

聴覚障害者のための講演・講義の音声 認識による字幕付与*

解説

河原達也, 秋田祐哉 (京都大学)**

43.72.Ne

1. はじめに

2015年3月に情報処理学会の全国大会を京都大学で開催する機会があり、新たな試みをいくつか行った。その一つが託児室の開設で、もう一つが聴覚障害者への情報保障である。いずれも関連学会の全国大会ではほとんど前例がなかった。聴覚障害者の情報保障というと、手話通訳を用意すればよいと思っている人が意外に多い。しかし、手話(視覚言語)でコミュニケーションを行うのは生まれつき聴覚に障害のあるろう者が中心で、はるかに多数の難聴者や中途失聴者は音声言語を中心にコミュニケーションを行い、手話を解さない人も多い。したがって情報保障には、音声言語を要約筆記・字幕化することが必要になる。しかし、要約筆記ボランティアの大半の方は「本職」があり、平日に複数人を手配するのが非常に困難であることがわかった。京都のような規模の都市でこのような状況であることに愕然とした。

2016年度から施行されている障害者差別解消法では、障害者の社会的障壁の除去について「必要かつ合理的な配慮」を行うことが義務づけられている(民間の場合は努力義務)。「その実施に伴う負担が過重でない」範囲でという留保がついているが、学会や大学で対応するにははなはだ心もとない状況である。

上記全国大会では、メインのイベント企画の動画配信も初めて実施した。最近、MOOC(Massive Open Online Courses)をはじめとして、教育関

係でも多くの動画コンテンツが作成・配信されるようになってきている。テレビ番組では国の行政指針もあり、生放送も含めてかなりの割合で字幕が付与されるようになったが、ネット配信のコンテンツでは字幕はほとんど皆無といってよい。米国ではオバマ前大統領による「[21世紀の通信と映像アクセシビリティ法\(CVAA: The 21st Century Communications and Video Accessibility Act\)](#)」の制定があったが、我が国ではそういう機運もない。動画の数は増える一方なのに、字幕付与のための人的資源や予算が決定的に不足している。

このように情報保障のニーズと現状には大きなギャップがある。著者らは、音声認識技術がこの状況を改善できるものと考えて研究開発を行っている[1][2]。

2. 音声認識を用いた字幕付与における課題

2.1 字幕付与の手段

音声をリアルタイムにテキスト化し、字幕付与を行う手段を表1にまとめる。ここでは、速記者などのプロが行う特殊な場面と、一般の講演会などでボランティアが行う場面を分けて記載している。このうち手書きの場合は、多数の人への画面表示に向いていない。また、通常書く速度は話す速度より大幅に遅いので、「2割要約」と揶揄されている。原稿テキスト送出(前ロール)は、事前に原稿を用意して読み上げる場合のみ可能で、限られた場面ではしか使えない。したがって現在、字幕付与に主に用いられているのは、タイプ入力である。一般のパソコン要約筆記の場合も、2名の連係入力で交代しながら行うので3~4名は必要となる。要約筆記ボランティアの養成と確保が課題となっている[3]。

* Captioning lectures using automatic speech recognition for hearing-impaired people

** Tatsuya Kawahara and Yuya Akita (Kyoto University, Kyoto 606-8501)

e-mail: kawahara@i.kyoto-u.ac.jp

表 1 音声のリアルタイムテキスト化・字幕付与の手段

	特殊な場面 (プロの職業) テレビ・議会・法廷など	一般の場面 (主にボランティア) 講演会や講義など	課題
手書き	速記 ...リアルタイム反訳不可	手書きノートテイク	会場に表示できない 情報量が少ない
事前原稿利用	原稿テキスト送付	原稿テキスト送付 (前ロール)	原稿を読む場合のみ可
タイプ入力	ソクタイプ	パソコン要約筆記 (PCテイク)	要員の確保
音声入力	復唱入力 (リスピーク)		要員の養成
直接音声認識	専用の音声認識	音声認識のカスタマイズ	精度の確保

