# Seq2Seq Modeling for Grammar Correction
## Emanuel Cortes
## Department of Computer Science, Stanford University

## Motivation

- Develop model to correct simple grammatical errors in sentences.
- Goal: To create a writing tool to help anyone learn grammar and improve their grammar skills when working on a writing activity .

## Data Preparation and Training Pipeline

Query Student Submissions (Source) with Correct Responses (Target)

⬇

Filter responses by levenshtein distance between source and target to encourage model to learn simple mistakes

⬇

Group responses based on length of source and target to avoid over-padding, referring each group as a "bucket"
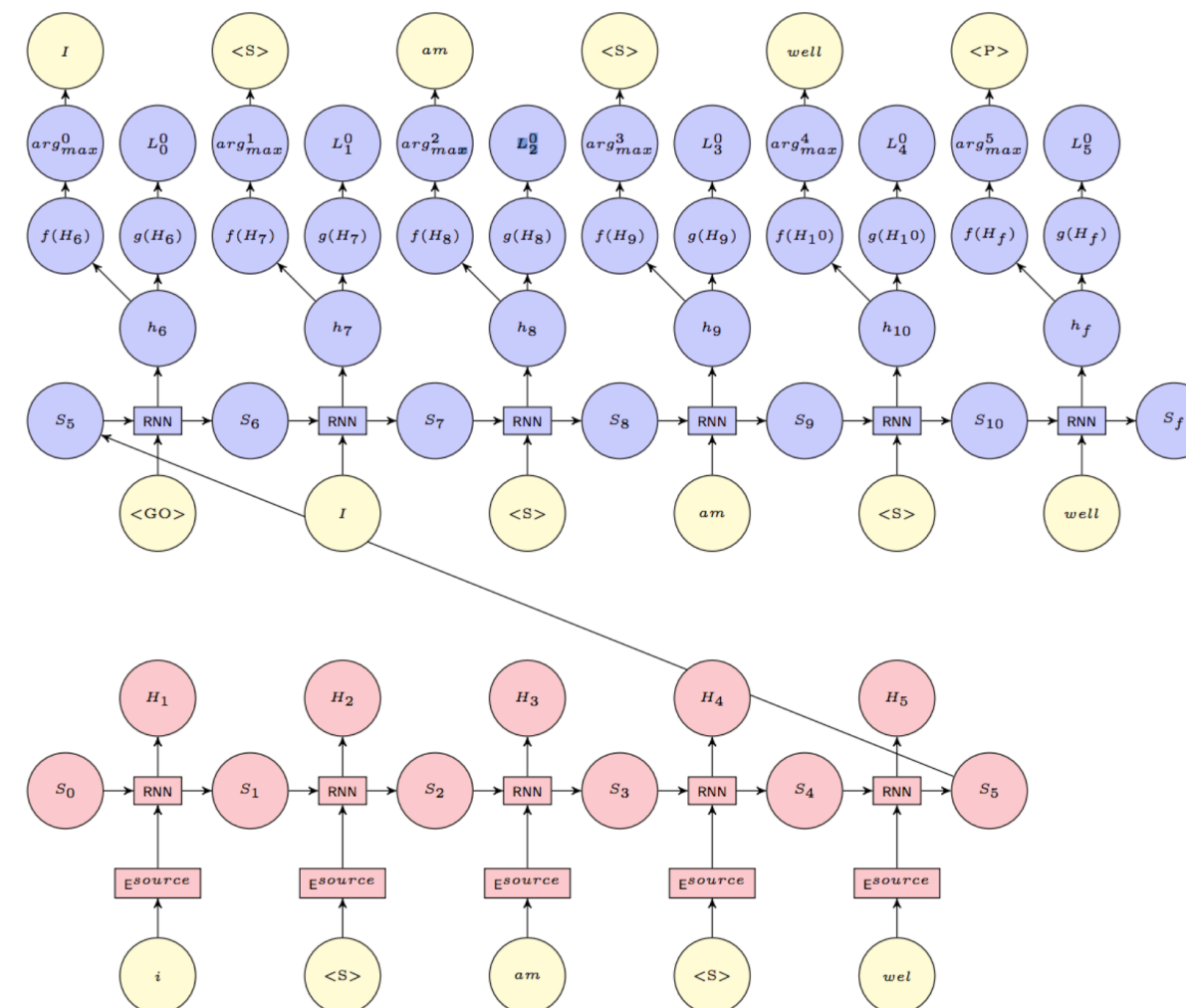
⬇

