# Eliot Kmiec

## Data Scientist

*About Me*

- Augustana College 2020
  - Biochemistry, Public Health
- Fields of Interest:
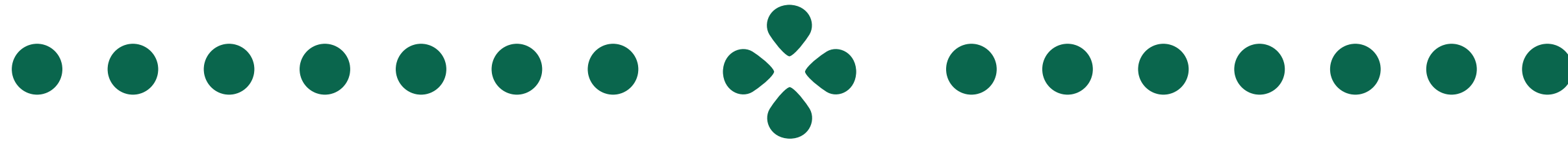  - Bioinformatics
  - Cloud Computing

# Why Automate Database Indexing?

- Subject headings are powerful indexing tools
- Labor intensive
- <u>Ongoing work</u> to automate MeSH for PubMed
- Saves Money, Streamlines Research
- ***Bring this technology to arXiv***

# Implementing Automated Indexing
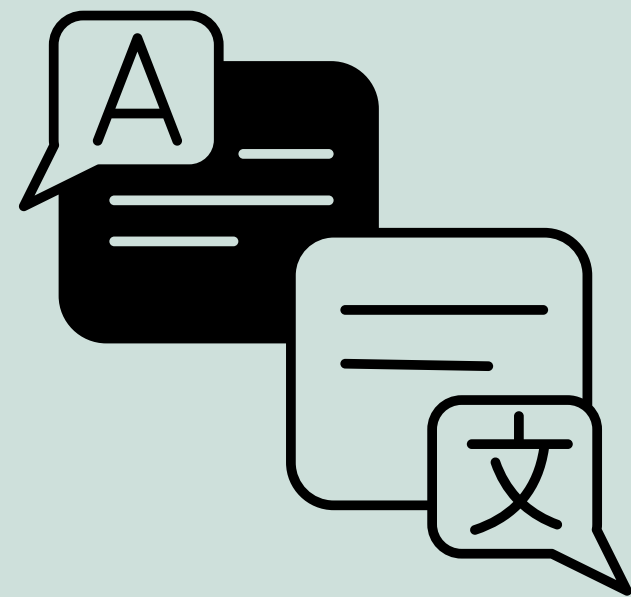
## Supervised Arm

Identifying the subject of an article

## Unsupervised Arm

Identifying new subjects for the database

# Approach to Modeling

NLP with RNNs

### Data Structure

- Abstracts : Subject Headings
- Available from <u>Kaggle</u>
- 1.5 mil abstracts

### Technical Hurdles

- Data Sparsity/High Dimensionality
- Numerous Categories/Subjects

### Model Architecture

- Recurrent Neural Networks
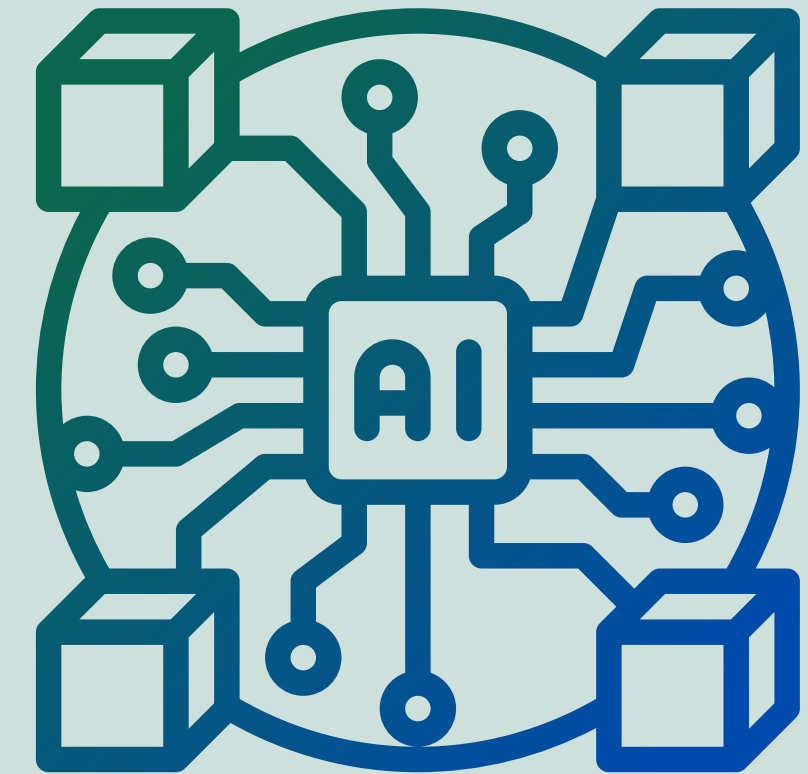- Retain more info VS. traditional models

