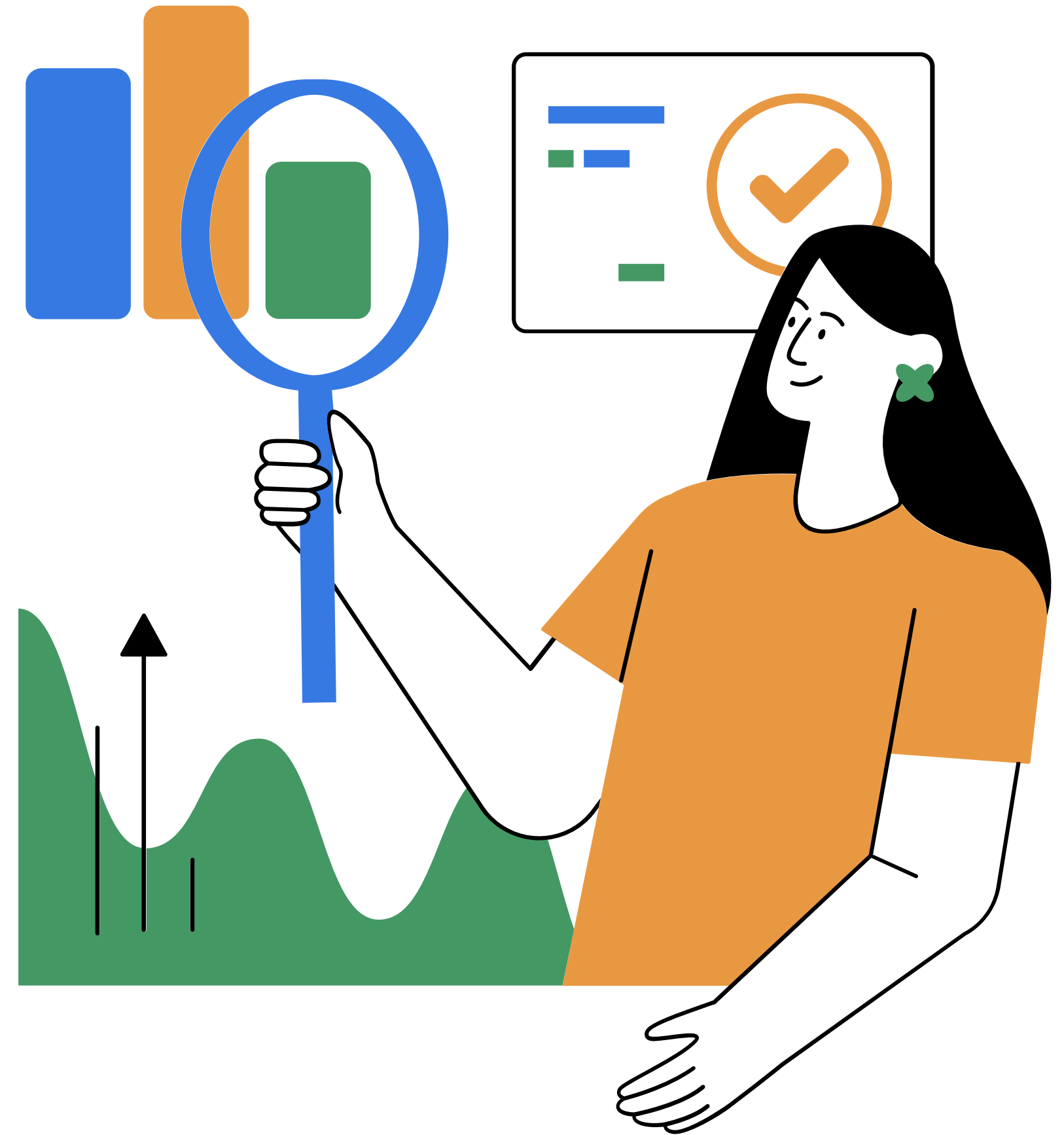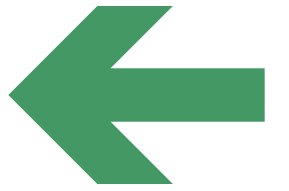# Leveraging Emotion Classification for Enhanced Content Intelligence

