

[招待講演] 音声認識の方法論の変遷と展望 —Acoustic-to-Word モデルを中心に—

河原 達也[†]

† 京都大学情報学研究科
〒606-8501 京都市左京区吉田本町
E-mail: †kawahara@i.kyoto-u.ac.jp

あらまし 音声認識の方法論は深層学習、特に End-to-End モデルの導入で大きく変わりつつある。本稿では、従来の方法論を概観し、End-to-End モデルに至るまでの変遷を述べる。単語単位の End-to-End モデルである Acoustic-to-Word モデルは、音響特徴量系列から単語列を直接求めるもので、音響モデルと言語モデルを内包し、発音辞書や複雑な認識プログラムを必要としない革新的な方式である。この方式の課題と解決法についても述べる。

キーワード 音声認識, End-to-End モデル, Acoustic-to-Word モデル

[Invited Talk] Review of Automatic Speech Recognition Methodology —Outlook of Acoustic-to-Word Model—

Tatsuya KAWAHARA[†]

† School of Informatics, Kyoto University
Sakyo-ku, Kyoto 606-8501, Japan
E-mail: †kawahara@i.kyoto-u.ac.jp

Abstract The methodology of speech recognition has been changing due to the introduction of deep learning, in particular end-to-end modeling. This article gives a brief overview of the conventional methodologies leading to the end-to-end models. Word-based end-to-end model, referred to as acoustic-to-word model, directly converts a sequence of acoustic features into a word sequence. This model contains acoustic and language models, and does not require a pronunciation lexicon and a complex decoding program. The problems of this new promising model and current solutions are also described.

Keywords Speech Recognition, End-to-End Model, Acoustic-to-Word Model

1. はじめに

音声認識の研究の歴史は長いですが、著者が研究を開始した約 30 年前（1980 年代後半）においても、基礎研究/SF の範疇であった。ちょうどその頃に、確率的な枠組みと統計的モデル[1]が隆盛し、その後加速度的に実用化まで展開された。さらに、この 10 年の間に、音声認識を取り巻く情勢は大きく変化している。

1 つは、認識性能の劇的な進展である。大規模データベース（ビッグデータ）の蓄積[2-3]と深層学習（ディープラーニング）の導入[3-6]により、パターン認識の多くの問題は解決されつつある。音声認識も、静かな環境で「普通に」発話すればほぼ正しく認識されるレベルになっている。2016 年には、IBM と Microsoft が Switchboard 電話会話音声認識のタスクで約 95% の認識率を達成し、「人間と同等認識精度を実現した」と発表した。「人間と同等」の能力については様々な議論があると思われるが、これらのシステムはかなり複雑な構成になっている。

これに対して、もう 1 つの流れとして、End-to-End の枠組みが盛んに研究されるようになった。音声認識では、音響特徴量系列から文字/単語列への系列間 (sequence-to-sequence) 写像学習として定式化される。特に、単語列を直接求める Acoustic-to-Word モデルはきわめて単純で、これまでの音声認識の枠組み・アーキテクチャを一変するものである。本稿では、この方式に至るまでの変遷を概観した上で、現在の課題と解決法についても述べる。

2. 音声認識の方法論の変遷

音声認識研究の歴史は 60 年以上に渡る。京都大学では 1960 年頃に単音節単位の認識を行う「音声タイプ」が構築されている[7]。その後、音声認識に有効な音響特徴量と、DP マッチングに代表される動的パターンのマッチング手法に関する基礎的な研究が行われた。これらは、テンプレートベースの方法であり、多数話者のバリエーションをモデル化するには不十分であった。

表 1 音声認識の方法論の変遷

第 1 世代	1950～1960 年代	ヒューリスティック
第 2 世代	1960～1980 年代	テンプレート (DP マッチング, オートマトン)
第 3 世代	1980～1990 年代	統計モデル (GMM-HMM, N-gram)
3.5 世代	1990～2000 年代	統計モデルの識別学習
第 4 世代	2010 年代	ニューラルネット (DNN-HMM, RNN)
4.5 世代	2015 年～	ニューラルネットによる End-to-End

