August 26th 2021

# Modergator:
## Automated Content Moderation for Hateful Memes, Speech and Text

Korbinian Koch, Skadi Dinter, Katrin Caragiuli

**17%** of German youth aged 18-24 have been **personally affected** by hate speech online*

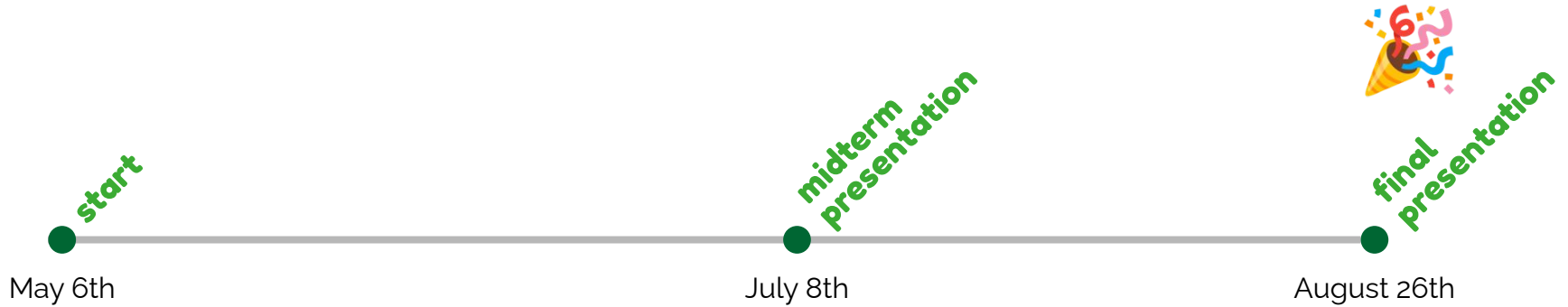*YouGov/Institut für Demokratie und Zivilgesellschaft (2019)

# The Mission



Providing safer social spaces for everyone by making group moderation for hateful content easier and faster.

# The Scope

- Master's project *Web Interfaces for Language Processing Systems*
- Team of 3⁽⁺⁾ students
- duration of ~4 months

start

midterm presentation

final presentation 🎉

May 6th       July 8th       August 26th

# The Team

**Korbinian**

*M.Sc. Intelligent Adaptive Systems*

bot
functionality &

voice
processing

design

(and more)

**Skadi**

*M.Sc. Informatics*

meme model +
api

interactivity and
feedback

target group
detection

(and more)

**Katrin**

*M.Sc. Informatics*

meme model +
api

target group
detection

meme
detection api

(and more)

# Modergator is ...

a collection of 6 APIs revolving around hate speech detection

**and**

a ready-to-use Telegram bot moderating hateful content in groups

| 1 | 2 | 3 |
| 4 | 5 | 6 |

modergator bot

*curated components*

*useable product*

# The Components

modergator bot

pyAPIstelegram-bot

Meme Detection
API

[0,1]

Contribution
from Niklas

OCR
API

Contribution
from Niklas

Text Hatefulness
API

[0,1]

Meme Hatefulness
API

[0,1]

Automatic Speech
Recognition
API

Target Group
API

[Race, Sexuality, ...]

Contribution
from Fabian

# The Features

Analyze voice messages

Automatic text classification: hateful, offensive or normal

Group members can discuss the classification(s)

GDPR-friendly opt-out

Detect affected target group(s)

Analyze memes

modergator bot

# Published Project



## 🐊 Modergator - Hate Detection for Text, Speech and Memes

Modergator is a Telegram bot able to moderate Telegram groups for hateful content.

Text messages are checked for whether they contain offensive and hateful speech, as well as the target groups that the speech is directed against (if there are any). Voice messages are transcribed and then handled the same way as a text messages.

Memes are also checked for hate which arises due to the combination of text and an image.