Tudor Pitulice
Noah Ivanisevic
Zakariae El Moumni
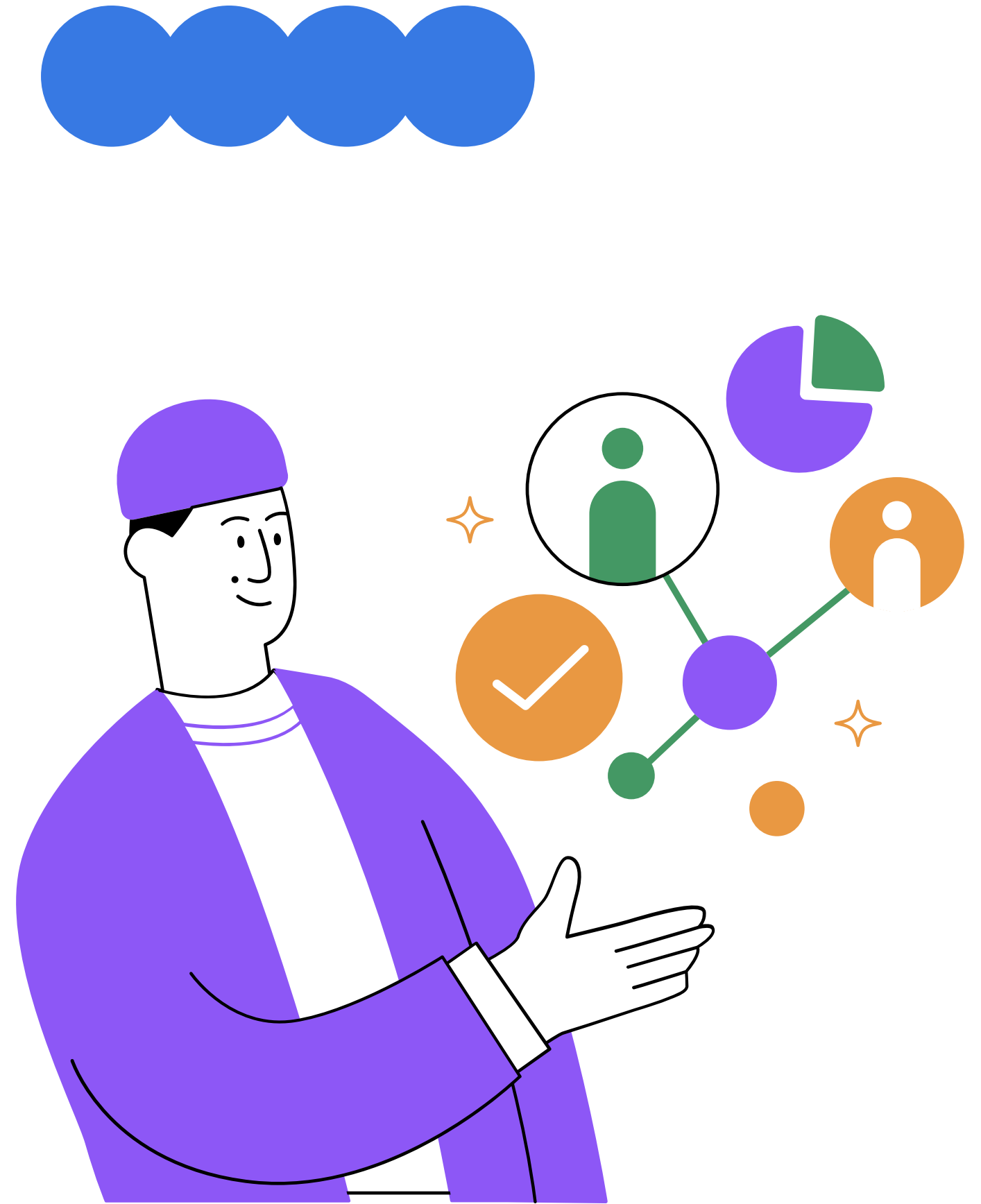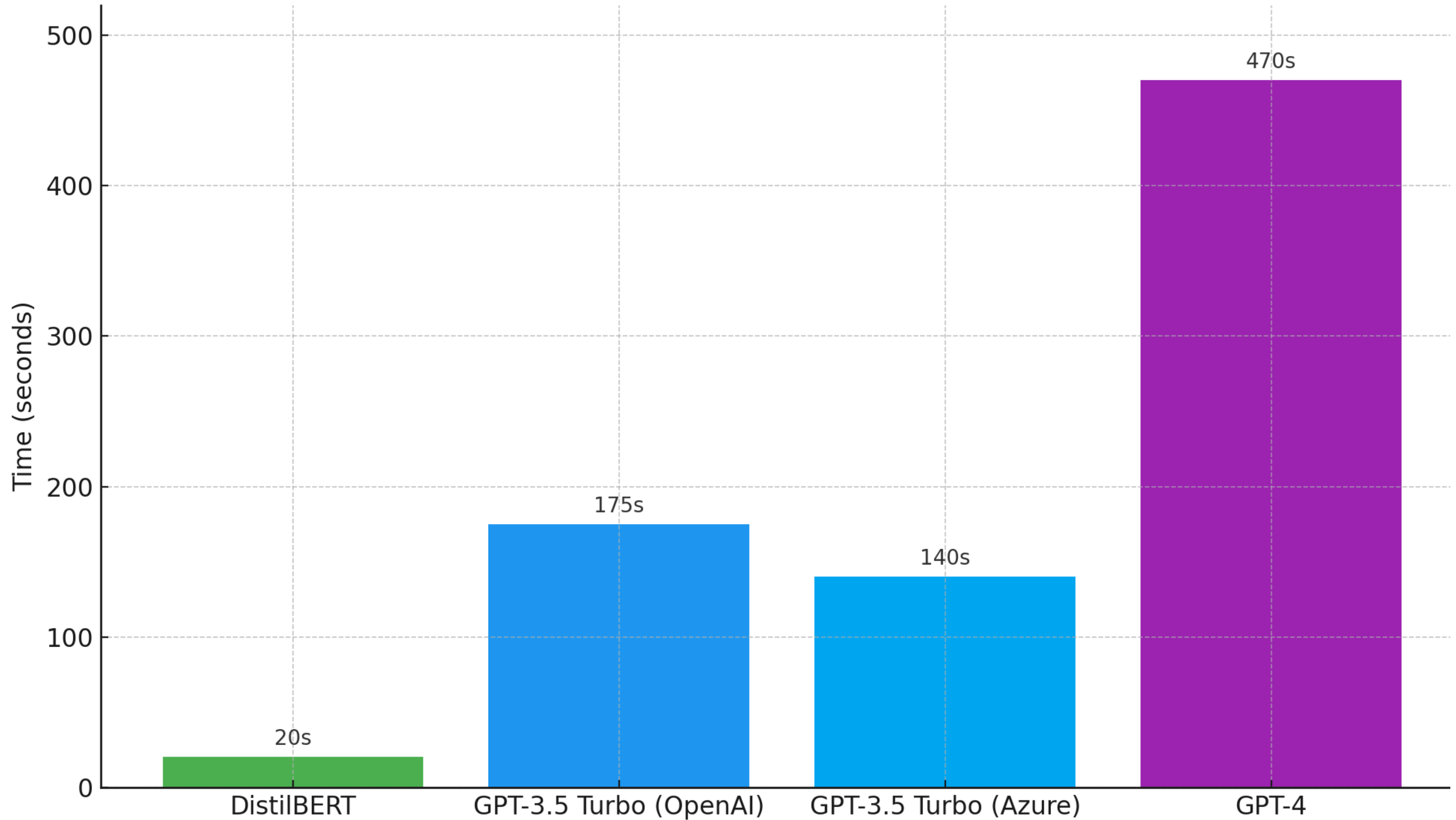
# Topic



