

英会話教育支援システムにおける 類似発音英単語の検索アルゴリズム

松野 了二・堤 豊

(熊本県立大学) (九州帝京短期大学)

著者らは、映画ソフトのクローズド・キャプション機能を用いた英会話教育支援システム「LETS」の研究開発を行っている。「LETS」では、映画ソフトから、クローズド・キャプション情報をパソコン上に取り込み、シナリオを作成し、そのシナリオから文単位で指定する場面を再生することができる。

今回、著者らは「LETS」において、シナリオ中の任意の文を抽出するための機能の一つとして、類似発音英単語に着目し、それを利用した検索機能を実装した。英語を母国語としない英語学習者にとって、発音をマスターすることは困難な問題の一つである。著者らは、注目する音素だけが異なり、残りの部分の発音は全く同じ2つの英単語を類似発音英単語として定義した。類似発音英単語を聞き比べることで、音素だけを聞き比べる場合よりも、聞き分ける能力が向上することが期待される。本論文では、「LETS」の概要と、類似発音英単語を検索する意義、および検索アルゴリズムの詳細、評価実験について述べ、有効性について議論する。

An Algorithm for Searching Minimal Pairs

Ryoji MATSUNO

Yutaka TSUTSUMI

(Prefectural University of Kumamoto) (Kyushu Teikyo Junior College)

We have been developing an English conversation education support system “LETS” for people of ESL/EFL, which makes use of video closed captioning (CC). “LETS” allows users to extract CC from any video movie so long as the movie has CC, to produce a scenario database, and to play back immediately any scene of the movie that are requested.

For students of ESL/EFL, it is a strong desire but sometimes difficult to acquire clear and correct pronunciation of English words, phrases and sentences. Using one feature of “LETS” mentioned above, they can search a variety of sentences that include words with similar pronunciation such as minimal pairs, and can listen to the contrasting of two similar words. As a result, it is anticipated to increase their ability to distinguish the difference between those two words, and it is also anticipated to decrease the peculiar accent coming from their own language.

In this paper, we will describe an overview of “LETS”, the significance of searching similar words, the details of our searching algorithm, and then discuss its effectiveness for language education.

1. はじめに

著者らは、クローズド・キャプション機能を用いた英会話教育支援システム「LETS」(A Language Education Tool for Speaking-up)の研究開発を行っている^[1]。クローズド・キャプション機能とは、米国で開発され、実用化されている聴覚障害者のための字幕機能である。これは通常のテレビやビデオなどの映像信号の空いている場所を利用して、字幕情報を入れておき、専用デコーダを通して再生すると、通常の映像の中に字幕が表示されるもので、米国のビデオソフトは、現在ほとんどのものが対応している。日本においても、映画の大半は米国製のものであり、かなりの映画ビデオが対応している。クローズド・キャプションを英会話教育に利用する試みは米国で活発に行われている^{[2][3]}。しかし、これ

らは単に映画をクローズド・キャプション付きで見せるだけの使い方である。英会話教育支援システム「LETS」を使えば、クローズド・キャプションで得られる英語文をパソコン上に取り込み、クローズド・キャプション・データベースを作成し、このデータベースから対応するビデオの画面を授業中に表示することができる。

この方法は次のような特長がある。

(1) 生きた英語を教えることができる。

英語教師が個人的に集め得るテキストは限られている。また、イディオムなどの最新情報も集めにくい。映画の中で使われている用語は、これらをカバーするには最適である。

(2) 不特定多数の話者による音声の聞き取りができる。

聞き取り能力の向上のためには、さまざまな話者による会話を聞くことが大切である。しかしながら、教室での教師対学生という教育では、教師個人の発音に慣れてしまい、教師以外の外国人の発音が聞き取れない、というケースも多い。メディアを使うことで、多数の話者による発音を聞く機会を与えることができる。

(3) 映像を見ながらの学習

映像を見ながら学習することで、口の動き、文脈などの情報も得ることができる。また、自分の知っている俳優の映像を見ながらの学習ということで楽しく学ぶことができる。

(4) 教材が豊富である

映画は、毎年数多く制作されており、すぐにビデオ化される。このビデオを教材として利用するため、量的には不足することは考えられない。質的には、英会話教育支援システム「LETS」にさまざまな機能を用意しそれを使うことで解決できるようにする。

