

## 未登録語のクラス依存サブワードモデルを用いた音声認識

谷垣 宏一 山本 博史 匂坂 芳典

ATR 音声翻訳通信研究所

〒619-0288 京都府相楽郡精華町光台 2-2

{tanigaki,yama,sagisaka}@itl.atr.co.jp

あらまし 本稿では、未登録語を含む音声の高精度な認識を可能とする言語モデルを提案する。単語のクラス **N-gram** をベースとする本言語モデルは、未登録語区間に対し、その語彙クラスの読みの統計的特徴を反映したサブワードモデルを用いる点を特徴とする。また、未登録語区間の認識結果として、クラスラベル付きの読みが与えられるため、後段の言語処理が容易になっている。本方式を日本人姓・名の両クラスに適用し検討を行った。日本人姓・名データの分析結果に基づき、サブワードモデルは、単語長（モーラ数）のガンマ分布と、自動獲得したサブワード単位の **N-gram** とによる統合モデルとして構築した。音声認識実験の結果、登録語として認識を行った場合とほぼ同等の精度で、未登録語の区間・読み・クラスを同定できることがわかった。

キーワード

未登録語, 未知語, 音声認識, サブワード, 統計的言語モデル, 固有名詞

## CLASS DEPENDENT SUBWORD-MODELS FOR OUT-OF-VOCABULARY WORDS RECOGNITION

*Koichi Tanigaki, Hirofumi Yamamoto, and Yoshinori Sagisaka*

ATR Interpreting Telecommunications Research Laboratories

2-2 Hikaridai, Seika-cho, Soraku-gun, Kyoto 619-0288 Japan

{tanigaki,yama,sagisaka}@itl.atr.co.jp

Abstract

A new language model is proposed for Out-Of-Vocabulary (OOV) words to cope with inevitable demands for the recognition of proper nouns not registered in the lexicon. Multiple subword models are created for each lexical class where OOVs are predicted, so that the models are embedded in a class **N-gram** language model. The efficiency of this modeling is evaluated on Japanese family names and personal names. Speech recognition experiments show that the proposed method achieves 70% recall accuracy for OOV Japanese names, where recall is defined as correct identification of readings, classes, and locations, simultaneously. This is nearly equal to the plausible upper bound of 73% achieved under in-vocabulary condition.

key words

out of vocabulary, continuous speech recognition, subword model, proper nouns

## 1. はじめに

近年、音声認識技術の進展に伴い、音声認識の大語彙タスクへの適用が盛んに行われている[4]。しかし、大語彙音声認識のパラダイムにおいても、未登録語の問題が完全に解決するわけではない。特に、人名などの固有名詞に関しては、すべてを網羅することが困難であるといった本質的な問題もある。一方で、固有名詞にはタスク達成上重要な情報であるものも多く含まれ、音声認識の実タスク上での運用を考える際、固有名詞の未登録語処理技術は重要な課題となる。

従来、連続音声認識における読みを含めた未登録語の検出方式としては、1)音素タイプライタ等のサブワードデコーダを併用する方式[9][11]、2)サブワードを擬似的な単語として言語モデルに組み込む方式[3][8][10]、が提案されている。しかし、1)の方式は、別のデコーダを駆動する必要があるため、処理量の観点で望ましくない。また、推定未知語区間の音響スコアには最尤音素系列のスコアが使われるため、語彙内単語系列仮説との統合には、ペナルティや閾値などのヒューリスティクスが絡む。一方、2)の方式は、デコーダの変更なしに実現できる利点がある。しかし、サブワード系列として得られる未登録語に対し有効な言語処理を行うためには、後処理として、認識語彙よりも大きな語彙による形態素解析などを要する。また、単語とサブワード、あるいは、サブワード間の N-gram 確率で、言語的特質を十分反映するモデル化ができるとは考えにくく、認識制約としての有効性に疑問が残る。