2.2 音声認識を用いた情報保障の試み

これに対して、音声認識を用いる試みも様々に行われている。テレビ番組の字幕付与や国会の会議録作成にも音声認識システムが導入されているが、これらは基本的に話し慣れた話者による発話であり、しかも膨大な学習データを用いて専用の音声認識システムを構築したもので、90%を上回る正解率を実現している。

一般の音声認識システムを講演や講義のような話し言葉に適用しても、十分な認識精度を確保するのは困難である。そこで、2000年代にいくつかの大学で講義の情報保障のために行われた試みは、復唱入力によるものであった[4]。しかし、復唱者の養成・確保は、パソコン要約筆記と同様に容易でない。

これに対して著者らは、2007年頃から講師の発話音声直接音声認識する方式の研究開発を行っている。講演の音声認識に関しては、我が国の『日本語話し言葉コーパス』(CSJ)プロジェクト[5]を初めとして、2000年頃から活発に研究が行われているが、聴覚障害者の字幕付与に使用された事例は(最近まで)ほとんどない。

米国では、このようなプラットフォームの開発と普及を行う Liberated Learning Consortium なども設立されたが、2013年に開催されたアクセシビリティの情報技術に関するトップカンファレンス ACM ASSETS で行われた Captioning Challenge においても、Google の音声認識を用いた方式は、プロのソクタイプや一般の PC 連携入力に及ばず、“readable or usable”な字幕を生

成できなかったと報告されている[6]。Google は YouTube でも自動音声認識による字幕付与のサービスを提供しているが、認識精度はまだ十分とはいえない[7]。

著者らは 2007 年から毎年『聴覚障害者のための字幕付与技術』シンポジウムを開催しているが、2009年のシンポジウムの第一著者(河原)の講演において、初めて音声認識を用いた字幕付与を行った。これは、一般の講演会で直接音声認識による情報保障を行った国内最初の事例と考えており、その後 4 節で述べるように様々な学会の研究会などに展開している。一般の講義については 2009年に実験を行ったが[8][9]、十分な情報保障ができるには至っていない。

近年、深層学習の導入によって、音声認識精度は格段に向上している。これに伴って、音声認識を用いた様々なシステム/サービスも登場している。例えば、聴覚障害者の日常生活におけるコミュニケーション支援のために、NICT ではスマートフォンアプリの「こえとら」、タブレットアプリの [Speech Canvas](#) などが開発されている。同様にコミュニケーション支援・会話のテキスト化アプリとして、[UD トーク](#)があり、こちらは教育機関への導入も進んでいる。ただし授業やゼミなどで使用する際には、話者は一定の質で音声入力するためにはかなり意識して発話する必要があるようである[10]。最近では、富士通の [LiveTalk](#) や東芝の [RECAIUS 音声ビューア](#) などのように会議や講演会での利用を掲げた製品も出ているが、本格的な利用例はまだ確認されていない。

2.3 講義・講演の字幕付与における課題

講義や講演の音声認識が容易でないのは、機械（音声認識装置）を意識しない自然な話し言葉であるためである。読上げ調の音声と比較して、音響的・言語的に様々な違いがあり、別のモデル化が必要と考えられる[11]。また、認識結果がフィードバックされないため、どう話せば認識される／されないがわからないという点もある。

ただし、講演・講義・会議などにも様々な発話スタイルがあり、その分類を表2に示す。原稿を読む場合は、通常の音声認識システムでもほとんど問題なく対応可能であるが、情報保障の観点からは原稿テキスト表出(前ロール)で十分である。MOOCにおいても通常は原稿を用意するので、字幕に転用可能である。予稿はあるが自然な発話である学会発表やTEDトークなどは、音声認識研究で扱われてきたタスクであり、最近では約90%の認識精度を達成している。