英会話教育支援システム「LETS」では、ビデオソフトの英文テキストを全文データベース化している。このデータベースの英文テキストと映像の間にリンクを張っており、英文テキストから映像を検索することができる。これを映像連動機能と呼ぶ。ユーザが英文テキストを指定するためには、時系列に並べられたシ

ナリオをブラウザを使って眺めながら指示する方法のほかに、英会話教育に役立つ、さまざまな検索方法が考えられる。われわれは、検索方法の1つとして、英単語間の発音の類似性 (minimal pairs) に着目し、その検索アルゴリズムを考案し、実装した。

本論文では、2章で英会話教育支援システム「LETS」の概要について述べ、3章で、英単語の発音の類似性検索の意義について述べる。4章で作成の方針について述べ、5章でアルゴリズムについて詳細を記述する。6章で評価実験、7章で考察を行う。

2. 英会話支援システム「LETS」の概要

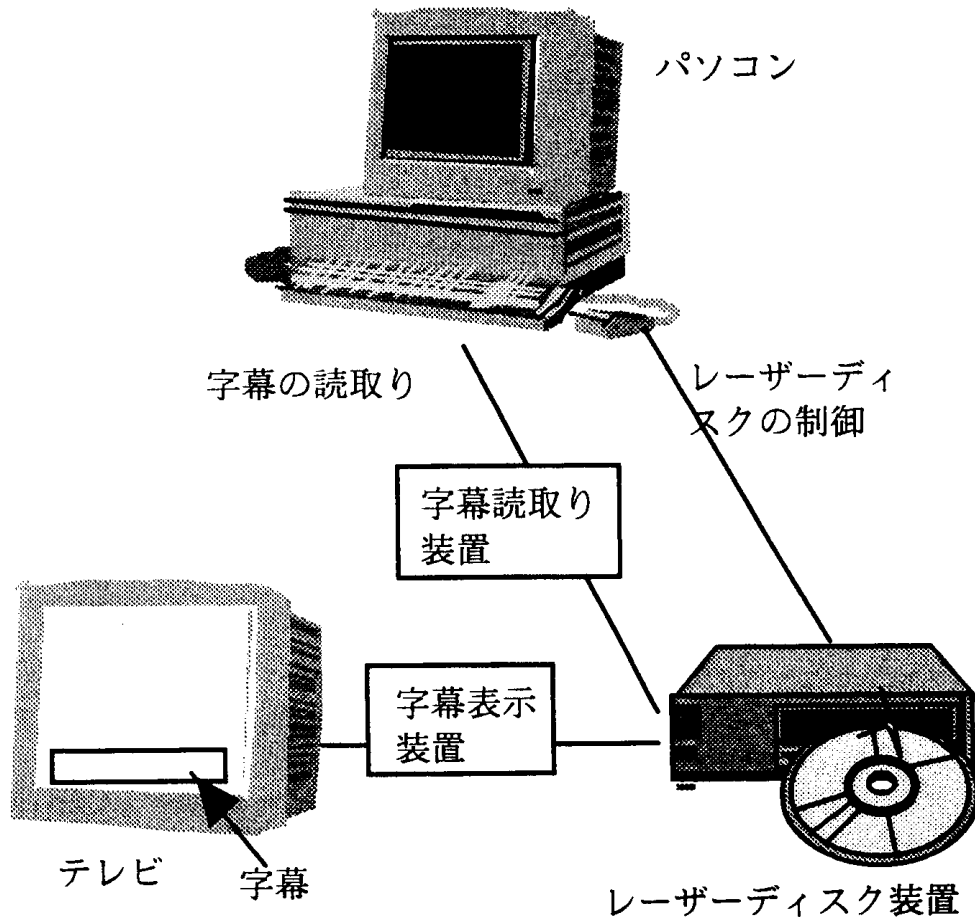


図1 システム構成

本章では、類似発音英単語の検索を実装するための環境である、英会話教育支援システム「LETS」の概要について述べる。

2. 1. ハードウェア構成

「LETS」のシステム構成を図1に示す。システムはこのように、学習教材を再生するためのレーザーディスク装置、字幕情報取り込み装置、字幕表示装置（クローズド・キャプション・アダプタ）そして、これらを制御するパソコンとから構成される。レーザーディスク装置とパソコン間および、字幕情報取り込み装置とパソコン間はRS232C回線を使って接続されている。

