



Advancing AI Incidents Classification: Leveraging LLMs with Strategic Prompting

Yian Chen, Lana Do, Liheng Yi, Ricardo Baeza-Yates and
John A. Guerra-Gomez

EasyChair preprints are intended for rapid dissemination of research results and are integrated with the rest of EasyChair.

February 11, 2025

Advancing AI Incidents Classification: Leveraging LLMs with Strategic Prompting

Yian Chen, Lana Do, Liheng Yi, Ricardo Baeza-Yates^[0000-0003-3208-9778], and
John A. Guerra-Gomez^[0000-0001-7943-0000]

Northeastern University, Oakland, CA 94613, USA

{chen.yian1, do.ng, yi.lih, r.baeza-yates, jguerra}@northeastern.edu

Abstract. This study examines the efficacy of large language models (LLMs), particularly GPT-4, in classifying AI incident reports documented in the AI Incidents Database (AIID), with the goal of enhancing our understanding and management of AI-related harm. The data of incident reports are all on news events that detail specific incidents related to AI technology that have resulted in harmful effect on our society. We explore the use of different prompting techniques on GPT-4 and assess the classification results of those incidents by subjective and objective evaluations. This work lays the groundwork for a comprehensive, automated classification framework for AI incident reporting, balancing LLM capabilities with the intricacies inherent in human judgment.

Keywords: Large Language Models · LLM Classification · Prompt Engineering · Responsible AI.

1 Introduction

The proliferation of artificial intelligence (AI) technologies has introduced unprecedented opportunities for innovation across diverse sectors. However, with these advancements come emerging concerns about AI failures. Instances of AI causing harm and sparking controversies underscore the critical need for effective classification frameworks to systematically categorize these incidents.

As part of a larger study of Responsible AI [1], this study primarily focuses on the development of a robust classification methodology to categorize a vast volume of technology-related incidents sourced from the AI Incidents Database (AIID)¹. The AIID serves as a repository of news articles documenting instances where AI technologies have caused harmful outcomes in the real-world.

Our objective is to streamline the process of the incidents classification for data visualization. However, the sheer volume of incidents—exceeding 600 cases—renders manual classification impractical. To address this challenge, we propose leveraging large language models (LLMs) to automate the classification process.

¹ <https://incidentdatabase.ai/>

Pre-trained transformer language models, such as GPT-3 and GPT-4, demonstrate great in-context learning capabilities through natural language task description [2]. Those LLMs exhibit breakthrough performances in various downstream tasks such as text-interpretation through semantic processing abilities [4] and medical diagnoses or theorem proving through its reasoning capabilities by general deduction [13]. We harness the power of these Artificial Intelligence Generated Content (AIGC) tools for our classification task by providing the data from AIID to the LLM and use various prompting techniques to refine the results. In our use-case, it is critical to ensure that the automated classifications are not only well-supported by the data given, but also closely aligned with human judgment and standards. This necessitates a validation process to ascertain the reliability and trustworthiness of the classification outcomes against established human benchmarks.

By utilizing LLM technology, we aim to develop a scalable and efficient approach to classify these AI-related incidents. Our efforts are geared towards providing insights into the prevalence, trends, and implications of AI-related harm, thereby contributing to the advancement of responsible AI researches.

AI engineering is an emerging field, and consequently, AI incidents documentation practices are still evolving. Our contributions to the field include:

- An overview of three prompting techniques we tested for classifying AI-related incidents by a predefined taxonomy of classes of harm, geographic location, societal impact, intended application, field of deployment, and others. We evaluate the three different techniques and the results and compare the LLM classifications to our human manual classifications.
- A classification of 17 incidents from AIID by LLM using various prompting techniques and another benchmark classification set by human manual efforts. Future efforts will aim to apply the most appropriate classification technique to the full 628 incidents on AIID and use those results for data visualization as part of the Responsible AI study.
- Evaluation and validation of the classification results by two domain experts, which assess the capability of the LLM in this classification task.

2 Related Work

2.1 AI Incidents Documentation

The documentation of AI failures has become increasingly crucial as AI technologies are more widely deployed in sensitive and critical contexts. To improve the quality of AI systems and prevent future incidents, the AI engineering community needs a comprehensive record of previous AI failures[10]. Various databases have been established to document AI incidents.

- **The AI and Algorithmic Incident and Controversies (AIAAIC) Repository** An independent initiative that logs global AI incidents and

controversies, managed by Charlie Pownall, and is used by over 60 universities and organizations ².

- **The AI Vulnerability Database (AVID)** The database focuses on cataloging proven high-level failure modes and specific vulnerabilities, providing valuable insights for AI engineers and auditors ³.
- **Awful AI** A curated list to track current harmful usage of AI, aiming to raise awareness about the dark side of AI technologies ⁴.
- **AI Incidents Database (AIID)** The most utilized database for reporting AI incidents ⁵. It encompasses a wide range of AI-related harms affecting physical health, safety, and social/political systems. This database also serves as the primary resource for our study.

Research using the records, including efforts to categorize the causes of AI failures [12] and a sector-based approach analysis of incidents [3], underscores the need for systematic approaches to classify and analyze these events. Despite the availability of databases cataloging numerous AI failures, similar incidents continue to occur. For instance, there are recurring instances where facial recognition systems disproportionately misidentify people of color, leading to wrongful accusations and arrests, despite widespread documentation and acknowledgment of bias in AI systems. Industry practitioners face significant challenges in effectively preventing and mitigating bias, often only identifying these serious issues after deployment [14].

These studies indicate that there is still much to learn and significant work to be done to develop a robust classification model capable of handling the complexity and variety of incident reports.

2.2 Large Language Models

With the advent of *pre-trained* language models (LM) with fixed architecture came a new paradigm shift in the learning of these models. LMs built on transformer architecture no longer rely on objective engineering nor fine-tuning to perform specific downstream tasks [9]. In the new paradigm, pre-trained LMs can solve various tasks with only textual *prompts* [9]. Previous works that have shown how LLMs can adapt and recognize tasks at inference time, also referred to as *in-context* learning [2]. These large models showed improved ability to learn a task from contextual information, where a model is guided by textual cues to perform a task. By processing an instruction or a handful of task examples, the model then predicts subsequent actions within the same context.

Specifically, in the domain of classification, LLMs have been applied with some success. For example, recent studies have deployed LLMs for classifying legal documents using detailed criteria [15] and have evaluated the GPT family of models for their effectiveness in biomedical reasoning and classification tasks

² <https://www.aiaaic.org/aiaaic-repository>

³ <https://avidml.org/>

⁴ <https://github.com/daviddao/awful-ai>

⁵ <https://github.com/responsible-ai-collaborative/aiid>

[6]. The flexibility of LLMs to parse and organize complex data has facilitated their adoption in various burgeoning fields. However, relying entirely on LLM responses can pose risks, such as generating misleading answers in scenarios like academic forums, where models may produce incorrect responses if they fail to understand a question [18]. To mitigate these risks, type-specific prompting has been proposed to enhance the performance of GPT models in classification tasks. Our research builds upon this foundation by utilizing LLMs to automate the classification of AI incidents according to a predefined taxonomy, employing diverse prompting techniques to enhance this process.

3 Background

Previously introduced in Section 2.1, the dataset used in our study is sourced from the AI Incidents Database (AIID) GitHub repository. This dataset was chosen due to its extensive collection of incidents, featuring over 1,000 archived reports from more than 600 contributors across various sectors. These incidents document real-world cases where the use of AI has resulted in harm, covering a range of issues from facial recognition and targeted advertising to collisions involving autonomous vehicles, among many others.

Before introducing our task setup, we present the main classification fields that our LLM aims to identify. The predefined taxonomy for this study encompasses classes of AI-related harm ranging from discrimination, disinformation, human incompetence, pseudoscience, copyright violation, mental health implications, to environmental impact. In addition to the taxonomy of classes of harm, the LLM should also identify other fields of classification germane to our study of Responsible AI.

3.1 Classification Fields Explained

To clearly define the scope and detail of the data the LLM analyzes in each incident, we outline the specific fields it is designed to classify:

- **Geographical Location (Country, State, City, Continent):** These fields determine an incident’s location, with the possibility of indicating "Worldwide" for events not confined to a single geographical area.
- **Company and Location (Company, Company City, Company State):** This identifies the organization responsible for the AI technology involved in the incident, including the precise location of its headquarters.
- **Affected Population:** This field aims to identify which demographic groups were impacted by the incident, shedding light on the societal segments most at risk from irresponsible AI usages.
- **Number of People Affected:** Both the actual and potential number of individuals impacted by the incident are quantified to gauge its severity and reach.

- **Taxonomic Classification of Classes of Harm (Classes, Subclasses, Sub-subclass):** Through a structured taxonomy, incidents are categorized into specific classes and subclasses of AI-related harm, enabling a systematic analysis of the types of irresponsible AI usage occurring. This taxonomy was collectively decided by the team and serve as labels or tags for the incident classification.
- **Application Area:** The specific domain of AI application involved in the incident is identified, such as healthcare, surveillance, or social media, providing insights into where AI poses the greatest risks.
- **Online:** Indicates whether the incident took place online, highlighting the digital nature of many AI-related ethical concerns.

Each classification field plays a crucial role in our comprehensive analysis of AI-related incidents from AIID, facilitating a nuanced understanding of how and where AI technologies are causing harm.

3.2 Taxonomy

We have adopted the comprehensive taxonomy developed by Dr. Ricardo Baeza-Yates as the classification framework for analyzing our dataset of incidents (see Figure 1). The taxonomy outlines the various classes of harm resulting from irresponsible AI use, forming a hierarchical structure where the main categories represent the broad types of harm, and the sub-classes offer a more detailed breakdown of these categories. For instance, 'Discrimination' serves as a main category that can be further dissected into 'Data Bias' or 'Algorithmic Bias.' Within 'Data Bias,' specifics such as gender, race, sexual orientation, economics, and others are identified as sub-subclass. This taxonomy will serve as the labels for the classes of harm in our classification.

