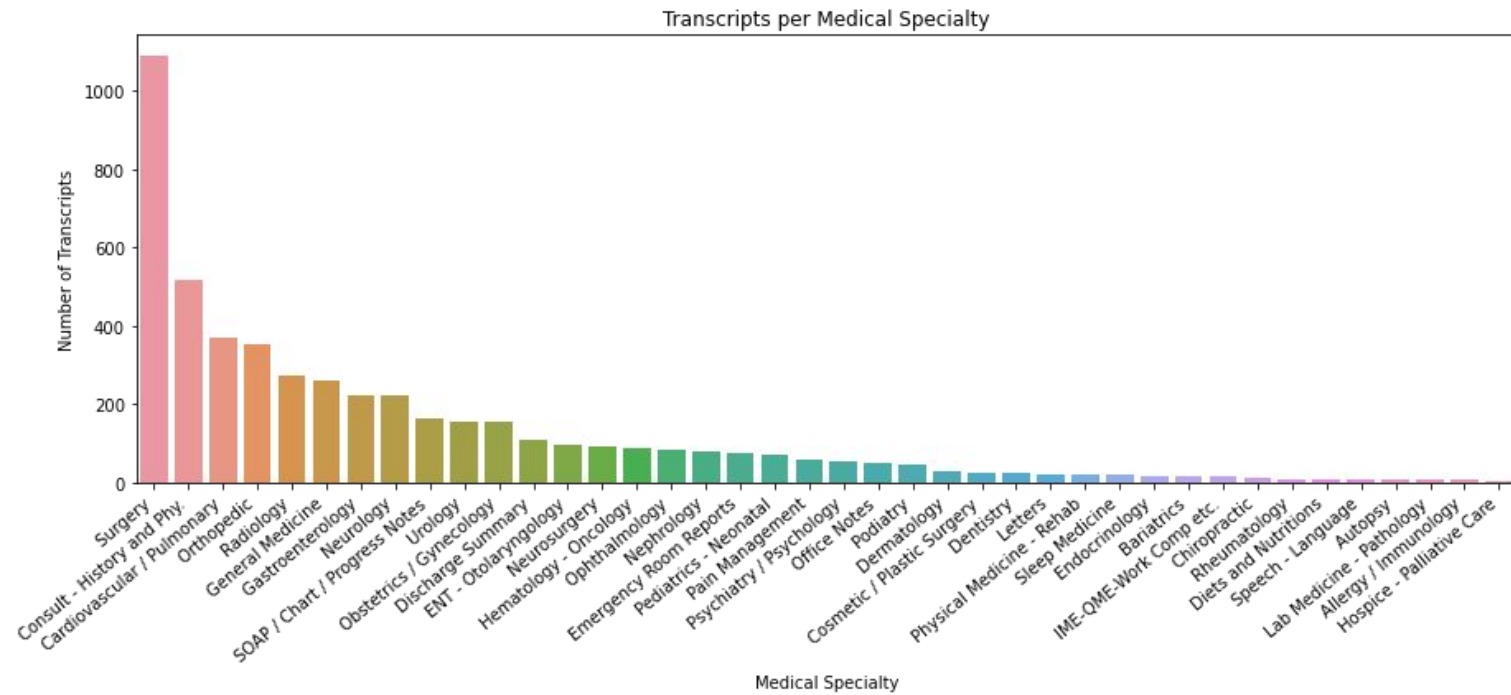# Medical Transcription Classification

Freya Gray
CS39AA
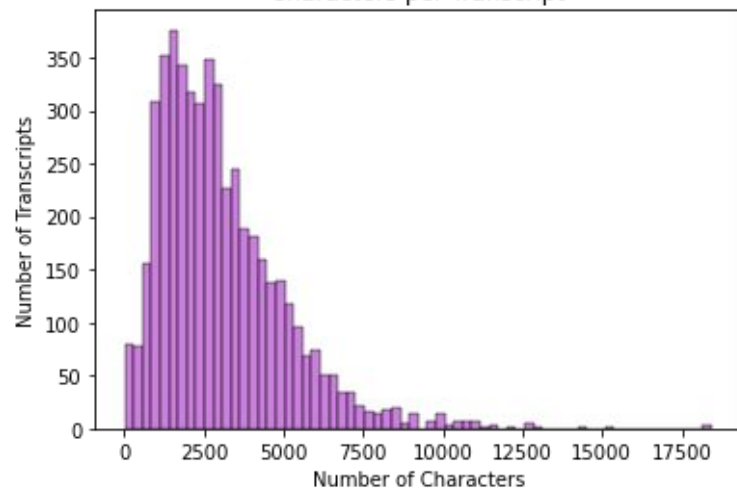
# Dataset

- Medical transcripts from mtsamples.com
- Transcripts from medical professionals after seeing a patient
- All names and dates have been changed or removed
- 40 different medical specialties
- 4966 transcripts
- Transcripts come in a wide range of formats and styles

Transcripts per Medical Specialty

## Characters per Transcript

Number of Transcripts — Number of Characters

## Words per Transcript

Number of Transcripts — Number of Words

## Sentences per Transcript

Number of Transcripts — Number of Sentences

## Average Word Length per Transcript

Number of Transcripts — Average Word Length

Most Common Words

Most Common Bigrams

Most Common Trigrams

# Baseline Model

- Reduced medical specialties to 18
- Support Vector Machine
- GridSearchCV to find best parameters
  - Best parameters: C = 1, kernel = linear
- Accuracy 23%

Confusion matrix (True rows × Predicted columns):