2. 2. ソフトウェア構成

パソコン上に作成したソフトウェアは、データベース作成用、検索用、表示用などのサブプログラム群を統合しており、メニューから選択することで、各機能を利用することができる。メインパネルを図2に示す。

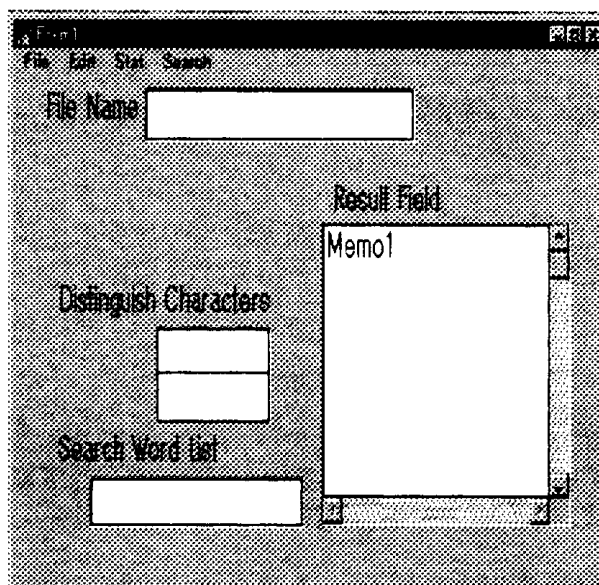


図2 メインパネル

2. 3. 運用

システムの運用は、次の2つの段階に分けられる。

(1) データベース作成

映画の字幕をすべてパソコン上に取り込み、データベース化する。

(2) 検索および表示

必要なキーワードを与えて、データベースから該当する文を検索する。検索された文をもとに映像を再生する。

データベース作成の結果はファイルとして保存することも可能であるので、必要なときに、いつでも検索することができる。また、教師は映像連動機能を使用しないならばレーザーディスクが接続されていないパソコンでも、利用することができる。従って、映画一本ごとに映画字幕データベースを保存したフロッピーディスクを用意しておき、教師が自宅で必要な表現をフロッピーディスクをもとに探し出し、教材として使える映画を検索する。という利用方法が考えられる。

以下に、英会話教育支援システム「LETS」の想定される用途を列挙する。

(1) 英会話教育の授業補助教材として利用する

本システムの開発の動機となった利用方法である。英会話教育において、学生に聞かせたい発音や、例文などを検索し、レーザーディスクで再生することで、教育効果を高めるものである。

(2) 英会話教育の授業の教材収集のために利用する

教師が自分で集めうる教材には限界がある。また、ある構文を教えたい時、例文を自分で作っても、しっくりこない、などの場合に利用する。

(3) 英会話学習のために学生が利用する

教師ではなく、学生が、自分の学習において利用する。授業中に利用する場合に比べて、自分で何度でも反復できるため、学習効果が上がる。

(4) 映画データベースを使った、英語の研究のために利用する

本システムを利用して、研究者が、映画データベースの統計的処理を簡単に行うことができる。

3. 類似発音英単語を検索する意義

ESL (English as a Second Language) や EFL (English as a Foreign Language) 学習において、困難なことのひとつに、発音の問題がある。母国語とは異なる母音、子音については、単独で聞く訓練をするだけでは修得することは難しい。これは、前後の音にひきずられて、注目している音が変化してしまうためである。このため、多数の例題を聞いて訓練する必要がある。この場合、もっとも有効な学習方法の1つが、間違いやすい2つの単語を聞き比べることである。類似発音英単語を検索することができ、間違いやすい英単語ペアを聞き比べることができれば、聞き取り能力は格段に向上することが期待できる。