一方、大学の講義やミーティングなども対象とした音声認識も行われているが、かなり低い認識精度(おおむね60-80%)しか報告されていない。これは、発話スタイルの問題に加えて、收音環境の問題もあると考えられる。一般の講演や講義では、部屋の音響条件や收音機器さらに拡声機器が多様であり、必ずしも音声認識システムが想定しているクリーンな音声が入力されない。音響的な知識のあまりないユーザによる利用も想定する必要がある。

次に、講演・講義の音声認識の形態・応用について図1に示す。大きく分けて、その場でリアルタイムに字幕表示する場合と録音・録画したコンテンツをオフライン処理する場合がある。

表2 講演・講義・会議などの発話スタイル

	話者一名	多人数
原稿を読む	スピーチ MOOC	議会の本会議
予稿あり、入念に準備	学会発表 TED など	議会の委員会 組織の会議
公共の場で考えながら発話	一般の講演・講義	ミーティング ゼミ
プライバシーのある空間	インタビュー	会話

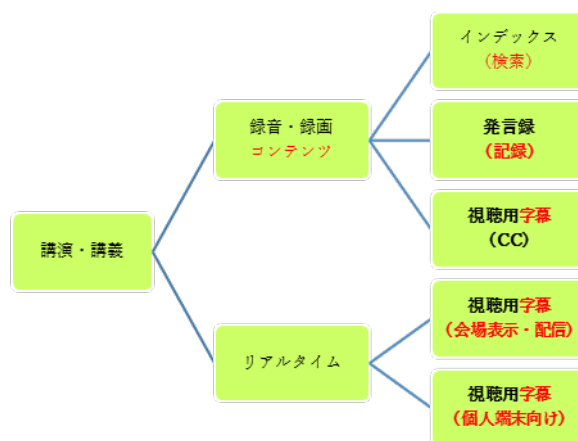


図1 講演・講義の音声認識の形態・応用

後者のうち、検索インデックスとしての利用であればキーワードが認識できるレベルでよく、一方、発言記録として利用するのであれば人手で修正・編集し、読みやすい文章にする必要がある。本稿で扱う情報保障用の字幕は、これら以外の場合に該当するが、視聴の形態によって求められる品質が異なってくる。録音・録画コンテンツの配信字幕(クローズドキャプション:CC)の場合は、人手で修正し、音声に忠実にする必要がある。リアルタイム字幕でも、会場のディスプレイに表示したり、ネット中継の配信に用いる場合は、人手で修正し、誤りをできるだけ出さないようにする必要がある。一方、字幕を個人端末のみに出す場合は、利用は個人が取捨選択するので、ベストエフォートでよいとも考えられる。

しかし、音声認識精度が80%を下回ると内容を把握するのは困難になるし、すべての誤りを訂正し字幕を効率よく生成するには90%以上が望ましい。ただし、90%の音声認識精度を実現するには、話し言葉の精緻なモデル化に加えて、一定水準以上の発話スタイルと音響条件が必要になる。

また、講演や講義の字幕において問題になるのは、専門用語である。大学レベル以上の専門用語は通常の音声認識のモデルでカバーされていないことが多い。専門用語の頻度や音声認識精度に与える影響はそれほど大きくないが、専門用語が正しく認識・表示されないと、情報保障の品質に重大な影響がある。専門用語をユーザが登録できる枠組みも必要であるが、予稿やスライドなどから自動的に獲得できることが望ましい。

3. 講演・講義コンテンツへの字幕付与

本節では、講演・講義を録画したコンテンツに対して字幕付与を行った事例について述べる。このような字幕では、誤りが無い完璧なものが要求され、タイプ入力することも容易であるので、それよりも効率的な優位性を実現する必要がある。

3.1 システムの構成

著者らは、講演や国会審議を対象とした音声認識の研究を進めてきたが[2]、そのような音声コンテンツに対して音声認識と字幕付与を行うサーバ (URL は以下) を構築している[12]。

