PodCastle:ポッドキャスト音声認識の ための集合知を活用した言語モデル学習

緒 方 淳^{†1} 後 藤 真 孝[†]

本稿では、ポッドキャスト音声認識の性能向上のための言語モデル学習手法について述べる。実環境音声であるポッドキャストは、その発話スタイルやトピックなどが多様であるため、従来のように特定タスクの事前コーパスに基づいて高精度な言語モデルを構築することはできない。そこで、本研究ではWebサービス「PodCastle」を通じて得られる集合知、すなわちエンドユーザによる音声認識誤りの訂正結果を活用した言語モデル学習手法を提案する。ポッドキャスト音声認識実験の結果、本学習システムが有効に働くことを確認した。

PodCastle: Collaborative Training of Language Models Based on Wisdom of Crowds for Podcast Transcription

Jun Ogata^{†1} and Masataka Goto^{†1}

This paper presents language modeling techniques for improving automatic transcription of podcasts. Since podcasts include various kinds of tasks and topics, accurate language modeling based on task-specific corpora is impractical. To solve this problem, we introduce collaborative training of language models on the basis of wisdom of crowds, i.e., podcast-speech transcripts annotated by anonymous users on our web service PodCastle. From our experimental results on actual podcast speech data, the effectiveness of the proposed language model training was confirmed.

†1 産業技術総合研究所

National Institute of Advanced Industrial Science and Technology (AIST)

1. はじめに

近年,音声認識技術の高度化,ならびに実環境音声データを対象とした大規模なコーパス $^{1)2}$ が整備され始めたことにより,アーカイビング,情報検索,要約,字幕化,翻訳など様々な「音声ドキュメント処理」に関する研究が精力的に進められるようになった.

音声情報検索は、このような音声認識技術のアプリケーションの1つとして重要視され、これまでにも様々な手法、システムが構築されてきた。しかしながら、現状の音声認識技術では、あらゆる音声データから検索に必要な索引情報を、精度よく抽出することが困難なこともあり、Google、Yahoo!等の代表されるテキストの検索のように日常的に利用されるには至っていない。

一方、最近では Web 上において、音声版のブログ (Weblog) ともいえる「ポッドキャスト」や動画共有サービス等が普及したことで、音声データに対する検索の重要性がより一層増してきたといえる。このような状況の中、我々は現在、ポッドキャストを音声認識によって自動的にテキスト化することで、それらをユーザが全文検索できるだけではなく、詳細な閲覧、編集も可能な音声情報検索システム「PodCastle³⁾⁴⁾⁵⁾」の開発を行っている。本研究では、そのようなシステムを実現、運用するために、実環境音声データであるポッドキャストの音声認識手法について検討を行っている。

従来,実環境音声の認識に関する研究においては,放送ニュース⁶⁾,学会講演¹⁾,大学講義⁷⁾,国会中継⁸⁾等を対象とした報告がなされており,いずれの場合においても,それぞれのタスク,ドメインにマッチしたコーパスを事前に用意し,音声認識器(音響モデル,言語モデル)を学習することで大きな改善が得られている.一方,本研究で対象とするポッドキャストは,その発話内容や録音環境などが多種多様であるという特徴を持っている.そのため,ポッドキャスト上で出現する全てのタスク,ドメインに対して,従来研究のようにコーパスを事前に構築することは現実的に不可能である.したがって,ポッドキャスト音声認識においては,事前コーパスに依存することなく,いかに高精度な音響モデル,言語モデルを構築,学習するかが性能向上への鍵となる.特に言語モデル(N-gram)は単純なモデル構造であるが故に,音響モデルに比べ,学習データにより強く依存する傾向があり⁹⁾,ポッドキャスト音声認識性能を劣化させる大きな要因となっている.

PodCastle では、最も特徴的な機能として、多数のユーザに認識誤りを訂正 (アノテーション) する協力をしてもらうことで、音声認識率をシステムの運用中に向上させる枠組みを採用している。こうすることで、検索サービスとしての質を向上させるだけでなく、音声認識



図 1 1 つのポッドキャストの構成 (RSS の例)

の基盤技術の向上も狙っている。本研究では、このような枠組みの一環として、PodCastle を通じて得られる集合知、すなわちユーザによる音声認識誤りの訂正結果を活用した言語モデル学習手法により、ポッドキャスト音声認識の性能向上をはかる。ここでは、集合知により形成される独自コーパスの利用形態により、2つの言語モデル学習アプローチについて提案し、それぞれの有効性について評価する。