Irresponsible AI Taxonomy

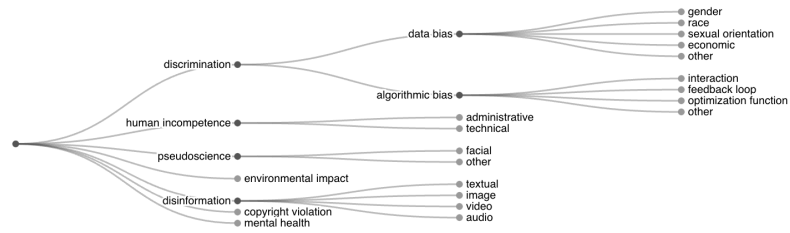


Fig. 1. Taxonomy for Classifying the Impacts of Irresponsible AI Utilization

Automation with Large Language Models (LLMs): We will program the LLMs to recognize and classify incidents according to the taxonomy. Their natural

language understanding will be leveraged to assess the incidents’ attributes and categorize them accordingly.

Human Classification: Simultaneously, a manual classification process will be conducted. This will involve a thorough review of each incident by the members of the team who will then classify them based on the detailed criteria set out in the taxonomy. This process will not only validate the automated classifications but also ensure that nuanced or ambiguous cases are properly categorized.

4 Curation Methodologies

4.1 Experiment Setup

GPT-4 Model For our project, we utilized the Chat Completions API ⁶ by OpenAI to leverage the capabilities of an advanced language model for classifying AI-related incidents documented in our dataset. The model of choice for our classification task was the GPT-4 model [11]. The GPT-4 language model represents an advancement over its predecessor, GPT-3, featuring significantly more parameters. GPT-4’s underlying technology is based on transformer architecture—specifically, attention mechanisms that emphasize different parts of the input data depending on the context. This model has been further refined with a technique called reinforcement learning from human feedback (RLHF), where human judgments help to fine-tune the model’s responses [5]. Although we specifically selected the GPT-4 model for this study, our future work may involve working with different models and evaluating their performances comprehensively.

Input Data From our dataset, we concentrated on eliminating redundant entries from the dataset. We merged all URLs linked to a specific incident under a singular incident ID, as the same incident might be reported by multiple news sources. This de-duplication process significantly streamlined the dataset, ensuring that each unique incident was catalogued as a single, exhaustive entry.

The result is a curated dataset devoid of redundancies, comprising 628 incidents⁷. Each incident is distinctly identified by an incident ID, consolidating all pertinent URLs under this unique identifier. This enables the classification to be based on the substantive content of the articles associated with each incident.

Next, we aggregated all available URL links, drawing reports from various news sources, social media posts, or official statement, for each incident within the AIID database. Each URL serves as a crucial source of information for specific incidents, offering additional context, details, and diverse perspective. Using Newspaper3k article scraping library⁸, we fetched the HTML content of each

⁶ <https://platform.openai.com/docs/api-reference/chat/create>

⁷ <https://observablehq.com/@irresponsible-ai/aiid-data-processing>

⁸ <https://newspaper.readthedocs.io/en/latest/>

web-page, parsed it to extract relevant article text, and filtered out noise or irrelevant content. The extracted text was concatenated into a single comprehensive string for each incident, which then fed into the LLM as the input. By providing a substantial corpus of text data from multiple sources, the model is able to gain exposure to diverse language patterns, enhancing its ability to accurately understand and classify incidents.

4.2 GPT-4 Configurations

A critical parameter in our use of the Chat Completions API was the setting of the *temperature*, which influences the model's output variability. The temperature parameter can range from 0 to 1, with lower values producing more deterministic and less varied outputs, and higher values encouraging more creativity and diversity in the generated text⁹. For the scope of our study, which centered around the straightforward tasks of text interpretation and classification, we adjusted the temperature setting from 0.5 to 0.3. This adjustment was made with the intention of minimizing the model's propensity for generating creative responses, thereby ensuring that the output remained closely aligned with the factual content of the articles and the classification taxonomy.

4.3 Prompting Strategies

Delimiters The first few prompts did not have clear delimiters so sometimes the model was unable to understand the prompt instructions. Later we used delimiters to clearly indicate distinct parts of inputs¹⁰. This approach significantly enhanced the model's ability to comprehend and respond to the prompts as intended.

```

1 ===== Start of Article Content =====
2                               {article_text}
3 ===== End of Article Content =====

```

Listing 1.1. Using Delimiters in Prompts

Persona-Based Prompting To guide the behavior of the language model, we employed a strategy of persona adoption, where the model is instructed to assume a specific role within its operational framework. For instance, we asked the model to operate as "a helpful assistant designed to classify news articles into specific categories," with the expected output format being JSON. This approach is designed to orient the model towards a particular behavior pattern¹¹

⁹ <https://platform.openai.com/docs/guides/text-generation/how-should-i-set-the-temperature-parameter>

¹⁰ <https://platform.openai.com/docs/guides/prompt-engineering/tactic-use-delimiters-to-clearly-indicate-distinct-parts-of-the-input>

¹¹ <https://platform.openai.com/docs/guides/text-generation/chat-completions-api>

Role Prompting for Enhanced Performance : The technique of role prompting leverages on the known sensitivity of Large Language Models (LLMs) to the construction of prompts. LLMs’ performance can significantly benefit from the specification of roles within prompts, improving by at least 20 percent over control prompts where no specific context is provided [20]. This sensitivity underscores the effectiveness of role prompting in enhancing the model’s task alignment and output accuracy.

We assigned an occupational role to our LLM. The initial prompt encouraged the LLM to conceptualize itself as part of a computer research institute tasked with categorizing incidents related to the irresponsible use of AI technology. Subsequently, we refined the prompting strategy by directly addressing our LLM as a computer researcher engaged in the same categorization task, eliminating the use of "imagine." This progression in prompt specificity was informed by evidence suggesting that LLMs respond with greater precision to direct role assignment, a reflection of their sensitivity to the nuances of prompt phrasing [20]. This iterative refinement in role prompting tailors the model’s focus and improves its classification performance.

Detailed Step Instructions In the refined prompts, we detailed the steps necessary for task completion by specifying the actions required to achieve a desired outcome¹². This modification improved GPT-4’s compliance with the given instructions, a significant enhancement over earlier versions of the prompts that lacked such explicit direction. Previously, the model occasionally deviated from the prescribed format or introduced its own categories, despite explicit instructions to adhere to the predefined taxonomy. The incorporation of clearly defined steps and rules in the prompts resulted in outputs that more accurately reflected the intended instructions, demonstrating the LLM’s increased ability to follow directions precisely.¹³

By adding step-by-step instructions refined the model’s performance, the model was able to adhere to our instructions more closely and maintained a more consistent output format, whereas in the former zero-shot prompt some of the output format deviated from our specifications.

4.4 Prompting Techniques

Zero-Shot and Few-Shot Prompting The work of Kojima et al. [8] has shown that these LLMs possess remarkable capabilities for reasoning without prior exposure to specific data (zero-shot reasoning) and for rapidly adapting to tasks when provided with a small number of targeted examples (few-shot learning). The complexity of the task and capability of the model could influence how useful it is to include some examples in the prompt, as this may offer extra context and direction, potentially improving model performance [5]. However,

¹² <https://platform.openai.com/docs/guides/prompt-engineering/tactic-specify-the-steps-required-to-complete-a-task>

¹³ See appendix A

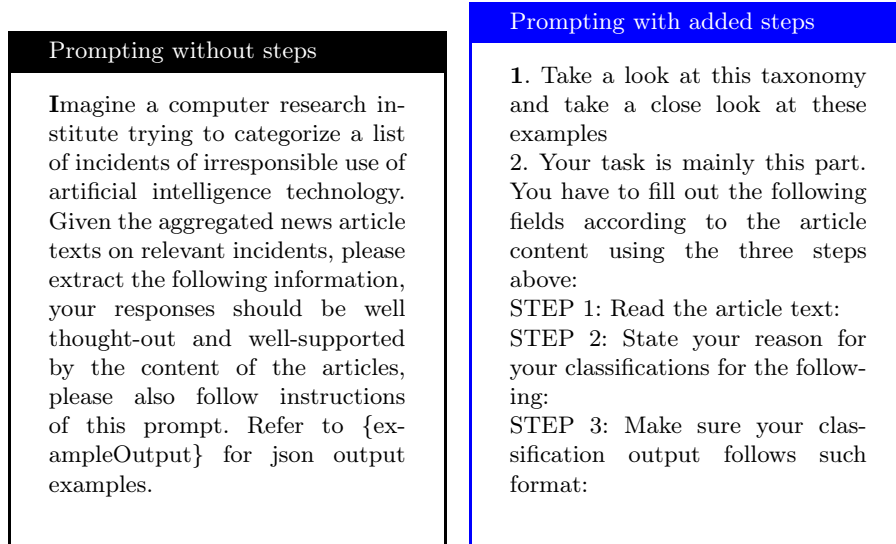


Fig. 2. Example of prompting without steps and prompting with added steps

it has also been shown that a well-designed zero-shot prompt, such as adding an intermediate reasoning step [8], could even lead to better performance and outperform few-shot prompts [5]. In our experiment, the few-shot learning with exemplars enabled GPT-4 to yield results that aligned more closely to the manual classifications¹⁴.

Chain of Thought Chain-of-Thought (CoT) prompting is introduced in the work of Wei et al. [16] that by eliciting an intermediate natural language reasoning step, the prompt can expand the set of diverse tasks that LLMs can perform successfully. Other research literature have shown that LLMs, combined with in-context learning (ICL) and chain-of-thought (CoT) prompting, are capable of various high-level reasoning tasks. [8] also introduced by simply adding "Let's think step by step" to a zero-shot prompt, the model outperforms standard zero-shot prompting. In our classification task, we set out with standard zero-shot prompting and also zero-shot prompting combined with Chain-of-Thought to elicit reasoning steps. We find that in the case of zero-shot prompting, there was not a big difference in results whether or not we added the reasoning steps. However, adding reasoning steps in few-shot learning prompts significantly improved the classification outcomes.

