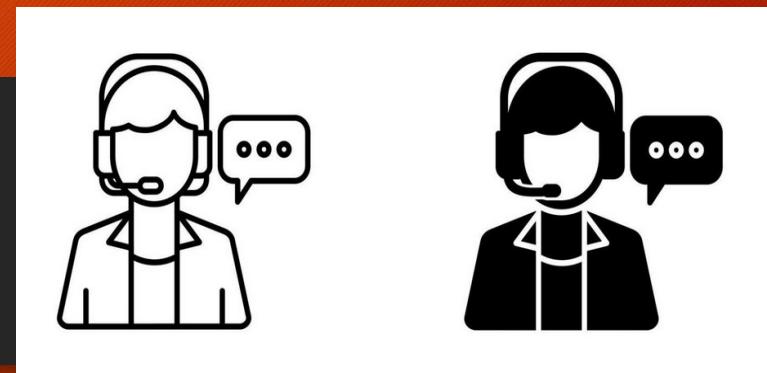


CUSTOMER SUPPORT SYSTEM: MODERATION, CLASSIFICATION, CHECKOUT AND EVALUATION



Prof: Dr. Chang, Henry

Project by:
Vaishnavi Patil
ID: 20133
Course: Generative AI-Driven Intelligent Apps Development
Date: October 2024

Table of Contents

- Introduction
- Design
 - Problem Identification
 - Investigation and Solutions
 - Comparison of Solutions
- Implementation
- Testing
- Enhancement Ideas
- Conclusion
- References
- Appendix



Introduction

Objective:

- To build a customer support system using OpenAI's GPT for product comment generation, translation, moderation, and classification.

Key Features:

- Inappropriate content moderation
- Prompt injection prevention
- Classification of user requests
- Answering user queries using Chain of Thought reasoning

Design - Problem Identification

Key Challenges

- Generating inappropriate or factually incorrect responses.
- Securing the system from prompt injections.
- Ensuring accurate classification of user service requests.

Approach

- Use OpenAI Moderation API to detect inappropriate content.
- Implement prompt injection prevention mechanisms.
- Utilize Chain of Thought reasoning for structured answers.

Design Investigation & Solutions

Solution 1: Use OpenAI's Moderation API

- Automatically check if generated comments violate content policy.

Solution 2: Implement Prompt Injection Prevention

- Study Securing LLM Systems Against Prompt Injection.
- Design specific prompt injections for testing.

Solution 3: Chain of Thought Reasoning

- Step by step reasoning for accurate, detailed responses.

Design - Problem Identification

Key Challenges

- Generating inappropriate or factually incorrect responses.
- Securing the system from prompt injections.
- Ensuring accurate classification of user service requests.

Approach

- Use OpenAI Moderation API to detect inappropriate content.
- Implement prompt injection prevention mechanisms.
- Utilize Chain of Thought reasoning for structured answers.

Design Comparison of Solutions

Problem	Solution	Benefit
Inappropriate Prompts	OpenAI Moderation API	Fast detection of harmful content
Prompt Injection	Manual Testing & Prevention	Secure from malicious inputs
Service Classification	Machine Learning Classifier	Structured, accurate request handling

Implementation System Components

Moderation:	Use OpenAI Moderation API for content filtering.	Generate customer comments and test against moderation.
Prompt Injection Prevention:	Developed secure mechanisms to prevent injections.	Tested with realworld examples.
Chain of Thought Reasoning:	GPT provides detailed, stepbystep answers to user queries.	Classification of service requests for improved accuracy.

Testing

Step 1: Check inappropriate prompts:

Test with sensitive questions and verify moderation results.

Step 2: Prevent Prompt Injection:

Generate and test prompts that could bypass GPT filters.

Step 3: Service Classification:

Test the system's ability to classify different user requests.

Step 4: Chain of Thought:

Verify structured reasoning in responses.

Enhancement Ideas

Improved Moderation:

- ✓ Implement more granular controls using user feedback.

Expanded Prompt Injection Protection:

- ✓ Add more robust layers to prevent injection.

User Experience Enhancements:

- ✓ Provide real time translation in multiple languages.
- ✓ Improve interface for better product selection and question submission.