これに対して、確率的なモデルを導入することにより解決が図られた。DP マッチングを拡張した形で隠れマルコフモデル(HMM)が導入され、その改良が様々に行われた。まず、HMM の各状態の音響特徴量のパターンを混合正規分布(GMM)でモデル化することが導入された。そして、これを最尤推定する代わりに、識別誤りを最小化するように学習(識別学習)する方法が検討された。2000 年代に実用化された音声認識システムは、GMM-HMM の識別学習に基づくものである。一方、言語モデルについては、単語の接続規則(文法)をオートマトンで記述したもから、その遷移を確率的なものにし、その確率をコーパスから最尤推定する N-gram モデルに移行した。

その後、深層学習の導入が進められた。音声認識にニューラルネットワークを用いることは 1990 年頃から研究されていたが、主流になったのは 2010 年以降である。音響モデルについては、GMM による確率計算をディープニューラルネットワーク(DNN)に置き換えた DNN-HMM が、言語モデルについては、リカレントニューラルネットワーク(RNN)を N-gram と併用するモデルが一般的になっている。最近では、RNN を発展させた LSTM (Long Short-Term Memory) のみで音声認識を行う End-to-End の枠組みが研究されている。

以上の変遷をまとめたのが表 1 である。世代の定義は古井[8]に従ったものであるが、第 4 世代以降は著者が追加したものである。

3. 音声認識の定式化とアーキテクチャ

音声認識は、音声の特徴量の時系列 X が与えられたときに、その単語列 W を同定する問題である[1]。すなわち、 $p(W|X)$ を最大化する W を求めればよい。ただし、 X も W も系列である点が、画像認識などの多くのパターン認識の問題と異なる点である。なお、 W については、単語列でなく文字列でもよいが、語彙の知識は認識の上できわめて有効であるし、単語がわからないと、検索・翻訳・対話などの後段の処理も意味をなさない。単語単位の系列が一般に用いられる。

3.1. ベイズ則に基づく階層モデル

従来は、この $p(W|X)$ をベイズ則で書き換えることにより定式化してきた。

$$p(W|X) = \frac{p(W)p(X|W)}{p(X)} \quad (1)$$

ここで W の同定において、分母 $p(X)$ は無関係なので、分子の 2 つの項のみを計算すればよい。 $p(W)$ は(その言語あるいは状況において)単語列 W が生成される先験的な確率 (=言語モデル確率) であり、 $p(X|W)$ は単語列 W から音声の特徴量 X が生成される確率 (=音響モデル確率) である。さらに、各単語は音素などのサブワード単位 S でモデル化され、単語と音素の関係は辞書で決定的に与えられる ($p(S|W) = \{1, 0\}$) ので、以下のようになる。

$$p(W)p(X|W) = \sum_S p(W)p(S|W)p(X|S) \approx \max p(W)p(X|S) \quad (2)$$

この定式化は、1990 年頃から標準的になり、その後四半世紀以上にわたって、世界中の言語において普遍的に用いられ、あらゆる教科書の冒頭に記述されてきた。この理由として、 $p(W|X)$ の確率分布を直接推定することが難しいこと、及びパターン認識処理で実現される音響モデルと自然言語処理で実現される言語モデルに分離した方が問題を扱いやすいことが挙げられる。

これに基づく音声認識システムのアーキテクチャを図 1 に示す。各モデルが階層的に構成され、各モデルで計算される確率を順次結合することによって $p(W|X)$ が計算される。これを本稿では階層モデルと呼ぶ。ここでは、各モデルは生成モデルとして定式化・学習される。

3.2. End-to-End モデル

これに対して図 2 に示すように、上記のモジュールを統合的にニューラルネットワークで構成して、 $p(W|X)$ あるいは $p(S|X)$ を直接推定する End-to-End の枠組みが研究されている。図 2 では、音響モデルを音素状態の認識部(DNN)と音素などのモデル(HMM)にさらに分けている。