<http://caption.ist.i.kyoto-u.ac.jp/>

利用者は、音声ファイルや映像ファイルを当該サーバにアップロードし、所定の手続きをすると、音声認識による書き起こしにタイムスタンプが付与されたファイル (SAMI や SRT など複数のフォーマット) が生成される。これらは、一般的な再生ソフトで字幕ファイルとして利用可能である。

音声認識には誤りが含まれる上に、話し言葉には言い淀みなども多いため、字幕として提示するには編集が必要である。また適当な位置での改行や句読点挿入も必要である。そのためのエディタ (図 2 参照) も上記サイトで提供している。

現在、想定しているコンテンツは以下の 3 種類であり、各々について音声認識のモデルが用意されている。

- 講演：学会や講義など大教室で 1 人で行う学術講演 (CSJ の学会講演データで構築)
- スピーチ：一般的な話題に関してゆっくり話すもの (CSJ の模擬講演データで構築)
- 討論：議会審議など公共の場で複数人で行う討論 (国会審議のデータで構築)

音声認識システムはすべて、Julius と DNN-HMM 音響モデル (系列識別学習) 及び単語 3-gram 言語モデルで構築している (学習データベースは上記参照)。

また、コンテンツに関連するテキスト (予稿やスライドなど) を同時にアップロードすることで、音声認識の言語モデルを適応することも可能である。

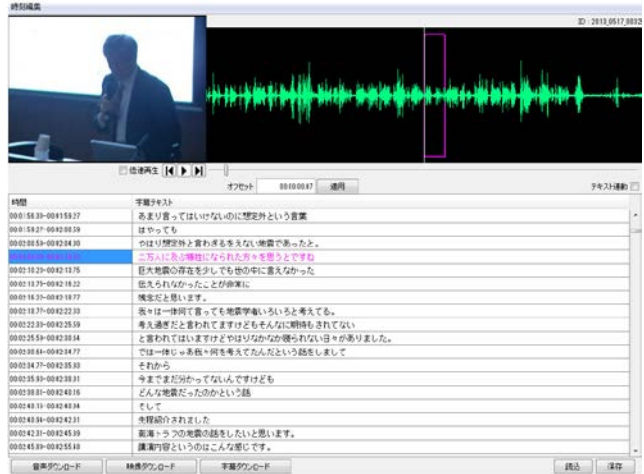


図 2 字幕編集用エディタ

3.2 京都大学 OCW 配信講義への字幕付与

京都大学 OCW (OpenCourseWare) では数千件の講演動画が配信されているが、このうち、CiRA (iPS 細胞研究所) の一般の方対象シンポジウムと「大震災後を考える」シンポジウムシリーズの講演に対して取り組んだ[13]。

関連するテキストを用いて音声認識の単語辞書や言語モデルを適応し、また音響モデルも話者に適応することで、65~88%程度の認識率が得られた。認識率のばらつきが大きい、最も影響するのは音響条件 (マイクとの距離や雑音・残響) であり、収録時に留意してもらい必要がある。この音声認識結果を編集して作成した字幕を付与・配信している (図 3 参照)。

2012年ノーベル生理学・医学賞 受賞
山中 伸弥 教授による講演「iPS細胞研究の進展と課題」

(2010年CiRA一般の方対象シンポジウム「iPS細胞研究の最前線」より)



図 3 京都大学 OCW の講演の字幕付与

3.3 放送大学の講義への字幕付与

放送大学で配信されている講義では、受講生の要望に応じて一部に字幕が付与されているが、テレビ番組の半数程度にとどまっている。そこで、インターネットで配信されている講義に対して、音声認識を用いた字幕付与に取り組んだ[14]。

講義はスタジオで収録されるので、音響条件もよい。すべてに台本が用意されているわけではないが、一般の講演や講義に比べると、はるかに発声は明瞭である。通常、科目毎に1冊のテキストが作成されている。

