

ゼロ資源言語の認識に向けたMAMLに基づく end-to-end 音声認識の検討

周 鋭^{1,a)} 伊藤 彰則^{1,b)} 能勢 隆^{1,c)}

概要: 高い精度の音声認識システムを開発するためには、大量のラベル付け音声データが必要である。しかし、世界中にある言語の大部分については、そのように多くの学習データを用意することはできない。このような言語の音声認識システムを開発するため、学習時に目標言語のデータを使用しない、または非常に少ないデータのみを使って開発する音声認識システムがゼロリソース音声認識システムである。本稿では、メタラーニングの学習方法である MAML とハイリソース言語データを用いて、モデルを事前学習する。そして、少数の目標言語のデータを使って、モデルを微調整する。20 分程度の目標言語データを使用して、MAML による学習を行ったところ 40 %程度の CER が得られた。

キーワード: ゼロリソース, メタラーニング, 音声認識

MAML-based End-to-End Speech Recognition for Zero-Resource Language Recognition

ZHOU RUI^{1,a)} ITO AKINORI^{1,b)} NOSE TAKASHI^{1,c)}

Abstract: High-precision speech recognition systems require large amounts of labeled speech data. However, most of the languages in the world do not have that much training data. A speech recognition system that uses no speech data, or uses very little data from such training target languages, is called a zero resource speech recognition system. In this paper, we pre-train a model using a training method called MAML and high-resource language data. Then we use a data of a minority language to fine-tune the model. When we use 20 minutes of the target language data, we obtained a character error rate of about 40%.

Keywords: Zero Resource, Meta Learning, Speech Recognition

1. はじめに

自動音声認識 (ASR) は、入力された音声信号から、音声とテキストを出力する技術である。音声認識技術は、私たちの日常生活の中で幅広く活用されている。例えば、スマートフォンで調べたいことを口に出して検索したり、チャットの際に文字を打ちたくない場合は直接音声で入力したり、自動車に音声で指示を出したりすることができる。音声認

識モデルには、2つのタイプがある。それは、End-To-End モデルと非 End-To-End モデルである。2010 年以前には、主に非 End-To-End モデルである GMM-HMM[1] が使われていた。しかし、ニューラルネットワークの発展により、GMM から DNN, すなわち DNN-HMM[2] への置き換えが進んだ。非 End-To-End モデルは、HMM と DNN, CNN, RNN のハイブリッドモデルであり、音響モデル、組合せ発音モデル、言語モデルの 3つのモデルを組み合わせたものである。これに対して、End-To-End モデルは、すべての音声認識ステップを1つのモデルで行うため、よりシンプルで、パラメータも少ない。

このような音声認識モデルを学習させる通常の方法で

¹ 東北大学工学研究科
Tohoku University Graduate School of Engineering

a) zhou.rui.p1@dc.tohoku.ac.jp

b) akinori.ito.a2@tohoku.ac.jp

c) nose@tohoku.ac.jp

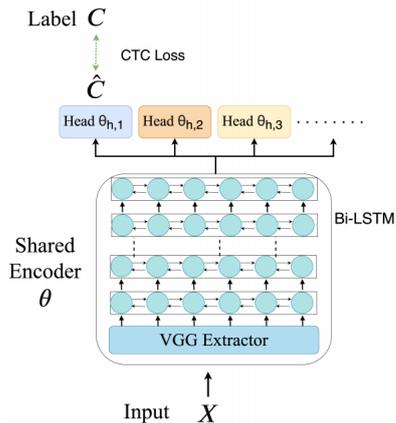


図 2: モデルの構造

Fig. 2 Structure of model

表 1: モデルのパラメータ
Table 1 Parameter of model

モデル	パラメータ
LSTM Layer	6
LSTM Cell	360
LSTM Dropout	0.1
VGG Dropout	0.1

3.2 CTC Loss

音声認識における音声は、一般的にテキストよりも長くなる。そのため、テキストに対応する音声の長さを決めることができないが、CTCはこの長さの問題を解決する。音声信号の長さを L 、テキストの文字数を U とし、テキストの全長が L となるように、 $L-U$ の数の空白トークンをランダムにテキストトークンに挿入する。このような挿入を行う場合のすべてのアライメント確率の和を使ってロスを計算する。 $\text{Loss} = -\log(P(Y|X))$ 、ここで X は認識される音声で、 Y は認識したテキスト結果で、 P は音声認識結果の確率である。

