

Multilingual Multi-Task Learning for Low-Resource Acoustic Modeling

Josh Meyer¹

¹University of Arizona

joshua.richard.meyer@gmail.com

Abstract

Index Terms: speech recognition, multi-task learning, acoustic modeling

1. Introduction

Previous work has shown that performance for a low-resource language on speech recognition can be improved by adding training data from another, resource-rich language. Typically, data from another language is added as a separate task in the Multi-Task Learning framework via an additional output layer (Caruana 1997).

The targets for this additional language have always been states of context-dependent triphones, defined by some tree clustering algorithm. This current research builds off the intuition that triphones encode information which is too fine-grained to be maximally useful for language-transfer. Using a higher-level of linguistic abstraction (eg. the monophone), we are able to better extract the kind of language-general information useful in training an acoustic model for some target language.

To put it another way, if we want to lower Word Error Rates for a language like Urdu, learning to distinguish different versions of English [th] in context is probably not very useful. A better way to make use of English data would be to focus on distinguishing more common linguistic contrasts like [p vs b]. In adding an additional language as an auxiliary task, it would be better to focus on distinctions which are robust and will transfer well to a new, target language.

The tasks are created by redefining the parameters of the HMM-GMM system used to bootstrap a DNN-hybrid system, such that the phonetic decision tree is cut short.

The target language is Kyrgyz, and the source language is English. Both data come from audiobooks, English being from LibriSpeech and Kyrgyz from the Bizdin.kg project.

2. Background Literature

Past attempts at MTL

3. DNN-Hybrid Training as Model Transfer

The standard DNN-Hybrid approach requires the GMM-HMM system to provide the labels for supervised training. This reliance of the DNN on GMM alignments is actually a form of model transfer, where the DNN is trained to perform the exact same classification as its GMM predecessor. The DNN not only learns the frame alignments from the individual GMMs, but also the DNN indirectly learns the structure of the phonetic decision tree used to define the tied-state system. This is because the output layer of the DNN is trained to predict targets which were defined via leaves of the decision tree.

Given that standard triphones encode very fine-grained information which may not help performance on a target lan-

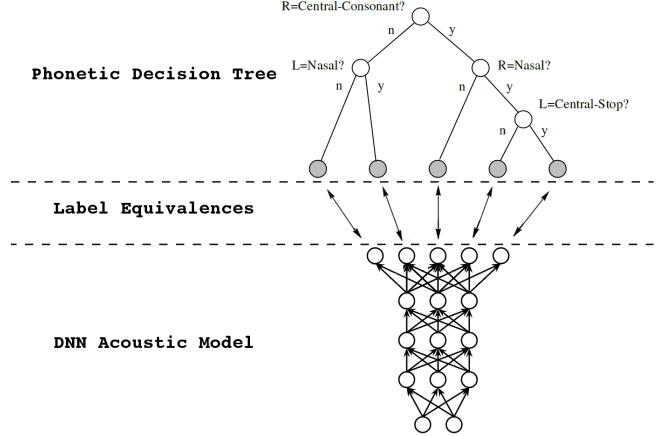


Figure 1: *GMM→DNN Model Transfer*

guage, the following experiments investigate GMM→DNN model transfer at a higher level for an additional source language.

4. Experiments

4.1. Data

Two speech corpora are used in the following experiments:

1. LibriSpeech selection (5 hours)
2. Kyrgyz audiobook section (1.5 hours)

4.2. Model Building

All models were build using the Kaldi `nn3` approach. The scripts used in this paper can be found at (XXX). These scripts are based on the offical repo multilingual Babel scripts here (XXX).

In GMM training, monophones were allotted a total of 1000 Gaussian components and trained over 25 iterations of EM. These monophones were then expanded into context-dependent triphones via a phonetic decision tree, with a maximum of 1000 leaves (ie. tied-states). The resulting leaves (state clusters) are then trained as context-dedendent triphones with a total of 2000 Gaussian components over 25 iterations of EM.

Given the alignments from the GMM-HMM models, a 5-layer, 200-dimensional TDNN is trained over 5 epochs of backprop on a single GPU instance.

Each auxiliary task is implemented as a separate output layer along with a penultimate hidden layer. All other hidden layers are trained in parallel.

4.2.1. English

Using the standard CMUDict phoneset, 3-state monophones are trained from a flat-start via EM training, and then expanded into triphones via Kaldi's decision tree.

4.2.2. Kyrgyz

The Kyrgyz phoneset can be found here (XXX). The transcriptions are basically a simple grapheme set, with predictable variation (allophony) explicitly encoded. For example, the letter K is pronounced as [q] in the context of back vowels, and in the context of front vowels it is pronounced [k]. The phoneme set includes both [q] and [k].

Decoding is performed with a bigram backoff language model was trained on a Wikipedia Kyrgyz dump, and contains, 103,998 unigrams and 56,6871 bigrams.

Because I am interested in the application of this MTL technique to low resource languages, I examine how much the performance boost from the source language changes with regards to the size of the training data for the target language.

Baseline

All the following architectures will be compared to the performance of the following baseline.

To account for any advantage multiple output layers may bring about, the baseline also contains two output layers, where the tasks are identical. In this way, random initializations in the weights and biases for each task are accounted for.

During testing, *only one* of the tasks (ie. the main task) is used. The additional tasks are dropped and the baseline Kyrgyz triphones are used in decoding. This highlights the purpose of the extra tasks: to force the learning of robust representations in the hidden layers during training. The tasks may in fact not be the best option for final classification; they serve as "training wheels" which are then removed once the net is ready.

