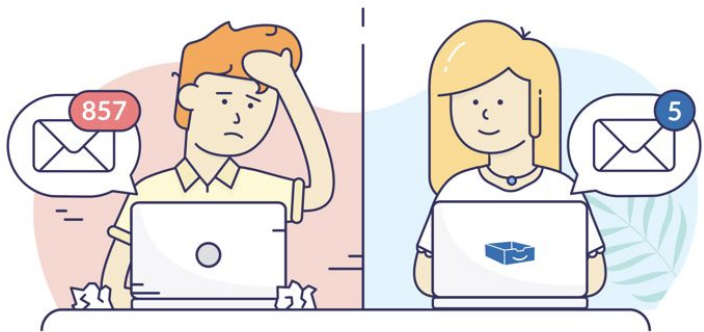# Introducing... ClutterCutter!

The smart, effective tool for your messy inbox

# Introduction and Motivation

❝

## The Problem

A messy, disorganized inbox makes life more complicated than it needs to be.

## Why do we need a solution?

- To save valuable time
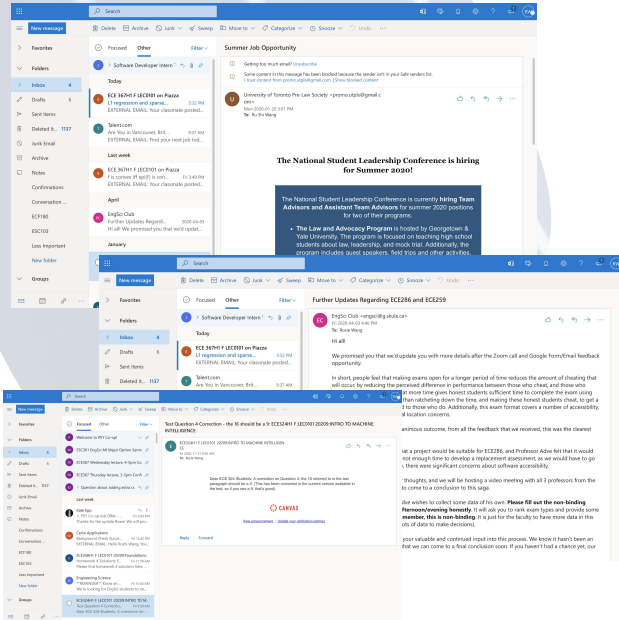- To make email search easier
- To increase the quality of life

# The Solution Process

"

1. Provide examples of how we sort past emails
2. Apply ML to learn the patterns
3. Apply this knowledge to make predictions
4. Automated sorting for future emails!

# Data & Data Processing



## Data Used

- Real-life emails from our own accounts

## Data Cleaning

- Standardized text
- Removed personal info for privacy reasons

## Amount of Data

- <u>Initial</u>: 150/class for a total of 750
- <u>Re-sampling</u>: Increased to ~200/class for ~1000 total
- <u>Data Augmentation</u>: Increased to 2850 total
  - Synonym Replacement
  - Back Translation

# Data Statistics



**8872 unique words** in 2850 emails!

Below are the 5 most frequently appearing "meaningful" words for each category

### Academics:
- 'course': 124, 'students': 121, 'final': 82, 'questions': 82, 'engineering': 81

### Alerts:
- 'order': 218, 'email': 169, 'account': 147, 'password': 71, 'information': 70

### Personal:
- 'practice': 45, 'think': 32, 'going': 20, 'other': 19, 'email': 19

### Professional:
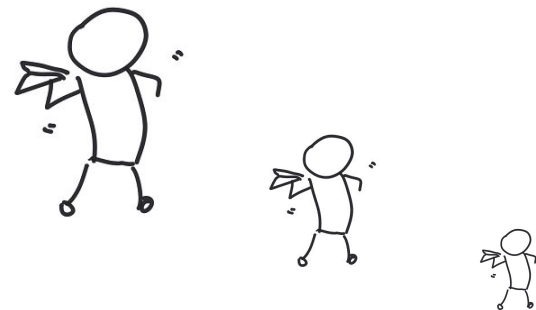- 'application': 177, 'interview': 145, 'position': 100, 'interest': 68, 'assessment': 59