**Currently, the bot can only understand English language.**

### 🎯 Key Features

The bot will:

- check texts, voice messages and memes for hate
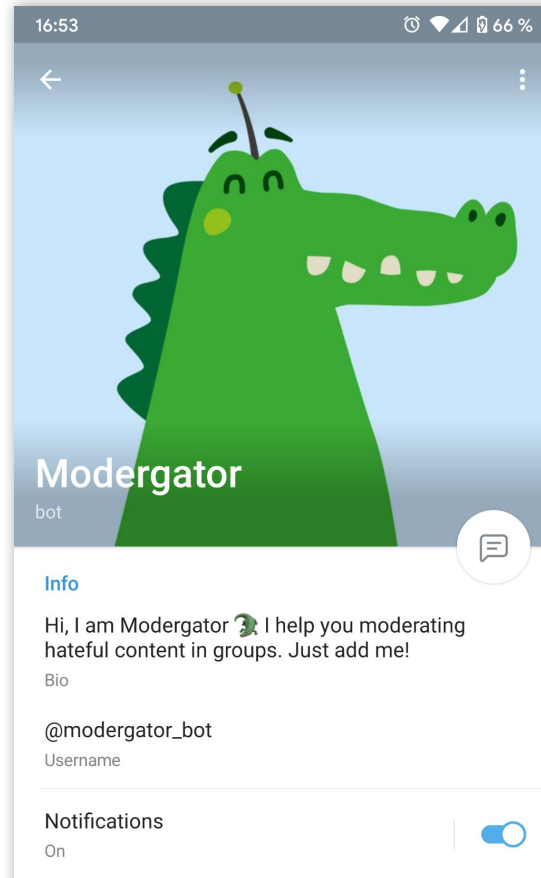- intervene if necessary

Group members can:

- dispute wrong classifications in a /poll
- optionally /optout of data processing (GDPR-compliant)

### 💡 How To Use

In order to interact with the bot, a Telegram account is needed. For instructions on how to create an account see: https://telegram.org/. To find
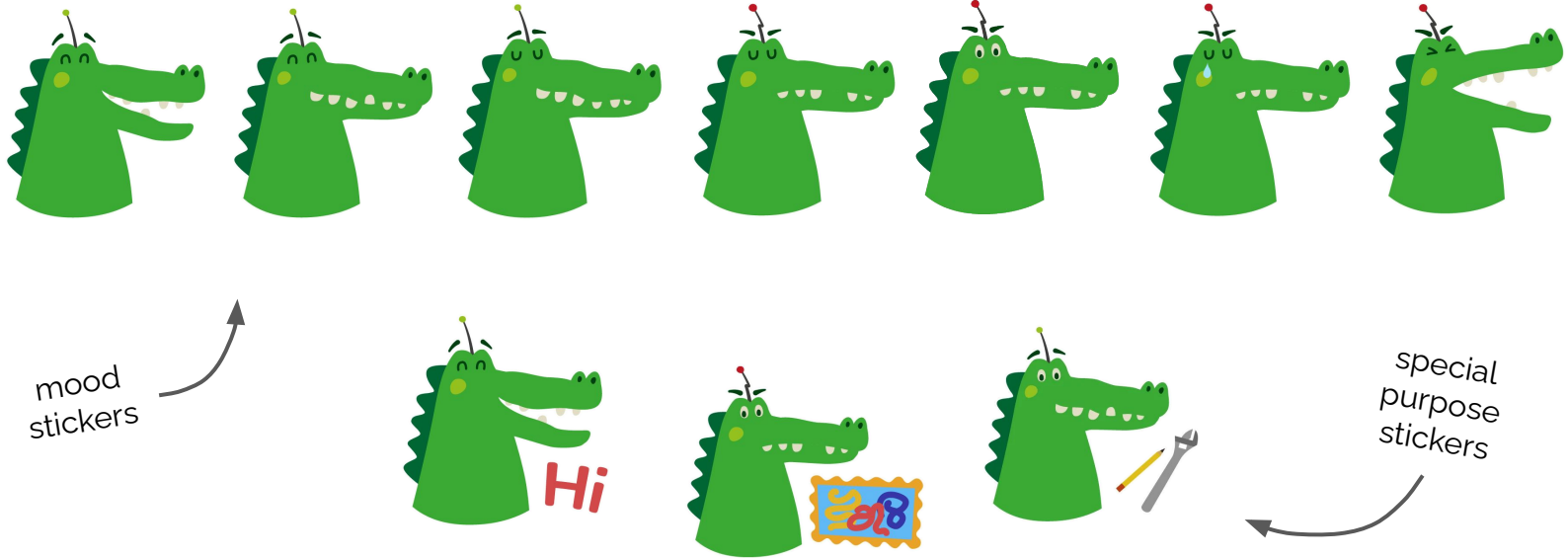
# Published Bot