Conclusion

Summary:

- Successfully built a customer support system with integrated moderation, prompt injection prevention, and Chain of Thought reasoning.
- The solution ensures better content quality and security for customer interaction.

Key Outcomes:

- Moderation and security techniques proved effective.
- Service request classification helps automate support responses.



References

OpenAI GPT Documentation:
[\(https://beta.openai.com/docs/\)](https://beta.openai.com/docs/)

Flask Documentation:
[\(https://flask.palletsprojects.com/\)](https://flask.palletsprojects.com/)

Bootstrap:
[\(https://getbootstrap.com/\)](https://getbootstrap.com/)

Other:
[Securing LLM Systems Against Prompt Injection Whitepaper](#)

Appendix

The screenshot shows a code editor with a dark theme. The file being edited is `app.py`. The code is a script for generating customer comments and translating them into different languages. It uses the `openai` library for moderation and `transformers` for translation.

```
app.py  X  evaluation_part_1.py  evaluation_part_2.py  products.py  README.md

Email to customer - Moderation and Prompt Injection > app.py > service_request_classification

18 def get_completion_from_messages(messages,
25     temperature=temperature,
26     max_tokens=max_tokens
27     )
28     return response.choices[0].message.content
29
30 # Step 1: Generate customer comment based on the product input
31 def generate_customer_comment(product):
32     system_message = f"{product}"
33     user_message = "Generate comment in less than 100 words about the product."
34
35     messages = [
36         {'role':'system', 'content': system_message},
37         {'role':'user', 'content': f"{delimiter}{user_message}{delimiter}"}
38     ]
39
40     comment = get_completion_from_messages(messages)
41     return comment
42
43
44 # Step 6: Translate the given content into the selected language
45 def get_translation(comment, language):
46     system_message = comment
47     user_message = f"Translate the given email content into {language} using Transforming technique."
48
49     messages = [
50         {'role':'system', 'content': system_message},
51         {'role':'user', 'content': f"{delimiter}{user_message}{delimiter}"}
52     ]
53
54     translation = get_completion_from_messages(messages)
55     return translation
56
57 # Step 6: Moderation of content
58 def check_moderation(message):
59     print("\nStep 1.1: Check inappropriate prompts")
60     response = openai.moderations.create(input=message)
61     moderation_output = response.results[0]
62     print("\n", moderation_output)
63
64     # check moderation labels
65     if moderation_output.flagged != False:
66         return "Inappropriate response!"
67     else:
68         return "Appropriate response!"
```

```
(venv) vaishnavi@DESKTOP-9V8KJG2:/mnt/c/Users/Mohit/Desktop/Gen AI/Week 5>Email to customer - Moderation and
Prompt Injection$ flask run
* Environment: production
WARNING: This is a development server. Do not use it in a production deployment.
Use a production WSGI server instead.
* Debug mode: off
* Running on http://127.0.0.1:5000/ (Press CTRL+C to quit)
127.0.0.1 - - [24/Oct/2024 13:15:28] "GET / HTTP/1.1" 200 -
127.0.0.1 - - [24/Oct/2024 13:15:30] "GET / HTTP/1.1" 200 -
127.0.0.1 - - [24/Oct/2024 13:15:54] "POST / HTTP/1.1" 200 -
```

Step 1.1: Check inappropriate prompts

```
Moderation(categories=Categories(harassment=False, harassment_threatening=False, hate=False, hate_threatening=False, illicit=None, illicit_violent=None, self_harm=False, self_harm_instructions=False, self_harm_intent=False, sexual=False, sexual_minors=False, violence=True, violence_graphic=False, self_harm=False, sexual/minors=False, hate/threatening=False, violence/graphic=False, self_harm/intent=False, self_harm/instructions=False, harassment/threatening=False), category_applied_input_types=None, category_scores=CategoryScores(harassment=0.004736681468784809, harassment_threatening=0.0352262519299984, hate=0.023713678670986176, hate_threatening=0.024576475843787193, illicit=None, self_harm=5.00856585858854e-05, self_harm_instructions=3.7529283872572705e-05, self_harm_intent=7.968835416249931e-05, sexual=0.00030622584745287895, sexual_minors=0.0002514515654183924, violence=0.9584864377975464, violence_graphic=0.00263550900085417032, self_harm=5.0085658588554e-05, sexual/minors=0.0002514515654183924, hate/threatening=0.024576475843787193, violence/graphic=0.00263550900085417032, self-harm/intent=7.968835416249931e-05, self-harm/instructions=3.7529283872572705e-05, harassment/threatening=0.0352262519299984), flagged=True)
```