### Promotions and Events:
- 'students': 275, 'engineering': 198, 'november': 151, 'student': 144, 'university': 143
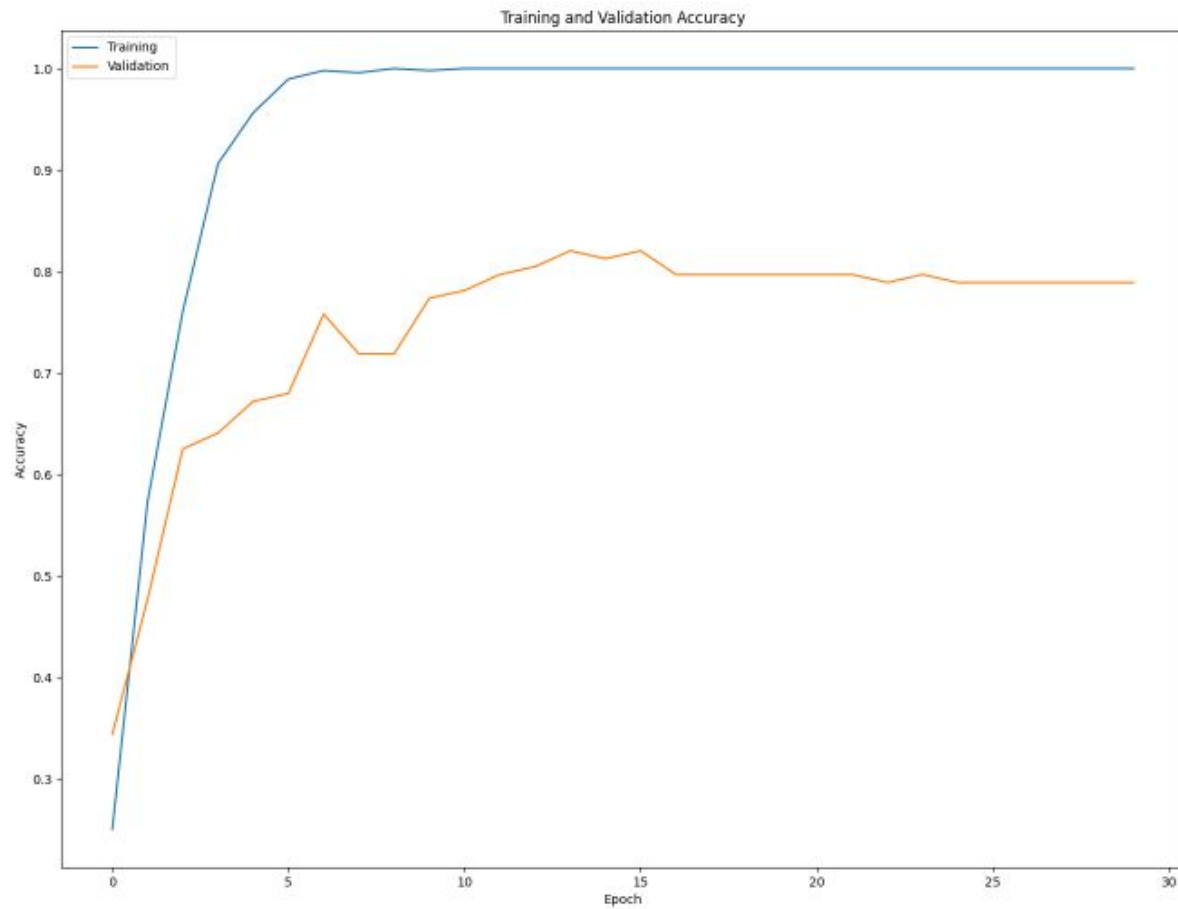
# Data Statistics cont.d...