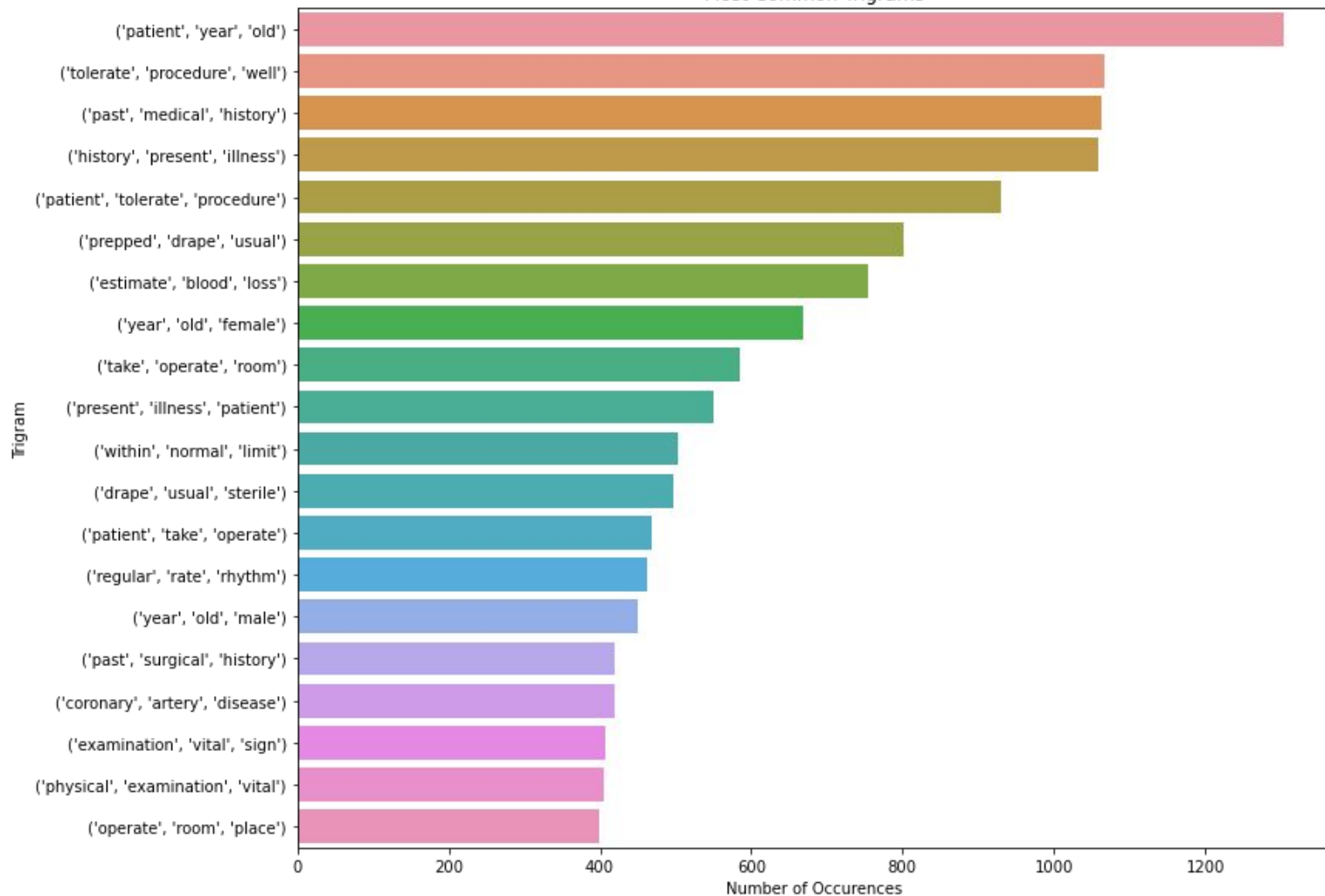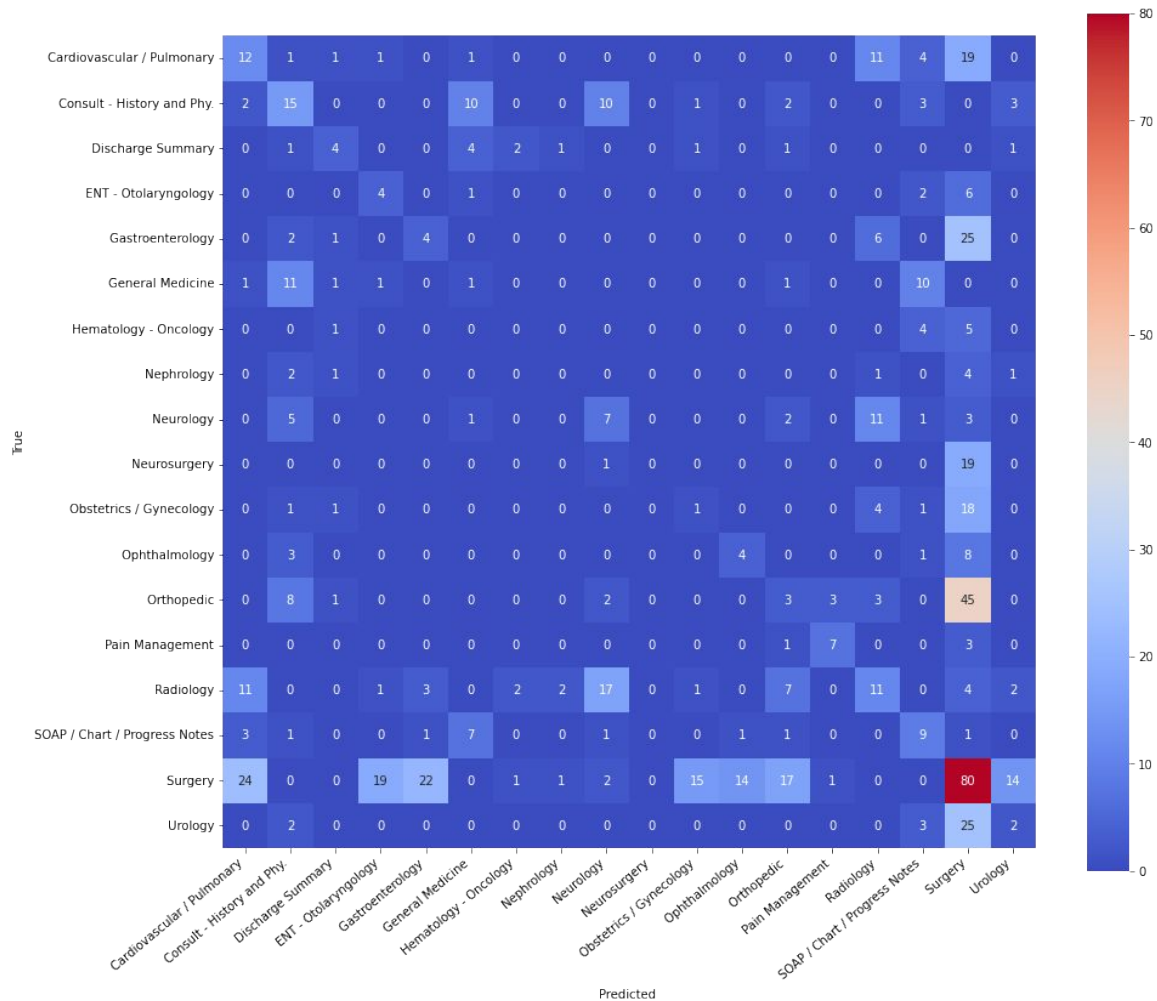| True \ Predicted | Cardiovascular / Pulmonary | Consult - History and Phy. | Discharge Summary | ENT - Otolaryngology | Gastroenterology | General Medicine | Hematology - Oncology | Nephrology | Neurology | Neurosurgery | Obstetrics / Gynecology | Ophthalmology | Orthopedic | Pain Management | Radiology | SOAP / Chart / Progress Notes | Surgery | Urology |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Cardiovascular / Pulmonary | 12 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 11 | 4 | 19 | 0 |
| Consult - History and Phy. | 2 | 15 | 0 | 0 | 0 | 10 | 0 | 0 | 10 | 0 | 1 | 0 | 2 | 0 | 0 | 3 | 0 | 3 |