4. 実装と実験

MAML の効果を証明するために、三つの実験を行った。三つの実験で用いる事前学習モデルは同じなので、まず事前学習のデータを示す。

表 2 は事前学習データの情報である。言語データは common voice コーパス [12] のものを使用した。これは九種類の言語で、キャラクターはアルファベットである。言語ごとに 2000 のタスクセットを作成したので、全部で 18000 のタスクセットを使用した。

4.1 大量の英語データを用いた実験

最初の実験は、異なる事前学習エポックに対応するモデルで、5000 個の英語ファイルを学習した結果である。英語ファイルは、Librispeech コーパス [13] の 360 時間の

表 2: 事前学習データ情報
Table 2 Pre-training data information

コーパス	時間 (h)	音声ファイル数 (個)	キャラクター
Chuvash	1.38	1211	86
Frisian	5.05	3707	77
Greek	2.12	1976	93
HakhaChin	0.64	815	58
Kyrgyz	2.28	1788	76
Maltese	2.44	1977	81
Sakha	2.34	1444	73
Slovenian	1.33	1348	65
Latvian	2.77	3198	80

表 3: 英語データの情報
Table 3 English data information

ファイル数	5000	3000	1000	500	100
時間 (h)	17.3	10.3	3.4	1.74	0.35

学習データから選択された。事前学習エポックは 3、微調整エポックは 10 である。評価指標は Train-Loss, Valid-Loss, Character Error Rate, Word Error Rate とした。図 3 は、5000 個の英語ファイルの微調整の結果を示している。epoch0 は事前学習モデルを用いない通常の学習方法、epoch1~epoch3 は事前学習エポックの数である。その結果、事前学習モデルは、通常の学習方法よりも早く低損失を実現し、文字誤り率、単語誤り率も低くなっていることがわかる。しかし、エポック数を増やしても、効果はあまり上がらない。

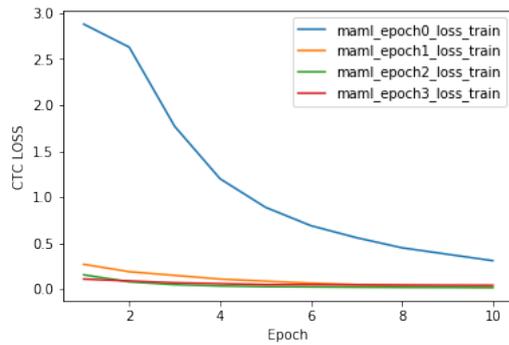
4.2 少量の英語データによる実験

前記の実験は、英語の微調整用ファイルを 5000 個使用した実験である。これに対し、リソースが少ない場合の効果を実証するため、微調整用の英語ファイルの数を減らして、実験を実施した。このとき使用した事前学習モデルは、8 エポックモデルである。英語ファイルの数は表 3 に示す通りである。英語ファイル数は 5000, 3000, 1000, 500, 100 であり、微調整エポック数は 10 である。評価指標は Train-Loss, Valid-Loss, Character Error Rate, Word Error Rate である。

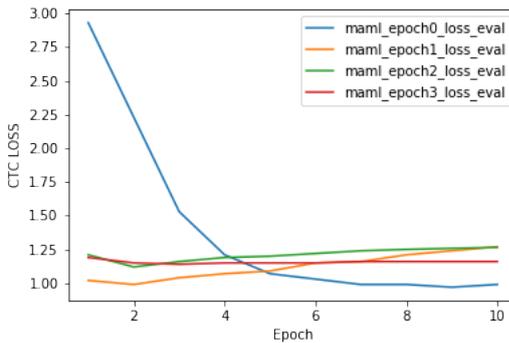
結果を図 4 に示す。結果から見ると、5000 ファイルと 3000 ファイルの場合の最終的な CER と WER はほぼ同じである。ファイルの数が減少するにつれて、CER と WER は増加する。しかし、一番少ない 100 個のファイルでも、CER が 28%, WER が 70% であり、通常の学習方法よりも良い結果となった。

4.3 少数言語の実験

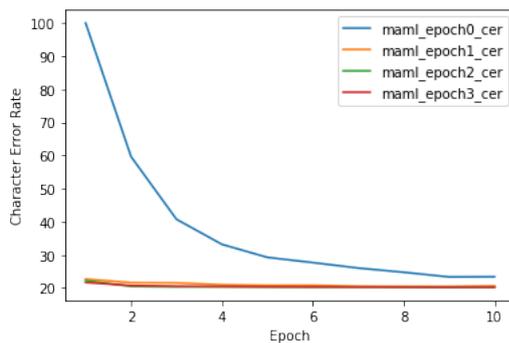
前節の英語実験の結果から、MAML を使うことにより、微調整データが非常に少ない場合でも、ある程度の認識性



