

# CIS545 Project Proposal: Style Transfer Exploration with Neural Net Models

*Peter Kong*

## 1 Background

Style transfer is the process of transforming the style of a text into a new style while retaining its original meaning. There is active research in style transfer - [Fidler, Hu, Prabhumoye]. Unfortunately, the non-neural metrics that these authors use to evaluate their style transfer models (like BLEU, ROUGE, and perplexity) offer no mechanism for separating style and content information from an input text.

Thus, I focus on neural net classifiers. In certain configurations, neural classifiers are also capable of accepting arbitrary, external data via injection into various hidden layers. (Hu) provides details of a wake/sleep variational auto encoder model that does this. In the future, a classifier built in this CIS 545 project may be able to be used as a component in a style transfer pipeline - it could accept a content vector as input and use the knowledge within its weight matrix to style that content.

## 2 Objective

To replicate existing best-in-class neural net style evaluators (classifiers) and possibly improve on their benchmarks.

## 3 Action Plan

### 3.1 Dataset

Although (Hu) used the IMDB review corpus, I prefer to use project Gutenberg's author-based text datasets since they are closer to "literary" English and offer greater diversity of style and content.

### 3.2 Data Wrangling

This project will require a significant amount of data wrangling and pre-processing. Gutenberg provides a large catalog of every entry in its database. This catalog must be parsed and filtered for relevant articles (I plan to choose a small set of authors and retrieve their work only).

Once the catalog is parsed, each article must be downloaded via a customized script. Article text must be sanitized for textual artifacts (example: Gutenberg-inserted boilerplate preambles) and tokenized.

Labels can be programmatically assigned, providing that we assume that each author is his/her own style.

### 3.3 Data Visualization

Visualizing the data will be important for generating insights during the exploration process. I plan to create visualizations from clustering algorithms and use metrics like document cosine similarity to find common clusters of documents.

### 3.4 Feature Engineering/Model

I plan to use Keras because of its straightforward interface and support for VAEs. Word embeddings seem like a useful tool to use as well, space and time permitting.

## 4 Related Work

(Ficler) "Controlling Linguistic Style Aspects in Neural Language Generation". Jessica Ficler and Yoav Goldberg, 2017. <https://arxiv.org/pdf/1707.02633.pdf>

(Hu) "Toward Controlled Generation of Text". Hu, et al. <https://arxiv.org/abs/1703.00955>

(Prabhumoye) "Style Transfer Through Back-Translation". Prabhumoye, et al. <https://arxiv.org/abs/1804.09000>