図1 従来の音声認識システムの構成

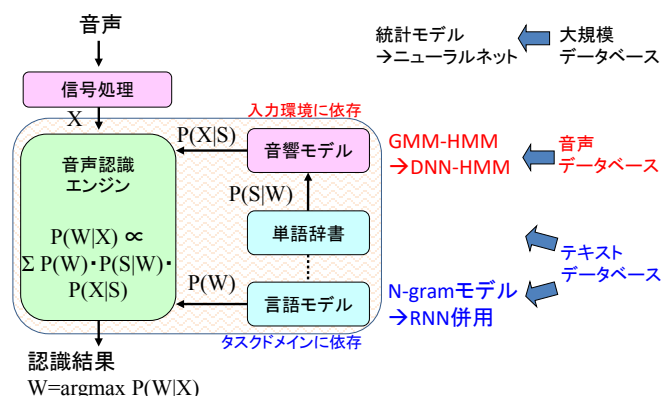
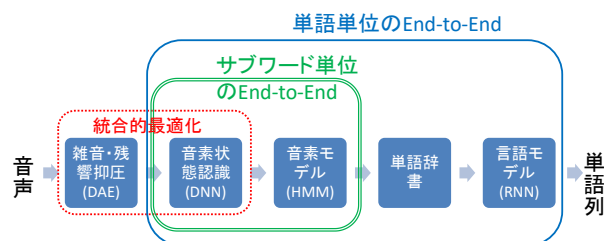


図2 End-to-End音声認識



この枠組みで最も早く実現されたのが、音素や文字などのサブワードを出力の単位とする LSTM を用いて、その出力系列を縮約する CTC (Connectionist Temporal Classification) である (5.1 節参照)。これとは別に、LSTM で入力系列をいったん符号化した後に、サブワード系列に復号化する注意機構モデル (5.2 節参照) も導入されている。これらのモデルはサブワード単位の言語モデルを暗黙的に包含し、 $p(S|X)$ を直接計算するものと捉えられる (図 2 の二重線枠部分; 5 章で詳述) が、単語辞書や言語モデルはこの後に適用する必要がある。

さらに最近では、これらのモデルを発展させて、単語を出力の単位としたモデル化も検討されている。これにより、 $p(W|X)$ を直接計算することができる (図 2 の実線枠部分; 6 章で詳述)。

これとは別に、雑音や残響抑圧のためのニューラルネットワーク (デノイジングオートエンコーダ (DAE) など) で構成) と音素状態認識の DNN を統合して学習することも研究されている (図 2 の点線枠部分)。

本稿ではフロントエンドは扱わず、音素状態認識以降の処理 (図 2 の実線枠内) について述べる。

4. 階層モデルにおける深層学習の導入

深層学習の音声認識への導入は当初、階層モデルの枠組みで実現された。最も直接的な方法は、HMM における各音素状態の GMM による確率計算を DNN に置き換えるもので、DNN-HMM ハイブリッドモデルと呼ばれる [9,10]。

4.1. DNN-HMM ハイブリッド音響モデル

DNN-HMM における DNN の入力音響特徴量であり、当該フレームの前後のフレームの周波数特徴量も結合して与える。出力層のノードは、HMM の各状態に対応付けられるが、一般的な音声認識ではトライフ

オンモデルの共有状態となる。これは、先行音素と後続音素の文脈を考慮したもので、クラスタリングを行っても数千個のオーダになる。音声認識で必要となるのは、式 (2) の $p(X|S)$ であるが、DNN で計算される出力確率は通常 Softmax 関数を経た事後確率 $p(S|X)$ の形となるので、事前確率 $p(S)$ で除して、HMM に渡す。

DNN-HMM の学習手順は以下の通りである。

- (1) GMM-HMM を学習する。これは、トライフォン状態のクラスタリングを含む。
- (2) 学習データベースをトライフォン HMM の状態でアライメントする。
- (3) アライメントされたデータ (音響特徴量と HMM の状態ラベルの対) を用いて、DNN を学習する。