本稿では、未登録語を含む音声の高精度な認識を可能とする、新しい言語モデルを提案する。本言語モデルは、単語のクラス N-gram と、未登録語認識用の複数のサブワードモデルから構成される。これらサブワードモデルは、各語彙クラスに依存して構築される。サブワードモデルのクラス依存化により、次の効果が期待できる。

- サブワードモデルの高精度化：モデル化対象を限定することで、読みの統計的特徴をより明確化することができる。更に、クラス固有のパラメータ制約を導入できるため、高精度なモデル化が期待できる。
- 検出区間の言語処理が可能：未登録語は、読みに加えクラスも同時に同定される。読みとクラスは、固有名詞の言語処理において必要十分な情報となるケースが多いものと思われる

以下では、固有名詞の下位クラスである、日本人姓・名の未登録語を対象を限定し、検討を行う。

## 2. 日本人姓・名データの分析

日本人の姓や名をサブワードの系列として眺めると、次の特徴を有することが容易に予想される。

- 長さに関する特徴：姓ではスズキ、サトウ、タカハシなど、名ではヒロシ、アキラ、イチロウなど、3ないし4モーラ長の名前が一般的と思われる。
- 並びに関する特徴：日本人の姓・名は、基本的に漢字で構成されており、姓ではヤマ、ムラ、ナカなど、名ではロウ、イチ、ヒロなど、高頻度の単位が存在するものと思われる。

こうした観点から、日本人姓・名の読みに関する統計的特徴を分析した。人名データとしては約 30 万人の著名人の名前を集録した人名リスト[2]を用いた。このデータから、漢字と平仮名のみで構成される姓・名を日本人名として抽出し、得られた姓 303,552 人分、名 295,148 人分を対象に分析を行った(表 1)。併せて比較のため、日本人姓・名以外の単語の特徴を分析する。データとしては、ATR 自然発話旅行会話データベース[4][6]より、日本人姓・名を除いた、のべ 1,155,183 単語を用いた。

単語の長さに関する統計を図 1 に示す。長さの単位としては、モーラ数を用いた。この結果から、日本人姓・名の長さが 3、4 モーラを中心に非常に偏った

表 1: モデルの訓練データ

単語総数 異なり語彙	日本人名		旅行会話
	姓	名	
	303,552	295,148	1,161,576
	19,018	20,413	13,453

日本人名の異なり語彙は、読みの異なり単語で評価し、漢字表記の違いは無視した

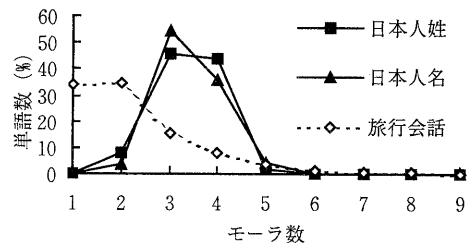


図 1: 単語の長さの分布

表 2: モーラ並びの偏り

二連鎖モーラ の種類 (頻度 上位 N 種類)	モーラ並びの被覆率(%)		
	日本人名 姓	日本人名 名	旅行会話
1	3.8	4.9	0.1
10	23.3	28.3	5.1
100	59.8	66.6	19.4
1000	84.3	82.4	35.6

頻度上位 N 種類の二連鎖モーラによる、モーラ並びの被覆率(%). 奇数長の単語があるため、被覆率が 100% になることはない。

分布を持つことが確認できる。3, 4 モーラを合わせると、姓・名ともにほぼ9割の人名が該当することになる。

次に、モーラの並びに関する統計を表 2 に示す。並びの偏りの指標として、頻度上位 N 種類のモーラ二連鎖による、モーラ並びの被覆率を調べた。例えば、日本人姓・名では、それぞれの高頻度 1000 種類のモーラ二連鎖だけで、姓・名におけるモーラ並びの 8 割以上が被覆される。