Auxiliary Tasks

The auxiliary tasks all related to the English language data from the LibriSpeech corpus. Investigating the intuition that labels generated by a standard triphone phonetic decision tree are not the best representation of data for transfer learning, the auxiliary tasks here investigate different levels in the decision tree's branches. By forcing the neural net to recognize higher levels in the tree, we will learn representations which are more abstract, and therefore more likely to be relevant multi-lingually.

4.3. Results

All results are performed on the same held-out section of Kyrgyz audiobook. The bigram language model, lexicon, and main-task decision tree are build into a standard decoding graph in the traditional Kaldi style.

5-layer, 500 dimensial hidden layer, 10 epoch results:

As we can see, the baseline model overfits to the training data after about 250 iterations. The performance for the final model (clearly would be beat out if we used early stopping) is 53.66% WER on decoding the held-out test data.

Above, we see the addition of triphones as an auxiliary task from the LibriSpeech corpus. Overfit to Atai after 500 iterations. The performance for the final model (clearly would be beat out if we used early stopping) is 49.85% WER on decoding the held-out test data.

Finally, monophones from the LibriSpeech corpus. Overfit to Atai after 500 iterations. The performance for the final model (clearly would be beat out if we used early stopping) is 51.32% WER on decoding the held-out test data.

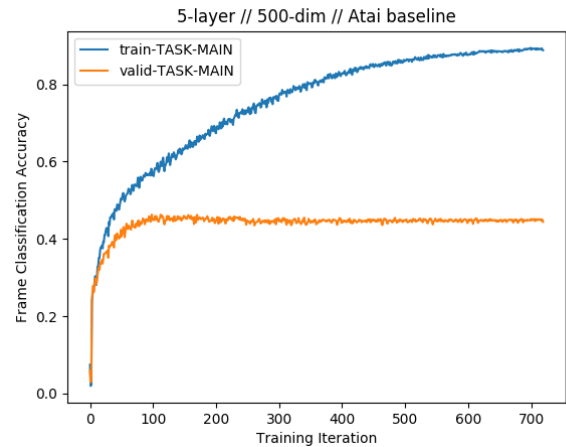


Figure 2: Baseline Model

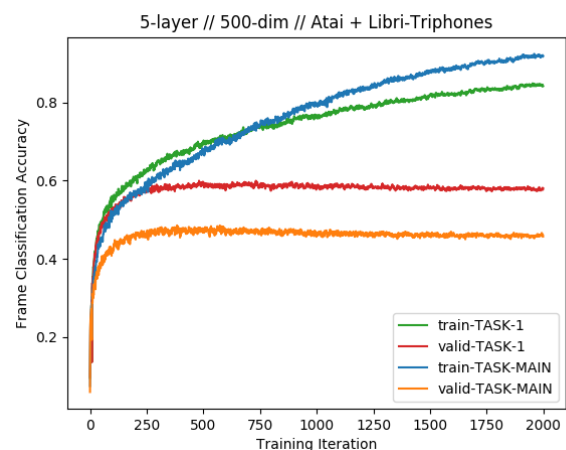


Figure 3: Aux Task == LibriSpeech Triphones

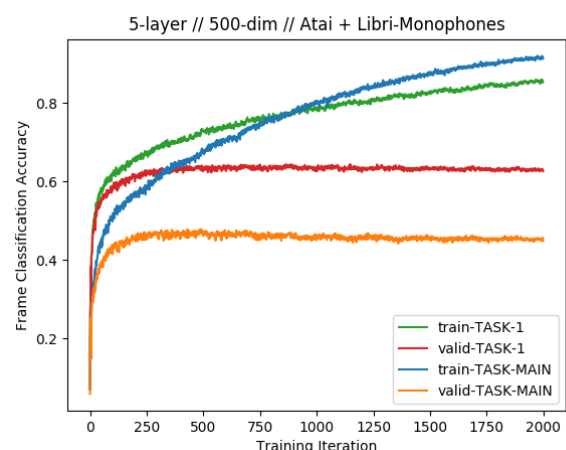


Figure 4: Aux Task == LibriSpeech Monophones

Here we are with the last model:

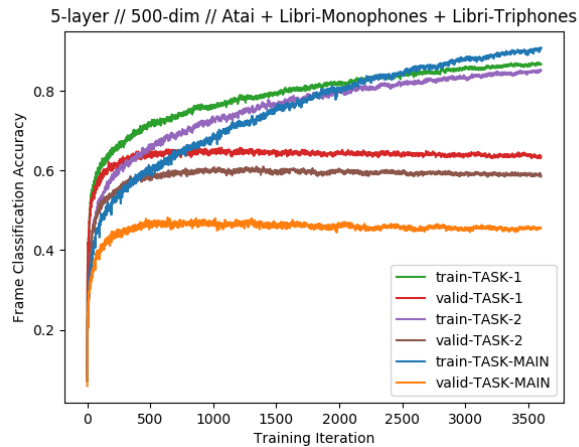


Figure 5: *Aux Task == LibriSpeech Monophones + Triphones*

Tasks	WER%
Atai (baseline)	
Atai + Libri-Triphones	
Atai + Libri-Monophones	
Atai + Libri-Monophones + Libri-Triphones	

Table 1: *Experimental Setup*

Takeaways from 500-dim // 5-epoch // 2-to-1 weighting:

1. Addition of LibriSpeech beats out Kyrgyz-only.
2. Triphones work better than monophones
3. Both languages / tasks overfit
4. best triphone model beats best baseline (ie. early stopping)
5. atai overfit slower with additional tasks
6. atai overfits with monophones earlier, I think because it's an easier task

5. Discussion

6. Conclusions

7. Acknowledgements