また、「LETS」で利用している映画ソフトのように学習教材が一般に普及しているものである場合、多数の話者による発音を聴き比べることができる。これも大いに学習に役立つと思われる。

4. アルゴリズム作成の方針

4. 1. 再現率と適合率

情報検索において、検索された結果が正解でない場合、誤りは次の2種類に分類される。

- (1) 検索しなければならないものを検索できなかった。これを第一種の誤りという。
- (2) 検索したものの中に意図するものでないものが含まれているとき。これを第二種の誤りという。

情報検索の分野ではこれらを再現率と適合率という言葉を用いて表現する。すなわち、

$$\text{再現率} = \frac{\text{検索した候補のうち、検索すべき対象の数}}{\text{検索されるべき数}}$$

$$\text{適合率} = \frac{\text{検索結果の中で意図したものの数}}{\text{検索した数}}$$

である。

類似発音英単語の検索においても、前述の再現率と適合率が適用される。ここで英会話教育支援システム「LETS」は、対話的に利用されるため、多少検索されたものに、意図しないものがあったとしても、利用者が排除することができる。従って、適合率は100%でなくてもよいが、再現率は100%であってほしい、ということになる。そこで、われわれのアルゴリズム作成の方針(1)は、**再現率を100%に保ったまま、できるだけ適合率を高くすることとする。**

4. 2. 速度

英会話教育支援システム「LETS」は、対話的に利用される。従って、利用者が待ち遠しくない時間で処理する必要がある。「LETS」で利用しているレーザーディスク装置のシーク時間が平均4秒である^[4]ため、これよりも遅くなければ実用としては利用できるものとする。すなわち、アルゴリズム作成の方針(2)は、**検索時間は、4秒以内を実現することとする。**

5. アルゴリズムの詳細

5. 1. 類似発音英単語の定義

「類似」という言葉はあいまいであり、人によって異なった許容度を持つことがあるので、類似発音英単語についても、ここで、厳密に定義しておく。

人によっては、melon-lemon のペアを類似しているという認識を持つかもしれない。しかし、ここでは、英会話の特に聞き取り能力向上のための学習を対象にしているので、次のように定義する。

[定義] 類似発音英単語とは、注目する音素だけが異なり、残りの部分はまったく同じ発音をする2つ以上の英単語をいう。

すなわち、like-lime は、類似発音英単語であるが、mile-lime は類似発音英単語ではない。

英語は基本的に表音文字であり、英単語のスペルは、ほぼ発音に対応する。しかし、一般によく知られているように、例外も多い。表1には、同じスペルで発音が違う例を、表2には、違うスペルで発音と同じ例を示す。しかし、われわれのアルゴリズム作成方針(1)より、多少、検索が失敗しても、構わないということがあげられる。多くの利用者に使ってもらうためには、検索失敗よりも、検索に時間がかかりすぎたり、あるいはシステム要件が厳しくなる方が問題である。

表1 同じスペルで発音が違う例

スペル	単語
gh	laugh-light
th	thing-thee
a	cat-car
ea	read(present)-read(past)

表2 違うスペルで同じ発音の例

スペル	単語
gh-f	laugh-loft
s-c	site-cite
o-ow	no-know
ph-f	phase-faze

5. 2. アルゴリズム

図3に類似発音英単語の検索アルゴリズムを示す。この図に示すように、アルゴリズムは、次の部分から構成される。

(1) 英文を単語リストにする

検索対象となる、映画のクローズド・キャプション英文を単語ごとに分け、出現した文番号と共に格納する。

(2) 単語リストをソートする

前項の単語リストを、アルファベット順に並べ替える。

(3) 重複した単語を削除する

前項で並べ替えられた単語リストの中で、重複して出現した単語については、1項目に縮約する。この際、文番号には、現れた文すべての番号を保存しておく。以上については、1つの映画ソフトについて、1回だけ行っておけばよい。これにより作成された、単語リストを単語リストデータベースと呼ぶことにする。