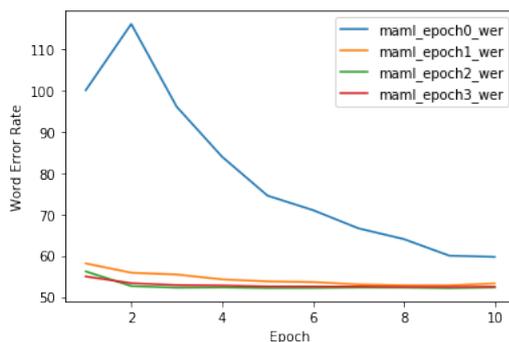
(a) Train-Loss



(b) Valid-Loss



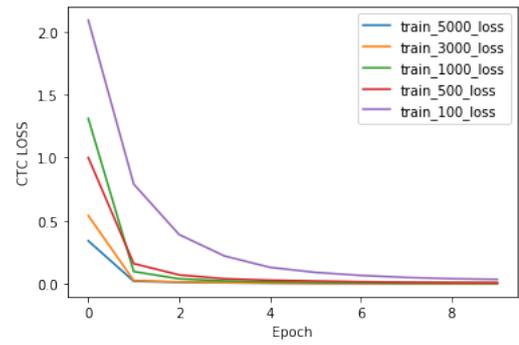
(c) Character Error Rate



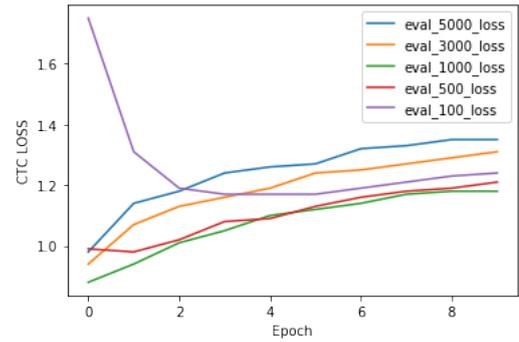
(d) Word Error Rate

図 3: 5000 英語微調整ファイル, 異なる epoch の結果
Fig. 3 5000 English fine-tuning files, different epoch results

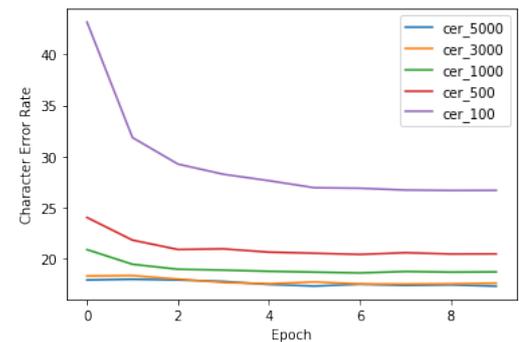
能が得られた. この結果を受けて, 少数言語の認識実験を行った. 表 4 は微調整用の少数言語の情報である. 実験に使用したデータの長さは 30 分程度であり, 通常の学習方法では全く認識をすることができない. 前節の英語実験と