## Data Splits

- Training = 0.64
- Validation = 0.16
- Test = 0.20
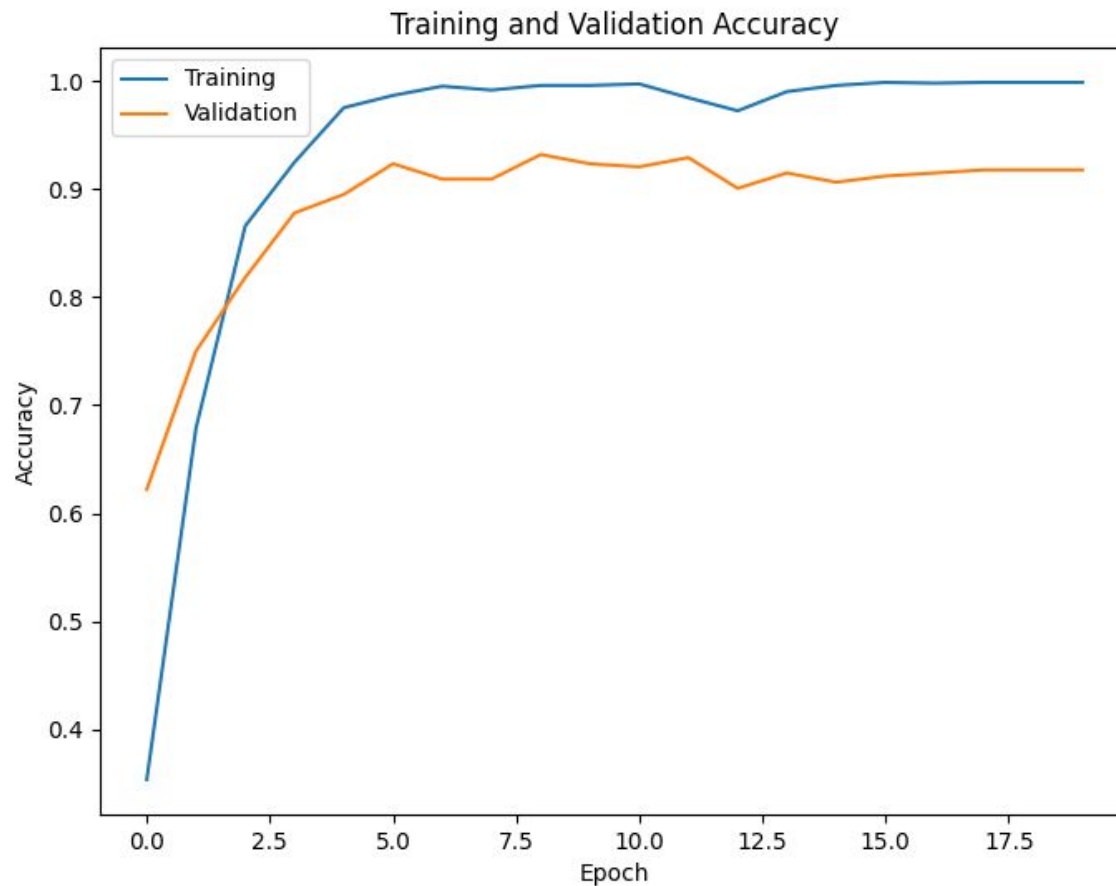
## Class Balance

| Labels | Train | Validation | Test |
|---|---|---|---|
| Academics | 355 | 89 | 111 |
| Alerts | 370 | 93 | 116 |
| Personal | 363 | 91 | 113 |
| Professional | 376 | 94 | 118 |
| Promotions and Events | 355 | 88 | 111 |

# Case 1: No Data Augmentation



Training and Validation Accuracy

# Case 2: 35% Data Augmentation



Training and Validation Accuracy

**Case 3:** 50% Data augmentation



Training and Validation Accuracy

# Our Model

## Vectorization

- Word embeddings using <u>Fasttext</u>

## Baseline

- Took <u>average meaning</u> of each email
- Fed to linear layer to produce 5 outputs

## RNN

- Single <u>GRU</u> unit
- Initial hidden dimension = 50
- Fed to linear layer to produce 5 outputs

class label

Output

$h_t$

$h_{t-1}$

RNN (GRU)

$X_t$

back propagation

- calculation of gradient & loss
- update parameters

input to RNN (GRU)