(4) キーとなるアルファベットを2個入力する

類似発音検索のキーとなる、2つのアルファベットを指定する。例えば、fire-fileを検索したければ、R-Lを指定する。

(5) キーが含まれる単語について、キーアルファベットを*に変換し、保存する。

単語リストデータベースで、キーアルファベットを含む単語について、キーアルファベットを*に変更して、出現リストに登録する。1つのアルファベットそれぞれに、別のリストに登録する。この際、変更する文字は、アルファベット以外の文字であれば、*である必要はないが、現在は、便宜上、*を利用している。例えば、単語リストデータベースにfileがあり、キーアルファベットの1つが、lであれば、fi*eとしてリストに登録する。

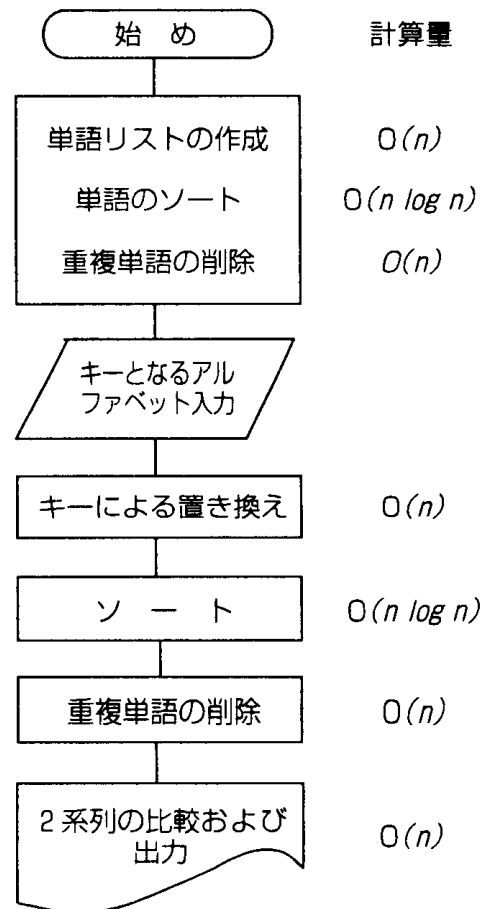


図3 アルゴリズム

(6) 前項の結果をソートする

前項の結果, 2つの出現リストが作成されるので, おのおのソートする。

(7) 重複した単語を削除する

前項でソートした2つの出現リストについて, 重複単語を削除する。この作業は, 前述したものとまったく同じ作業である。

(8) キーごとの前項の結果を比較し, 共通する項目を出力する

前項で修正された2つの出現リストを順に比較する。同一項目があれば, 類似発音の単語として出力する。先の例であれば, 単語リストデータベースに file, fire があれば, 出現リストには両方とも, fi*e が含まれることとなり, 類似発音の単語として検索されることになる。

5. 3. 計算量の計算

評価実験において, 実際の計算時間を計測するが, ここでは, 理論的な計算量について述べる。理論的な計算量が重要な理由は, 検索対象の映画字幕の文数が増えたときに, 実用的な時間で検索ができることを確認するためである。前節で述べた各部分について, データベース中の単語数を n としたときの計算量を図 3 に示す。

- (1) 単語リストを作成するのに要する時間は, 単語数に比例するため, $O(n)$ で示される。
- (2) 単語をソートするのに要する時間は, quick sort を使えば, $O(n \log n)$ である^[5]。
- (3) 重複単語を削除するのに要する時間は, 順序リストにおいては, $O(n)$ である。これにより削除された後の項目数を n_1 とする。
- (4) 入力時間は無視するものとする
- (5) 文字の検索にかかる時間は単語数に比例する。キーとなる文字 2 種類について検索を行うため, $O(2n_1)$ の時間がかかる。
- (6) ソートは前述の通り, $O(2n_1 \log n_1)$ であるが, ここで, n_1 は, 前項により大幅に減っているため, これを n_2 として別の数とする。従って, この部分でかかる時間は, $O(2n_2 \log n_2)$ である。
- (7) 重複項目の削除に要する時間は, $O(2n_2)$ である。

