する。

3.1 評価方法

被験者に対して「最近のニュースを検索するためにキーワードを発声して下さい」という要求を行い、その結果を評価元データとした。延べ132名の被験者から616キーワードを収集した。なお、本研究の目的を考慮して、被験者にはキーワードを記述するのではなく発声するよう指示し、それを録音したもののからテキストに起こしたものを評価元データとして利用した。

今回は提案手法の評価を主眼とするために音声認識率を100%として、発声されたキーワードと完全一致するキーワードが認識語彙の中にあっただうかを評価した。なお、音声認識の特性上、キーワードの「ヨミ」が認識対象となるため、被験者発声か

表 7: 削減後の認識語彙

見出し	観光誘致, つばめ, 南九州, 九州新幹線
本文	割引運賃, 開業, 日本旅行, 福岡, 現在, 韓国, 最短, 鹿児島市, 鹿児島間, 動き, 熊本, JR九州, 定期高速船, 日本航空システム, ショットバー, 増便, 鹿児島中央駅構内, 割引券, 石原進社長, 部分開業, つばめ, 九州, 協力, 特産品, 乗り継ぎ, 価格競争, 鹿児島線, 南九州, 博多, 観光需要, 部分, 各地, 釜山, 熊本県, 九州新幹線, 地元自治体, 温泉, 西鹿児島, 十分, 情報発信, 時間短縮, 八代, 開店, 対抗, 記念乗車ツアー, 全便, 観光客誘致, 開業日, ゴールデンウィーク期間中, 鹿児島両県, 激化, 大幅, キャンペーン, 値下げ, 活発, ため投資効果, 期待, 半分, 合同, 掘り起こし, PR, JAL, 鹿児島中央, 開業区間, 在来線, 旅行代理店, ツアー, 近畿日本ツーリスト, 韓国新幹線, 焼酎

ら起こしたヨミと、キーワード抽出の際に茶筌が出力するヨミとが完全一致する場合に「一致した」と判断する。従って「松井稼頭央」などの難読キーワードについては、キーワードとして抽出されていても茶筌が出力するヨミと一致しないため、本評価では「一致せず」と判断される。本課題については、頻繁に出現するキーワードのヨミを適宜メンテナンスすることで、比較的容易に解決可能であると考えられることから本論文では言及しない。

本論文では評価元データと一致するキーワードを検索する対象記事として、記事登録時間が分単位で表示されている毎日新聞社の Web サイト [6] の記事を用いた。被験者がキーワードを発声した時間から、6,12,18,24,30,36,42,48 時間前までさかのぼった期間に登録された記事を検索対象として評価した。すなわち「24 時間」であれば、被験者発声時間から 24 時間前までに登録された記事を検索対象として評価した。

3.2 評価尺度・比較対象

評価尺度

そもそも認識語彙が存在しなければ音声認識が不可能であるため、第一の評価尺度として、発声キーワードと一致する認識語彙が存在するかどうか、という Hit 率を採用した。また、認識語彙数と認識率がトレードオフの関係にあることから、認識語彙数を第二の評価尺度とした。

比較対象

提案手法によって抽出された認識語彙（複合名詞から数名詞・短名詞をのぞいたもの）による Hit 率および認識語彙数と以下のものをそれぞれ比較評価する。

(1) 検索キーワードに一致する記事が存在するか、という観点から、インターネット上の検索エンジンなどで多く利用されるテキストベースの検索結果を Hit 率の比較対象とする。本論文ではシソーラスなどは用いずに、単純に見出しおよび本文の部分一致全文検索を行った結果と比較した。

(2) 「その記事内容をよく表しているのが見出しである」という仮定から、見出しから抽出されたキーワードが検討されることがある。そこで、見出しから抽出した名詞による認識語彙を用いて、特に Hit 率について比較評価する。

(3) 本提案手法のキーワード抽出手法の有効性を評価するため、見出しおよび本文から抽出した名詞を認識語彙とした場合の Hit 率および認識語彙数について比較評価する。