- ***Topic:*** *"The DNA of outstanding content"*

- ***Domain:*** Media and content analytics tool

- ***Use case:***
  - Predicting the emotion from an audio file for further analysis

- ***Benefits***
  - Privacy concerns
  - Faster predictions for the users
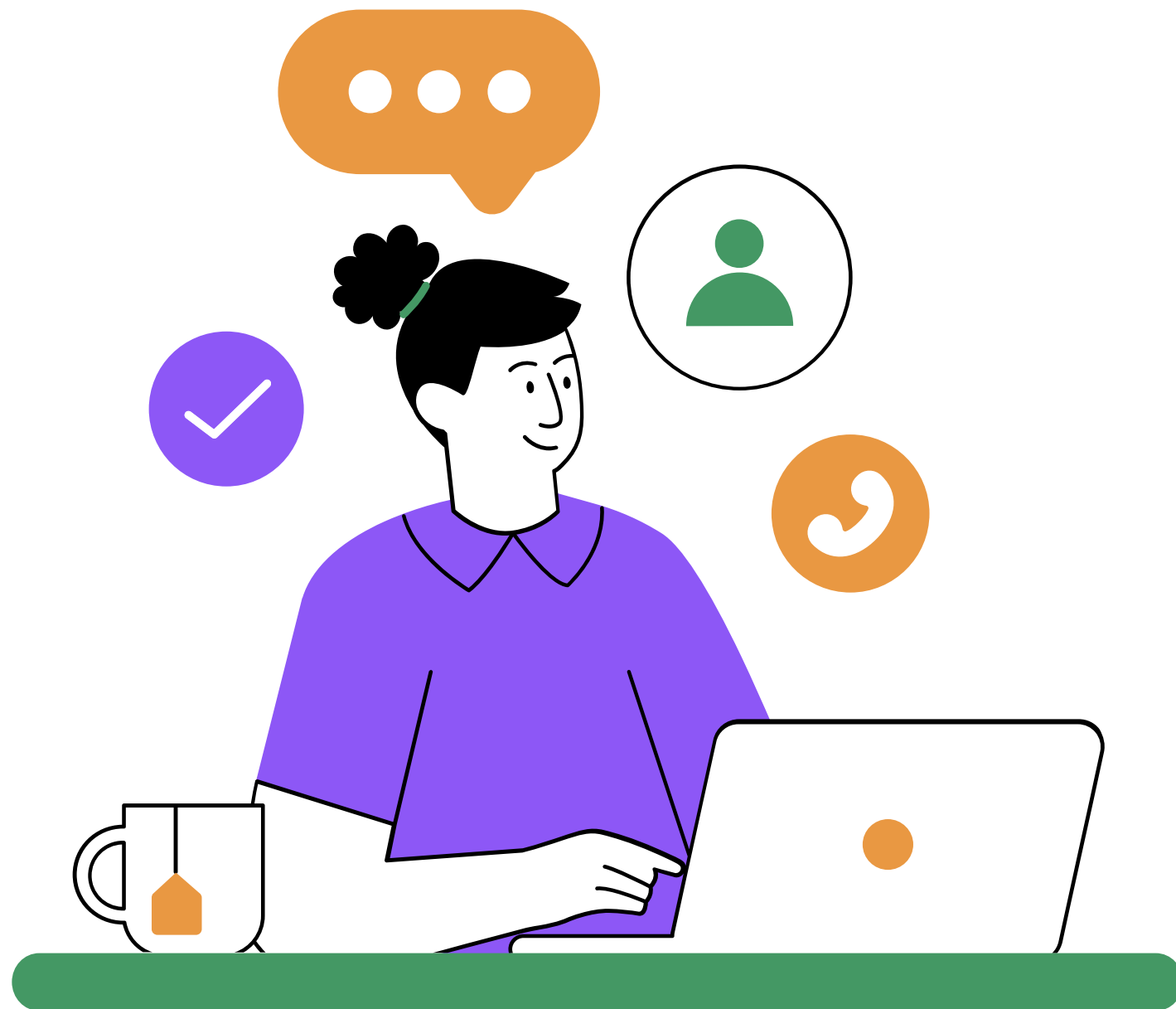
# Business Value