| Discharge Summary | 0 | 1 | 4 | 0 | 0 | 4 | 2 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 |
| ENT - Otolaryngology | 0 | 0 | 0 | 4 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 6 | 0 |
| Gastroenterology | 0 | 2 | 1 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 6 | 0 | 25 | 0 |
| General Medicine | 1 | 11 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 10 | 0 | 0 |
| Hematology - Oncology | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 5 | 0 |
| Nephrology | 0 | 2 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 4 | 1 |
| Neurology | 0 | 5 | 0 | 0 | 0 | 1 | 0 | 0 | 7 | 0 | 0 | 0 | 2 | 0 | 11 | 1 | 3 | 0 |
| Neurosurgery | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 19 | 0 |
| Obstetrics / Gynecology | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 4 | 1 | 18 | 0 |
| Ophthalmology | 0 | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 0 | 0 | 0 | 1 | 8 | 0 |
| Orthopedic | 0 | 8 | 1 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 3 | 3 | 3 | 0 | 45 | 0 |
| Pain Management | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 7 | 0 | 0 | 3 | 0 |
| Radiology | 11 | 0 | 0 | 1 | 3 | 0 | 2 | 2 | 17 | 0 | 1 | 0 | 7 | 0 | 11 | 0 | 4 | 2 |
| SOAP / Chart / Progress Notes | 3 | 1 | 0 | 0 | 1 | 7 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 9 | 1 | 0 |
| Surgery | 24 | 0 | 0 | 19 | 22 | 0 | 1 | 1 | 2 | 0 | 15 | 14 | 17 | 1 | 0 | 0 | 80 | 14 |
| Urology | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 25 | 2 |