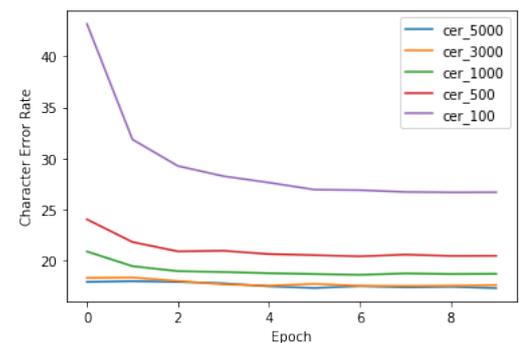
(a) Train-Loss



(b) Valid-Loss



(c) Character Error Rate

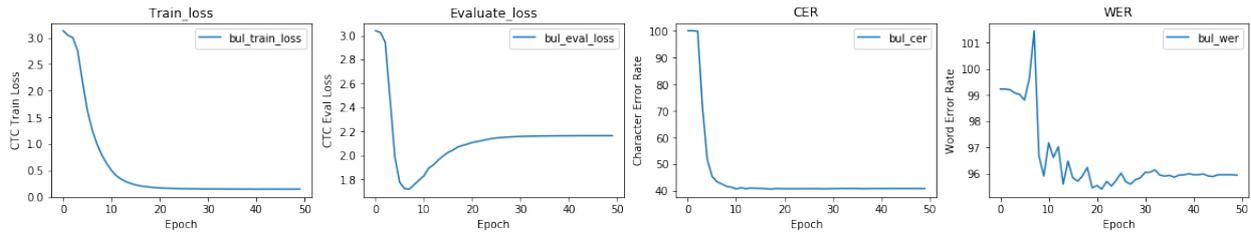


(d) Word Error Rate

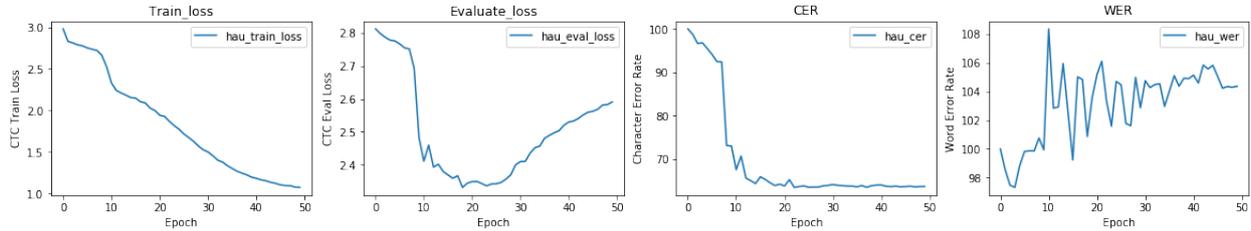
図 4: 異なる微調整英語ファイル数の結果
Fig. 4 Results of different number of fine tune English files

同じように 8 epoch の事前学習モデルを使用して, この少数言語で微調整を行う. 微調整の epoch 数は 50 とし, 評価指標は前節と同じとした.

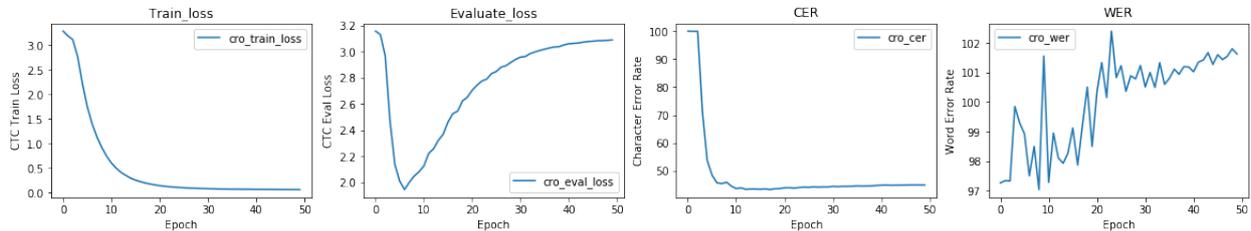
図 5 はそれぞれの言語の結果である. 30 分程度の学習



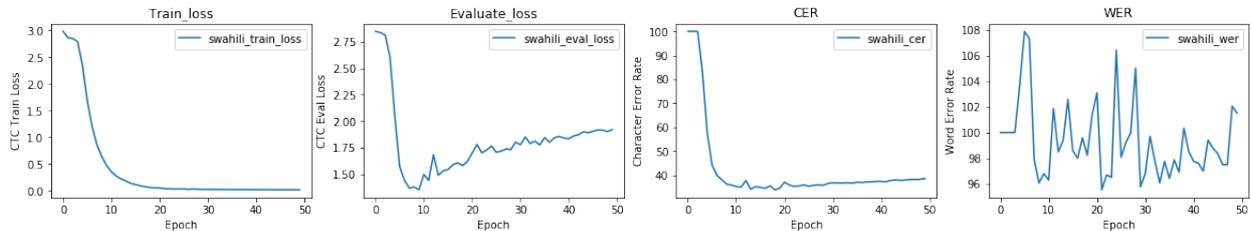
(a) Bulgarian 語の結果



(b) Hausa 語の結果



(c) Croatian 語の結果



(d) Swahili 語の結果

図 5: 少数言語実験の結果

Fig. 5 Results of minority language experiments

表 4: 微調整用少数言語のデータ情報

Table 4 Minority language data information for fine tuning

コーパス	時間 (h)	音声ファイル数 (個)	文字種
Bulgarian	0.44	200	66
Croatian	0.35	150	77
Hausa	0.24	200	55
Swahili	0.25	150	39

表 5: ゼロリソース言語の実験結果

Table 5 Experiment results of zero resource language

言語	CER(%)	WER(%)	AVG.LOSS
English	40.66	84.7	2.93
Bulgarian	93.16	99.9	5.63
Croatian	90.77	99.9	5.53
Swahili	90.2	100	4.94
Hausa	93.6	100	6.28

データにより, 50%以下の CER が得られた. 通常の学習方法では全く認識できないので, MAML の効果があったことが分かる.

