

### 3.5.1 ユニバーサルコミュニケーション研究所 音声コミュニケーション研究室

室長 堀 智織 ほか 17 名

#### 音声言語コミュニケーションシステムのための音声認識、音声合成、対話制御技術の研究

##### 【概要】

本研究室では、人間にとって自然で簡便な情報伝達手段である音声によるコミュニケーションを用いた音声対話・音声翻訳システムを実現するため、音声認識、音声合成、対話処理の研究開発を行っている。さらに、インターネット上の音声を含むマルチメディアデータに対する情報検索を実現するため、高速な音声インデキシング技術および多言語字幕付与の研究開発を行っている。今年度、音声認識を行うために必要な多言語の学習データを効率的に収集し、英語・中国語音声のニュース音声を対象とした高精度モデルの構築を行った。さらに、英語講演音声認識を対象とした競争型国際ワークショップでは2年連続で認識性能が首位となった。また、多言語音声翻訳技術の研究開発を目的とした23カ国28研究機関(平成26年3月末)から構成される国際研究共同体U-STARを主導し、ネットワーク型多言語音声翻訳の実証実験を行った。その結果、実証実験で取得された音声データを用いて、タイ語の音声認識性能を単語正解精度が30%から60%に改善した。

##### 【平成25年度の成果】

##### ● Web上の動画から音声特徴データを収集

Web上には膨大な量の音声付き動画データがあり、それらの大量の音声データから音響モデルを学習することにより認識性能を改善することが可能である(図1)。今年度は、Web上の音声5,000時間の収集目標のうち、中国語約800時間、英語約6,000時間の音声データから音響モデルを学習するための音響特徴量を抽出した。その特徴量を用いて中国語ニュースの音声認識の単語正解率を77.2%(平成24年度59.4%)、英語ニュースでは82.9%(平成24年度63.8%)に向上させた。今後は、同様の収集システムを用いて、日本語、英語、中国語など多言語音声データを収集し、さらに音声認識性能の改善を行う。

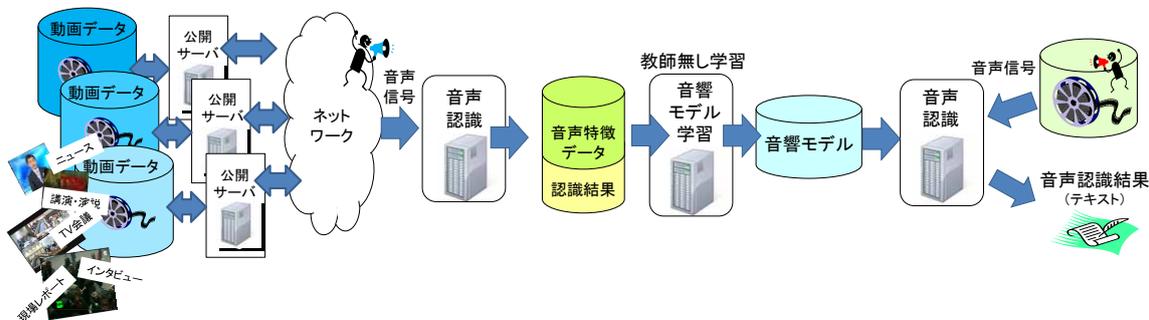


図1 Web上の大量音声データを用いた音声認識システムの構成

##### ●国際連携 U-STAR による多言語音声認識の研究加速

NICTはアジア・ヨーロッパの音声・言語の研究機関6(23カ国28機関)から成る国際研究共同体U-STAR(<http://www.ustar-consortium.com/>)を主導し、2010年にNICTが国際標準化(ITU-T勧告書F.745およびH.625に準拠)したネットワーク型音声翻訳通信プロトコルを用いて各加盟機関の音声翻訳サーバを相互接続し、ネットワーク型多言語音声翻訳システムを開発した。平成24年7月から継続的に公開している音声翻訳アプリVoiceTra4U(図2)では、17言語の音声認識と14言語の音声合成を実現した。本実証実験を通して収集された実利用データを用いて、タイ語の単語正解率を約60%(平成24年度約30%)に大幅に改善した。