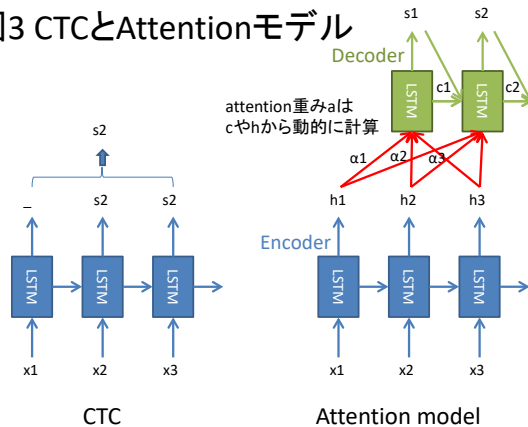
HMM の状態遷移確率は (1) で推定され、各状態の事前確率は (2) で推定されるものを用い、出力確率を DNN で計算する。DNN の学習を行うための音素状態ラベルは、GMM-HMM によって作成されていることに留意する。

各種の音声認識タスクにおいて、DNN-HMM が GMM-HMM を大きくしのぐ (誤り率で概ね 20~30% の削減) ことが示され、普遍的に用いられるようになった。その後も、CNN や ResNet、さらに LSTM などを多層に積み重ねるなどの改良が行われている。

4.2. RNN 言語モデル

音響モデルだけでなく、言語モデルにおいてもニューラルネットワークの導入が進められており、RNN が一般に用いられている。ただし、入力単語を少ないノードの数値データ (分散表現) に射影する層を別途用意する。中間層はこれと履歴を符号化したものと捉えられ、N-gram モデルと比べて非常に長い履歴を考慮することができる。ただし、N-gram モデルの方が低頻度語のスムージングが効果的に行えることもあり、N-gram モデルと併用・線形補間する場合が多い。

図3 CTCとAttentionモデル



RNNを発展させたものとして、種々のゲートを導入したLSTM (Long Short-Term Memory) がある。LSTMでは、内部メモリで1つ前の状態(分散表現)を記憶しつつ、入力・記憶の利用及び出力の各過程においてゲートを設定し、情報の制御を行っている。このゲートの値(2値に近い)も入力と1つ前の出力を元に決まるように学習される。LSTMは次節以降に述べるEnd-to-Endモデルの基盤となっている。

5. End-to-Endモデルによる系列写像学習

階層モデルのように多数の処理を経るのではなく、入力から出力への写像を直接(統合的に)学習するものをEnd-to-Endモデルと呼ぶ。特に、音声認識では、入力が音声の時間フレーム長 T で規定される音響特徴量の系列、出力が音素・文字・単語などの記号列であり¹、このように長さの異なる時系列間(sequence-to-sequence)の写像を行うものをEnd-to-Endモデルと呼ぶ。音声合成(単語列→音響特徴量系列)や機械翻訳(単語列→単語列)でも同様の定式化が行える。

音声認識におけるEnd-to-Endモデルは、HMMを一切用いずに実現する。その方法として、現在主に採用されているのが以下の2つである。両者の概念を図3に示す。

5.1. Connectionist Temporal Classification (CTC)

HMMを用いることなく、ニューラルネットワークのみで時系列パターンを分類しようと定式化されたのがCTCである[11]。CTCでは通常、音素単位や文字単

¹ 「真のEnd-to-End」は音声波形を入力とすべきとする向きもあるが、本稿では音響特徴量系列を入力 X とする。ヒトの蝸牛でも周波数分析を行っている一方で、我々が通常扱う音声波形はマイクロフォンで入力されたものである。

位のLSTMが用いられる。ただし、どの音素・文字でもない空白記号($_$)を導入し、各音素・文字の間に挿入する。入力時間フレーム毎にこれらの記号が出力され、これを集積する。この際に、空白記号を消去し、時間的に連続した同一の出力記号を1つに縮約する操作を行う(図3左)。

これを確率的に定式化すると、各フレームの音響特徴量 x_t に対して各記号 s_t の事後確率を計算し、同じ記号列 S に縮約されるものの総和を計算することになる。

$$p(S|X) = \sum_{g \rightarrow S} p(g|X) = \sum_{g \rightarrow S} \prod_{t=1}^T p(s_t | x_t) \quad (3)$$

