

1 What I'm doing & Why it's Extremely Cool

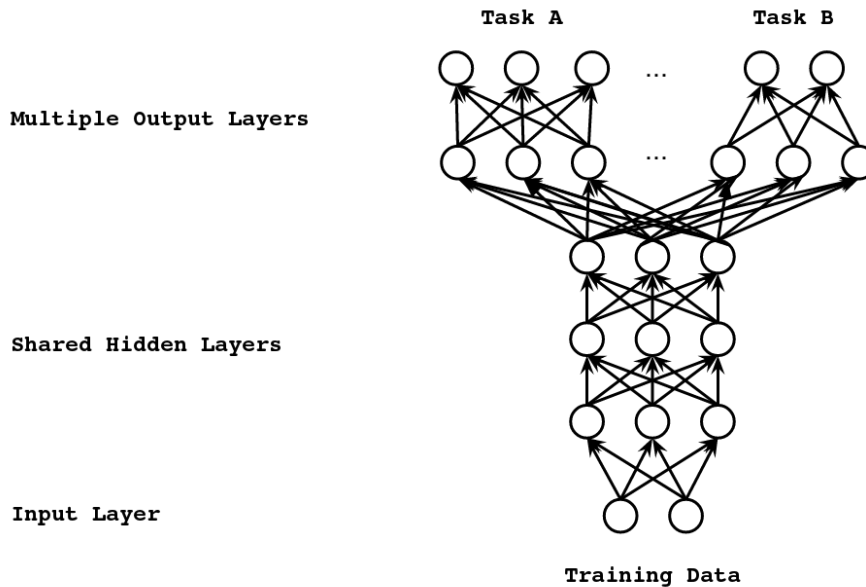
I'm creating Neural Nets which better classify data from unseen conditions, without any explicit adaptation of model parameters or data transformations. In my dissertation, the neural nets of interest are acoustic models for automatic speech recognition, and the unseen data conditions are (1) **new noise**, (2) a **new speaker**, or (3) a **new language**. For each of these three cases, I have two conditions of data. I train the model on one data condition, and then test the on other.

1. For the **Noise** condition, I train on clean audio and test on noisy audio.
2. For the **Speaker** condition, I train on a set of speakers, and I test on an unseen speaker.
3. For the **Language** condition, I train on (mostly) English, and I test exclusively on Kyrgyz. I must include *some* Kyrgyz in training because the DNN+HMM pipeline requires the testing phonemes defined in the model.

My approach is extremely cool because it doesn't require tons of data, it is not specific to some dataset, and it can be used to train any Neural Net (not just for speech recognition).

In the past people have dealt with new data by (1) adapting an existing model to the new data, or (2) normalizing the new data to look more like the old training data. More recently, Multi-Task Learning (MTL) has been found to produce models which are better at handling new data, because these models inherently learn more robust features which generalize to unseen domains.

MTL training for neural nets works by training a set of hidden layers to perform multiple tasks (multiple output layers). Here is an example of MTL architecture for an acoustic model from Heigold 2013.



This is the same architecture that will be used in my dissertation, but obviously I will have more layers and more nodes per layer. The number of input nodes on my neural nets corresponds to the dimensionality of audio features, and the number of nodes on each output layer corresponds to the number of phonemes (i.e. monophones or triphones) I've defined for the language.

In this dissertation (following the MTL tradition), each output softmax layer will be referred to as a **Task**.

2 Overview of Speech Recognition

This section contains an overview of the training and testing procedures for standard automatic speech recognition (ASR) pipelines. The overview will provide the reader with a technical grounding in ASR, so that the rest of the dissertation will have some point of reference.

- **Gaussians + HMMs**
- **Neural Nets + HMMs**
- **end-to-end Neural Nets**

3 Overview of Multi-Task Learning

A task (in classification) is a set of (data,label) pairs.

Most machine learning training uses single-task learning (e.g. classifying an image as a digit [0-9]).