Tree of Thoughts The Tree of Thoughts (ToT) concept was introduced in Yao et al. [17], which enables LLMs to perform deliberate problem-solving by

¹⁴ See appendix B

maintaining and exploring a tree of thought steps. ToT frames problem-solving as a search over a tree of coherent language sequences ("thoughts") that serve as intermediate steps toward the solution, allowing exploration of multiple reasoning paths. This process bears resemblance to the thought and reasoning process typical to humans in that it encourages the LLM to plan and even backtrack in its reasoning to assess different paths of thoughts. It has been demonstrated that the ToT framework significantly outperforms standard input-output (IO) prompting and Chain-of-Thought prompting on certain reasoning and creative writing tasks [17].

Our classification task’s prompting methodology draws inspiration from these research insights. We present our observations from employing various prompting strategies with modifications and ensemble approaches that combine select techniques.

5 Output

5.1 Zero-Shot Learning

In the zero-shot learning results, GPT-4 demonstrates proficiency in pinpointing geographic and location information within articles with considerable accuracy. However, it occasionally misidentifies the city associated with a company. For instance, it incorrectly recognized Stanford University’s location as "Stanford" instead of the correct "Palo Alto" ¹⁵. The "Affected population" category permitted a broader range of responses, given that it wasn’t constrained to a predetermined set of labels. The model was successful in this field of classification by recognizing the impacted group of people. Additionally, GPT-4 was able to identify the hierarchy of harm classes along with their subcategories. Nevertheless, its categorization at times varied from those done manually ¹⁶.

5.2 Few-Shot Learning with Chain-of-Thought Technique

In our few-shot learning prompt, we included two examples in the prompt featuring article texts along with their example classification results. We also integrated the Chain of Thought (CoT) approach by providing a supporting rationale for each classification outcome. It has been shown that LLMs are capable of producing sequences of reasoning when provided with examples that showcase this process in the context of few-shot prompting [16]. The purpose of these examples is to guide the models in replicating a similar style of reasoning. By enumerating reasoning step examples, the model is guided to arrive at a particular answer by extracting information from the article. Eliciting reasoning steps from the model also allows it to correct its own answers [16]. Our findings reveal that GPT-4 was generally accurate in identifying geographic and location information, with the

¹⁵ <https://visit.stanford.edu/basics/>

¹⁶ See appendix E for the full result on the 17 incidents and appendix F for table comparisons

exception that the company city was sometimes misidentified. Figure 3 shows that location categories scored highest in accuracy when compared to manual classifications. GPT-4 also reliably classified "Affected population" and "Area of AI Application," which are categories that allow for a wider range of answers without fixed labels. With reasoning guidance provided in the examples, GPT-4's identification of the classes of harm were more consistent with manual classifications than with the results from zero-shot learning ¹⁷.

5.3 Tree of Thought Framework with Few-Shot Learning and Chain-of-Thought Technique

In addition to zero-shot and few-shot learning with CoT reasoning, we used an ensemble approach that combined few-shot learning and CoT reasoning with the Tree of Thoughts (ToT) framework. In the prompt, we asked GPT-4 to act as three experts in order to branch out into three reasoning and classification paths. The prompt guides the model to create a breadth-first search of the tree of probable classifications and ask the experts to vote on the classification is the most promising and well-supported by the article text and use that as the final result. Reference [17] mentions the "self-reflection" approach, wherein language models can provide feedback on the outputs they generate. The ToT framework leverages this capability to get the model to evaluate and select among its generated reasoning and outcome branches. We find GPT-4 performs generally accurate classification for geographic and location information, but again with exception of some misidentified company cities. The model also successfully identified the "Affected population" and "Area of AI Application" fields. The classifications for "Affected Population" provided by the ToT method included more detail compared to those from the Few-Shot Learning with Chain of Thought (CoT) approach. This variation could be attributed to the use of different sets of examples for the "Affected Population" category, with one set featuring concise, single-word examples and others offering descriptive examples of the affected demographic groups. It is uncertain which style or format the LLM will choose to emulate when provided with several different examples. When incorporating the Tree of Thoughts (ToT) into our prompts, the overall classification outcomes ¹⁸ were similar to those from few-shot learning with CoT, displaying only minor differences in certain areas.

6 Evaluation

Evaluation of the efficacy of prompt methods in LLMs are mainly divided to subjective and objective categories. Subjective evaluations generally rely on human evaluators to assess the quality of the generated content. Such evaluation methods are subjective and can be prone to inconsistencies [5]. In some cases, human

¹⁷ See appendix C for full result on 17 incidents

¹⁸ See appendix D

evaluation is the better way to evaluate the accuracy and quality of LLM generated results because evaluators (typically experts) can assess the results more comprehensively with more accurate feedback [4]. On the other hand, objective evaluation, or automatic evaluation methods include the use of machine learning algorithms to score the quality of the text generated by the LLMs. However, these automated metrics sometimes fail to fairly capture the full semantic considerations and cultural context in generated text as well as assessment results of human evaluators.

6.1 Objective Evaluation

We measured the accuracy score of the classification results from LLM from the manual results (used as the true labels). This evaluation metric focuses strictly on the literal exactness of the LLM’s outputs, but it does not consider semantic meanings or contextual nuances. Consequently, only single-word fields such as "State", "Company State", "City", tend to score high with this metric. We find that the few-shot learning with ToT and CoT prompt results had the highest matches for these fields, as shown in figure 3. This suggests that few-shot learning with ToT and CoT has the closest reasoning process to our manual classifications. Other fields such as "Affected population", "Area of AI application", and the classes of harm are expected to have lower scores due to the literal nature of this accuracy measurement. In addition to accuracy, we used BERTScore [19], which measures semantic similarity between pairs of text sequences. It provides a more nuanced evaluation compared to taking the accuracy score or other traditional evaluation metrics that rely more on matching the exactness of n-grams. The F1 scores are relatively high across all fields for each prompting technique as seen in figure 4, suggesting that the BERT model is able to capture the similarity between the evaluated text and the reference. The few-shot with ToT and CoT results tends to have the highest or near-highest F1 score in almost every field compared to the other sources, which may suggest that this approach is better at aligning with the manual evaluation. Although BERTScore is better at capturing semantic nuances, it still does not take on full contextual and semantic meanings, which could account for the lower F1 scores across the classification fields for some categories.

6.2 Subjective Evaluation

Since our classification task is domain-specific to the study of Irresponsible AI and its impact, we used human subjective evaluation from domain experts to assess the quality of GPT-4’s output. Building on the experimental framework outlined in Section 4, our evaluation method involved a blind review by two independent domain experts. Each expert assessed the classifications—both human and LLM-generated—across various dimensions such as geographic location, affected population, incident impact, application of AI use, and the classes of harm. The reviewers were unaware of which incidents have been classified by GPT-4 and which fields were manually classified by the team. The classifications

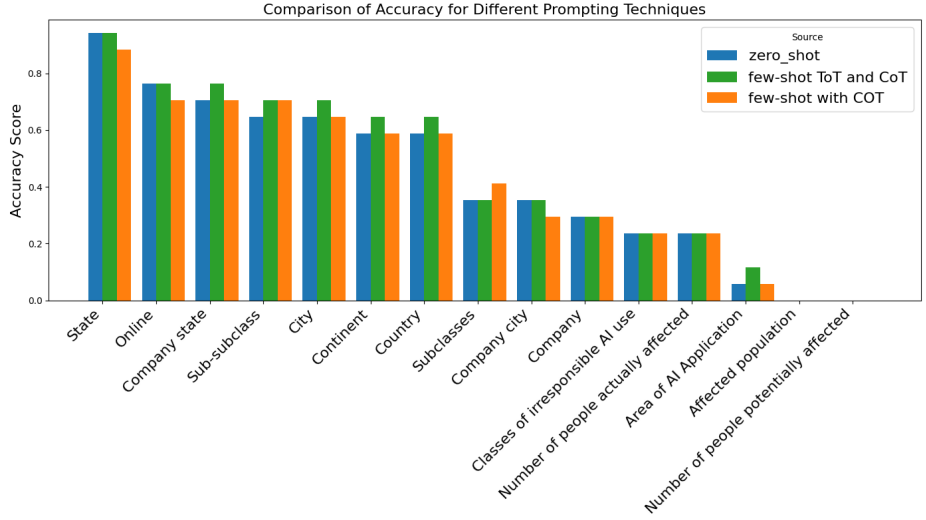


Fig. 3. Accuracy score of GPT-4 results. The accuracy scores here are measured by the exactness of literals using the manual classification results as the true label.

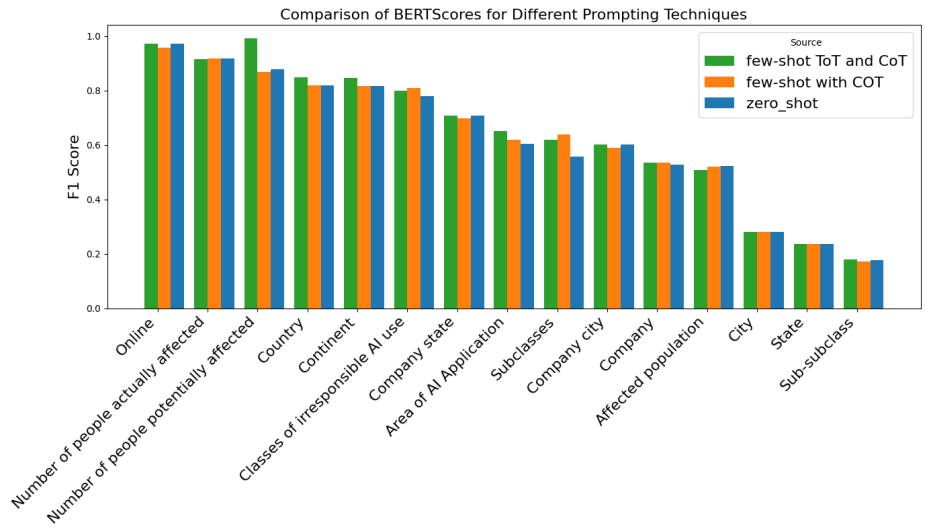


Fig. 4. F1 scores for different prompting technique results evaluated by the BERTScore evaluator

were scored on a scale from 1 to 10, based on accuracy of identification of each field. The average scores of the LLM classifications on 11 incidents achieved a slightly higher average score (at 7.73) than the average score of human manual classifications on 9 incidents (at 7.17). For both reviewers, the median scores for LLM classifications appear higher than for human classifications as can be seen in figures 5 and 6. This suggests that, on average, both reviewers rated LLM classifications more favorably than human ones. In both ratings, the variability of LLM classification ratings are generally wider than human classifications, indicating the reviewers both perceived some level of inconsistency in performance of the LLM results. Reviewer feedback on human classifications often cited issues such as typographical errors, incorrect identification of places, and lack of consensus regarding the labeling of harm classes, highlighting that human-generated results can also be prone to errors. On the other hand, the feedback for LLM classifications commonly mentioned the absence of predicted potential impacts, suggesting that the model may not have discerned certain information from the text or the information was not present for the model to make a judgment.