### 3. 日本人姓・名のサブワードモデルに基づく言語モデル

2章で得られた知見に基づき、日本人姓・名クラスのサブワードモデルに基づく言語モデルを構築する。

また、デコーディングの観点から、言語モデルは、近年広く用いられている N-gram 形式で扱えることが望ましい。3.4節では、本サブワードモデルを単語 N-gram 形式で実装する方法について述べる。

#### 3.1. 未登録語を含む単語系列のモデル化

提案する言語モデルのベースとなるのは、単語のクラス N-gram である。クラス N-gram では、単語系列  $W$  の言語的尤度  $\hat{p}(W)$  が一般に次式(1)で与えられる。

ただし、 $w_t$  は  $W$  の  $t$  番目の単語であり、 $c^{w_t}$  は  $w_t$  の語彙クラスを表すものとする。

$$\hat{p}(W) = \prod_t p(w_t | c^{w_t}) \cdot p(c^{w_{t+1}}, \dots, c^{w_{|W|}} | c^{w_t}) \quad (1)$$

ところで、単語  $w$  には認識語彙にない未登録語が含まれている。これら未登録語の生起確率を読み込みの統計的特徴に基づいて推定するとき、上式(1)中のクラス内単語 1-gram 確率  $p(w | c^w)$  は次式(2)により与えられる。ただし、 $M^w$  は  $w$  のモーラ系列を表す。

$$p(w | c^w) = \begin{cases} \text{if } w \in \text{Vocabulary} \\ (1 - p(OOV | c^w)) \cdot p(w | c^w, \text{inVoc}) \\ \text{otherwise} \\ p(OOV | c^w) \cdot p(M^w | c^w, OOV) \end{cases} \quad (2)$$

上式(2)において、 $p(OOV | c^w)$  は、クラスから未登録語が生起する確率であり、文献[7]等の手法で推定できる。本稿では、限られた評価セット上でサブワードモデル  $p(M^w | c^w, OOV)$  の有効な評価を行うことに主眼を置き、次式(3)および(2)によるモデル化を行う。

すなわち、未登録語の生起は予め規定したいくつかのクラス（クラスの集合を  $C^{OOV}$  とする）のみに許すこととし、これらクラスからの単語生起は全てサブワードモデルで説明することとする（登録語を作らない）。

$$p(OOV | c^w) = \begin{cases} 0 & \text{if } c^w \notin C^{OOV} \\ 1 & \text{if } c^w \in C^{OOV} \end{cases} \quad (3)$$

$$\therefore p(w | c^w) = \begin{cases} \text{if } c^w \notin C^{OOV} \\ p(w | c^w, \text{inVoc}) \\ \text{if } c^w \in C^{OOV} \\ p(M^w | c^w, OOV) = p(M^w | c^w) \end{cases} \quad (2')$$

#### 3.2. 日本人姓・名のサブワードモデル

2章で述べたように、日本人姓・名の読みには、モーラ長、およびモーラ並び、それぞれに関して特徴的な傾向が見られた。したがって、式(2)'のサブワードモデル  $p(M^w | c^w)$  は、次式(4)のように展開することにより、高精度なモデル化が可能と考えられる。ただし、 $len(M^w)$  は  $w$  のモーラ長を表す。

$$p(M^w | c^w) = p(len(M^w) | c^w) \cdot p(M^w | c^w, len(M^w)) \quad (4)$$

上式(4)の  $p(len(M) | c)$  は、日本人姓または名クラス  $c$  において、長さ  $len(M)$  の単語が生起する確率である。本稿では、その確率分布が次式(5)で与えられるガンマ分布に従うことを仮定する。ただし、 $\alpha$ 、 $\lambda$  はクラス  $c$  に依存するパラメータであり、モーラ長の平均と分散より定まる。

$$g(x) = \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x} \quad (x \geq 0) \quad (5)$$

$$\text{ここで } \Gamma(\alpha) = \int_0^\infty x^{\alpha-1} e^{-x} dx$$

一方、(4)式の  $p(M^w | c^w, len(M^w))$  は、クラス  $c^w$  において長さ  $l = len(M^w)$  のモーラ並びが  $M^w = m_1^w, m_2^w, \dots$  となる確率であり、次式(6)のサブワード N-gram によりモデル化する。ただし、 $U = u_1, u_2, \dots$  は後述の手法 (3.3節) で自動獲得した