2. 集合知により生成されるポッドキャストコーパス

2.1 ポッドキャストの構成

ポッドキャストは音声版のブログ (Weblog) ともいえる Web 上のコンテンツであり、個人の日記から大学の授業、ニュースまで多岐に渡る内容の音声データが日々配信されている。ポッドキャストは、一連のエピソードと呼ばれる音声データ (MP3 ファイルなど) に加え、その流通を促すためにブログなどで更新情報を通知するために用いられているメタデータ RSS が付与されている (図 1). PodCastle では、ユーザ側で任意のポッドキャスト (RSS) を登録することができ、登録したポッドキャストで新たなエピソードが配信される毎に、音声認識が自動的に開始されるようになっている。

2.2 訂正状況の分析とポッドキャストコーパス

PodCastle では、図2に示すように競合候補のリストという形で訂正インタフェース¹⁰を提供している。ユーザは本インタフェースを通じて、認識誤りが見つかれば、「候補選択」、「タイプ入力」のいずれかの手段で訂正を行う。現在の状況として、サイトに登録されているいずれかの音声に対して、ほぼ毎日のペースで訂正がなされており、ポッドキャストによっ



図 2 音声訂正インタフェース

表 1 PodCastle における各データ量

登録済みポッドキャスト数	621
エピソード (mp3 音声ファイル)	63278
訂正されたエピソード数	2022

ては全エピソードのほぼ全ての認識誤りが訂正されているものもある。ここで、PodCastle の 2009 年 11 月 15 日時点における各種データ量を表 1 に示す。「訂正されたエピソード数」とは、1 つのエピソード中で訂正が 1 カ所でも行われているものをカウントしている。また、訂正された 2022 のエピソードにおける訂正単語数 *1 、全訂正のうち各訂正手段がどの程度利用されたかの内訳を表 2 に示す。ボランティアベースにも関わらず、2009 年 11 月 15 日時点までに 44 万単語もの訂正が得られている。訂正手段の内訳としては、候補選択による訂正がより多く利用されていた。実際の利用状況を確認すると、1-best の認識結果に誤りが多くても、競合候補中に本来の正解がある程度含まれていると、候補選択による訂正インタフェースが積極的に利用される傾向にあった。ただし、現状ではそもそも競合候補にも正解がほとんど含まれないような認識困難なポッドキャストも多く存在し (例えば芸能人の会話番組など)、そのようなデータではタイプ入力が主な訂正手段となっている。

このように、実際にユーザが訂正することで生成された書き起こし、並びにそれらに対応 する音声データによりポッドキャストコーパスが構成される。前述したようにポッドキャス

^{*1} ここでの「単語」は訂正インタフェース (図 2) における区間を表す. ただし訂正時には 1 つの区間に対して複数単語、もしくはフレーズが入力されることもあるため、全てが厳密な単語に対応するわけではない.

表 2 訂正に関する分析

総訂正単語数	440570
平均訂正単語数/エピソード	218
訂正手段 (候補選択)	227293
訂正手段 (タイプ入力)	213277

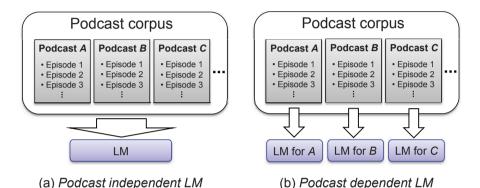


図 3 2 つの言語モデル学習アプローチ (LM: Language Model). 実際は、各言語モデルは他タスク言語モデルと 混合される.

トは多種多様な音声データであるため、本コーパスはタスクに非依存な特性を持っている. また最も重要なポイントとして、Web サービスの利用に伴って音声、テキストデータともに日々蓄積され、コーパスの規模が徐々に拡大していくことが挙げられる.

3. 集合知を活用した言語モデル学習

前章で述べたポッドキャストコーパスを利用した言語モデルの学習について説明する.ここでは、ポッドキャスト音声認識における言語モデルの汎用性と適応性の観点から以下で述べる2つの学習アプローチを導入する.