In some classification output, we have seen that GPT-4 is capable of producing very consistent output that is very close to human classification results¹⁹. Other cases show that the results varied with different prompts. The outcome of the assessment suggests that LLMs can outperform human efforts in classification tasks in some cases, as shown by the higher average scores. However they are not without limitations, such as failing to predict certain outcomes from information gleaned from article texts. The variability and errors in human classifications also demonstrate the complexities and challenges that we often see in the manual categorization processes. This comparison suggests that while LLMs show promise as efficient tools for classification, combining their strengths with human oversight could offer a more robust approach.

It is worth noting that this experiment represents a small subset of the 628 incidents documented in the AIID, meaning the findings might not sufficiently reflect the broader performance capabilities of LLM and human classifications. We plan to experiment in small batches before we proceed with a suitable approach that could successfully classify all the incidents. Additionally, the variability in evaluations among expert reviewers might indicate the need for a more robust evaluation design, since it is inevitable for subjective evaluations to have wide variability because of differences in opinions.

7 Conclusions

Our study’s focus was to enhance classification methods for the numerous AI incidents documented in the AI Incidents Database, aiming to improve understanding and management of AI-related harm.

The study leveraged the advanced capabilities of large language models (LLMs) like GPT-4 to classify incidents more efficiently than manual methods. Our findings indicated that while LLMs can sometimes outperform human classification

¹⁹ See appendix F

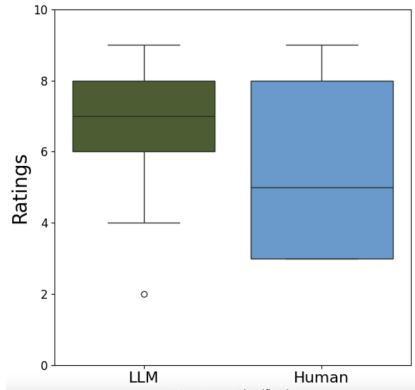


Fig. 5. Reviewer 1 ratings: distribution of reviewer 1's ratings on incidents classification from a blind review

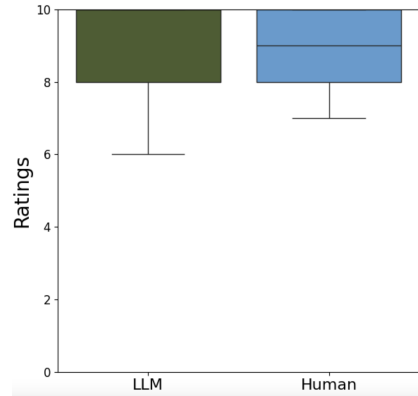


Fig. 6. Reviewer 2 ratings: distribution of reviewer 2's ratings on incidents classification from a blind review

in terms of average scoring, they are not infallible, occasionally misidentifying locations or missing potential impacts. This underscores the need for human oversight to ensure the highest quality and reliability in classification.

The investigation into prompt engineering, spanning zero-shot to few-shot learning with Chain of Thought and Tree of Thoughts frameworks, revealed nuances in model performance. Combining Tree of Thoughts method with few-shot prompts and explicit reasoning steps offered results most closely aligned with manual classifications. There are still known limitations to LLMs such as GPT-4. Most notably, like its predecessors, it still lacks complete reliability, sometimes producing fabricated information (known as "hallucinations") or committing logical errors [11]. Hallucinations occur because the model might not have enough supporting evidence in its training data for its responses, or it could generalize patterns too broadly to produce coherent output [7]. In some instances, we observed GPT-4 exhibiting this behavior where it could not accurately identify specific classification fields from the text of an article. For instance, it mistakenly identified the city of Stanford University as "Stanford" rather than "Palo Alto." In another case, it incorrectly classified "Silicon Valley" as the city where Chai Research ²⁰ is headquartered, when it is actually located in Palo Alto ²¹.

7.1 Future Work

Moving forward, we plan to extend our methodology to the entire corpus of incidents in AIID and apply a more robust and effective approach to classify all the news incidents with high accuracy, which will be corroborated by experts in the field. Our aim is to harness these classifications to create data visualizations

²⁰ www.chai-research.com

²¹ See appendix C

that will contribute to the ongoing discourse on Responsible AI, providing greater insight into AI-related harm trends and influencing the development of more responsible AI systems.

This work serves as a step toward understanding and documenting AI failures and establishing a robust, automated classification framework to handle the complexity of AI incident reports. It highlights both the capabilities and limitations of LLMs in this domain and points to a future where collaboration between human expertise and AI could offer robust solutions for managing AI’s societal impacts.

Acknowledgments. We thank Khoury West Coast Research Program for funding this research and the team behind the AI Incidents Database (AIID) for publishing their data.

References

1. Baeza-Yates, R.: LECTURE HELD AT THE ACADEMIA EUROPAEA BUILDING BRIDGES CONFERENCE 2022: An Introduction to Responsible AI. *European Review* **31**(4), 406–421 (8 2023). <https://doi.org/10.1017/S1062798723000145>, <https://www.cambridge.org/core/journals/european-review/article/lecture-held-at-the-academia-europaea-building-bridges-conference-2022/C61F82D2FF64FFB933282F5E65CDB22A>
2. Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D.M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., Mccandlish, S., Radford, A., Sutskever, I., Openai, D.A.: *Language Models are Few-Shot Learners* (2020)
3. Burema, D., Debowski-Weimann, N., Von Janowski, A., Grabowski, J., Maftai, M., Jacobs, M., Van Der Smagt, P., Benbouzid, D.: A sector-based approach to AI ethics: Understanding ethical issues of AI-related incidents within their sectoral context. *AIES 2023 - Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society* pp. 705–714 (8 2023). <https://doi.org/10.1145/3600211.3604680>, <https://dl.acm.org/doi/10.1145/3600211.3604680>
4. Chang, Y., Wang, X.U., Yi, X., Wang, Y., Ye, W., Yu, P.S., Chang, Y., Wang, X., Wu, Y., Yi, X., Xie, X., Yang, .L., Wang, C., Zhang, Y., Wang, Y., Ye, W., Wang, X., Wang, J., Wu, Y., Yang, L., Zhu, K., Chen, H., Wang, C., Zhang, Y., Chang, Y., Yang, Q., Xie, X.: A Survey on Evaluation of Large Language Models. *ACM Trans. Intell. Syst. Technol* **15**(3), 39 (2024). <https://doi.org/10.1145/3641289>, <https://github.com/MLGroupJLU/LLM-eval-survey>
5. Chen, B., Zhang, Z., Langrené, N., Zhu, S.: Unleashing the potential of prompt engineering in Large Language Models: a comprehensive review (10 2023), <http://arxiv.org/abs/2310.14735>