- Greater security during unstable times

- Lower operational cost compared to commercial API's

- Faster prediction times

Prediction Time for 5,000 Sentences by Model

# Data Characteristics

**GoEmotions dataset:**
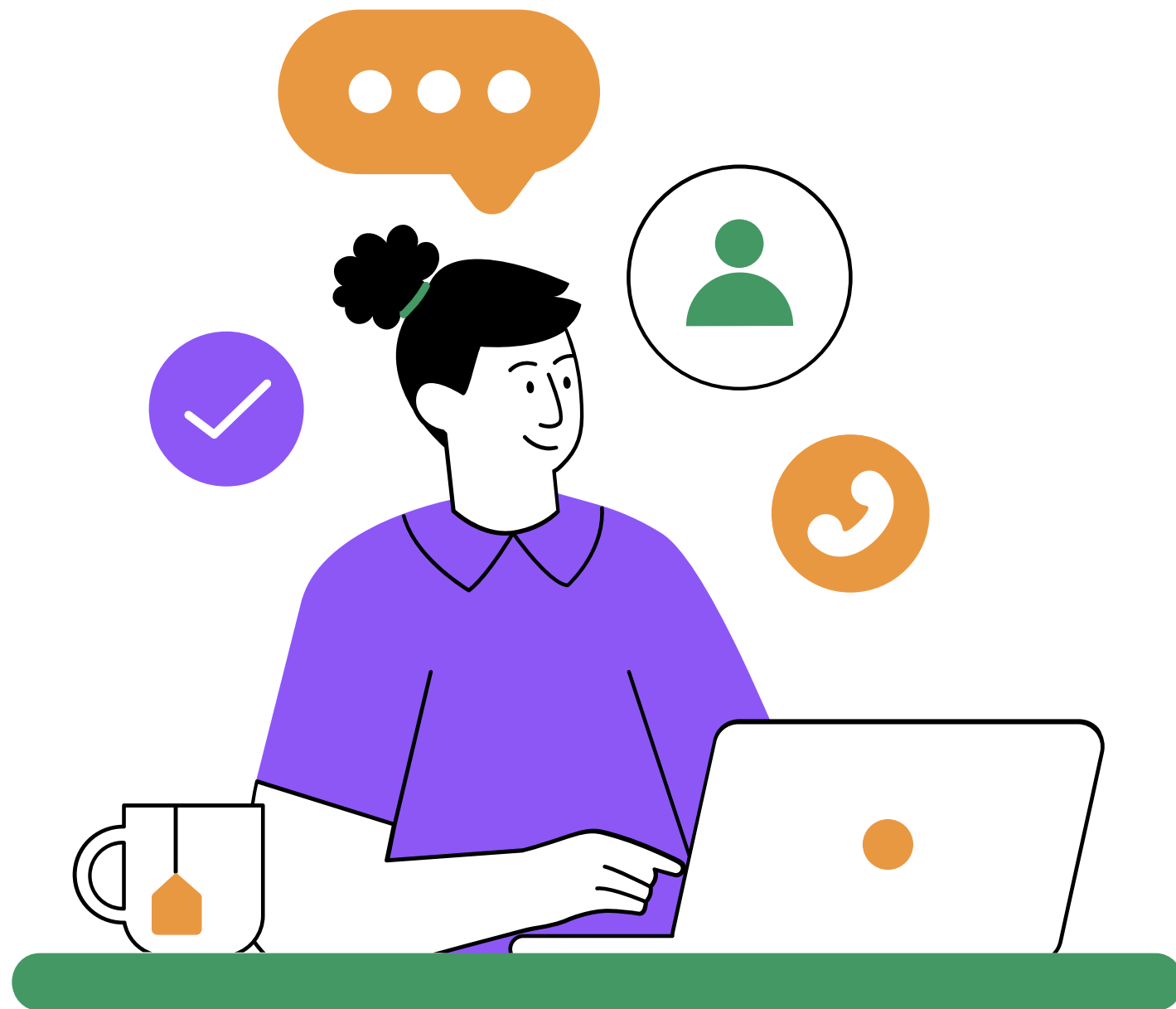- 58.000 sentences