*The Good*

- RNNs scale well
- Accurate even with few examples

*The Bad*

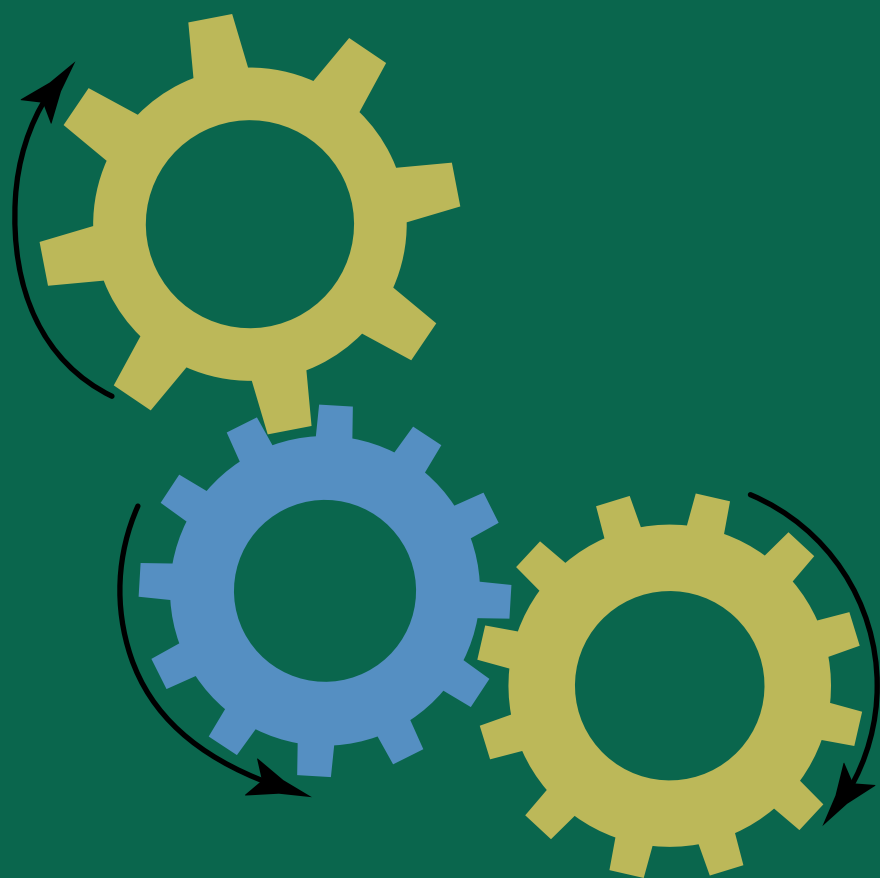- Lots of parameters
- Lots of resources to train

*The Ugly*

- Prone to gradient instability
- Initial models had "exploding gradient"