6. Chen, S., Li, Y., Lu, S., Van, H., Aerts, H.J.W.L., Savova, G.K., Bitterman, D.S.: Evaluating the ChatGPT family of models for biomedical reasoning and classification. *Journal of the American Medical Informatics Association* (4 2024). <https://doi.org/10.1093/JAMIA/OCAD256>
7. Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., Ishii, E., Bang, Y.J., Madotto, A., Fung, P.: Survey of Hallucination in Natural Language Generation. *ACM Computing Surveys* **55**(12) (3 2023). <https://doi.org/10.1145/3571730>, <https://dl.acm.org/doi/10.1145/3571730>
8. Kojima, T., Gu, S.S., Reid, M., Matsuo, Y., Iwasawa, Y.: Large Language Models are Zero-Shot Reasoners (5 2022), <http://arxiv.org/abs/2205.11916>
9. Liu, P., Yuan, W., Fu, J., Jiang, Z., Hayashi, H., Neubig, G.: Pre-train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing (2021)
10. Mcgregor, S.: Preventing Repeated Real World AI Failures by Cataloging Incidents: The AI Incident Database (2021), www.aaai.org
11. OpenAI, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F.L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., Avila, R., Babuschkin, I., Balaji, S., Balcom, V., Balthescu, P., Bao, H., Bavarian, M., Belgum, J., Bello, I., Berdine, J., Bernadett-Shapiro, G., Berner, C., Bogdonoff, L., Boiko, O., Boyd, M., Brakman, A.L., Brockman, G., Brooks, T., Brundage, M., Button, K., Cai, T., Campbell, R., Cann, A., Carey, B., Carlson, C., Carmichael, R., Chan, B., Chang, C., Chantzis, F., Chen, D., Chen, S., Chen, R., Chen, J., Chen, M., Chess, B., Cho, C., Chu, C., Chung, H.W., Cummings, D., Currier, J., Dai, Y., Decareaux, C., Degry, T., Deutsch, N., Deville, D., Dhar, A., Dohan, D., Dowling, S., Dunning, S., Ecoffet, A., Eleti, A., Eloundou, T., Farhi, D., Fedus, L., Felix, N., Fishman, S.P., Forte, J., Fulford, I., Gao, L., Georges, E., Gibson, C., Goel, V., Gogineni, T., Goh, G., Gontijo-Lopes, R., Gordon, J., Grafstein, M., Gray, S., Greene, R., Gross, J., Gu, S.S., Guo, Y., Hallacy, C., Han, J., Harris, J., He, Y., Heaton, M., Heidecke, J., Hesse, C., Hickey, A., Hickey, W., Hoeschele, P., Houghton, B., Hsu, K., Hu, S., Hu, X., Huizinga, J., Jain, S., Jain, S., Jang, J., Jiang, A., Jiang, R., Jin, H., Jin, D., Jomoto, S., Jonn, B., Jun, H., Kaftan, T., Kaiser, , Kamali, A., Kanitscheider, I., Keskar, N.S., Khan, T., Kilpatrick, L., Kim, J.W., Kim, C., Kim, Y., Kirchner, J.H., Kiros, J., Knight, M., Kokotajlo, D., Kondraciuk, , Kondrich, A., Konstantinidis, A., Kosic, K., Krueger, G., Kuo, V., Lampe, M., Lan, I., Lee, T., Leike, J., Leung, J., Levy, D., Li, C.M., Lim, R., Lin, M., Lin, S., Litwin, M., Lopez, T., Lowe, R., Lue, P., Makanju, A., Malfacini, K., Manning, S., Markov, T., Markovski, Y., Martin, B., Mayer, K., Mayne, A., McGrew, B., McKinney, S.M., McLeavey, C., McMillan, P., McNeil, J., Medina, D., Mehta, A., Menick, J., Metz, L., Mishchenko, A., Mishkin, P., Monaco, V., Morikawa, E., Mossing, D., Mu, T., Murati, M., Murk, O., Mély, D., Nair, A., Nakano, R., Nayak, R., Neelakantan, A., Ngo, R., Noh, H., Ouyang, L., O'Keefe, C., Pachocki, J., Paino, A., Palermio, J., Pantuliano, A., Parascandolo, G., Parish, J., Parparita, E., Passos, A., Pavlov, M., Peng, A., Perelman, A., Peres, F.d.A.B., Petrov, M., Pinto, H.P.d.O., Michael, Pokorny, Pokrass, M., Pong, V.H., Powell, T., Power, A., Power, B., Proehl, E., Puri, R., Radford, A., Rae, J., Ramesh, A., Raymond, C., Real, F., Rimbach, K., Ross, C., Rotsted, B., Roussez, H., Ryder, N., Saltarelli, M., Sanders, T., Santurkar, S., Sastry, G., Schmidt, H., Schnurr, D., Schulman, J., Selsam, D., Sheppard, K., Sherbakov, T., Shieh, J., Shoker, S., Shyam, P., Sidor, S., Sigler, E., Simens, M., Sitkin, J., Slama, K., Sohl, I., Sokolowsky, B., Song, Y., Staudacher, N., Such, F.P., Summers, N., Sutskever, I., Tang, J., Tezak, N., Thompson, M.B., Tillet,

- P., Tootoonchian, A., Tseng, E., Tuggle, P., Turley, N., Tworek, J., Uribe, J.F.C., Vallone, A., Vijayvergiya, A., Voss, C., Wainwright, C., Wang, J.J., Wang, A., Wang, B., Ward, J., Wei, J., Weinmann, C., Welihinda, A., Welinder, P., Weng, J., Weng, L., Wiethoff, M., Willner, D., Winter, C., Wolrich, S., Wong, H., Workman, L., Wu, S., Wu, J., Wu, M., Xiao, K., Xu, T., Yoo, S., Yu, K., Yuan, Q., Zaremba, W., Zellers, R., Zhang, C., Zhang, M., Zhao, S., Zheng, T., Zhuang, J., Zhuk, W., Zoph, B.: GPT-4 Technical Report (3 2023), <http://arxiv.org/abs/2303.08774>
12. Pittaras, N., McGregor, S.: A taxonomic system for failure cause analysis of open source AI incidents (11 2022), <https://arxiv.org/abs/2211.07280v1>
 13. Saparov, A., Pang, R.Y., Padmakumar, V., Joshi, N., Kazemi, S.M., Kim, N., He, H.: Testing the General Deductive Reasoning Capacity of Large Language Models Using OOD Examples (5 2023), <http://arxiv.org/abs/2305.15269>
 14. Turri, V., Dzombak, R.: Why We Need to Know More: Exploring the State of AI Incident Documentation Practices. AIES 2023 - Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society pp. 576–583 (8 2023). <https://doi.org/10.1145/3600211.3604700>, <https://dl.acm.org/doi/10.1145/3600211.3604700>
 15. Wei, F., Keeling, R., Huber-Fliflet, N., Zhang, J., Dabrowski, A., Yang, J., Mao, Q., Qin, H.: Empirical Study of LLM Fine-Tuning for Text Classification in Legal Document Review. Proceedings - 2023 IEEE International Conference on Big Data, BigData 2023 pp. 2786–2792 (2023). <https://doi.org/10.1109/BIGDATA59044.2023.10386911>
 16. Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q., Zhou, D.: Chain-of-Thought Prompting Elicits Reasoning in Large Language Models (1 2022), <http://arxiv.org/abs/2201.11903>
 17. Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T.L., Cao, Y., Narasimhan, K.: Tree of Thoughts: Deliberate Problem Solving with Large Language Models (5 2023), <http://arxiv.org/abs/2305.10601>
 18. Zhang, P., Jaipersaud, B., Ba, J., Petersen, A., Zhang, L., Zhang, M.R.: Classifying Course Discussion Board Questions using LLMs. Annual Conference on Innovation and Technology in Computer Science Education, ITiCSE **2**, 658 (6 2023). <https://doi.org/10.1145/3587103.3594202>
 19. Zhang, T., Kishore, V., Wu, F., Weinberger, K.Q., Artzi, Y.: BERTScore: Evaluating Text Generation with BERT (4 2019), <http://arxiv.org/abs/1904.09675>
 20. Zheng, M., Pei, J., Jurgens, D.: Is "A Helpful Assistant" the Best Role for Large Language Models? A Systematic Evaluation of Social Roles in System Prompts. arxiv.orgM Zheng, J Pei, D JurgensarXiv preprint arXiv:2311.10054, 2023•[arxiv.org](http://arxiv.org/abs/2311.10054) (11 2023), <http://arxiv.org/abs/2311.10054>

Appendices

A

Prompt with clear steps

```

1 STEP 1: Read the article text:
2 ===== Start of Article Content =====
3 {article_text}
4 ===== End of Article Content =====
5 STEP 2: State your reason for your classifications for the
   following:
6 =====Classification Fields=====
7 - Country (output "Worldwide" if the incident happened across
   multiple countries):
8 - State (if not applicable leave blank):
9 - City (if not applicable leave blank):
10 - Continent (output "Worldwide" if the incident happened
   across multiple countries):
11 - Company (i.e. the company that developed the technology
   involved in this incident):
12 - Company city (the city where the headquarters of this
   company is located. If the company recently moved
   headquarters, please use the location of the new
   headquarter):
13 - Company state (the state of the company city, if applicable
   , if not leave blank):
14 - Affected population (think about which groups of people are
   directly affected by the incident in the article.)
15 - Number of people actually affected (let's check the number
   of people directly affected according to the article.
   Give a total number. If unknown output 'Unknown'):
16 - Number of people potentially affected (let's check the
   article text to see if this information is provided or
   suggested, if not you may ouput 'Unknown'):
17 - Classes of irresponsible AI use (please follow the rules
   and refer to this taxonomy:
18 “‘taxonomy classes
19     {taxa.classes}
20 “‘
21 Rule1: There could be more than one classes the article
   classifies as.
22 Rule2: DO NOT create your own class, adhere strictly to the
   provided list.

```

```

23 - Subclasses (please follow the rules and refer to this
      taxonomy structure '<class>:[<subclass>]'):
24   "'taxonomy subclasses
25     {taxa.subclasses}
26   "'
27 =====Classification Fields=====
28 STEP 3: Make sure your classification output follows such
      format:
29   "'THIS IS AN EXAMPLE"'
30 {example_output_id_6}
31   "'END OF EXAMPLE"'

```

Listing 2.1. Detailed Instructions with Steps in Prompts

B

Manual classification results

The following JSON of 17 incidents are the results of our manual classifications that was discussed and decided among three people after reading the articles of the incidents. This was used as a benchmark for comparison with the LLM generated results.

```

1  [
2  [
3    {
4      "id": 1,
5      "Country": "Worldwide",
6      "State": "",
7      "City": "",
8      "Continent": "Worldwide",
9      "Company": "Google LLC",
10     "Company city": "Mountain View",
11     "Company state": "California",
12     "Affected population": "Children on Youtube",
13     "Number of people actually affected": "unknown",
14     "Number of people potentially affected": "unknown",
15     "Classes of irresponsible AI use": "human incompetence,
16     mental health,copyright",
17     "Subclasses": "technical,administrative",
18     "Sub-subclass": "",
19     "Area of AI Application": "content filtering",
20     "Online": "yes"
21   },
22   {
23     "id": 2,
24     "Country": "United States",
25     "State": "New Jersey",
26     "City": "Robbinsville",
27     "Continent": "North America",