(4) 本提案手法の認識語彙数低減手法の有効性を評価するため、見出しおよび本文から抽出した複合名詞から数名詞、短名詞を削除しない認識語彙の場合の認識語彙について比較評価する。また、認識語彙削減によって Hit 率がどの程度減少するかについても評価する。

3.3 評価結果

提案手法を評価した結果について以下に説明する。評価対象期間には計 5371 件の記事が存在した。また、検索対象時間ごとの平均検索対象記事件数を表 8 に示す。

表 8: 検索対象時間毎の平均検索対象記事数

検索対象時間	平均検索対象記事数
6	43.53
12	68.62
18	104.52
24	186.40
30	228.31
36	250.47
42	282.42
48	362.97

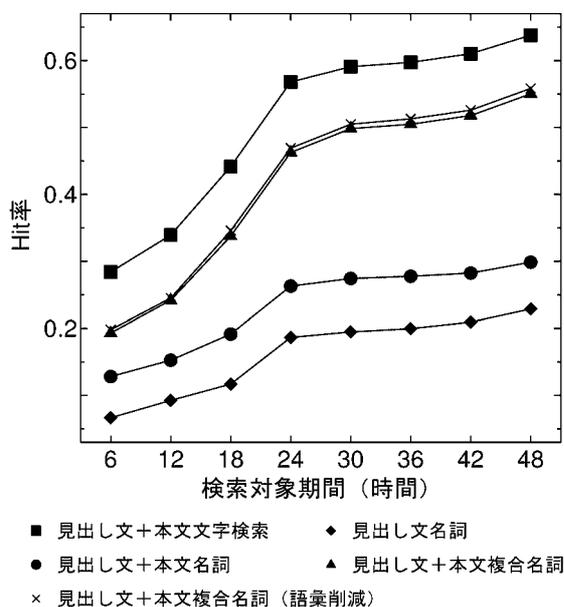


図 1: 各種抽出方法による Hit 率の評価結果

3.3.1 Hit 率

各種キーワード抽出手法による Hit 率の評価結果を図 1 に示す。横軸に検索対象とした期間、縦軸に Hit 率を表す。

提案手法によって抽出された認識語彙は、見出し文や見出し文+本文から抽出した名詞を認識語彙とした場合と比較して 50%から 190%、平均 117%の Hit 率向上が確認され、利用者は意味あるキーワードとしての複合名詞を発声する可能性が高い、という仮定を裏付ける結果となった。個別に見ると、見出し文の名詞に対しては、140%から 190%、平均 160%の Hit 率向上であり、見出し文+本文の名詞に対しては、50%から 84%、平均 74%の Hit 率向上であった。被験者が意味をなしにくい名詞を発声す

る可能性が低いと共に、見出し文では新聞特有の記述が見られ、音声認識という観点では被験者の発声するヨミと乖離する傾向が強いことが原因として考えられる。

一方、テキストベースの全文検索と比較すると、14%から 32%、平均 20%程度の Hit 率下落にとどまった。また、提案手法のうち認識語彙数低減手法を用いない複合名詞と比較した場合には、1.3%から 2.5%、平均 1.7%の Hit 率下落に過ぎず、数名詞および短い単語は発声されにくいという仮定が検証されたと考えられる。

図 1 によれば、6 時間前から 24 時間前にかけて Hit 率が急激に上昇し、その後横ばいになることから、一般的に「最近のニュース」といった場合に約 1 日前までのニュースをヒトは思い浮かべることが示されている。従って、最近のニュース検索をするサービスを考える際には、発声時から 24 時間前までのニュース記事を検索対象とすることで十分であるとも考えられる。