音声認識には、3.1節で紹介した「講演モデル」を用いるが、科目毎に教科書テキストを追加混合することで適応を行う。これにより専門用語の追加も自動的に行われるが、正しい読みが付与されているか確認が必要である。

以下のラジオ講義を対象に、音声認識を行ったものを編集することで字幕テキストを作成する実験を行った。

- (1) 心理臨床の基礎
- (2) リスク社会のライフデザイン
- (3) CGと画像処理の基礎

講義はいずれも各回45分で15回あるが、(2)については台本のある3回分を除いた。また、(3)については本実験では7回分のみを用いるが、これらには台本がある。

平均の音声認識率を表1に示す。これは、最終的に生成された字幕テキストに対する文字単位の正解率である。この編集に要した時間と実時間比も表3に示す。この編集作業は、1名の作業者が図2のエディタを用いて行ったものである。この後、実時間と同程度の確認作業を行っている。また講義(1)(2)について、各回の認識率と編集時間の関係をプロットしたものを図4・5に示す。

これらから音声認識率と編集時間の間に高い(0.5~0.6)相関があることがわかる。以前放送大学で行った実験では、音声認識を用いずに放送大学の講義を書き起こした場合、実時間の平均5.3倍(約4時間)程度と報告されている[15]。これをグラフ中の直線と重ねると、87%程度の認識率の場合に相当する。これから、音声認識率が87%以上の場合にその効果があることが示唆される。また、93%になると1/3以上の時間短縮効果が示され、かなり優位性があるといえる。

表3 講義の音声認識率と編集時間

	講義数	認識率	編集時間	実時間比
(1)	15	90.8%	3時間 16分	4.4
(2)	12	88.5%	3時間 46分	5.0
(3)	7	94.4%	2時間 52分	3.8

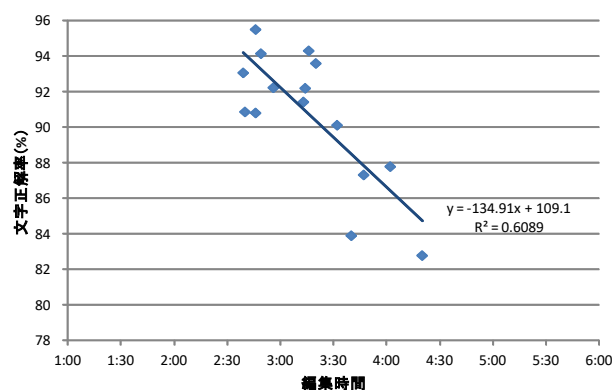


図4 音声認識率と編集時間の相関 (「心理臨床の基礎」)

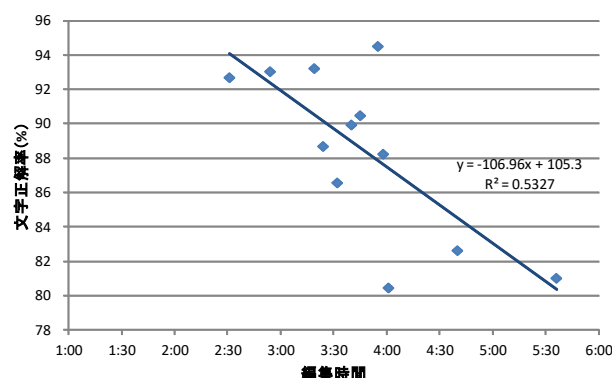


図5 音声認識率と編集時間の相関 (「リスク社会のライフデザイン」)

前記の3つのラジオ講義のインターネット配信しているコンテンツについては、2016年度から字幕付与されている。その際に、講義で用いている図やグラフを静止面の形で貼り付けることで、通常のラジオ講義と比べて理解がしやすいようにしている(図6参照)。