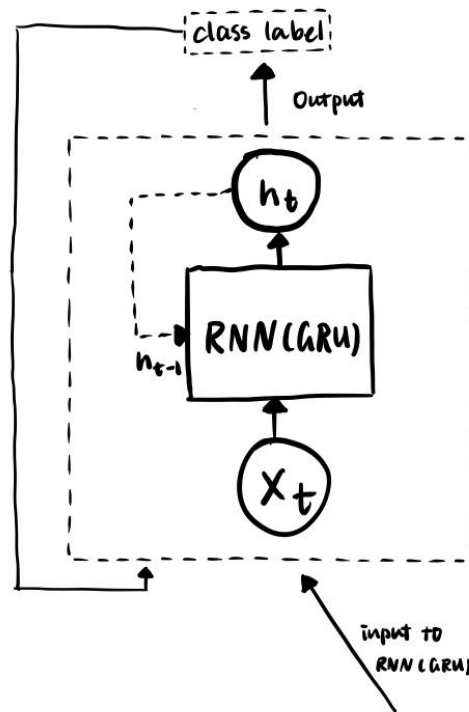raw email samples

preprocessing

labelled .txt files
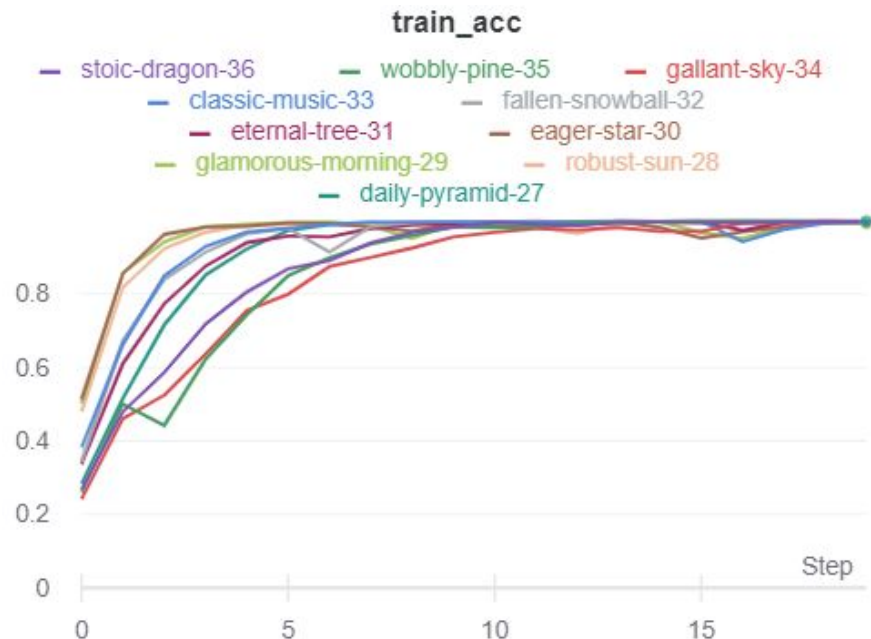
data augmentation

more labelled .txt files

