# Project Elimi-'Hate'

Chua Soon Poh
DSI-32

# Table of contents

Problem

# What is Hate Speech?

"**any kind of communication** in speech, writing or behaviour, that **attacks** or uses **pejorative** or **discriminatory** language with reference to a person or a group on the basis of **who they are**, in other words, based on their religion, ethnicity, nationality, race, colour, descent, gender or other identity factor."

UNITED NATIONS

"

# Social media provides a global megaphone for hate."

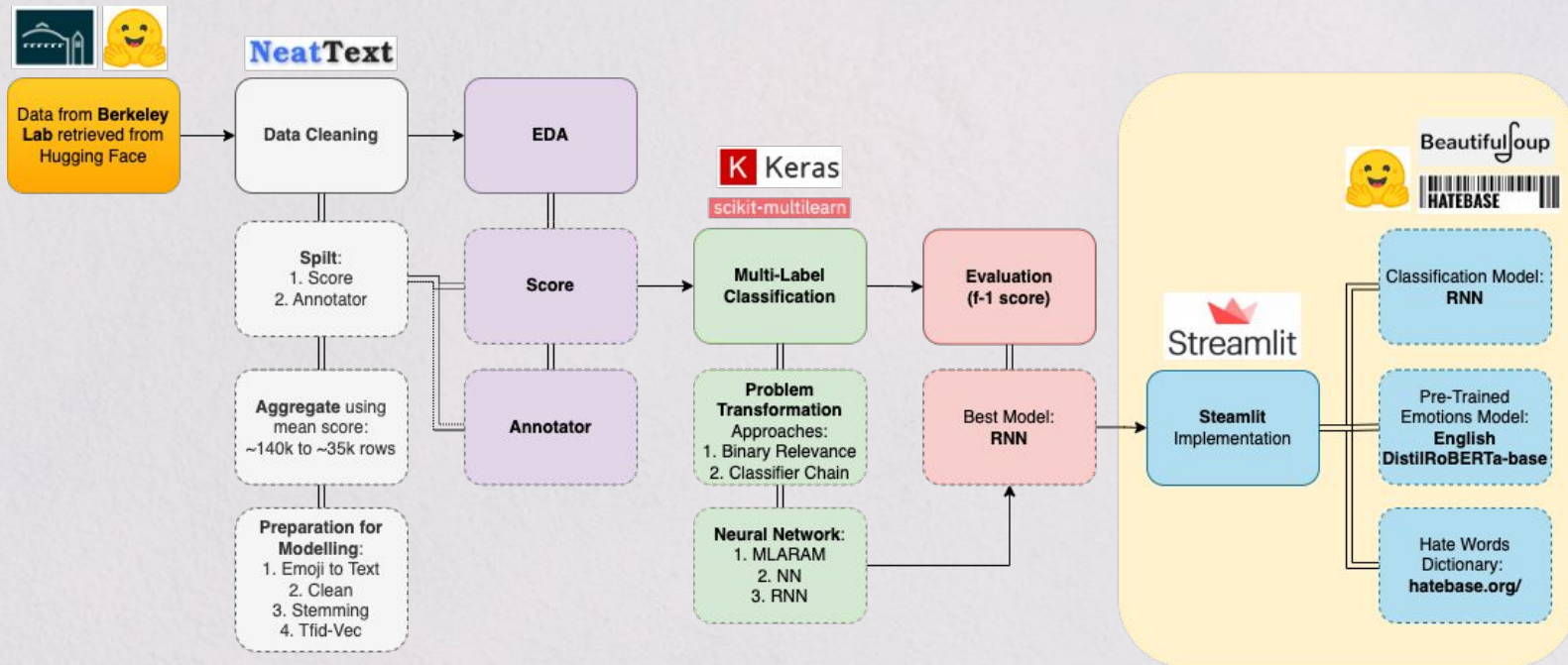**ANTÓNIO GUTERRES**, *United Nations Secretary-General, 2021*
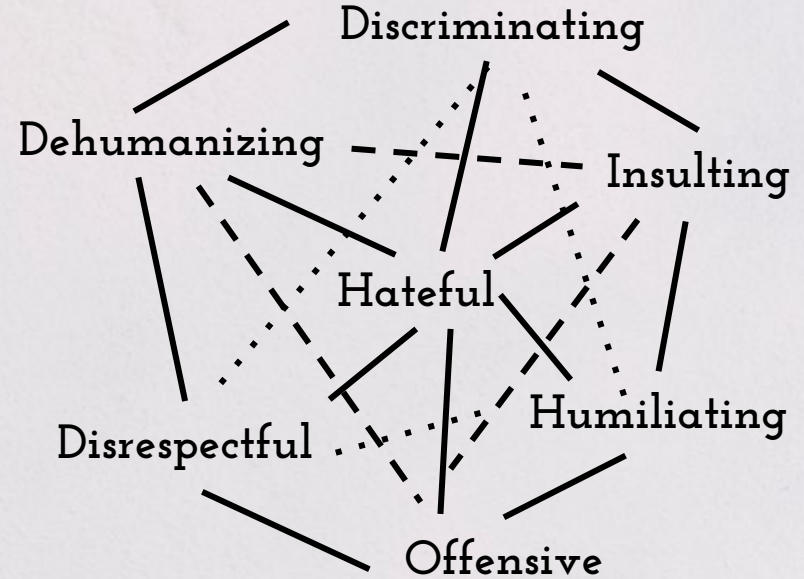
# Project Overview