# Design and UX: Name And Logo

modergator = portmanteau of moderator + alligator

# Design and UX: Stickers

Custom Telegram Stickers to provide engaging user experience and create emotional connection



mood stickers

special purpose stickers

Hi

# Technologies used

`python-telegram-bot`

# Telegram Bot

- programmed in Python
- defines all commands such as
  - /help
  - /optout
  - ...
- and accesses the other APIs
- responsible for receiving and filtering messages
- as well as sending answers

```python
def main() -> None:
    """Start the bot."""
    # Create the Updater and pass it your bot's token.
    updater = Updater(TOKEN)

    # Get the dispatcher to register handlers
    dispatcher = updater.dispatcher

    # on different commands - answer in Telegram
    dispatcher.add_handler(CommandHandler("start", start_command))
    dispatcher.add_handler(CommandHandler("about", about_command))
    dispatcher.add_handler(CommandHandler("help", help_command))
    dispatcher.add_handler(CommandHandler("optout", optout_command))
    dispatcher.add_handler(CommandHandler("joke", joke_command))
    dispatcher.add_handler(CommandHandler("howto", howto_command))
    dispatcher.add_handler(CommandHandler("poll", poll_command))
    dispatcher.add_handler(CommandHandler("optin", optin_command))
    dispatcher.add_handler(CommandHandler("debug", debug_command))
    dispatcher.add_handler(CommandHandler("goodvibes", goodvibes_command))
    dispatcher.add_handler(PollAnswerHandler(receive_poll_answer))
    dispatcher.add_handler(MessageHandler(Filters.poll, receive_poll))
    dispatcher.add_handler(MessageHandler(Filters.status_update.new_chat_members, welcome_message))

    # on non command i.e message - echo the message on Telegram
    dispatcher.add_handler(MessageHandler(Filters.text & ~Filters.command, handle_text))
    dispatcher.add_handler(MessageHandler((Filters.photo | Filters.document.category('image')) & ~Filters.comman
    dispatcher.add_handler(MessageHandler(Filters.voice & ~Filters.command, handle_voice))
```

# Target Detection

- Based on HateXplain data set (https://github.com/hate-alert/HateXplain)
- Built a model and training pipeline
- Classify input into one (or more) target groups out of 24
- Evaluation results:
  - F1: , 0.058, Precision: 0.3, Recall: 0.032
- Immense help from Fabian

> I think that this text message is offensive. Please be nice and stick to the community guidelines.
>
> Your hate was probably directed towards the following group(s): Women.

# Meme Detection

- Contribution from Niklas
- We built the meme-detection-api
- In the bot the input image is classified: is it a meme?
  - If False: do nothing with the image
  - If True: hand the image over to the ocr-api and meme-model-api to classify

```python
398    def detect_meme(url):
399        print("Start Meme Detection")
400        params = {"url": url}
401        r = requests.get(url=f"http://127.0.0.1:{PORTDICT['meme-detection-api']}/classifier", params=params)
402        is_meme = r.json()["result"]
403        print("is_meme: ", is_meme)
404        return is_meme
```
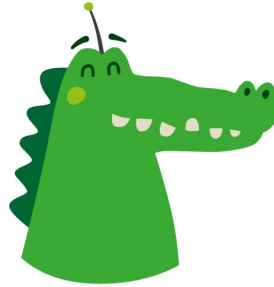
# Documentation

- **Swagger**
  - Build with flask-apispec
  - Automatically builds Swagger documentation
  - Documentation can be used to test functionality
  - Later shown in demo

- **Readme**
  - Detailed instructions both for the end-user and anyone hosting an instance of the bot
  - Includes links to files that need to be downloaded

Demo

# The End.

Thank you for your attention.