(8) 2系列の順序リストの比較は、 $O(n_2)$ で可能である。

以上のうち、検索時に毎回計算しなければならない項目は、(4)から(8)である。ここで、 n と n_1 、 n_2 の関係を考えてみる。延べ単語数 n に対する単語の種類数 n_1 の割合は、われわれが調査したところ、図4のように、対数的に減少する。従って、 $O(n_1) = O(\log n)$ 。文字が等確率で出現し、単語が平均6文字から構成されているとすると、ある文字が単語中に含まれている確率は

$$1 - \left[\frac{25}{26} \right]^6$$

であり、これは約0.2である。従って、 $O(n_2) = O(0.2n_1)$ 。

以上をまとめると、

$$(5)'O(2 \log n)$$

$$(6)'O(0.4 \log n (\log 0.2 + \log(\log n)))$$

$$(7)'O(0.4 \log n)$$

$$(8)'O(0.2 \log n)$$

となる。検索時に毎回計算しなければならない部分の計算量 S は、

$$S = (5)' + (6)' + (7)' + (8)'$$

$$= O(2 \log n + 0.4 \log n \cdot \log(\log n))$$

となり、これは $O(n)$ よりも小さい。従って、処理したいテキストの量が100倍になっても、処理時間は、約12倍にすぎないことがわかる。

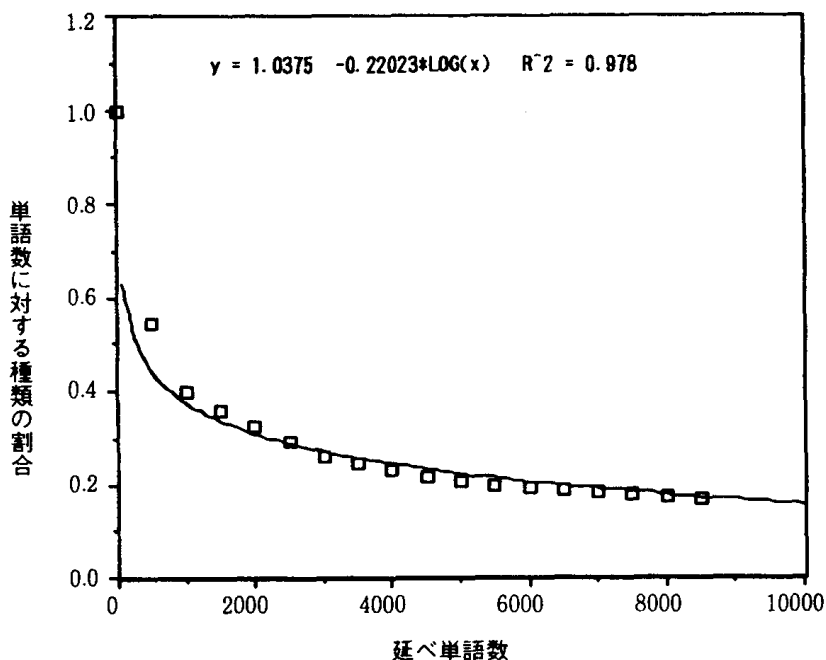


図4 ある映画ソフトにおける単語数に対する種類の割合

6. 評価実験

類似発音英単語が、期待通りに検索できることを確認するために、英会話教育支援システム「LETS」に組み込み、評価実験を行った。評価実験は、ユーザの思い通りの検索ができるかどうかを調べるための実験と、ユーザが待ち遠しくない時間で結果が検索されるかどうかを調べるための速度実験の2種類を行った。

6. 1. 再現率と適合率

映画のビデオソフト3本に関して、「LETS」でシナリオを読み込み、数種類の類似発音英単語の検索をした、検索結果を表3に示す。表で示されるように、任意の映画1本で、ユーザが期待した類似発音英単語が含まれる確率、すなわち、再現率は、

$$\text{再現率} = \frac{16}{21} = 76.2\%$$

となる。

検索された結果のうち、意図しないものは、含まれていないと考えられるので、適合率は100%である。

ビデオ(1)に関しては、再現率は100%である。逆にビデオ(2)は、再現率が57%である。これは、ビデオに含まれる台詞の量や単語の種類が影響している。これらのことから、再現率は、教材に適したビデオを選ぶことにより、高くすることができる。また、システムを拡張し、複数の映画ソフトにまたがって検索することができれば、再現率を高くすることができる。