# Using RNNs for NLP

Training Hurdles

# Model Training & Gradient Instability



Loss

Accuracy

Epochs

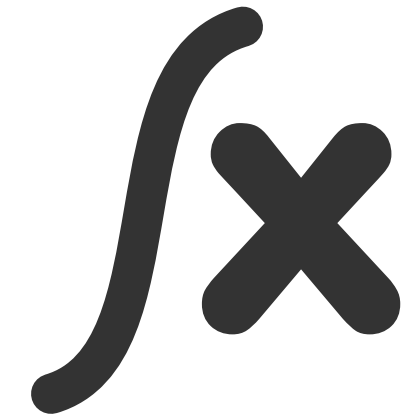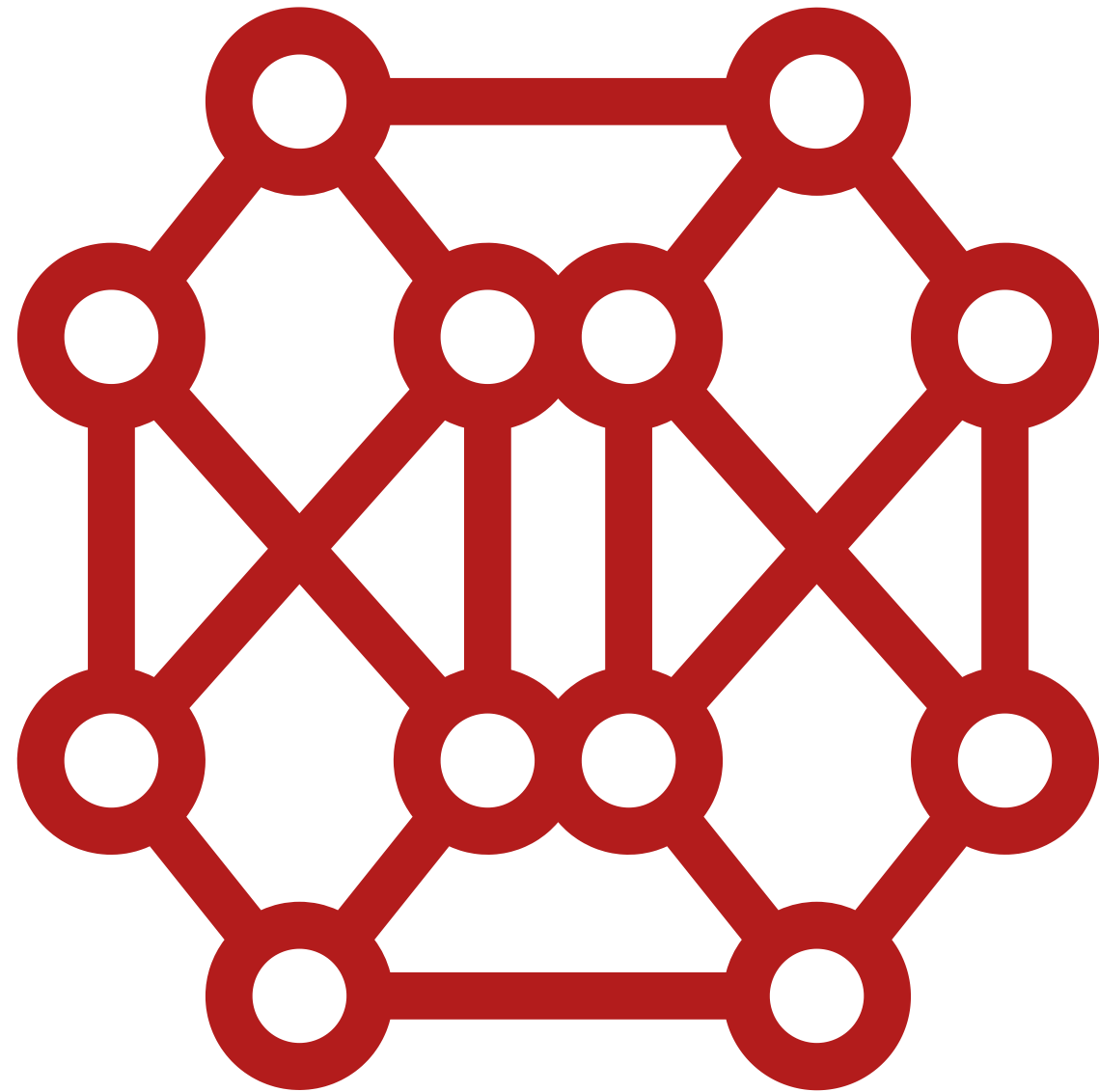| Accuracy | Loss | AUC Score |
|:---:|:---:|:---:|
| 51.7% | 0.541 | 0.395 |

# _Test Scores and Remaining Issues_

- 10-classes with 18% dummy model accuracy
- Varying degrees of convergence for most models
- MeSH studies have shown near 80% accuracy

# Next Steps

- Parallelize training on GPU cores
- Further work on stabilizing the gradient
- Expand model training to all subjects

# Questions?

Project Github

https://github.com/ek775

Eliot Kmiec