**Dataset**

1. Increase empathy towards others in the social media space
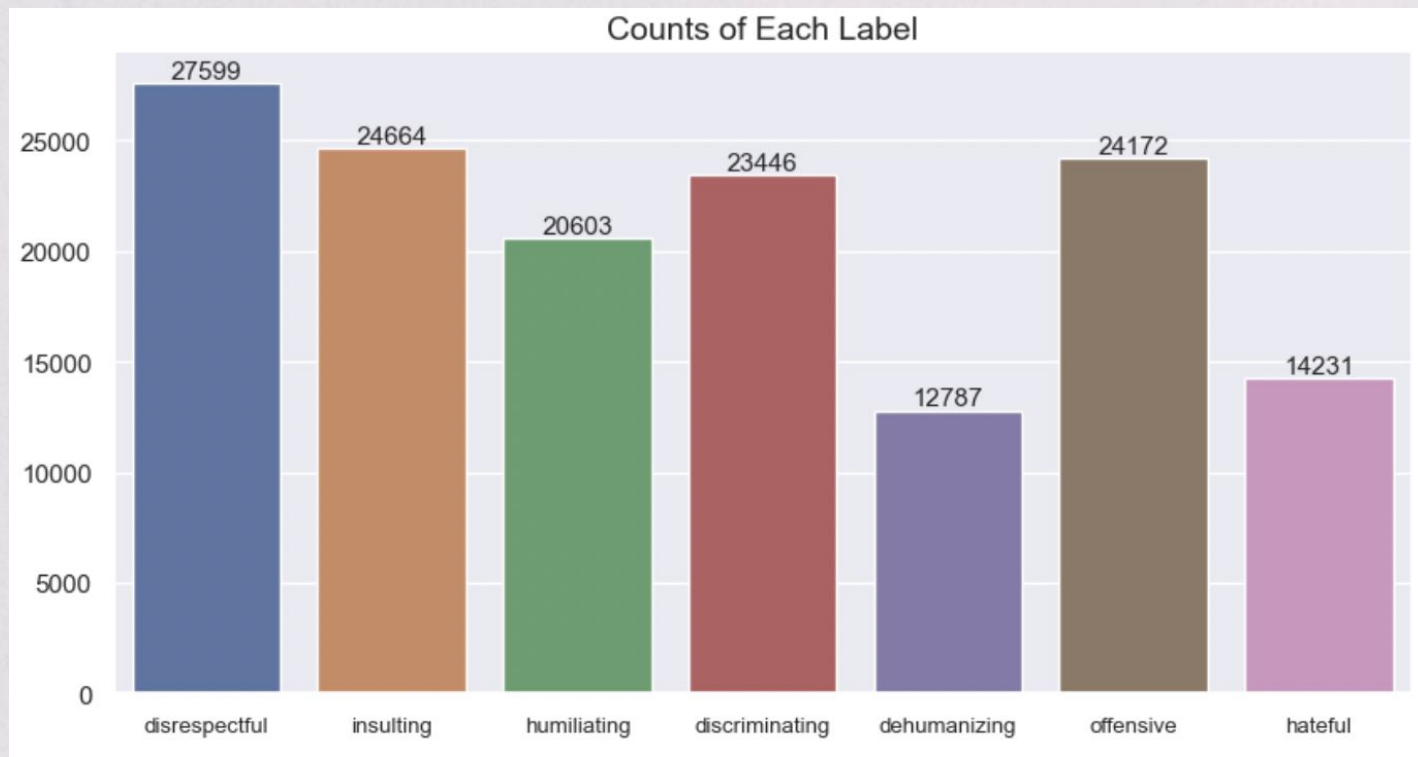2. Cultivate better awareness of unhelpful written expressions