In Mutli-Task learning, we learn multiple tasks which share useful information. An example of a non-useful auxillary task is classifying a number as greater or less than 7, when the main task is to identify the identity of a single digit.

4 Background Literature

Here I will cover the literature relevant to working with small (or completely new) datasets. There are two main approaches, (1) adapt a model from one training dataset to a new, smaller dataset; (2) create a model that is robust enough to handle data from multiple domains.

- **Model Adaptation:** (e.g. **Speaker; Language**)
- **Model Robustness:** (e.g. **Noise; Channel**)

5 Experiments

This section contains the main contributions of the dissertation research.

This dissertation investigates training methods for acoustic modeling in the Neural Net + HMM ASR pipeline.

I aim to produce acoustic models which perform better (i.e. lower Word Error Rates) on datasets which are not similar the original training dataset.

I investiage the effectiveness of different **Tasks** (eg. linguistic **Tasks** vs machine learning **Tasks**) in a Multi-task Learning framework.

5.1 Data

I am creating acoustic models which generalize well to new data. To measure how well the models generalize, I use a set of speech corpora which exhibit some interesting differences between training and testing data. These differences between corpora exemplify the typical challenges faced in speech recognition generalization.

The training and testing data will differ in either (1) the recording **noise** conditions, (2) who the **speaker** is, or (3) what **language** the speaker is using. The following table shows which data sets are used for each audio condition.

AUDIO CONDITION	CORPUS		
	Train		Test
	Noise	TIDIGITS	Aurora 5
	Speaker Language	LibriSpeech-A LibriSpeech	LibriSpeech-B Kyrgyz Audiobook

Table 1: Speech Corpora

5.2 Model Training Procedure

This dissertation investigates the creation of new tasks for MTL, either using (1) linguist-expert knowledge, (2) ASR Engineer-expert knowledge, or (3) general Machine Learning knowledge.

The former two knowledge sources are useful for building acoustic models, but not much else. On the other hand, the final knowledge source (general machine learning concepts) can be applied to *any* classification problem.

The three knowledge sources will be abbreviated as such:

- (LING) **Linguistic Knowledge**
- (ASR) **Traditional Speech Recognition Pipeline**
- (ML) **General Machine Learning**

Each of these categories contains a wealth of ideas, but I will consolidate each into three experiments. With three experiments for each knowledge source, my dissertation will contain nine (9) experimental conditions (for each audio condition).

Specifically, I will use the following concepts to create new tasks to be used in MTL training:

EXPERIMENTS	KNOWLEDGE SOURCE		
	LING	ASR	ML
	voicing	monophones	k-means
	place	1/2 triphones	random forests
	manner	3/4 triphones	bootstrapped resamples

Table 2: Experimental Setup

Each of these tasks will be added to a Baseline model. More specifically, the Baseline model will be a Neural Net with a single output layer (Task A), and the tasks above will be added as a second task (Task B). You can think of the tasks as simply a new set of labels for the existing data set. For example, when the LING task of VOICING is used, any audio segment labeled [b] will be assigned the new label **voiced**.

When these experiments will be applied to each of the three audio conditions, we get the following 30 experiments:

Data Condition	Train Data	Test Data	MTL Training Tasks	Num. Exps
NOISE	TIDIGITS	AURORA 5	Baseline	1
			Baseline + LING	3
			Baseline + ASR	3
			Baseline + ML	3
SPEAKER	LIBRISPEECH-A	LIBRISPEECH-B	Baseline	1
			Baseline + LING	3
			Baseline + ASR	3
			Baseline + ML	3
LANGUAGE	LIBRISPEECH + KYRGYZ-A	KYRGYZ-B	Baseline	1
			Baseline + LING	3
			Baseline + ASR	3
			Baseline + ML	3
				30

Table 3: Experimental Setup

5.3 Task Creation Specifics

1. 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** is used. The additional **Tasks** are dropped and the **Baseline 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.