3.1 ポッドキャスト非依存言語モデリング

PodCastle においては、様々なタスク、ドメインの書き起こしデータが収集される。本アプローチは収集されたポッドキャストコーパス全体を利用することで、様々なポッドキャストに対して一定の性能が得られる、汎用的な言語モデルの構築を目指すものである(図 3-(a)). なお、一般的に音声認識のための統計的言語モデルにおいては、幅広いタスクの膨

大なテキストを学習に利用するよりは、小/中規模でもよいので該当するタスクに特化した テキストのみで学習する方が性能/効率ともによいとされる¹¹⁾. ただしポッドキャスト音声 認識では、事前に対象音声のトピックや発話スタイル、語彙などの特定ができないため、(少 なくとも初期段階の認識においては)このような汎用性を目指すアプローチが有効に働くと 考えられる. また、後述するポッドキャスト依存言語モデリングにおける初期モデルとして も有効となる. システムの概要を以下に示す.

まず、事前に複数のコーパスを用いてそれぞれのタスク、ドメインに依存した言語モデル (要素モデルと呼ぶ)を構築する。本研究では、日本語話し言葉コーパス (CSJ)、Web ニュースコーパス、そして前節で述べたポッドキャストコーパスを利用することで3種類の要素モデルを学習する。CSJ では主に話し言葉への対応、Web ニュースコーパスでは幅広いトピック、最新の単語への対応をそれぞれねらっている。一方、ポッドキャストコーパスにおいては、様々なタスクにおける音声データの書き起こしが得られるため、幅広いトピック、講演口調以外の様々な発話スタイルや会話等への対応などが期待できる。

次に、得られた要素モデルをモデルベースで混合することで、最終的な言語モデルを構築する。ここで混合手法としては線形補間法を用いる。線形補間では、例えば2つの要素モデルを LM_1 , LM_2 とすると下記に示す確率を持つN-gram を生成する。

$$p(w|h) = \lambda p_1(w|h) + (1 - \lambda)p_2(w|h) \tag{1}$$

ここで、 p_1 は LM_1 、 p_2 は LM_2 における N-gram 確率を表す。 λ は補間パラメータ (重み) であり、一般的に評価セットと同一タスクのヘルドアウトセットを用いて最適化する。線形補間法では、この補間パラメータをいかに最適に設定するかが重要となる。本研究では、現段階までに蓄積されたポッドキャストコーパスの汎用性を評価するために、様々なトピック、発話スタイルのポッドキャストから構成される開発セット (development set) を補間パラメータの最適化に用いる。最適化には、開発セットのパープレキシティが最小となるように、EM アルゴリズムによる繰り返し推定を行う。

3.2 ポッドキャスト依存言語モデリング

ポッドキャスト音声認識の更なる性能向上につなげるために、ポッドキャストにおける 個々のタスク、ドメインに特化した(適応された)言語モデルの構築を目指す.

本研究では、その1つとして PodCastle に登録されている個々のポッドキャストごとに言語モデルの学習を行う、ポッドキャスト依存言語モデリングを導入する(図 3-(b)). この理由として、同一ポッドキャスト中の各エピソードは、同じ言語的特性(トピック、発話スタイル等)を持っている可能性が高いことが挙げられる. また図1に示すポッドキャスト

情報処理学会研究報告 IPSJ SIG Technical Report

(RSS) の仕組みにより、認識対象となる各音声ファイルがどのポッドキャストのエピソードなのか、すなわち各音声ごとにどの言語モデルを認識時に適用すべきか自明であるという利点もある。ポッドキャスト依存言語モデリングの流れを以下に示す。

- (1) あるポッドキャスト A において、訂正がなされたエピソードの全書き起こしを利用して、ポッドキャスト A の要素モデルを学習する.
- (2) ポッドキャスト A の要素モデルとベースライン言語モデル (ここでは CSJ と Web ニュースを混合したモデル) を線形補間する. この段階では初期パラメータとして $\lambda=0.5$ と設定する.
- (3) 対象となる音声データ (新たなエピソード) に対して, (2) で構築した言語モデルにより音声認識を行い、認識仮説を生成する。
- (4) 認識仮説を開発セットとして、EMアルゴリズムにより補間パラメータの最適化を行う (対象音声データに特化した補間パラメータ).
- (5) 最適化された補間パラメータにより再度線形補間を行い、ポッドキャスト A 依存の言語モデルを構築する.

4. 実 験

前章で述べた2種類の言語モデリングアプローチに対してそれぞれ評価実験を行う.

4.1 ポッドキャスト非依存言語モデリングの評価

本実験で利用する各学習セット、評価セット、開発セットについて表 3 にまとめる. 評価セットは 7 つのポッドキャストで構成されており、主に、ニュース、経済コラム、対談、雑談、4 種類に内容的に分類できる. 経済コラムは、発話スタイルとしては学会講演に近いが、雑音 (背景音楽)を含んでいる. また、対談、雑談については日常会話と同等のカジュアルな発話スタイルであり、発話速度も速いため、非常に認識が困難なタスクとなっている. 開発セットは、同じ 7 つのポッドキャストのうち、評価セットに含まれないエピソードで構成される.