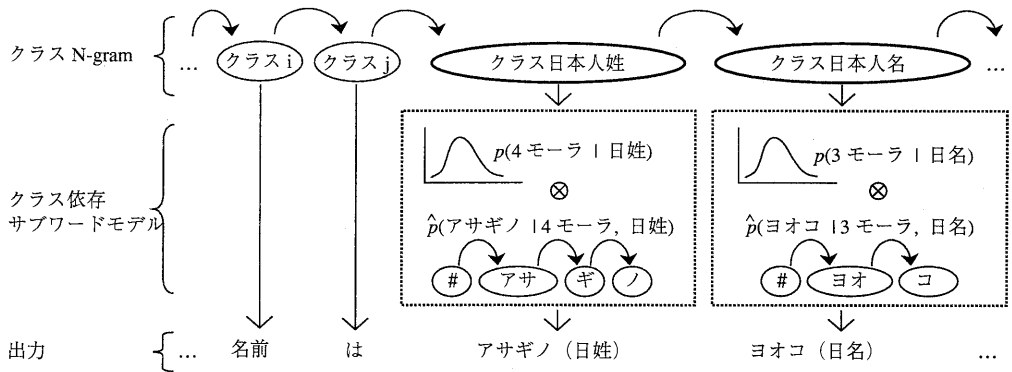


図2: クラス依存サブワードモデルに基づく言語モデル

サブワード単位（モーラまたはモーラ連鎖）の系列である。また、式中の N-gram には終端記号への遷移を含まない。

$$\begin{aligned}
 & p(M^w | c^w, \text{len}(M^w)) \\
 &= \prod_{i=1}^{\text{len}(M^w)} p(m_i^w | m_1^w, \dots, m_{i-1}^w, c^w) \\
 &= \sum_{U, s.t. U \text{ eq } M^w} \prod_j p(u_j | u_1, \dots, u_{j-1}, c^w) \quad (6) \\
 &\cong \max_{U, s.t. U \text{ eq } M^w} \prod_j p(u_j | u_1, \dots, u_{j-1}, c^w) \\
 &\cong \max_{U, s.t. U \text{ eq } M^w} \prod_j p(u_j | u_{j-n+1}, \dots, u_{j-1}, c^w)
 \end{aligned}$$

以上述べてきた提案言語モデルにおいて、「... あさぎ野 陽子 と ...」が出力される例を図2示す。例では、日本人姓・名クラスの単語「あさぎ野」、「陽子」に対し、クラスラベル付きモーラ系列「アサギノ（日姓）」、「ヨオコ（日名）」が出力される。本モデルでは、日本人姓・名の生起に対し、次の3レベルから言語的制約をかける。

【3レベルの言語的制約】

1. 単語間制約：単語のクラス N-gram を用い、単語コンテキストにおいて日本人姓・名（クラス）が生起する尤度を評価する。サブワードによる姓・名のモデル化は下位の階層に隠蔽されるため、登録語系列のモデル化には悪影響を及ぼさない。
2. 姓・名区間の継続長制約：姓、名それぞれのモーラ長に関するガンマ分布を用い、区間の姓・名らしさを評価する。この制約により、不当に短い、もしくは長いモーラ系列の湧き出しを防ぐことができる。
3. サブワードの並び制約：モーラとモーラ連鎖を単位とする姓・名のサブワード N-gram を用いる。モーラ連鎖を単位とすることで、N-gram の高精

