# Multimodal Moderator

AI content moderation for text and images

# Problem

## Single modalities

Content moderators can process text or images, but not both. Separate applications are needed to moderate multiple content types.

## Binary output

Content moderators produce binary benign or toxic output. Most don't provide explanation about why certain content is not appropriate.

# Solution

The Multimodal Moderator

- Single application that can check if text or image is appropriate or not
- Can "understand" the message and provide explanation about why certain content is not appropriate

# Technology

## Zapier

We used Zapier with Discord as the front end. Zapier is a no-code solution to connect with an AI model, and does not require a constantly running program to connect with Discord.

## GPT-4 Vision

GPT-4 Vision hosted on Clarifai platform is the engine that can check if text or image is appropriate. It sends back a narrative to explain the reasoning.
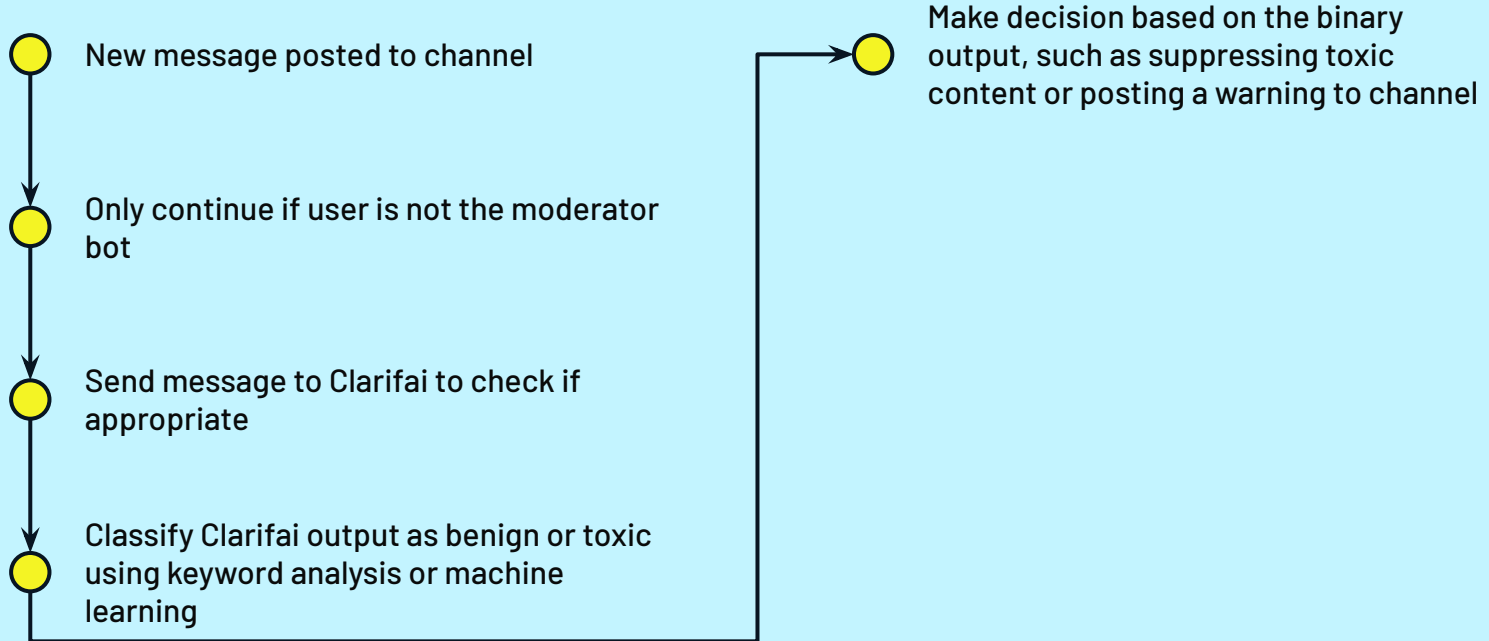
# Workflow

# Future State

New message posted to channel

Only continue if user is not the moderator bot

Send message to Clarifai to check if appropriate

Classify Clarifai output as benign or toxic using keyword analysis or machine learning

Make decision based on the binary output, such as suppressing toxic content or posting a warning to channel

# Demo



**Discord - https://discord.gg/qsw2bJqJ**

# Team

## Don

AI Developer

## Vesselin

Entrepreneur Scientist

# Thanks!

**Thanks to Lablab and Clarifai for hosting.
Do you have any questions?**
https://discord.gg/qsw2bJqJ