また、Hit 率の絶対値について、テキストベースの全文検索であっても 6 割程度の Hit 率である理由として、“酒鬼薔薇聖斗”や“松井稼頭央”などの難読キーワードや“横須賀の桜の開花予想”や“明日の天気”などの元々ニュース記事に存在しないキーワードが発声された場合が比較的多かったことが上げられる。サービスとしては、キーワード例を示すなどしてユーザを利用することで回避できる可能性が高いと考えられる。

3.3.2 認識語彙数

音声認識率とトレードオフの関係にあることから、可能な限り少数であることが望ましい認識語彙数について図 2 に示す。横軸に検索対象とした期間、縦軸に認識語彙数を表す。

提案した認識語彙数削減手法を用いない複合名詞に対して、34%から 36%、平均 34%の認識語彙数削減の効果があることが示されている。一方、見出し文と本文の名詞を認識語彙とした場合と比較すると、6 時間から 24 時間程度は提案手法の方が認識語彙数が平均 11%下回っており、24 時間以降は提案手法が逆に平均 11%上回っている。従って、提案した認識語彙数削減手法により、平均 34%程度の認識語彙数削減の効果が見込まれ、その結果、認識語彙の重複が多く存在する見出し文+本文の名詞の場合と

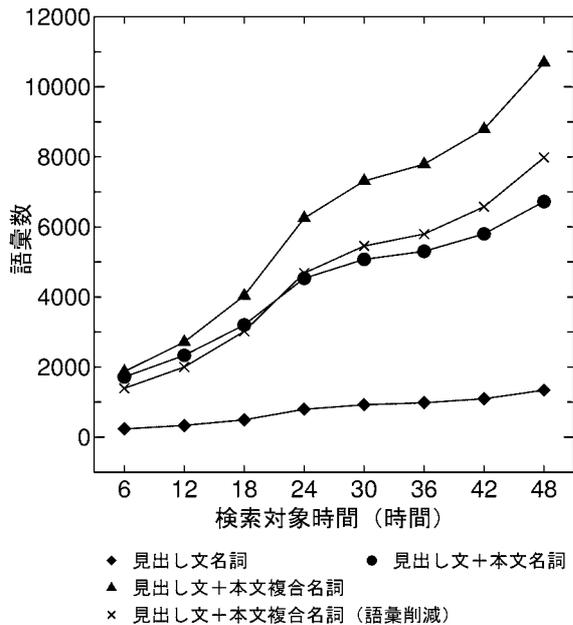


図 2: 各種抽出方法による認識語彙数の評価結果

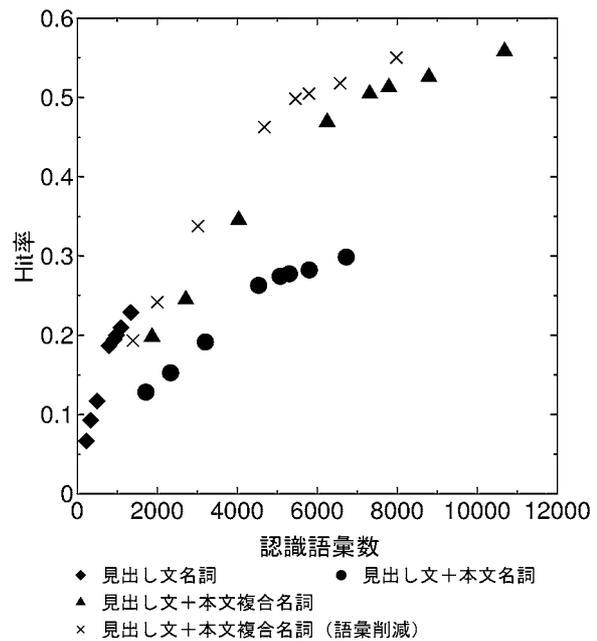


図 3: 認識語彙数当たりの Hit 率比較

同様の認識語彙数まで削減することが可能となっている。

3.3.3 認識語彙数当たりの Hit 率

本論文の評価尺度として設定してはいないが、認識語彙数当たりの Hit 率について検証した結果を図 3 に示す。横軸に認識語彙数、縦軸に Hit 率を表す。少ない認識語彙数で高い Hit 率を示すことが好ましいと考えられるため、グラフの左上側が優れていると考えられる。