度が期待できる。（モデル化単位とするモーラ連鎖は、次節で述べる繰り返し学習において自動的に獲得する。）

3.3. サブワードモデルの学習

表1の人名リストをもとに、姓クラス、および名クラスのサブワードモデルを構築する。以下では、個人名は等確率で出現するとし、各姓（名）の観測頻度として人名リスト中の同姓（名）者の人数を用いることとする。サブワード N-gram には、初期単位セットとして単一モーラのみを与え、後述の繰り返し学習において、逐次的に新たなモーラ連鎖を単位セットに追加していく。これら単位候補となるモーラ連鎖には頻度による予備選択を施すことで、学習の効率化を図った。モデルの構築手順は次のようになる。

【サブワードモデルの構築手順】

1. モーラ長のガンマ分布を推定する
2. サブワード N-gram の単位候補となる、高頻度のモーラ連鎖を抽出する
3. サブワード N-gram の初期単位セットとして、単一モーラを与える
4. 単位候補となる各モーラ連鎖に対し、そのモーラ連鎖を現在の単位セットに追加したときの平均尤度（式(4)）を計算する
5. 平均尤度が最大となるモーラ連鎖を現在の単位セ

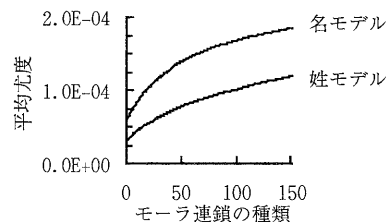


図3: モーラ連鎖の単位化による平均尤度向上

ットに追加する

#### 6. 4へ

図3に繰り返し学習における平均尤度(式(4))の変化を示す。単位候補とするモーラ連鎖は、頻度が100以上のものとした。表1のデータからは、姓モデルで1,829種類、名モデルで1,660種類のモーラ連鎖が単位候補となる。サブワードN-gramはN=2とし、1-gram, 0-gramを用いた削除補間法で補間した。図3に示すように、モーラ連鎖をサブワード単位として追加していくことで、学習データに対する平均尤度は単調に上昇する。モーラ連鎖を150個追加したモデルの平均尤度は、モーラ連鎖を用いないモデルに比べ、姓モデルで3.9倍、名モデルで3.2倍となった。サブワードモデルを単語1-gramとみなすと、単語の学習セットパーレキシティは姓モデルで74%、名モデルで69%改善されることになる。

### 3.4. 単語 N-gram 形式による実装

前節までに述べたサブワードモデルは、以下に述べる方法により、近似なく、クラス N-gram の形式で扱うことができる。そのため、言語モデルとしてクラス N-gram を扱うことが可能なデコーダであれば、デコーダの変更無しに、本方式による未登録語の認識が可能となる。ただし、極端に長い未登録語(本稿では、10モーラ以上の姓・名)が認識対象とならないことが条件となる。

サブワード N-gram で単位として用いるモーラおよびモーラ連鎖は、擬似的な単語として扱い、認識辞書、およびクラス N-gram に組み込む。その際、各サブワード単位は以下のラベル付けによる展開を行い、ラベル違いの同一サブワード単位を複数生成する。ラベルは、a)クラス、b)単語内での開始モーラ位置、c)単語終端か否か、の3項組みである。a)のクラスは、本稿では日本人姓、または日本人名の何れかである。b)の開始モーラ位置による展開は、表1に出現する最長の姓、名に合わせ、ともに終端位置が9モーラまでとなるようにした。c)で単語終端ラベルを付与したサブワード単位には、読みの終端にポーズが入ることを許容する。

ラベル付きサブワードは、その遷移に次の制約を受ける。i)登録語(のクラス)からラベル付きサブワードへの遷移は、ラベル付きサブワードの開始モーラ位置が1の場合のみ許される。逆に、ii)ラベル付きサブワードから登録語(のクラス)への遷移は、ラベル付きサブワードに単語終端ラベルが付与されている場合のみ許される。iii)ラベル付きサブワード間の遷移は、単語内でのモーラ位置が接続し、かつ同クラスに属する場合のみ許される。

