**Name: Bilal Butt**

**Intern ID: DHC-356**

## Task 2: Multi Label Emotion Recognition

a. **Dataset Preprocessing Steps**

- Loaded GoEmotions using HuggingFace datasets.
- Tokenized text with BERT tokenizer.
- Converted label lists to 28-length multi-hot vectors.
- Removed variable-length arrays to avoid batching errors.

b. **Model Selection and Rationale**

- Selected **BERT base uncased** due to its proven performance in NLP tasks.
- Fine-tuned on multi-label classification using BertForSequenceClassification.
- Used problem_type="multi_label_classification" to handle multiple emotions per sentence.

c. **Challenges Faced and Solutions**

- **fsspec '**' error:** solved by upgrading datasets and fsspec.
- **RuntimeError:** tensor stack mismatch: resolved using a custom DataCollator.
- **Labels mismatch:** resolved using manual one-hot encoding per label list.

d. **Results with Visualizations and Interpretations**

- Achieved reasonable performance on the test set using a small data subset.
- **Micro-F1 ≈ 0.56**
- **Hamming Loss ≈ 0.21**
- The classification report shows that the model could detect multiple emotions per sentence, especially on common emotion labels like "joy", "sadness", or "anger".
- The label distribution chart highlights imbalanced classes which can be improved with full dataset training.

# Emotion distribution bar chart