Vectorization

word embeddings

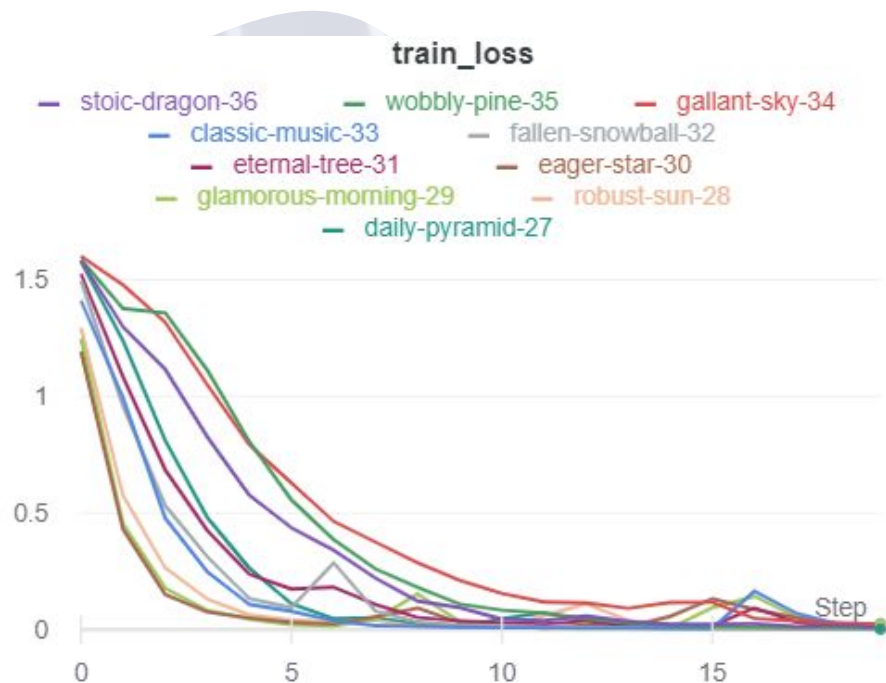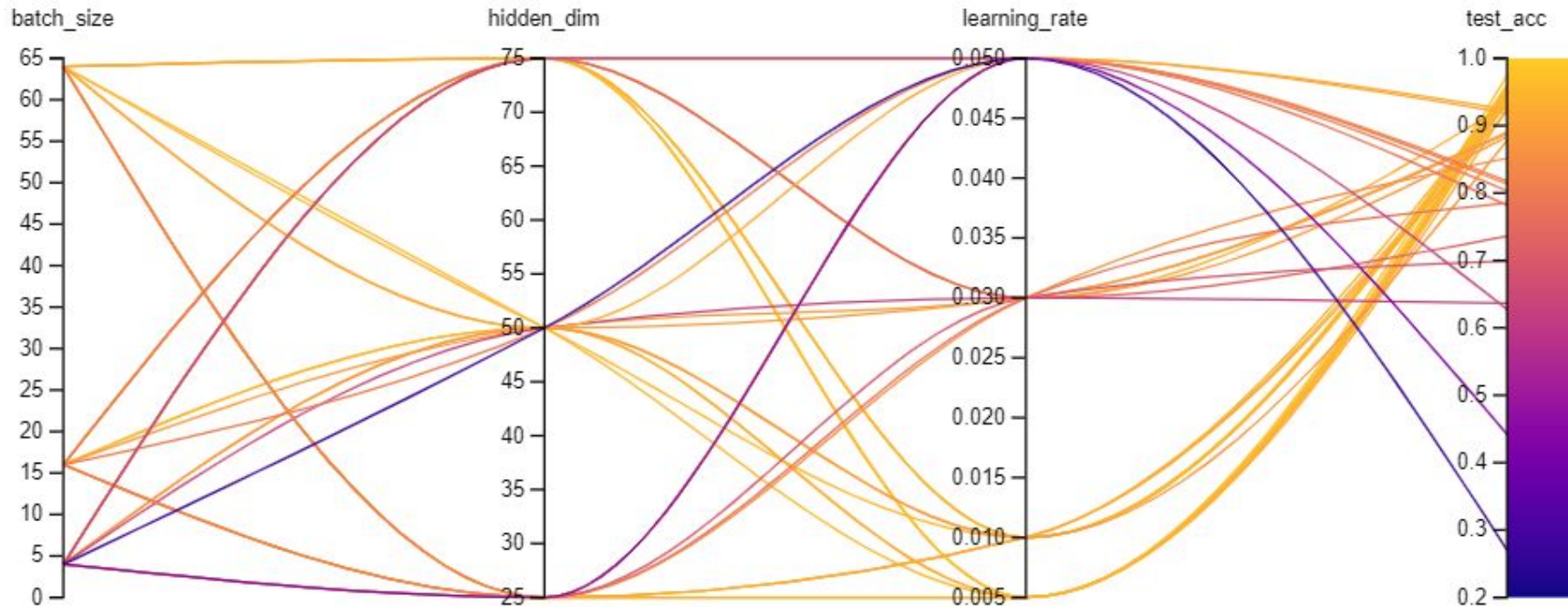# Hyperparameter Search: Loss and Accuracy Curves