- Reddit comments

**Apple ICloud IMessage dataset:**
- 400.000 sentences
- Messages from Apple users

**Test sets from other groups:**
- 1.800 sentences (2 test sets)
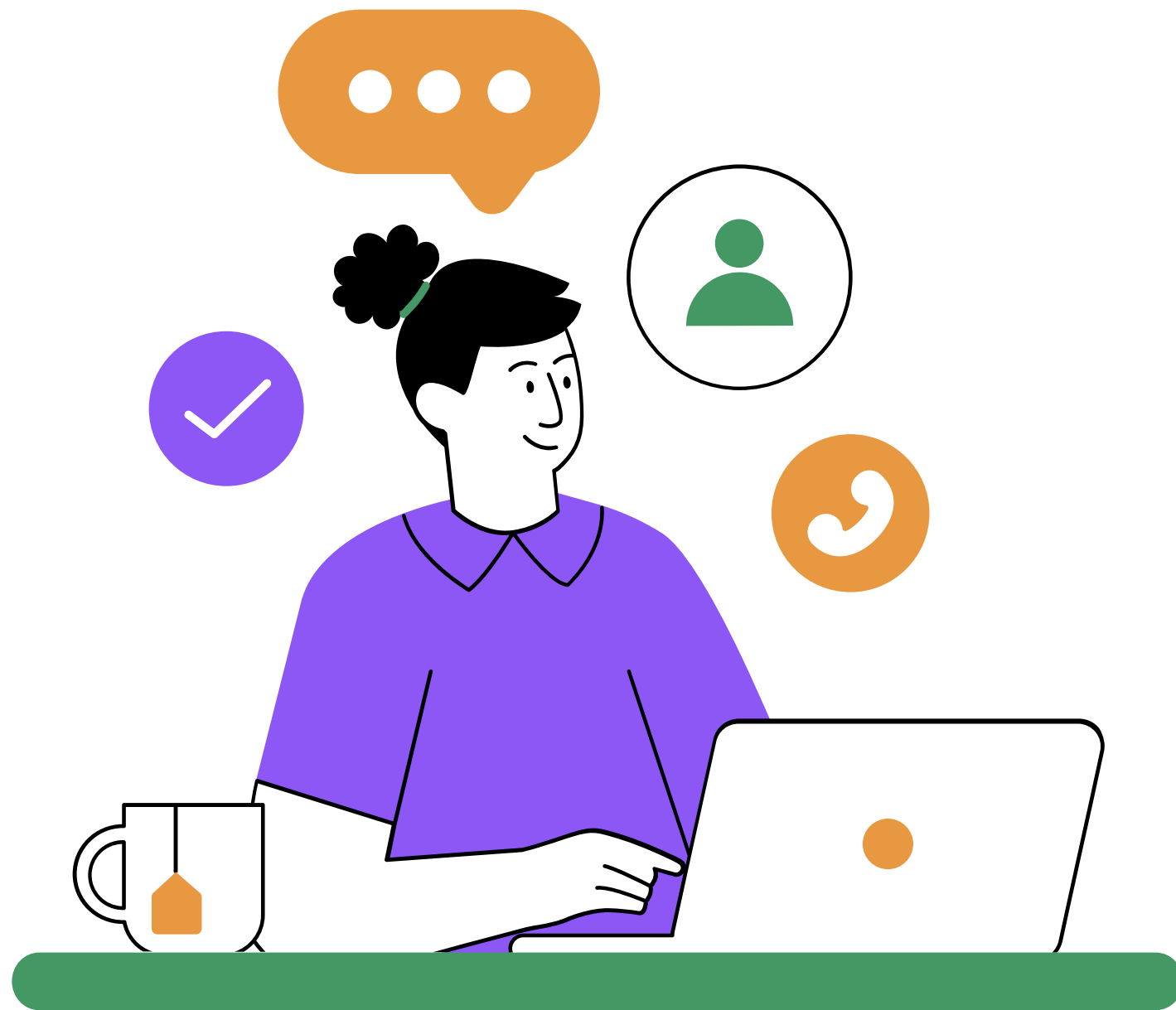- contextual diversity
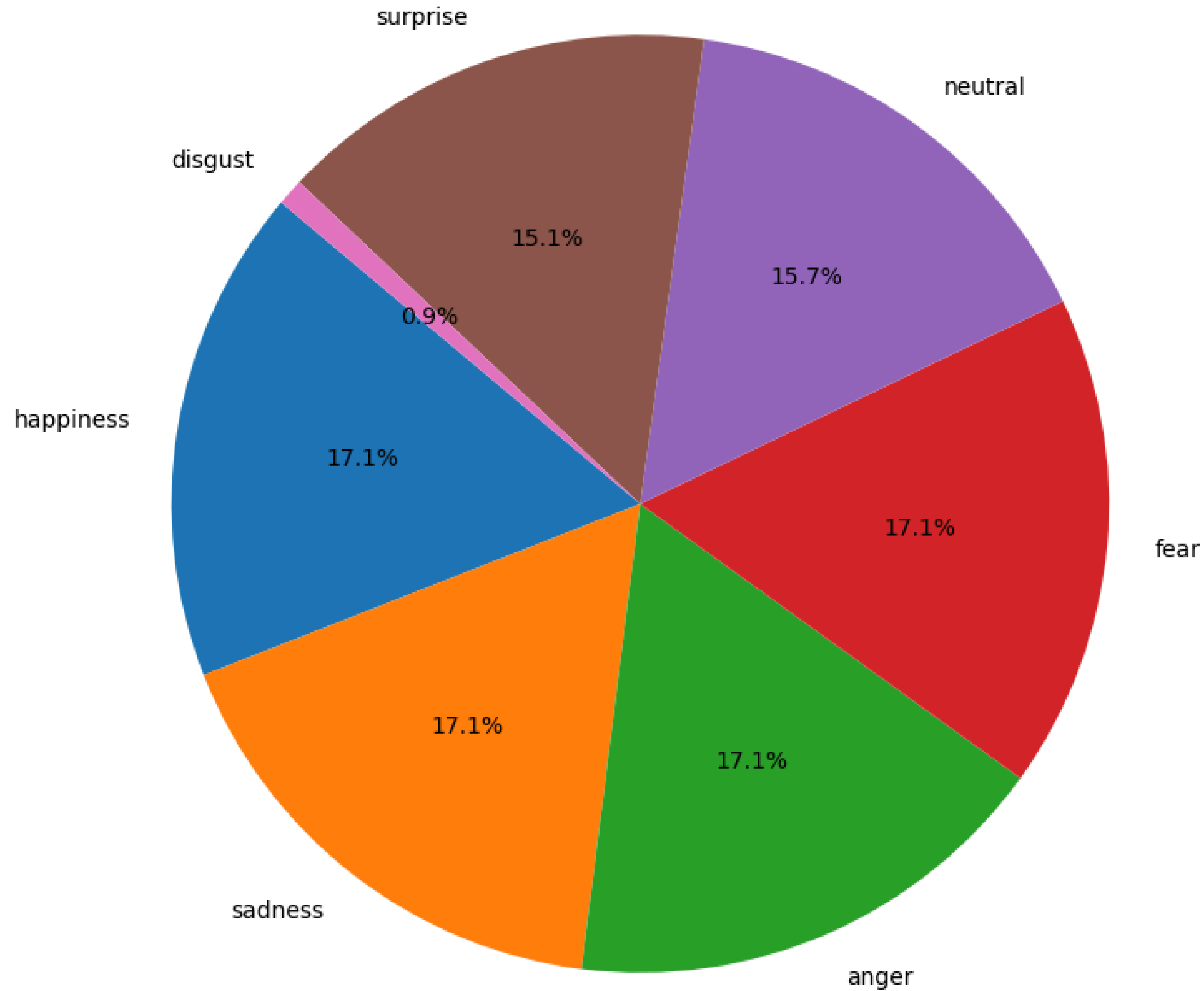
# Data Characteristics