本実験で構築した全ての言語モデルは 3-gram である。 CSJ は,日本語話し言葉コーパス (CSJ) の講演書き起こし,Webnews は 2006 年 8 月~2009 年 11 月における Yahoo!ニュースの記事である.PodcastALL はポッドキャストコーパス中のデータであり,評価セットのポッドキャストは含んでいない.PodcastALL はユーザによる訂正箇所以外は全て音声認識 結果となっている.ここでは比較として,ユーザによる訂正がなされた発話 (文) のみを学習に利用する学習セット PodcastUSER も考慮する.各モデルの評価には,評価セットに対

表 3 学習、評価セット各データ量

	文数	単語数	ポッドキャスト数	エピソード数
CSJ	389,309	7,043,529	-	-
Webnews	18,190,294	456,017,101	-	-
PodcastALL	412,047	4,350,825	253	1,984
PodcastUSER	83,672	1,467,596	253	1,984
評価セット	2,088	49,634	7	18
開発セット	1,728	35,346	7	16

表 4 各言語モデルの評価

	PP	APP(#OOVs)	WER
CSJ	153.94	245.46(3455)	52.22%
Webnews	548.36	583.80(594)	58.20%
CSJ + Webnews	180.49	189.23(459)	47.86%
+PodcastHypo	155.82	161.65(371)	47.04%
+ PodcastUSER	155.74	161.81(383)	46.63%
+PodcastALL	151.11	156.33(348)	46.81%

するパープレキシティ(PP) と補正パープレキシティ(APP), 単語誤り率 (WER) を用いた. 補正パープレキシティは, 評価データ中の未知語 (OOV) 出現率を考慮したパープレキシティであり, 下記の式で表される $^{12)13}$).

$$APP = (P(w_1 \dots w_n)m^{-n_u})^{-\frac{1}{n}}$$
 (2)

ここでn は評価セットの総単語数, n_u,m は評価セット中の未知語の総数, 種類数をそれぞれ示す.

本実験では、音声認識システムとしてはシンプルなワンパスのデコーディングを用いており、音響モデルの教師なし適応等による多段処理は行っていない.

4.1.1 実験結果

各言語モデルの評価を表 4 に示す。まず上段の CSJ,Webnews の要素言語モデルについて考察する。CSJ,Webnews を単独で用いた場合,PP,OOV に関してそれぞれ一長一短があることがわかる。Webnews は話し言葉でなく基本的に書き言葉であるため,パープレキシティ自体は大きくなるが未知語を劇的に減らすことができている。これら 2 つを線形補間したモデル CSJ+Webnews は,それぞれの特徴を表現したモデルとなっており,APP,WER ともに大きく削減できていた。

表中、下段3つはCSJ+Webnewsとそれぞれ線形補間したモデルの結果である.Pod-

表 5 評価セット

ID	カテゴリ	エピソード数	時間 (sec.)
Α	ニュース	2	2282.56
В	雑談 (独話)	2	2845.26
C	コラム	2	846.76

castHypo は、ユーザの訂正結果を含まず、全て認識結果を学習テキストとしたモデルである。ただし、ここでの認識結果は、PodCastle の実際の運用上で得られた認識結果であり、本実験と同一の認識システム、モデルを利用して得られたものではない。ポッドキャストコーパスの要素モデルを線形補間することで、いずれの場合においても CSJ+Webnews と比較して性能が改善していることがわかる。評価セット中の全てのポッドキャストに対して改善が得られ、特に対談や雑談など、CSJ+Webnews ではカバーできない発話スタイルや言い回しに対する改善が大きかった。パープレキシティ、補正パープレキシティについてはPodcastALL が最もよい性能を示したものの、WER についてはユーザからの訂正を含む発話のみを用いた PodcastUSER が最も削減率が大きい結果となった。これは PodcastALL では、PodcastUSER が最も削減率が大きいお果となった。これは PodcastALL では、PodcastUSER と比較して、単純にデータ量が多いため PP、APP ではより良い値を示したものの、誤認識単語が言語モデル学習の際に考慮されるため、その結果 WER が高くなったと考えられる。これに対しては今後、認識仮説の信頼度等の基準による発話選択処理を導入することにより更なる向上が見込める。