# Summary Sweep of Hyperparameters

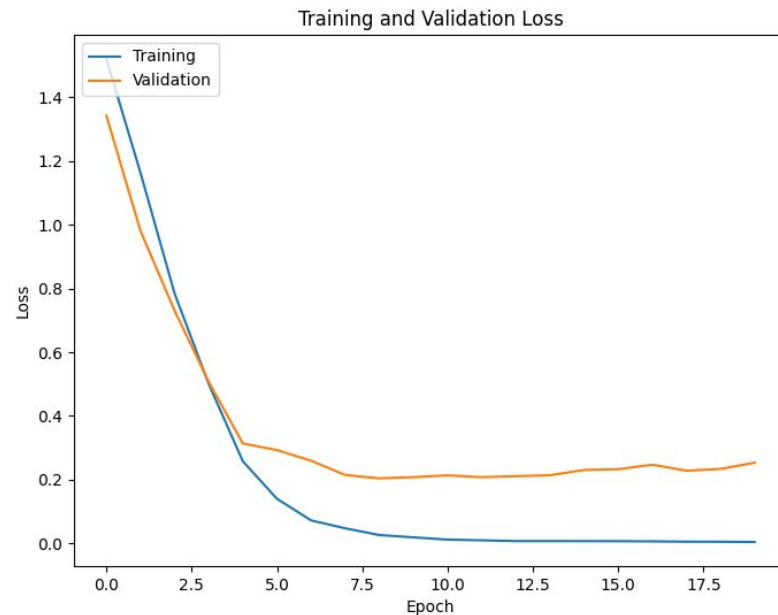# Best Model

**Test Accuracy:** **96.65%**

**Hyperparameters:** Batch size - 32, Learning rate - 0.005, Hidden dimension - 75

# Video Demo

# Video Summary

## From the video we saw

- Ex 1: Question from Piazza was correctly predicted as <u>Academics</u>
- Ex 2: Ad for job opportunities was predicted as <u>Promotions & Events</u>

These examples were from new emails the model has never seen!

# Visualizing Errors

- The **poorest performing** class was <u>Promotions & Events</u>
  - Actual <u>P&E</u> got confused for <u>Academics</u>
  - Most samples falsely labeled as <u>P&E</u>

# Examples of Incorrect Predictions 1

**Prediction: Promotions & Events**
**Ground Truth: Academics**



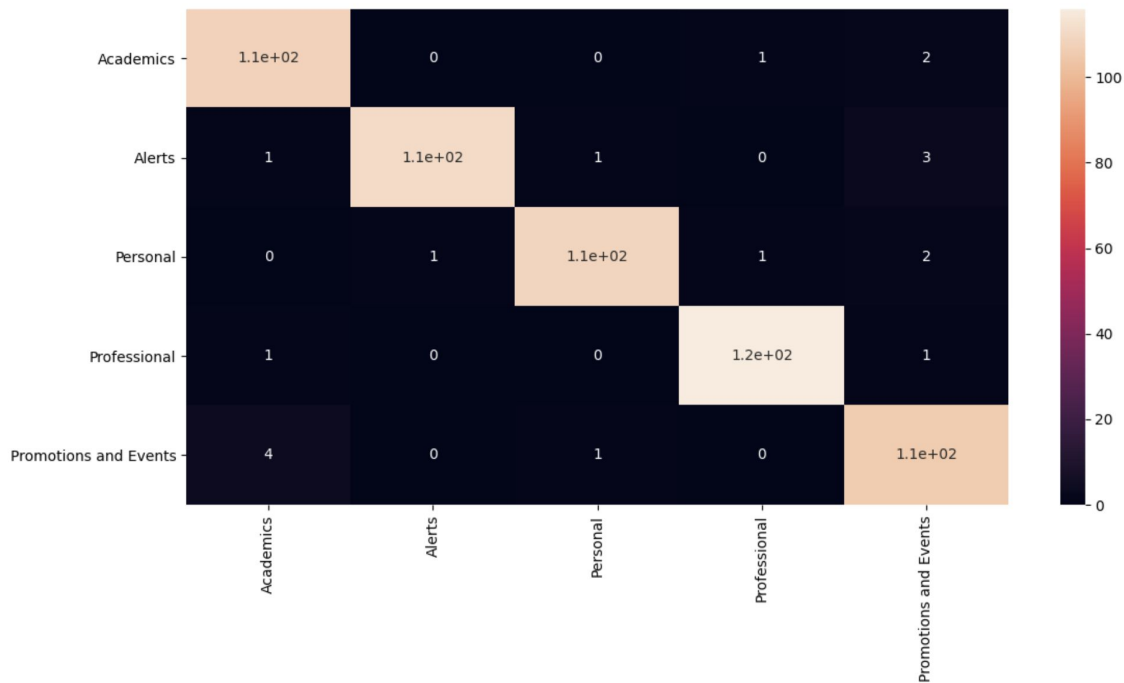Course Evaluations: Closing Soon! Complete Yours Today! (Engineering)

CE Course Evaluations <course.evaluations@utoronto.ca>
Sun 2020-11-29 1:25 AM
To:

**Make your voice heard!**

Hello,

Now is your opportunity to evaluate your courses and instructors. We appreciate you taking the time to provide comments to your instructor, as they can use this feedback to reflect upon their own work in these challenging times. Course evaluations are only open for a limited time, so complete yours today!