```

```
27     "Company": "Amazon",
28     "Company city": "Seattle",
29     "Company state": "Washington",
30     "Affected population": "Amazon Workers and Families",
31     "Number of people actually affected": 54,
32     "Number of people potentially affected": 80,
33     "Classes of irresponsible AI use": "human incompetence",
34     "Subclasses": "technical",
35     "Sub-subclass": "",
36     "Area of AI Application": "robotics",
37     "Online": "no"
38 },
39 {
40     "id": 3,
41     "Country": "Indonesia",
42     "State": "",
43     "City": "Bali",
44     "Continent": "Southeast Asia",
45     "Company": "Boeing",
46     "Company city": "Arlington",
47     "Company state": "Virginia",
48     "Affected population": "Lion Air Jet Passengers and
49     Families",
50     "Number of people actually affected": 189,
51     "Number of people potentially affected": "unknown",
52     "Classes of irresponsible AI use": "human incompetence",
53     "Subclasses": "technical",
54     "Sub-subclass": "",
55     "Area of AI Application": "airspeed indicator",
56     "Online": "no"
57 },
58 {
59     "id": 4,
60     "Country": "United States",
61     "State": "Arizona",
62     "City": "Tempe",
63     "Continent": "North America",
64     "Company": "Uber",
65     "Company city": "San Francisco",
66     "Company state": "California",
67     "Affected population": "Pedestrians",
68     "Number of people actually affected": 1,
69     "Number of people potentially affected": "unknown",
70     "Classes of irresponsible AI use": "human incompetence",
71     "Subclasses": "technical",
72     "Sub-subclass": "",
73     "Area of AI Application": "autonomous driving",
74     "Online": "no"
75 }
```

```

76     "id": 5,
77     "Country": "United States",
78     "State": "",
79     "City": "",
80     "Continent": "North America",
81     "Company": "",
82     "Company city": "",
83     "Company state": "",
84     "Affected population": "Patients with robotic procedures",
85     "Number of people actually affected": 1535,
86     "Number of people potentially affected": "unknown",
87     "Classes of irresponsible AI use": "human incompetence",
88     "Subclasses": "technical",
89     "Sub-subclass": "",
90     "Area of AI Application": "robotics ",
91     "Online": "no"
92 },
93 {
94     "id": 6,
95     "Country": "United States",
96     "State": "",
97     "City": "",
98     "Continent": "North America",
99     "Company": "Microsoft",
100    "Company city": "Seattle",
101    "Company state": "Washington",
102    "Affected population": "Tweeter users",
103    "Number of people actually affected": "unknown",
104    "Number of people potentially affected": "unknown",
105    "Classes of irresponsible AI use": "discrimination,
106    disinformation",
107    "Subclasses": "data bias,algorithmic bias,textual",
108    "Sub-subclass": "race,gender,feedback loop",
109    "Area of AI Application": "chatbot",
110    "Online": "yes"
111 },
112 {
113     "id": 9,
114     "Country": "United States",
115     "State": "New York",
116     "City": "New York",
117     "Continent": "North America",
118     "Company": "",
119     "Company city": "New York",
120     "Company state": "New York",
121     "Affected population": "NYC teachers",
122     "Number of people actually affected": "unknown",
123     "Number of people potentially affected": "unknown",

```

```

123     "Classes of irresponsible AI use": "discrimination,
124     disinformation",
125     "Subclasses": "data bias,algorithmic bias,textual",
126     "Sub-subclass": "other,other",
127     "Area of AI Application": "data prediction",
128     "Online": "no"
129 },
130 {
131     "id": 10,
132     "Country": "United States",
133     "State": "",
134     "City": "",
135     "Continent": "North America",
136     "Company": "UKG",
137     "Company city": "Lowell",
138     "Company state": "Massachusetts",
139     "Affected population": "Starbucks braristas",
140     "Number of people actually affected": "130,000",
141     "Number of people potentially affected": "unknown",
142     "Classes of irresponsible AI use": "human incompetence",
143     "Subclasses": "technical",
144     "Sub-subclass": "",
145     "Area of AI Application": "job scheduling",
146     "Online": "no"
147 },
148 {
149     "id": 11,
150     "Country": "United States",
151     "State": "Florida",
152     "City": "Broward County",
153     "Continent": "North America",
154     "Company": "Northpointe Bank",
155     "Company city": "Grand Rapids",
156     "Company state": "Michigan",
157     "Affected population": "Defendants in Broward County,",
158     "Number of people actually affected": "unknown",
159     "Number of people potentially affected": "unknown",
160     "Classes of irresponsible AI use": "discrimination,
161     pseudoscience",
162     "Subclasses": "data bias,facial",
163     "Sub-subclass": "race",
164     "Area of AI Application": "predictive algorithms - COMPAS
165     ",
166     "Online": "no"
167 },
168 {
169     "id": 13,
170     "Country": "United States",
171     "State": "",
172     "City": "",

```



```

170     "Continent": "North America",
171     "Company": "Google LLC",
172     "Company city": "Mountain View",
173     "Company state": "California",
174     "Affected population": "Minorities and groups like woman,
    gay, etc.",
175     "Number of people actually affected": "unknown",
176     "Number of people potentially affected": "unknown",
177     "Classes of irresponsible AI use": "discrimination",
178     "Subclasses": "data bias",
179     "Sub-subclass": "gender,race,sexual orientation,other",
180     "Area of AI Application": "predictive algorithm",
181     "Online": "yes"
182 },
183 {
184     "id": 14,
185     "Country": "United States",
186     "State": "",
187     "City": "",
188     "Continent": "North America",
189     "Company": "Google LLC",
190     "Company city": "Mountain View",
191     "Company state": "California",
192     "Affected population": "Minorities and groups like woman,
    gay, etc.",
193     "Number of people actually affected": "unknown",
194     "Number of people potentially affected": "unknown",
195     "Classes of irresponsible AI use": "discrimination",
196     "Subclasses": "algorithmic bias,data bias",
197     "Sub-subclass": "feedback loop,race,gender",
198     "Area of AI Application": "NLP analysis",
199     "Online": "yes"
200 },
201 {
202     "id": 451,
203     "Country": "United States",
204     "State": "",
205     "City": "",
206     "Continent": "North America",
207     "Company": "Stability AI",
208     "Company city": "Houston",
209     "Company state": "Texas",
210     "Affected population": "Getty Company",
211     "Number of people actually affected": "unknown",
212     "Number of people potentially affected": "unknown",
213     "Classes of irresponsible AI use": "copyright violation",
214     "Subclasses": "",
215     "Sub-subclass": "",
216     "Area of AI Application": "image generation",
217     "Online": "yes"

```

```
218 },
219 {
220   "id": 382,
221   "Country": "United Kingdom",
222   "State": "",
223   "City": "London",
224   "Continent": "Europe",
225   "Company": "Facebook",
226   "Company city": "Menlo Park",
227   "Company state": "California",
228   "Affected population": "instagram users",
229   "Number of people actually affected": 1,
230   "Number of people potentially affected": "unknown",
231   "Classes of irresponsible AI use": "mental health",
232   "Subclasses": "",
233   "Sub-subclass": "",
234   "Area of AI Application": "social media",
235   "Online": "yes"
236 },
237 {
238   "id": 505,
239   "Country": "Belgium",
240   "State": "",
241   "City": "",
242   "Continent": "Europe",
243   "Company": "EleutherAI",
244   "Company city": "New York",
245   "Company state": "New York",
246   "Affected population": "Chai app users",
247   "Number of people actually affected": 1,
248   "Number of people potentially affected": "unknown",
249   "Classes of irresponsible AI use": "mental health",
250   "Subclasses": "",
251   "Sub-subclass": "",
252   "Area of AI Application": "chatbot",
253   "Online": "yes"
254 },
255 {
256   "id": 167,
257   "Country": "United States",
258   "State": "",
259   "City": "",
260   "Continent": "North America",
261   "Company": "Standford researchers",
262   "Company city": "Palo Alto",
263   "Company state": "California",
264   "Affected population": "LGBT group",
265   "Number of people actually affected": "unknown",
266   "Number of people potentially affected": "unknown",
```

```

267   "Classes of irresponsible AI use": "pseudoscience,
      discrimination",
268   "Subclasses": "facial,data bias",
269   "Sub-subclass": "sexual orientation",
270   "Area of AI Application": "behavioral modeling",
271   "Online": "yes"
272 },
273 {
274   "id": 0,
275   "Country": "United States",
276   "State": "",
277   "City": "",
278   "Continent": "North America",
279   "Company": "",
280   "Company city": "",
281   "Company state": "",
282   "Affected population": "General Human",
283   "Number of people actually affected": "unknown",
284   "Number of people potentially affected": "unknown",
285   "Classes of irresponsible AI use": "environmental impact
      ",
286   "Subclasses": "",
287   "Sub-subclass": "",
288   "Area of AI Application": "AI training",
289   "Online": "no"
290 },
291 {
292   "id": 39,
293   "Country": "United States",
294   "State": "",
295   "City": "",
296   "Continent": "North America",
297   "Company": "",
298   "Company city": "",
299   "Company state": "",
300   "Affected population": "online audience",
301   "Number of people actually affected": "unknown",
302   "Number of people potentially affected": "unknown",
303   "Classes of irresponsible AI use": "disinformation",
304   "Subclasses": "video",
305   "Sub-subclass": "",
306   "Area of AI Application": "deepfake video generation",
307   "Online": "yes"
308 }
309 ]

```

Listing 2.2. Manual classification results of 17 incidents handpicked by the team. The 17 incidents covered all the classes of harm of our taxonomy.

C

Few-Shot Learning Classification Results

These are the classification results of 17 incidents from GPT-4 with two-shot prompting with Chain-of-Thought (CoT) technique.

```

1 {
2   "1": {
3     "Country": "Global",
4     "State": "N/A",
5     "City": "N/A",
6     "Continent": "Global",
7     "Company": "Google LLC",
8     "Company city": "Mountain View, CA",
9     "Affected population": "Children",
10    "Number of people actually affected": "Unknown",
11    "Number of people potentially affected": "Unknown",
12    "Class of irresponsible AI use": [
13      "Disinformation",
14      "Human Incompetence",
15      "Mental Health",
16      "Other"
17    ],
18    "Subclasses": [
19      "Disinformation-> Video",
20      "Human Incompetence-> Technical",
21      "Mental Health-> Psychological Impact",
22      "Other-> Inappropriate Content"
23    ],
24    "Sub-subclass": [],
25    "Area of AI Application": "Video Sharing/Streaming",
26    "Online": "Yes"
27  },
28  "2": {
29    "Country": "United States",
30    "State": "New Jersey",
31    "City": "Robbinsville",
32    "Continent": "North America",
33    "Company": "Amazon",
34    "Company city": "Seattle, WA",
35    "Affected population": "Amazon workers",
36    "Number of people actually affected": "24",
37    "Number of people potentially affected": "54",
38    "Class of irresponsible AI use": [
39      "Human Incompetence"
40    ],
41    "Subclasses": [
42      "Human Incompetence-> Technical"
43    ],
44    "Sub-subclass": [],