**Preprocessing steps:**

- mapping of the emotions into the 6 emotions + neutral
- handling missing values
- balancing the data distribution
- data augmentation targeting minority classes
  - synonym replacement
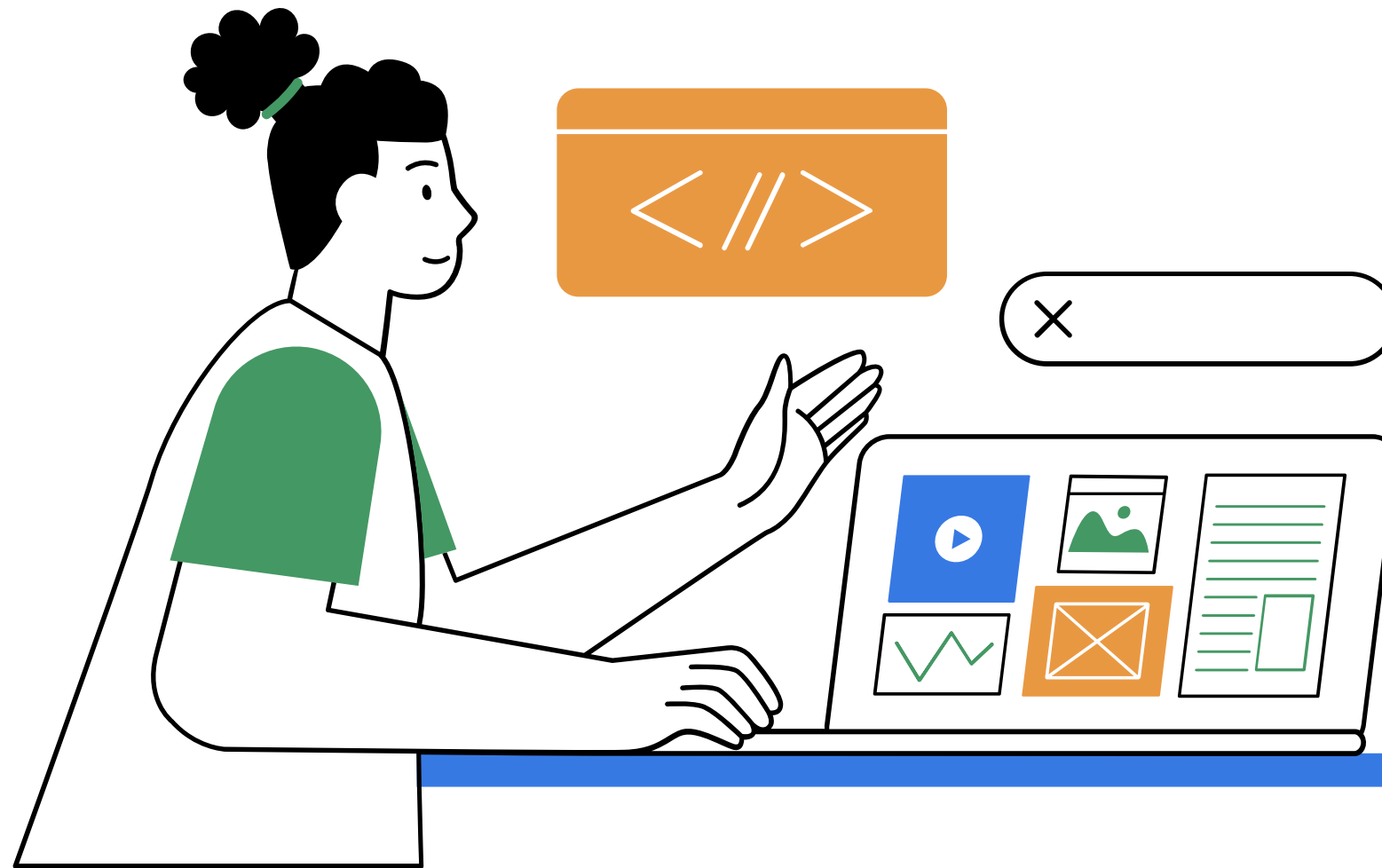  - random word swap
  - random deletion

# Data Limitations

- Contextual gap between training and testing sets

- Underrepresentation of some emotions
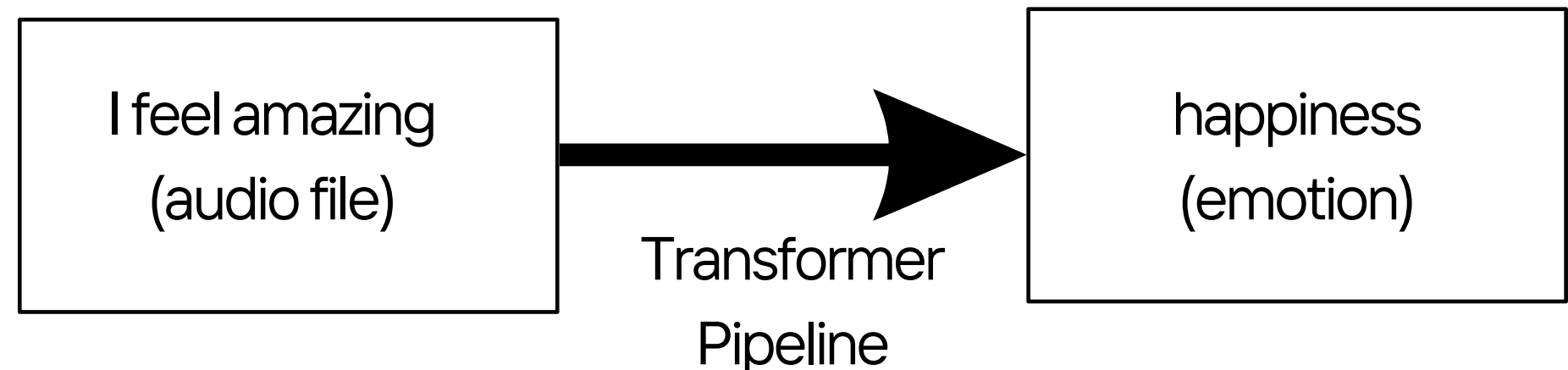
Emotion Distribution in Training Set