# Multi-Labels?

**Dataset**

Disrespectful

↓

Insulting

↓

Offensive

↓

Discriminating

↓

Humiliating

↓

Dehumanizing

↓

Hateful

Discriminating

Dehumanizing

Insulting

Hateful

Disrespectful

Humiliating

Offensive

# Multi-Labels

Counts of Each Label

# Multi-Labels

Correlation between Labels

# Results

**Modelling**

| Model | Precision Macro Avg | Recall Macro Avg | F-1 Macro Avg |
|---|---|---|---|
| Binary Relevance & Naives Bayes | 0.73 | 0.74 | 0.70 |
| Classifier Chain & Naives Bayes | 0.62 | 0.93 | 0.74 |
| MLARAM | 0.62 | 0.51 | 0.53 |
| LSTM NN with GLOVE Embedding | <u>0.71</u> | <u>0.82</u> | <u>0.76</u> |

**Best Model:** LSTM NN with GLOVE Embedding

Modelling

|                  | precision | recall | f1-score |
|------------------|-----------|--------|----------|
| disrespectful    | 0.82      | 0.93   | 0.87     |
| insulting        | 0.77      | 0.92   | 0.84     |
| offensive        | 0.74      | 0.89   | 0.81     |
| discriminating   | 0.71      | 0.87   | 0.78     |
| humiliating      | 0.72      | 0.84   | 0.78     |
| hateful          | 0.66      | 0.72   | 0.69     |
| dehumanizing     | 0.58      | 0.55   | 0.57     |
|                  |           |        |          |
| micro avg        | 0.73      | 0.85   | 0.79     |
| macro avg        | 0.71      | 0.82   | 0.76     |

# Demo Part 1

**Streamlit API**

https://tinyurl.com/Elimi-Hate-Demo



Write your post below to check if:

1. **The negative emotion that your post may contain;**

2. **Your post is potentially disrespectful, insulting, offensive, discriminating, humiliating, hateful or dehumanizing towards others;**

3. **Your post contains any hate words from https://hatebase.org.**

🟣 j-hartmann
**/emotion-english-distilroberta-base**

Best Model: LSTM NN with GLOVE Embedding

HATEBASE

**Streamlit API**

Jonathan

DSI-32

Coon

General

Rube

Salad

Wink

Assembly

Gooks

Pain

Gimps

Lubra

Ho

Snowflake

Libtard

Pikey

Suntan

Thirty-Six

Yoga

Deep

**Demo Part 2**

Streamlit API

Jonathan

DSI-32

Coon

General

Rube

Salad

Wink

Assembly

Gooks

Pain

Gimps

Lubra

Ho

Snowflake

Libtard

Pikey

Suntan

Thirty-Six

Yoga

Deep

# Results from HateBase



rube

English A rural person.