## 4. 音声認識実験

提案手法の有効性を確認するため、音声認識実験を行った。以下では、二種類の言語モデルを用いて比較評価を行う。両言語モデルは、共通のベースモデルとして、表1の旅行会話データのみから生成したクラス N-gram を用いる。このベースモデルに対し、それぞれの方法で日本人姓クラス、および名クラスのクラス内単語 1-gram を置換する。

#### 【評価を行う言語モデル】

- 提案方式：日本人姓、名クラスの単語 1-gram として、姓、名それぞれのサブワードモデルを用いる。サブワード N-gram で単位として用いるモーラ連鎖は、特に断らない限り 150 個の場合を評価する。認識語彙は、日本人姓・名以外の単語 12,755 単語+サブワードで構成し、登録語の日本人姓・名は作らない。
- 登録語方式：日本人姓、名クラスの単語 1-gram として、表1の人名データによる単語 1-gram を用いる。認識語彙は、日本人姓・名以外の単語 12,755 単語+日本人姓・名 39,431 単語となる。この方式は、評価セット中のほぼ全人名をカバーする語彙を持つこと、また、提案方式がサブワードモデルの最尤推定に用いる人名データを単語 1-gram として直接用いることから、概ね提案方式による認識精度の上限値を与えるものと考えられる。

これら二方式の音声認識率を、以下の基準により評価する。

#### 【音声認識率の評価基準】

- 単語認識率：評価データに出現する全単語の認識率を評価する。日本人姓・名は、クラス(「日姓」または「日名」、読み(モーラ並び)、位置(DPによる対応付け)、が全て正しい場合のみを正解とする。ただし、読みに関し、明らかに等価な長音(ヨウコとヨオコ)は手作業で修正・評価した。
- 姓・名单語の再現率・適合率：単語認識率評価時の DP マッチに基づき、日本人姓・名のみの再現率と適合率を評価する。

評価セットには、旅行会話ドメインの 42 片側会話 4,990 単語を用いた。評価セットに出現する日本人名は、姓、名、合わせて 70 単語(異なり単語数 52)である。うち、表1の人名リストにも出現しない姓は 3 単語(アサギノ 1 単語、チンザイ 2 単語)であった。

### 4.1. 音声認識率

表3に提案方式、および登録語方式の音声認識率

を示す。提案方式では、未登録語である姓・名を、登録語として認識した場合とはほぼ同等の精度で認識できた。

予想に反し、提案方式の単語認識率が登録語方式を上回った理由の一つとして、以下が挙げられる。音響尤度の低い一部の姓・名に対し、提案方式では読み誤りはあるものの区間が正しく検出され、結果、前後の単語にまで認識誤りを誘発することが少なかったと考えられる。このことは、表 4 に示す読み誤りを無視した姓・名区間の再現率・適合率において、提案方式が優れていることから裏付けられる。

図 4 に、サブワードモデルで N-gram 単位に用いるモーラ連鎖数と姓・名单語再現率との関係を示す。単位化するモーラ連鎖を増やすことで、モデルによる姓・名の尤度が上がり、再現率が改善されるものと思われる。これは、3.2 節で述べた学習セットに対する平均尤度の改善傾向と合致する。