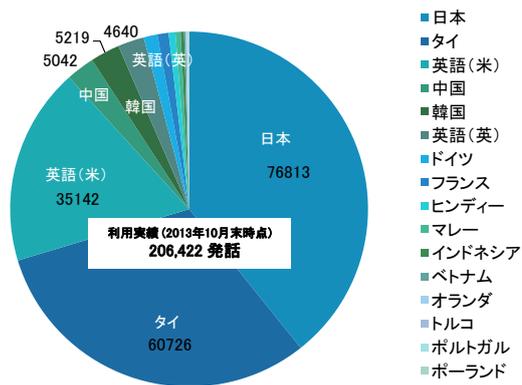


図2 U-STAR 実証実験ログデータ

●評価型国際ワークショップで2年連続首位獲得

実世界の大規模な語彙を実時間で高精度に認識する新手法として、重み付き有限状態トランスデューサ(WFST)に基づく大語彙連続音声認識システムを研究開発し、高速かつ高精度な認識を実現している。提案手法に基づく音声認識システムは、英語講演音声 TED(図3)に対し、話し終わると同時に結果を出力する実時間音声認識の条件で、単語正解精度80%という高精度な書き起こしを生成することができた。さらに、より長い認識時間をかけることにより、90%の単語正解精度を達成することができた。本システムを用いて英語講演音声認識を対象とした競争型国際ワークショップ IWSLT に参加し、音声認識性能で世界第1位を2年連続で獲得した。日本語の音声認識だけでなく、英語、中国語の音声認識で高性能であったことから、Web上にある多言語の音声データに対するリアルタイムインデキシングの研究に本システムを適用することが有効と考えられる。



http://www.ted.com/

参加組織	2012	参加組織	2013
NICT	12.1	NICT	13.5
KIT-NAIST	12.4	KIT	14.4
KIT	12.7	MITLL-AFRL	15.9
MITLL	13.3	RWTH	16.0
RWTH	13.6	NAIST	16.2
UEDIN	14.4	UEDIN	22.1
FBK	16.8	FBK	23.2

図3 英語講演音声 TED の音声認識

●因子分解 RNN 言語モデルによる性能改善

近年、言語モデルの性能改善に寄与する手法として、再帰的な接続を持つニューラルネットワークにより、文全体における単語間の依存性を推定するリカレントニューラルネットワーク(RNN)が提案された。本研究室では、単語の表層に留まらず、品詞、語幹などを考量した因子分解 RNN (fRNN) 言語モデルを提案し、IWSLT における英語講演音声認識で約1%の性能改善を果たし、首位獲得に貢献した。

- IWSLT (<http://htc.cs.ust.hk/iwslt/>)
- KIT :カールスルーエ工科大学(ドイツ)  
(Facebookに採用されたCMUのエンジンと同等)
- MITLL/AFRL:マサチューセッツ工科大学リンカーン研究所  
/空軍研究所(アメリカ)
- RWTH :アーヘン工科大学(ドイツ)
- NAIST :奈良先端科学技術大学院大学(日本)
- UEDIN :エディンバラ大学(イギリス)
- FBK :ブルーノ・クessler財団 研究所(イタリア)

●多言語音声コミュニケーション技術の事業化

今年度研究開発された上記研究成果に基づく音声認識システムを株式会社フィートに商用リリースし、au「おはなしアシスタント」の多言語音声翻訳技術に採用され、NICTの音声認識技術が音声コミュニケーションシステムの普及に大きく貢献した。

●学術的な成果

学術論文誌7本、トップレベルの国際学会(採択率20%以下)12本、ほか国際会議に15本の研究成果を発表し、学会において活発な研究発表を行うだけでなく、U-STARにおける主導的な役割と競争型国際ワークショップにおいて首位を獲得した技術力により、NICTの世界的なプレゼンスを高めた。

特記事項

ASTAP(アジア・太平洋電気通信標準化機関)の音声・自然言語処理専門家グループを主導し、アジア・環太平洋地域における多言語音声コミュニケーション技術の研究開発を推進した。

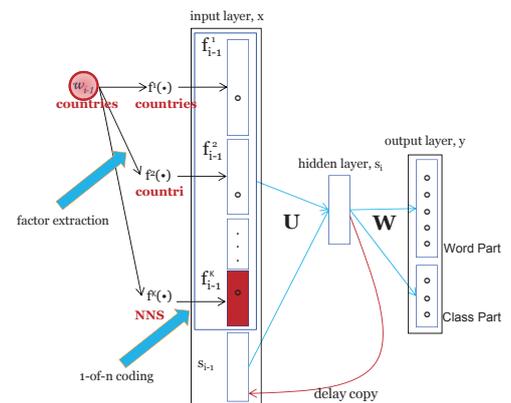


図4 因子分解リカレントニューラルネットワーク言語モデル