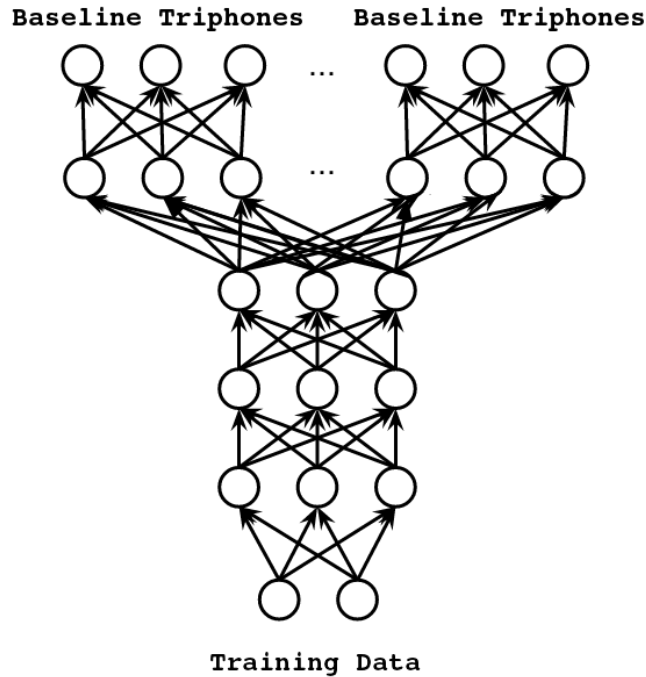


Figure 1: Baseline

2. LING

All of the linguistic knowledge tasks view the phoneme as a bundle of features.

Using standard features from articulatory phonetics (voicing, place, and manner), the following tasks generate labels for each data point by collapsing the given phoneme labels along one of these three dimensions.

All information from one dimension is removed from the labeled data. This forces the classifier to rely on audio signal features which do not relate to that dimension. The DNN must project the input data into a new space for classification, using only information from the other two dimensions.

(a) VOICING

Voicing information is removed from the data labels.

Speaker Robustness Experiments: The training data is a 4.5 hour subset of Librispeech, with mixed speakers, men and women. The testing data is 30 minutes of speech from two speakers (one man one woman).

First, two separate GMM-HMM models are trained on the training data. The first GMM-HMM model uses the standard CMUDict phoneset (39 phones + stress variants).

From this standard phoneset, the normal 3-state monophones are trained from a flat-start via EM training. A total of XXX states are trained with a total of XXX Gaussian components over XXX iterations of EM. These monophones are then expanded into context-dependent triphones via a phonetic decision tree, with a maximum of XXX leaves. The resulting leaves (state clusters) are then trained with XXX Gaussian components over XXX iterations of EM. The final model achieves a WER of XXX on the testing data.

Phoneme	Example	Translation
AA	odd	AA D
AE	at	AE T
AH	hut	HH AH T
AO	ought	AO T
AW	cow	K AW
AY	hide	HH AY D
B	be	B IY
CH	cheese	CH IY Z
D	dee	D IY
DH	thee	DH IY
EH	Ed	EH D
ER	hurt	HH ER T
EY	ate	EY T
F	fee	F IY
G	green	G R IY N
HH	he	HH IY
IH	it	IH T
IY	eat	IY T
JH	gee	JH IY
K	key	K IY
L	lee	L IY
M	me	M IY
N	knee	N IY
NG	ping	P IH NG
OW	oat	OW T
OY	toy	T OY
P	pee	P IY
R	read	R IY D
S	sea	S IY
SH	she	SH IY
T	tea	T IY
TH	theta	TH EY T AH
UH	hood	HH UH D
UW	two	T UW
V	vee	V IY
W	we	W IY
Y	yield	Y IY L D
Z	zee	Z IY
ZH	seizure	S IY ZH ER

Table 4: CMUDict Phoneset

The second GMM-HMM model trained differs from the first model in its set of initial phones. Instead of building monophones (and then triphones) from the standard CMUDict, this **-Voicing** model collapsed all voicing information from the phonetic dictionary (i.e. the lexicon file).