例えば、以下の記号系列はすべてhaiを表現するものとしてまとめられる。

_h____a_____i_____
 _hh____aa_____ii_____
 _h____aaaa_____iii_____

モデル学習の際には、正解記号系列 S (上記の例ではhai)に縮約されるすべての系列の尤度の総和を求める。この尤度計算は、HMMの尤度計算と同様に、前向き・後向きアルゴリズムで効率的に実現できる。この対数尤度を基に、LSTMの各パラメータを更新する学習則が導出される。

なお、CTCではLSTMにより時間フレーム間の依存性はモデル化されるが、式(3)においては記号間の関係は独立に扱われている。

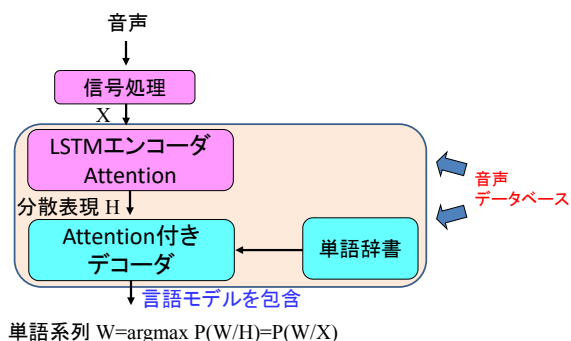
5.2. 注意機構(Attention)モデル

注意機構モデルは、正確には注意機構付きエンコーダデコーダ(encoder-decoder)モデルであり、符号部と復号部から構成される(図3右)[12]。

エンコーダでは、入力フレーム毎に音響特徴量をLSTMにより別の数値ベクトル(分散表現) h_t に符号化する。デコーダは、この符号化された分散表現 h_t の系列(長さ T)を入力として、音素・文字などの記号 s_t の系列(長さ L)を順次予測する。その際に、例えば最初の方の音素は音声の最初の方の情報を用いるのが有用であるので、重み α_t を付ける。これが注意機構であり、この重み自体も動的に計算され、重みを計算するパラメータは統合的に学習される。 c_t はLSTMの内部メモリ状態であり、以下のように更新される。

$$\begin{aligned} c_t &= LSTM(c_{t-1}, g_t, s_{t-1}) \\ g_t &= \sum_i \alpha_{i,t} \cdot h_i \\ \alpha_{i,t} &= \exp(e_{i,t}) / \sum_{i'} \exp(e_{i',t}) \\ e_{i,t} &= f(c_{t-1}, h_i, \alpha_{i-1}) \end{aligned} \quad (4)$$

図4 単語単位のEnd-to-End音声認識



その上で、記号 s_t は内部状態 c_t (と g_t) に基づいて (one-hot ベクトルとして) 計算される。デコーダは、文頭(sos)記号から予測を開始し、文末(eos)記号が出力されると終了する。デコーダ LSTM は次の記号を予測する際に、直前の状態 c_{t-1} に加えて、直前の記号 s_{t-1} を用いており、言語モデルの機能を包含している。

エンコーダ・デコーダ及び注意機構の学習は、正解記号系列と予測記号系列のクロスエントロピーに基づいて統合的に行われる。

以上をまとめると、入力音響特徴量系列 X をいったん分散表現の系列 H に変換した上で、 $p(S|H)=p(S|X)$ が最大となる記号系列 S を出力するモデルであるので、 H を音素状態の尤度に代わるものと捉えれば、エンコーダは音響モデル、デコーダが言語モデルに対応すると考えられる。

また、注意機構は音声と記号の対応付けを行うものと捉えられる。このアライメントを正確に行うために、順方向の LSTM のみでなく、発話末から逆方向の LSTM を併用した双方向 LSTM を用いるのが一般的であり、エンコーダではこれを多層用いることが多い。これは、CTC においても同様である。

6. 単語単位の End-to-End (Acoustic-to-Word) モデル

End-to-End 音声認識に関する研究の大半は、音素や文字を単位として行われていた。語彙の制約や単語単位の言語モデルは、後処理として適用する必要がある。これは、式(2)の $p(X|S)$ を $p(S|X)$ で代用し、 $p(W)$ を組み合わせる過程に相当する。

認識精度の点において、End-to-End モデルは従来の DNN-HMM ハイブリッドシステムを上回る結果はあまり得られていないが、単語辞書や言語モデルをデコーディングの早い段階で適用していないことが大きな原因と考えられる。