表 8: 音声認識率

認識率(%)	提案方式	登録語方式
単語認識率	87.51	87.30
姓・名单語再現率	70	73
適合率	67	75

姓・名は、読み・クラス・区間が全て正しい場合のみ正解として評価。登録語方式の認識率は、概ね提案方式の上限値に相当すると考えられる。

表 4: 姓・名单語の区間検出率

認識率(%)	提案方式	登録語方式
姓・名区間再現率	87	80
適合率	84	82

姓・名のクラス・区間が正しい場合を、正解として評価（音響尤度の影響が強い読み誤りは無視する）

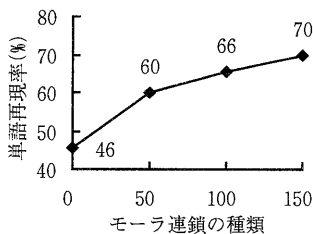


図 4: 姓・名再現率にみるモーラ連鎖の単位化効果

#### 4.2. 希有な姓・名に対する音声認識率

本稿で提案するサブワードモデルの利点は、事前に予測できない希有な単語も正しく認識できる可能性があることにある。本節では、そうした希な姓・名を模倣的に作り出すことで、提案方式の評価を行う。

評価セットには、52 種類の日本人姓・名が出現す

る。そこで、これらの単語と同じ読みを持つ全ての姓・名を表 1 の学習データから削除した後、前節と同様に提案方式と登録語方式による音声認識率の比較実験を行った。

表 5 に結果を示す。提案方式では、学習データに存在しない姓・名を与えても、31%の再現率で、その読み・クラス・区間を正しく認識できた。結果、単語認識率でも登録語方式を 0.58 ポイント上回った。

表 5 希有な姓・名入力時の音声認識率

認識率(%)	提案方式	登録語方式
単語認識率	86.66	86.08
姓・名单語再現率	31	6
適合率	36	8

訓練に用いる姓・名データから、評価セットに出現する姓・名と同じ読みを持つエントリを全て削除して実験。姓・名は、読み・クラス・区間が全て正しい場合のみ正解として評価。登録語方式の再現率・適合率が 0%にならないのは、形態素の不備により、一部の姓が「普通名詞」になっていたため。

#### 5. おわりに

音声認識における未登録語の問題に対処するため、クラス依存サブワードモデルに基づく言語モデルを提案した。また、本言語モデルを日本人姓・名の未登録語に適用し、評価を行った。音声認識実験の結果、本方式では未登録語の姓・名を、登録語として認識した場合とはほぼ同等の精度で認識できることがわかった。今後、日本人姓・名以外の語彙クラスにおいても、検討を行う必要がある。

#### 参考文献

- [1] 山本 他, “接続の方向を考慮した多重クラス複合 N-gram 言語モデル”, 信学技報, SP98-102, pp. 49-54, 1998.
- [2] 日外アソシエーツ, “30 万人よみ方書き方辞典”, ISBN4-8169-7020-7, 1993.
- [3] D. Klakow et al., “OOV-detection in large vocabulary system using automatically defined word-fragments as fillers”, *Proc. Eurospeech99*, pp. 49-52, 1999.
- [4] J. Davenport et al., “Toward realtime transcription of broadcast news”, *Proc. Eurospeech99*, pp. 651-654, 1999.
- [5] T. Morimoto et al., “Speech and language database for speech translation research”, *Proc. ICSLP94*, pp. 1791-1794, 1994.
- [6] A. Nakamura et al., “Japanese speech databases for robust speech recognition”, *Proc. ICSLP96*, pp. 2199-2202, 1996.
- [7] 政瀧 他, “品詞および可変長形態素列の複合 N-gram を用いた形態素解析”, 言語処理学会誌「自然言語処理」, Vol. 6, No. 2, pp. 41-57, 1999.
- [8] 内山 他, “仮名文字と連語登録を併用した統計的言語モデル”, 信学技報, SP99-38, pp. 87-94, 1999.
- [9] 甲斐 他, “単語 N-gram 言語モデルを用いた音声認識システムにおける未知語・冗長語の処理”, 情処論, Vol. 40, No. 4, 1999.
- [10] 伊藤 他, “被覆率を重視した大語彙連続音声認識用統計的言語モデル”, 音講論集, 2-1-7, pp. 65-66, 1999.3.
- [11] 伊藤 他, “連続音声認識における未知語の扱い”, 信学技報, SP91-96, pp. 41-47, 1991.