4.4 ゼロリソースの実験

前の実験では目標言語のデータを少量使用して微調整を行った. 最後の実験では事前学習モデルだけを使用して,

微調整なしに目標言語の音声を認識する. 言語は四種類の少数言語と英語を使用した. 結果を表 5 に示す.

ゼロリソースの実験結果から, 微調整をしない場合, 英語の CER だけが 40%であったものの, 他の言語はあまり認識されていない. また, LOSS も大きくなっている.

5. おわりに

本論文では、まず、ゼロリソースの音声認識手法であるメタ学習を紹介した。その後、メタ学習を実装し、実験を行った。結論は3つある。第一に、メタ学習は、事前に学習したモデルを用いることで、高リソース音声認識の効果を向上させることができる。第二に、事前学習済みモデルを用い、非常に少ないデータで微調整を行うことで、ある程度の認識が可能となる。わずか30分程度の音声データで、50%以下の文字誤り率を得ることができた。最後に、対象言語での微調整を行わず、事前学習済みモデルのみを用いた場合、少数言語の認識性能はよくなかった。今後は、さらに性能を向上させるため、言語間の文字表記の統一、より少ない微調整データによる認識方法の検討などを行っていく。

参考文献

- [1] Jaitly, N., Nguyen, P., Senior, A., & Vanhoucke, V. (2012). Application of pretrained deep neural networks to large vocabulary speech recognition.
- [2] Li, J., Yu, D., Huang, J. T., & Gong, Y. (2012, December). Improving wideband speech recognition using mixed-bandwidth training data in CD-DNN-HMM. In 2012 IEEE Spoken Language Technology Workshop (SLT) (pp. 131-136). IEEE.
- [3] Vilalta, R., & Drissi, Y. (2002). A perspective view and survey of meta-learning. *Artificial intelligence review*, 18(2), 77-95.
- [4] Finn, C., Abbeel, P., & Levine, S. (2017, July). Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning* (pp. 1126-1135). PMLR.
- [5] Chan, W., Jaitly, N., Le, Q., & Vinyals, O. (2016, March). Listen, attend and spell: A neural network for large vocabulary conversational speech recognition. In 2016 IEEE international conference on acoustics, speech and signal processing (ICASSP) (pp. 4960-4964). IEEE.
- [6] Paterlini-Brechot, P., & Benali, N. L. (2007). Circulating tumor cells (CTC) detection: clinical impact and future directions. *Cancer letters*, 253(2), 180-204.
- [7] Graves, A. (2012). Sequence transduction with recurrent neural networks. arXiv preprint arXiv:1211.3711.
- [8] Zhang, Q., Lu, H., Sak, H., Tripathi, A., McDermott, E., Koo, S., & Kumar, S. (2020, May). Transformer transducer: A streamable speech recognition model with transformer encoders and rnn-t loss. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 7829-7833). IEEE.
- [9] Chiu, C. C., & Raffel, C. (2017). Monotonic chunkwise attention. arXiv preprint arXiv:1712.05382.
- [10] Beckmann, P., Kegler, M., Saltini, H., & Cernak, M. (2019). Speech-vgg: A deep feature extractor for speech processing. arXiv preprint arXiv:1910.09909.
- [11] Graves, A., & Schmidhuber, J. (2005). Framewise phoneme classification with bidirectional LSTM and other neural network architectures. *Neural networks*, 18(5-6), 602-610.
- [12] Ardila, R., Branson, M., Davis, K., Henretty, M., Kohler, M., Meyer, J., ... & Weber, G. (2019). Common voice: A massively-multilingual speech corpus. arXiv preprint arXiv:1912.06670.
- [13] Panayotov, V., Chen, G., Povey, D., & Khudanpur, S. (2015, April). Librispeech: an asr corpus based on public domain audio books. In 2015 IEEE international conference on acoustics, speech and signal processing (ICASSP) (pp. 5206-5210). IEEE.