```

```

45     "Area of AI Application": "warehouse automation,
robotics",
46     "Online": "No"
47 },
48 "3": {
49     "Country": "Indonesia",
50     "State": "",
51     "City": "Jakarta",
52     "Continent": "Asia",
53     "Company": "Boeing",
54     "Company city": "Chicago, IL",
55     "Affected population": "Airline passengers and crew",
56     "Number of people actually affected": "189",
57     "Number of people potentially affected": "Unknown",
58     "Class of irresponsible AI use": [
59         "Human Incompetence"
60     ],
61     "Subclasses": [
62         "Human Incompetence-> Technical"
63     ],
64     "Sub-subclass": [],
65     "Area of AI Application": "Aviation safety systems",
66     "Online": "No"
67 },
68 "4": {
69     "Country": "United States",
70     "State": "Arizona",
71     "City": "Tempe",
72     "Continent": "North America",
73     "Company": "Uber",
74     "Company city": "San Francisco, CA",
75     "Affected population": "Pedestrians",
76     "Number of people actually affected": "1",
77     "Number of people potentially affected": "Unknown",
78     "Class of irresponsible AI use": [
79         "Human Incompetence",
80         "Technical Failure"
81     ],
82     "Subclasses": [
83         "Human Incompetence-> Monitoring Failure",
84         "Technical Failure-> Detection Failure"
85     ],
86     "Sub-subclass": [],
87     "Area of AI Application": "autonomous vehicles",
88     "Online": "No"
89 },
90 "5": {
91     "Country": "United States",
92     "State": "Not specified",
93     "City": "Not specified",

```

```

94     "Continent": "North America",
95     "Company": "Not specified",
96     "Company city": "Not specified",
97     "Affected population": "Patients undergoing robotic
surgery",
98     "Number of people actually affected": "144 deaths,
1391 injuries",
99     "Number of people potentially affected": "More than
1.7 million (number of robotic procedures carried out
between 2007 and 2013)",
100     "Class of irresponsible AI use": [
101         "Human Incompetence",
102         "Technical Failure"
103     ],
104     "Subclasses": [
105         "Human Incompetence-> Operator error",
106         "Technical Failure-> Device malfunction",
107         "Technical Failure-> System error"
108     ],
109     "Sub-subclass": [
110         "Device malfunction-> Electrical sparking",
111         "Device malfunction-> Uncontrolled movement",
112         "Device malfunction-> Broken pieces",
113         "System error-> Loss of video feed"
114     ],
115     "Area of AI Application": "robotic surgery",
116     "Online": "No"
117 },
118 "6": {
119     "Country": "Global",
120     "State": "Not specified",
121     "City": "Not specified",
122     "Continent": "Not specified",
123     "Company": "Microsoft",
124     "Company city": "Redmond, WA",
125     "Affected population": "Twitter Users, Online
Community",
126     "Number of people actually affected": "Unknown",
127     "Number of people potentially affected": "Millions (
Twitter user base and wider online community)",
128     "Class of irresponsible AI use": [
129         "Human Incompetence",
130         "Disinformation"
131     ],
132     "Subclasses": [
133         "Human Incompetence-> Lack of foresight",
134         "Disinformation-> Racial",
135         "Disinformation-> Political",
136         "Disinformation-> Gender"
137     ],

```

```

138     "Sub-subclass": [],
139     "Area of AI Application": "Social media interaction,
Conversational AI",
140     "Online": "Yes"
141 },
142 "7": {
143     "Country": "Global",
144     "State": "Global",
145     "City": "Global",
146     "Continent": "Global",
147     "Company": "Wikipedia",
148     "Company city": "San Francisco, CA",
149     "Affected population": "Wikipedia Users",
150     "Number of people actually affected": "Unknown",
151     "Number of people potentially affected": "Millions",
152     "Class of irresponsible AI use": [
153         "Human Incompetence"
154     ],
155     "Subclasses": [
156         "Human Incompetence-> Technical"
157     ],
158     "Sub-subclass": [],
159     "Area of AI Application": "Online encyclopedia
editing",
160     "Online": "Yes"
161 },
162 "8": {
163     "Country": "United States",
164     "State": "California",
165     "City": "San Francisco",
166     "Continent": "North America",
167     "Company": "Uber Technologies Inc.",
168     "Company city": "San Francisco, CA",
169     "Affected population": "General Public",
170     "Number of people actually affected": "Unknown",
171     "Number of people potentially affected": "Residents
and visitors of San Francisco",
172     "Class of irresponsible AI use": [
173         "Human Incompetence",
174         "Other"
175     ],
176     "Subclasses": [
177         "Human Incompetence-> Technical",
178         "Other-> Safety"
179     ],
180     "Sub-subclass": [],
181     "Area of AI Application": "autonomous vehicles",
182     "Online": "No"
183 },
184 "9": {

```

```

185     "Country": "United States",
186     "State": "New York",
187     "City": "New York City",
188     "Continent": "North America",
189     "Company": "New York City Department of Education",
190     "Company city": "New York, NY",
191     "Affected population": "Teachers",
192     "Number of people actually affected": "12000",
193     "Number of people potentially affected": "Unknown",
194     "Class of irresponsible AI use": [
195         "Human Incompetence",
196         "Other"
197     ],
198     "Subclasses": [
199         "Human Incompetence-> Technical",
200         "Other-> Misuse of Data"
201     ],
202     "Sub-subclass": [],
203     "Area of AI Application": "Educational assessment",
204     "Online": "Yes"
205 },
206 "10": {
207     "Country": "United States",
208     "State": "Various",
209     "City": "Various",
210     "Continent": "North America",
211     "Company": "Starbucks",
212     "Company city": "Seattle, WA",
213     "Affected population": "Starbucks Workers",
214     "Number of people actually affected": "Unknown",
215     "Number of people potentially affected": "130,000 (
total number of Starbucks baristas in the U.S.)",
216     "Class of irresponsible AI use": [
217         "Human Incompetence"
218     ],
219     "Subclasses": [
220         "Human Incompetence-> Technical"
221     ],
222     "Sub-subclass": [],
223     "Area of AI Application": "workforce management,
scheduling",
224     "Online": "No"
225 },
226 "11": {
227     "Country": "United States",
228     "State": "Florida",
229     "City": "Broward County",
230     "Continent": "North America",
231     "Company": "Northpointe (now part of Equivant)",
232     "Company city": "Canton, OH",

```



```

233     "Affected population": "Criminal Defendants",
234     "Number of people actually affected": "Unknown",
235     "Number of people potentially affected": "Unknown",
236     "Class of irresponsible AI use": [
237         "Discrimination",
238         "Human Incompetence"
239     ],
240     "Subclasses": [
241         "Discrimination-> Data bias",
242         "Discrimination-> Algorithmic bias",
243         "Human Incompetence-> Technical"
244     ],
245     "Sub-subclass": [
246         "Data bias-> Race",
247         "Data bias-> Gender",
248         "Algorithmic bias-> Race",
249         "Algorithmic bias-> Gender"
250     ],
251     "Area of AI Application": "criminal risk assessment",
252     "Online": "Yes"
253 },
254 "12": {
255     "Country": "",
256     "State": "",
257     "City": "",
258     "Continent": "",
259     "Company": "",
260     "Company city": "",
261     "Affected population": "",
262     "Number of people actually affected": "",
263     "Number of people potentially affected": "",
264     "Class of irresponsible AI use": [],
265     "Subclasses": [],
266     "Sub-subclass": [],
267     "Area of AI Application": "",
268     "Online": ""
269 },
270 "13": {
271     "Country": "Global",
272     "State": "N/A",
273     "City": "N/A",
274     "Continent": "Global",
275     "Company": "Google LLC",
276     "Company city": "Mountain View, CA",
277     "Affected population": "Online Users",
278     "Number of people actually affected": "Unknown",
279     "Number of people potentially affected": "Millions",
280     "Class of irresponsible AI use": [
281         "Discrimination"
282     ],

```

```

283     "Subclasses": [
284         "Discrimination-> Data bias",
285         "Discrimination-> Algorithmic bias"
286     ],
287     "Sub-subclass": [
288         "Data bias-> Race",
289         "Data bias-> Sexual Orientation",
290         "Data bias-> Economic",
291         "Algorithmic bias-> Gender",
292         "Algorithmic bias-> Race"
293     ],
294     "Area of AI Application": "content moderation, online
discussion forums",
295     "Online": "Yes"
296 },
297 "14": {
298     "Country": "Global",
299     "State": "Global",
300     "City": "Global",
301     "Continent": "Global",
302     "Company": "Google LLC",
303     "Company city": "Mountain View, CA",
304     "Affected population": "Users of Google's Cloud
Natural Language API, including developers and indirectly
, religious and ethnic minorities",
305     "Number of people actually affected": "Unknown",
306     "Number of people potentially affected": "Unknown",
307     "Class of irresponsible AI use": [
308         "Discrimination"
309     ],
310     "Subclasses": [
311         "Discrimination-> Data bias",
312         "Discrimination-> Algorithmic bias"
313     ],
314     "Sub-subclass": [
315         "Data bias-> Race",
316         "Data bias-> Sexual Orientation"
317     ],
318     "Area of AI Application": "Natural Language
Processing, Sentiment Analysis",
319     "Online": "Yes"
320 },
321 "15": {
322     "1": {
323         "Country": "Global",
324         "State": "",
325         "City": "",
326         "Continent": "",
327         "Company": "Amazon",
328         "Company city": "Seattle, WA",

```

```

329         "Affected population": "Authors and readers of
gay and lesbian-themed books",
330         "Number of people actually affected": "Unknown",
331         "Number of people potentially affected": "Unknown",
",
332         "Class of irresponsible AI use": [
333             "Human Incompetence"
334         ],
335         "Subclasses": [
336             "Human Incompetence-> Technical"
337         ],
338         "Sub-subclass": [],
339         "Area of AI Application": "Online retail, Sales
ranking system",
340         "Online": "Yes"
341     }
342 }
343 }