Predicted

True

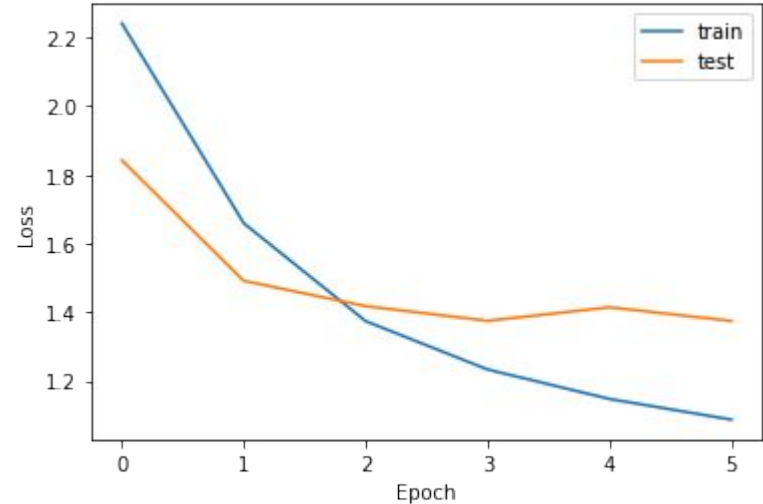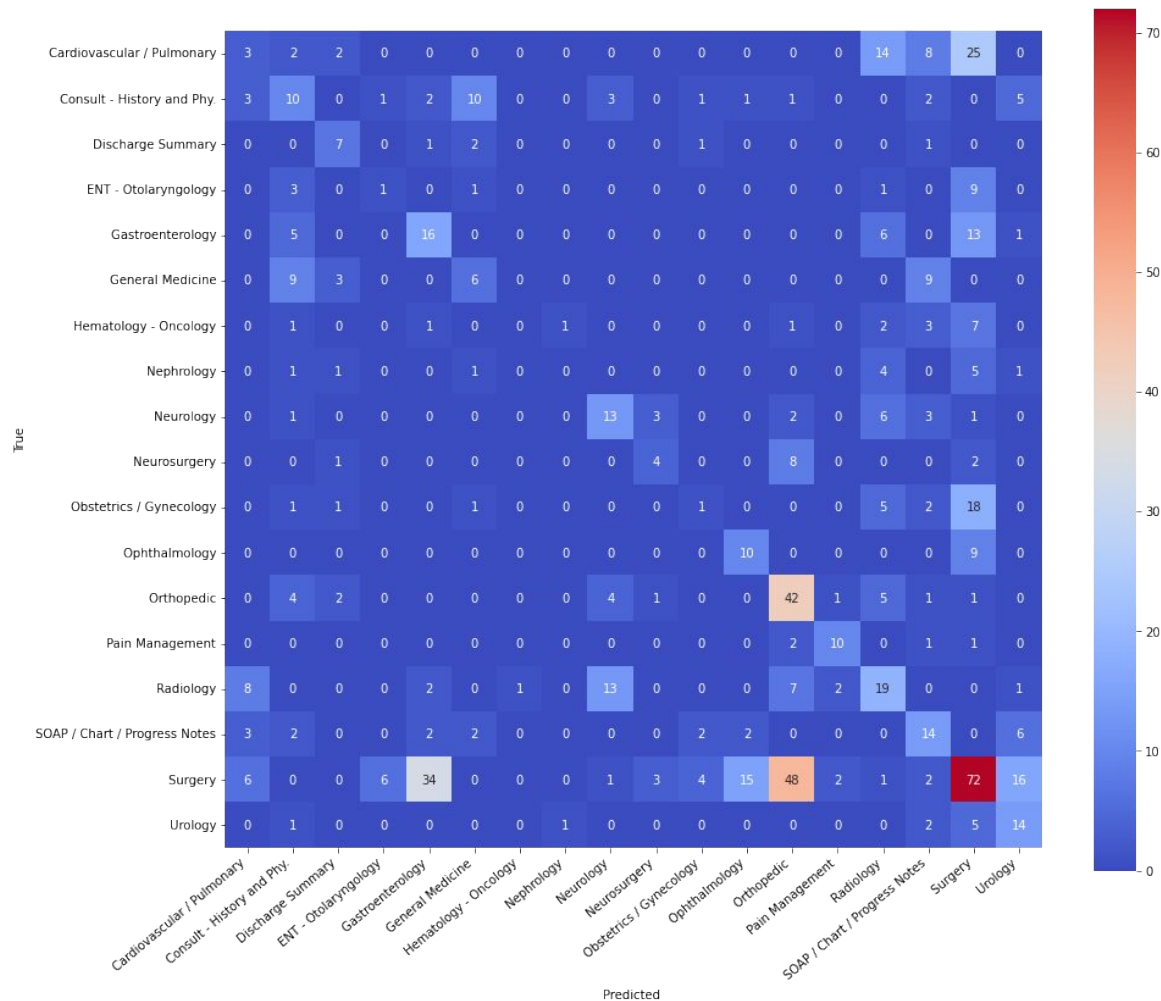Accuracy:0.2332859174964438
Precision:0.20922610493758473
Recall:0.2332859174964438
F1 score:0.21742526566778728