A unique model will be trained on each available bucket of data

⬇

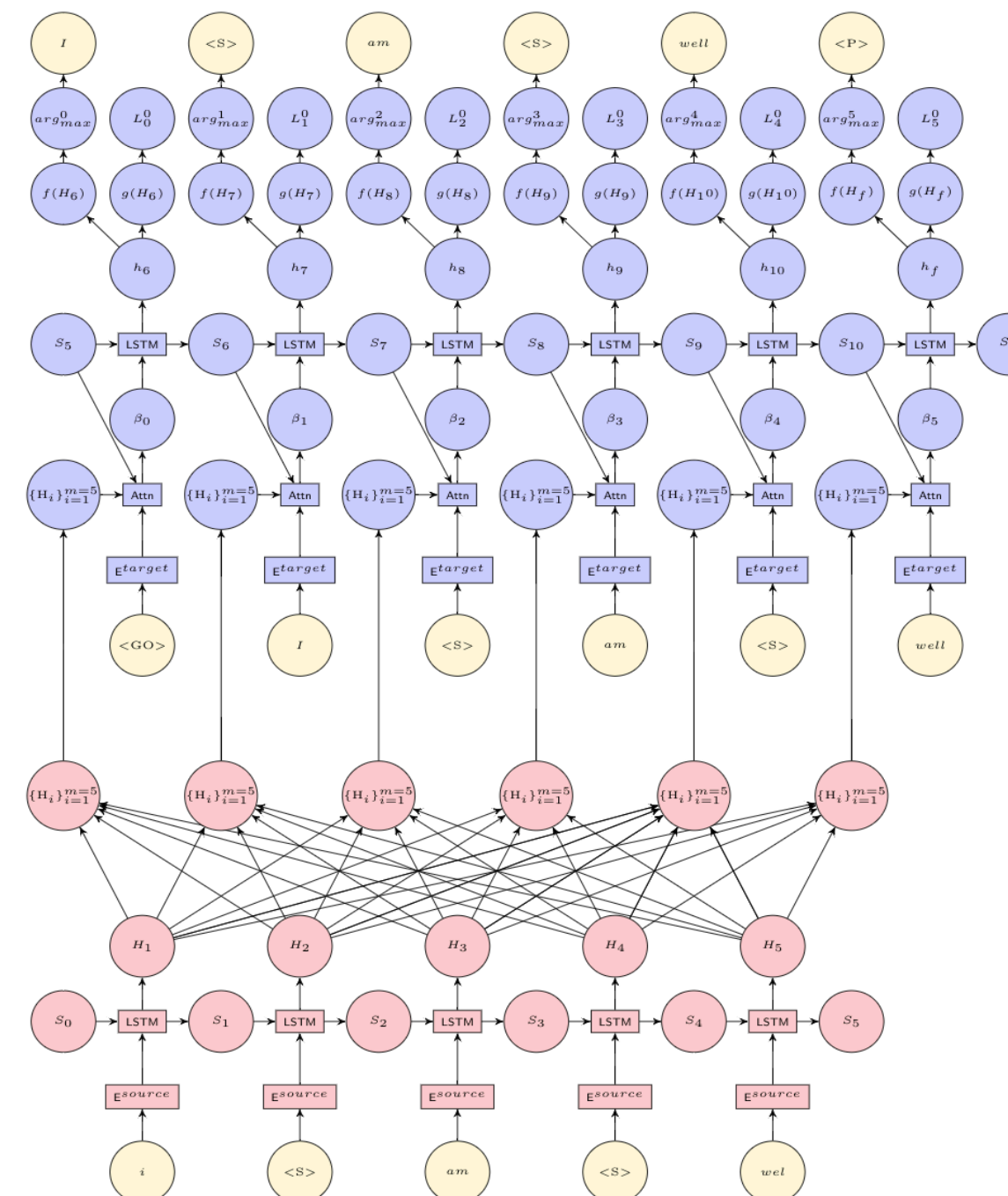For each bucket of data, optimize the loss in each model (across all buckets)

