



表1 従来手法と提案手法の認識実験結果

モデル	CER
(1) TM-GMM-HMM	52.4 %
(2) TM-DNN-HMM	49.3 %
(3) LSTM-GMM-HMM	39.7 %
(4) TM-LSTM-DNN-HMM	45.4 %
(5) LSTM-DNN-HMM-KD	39.5 %

力分布を用いる。この学習方法は知識蒸留と呼ばれ、強い正則化効果があると報告されており[5]、過学習防止に効果が期待できる。

### 3 認識実験

#### 3.1 実験条件

評価データとして咽喉マイクで収録した男性話者10名による自由発話の音声(約20分)を使用し、文字誤り率(CER)による評価を行った。教師モデルの学習に用いる接話マイク音声には日本語話し言葉コーパス(CSJ)より約240時間を使用した。咽喉マイクの音響モデルの学習に用いる接話マイクと咽喉マイクのパラレルデータには簡易防音室で同期録音した男性話者8名による約3時間の音素バランス文読み上げ音声を使用した。言語モデルにはCSJの書き起こし文から生成した3-gram言語モデルを用いた。

教師モデルの入力特徴量はfMLLRを適用した40次元のFBANKの前後5フレームを結合した440次元の特徴量である。教師モデルはStacked Denoising Autoencoderによる事前学習を行った後、ファインチューニングを行った。咽喉マイクの音響モデルの入力特徴量はfMLLRを適用した40次元のFBANKである。特徴マッピングのLSTMのルックバック数は7とした。DNNのノード数は図1に示す通りである。

#### 3.2 実験結果

まず、従来手法と提案手法の実験結果を表1に示す。ここで表中のモデル(1)は約3時間の咽喉マイク音声のMFCCで学習したGMM-HMMを音響モデルとしたシステム、(2)は約3時間の咽喉マイク音声のFBANKとそのアライメントで学習したハイブリッド方式の音響モデルを有するシステム、(3)は先行研究[3]と同様に接話マイクのBNFで学習したGMM-HMMを音響モデルとし、咽喉マイク音声のFBANKから接話マイク音声のBNFへLSTMによってマッピングした特徴量を入力とするシステムである。(4)は提案手法(図1)と同じDNNの構造を有するが、アライメントから生成した正解ラベルを用いて学習されている。(5)は提案手法のシステムである。(2)のDNNのノード数は[440: 1024: 1024: 1024: 42: 1024: 出力層]とした。提案手法(5)は従来の(1)と(2)のシステムと比較してCERの大幅な削減を達成でき、先

表2 提案手法の初期化方法に関する認識実験結果

前段初期化	後段初期化	CER
Random	Random	41.0 %
FM-LSTM	Random	41.3 %
Random	CM-Hybrid	39.3 %
FM-LSTM	CM-Hybrid	39.5 %

行研究[3]の手法(3)と比べて性能が向上した。また、(5)は同じDNNの構造を有する(4)のシステムよりも精度が高く、知識蒸留によって学習がうまく進んだと考えられる。

次に、提案手法の初期化方法に関する実験を行った。実験結果を表2に示す。ここで表中のRandomはDNNのパラメータをランダムに初期化すること、FM-LSTMは特徴マッピング用LSTMのパラメータで初期化すること(図1-③)、CM-Hybridは接話マイク音声で学習したハイブリッド方式の音響モデルのDNNの後段のパラメータで初期化すること(図1-④)を表す。結果からDNNの後段を接話マイクのハイブリッド方式の音響モデルのDNNのパラメータで初期化する方法はランダムに初期化するよりも精度が高く、有効である可能性が示された。しかしながら、今回の実験ではDNNの前段を特徴マッピング用LSTMで初期化する方法の有効性は認められなかった。特徴マッピング用LSTMの学習データが少なく、パラメータの学習が不十分で有効な初期値と成り得なかった可能性が考えられる。

### 4 おわりに

本研究ではDNNの前段を咽喉マイクのFBANKから接話マイクのBNFへのマッピングを学習したLSTM、後段を接話マイクのBNFから音素を識別するよう学習したDNNのパラメータで初期化した後、パラレルデータを用いて知識蒸留に基づきネットワーク全体を最適化する手法を提案し、従来手法と比較して咽喉マイク音声の認識精度を改善することができた。しかしながら、今回の実験ではDNNの前段を特徴マッピング用のLSTMのパラメータで初期化する方法の有効性は認められなかった。今後の課題としてモデルのハイパーパラメータの調整や学習データの追加などが挙げられる。

#### 謝辞

本研究の一部は科研費(16H01817, 16K01543)の助成を受けた。

#### 参考文献

- [1] Dupont, S. et al., *In Proc. of Robust 2004*, 2004
- [2] Shengke, L. et al., *In Proc of NCSP*, pp. 363-366, 2018.
- [3] 鈴木他, SLP, Vol.2018-SLP-123, pp.1-6, 2018
- [4] 神田, 日本音響学会誌, 第73巻, 第1号, pp.31-38, 2017
- [5] Hinton, G. et al., arXiv:1503.02531, 2015