B P --> P
 CH JH --> CH
 D T --> T
 DH TH --> TH
 F V --> F
 G K --> G
 S Z --> S
 SH ZH --> SH

(b) PLACE

All place information is removed from the data labels.

F TH SH S HH --> F	voiceless fricatives
V DH Z ZH --> V	voiced fricatives
P T K --> P	voiceless plosives
B D G --> B	voiced plosives
M N NG --> N	voiced nasal
L R --> R	voiced laterals
Y W --> Y	voiced approximants

(c) MANNER

All manner information is removed from the data labels.

B M V W --> W	voiced labials
P F --> P	voiceless labials
D Z --> D	voiced alveolar
N L R --> R	voiced alveolar2
T S --> T	voiceless alveolar
ZH JH --> JH	voiced postalveolar
SH CH --> CH	voiceless postalveolar
NG G --> G	voiced velar

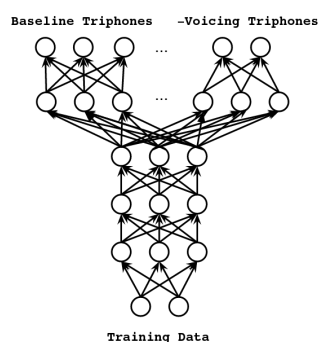


Figure 2: -Voicing

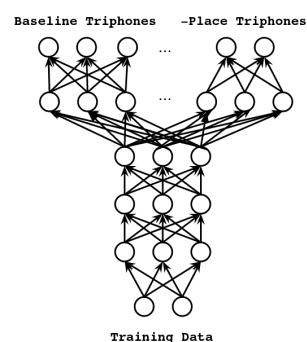


Figure 3: -Place

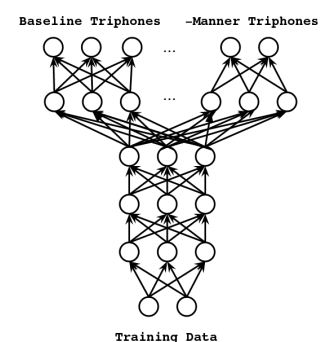


Figure 4: -Manner

3. ASR

All the following tasks relate to the structure of the phonetic decision tree used in the traditional ASR pipeline to cluster context-dependent triphones. In GMM training the leaves of the decision tree are then assigned their own Gaussians, and in DNN training the same leaves are used as labels during training via backprop.

The main intuition behind these experiments is that in using the decision tree labels as targets for the DNN classifier, we are performing model transfer. The decision tree and its associated Gaussians perform classification, and we are merely training a DNN to perform that same task. So, the decision tree can be thought of as a single task for the DNN to learn.

However, the DNN only sees the leaves of the decision tree. It doesn't see any of the branches, or any of its wonderful linguistic structure. So, in order to force the DNN to learn the information hidden in the decision tree, the following tasks are like cross-sections of the tree, slicing it from leaves up. The DNN then has to learn how to read these cross-sections, and how to map data onto each layer.

If we slice the tree at the roots, we have the MONOPHONES. If we slice down half-way ($1/2$ TRIPHONES), we have more contextual information than monophones but less than full triphones. If go a little farther down ($3/4$ TRIPHONES), we get even more context, but less general information about the original phoneme.

(a) monophones

When we chop the tree at the roots.

(b) $1/2$ triphones

Chop the tree half-way down.

(c) $3/4$ triphones

Chopping a little further.