**Online evaluations are currently available for the following courses:**

| COURSE | INSTRUCTOR | EVALUATION END |
|---|---|---|
| FOUNDATIONS OF COMPUTING ECE358H1-F-LEC0101 | Veneris | December 10, 2020 |
| INTRO TO MACHINE INTELLIGENCE ECE324H1-F-LEC0101 | Rose | December 10, 2020 |
| MATRIX ALGEBRA & OPTIMIZATION ECE367H1-F-LEC0101 | Draper | December 10, 2020 |
| SIGNAL ANALYSIS &COMMUNICATION ECE355H1-F-LEC0101 | Liang | December 10, 2020 |

Enter email text here:
Hello, ____. Now is your opportunity to evaluate your courses and instructors. We appreciate you taking the time to provide comments to your instructor, as they can use this feedback to reflect upon their own work in these challenging times. Course evaluations are only open for a limited time, so complete yours today! Online evaluations are currently available for the following courses: COURSE INSTRUCTOR EVALUATION END FOUNDATIONS OF COMPUTING ECE358H1-F-LEC0101 Veneris December 10, 2020 INTRO TO MACHINE INTELLIGENCE ECE324H1-F-LEC0101 Rose December 10, 2020 MATRIX ALGEBRA & OPTIMIZATION ECE367H1-F-LEC0101 Draper December 10, 2020 SIGNAL ANALYSIS &COMMUNICATION ECE355H1-F-LEC0101 Liang December 10, 2020 EVALUATE MY COURSES ››
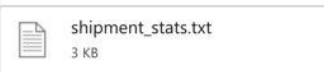
Model rnn predicts that this email is: Promotions and Events

# Examples of Incorrect Predictions 2

**Prediction: Promotions & Events**
**Ground Truth: Professional**



Assignment - shipment stats                                    📎 1 ∨

shipment_stats.txt
3 KB

**EXTERNAL EMAIL:**
Hi ▮▮▮▮

It was a pleasure talking to you today.

As promised during our chat, we would like to ask you to do the assignment. This will allow us to assess your technical skills in a "real-life" setting. We do not believe in asking candidates to code on the whiteboard under the pressure of the spotlight during the interview, we believe that candidates will produce the best code in a calm and relaxed environment (just like day to day in the office).  For this, we are asking you to spend some of your personal time and we appreciate that.

Please see the assignment attached. We would like to ask you to refrain from posting the assignment or solution on the internet.

When you are coding, do the research on SOLID (especially SOL), TDD, OO or functional design patterns and most importantly "Clean Code" and modularity. All of the items that I stated will contribute to making the software easy to change and understand.

We won't accept the assignment without passing tests. Most languages have a number of great unit testing frameworks.

Thanks and we hope you enjoy it. The assignment represents a snippet of one of the real problems that ▮▮▮▮ had to solve.

Enter email text here:
Hi ▮▮▮▮   It was a pleasure talking to you today.  As promised during our chat, we would like to ask you to do the assignment. This will allow us to assess your technical skills in a "real-life" setting. We do not believe in asking candidates to code on the whiteboard under the pressure of the spotlight during the interview, we believe that candidates will produce the best code in a calm and relaxed environment (just like day to day in the office).  For this, we are asking you to spend some of your personal time and we appreciate that.  Please see the assignment attached. We would like to ask you to refrain from posting the assignment or solution on the internet.  When you are coding, do the research on SOLID (especially SOL), TDD, OO or functional design patterns and most importantly "Clean Code" and modularity. All of the items that I stated will contribute to making the software easy to change and understand.  We won't accept the assignment without passing tests. Most languages have a number of great unit testing frameworks.  Thanks and we hope you enjoy it. The assignment represents a snippet of one of the real problems that had to solve.  Kind regards

Model rnn predicts that this email is: Promotions and Events

# Ignore this

E

**ECE 367H1 F LEC0101 on Piazza** <no-reply@piazza.com>
Thu 2020-12-03 4:00 PM
To: ██████████

**EXTERNAL EMAIL:**
Your classmate posted a new Followup.

Hello ████████

Great discussion on Linear Programming and its application! You mentioned that LP can be extended in to a variety of fields and I wonder if maximizing profit/minimizing cost etc. financial problems can also be addressed with LP. Since they also share a linear objective function which is to be optimized (maximize or minimize) subject to a certain number of constraints. Thanks!