Class

libtard

English Contraction of liberal and retard

Disability

wink                                                          Updates

English A Caucasian who emulates, or associates with, Chinese people. NON-HATEFUL MEANING A deliberate blink of one eye

Ethnicity

Targeted Groups

Ethnicities Chinese

Streamlit API

# Design Thinking

**Streamlit API**

**Step 1: Run Pre-Trained Emotion Model**

Anger 🤬
Disgust 🤢
Fear 😨
Sadness 😭

Surprise 😲
Joy 😃
Neutral 😐

**Step 2: Run Best Model**

**Step 3: Loop Hatebase Dictionary**

|  | precision | recall |
|---|---|---|
| disrespectful | 0.82 | 0.93 |
| insulting | 0.77 | 0.92 |
| offensive | 0.74 | 0.89 |
| discriminating | 0.71 | 0.87 |
| humiliating | 0.72 | 0.84 |

1. Increase empathy towards others in the social media space
2. Cultivate better awareness of unhelpful written expressions

# Limitations

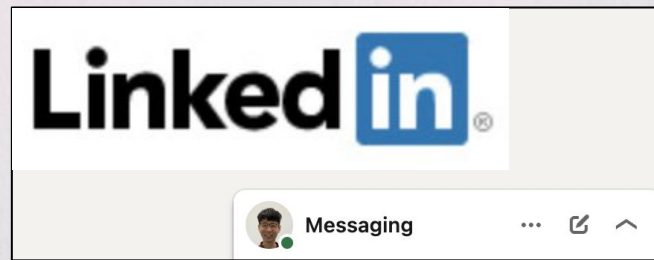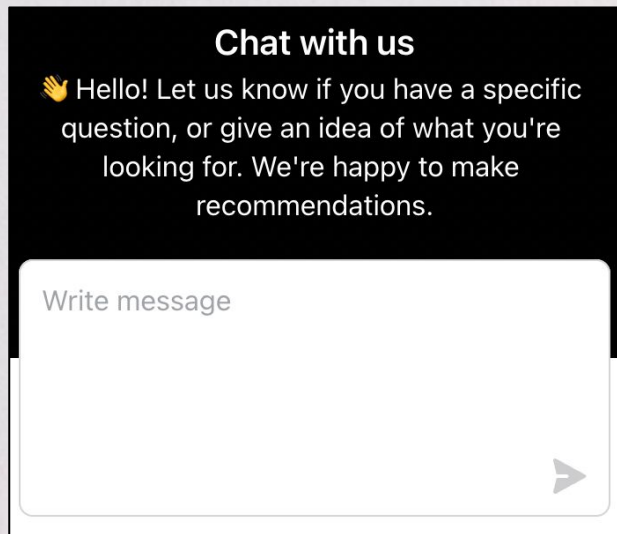**Conclusion**

1.   Nature of Dataset

2.   Nature of Language



3.   Testing of Streamlit API
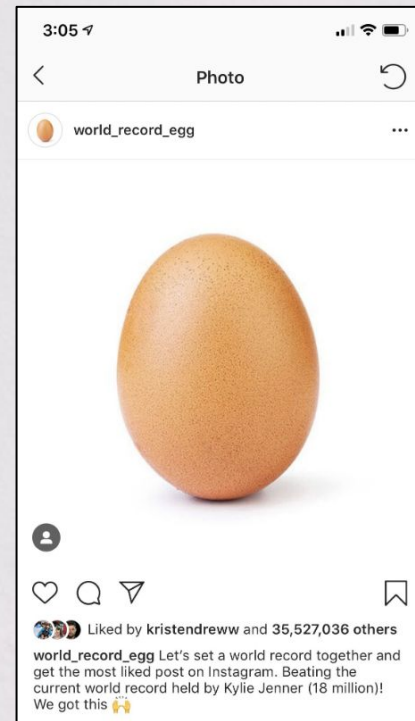
¯\\_(ツ)_/¯

# Future Work

1. Make the API more accessible



**Chat with us**
👋 Hello! Let us know if you have a specific question, or give an idea of what you're looking for. We're happy to make recommendations.

Write message



Linked in.

Messaging

# Future Work

2. Incorporate Images

# Deep Thanks!

Chua Soon Poh
DSI-32