表3 類似発音英単語の検索実験結果

スペル	ビデオ(1)	ビデオ(2)	ビデオ(3)
L-R	deal-dear	real-rear	blwon-brown, deal-dear older-order
L-WR	long-wrong		
S-TH	see-thee		pass-path
T-CH	out-ouch, gotta-gochoa too-choo	beat-beach, eat-each	eat-each
F-H	fire-hire, of-oh, foot- hoot	of-oh	of-oh
I-E	fini-fine, hi-he, hire- here, will-well	till-tell, will-well	bit-bet, hi-he, sit-set, sixsex, till-tell, will-well

ビデオ(1)：アラジンのチーム・ジェニー（ウォルト・ディズニー社）

ビデオ(2)：ダンボ（ウォルト・ディズニー社）

ビデオ(3)：マジソン郡の橋（タイム・ワーナー社）

6. 2. 速度

アルゴリズム作成の方針(2)を満たしているかどうか、検索時間を測定した。測定結果を表4に示す。ここで、実験に使用したコンピュータは、CPUがi486DX2/50MHzであり、主記憶20MB、ハードディスク540MBである。検索対象の映画ソフトは11545語、2444文からなっている。

現在のアルゴリズムでは、データはすべて主記憶上で処理しているので、CPUの速度が向上すれば、反比例して処理時間が短くなることが予想される。表からわかるように、最大でも2.6秒の検索時間で処理できる。従って、アルゴリズム作成の方針(2)は満たしている。

表4 検索時間

スペル	時間(秒)
B-V	0.7
M-N	1.7
S-TH	1.8
L-R	2.6
E-IR	2.2
平均	1.8

7. 考 察

7. 1. 実用性

評価実験で述べたとおり、速度に関しては、十分実用的であることが確認できた。しかしながら、システムを拡張して、複数の映画ソフトを扱えるようにした場合には、検索速度は十分であるとは言えない。例えば、映画100本をデータベースに格納できた場合、検索速度は現在の10倍以上かかることになる。現在のハードウェア構成をそのまま利用するとすれば、検索に30秒近くかかることになり、実用的とは言えない。パソコンの性能を上げることや、一度検索したものをキャッシュに貯えるなどの対応が必要となる。

検索結果に関しては、1本の映画ソフトだけでは、含まれる文の数が少ないため、検索できないような類似発音もある。複数のソフトにまたがって検索ができれば、この問題は解決されると考えられる。このためには、複数のソフトを同時にセットできるような機器が必要となる。

7. 2. 用途

本節では、類似発音英単語の検索の用途について述べる。本アルゴリズム作成の動機である、英会話教育支援システム「LETS」上の検索手段として有効であることは、前述の通りであるが、それ以外にも、次のような用途が考えられる。

(1) スペルチェックの補助

英単語のスペルチェック機能は、辞書に登録されている単語については指摘しない。すなわち、fly と書くべきところを fry と書いても指摘しない。本アルゴリズムでは、このように類似した単語もチェックできるので、英作文を行うときに、スペルチェックの後で本アルゴリズムを使い、チェックすれば、外国人特有の間違いなどについて、役に立つ可能性がある。

(2) 派生語の調査

-tion と -ly を指定することで、同一単語から派生した名詞と副詞を検索したり、re-, un-などを指定して、派生語を検索することができる。

7. 3. 授業への適用

類似発音英単語の検索機能を授業で利用する場合、考慮しなければならない点を考えてみる。

(1) 映画ソフトの選定

映画ソフトの選定にあたっては、授業で利用するための映画ソフトととして適当なものであるかどうか最優先となるが、そのほかに考慮すべきこととして、クローズド・キャプションのデータが実際の台詞にきちんと対応しているかどうかや、台詞の多さ、俳優の発音なども考慮することが望ましい。