提案手法である「見出し文+本文複合名詞(語彙削減)」が縦軸である Hit 率の上方までのびていると共に左上の方に位置しており、認識語彙数当たりの Hit 率の高い効率の良い手法であると考えられる。

なお、縦軸・横軸の重みを 1 対 1 として考えるならば、「見出し文名詞」の効率が最も高いと言える。見出し文は、その記事の内容を端的に示す文章であることを表していると考えられるが、抽出されるキーワード数の絶対数が少ないため、Hit 率が伸びないと言える。

4 まとめ

本論文では、モバイル環境をターゲットとして、ニュース記事を音声認識にて検索する際の音声認識語彙をニュース記事から抽出する手法について提案

し、被験者評価を通して提案手法の有効性を示した。

提案手法は、「サービス利用者はいくつかの連続した名詞からなるある一定の意味を持つ複合名詞を発声する可能性が高い」という仮定に基づき、複合名詞を利用することで Hit 率を向上させることを狙った。さらに、認識率とトレードオフの関係となる認識語彙数を削減するために、数字を含む複合名詞および 2 モーラ以下の短い複合名詞を削減することとした。数字を含む名詞は、検索結果として知りたい情報であって、手元に正確な数字のないモバイル環境においては検索をするために使う情報ではない、と考え削除対象とした。また、2 モーラ以下の短い複合名詞については、複合名詞採用の裏返しの理由となる、ある一定の意味を持ちにくい短い単語は記事を特定するためのキーワードとして利用される可能性は低い、と考え削除対象とした。

被験者発声による計 616 キーワードに対する評価結果において、見出し文の名詞を認識語彙とした場合に対して平均 160%、見出し文と本文の名詞を認識語彙とした場合に対して平均 74% の Hit 率向上が観測された。また、目標値としたテキストベースの全文検索に対しては平均 20% の Hit 率下落となった。提案手法による認識語彙数削減の Hit 率への影響は平均 1.7% の Hit 率下落に過ぎず、不要な認識語彙を効率よく削除していると考えられる。一方、認識語彙数の評価において、見出し文+本文の複合

名詞を認識語彙とした場合から平均 34%の認識語彙数削減を果たし、見出し文+本文の名詞を認識語彙とした場合と同等の認識語彙数に抑えられたことが確認された。また、認識語彙数当たりの Hit 率という効率性を見た際には、提案手法は他の手法と比較して Hit 率が高い中で効率が良い、という方向性が示された。

今後はさらなる Hit 率向上、認識語彙数削減を狙って、シソーラスの導入や認識語彙重要度によるふるい落としなどについて検討を進める予定である。

参考文献

- [1] 増井 俊之, “携帯端末のテキスト入力手法”, ヒューマンインタフェース学会誌, Vol.4, No.3, 2002.
- [2] 荒金 陽助, 久保田 浩司, “自動車内における情報入出力インタフェースの現状と課題”, 情報処理学会マルチメディア, 分散, 協調とモバイル (DICOMO2003) シンポジウム, 7F, pp.661-664, 2003.
- [3] 岩崎 知弘, 難波 利行, 石川 泰, “カーナビゲーション用音声インタフェース技術”, 自動車技術, Vol.57, No.2, pp.65-70, 2003.
- [4] NTT Communications, “V ポータル Web ページ”, <http://www.ntt.com/v-portal/>.
- [5] 奈良先端科学技術大学院大学 情報科学研究科 自然言語処理学講座, “形態素解析システム茶釜 FrontPage”, <http://chasen.aist-nara.ac.jp/hiki/ChaSen/>
- [6] 毎日新聞社, “Mainichi Interactive”, <http://www.mainichi-msn.co.jp/>