これに対して、1～2年前から単語を単位とした

End-to-End (Acoustic-to-Word) モデルの研究も行われている[13]。これは、語彙や言語モデルをすべて包含して、一気に $p(W|X)$ を求めるもので、ニューラルネットワークのみで音声認識処理が完結し、発音辞書や複雑なプログラムである音声認識デコーダを一切必要としない。3.1 節で述べた従来の標準的な枠組み・アーキテクチャを一新するものである。これを図4に示す。

ただし、音素や文字に比べて、単語のエントリ数は圧倒的に多く、また出現頻度の偏りも大きいため、学習が容易でない。実際はかなり大規模な音声データベースがないと構築が困難である。一方で、音声と書き起こしのパラレルデータからしか学習できず、テキストのみのデータを活用できないという問題もある。さらに、サブワード単位の認識と比較して、音声データベース中に出現しない未知語を認識できない(単語辞書に後で追加することもできない)という実用的に重要な問題もある。

これらの問題への解決策を以下に述べる。

6.1. 文字モデルとのハイブリッド・マルチタスク学習

上記の問題に対して著者らは、文字単位のモデルと併用することで解決を図っている[14]。単語単位のモデルについては、注意機構モデルの方が言語モデルを直接的に包含するので望ましい。この注意機構モデルの正則化の効果期待して、文字単位のモデルは CTC を採用し、エンコーダ部分を共有する構成とした。その上で、両者のマルチタスク学習を行う。また、単語単位のモデルで未知語を検出した際には、文字単位のモデルの認識結果を用いて復元する。

『日本語話し言葉コーパス』(CSJ)を用いて、約2万語彙のシステムを実装・評価したところ、学会講演・模擬講演ともに、DNN-HMM ハイブリッドモデルを認識率で上回ることができた。また、認識に要する時間の実時間比は約0.03でDNN-HMMに比べて30分の1以下となった。発話終了後でないで処理を開始できないが、この処理速度であれば問題ないと考えられる。

6.2. 言語モデルフュージョン

テキストのみのデータから学習される言語モデルについても、デコーディング時に統合することは比較的容易である。これは、shallow fusion と呼ばれる方法[15]で、Acoustic-to-Word モデルで得られる確率 $p(W|X)$ に、外部言語モデルによる確率 $p(W)$ を対数次元で重み付きで足し合わせるものである。注意機構モデルのデコーダにおいてビームサーチを導入する必要があるが、ビーム幅はかなり小さい値(10以下)でも十分である。さらに、End-to-End モデル学習時にも外部言語モデルを統合する方法(cold fusion)[16]も検討されている。た

文 献

だし、これらの方法では Acoustic-to-Word モデルにおける未知語は登録・認識されない。

6.3. エンコーダ転移による適応学習

Acoustic-to-Word モデルの学習には、少なくとも数百時間の学習データが必要となるが、あらゆるタスクドメインでこれだけのデータを収集するのは容易でない。そこで、CSJ のような大規模コーパスで学習したモデルを転移学習によって適応することが考えられる。

その際に音響モデルに対応するエンコーダは固定して、言語モデルに対応するデコーダのみを追加学習することで、適応データがそれほど大きくなくても効果的に適応することができる[17]。この方法では、適応先の語彙も追加されるが、音声データ(と書き起こし)が必要である。

6.4. 音声合成によるテキストからの学習データ生成

上記の2つの方法における問題点を解決し、学習データを効率的に増強するために、音声合成を用いる方法を検討している。最新の音声合成は End-to-End モデルに基づいており、構築が比較的容易な上に、品質も高い。音声認識の学習データには音響特徴量で十分で、音声波形を生成する必要はない。

単一話者の音声合成でも、上記のエンコーダ転移学習と併用することにより、テキストデータのみから語彙の追加を含めたドメイン適応が可能になり、しかも言語モデルフュージョンに比べて効果的である[18]。

ただし音声認識では、多様な話者のパターンをモデル化する必要があるため、音声合成において話者情報を含めることで多数話者の音声を出力できるように拡張している。これにより、テキストデータのみから音声認識学習用のパラレルデータをいくらかでも生成することが原理的に可能になる。