また、2016年度から開設されたオンライン授業の多くの番組で、音声認識を用いた枠組みで字幕が付与されている。

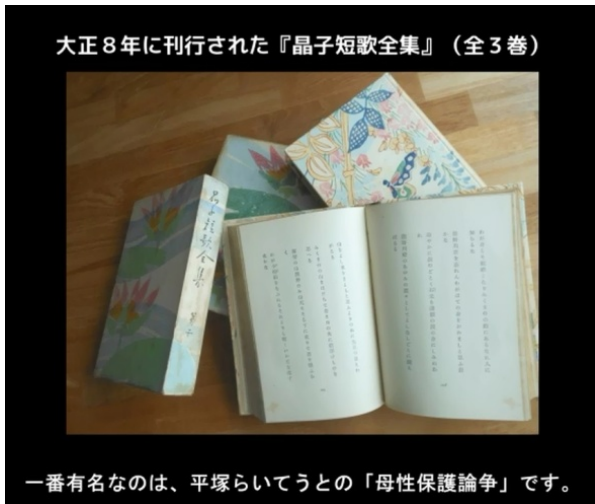


図6 放送大学の講義への字幕付与の例
 (「特別講義メディアと与謝野晶子」©放送大学)

4. リアルタイム字幕付与

本節では、講演の会場でリアルタイムに字幕付与を行った事例について述べる。このような字幕では、一定の正確性ととも時間に遅れも重要になる。情報保障の品質の観点からは、認識誤りの修正や不要な箇所の削除などを行う編集が必要となる。しかし、パソコン連係入力に対する効率性・優位性を目指して、1名で編集する枠組みを採用している (図7 参照)。

4.1 システムの構成

音声認識には、3.1 節で紹介した「講演モデル」を使用し、言語モデル及び単語辞書は講演予稿やスライドのテキストを用いて話題への適応を行う。これらのテキストが利用可能となるのは一般的に講演開催の直前で、適応や調整の作業に時間的な余裕がないため、適応には単純なテキストの線形補間手法を用いている。

講師の音声は、会場の拡声機器 (PA) から分配して編集端末 (PC) に入力することを想定しているが、これが困難な場合は独自にマイクを設置して入力する。入力された音声は、会場の端末では音声認識を行わず、音声区間が切り出されるたびにそのデータを 3.1 節の字幕サーバにインターネット (LAN または 4G モバイル回線) 経由で送信する。サーバ側で音声認識を行い、その結果を音声区間ごとに端末で逐次的に取得し、字幕としての編集と出力を行う。

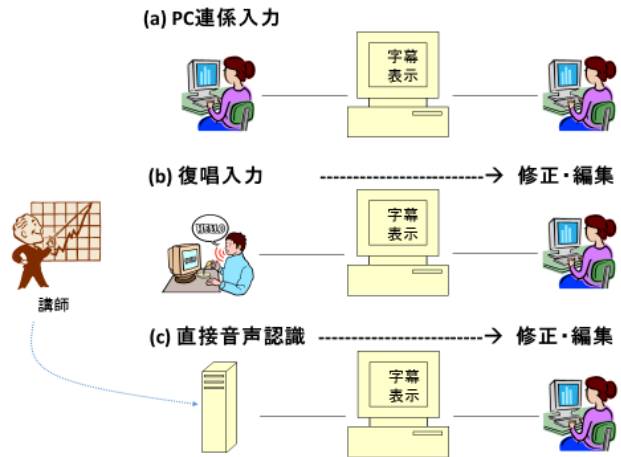


図7 PC 連係入力と音声認識による字幕付与

字幕の提示には、パソコン要約筆記で一般的に用いられているソフトウェア IPtalk を使用する。音声認識結果は IPtalk の「確認修正パレット」に取り込み、ここで修正・編集作業を行うこととした。これは、PC 連係入力に慣れている人の利便性を考慮したものである。音声認識システム (Julius) と IPtalk を接続するソフトは以下の Web サイトで公開している。