**Models:**

- **DistilBERT (transformer)**
  - generally good
  - requires smaller training sets
  - faster version of BERT
  - retains 94% of the performance of BERT
- **DistilROBERTA (transformer)**
  - better on informal speech
  - requires more data
  - slower than DistilBERT
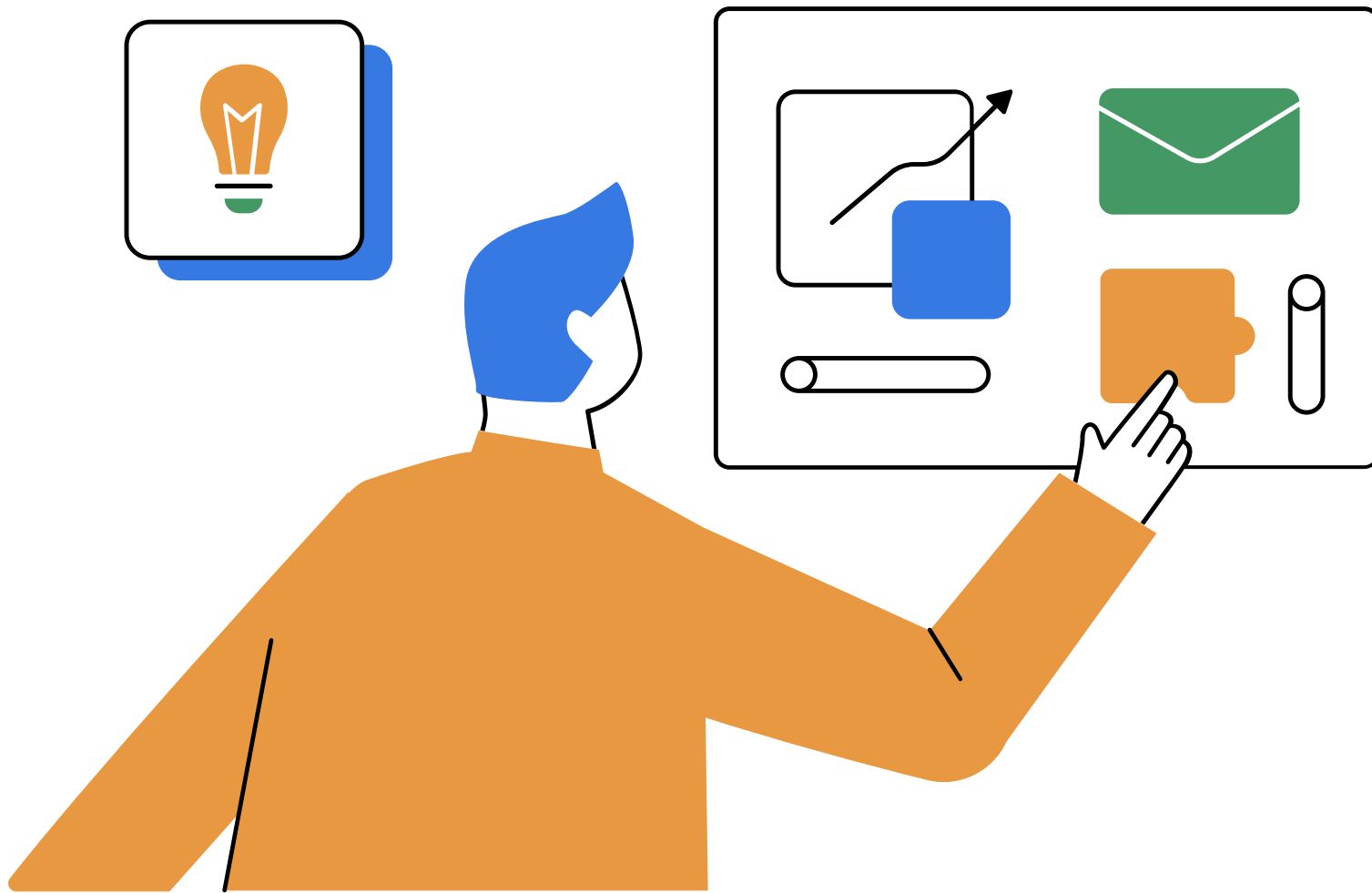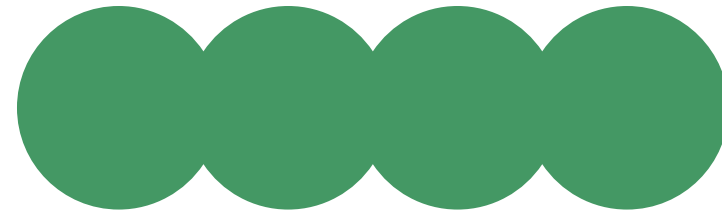  - better performance than DistilBERT

# Models

| I feel amazing (audio file) | → Transformer Pipeline | happiness (emotion) |

# Models Performance

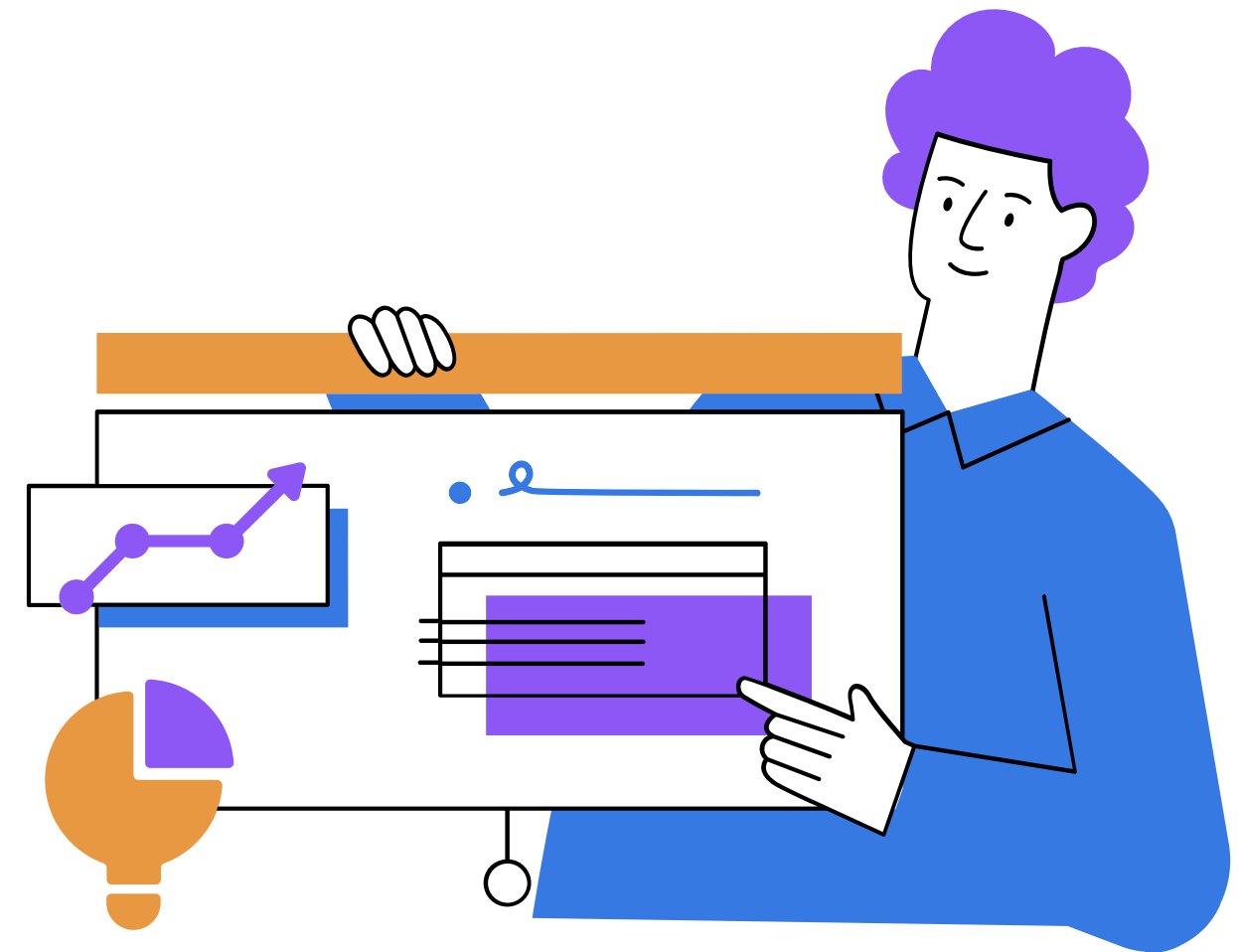|  | DistilBERT | DistilROBERTA |
|---|---|---|
| Accuracy | 0.73 | 0.48 |
| F1 score | 0.65 | 0.49 |
| Precision | 0.71 | 0.58 |
| Recall | 0.73 | 0.48 |