7. おわりに

本稿では、音声認識の方法論を概観し、End-to-End モデル、特に Acoustic-to-Word モデルに至るまでの変遷を述べた。つい数年前まで、図1に示す階層モデルが絶対的なものと考えられていた。しかし、人間の言語獲得とも比べながらあらためて考えると、音響モデルと言語モデルを独立にすることも、サブワード単位のモデル化も不自然といえる。

しかしながら、Acoustic-to-Word モデルにおいても、常にパラレルデータを必要とする点、発話の終端を検出しないと認識を行えない点、新たな語彙を追加できない点など、不自然な点がある。これらの点を解決することが次の段階につながると期待される。

- [1] 河原達也 編著. [音声認識システム \(改訂2版\)](#). オーム社, 2016.
- [2] 河原達也. 音声認識技術の現状と将来展望. 電気学会誌, Vol.133, No.6, pp.364--367, 2013.
- [3] 河原達也. 音声認識技術. 電子情報通信学会誌, Vol.98, No.8, pp.710--717, 2015.
- [4] 篠田浩一. 音声言語処理における深層学習: 総説. 日本音響学会誌, Vol.73, No.1, p. 25-30, 2017.
- [5] 神田直之. 音声認識における深層学習に基づく音響モデル. 日本音響学会誌, Vol.73, No.1, p. 31-38, 2017.
- [6] 増村亮. 深層学習に基づく言語モデルと音声言語理解. 日本音響学会誌, Vol.73, No.1, p. 39-46, 2017.
- [7] T.Sakai and S.Doshita. The Phonetic Typewriter. Proc. IFIP Congress 62, pp.445-450, 1962.
- [8] S.Furui. Selected Topics from 40 Years of Research on Speech and Speaker Recognition. Proc. INTERSPEECH, pp.1-8, 2009.
- [9] G.Hinton, L.Deng, Y.Dong, G.E.Dahl, A.Mohamed, N.Jaitly, A.Senior, V.Vanhoucke, P.Nguyen, T.N.Sainath and B.Kingsbury. Deep neural networks for acoustic modeling in speech recognition. IEEE Signal Processing Magazine, Vol.29, No.6, pp. 82-97, 2012.
- [10] D.Yu and L.Deng. Automatic Speech Recognition – A Deep Learning Approach. Springer, 2015.
- [11] A.Graves and N.Jaitly. Towards End-to-End speech recognition with recurrent neural networks. Proc. ICML, 2014.
- [12] J.Chorowski, D.Bahdanau, D.Serdyuk, K.Cho, and Y.Bengio. Attention-based models for speech recognition. Proc. NIPS, 2015.
- [13] H.Soltau, H.Liao, and H.Sak. Neural speech recognizer: acoustic-to-word LSTM model for large vocabulary speech recognition. Proc. INTERSPEECH, pp.3707-3711, 2017.
- [14] S.Ueno, H.Inaguma, M.Mimura, and T.Kawahara. Acoustic-to-word attention-based model complemented with character-level CTC-based model. In Proc. ICASSP, pp.5804-5808, 2018.
- [15] A.Kannan, Y.Wu, P.Nguyen, T.Sainath, Z.Chen and R.Prabhavalkar. An analysis of incorporating an external language model into a sequence-to-sequence model. Proc. ICASSP, pp.5824-5827, 2018.
- [16] A.Sriram, H.Jun, S.Satheesh and A.Coates. Cold fusion: training seq2seq models together with language models, Proc. INTERSPEECH, pp.387-391, 2018.
- [17] S.Ueno, T.Moriya, M.Mimura, S.Sakai, Y.Yamaguchi, Y.Aono, and T.Kawahara. Encoder transfer for attention-based acoustic-to-word speech recognition. In Proc. INTERSPEECH, pp.2424--2428, 2018.
- [18] M.Mimura, S.Ueno, H.Inaguma, S.Sakai, and T.Kawahara. Leveraging sequence-to-sequence speech synthesis for enhancing acoustic-to-word speech recognition. In Proc. IEEE SLT, 2018.