Best,

████████

```
Enter email text here:
Hello, Great discussion on Linear Programming and its application! You mentioned that
LP can be extended in to a variety of fields and I wonder if maximizing profit/minimiz
ing cost etc. financial problems can also be addressed with LP. Since they also share
a linear objective function which is to be optimized (maximize or minimize) subject to
a certain number of constraints. Thanks!   Best,

Model rnn predicts that this email is: Academics
```

# Ignore this too

Some content in this message has been blocked because the sender isn't in your Safe senders list.
I trust content from ▓▓▓▓▓▓▓▓▓▓▓▓ | Show blocked content

CB

▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓
ms.com>
Mon 2020-11-30 10:39 AM
To:▓▓▓▓▓▓

**EXTERNAL EMAIL:**

Hello ▓▓▓ ,

Thank you for your interest in our ▓▓▓ Operations: WHS Intern 2021 United States ▓▓▓▓▓ !

There are many bright and talented candidates like you in the marketplace, making our job all that much harder. With that said, we regret to inform you that we will be unable to extend the interview process at this time. We appreciate the time you invested in pursuing an opportunity with us.

While we know that this news is disappointing, we encourage you to apply to other roles on our careers website that look of interest to you and match your skillset. We acknowledge that your skills and interests are constantly developing and changing, and we would like to let you know that you are eligible to apply for this role again after 6 months.

We wish you all the best in your search and in your future endeavors!

Thank you for your interest,

▓▓▓▓ Recruiting Team

Enter email text here:
Thank you for your interest in our ▓▓▓ Operations: WHS Intern 2021 United States (1
▓▓▓▓▓ )!    There are many bright and talented candidates like you in the marketplace,
making our job all that much harder. With that said, we regret to inform you that we w
ill be unable to extend the interview process at this time. We appreciate the time you
invested in pursuing an opportunity with us.    While we know that this news is disapp
ointing, we encourage you to apply to other roles on our careers website that look of
interest to you and match your skillset. We acknowledge that your skills and interests
are constantly developing and changing, and we would like to let you know that you are
eligible to apply for this role again after 6 months.    We wish you all the best in yo
ur search and in your future endeavors!    Thank you for your interest, ▓▓▓▓ Recrui
ting Team

Model rnn predicts that this email is: Professional

# Discussion

We have <u>high test accuracy</u>... but:

- Data source (¾ Engineering + ¼ Artsci = 100% UofT)
- Only tested on emails in our accounts
- **May not generalize!**

<u>Keyword frequency</u> implies category:

Recall the most common words in each category

- "Course" → Academics
- "Order" → Alerts
- "Application" → Professional
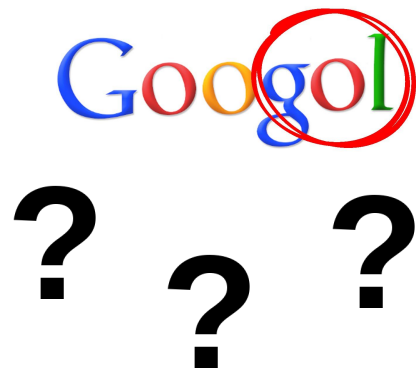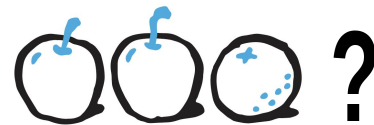- "Students" → Promotion & Events
- "Practice" → Personal

# What we learned

- Preparing raw data takes **time**!
  - Need enough
  - Manual cleaning
- Need to be **VERY clear** on the scope of the labels

# Things to Consider

- **Generalization** for all users?
- User choice for **labels**?
- Managing **typos**?
- Incorporation of **decision heuristics** into NN design?

# Any Questions?