

My Research Passion

Everyone should have access to speech technology in their native language.

This is the conviction that drives my research. I want the world to be a place where someone who's born blind doesn't have to settle for a lesser education because of lack of audiobooks. I've worked with people in that exact situation to build a Kyrgyz speech synthesizer, because they want more opportunities for other blind Kyrgyz speakers. I've played a small role to make speech synthesis available for free, but there is a hard limit as to what I can accomplish on my own.

The Google AI residency is the next step to accomplish my goals. I want to publish influential work and get more people caring about low-resource languages. With my knowledge of linguistics and passion for machine learning, the Google AI team are the perfect place for me to flourish and share my research with others.

Collaborating with researchers at Google, I will help create new approaches and algorithms, testing theories, and diving deeper into multi-task learning. In particular, I have many ideas for multilingual experiments, extracting information from large datasets to transfer knowledge to smaller domains. The competitive (yet collaborative) research atmosphere at Google excites me, and I will bring my enthusiasm and dedication with me.

Overcoming Obstacles

Working with low-resource languages, my thesis topic is Multi-Task Learning for deep neural net acoustic modeling in Automatic Speech Recognition (ASR).

Acoustic models in ASR are statistical classifiers, accepting as input some speech signal (ie. time chunks of an audio file), and returning probabilities over phonetic transcriptions. Using Multi-Task Learning, my research accomplishes this audio \rightarrow phonetic transcription with a feedforward neural net with multiple output layers. Each output layer represents a different task the net is required to perform. If two tasks are related, they will contain overlapping information. Learning related tasks in parallel should improve performance on any one of the tasks, because the weights in the net will be biased towards general, task-independent representations of the data.

The Multi-Task approach offers an elegant way to exploit small datasets, if you can come up with the right tasks. My early experiments show that by adding tasks along dimensions relevant to linguistics (place, voicing, manner), we can lower Word Error Rates for smaller data sets. I spent hours carefully defining each task, consulting both the machine learning and linguistic literature. Via trial-and-error, I eventually found tasks that showed improvement over my baseline. However, defining each task by hand is not a scalable approach. Since then, I've devoted my research to finding scalable, automatic solutions to task creation.

You can view Multi-Task classifier as a kind of ensemble model, and once you do, you open up a flood of inspirations from past work. Looking for automated ensemble model approaches, I came across the Random Forest.

The Random Forest is trained not by combining different architectures or different labels, but by merely resampling the data. Each subsample has its own

decision plane, specific to the data and noise in that particular sample. By averaging the votes of all trees, the noise is ignored, but the generalities remain. Taking inspiration from Random Forests, my current research investigates Multi-Task neural nets, where the different tasks are independent subsets of the data. I use the Kaldi Speech Recognition toolkit to run experiments.

Typically, I perform my experiments in this way: After I've defined a new training procedure, I run a toy model on a small subset of the data, to quickly make sure there's no bug in the code. Then, I run some practice models with fewer parameters and epochs on a larger subset of the data, to get a ball-park idea of model performance. Finally, I define an architecture with a real number of parameters and epochs on all the data, and set a loop to train multiple versions of the same model (to get a standard deviation for the differences in random initializations).

Random bootstraps of the data offer an easy way to average out misleading noise, but there's a catch. Typically, Random Forests are trained with hundreds of trees, if not thousands. Training decision trees takes a lot less time than training a neural net via backprop, and as such, random multi-task generation does have a down-side.

A better approach would be to automatically generate fewer tasks, but be sure that each task is more meaningful than just a random chunk of the data. My linguist-crafted tasks were all built by forcing the net to learn commonalities in the data at a level higher than the labels themselves. A corollary in an image recognition task like CIFAR-100 would be if I worked with a zoologist to identify not only species, but genres, and families, etc. A label-level task would be distinguishing wolves and coyotes, while a level higher would be distinguishing canines and felines, therefore understanding that wolves and coyotes share common features.

Instead of working with a domain-expert like a zoologist or linguist, it would be better to automatically discover those higher level groups automatically. This brings me to my next planned experiment (after I get Random Neural Forests working): using clustering methods like k-means to generate higher-level labels for the data, and using those labels as a second task.

The open-endedness of research does not phase me. I know my goal (make speech technology for all languages), and I know it is so big that I can happily spend a career chipping away at it, and Google AI is the perfect place to start chipping.

Thank you for your time and consideration.
