

ロボット対話における深層学習を用いたセミブラインド音声強調

和気 雅弥[†] 坂東 宜昭[‡] 三村 正人[‡] 糸山 克寿[‡] 吉井 和佳[‡] 河原 達也[‡]
[†] 京都大学 工学部情報学科 [‡] 京都大学 大学院情報学研究科 知能情報学専攻

1. はじめに

自然なロボット対話の実現には、人間同士の自然な対話で起こりうる状況に対応することが重要である。その一つに、一方が発話している最中に他方からの発話が割り込むバリエーションという状況がある。ロボット対話では、ロボットが発話している最中に人間からの発話が割り込むという状況である。このとき、マイクロホンの観測信号は両発話の混合音であるのでそのままの音声認識は困難であり、人間の発話を強調する必要がある。

本稿で扱う問題設定はセミブラインド音源分離と呼ばれ、除去の対象となるロボットの音源信号を利用できる。従来手法には、音声の混合モデルを基にして確率的に音声強調を行う手法 [1] や、独立成分分析を用いる Semi-blind ICA (SB-ICA) [2] といった手法が挙げられる。

本稿では、図 1 のような残響のある環境下でのセミブラインド音源分離を考える。観測された混合音からロボットの音源信号を利用して人間の発話を強調するために畳み込みニューラルネットワーク (CNN) を用いたセミブラインド音声強調 (SB-CNN) を提案する。CNN はネットワーク内で局所的な特徴の抽出を学習させる層を用い、その層を重ねることで複雑な特徴抽出を可能とし、特に画像認識の分野で大きな成果を挙げている [3]。音声においても、スペクトルを二次元の画像とみなすことが可能であり、音声の特徴を抽出することで音声の空間モデルを用いずに音源位置などに頑健な音声強調を行う。

2. 畳み込みニューラルネットワーク

本稿では、以下の音声強調の問題を扱う。

- 入力 1. マイクロホンの観測スペクトル X_{tf}
 2. ロボットの音源スペクトル S_{tf}^R

出力 人間の音源スペクトル S_{tf}^U

上記の t, f はそれぞれ時間フレームと周波数ピンを表す。このときロボット音源、人間音源とマイクロホン間の伝達関数をそれぞれ H_R, H_U とすると、観測スペクトルと各音源スペクトルとの関係は以下の式で表せる。

$$X_{tf} = H_R S_{tf}^R + H_U S_{tf}^U$$

また、本稿で用いるネットワークは複素数を扱わないので、各スペクトルの振幅をネットワークの入出力に用い、時間領域への変換に際してマイクロホン入力の位相を用い、逆短時間フーリエ変換の後に窓関数の影響が及ばないように重みをかけて調整している。

図 2 に提案法で用いる CNN の概形を示す。入力は、マイクロホンの観測とロボットの音源のスペクトルを前後 5 フレーム、各 11 フレームずつとする。前後のフレームも入力に加えることで、時間軸方向での音声の特徴が抽出され、残響除去の効果が期待できる。また、CNN の最初の畳み込み層のフィルタサイズを周波数軸方向では音声の基本周波数より大きくとることで、音声の周波数軸方向に持つ調波構造が考慮されることが期待される。

Semi-blind speech enhancement using deep learning in human-robot interaction: Masaya Wake, Yoshiaki Bando, Masato Mimura, Katsutoshi Itoyama, Kazuyoshi Yoshii, Tatsuya Kawahara (Kyoto Univ.)

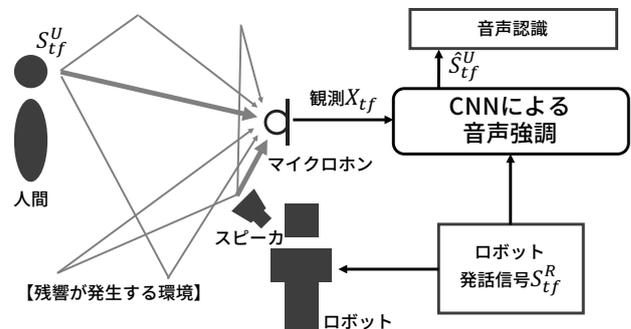


図 1: 状況設定及びデータの流れ

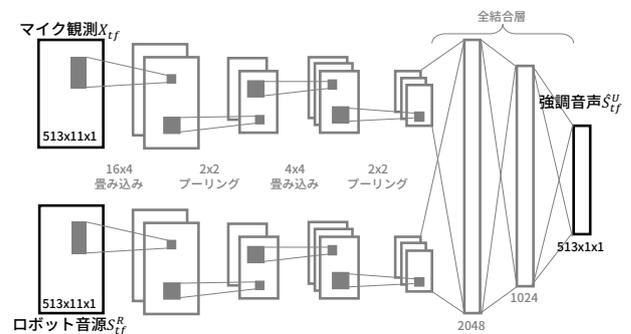


図 2: 用いる CNN の概形

提案法では最初の畳み込み層のフィルタサイズは、時間軸方向は 40 ミリ秒、周波数方向は 250Hz とするよう設定した。

プーリング層や全結合層での活性化関数には、[3] で高速な学習が可能と報告されている Rectified Linear Units (ReLU) を用いる。出力層は振幅スペクトルの各周波数ピンごとの値とし、非負にするため出力層にも ReLU を用いる。

3. 評価実験

CNN を用いた提案法の有効性を示すため、数値実験により評価を行った。

3.1 実験条件

実験のためのデータセットには、新聞記事読み上げコーパス (ASJ-JNAS) [4] の音素バランス 503 文を用い、人間の音声として男性話者、ロボットの音声として女性話者の発話データを利用した。学習データには各 60 話者 3012 発話、評価データには学習データと異なる各 4 話者 200 発話を用い、サンプリングレートは 16kHz、混合時の信号対雑音比 (SNR) は 0dB に統一した。

図 3 に示す部屋の中央にマイクロホンを置き、音源は伝達関数を測定した 40 箇所のうち、学習データにおいては白丸の 36 箇所、評価データにおいては黒丸の 4 箇所の中からランダムに配置する。またマイクロホン、音源の高さはそれぞれ 115cm、110cm である。伝達関数の有効長は 16384 サンプル (1024 ミリ秒) とした。スペクトルの窓幅は 1024 サンプル (64 ミリ秒)、窓のシフト幅は 160 サンプル (10 ミリ秒) とし、最初の畳み込み層の