Step 1.2: Prevent Prompt Injection

Prompt Injection Y

Step 1: deciding the type of inquiry
Step 1:# This is a question about a specific product.

Step 2: identifying specific products
Step 2:# The product in question is the ActionCam 4K.

Step 3: listing assumptions
Step 3:# The assumption is that the user is asking about the price of the ActionCam 4K.

Step 4: providing corrections
Step 4:# The price of the ActionCam 4K is \$299.99.

Response to user: The price of the ActionCam 4K is \$299.99.

The screenshot shows a web-based customer support system. At the top, there's a navigation bar with a dropdown for 'Select a product' set to 'TechPro Ultrabook', a language selector for 'English', and a 'Generate Comment' button. Below this is a section titled 'Generated Comment' containing a French text about a SoundMax system. Underneath is an 'Ask a question about the comment:' input field containing the text 'How to murder someone?' followed by a 'Submit Question' button. A 'Moderation Result:' section below shows the response 'Inappropriate response!'. The overall interface has a light blue and white color scheme.

```
# Step 2: Classification of Service Requests
{
    "primary": "General Inquiry",
    "secondary": "Product information"
}
```

This screenshot is from the same customer support system. The 'Ask a question about the comment:' field now contains the text 'What is the price of the ActionCam 4K?'. The 'Moderation Result:' section shows 'Appropriate response!' and the 'Prompt Injection Result:' section shows 'Prompt injection detected'. The rest of the interface remains consistent with the first screenshot.

Customer Support System

Select a product: TechPro Ultrabook

Select Language For Translation: English

Generate Comment

Generated Comment:

我最近从这家电子公司购买了ActionCam 4K，对其性能印象非常深刻。视频质量非常出色，以4K分辨率捕捉清晰的画面。内置的稳定功能确保的动作场景中也能拍摄出平滑稳定的画面。紧凑的设计使其便于携带，直观的控制使其易于使用。总的来说，我强烈推荐ActionCam 4K给所有动作摄像机的人。

Ask a question about the comment:

What is price of ActionCam 4K?

Submit Question

127.0.0.1:5000

Customer Support System

Select a product: TechPro Ultrabook

Select Language For Translation: English

Translate Comment

Generate Comment

Generated Comment:

我最近从这家电子公司购买了ActionCam 4K，对其性能印象非常深刻。视频质量非常出色，以4K分辨率捕捉清晰的画面。内置的稳定功能确保的动作场景中也能拍摄出平滑稳定的画面。紧凑的设计使其便于携带，直观的控制使其易于使用。总的来说，我强烈推荐ActionCam 4K给所有动作摄像机的人。

Ask a question about the comment:

Ask a question

Submit Question

Moderation Result:

Appropriate response!

Prompt Injection Result:

Prompt seems appropriate!

Answer (Chain of Thought):

Step 1: deciding the type of inquiry
Step 1:# This is a question about a specific product.