4.2 ポッドキャスト依存言語モデリングの評価

次に、ポッドキャスト依存言語モデリングの有効性を音声認識実験により評価する.本実験では3種類のポッドキャストを評価セットとして用いた(表 5). Aは、読み上げ音声であるが、一部の区間に背景音楽が存在する. Bは、女性声優による雑談(一般的なラジオ番組)であり、内容はエピソードごとに様々である. Cは、男性芸能人によるコラムであり、内容はエピソードごとに様々である. B、Cのデータは、ともに自由発話音声である.以上の3つのポッドキャストは、実際にユーザから比較的多くの訂正がなされている.本実験では、言語モデル学習セットとして、これらの各ポッドキャストにおいて訂正が1箇所でもなされた全てのエピソードを用いた(表 6).ただし、評価セットのエピソードは学習セットには含んでいない.

音声認識には PodCastle 音声認識システムを用いた 14)。音響モデルは、日本語話し言葉 コーパス (CSJ) から学習された triphone モデルである。ベースラインとなる言語モデルは、Web キーワードベースの N-gram 15)であり、Web ニューステキスト、CSJ の講演書き起

表 6 学習セット

ID	エピソード数	単語数
Α	67	270,447
В	56	283,414
С	30	39,098

表 7 各手法での WER(%)

ID	ベースライン	訂正なし	訂正あり
Α	16.88%	16.24%	14.49%
В	30.98%	28.52%	24.61 %
С	35.16%	33.22%	$\boldsymbol{26.24\%}$

こしを用いて学習したものである。また、評価用音声データの各エピソードごとに、繰り返し教師なし MLLR 適応を行っている。本音声認識システムのデコーディングは段階的探索に基づいている。まず、bigram を用いた N-best デコーディングにより単語グラフを生成する。次に、trigram を用いて、生成された単語グラフを、trigram 制約の単語グラフに拡張する。最後に、trigram 制約の単語グラフに対して、consensus デコーディング 16)を行い、confusion network 中の最大候補を最終の認識結果とした。

4.2.1 実験結果

実験結果を表 7 に示す. ここでは比較として、ポッドキャスト依存言語モデル学習を行わないベースライン認識システム (「ベースライン」)、ユーザによる訂正結果を用いないポッドキャスト依存言語モデル学習 (「訂正なし」)の結果も併せて示す.「訂正なし」では、音声認識結果のテキストのみを用いて言語モデル学習を行う. 実験結果より、ポッドキャスト依存言語モデリングによっていずれのポッドキャストに対しても性能改善が得られていることがわかる. 特に C のポッドキャストに対する改善が大きいが、これは C の話者が独特の発話スタイル、言い回しを持っているためであった (C の番組内容はクイズ形式のコラムであり、他のポッドキャストと比べても異質な内容であった). B に関しては、内容はエピソードにおいて様々であるが、独自の言い回しや冒頭、終了部分での特定のフレーズ、曲名やアルバム名等の専門用語に対して特に改善がみられた. A は基本的に読み上げ音声のニュース番組であり、言語的要因による改善効果は限定的である. 冒頭や終了時の背景音楽が存在する区間以外は、ベースラインにおいても高い認識性能が得られているため、訂正ありと訂正なしとの差は B、C ほど大きくはなかった.

5. おわりに

本稿では、ポッドキャスト音声認識を改善するための言語モデル学習手法について検討した. ポッドキャストのように、幅広いタスク、多様な言語的特性を持つ音声データに対し、高精度な言語モデルを学習することは従来困難であった. それに対し、本研究では、Webサービスを通じて得られる集合知 (認識誤りの訂正)の利用した 2 種類の言語モデル学習アプローチを導入することで、性能を改善できることを示した. なお、本報告においては、ポッドキャスト非依存言語モデリング、ポッドキャスト依存言語モデリングはそれぞれ独立した実験にて評価を行ったが、これら 2 つのアプローチは併用して更なる改善を得ることも可能である.