<http://sap.ist.i.kyoto-u.ac.jp/jimaku/julius2iptalk.html>

本システムでは、認識結果からフィルターを自動的に削除し、また典型的な文末表現に限って句点を自動挿入したうえで IPtalk に入力している。

なお、この音声認識システムは最近の PC (2.5GHz 程度のマルチコア CPU; GPU は不要) で、スタンドアロンで実行可能である。そのため「話し言葉音声認識キット」及び「講演音声認識キット」を Julius の Web サイトで公開している。

<http://julius.osdn.jp/index.php?q=dictation-kit.html>

また、IPtalk にも、これらのキットを使用して音声認識を行い、前記の接続ソフトを用いて、字幕を作成できる機能が搭載された。この機能の詳細については、IPtalk の Web サイトを参照されたい。

http://www.geocities.jp/shigeaki_kurita/

4.2 情報処理学会 SIG-AAC での運用

実際の学術講演に対する、前節の音声認識システム(サーバ)を用いたリアルタイム字幕付与の試みを、先述の『[聴覚障害者のための字幕付与技術](#)』シンポジウムや情報処理学会アクセシビリティ研究会(SIG-AAC)の一部の講演で実施している。作業員(1名)により、音声認識結果をできる限り編集したものを字幕として会場の聴衆に提示する。情報処理学会の研究会では講演の映像をインターネットで配信しており、SIG-AACの講演では、作成した字幕は映像と合わせて配信も行われる。

一例として、2016年2月のAAC研究会において作成した字幕の概要を表4に示す。3件の講演について、音声認識を用いてリアルタイム字幕を提供した。合計で2時間8分の講演時間に対して1名の作業員で編集作業を行い、32,601文字の字幕を送出した。編集した文字数は5,116文字で、このうち半数は音声認識結果の削除であった。講演時間1分あたりでは40文字の編集である。ただし、IPtalkでは不要な認識文は一括削除でき、編集の半数が削除であるので、実際の作業負荷はかなり小さい。なお、完全に正確な書き起こしを作成していないため、実際の音声認識の精度は明らかではないが、字幕の92%は音声認識結果をそのまま用いている(残りの8%は表4の置換・挿入にあたる)。

表4 音声認識を用いて作成した字幕の概要
(情報処理学会 SIG-AAC・2016年2月)

作業員	1名
字幕送出文字数	32,601
音声認識文字数	33,697
編集文字数	5,116 (16%) 内訳: 置換 1,190 (3.7%) 挿入 1,415 (4.3%) 削除 2,511 (7.7%) (字幕送出文字数に対する割合)

5. おわりに

実際に字幕が付与されたコンテンツを視聴すると、字幕がないもの比べて理解が深まる印象がある。これは、聴覚(音声)と視覚(文字)の相乗効果に加えて、日本語では表意文字の漢字を用いている効果もある。例えば、「カンサイボウ」と聞くだけより、「幹細胞」と見る方が概念をイメージしやすい。一方で、毎年『[聴覚障害者のための字幕付与技術](#)』シンポジウムで、聴覚障害者の方にも意見を伺うと、「表示のタイミングが早くて情報量も多くて良い」という意見もあれば、「五月雨式に文字がたくさん表示されても理解がつかない」という意見も多い。字幕は単に発話を文字にすればよいというものでなく、わかりやすいようにレイアウトする必要がある。しかしこれは、その場で冗長さを含めてやりとりする音声言語とそもそも遠隔に伝えるための文字言語の本質的な違いに起因する問題であるとも考えられる。将来的にはこのような点も考慮して字幕付与を高度化する検討もしていきたい。

謝辞

本稿で記載した内容の多くは、『[聴覚障害者のための字幕付与技術](#)』シンポジウムを通じて得たものである。参加・講演頂いた聴覚障害者、要約筆記者、教育関係者、速記者、ICT研究者などの多数の方々に感謝します。