# Ethical Considerations



- Different accents from different races
- Data Bias that might include
  - racism
  - xenophobia
  - sexism
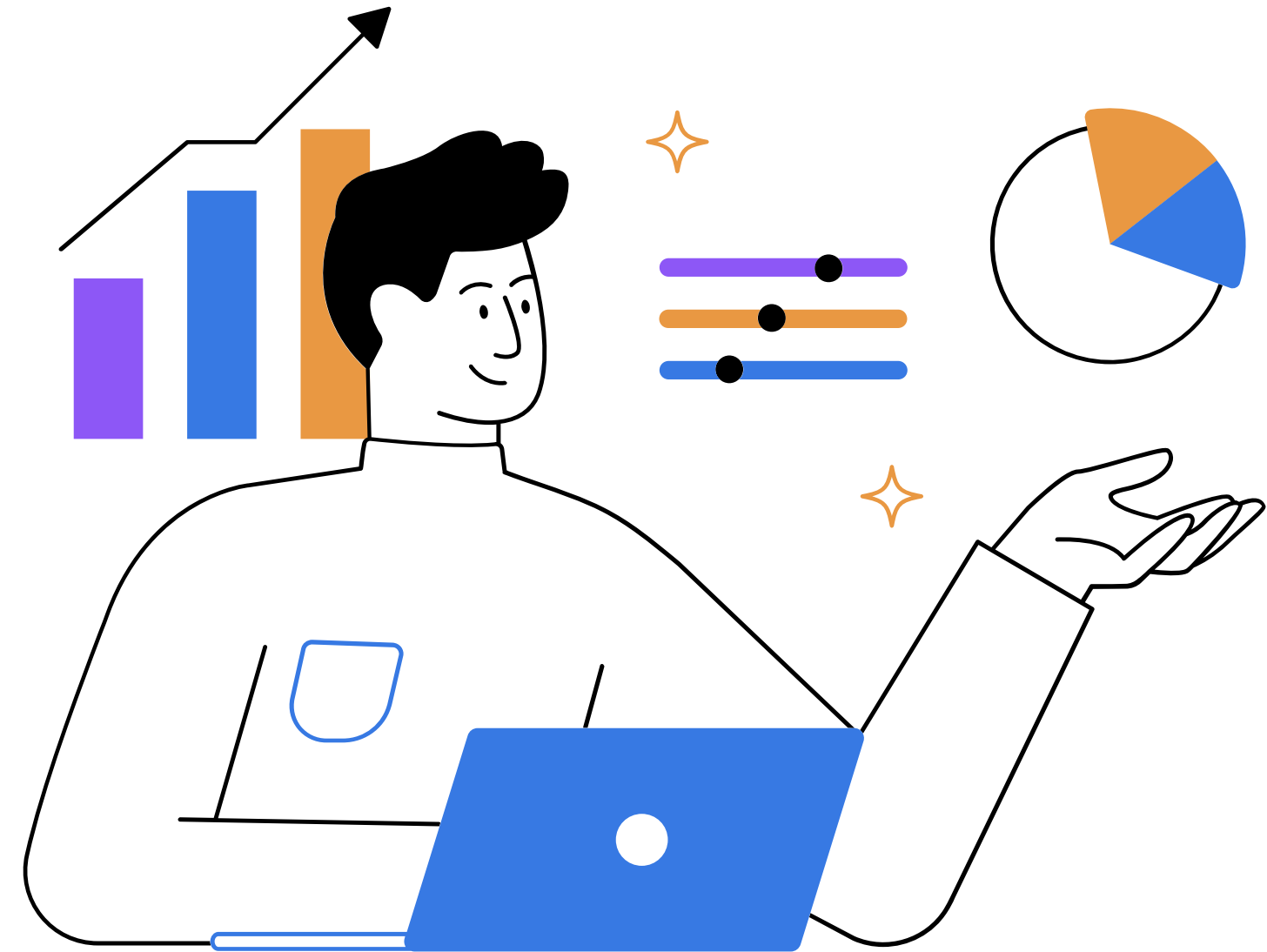  - demographics
  - Generational slang

# Limitations for Implementation

- Availability of data
- Quality of data
- Domain adaptation needed for media-specific texts
- Transcription issues
- Multilingual support not included
- Maximum file size limit: 5 GB

# Next Steps

- Gathering more data
- Model optimization (reach ChatGPT API performance)
- Deployment in production

Thank You