我々は以前の報告で、4.2節の評価実験と同等の条件にて音響モデル学習の性能改善に関する実験を行った(集合知に基づくポッドキャスト依存音響モデル学習)¹⁷⁾. 今回の実験は全体の傾向として、音響モデル学習の性能改善に比べて数値的には低い結果となっている。ただし、言語的要因による改善では、個々のポッドキャスト独特のトピックやキーワード、前述したような専門用語など音声データ中の重要箇所に直接影響するため、改善による意義は大きいと考えられる。また、PodCastle においては、ユーザがあるエピソードの訂正時に入力した特徴的な単語が、本研究の言語モデル学習によって別エピソードで正しく認識できるようになることで、訂正してくれたユーザに対してより良い印象を与えることができる。加えて、それによってユーザ自身がシステムへの貢献を実感することで、さらなる訂正の促進につながる可能性もある。

今後は、より大規模な実験を行うとともに、音響モデル学習との併用によるさらなる性能 改善についても調査していく.また、より高度な言語モデル補間手法¹⁸⁾、発話選択処理の 導入なども行う予定である.

参考文献

- 1) 河原達也:『日本語話し言葉コーパス』を用いた音声認識の進展,第3回話し言葉の科学と工学ワークショップ講演予稿集 (2004).
- 2) Takeda, K., Fujimura, H., Itou, K., Kawaguchi, N., Matsubara, S. and Itakura, F.: Construction and Evaluation of a Large in-car Speech Corpus, *IEICE Transactions on Information and Systems*, Vol.E88-D, no.3, pp.553–561 (2005).
- 3) 緒方 淳,後藤真孝, 江渡浩一郎: PodCastle: ポッドキャストをテキストで検索, 閲覧, 編集できるソーシャルアノテーションシステム, WISS 2006 論文集, pp.53-58

(2006).

- 4) Goto, M., Ogata, J. and Eto, K.: PodCastle: A Web 2.0 Approach to Speech Recognition Research, *Proc. of Interspeech 2007*, pp.2397–2400 (2007).
- 5) 後藤真孝,緒方 淳, 江渡浩一郎: PodCastle: ユーザ貢献により性能が向上する音声情報検索システム,人工知能学会論文誌, Vol.25, No.1, pp.104-113 (2010).
- 6) 今井 亨, 小林彰夫, 佐藤庄衛, 本間真一, 奥 貴裕, 都木 徹: 放送用リアルタイム字幕制作のための音声認識技術の改善, 第2回音声ドキュメントワークショップ講演論文集, pp.113-120 (2008).
- 7) 小暮 悟, 西崎博光, 土屋雅稔, 富樫慎吾, 山本一公, 中川聖一: 日本語講義音声コンテンツコーパスの構築と講義音声認識手法の検討, 第2回音声ドキュメントワークショップ講演論文集, pp.7-14 (2008).
- 8) Akita, Y., Mimura, M. and Kawahara, T.: Automatic Transcription System for Meetings of the Japanese National Congress, *Proc. of Interspeech 2009* (2009).
- 9) Lefevre, F., Gauvain, J.-L. and Lamel, L.F.: Genericity and portability for task-independent speech recognition, *Computer Speech & Language*, Vol.19, pp.345–363 (2005).
- 10) 緒方 淳,後藤真孝:音声訂正:選択操作による効率的な誤り訂正が可能な音声入力 インタフェース,情処学論, Vol.48, No.1, pp.375-385 (2007).
- 11) 伊藤彰則:音声認識における言語モデル, 日本音響学会誌, Vol.66, No.1, pp.32-35 (2010).
- 12) Ueberla, J.: Analysing a simple language model some general conclusions for language models for speech recognition, Computer Speech & Language, Vol.8, No.2, pp.153–176 (1994).
- 13) 中川聖一,赤松裕隆:未知語を含む文集合のパープレキシティの算出法-新補正パープレキシティー,日本音響学会講演論文集 (1998).
- 14) Ogata, J., Goto, M. and Eto, K.: Automatic Transcription for a Web 2.0 Service to Search Podcasts, *Proc. of Interspeech 2007*, pp.2617–2620 (2007).
- 15) 緒方 淳, 松原勇介, 後藤真孝: PodCastle: 集合知に基づく Web キーワードを考慮 した言語モデリング, 日本音響学会講演論文集 (2008).
- 16) Mangu, L., Brill, E. and Stolcke, A.: Finding Consensus in Speech Recognition: Word Error Minimization and Other Applications of Confusion Network, *Computer Speech & Language*, Vol.14, No.4, pp.373–400 (2000).
- 17) Ogata, J. and Goto, M.: PodCastle: Collaborative Training of Acoustic Models on the Basis of Wisdom of Crowds for Podcast Transcription, *Proc. of Interspeech* 2009, pp.1491–1494 (2009).
- 18) Hsu, B.-J.P.: Generalized linear interpolation of language models, $Proc.\ ASRU$ (2007).