(2) 検索速度

前章で述べた通り、検索速度は、最大3秒ほどかかる、従って、検索をより速く行うためには、コンピュータをより性能のよいものを使うことが必要である。また、検索した結果をビデオテープに録画したり、あるいはコンピュータ上にMPEGなどのファイルとして記録するなどの工夫をすることも可能である。

7. 4. 他方式との比較

われわれが提案しているような類似発音の英単語を検索するために、手作業で行うとすると、映画ソフトに出てきた単語のリストをすべて手元に置いて、いちいちチェックしていかなければならず、大変な手間がかかる。また、映画ソフトを検索するために、映画ソフトの上映時間分だけチェックするための時間がかかることになる。

従来、本論文で目指しているような授業を行おうとすると、準備に膨大な時間がかかり、実現することは困難であった。また、膨大な時間をかけて準備しても、提示できる例が少なく学習に効果を発揮するには至らなかった。しかし、「LETS」上で類似発音英単語の検索機能を使うことで、大量の例を提示できることになり、学習に大いに役立つと思われる。

7. 5. 改良点

スペル通りに発音しないような単語や、スペルにより一意に発音が決まらないような単語については、ユーザの意図通りには検索ができない。これらについて

は、単語毎に発音辞書を用意しておくことで解決することが可能である。しかし、辞書を検索する時間や、辞書を格納しておく領域の問題がある。速度の問題については、前述した通り、映画ソフト1本を扱うには、十分であるが、今後「LETS」で複数の映画ソフトを扱うようにすれば、何らかの改良をしなければならない。LやRなど頻繁に検索されるスペルについて、あらかじめ検索をしておくことで、対応することができると考えている。

8. まとめ

以上、類似発音の英単語検索アルゴリズムとその実装について述べた。アルゴリズムは、実用上問題がない速度で実装できるため、これを組み込んだ英会話教育支援システム「LETS」の有効な検索手段として利用できることが確認できた。

検索キーの指定方法も、アルファベット2文字を入力するという単純な方式を採用し、使いやすくした。発音のための特別な辞書をコンピュータ上に必要とせず、単純なアルゴリズムで、実用的な精度で検索できることを確認した。

また、機能面からは、英会話教育支援システム「LETS」の映像連動機能を利用することで、英会話学習に大いに役立つことが期待できる。

著者らは、自然言語文について、構文的な類似性^[6]、および意味的な類似性^[7]の両面から研究を行っており、今後、文の例示による検索なども英会話教育支援システム「LETS」に組み込んでいくことを計画している。これにより、クローズド・キャプションを用いた英会話教育の利用がさらに効率的に行えるようになる。

考察でも述べたように、検索の結果の映像をコンピュータに取り込み、コンピュータ上で再生することにより、素早く学生に提示できるようになる。この機能についても、今後研究を進めていく予定である。

謝 辞

本研究を進めていくにあたり、九州帝京短期大学ネットワーク研究室の諸君には議論を通じ、有益なコメントをいただいた。ここに記して感謝する。

参考文献

- [1] 堤, 松野:“映画字幕データベースの作成と英会話教育への適用と試み”, 電子情報通信学会技術研究報告, 教育工学, ET96-114~136, PP.73~80(1997),
- [2] Peter Shea: “Video Captioning and Language Learning”, A Review of the Literature with Implications for Multimedia Design Paper (1995).
- [3] Michal Kurmanowicz: “Video Indexing via Closed Captioning”, Report of Northeast Parallel Architectures Center at Syracuse University (1996).
- [4] “レーザディスクプレーヤ CLD-E505 ユーザーズマニュアル”, パイオニア株式会社 (1995).
- [5] Niklaus Wirth:“アルゴリズム+データ構造=プログラム”, 片山卓也訳, 日本コンピュータ協会 (1979).
- [6] 隅田, 堤:“翻訳支援のための類似用例の実用的検索法”, 電子情報通信学会論文誌, D-II, Vol, J74-D-II, No. 10, pp. 1437—1447 (1991).
- [7] 堤, 牛島:“電子メールを用いた日本語文による質問応答システムにおける類似質問の抽出法”, 情報処理学会自然言語処理研究会報告 NL-117-22 (1997).