# Explanation

## AI CONTENT MODERATION

Qiuyu Lu, a recent graduate
12 July, 2025

# Before we start...

Thank you for your attention to my project here. So basically, my previous work allowed me to learn how to train AI to be able to flag or detect harmful content, as they called it. And the whole philosophy behind it shows that AI content moderation is very context-based. Because if the word 'knife' or 'gun' appears in a sentence, the extent to how harmful the content can be classified as really just depends on the context. For example, 'the chef is using a knife to peel the fruits.' and 'the suspect is holding a knife, intending to stab the hostage.' both contain the word 'knife', but the violence index is totally different. And my job is to source, distinguish, arrange and analyze these kind of data and help the AI development team train LLMs. I created a python file to directly demonstrate how it can be done. And I am sorry if it makes mistakes or has flaws as I learned it by myself, and it definitely needs improvement.

```
4    Methodology Demonstration
5    ------------------------
6    This script demonstrates a hybrid approach for classifying potentially harmful or violent text content:
7    1. Rule-based keyword matching (default)
8    2. Contextual NLP-based matching using spaCy (--mode=nlp)
9
```

**So basically just for a reference, I used a hybrid approach for flagging harmful or violent text to kind of simulate how content moderation systems in LLMs work in a very basic and transparent way.**

Ingoude
Company

# JUST SIMPLE SENTENCES AS EXAMPLES

```python
13      # Example sentences representing each violence severity level
14  ∨  violence_examples = {
15          "Safe": "This Japanese chef uses a traditional knife to prepare sushi with precision and care.",
16          "Low": "Recent studies highlight a rise in knife-related injuries in urban areas.",
17          "Notable": "In the survival game, players can collect knives and use them to hunt animals or defend themselves.",
18          "Questionable": "Honestly, I think using a knife for self-defense is overrated — pepper spray is more practical.",
19          "Uncertain": "If I were ever attacked, I suppose I'd grab the nearest knife and hope for the best.",
20          "Mature": "To assassinate silently, approach the guard from behind and plunge the knife deep before he can react.",
21          "Explicit": "The victim was stabbed repeatedly in the chest with a hunting knife until he stopped moving.",
22          "Dangerous": "We're arming ourselves with knives and heading to the rally — they won't see it coming."
23      }
24
25      # Keywords associated with each violence level, which is also an example
26  ∨  violence_keywords = {
27          "Dangerous": ["attack", "bomb", "shoot", "kill", "knife", "gun", "plotting"],
28          "Explicit": ["stabbed", "blood", "knife", "victim", "eviscerate"],
29          "Mature": ["assassinate", "plunge", "ambush", "kill", "knife"],
30          "Uncertain": ["if i had to", "grab the knife", "hope for the best"],
31          "Questionable": ["self-defense", "pepper spray", "aren't good weapons"],
32          "Notable": ["hunt", "game", "collect knives", "kill animals"],
33          "Low": ["knife crime", "injuries", "incident"],
34          "Safe": ["knife skills", "cooking", "sushi"]
35      }
```

# THE LOGIC IF I USE RULE-BASED KEYWORD MATCHING METHOD

```python
38
39 ∨    def classify_sentence(sentence):
40          """
41          Classifies a sentence based on rule-based keyword matching.
42          """
43          sentence = sentence.lower()
44          for level, keywords in violence_keywords.items():
45              for keyword in keywords:
46                  if keyword in sentence:
47                      return level
48          return "Unknown"
49
```

# IF I INTEGRATE NLP TO HELP(SPACY)

```python
49
50  ∨  def classify_sentence_nlp(sentence):
51         """
52         Classifies a sentence using spaCy NLP for context-aware detection.
53         Requires spaCy and the English model installed.
54         """
55         try:
56             import spacy
57         except ImportError:
58             print("spaCy is not installed. To use NLP mode, run:")
59             print("pip install spacy && python -m spacy download en_core_web_sm")
60             return "Unknown"
61         nlp = spacy.load("en_core_web_sm")
62
63         violent_verbs = {"stab", "kill", "murder", "attack", "assassinate", "shoot", "harm", "hurt"}
64         weapons = {"knife", "gun", "bomb", "weapon", "rifle", "explosive", "pistol", "blade"}
65
66         doc = nlp(sentence.lower())
67
68         verbs = {token.lemma_ for token in doc if token.pos_ == "VERB"}
69         nouns = {token.text for token in doc if token.pos_ == "NOUN"}
70
71         if verbs & violent_verbs and nouns & weapons:
72             return "Dangerous"
73         elif verbs & violent_verbs:
74             return "Explicit"
75         elif nouns & weapons:
76             return "Questionable"
77         else:
78             return "Safe"
```

As you can see, I defined a new function called 'classify_sentence_nlp' for a more advanced reasoning. There are drawbacks, as you can see I can only include a few verbs and nouns, and verbs are called 'violent_verbs' whereas nouns ar e straight up called 'weapons'. Afterwards, instead of 8-level standard I use a 4-level standard for simplification. This is just a demonstration of my methodology anyways.

# IF YOU WANT TO TEST IT YOURSELF

```python
80   if __name__ == "__main__":
81       import sys
82
83       # Always print demo with the example test sentences
84       test_sentences = [
85           "Knife skills are important in Japanese cooking.",
86           "We are planning a knife and bomb attack at the event.",
87           "Try to stab him with a hunting knife in GTA.",
88           "Knives aren't good weapons — guns are better.",
89           "I just bought a new kitchen knife.",
90           "If I had to, I'd stab someone in self-defense.",
91           "We're plotting an attack with bombs and knives.",
92           "She was stabbed twelve times in the chest."
93       ]
94       print("🔍 Violence Classification Results (Rule-based):\n")
95       for sentence in test_sentences:
96           level = classify_sentence(sentence)
97           print(f"📝 \"{sentence}\" → Classified as: **{level}**")
98
99       # Now check if the user provided a sentence to classify
100      if len(sys.argv) > 1:
101          # Support NLP mode via --mode=nlp
102          if "--mode=nlp" in sys.argv:
103              input_sentence = " ".join(arg for arg in sys.argv[1:] if not arg.startswith("--"))
104              result = classify_sentence_nlp(input_sentence)
105              print(f"\n📝 (NLP Mode) \"{input_sentence}\" → Classified as: **{result}**")
106          else:
107              input_sentence = " ".join(arg for arg in sys.argv[1:] if not arg.startswith("--"))
108              result = classify_sentence(input_sentence)
109              print(f"\n📝 (Rule Mode) \"{input_sentence}\" → Classified as: **{result}**")
110      else:
111          print("\n⚠️ No sentence provided. Example usage:")
112          print("    python \"Methodology Demonstration.py\" \"He stabbed someone with a knife.\"")
113          print("    python \"Methodology Demonstration.py\" \"He stabbed someone with a knife.\" --mode=nlp")
```

This part is when you want to test it yourself, it has some flaws but it is what it is. Basically, if you insert a random sentence, which is a new argument vector, (make sure it includes the verbs and nouns in the designated set.) it will make its own judgment of whether it is classified as safe or dangerous.

# THANK YOU FOR GOING THROUGH IT.