# RoBERTa

- Reduced medical specialties to 18
- Model limited max sequence length
- Just transcriptions
- Accuracy: 48.06%
- Loss:1.0876
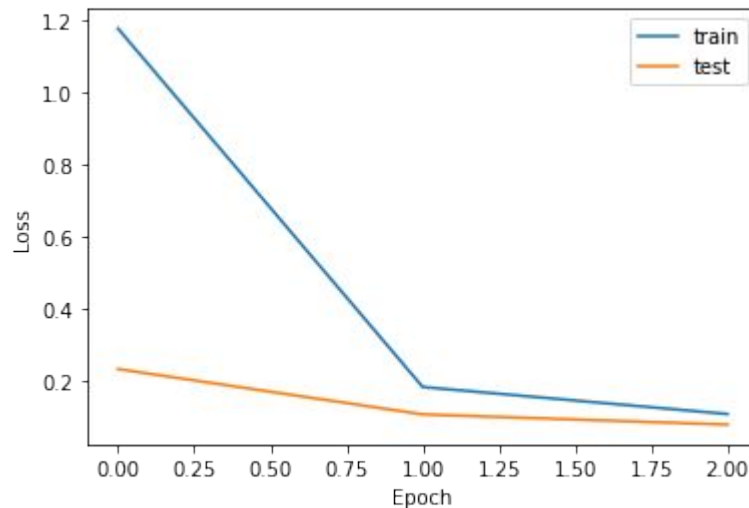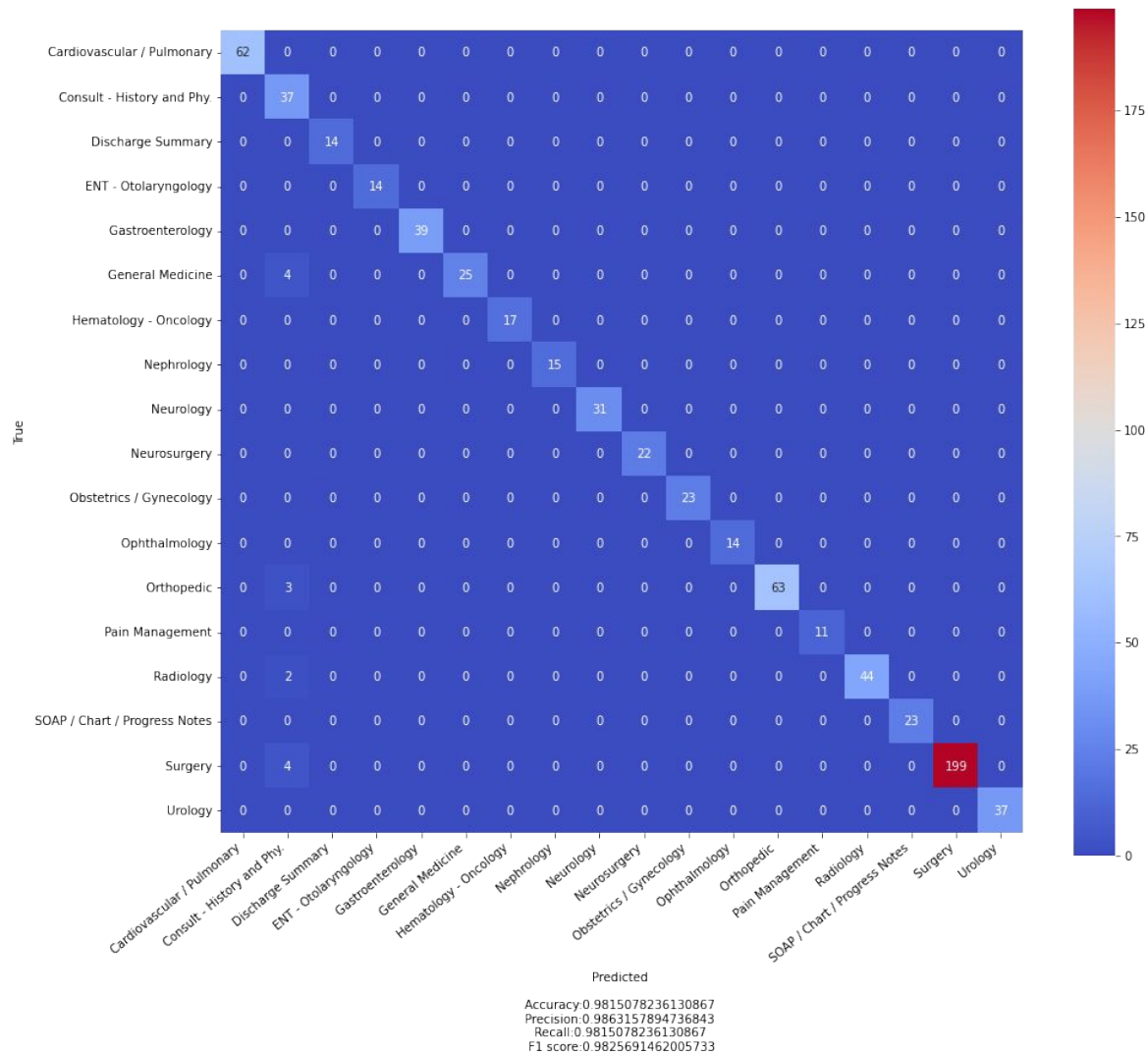- Validation Accuracy: 34.42%
- Validation Loss: 1.3746