放送大学の字幕付与に関しては、同大学の広瀬洋子教授との共同研究によるものである。

SIG-AACにおける字幕付与に関しては、前主査の平賀瑠美教授(筑波技術大学)をはじめとする関係者の協力を頂いた。IPtalkについては開発者の栗田茂明氏(日本遠隔コミュニケーション支援協会)に、Julius認識キットについては李晃伸教授(名古屋工業大学)に多大な協力を頂いた。深く感謝します。

文献

- [1] 河原達也. 聴覚障害学生支援の最先端 ―音声認識による字幕付与技術. 嶺重慎, 広瀬浩二郎(編), 知のバリアフリー, 第4章, pp. 109–122. 京都大学学術出版会, 2014.
- [2] 河原達也. 話し言葉の音声認識の進展 ―議会の会議録作成から講演・講義の字幕付与へ―. メディア教育研究, Vol. 9, No. 1, pp. 1–8, 2012.

- [3] 吉川あゆみ, 太田晴康, 白澤麻弓: 大学ノートテイク入門, 人間社, 2001.
- [4] 中野聡子, 牧原功, 金澤貴之, 中野泰志, 新井哲也, 黒木速人, 井野秀一, 伊福部達. 音声認識技術を用いた聴覚障害者向け字幕提示システムの課題 —話し言葉の性質が字幕の読みに与える影響— 電子情報通信学会論文誌. Vol. J90-D2, No. 3, pp.808-814, 2007.
- [5] 前川喜久雄. 『日本語話し言葉コーパス』の概観. 国立国語研究所, 2004.
- [6] R. S. Kushalnagar. ASSETS 2013 Captioning Challenge, ACM SIGACCESS Newsletter, Issue 108, January 2014.
- [7] H. Liao, E. McDermott, and A. Senior. Large scale deep neural network acoustic modeling with semi-supervised training data for YouTube video transcription. Proc. IEEE-ASRU, 2013.
- [8] 勝丸徳浩, 河原達也, 秋田祐哉, 森信介, 山田篤. 講義音声認識に基づくノートテイクシステム. 電子情報通信学会技術研究報告, SP2009-53, WIT2009-59, 2009.
- [9] T. Kawahara, N. Katsumaru, Y. Akita, and S. Mori. Classroom note-taking system for hearing impaired students using automatic speech recognition adapted to lectures. In Proc. INTERSPEECH, pp. 626--629, 2010.
- [10] 松崎丈. 音声認識アプリを活用した支援システムの構築に関する検討. 宮城教育大学情報処理センター研究紀要. No. 24, pp. 3-8, 2017.
- [11] S. Furui and T. Kawahara. Transcription and distillation of spontaneous speech. In J. Benesty, M. M. Sondhi, and Y. Huang, editors, Springer Handbook on Speech Processing and Speech Communication, chapter 32, pp. 627--651. Springer, 2008.
- [12] 秋田祐哉, 三村正人, 河原達也. 音声認識を用いた講義・講演の字幕作成・編集システム. 情報処理学会研究報告, SLP-108-2, 2015.
- [13] 秋田祐哉, 河原達也. オープンコースウェアの講演を対象とした音声認識に基づく字幕付与. 日本音響学会研究発表会講演論文集, 2-9-9, 春季 2013.
- [14] 河原達也, 秋田祐哉, 広瀬洋子. 自動音声認識を用いた放送大学のオンライン授業に対する字幕付与. 情報処理学会研究報告, AAC-2-5, 2016.
- [15] 長妻令子, 福田健太郎, 柳沼良知, 広瀬洋子. クラウドソーシングを活用した効率良い字幕作成手法. 電子情報通信学会技術報告, WIT2012-25, 2012.
- [16] S. Li, Y. Akita, and T. Kawahara. Semi-supervised acoustic model training by discriminative data selection from multiple ASR systems' hypotheses. IEEE/ACM Trans. Audio, Speech & Language Process., Vol. 24, No. 9, pp. 1524--1534, 2016.