# Disease Identification From Literature
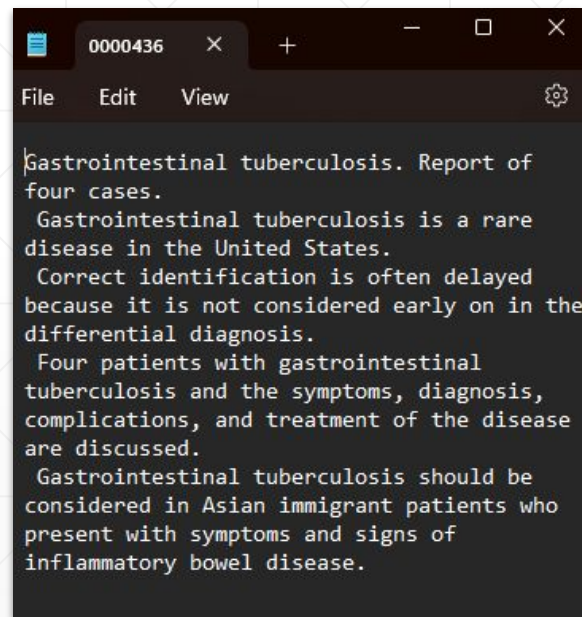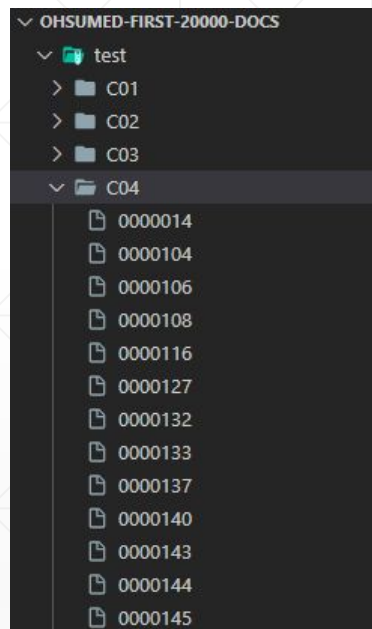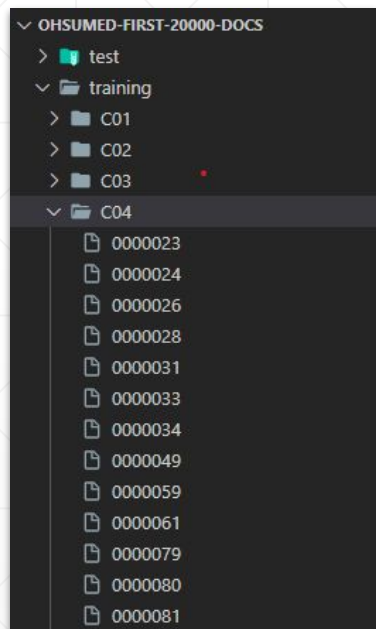
## Project Proposal

**Presented by:**

- **Asif Shahriar # 1805040**
- **Debojit Pandit # 1805042**

# Project Idea

- There are 23 categories of diseases: disease-categories

- Task: Given a medical literature/document, attribute it to one or more of these categories

- NLP & Multi-label classification problem

- Kaggle link

- Papers:

  - GCN-paper-1

  - GCN-paper-2

# Dataset

- **OHSUMED** [ohsumed-dataset](#) : consists of abstracts from medical journals

# Data Preprocessing

- Current representation: document 0008272 is contained in two directories: C07 & C15 -> two labels

- Need to create a csv file , with entries like this:

  - 0008272    0000001000000010000000

# Proposed Experiments

- We found papers using GCN models on Ohsumed dataset

- GCN-paper-1

- GCN-paper-2

**Table 4.** Test Accuracy (%) on text classification datasets. The numbers are averaged over 10 runs.

| Dataset | Model | Test Acc. ↑ | Time (seconds) ↓ |
|---------|-------|-------------|------------------|
| 20NG    | GCN   | $87.9 \pm 0.2$ | $1205.1 \pm 144.5$ |
|         | SGC   | $88.5 \pm 0.1$ | $19.06 \pm 0.15$ |
| R8      | GCN   | $97.0 \pm 0.2$ | $129.6 \pm 9.9$ |
|         | SGC   | $97.2 \pm 0.1$ | $1.90 \pm 0.03$ |
| R52     | GCN   | $93.8 \pm 0.2$ | $245.0 \pm 13.0$ |
|         | SGC   | $94.0 \pm 0.2$ | $3.01 \pm 0.01$ |
| Ohsumed | GCN   | $68.2 \pm 0.4$ | $252.4 \pm 14.7$ |
|         | SGC   | $68.5 \pm 0.3$ | $3.02 \pm 0.02$ |
| MR      | GCN   | $76.3 \pm 0.3$ | $16.1 \pm 0.4$ |
|         | SGC   | $75.9 \pm 0.3$ | $4.00 \pm 0.04$ |

| Model | 20NG | MR | Ohsumed |
|-------|------|-----|---------|
| CNN | $0.8215 \pm 0.0052$ | $0.7775 \pm 0.0007$ | $0.5844 \pm 0.0106$ |
| LSTM | $0.7318 \pm 0.0185$ | $0.7768 \pm 0.0086$ | $0.4927 \pm 0.0107$ |
| Graph-CNN | $0.8142 \pm 0.0032$ | $0.7722 \pm 0.0027$ | $0.6386 \pm 0.0053$ |
| Fast-GCN | **OOM** | $0.7510 \pm 0.0021$ | $0.5441 \pm 0.0081$ |
| Text-GCN | $0.8634 \pm 0.0009$ | $0.7674 \pm 0.0020$ | $0.6836 \pm 0.0056$ |
| Text-GNN | - | - | $0.6940 \pm 0.0060$ |
| WGCN | $\mathbf{0.8885} \pm 0.0012$ | $\mathbf{0.7794} \pm 0.0010$ | $\mathbf{0.6962} \pm 0.0024$ |

- We will use models like **LSTM-CNN** and **BERT** for better performance