⬇

Learn mappings from word to vector representation and learn to predict the correct sentence given an incorrect sentence.
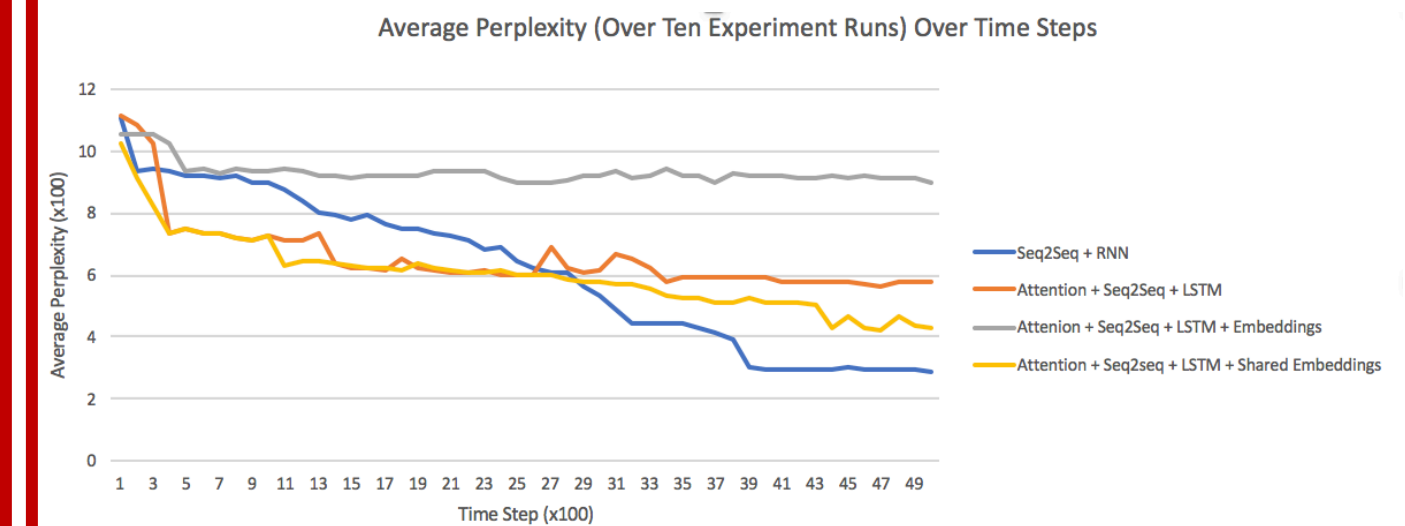
## Models

1.) Baseline:



2.) Attention Seq2Seq:



## Experimental Results



Average Perplexity (Over Ten Experiment Runs) Over Time Steps

- All average BLEU scores are near zero.
- For attention seq2seq + LSTM + shared embeddings, model is able to correct capitalization mistakes for smaller sentences.

## Discussion

- Attention seq2seq model with embeddings may need to be trained for a longer period of time to achieve better performance.
- The model seems to train faster when sharing embeddings between encoder and decoder.
- With limited gpu/cpu resources, I need to balance expressive power of model with time complexity.

## Future Work

- Less Complexity: Implement Sampled Softmax and Share Variables Across All Buckets.