Confusion matrix (True vs Predicted):

| True \ Predicted | Cardiovascular / Pulmonary | Consult - History and Phy. | Discharge Summary | ENT - Otolaryngology | Gastroenterology | General Medicine | Hematology - Oncology | Nephrology | Neurology | Neurosurgery | Obstetrics / Gynecology | Ophthalmology | Orthopedic | Pain Management | Radiology | SOAP / Chart / Progress Notes | Surgery | Urology |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Cardiovascular / Pulmonary | 3 | 2 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 14 | 8 | 25 | 0 |
| Consult - History and Phy. | 3 | 10 | 0 | 1 | 2 | 10 | 0 | 0 | 3 | 0 | 1 | 1 | 1 | 0 | 0 | 2 | 0 | 5 |
| Discharge Summary | 0 | 0 | 7 | 0 | 1 | 2 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| ENT - Otolaryngology | 0 | 3 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 9 | 0 |
| Gastroenterology | 0 | 5 | 0 | 0 | 16 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 6 | 0 | 13 | 1 |
| General Medicine | 0 | 9 | 3 | 0 | 0 | 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 9 | 0 | 0 |
| Hematology - Oncology | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 2 | 3 | 7 | 0 |
| Nephrology | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 0 | 5 | 1 |
| Neurology | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 13 | 3 | 0 | 0 | 2 | 0 | 6 | 3 | 1 | 0 |
| Neurosurgery | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 0 | 0 | 8 | 0 | 0 | 0 | 2 | 0 |
| Obstetrics / Gynecology | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 5 | 2 | 18 | 0 |
| Ophthalmology | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 10 | 0 | 0 | 0 | 0 | 9 | 0 |
| Orthopedic | 0 | 4 | 2 | 0 | 0 | 0 | 0 | 0 | 4 | 1 | 0 | 0 | 42 | 1 | 5 | 1 | 1 | 0 |
| Pain Management | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 10 | 0 | 1 | 1 | 0 |
| Radiology | 8 | 0 | 0 | 0 | 2 | 0 | 1 | 0 | 13 | 0 | 0 | 0 | 7 | 2 | 19 | 0 | 0 | 1 |
| SOAP / Chart / Progress Notes | 3 | 2 | 0 | 0 | 2 | 2 | 0 | 0 | 0 | 0 | 2 | 2 | 0 | 0 | 0 | 14 | 0 | 6 |
| Surgery | 6 | 0 | 0 | 6 | 34 | 0 | 0 | 0 | 1 | 3 | 4 | 15 | 48 | 2 | 1 | 2 | 72 | 16 |
| Urology | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 5 | 14 |

Accuracy:0.3442389758179232
Precision:0.31842947125638177
Recall:0.3442389758179232
F1 score:0.3199623527339034

# RoBERTa

- Keywords
- Accuracy: 98.29%
- Loss: .1126
- Validation Accuracy: 98.58%
- Validation Loss: .0667

Accuracy:0.9815078236130867
Precision:0.9863157894736843
Recall:0.9815078236130867
F1 score:0.9825691462005733

# Conclusions

- Pre-trained model outperformed the baseline model
- Overlap between specialties seems to have increased errors
- Possible Improvements
  - Different Model
  - Better data cleaning
  - Classify transcripts differently

# Sources

- [Dataset](#)
- [Hugging Face Guide](#)
- [Text Cleaning](#)
- Data Exploration: [link 1](#), [link 2](#)