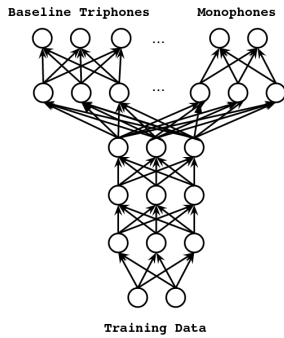


Figure 5: Monophones

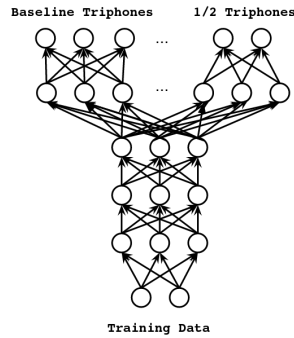


Figure 6: $1/2$ Triphones

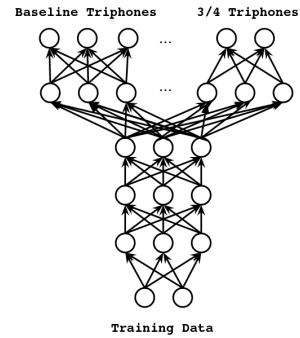


Figure 7: $3/4$ Triphones

4. ML

The following tasks do not make use of any linguistic knowledge or any part of the ASR pipeline. The only things needed to perform these tasks is labeled data.

The two approaches above use linguistics or the ASR pipeline to force a DNN to learn structure about the data, because that information is useful for classification.

We typically do not have this kind of *a priori* information about the datasets we use in Machine Learning. Therefore, an interesting problem is how to learn this structure in a data set when we don't have access to that expert knowledge.

The following tasks force the DNN to learn structure in the data without any knowledge about that structure. In order to do so, I make the assumption that the data does in fact have heirarchical relations. That is, I assume the (data,label) pairs were generated by something like a phonetic decision tree, and I try to recover the structure of that tree.

(a) k-means

Standard k-means on the data, with the caveat that labels cannot be split across clusters. A first round of clustering is performed, and then all data from the same original label are shifted to the cluster with the most data points from that label. Then, centroids are recalculated, and data is re-clustered. This adapted k-means should find related data points in the same clusters. If k-means is working, we would expect to be able to recover phonemes (monophones) from the labeled triphone data.

(b) random forest

In another attempt to cluster triphones along phonetic categories, the random forest procedure works as follows: (1) take a random subset of the labels, (2) train a random forest with all data points associated with those labels, (3) re-classify all the rest of the data with the new random forest. In this way, we will reduce the number of labels (eg. out of 2,000 triphone labels I choose 500), and classify unseen data as its closest seen data point.

(c) bootstrapped resamples

In this approach, new labels are not generated at all. The separate tasks for the DNN are just different samples of the data.

Some sub-samples may exhibit a more useful decision plane than others, and if we randomly subsample for multiple tasks, the different decision planes will all have something in common. The individual peculiarities of one sub-sample will have to be ignored for the DNN to perform well on all tasks.

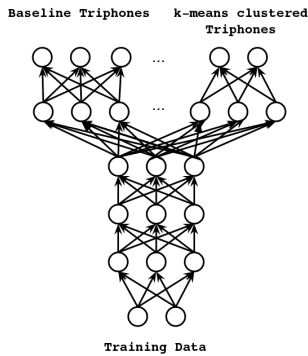


Figure 8: k-means

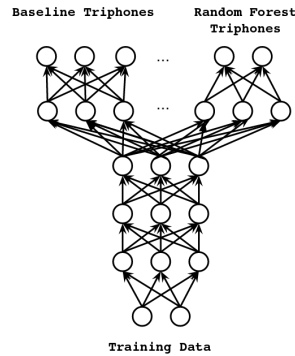


Figure 9: 1/2 Triphones

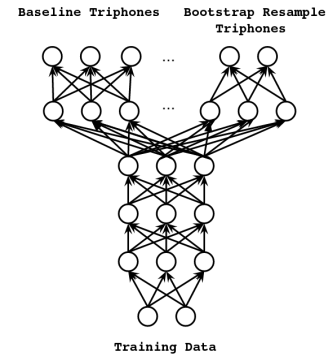


Figure 10: 3/4 Triphones