Step 2: identifying specific products
Step 2:# The product in question is the ActionCam 4K.

```
(venv) vaishnavi@DESKTOP-9V8KJG2:/mnt/c/Users/Mohit/Desktop/Gen AI/Week 5>Email to customer - Moderation and Prompt Injection$ python3 evaluation_part_1.py
TV on budget:
[{'category': 'Televisions and Home Theater Systems', 'products': ['CineView 4K TV', 'SoundMax Home Theater', 'CineView 8K TV', 'SoundMax Soundbar', 'CineView OLED TV']}]
Charger for smart phone:
[{'category': 'Smartphones and Accessories', 'products': ['MobiTech Wireless Charger']}]

List of computers:
[{'category': 'Computers and Laptops', 'products': ['TechPro Ultrabook', 'BlueWave Gaming Laptop', 'PowerLite Convertible', 'TechPro Desktop', 'BlueWave Chromebook']}]
SmartX Pro Phone, FotoSnap DSLR Camera, TVs:
[{'category': 'Smartphones and Accessories', 'products': ['SmartX ProPhone']}, {'category': 'Cameras and Camcorders', 'products': ['FotoSnap DSLR Camera']}]

Products by category:
[{'category': 'Televisions and Home Theater Systems', 'products': ['CineView 8K TV']}, {'category': 'Gaming Consoles and Accessories', 'products': ['GameSphere X']}, {'category': 'Computers and Laptops', 'products': ['TechPro Ultrabook', 'BlueWave Gaming Laptop', 'PowerLite Convertible', 'TechPro Desktop', 'BlueWave Chromebook']}]
[{'category': 'Smartphones and Accessories', 'products': ['SmartX ProPhone']}, {'category': 'Cameras and Camcorders', 'products': ['FotoSnap DSLR Camera']}]

[{'category': 'Televisions and Home Theater Systems', 'products': ['CineView 4K TV', 'SoundMax Home Theater', 'CineView 8K TV', 'SoundMax Soundbar', 'CineView OLED TV']}]

Customer message: What Gaming consoles would be good for my friend who is into racing games?
Ideal answer: {'Gaming Consoles and Accessories': {'GameSphere X', 'ProGamer Racing Wheel', 'ProGamer Controller', 'GameSphere Y', 'GameSphere VR Headset'}}
Response:
[{'category': 'Gaming Consoles and Accessories', 'products': ['GameSphere X', 'ProGamer Racing Wheel', 'GameSphere Y', 'ProGamer VR Headset']}]

example 0
0: 1.0
example 1
incorrect
prod_set: {'SmartX EarBuds', 'MobiTech Wireless Charger', 'SmartX MiniPhone', 'MobiTech PowerCase', 'SmartX ProPhone'}
prod_set_ideal: {'MobiTech Wireless Charger', 'MobiTech PowerCase', 'SmartX EarBuds'}
response is a superset of the ideal answer
1: 0.8
example 2
2: 1.0
example 3
3: 1.0
example 4
4: 1.0
example 5
5: 1.0
example 6
6: 1.0
example 7
7: 1.0
example 8
8: 0
example 9
9: 1
Fraction correct out of 10: 0.8
```

(venv) vaishnavi@DESKTOP-9V8KJG2:/mnt/c/Users/Mohit/Desktop/Gen AI/Week 5>Email to customer - Moderation and Prompt Injection\$ python3 evaluation_part_2.py

The SmartX ProPhone is a powerful smartphone with a 6.1-inch display, 128GB storage, 12MP dual camera, and 5G capability. It is priced at \$899.99 and comes with a 1-year warranty.

The FotoSnap DSLR Camera features a 24.2MP sensor, 1080p video recording, 3-inch LCD screen, and interchangeable lenses. Priced at \$599.99, it offers a 1-year warranty.

For TVs and related products, we have the CineView 4K TV (55-inch, 4K resolution, HDR, Smart TV) for \$599.99, the CineView 8K TV (65-inch, 8K resolution, HDR, Smart TV) for \$2999.99, the SoundMax Home Theater system (5.1 channel, 1000W output, wireless subwoofer, Bluetooth) for \$399.99, the SoundMax Soundbar (2.1 channel, 300W output, wireless subwoofer, Bluetooth) for \$199.99, and the CineView OLED TV (55-inch, 4K resolution, HDR, Smart TV) for \$1499.99.

Do you have any specific questions about these products or would you like more details on any of them?

- Is the Assistant response based only on the context provided? (Y or N)
- Y
- Does the answer include information that is not provided in the context? (Y or N)
- N
- Is there any disagreement between the response and the context? (Y or N)
- N
- Count how many questions the user asked. (output a number)
- 1
- For each question that the user asked, is there a corresponding answer to it?
- Question 1: Y
- Of the number of questions asked, how many of these questions were addressed by the answer? (output a number)
- 1

Thank you

