

# Evaluating and Developing Methods of Generating Code-Switched Data 07-400, Spring 2022

Abhishek Vijayakumar,  
under Prof. Alan W. Black, Language Technologies Institute  
<https://inkyubeytor.github.io/category/code-switching.html>

March 14, 2022

## 1 Major Changes

The primary goal of this project has not changed since the previous milestone.

## 2 What You Have Accomplished Since Your Last Meeting

Since our last meeting, I have determined the learning rate range in which I can train models effectively. I have trained several models with different learning rates and looked at their performance on part of speech tagging to determine at what learning rates training breaks down.

I have also identified a problem with the sentiment analysis task (see Surprises) from the GLUECoS benchmark.

## 3 Meeting Your Milestone

My primary goal was to identify the hyperparameters I should be using for efficient pre-training to actually have an impact on model parameters.

## 4 Surprises

It seems that regardless of training scheme, the model is producing outputs with most outputs belonging to a single class. I believe this applies to the published benchmark numbers as well, indicating that the benchmark number is not a good measure of performance. I have not been able to train a model to overcome this problem.

## 5 Looking Ahead

My goal for the next two weeks is to achieve a statistically significant performance improvement on an evaluation task via some form of pretraining from a baseline model, such as DistilBERT. This is the first step towards providing a metric for how useful generated code-switched data is for pretraining models.

## **6 Revisions to Your Future Milestones**

My future milestones may undergo revisions based on whether I can meet my new next milestone.

## **7 Resources Needed**

No further resources are needed for this project at this time.