```

Listing 2.3. Few-shot with CoT classification results of 17 incidents

D

Tree of Thought Framework Classification Results

The Tree of Thoughts framework allows the LLM to explore various reasoning paths and collectively determine the most well-supported classification based on the text of the article text.

```

1 {
2     "1": {
3         "Country": "Worldwide",
4         "State": "",
5         "City": "",
6         "Continent": "Worldwide",
7         "Company": "YouTube",
8         "Company city": "San Bruno",
9         "Company state": "California",
10        "Affected population": [
11            "Student",
12            "Online Female Population",
13            "Twitter Users"
14        ],
15        "Number of people actually affected": "Unknown",
16        "Number of people potentially affected": "Unknown",
17        "Classes of irresponsible AI use": [
18            "Human Incompetence",
19            "Disinformation",
20            "Mental Health",
21            "Other"

```

```

22     ],
23     "Subclasses": {
24         "Human Incompetence": [
25             "Technical"
26         ],
27         "Disinformation": [
28             "Textual",
29             "Image",
30             "Video"
31         ],
32         "Mental Health": [],
33         "Other": []
34     },
35     "Sub-subclass": [],
36     "Area of AI Application": "content filtering",
37     "Online": "yes"
38 },
39 "2": {
40     "Country": "United States",
41     "State": "New Jersey",
42     "City": "Robbinsville",
43     "Continent": "North America",
44     "Company": "Amazon",
45     "Company city": "Seattle",
46     "Company state": "Washington",
47     "Affected population": [
48         "Warehouse Workers"
49     ],
50     "Number of people actually affected": "24",
51     "Number of people potentially affected": "54",
52     "Classes of irresponsible AI use": [
53         "Human Incompetence"
54     ],
55     "Subclasses": {
56         "Human Incompetence": [
57             "Technical"
58         ]
59     },
60     "Sub-subclass": [],
61     "Area of AI Application": "Warehouse automation",
62     "Online": "no"
63 },
64 "3": {
65     "Country": "Indonesia",
66     "State": "",
67     "City": "Jakarta",
68     "Continent": "Asia",
69     "Company": "Boeing",
70     "Company city": "Chicago",
71     "Company state": "Illinois",

```

```

72     "Affected population": [
73         "Air Travelers"
74     ],
75     "Number of people actually affected": "189",
76     "Number of people potentially affected": "Unknown",
77     "Classes of irresponsible AI use": [
78         "Human Incompetence",
79         "Other"
80     ],
81     "Subclasses": {
82         "Human Incompetence": [
83             "Technical"
84         ],
85         "Other": []
86     },
87     "Sub-subclass": {
88         "Technical": []
89     },
90     "Area of AI Application": "Aviation safety and
control systems",
91     "Online": "no"
92 },
93 "4": {
94     "Country": "USA",
95     "State": "Arizona",
96     "City": "Tempe",
97     "Continent": "North America",
98     "Company": "Uber",
99     "Company city": "San Francisco",
100    "Company state": "California",
101    "Affected population": [
102        "Pedestrians"
103    ],
104    "Number of people actually affected": "1",
105    "Number of people potentially affected": "Unknown",
106    "Classes of irresponsible AI use": [
107        "Human Incompetence",
108        "Other"
109    ],
110    "Subclasses": {
111        "Human Incompetence": [
112            "Technical"
113        ],
114        "Other": []
115    },
116    "Sub-subclass": [],
117    "Area of AI Application": "Autonomous Vehicles",
118    "Online": "no"
119 },
120 "5": {

```

```

121     "Country": "United States",
122     "State": "",
123     "City": "",
124     "Continent": "North America",
125     "Company": "Not specified",
126     "Company city": "Not specified",
127     "Company state": "",
128     "Affected population": [
129         "Patients undergoing robotic surgery"
130     ],
131     "Number of people actually affected": "1535",
132     "Number of people potentially affected": "1.7 million
",
133     "Classes of irresponsible AI use": [
134         "Human Incompetence",
135         "Technical"
136     ],
137     "Subclasses": {
138         "Human Incompetence": [
139             "Technical"
140         ]
141     },
142     "Sub-subclass": [],
143     "Area of AI Application": "Robotic surgery",
144     "Online": "no"
145 },
146 "6": {
147     "Country": "Worldwide",
148     "State": "",
149     "City": "",
150     "Continent": "Worldwide",
151     "Company": "Microsoft",
152     "Company city": "Redmond",
153     "Company state": "Washington",
154     "Affected population": [
155         "Twitter Users",
156         "Online Population"
157     ],
158     "Number of people actually affected": "Unknown",
159     "Number of people potentially affected": "Unknown",
160     "Classes of irresponsible AI use": [
161         "Human Incompetence",
162         "Disinformation",
163         "Other"
164     ],
165     "Subclasses": {
166         "Human Incompetence": [
167             "Technical"
168         ],
169         "Disinformation": [

```

```

170         "Textual"
171     ],
172     "Other": []
173 },
174     "Sub-subclass": [],
175     "Area of AI Application": "Social Media Interaction",
176     "Online": "yes"
177 },
178 "7": {
179     "Country": "Worldwide",
180     "State": "",
181     "City": "",
182     "Continent": "Worldwide",
183     "Company": "Wikipedia",
184     "Company city": "San Francisco",
185     "Company state": "California",
186     "Affected population": [
187         "Online Users"
188     ],
189     "Number of people actually affected": "Unknown",
190     "Number of people potentially affected": "Unknown",
191     "Classes of irresponsible AI use": [
192         "Human Incompetence",
193         "Other"
194     ],
195     "Subclasses": {
196         "Human Incompetence": [
197             "Technical"
198         ],
199         "Other": []
200     },
201     "Sub-subclass": [],
202     "Area of AI Application": "Content management and
editing",
203     "Online": "yes"
204 },
205 "8": {
206     "Country": "United States",
207     "State": "California",
208     "City": "San Francisco",
209     "Continent": "North America",
210     "Company": "Uber",
211     "Company city": "San Francisco",
212     "Company state": "California",
213     "Affected population": [
214         "General Public",
215         "Bicyclists",
216         "Pedestrians"
217     ],
218     "Number of people actually affected": "Unknown",

```

```

219     "Number of people potentially affected": "Unknown",
220     "Classes of irresponsible AI use": [
221         "Human Incompetence",
222         "Other"
223     ],
224     "Subclasses": {
225         "Human Incompetence": [
226             "Technical"
227         ],
228         "Other": []
229     },
230     "Sub-subclass": {
231         "Technical": []
232     },
233     "Area of AI Application": "Autonomous Vehicles",
234     "Online": "no"
235 },
236 "9": {
237     "Country": "United States",
238     "State": "New York",
239     "City": "New York City",
240     "Continent": "North America",
241     "Company": "New York City Department of Education",
242     "Company city": "New York City",
243     "Company state": "New York",
244     "Affected population": [
245         "Student",
246         "Teachers"
247     ],
248     "Number of people actually affected": "Unknown",
249     "Number of people potentially affected": "Unknown",
250     "Classes of irresponsible AI use": [
251         "Human Incompetence",
252         "Other"
253     ],
254     "Subclasses": {
255         "Human Incompetence": [
256             "Administrative"
257         ],
258         "Other": []
259     },
260     "Sub-subclass": [],
261     "Area of AI Application": "Educational assessment",
262     "Online": "yes"
263 },
264 "10": {
265     "Country": "United States",
266     "State": "",
267     "City": "",
268     "Continent": "North America",

```



```

269     "Company": "Starbucks",
270     "Company city": "Seattle",
271     "Company state": "Washington",
272     "Affected population": [
273         "Retail Workers",
274         "Single Parents",
275         "Students"
276     ],
277     "Number of people actually affected": "Unknown",
278     "Number of people potentially affected": "Unknown",
279     "Classes of irresponsible AI use": [
280         "Human Incompetence",
281         "Discrimination"
282     ],
283     "Subclasses": {
284         "Human Incompetence": [
285             "Technical"
286         ],
287         "Discrimination": [
288             "Data bias"
289         ]
290     },
291     "Sub-subclass": {
292         "Technical": [],
293         "Data bias": [
294             "Economic",
295             "Other"
296         ]
297     },
298     "Area of AI Application": "Workforce management",
299     "Online": "yes"
300 },
301 "11": {
302     "Country": "United States",
303     "State": "Florida",
304     "City": "Broward County",
305     "Continent": "North America",
306     "Company": "Northpointe (now part of Equivant)",
307     "Company city": "Canton",
308     "Company state": "Ohio",
309     "Affected population": [
310         "Black Defendants",
311         "White Defendants"
312     ],
313     "Number of people actually affected": "Unknown",
314     "Number of people potentially affected": "Unknown",
315     "Classes of irresponsible AI use": [
316         "Discrimination"
317     ],
318     "Subclasses": {

```

```

319         "Discrimination": [
320             "Data bias",
321             "Algorithmic bias"
322         ]
323     },
324     "Sub-subclass": {
325         "Data bias": [
326             "Race",
327             "Gender"
328         ],
329         "Algorithmic bias": [
330             "Feedback loop",
331             "Optimization function"
332         ]
333     },
334     "Area of AI Application": "Criminal justice risk
assessments",
335     "Online": "yes"
336 },
337 "12": {
338     "Country": "",
339     "State": "",
340     "City": "",
341     "Continent": "",
342     "Company": "",
343     "Company city": "",
344     "Company state": "",
345     "Affected population": [],
346     "Number of people actually affected": "",
347     "Number of people potentially affected": "",
348     "Classes of irresponsible AI use": [],
349     "Subclasses": {},
350     "Sub-subclass": [],
351     "Area of AI Application": "",
352     "Online": ""
353 },
354 "13": {
355     "Country": "United States",
356     "State": "",
357     "City": "",
358     "Continent": "North America",
359     "Company": "Google's Alphabet",
360     "Company city": "Mountain View",
361     "Company state": "California",
362     "Affected population": [
363         "Online Population"
364     ],
365     "Number of people actually affected": "Unknown",
366     "Number of people potentially affected": "Unknown",
367     "Classes of irresponsible AI use": [

```

```

368         "Discrimination",
369         "Human Incompetence"
370     ],
371     "Subclasses": {
372         "Discrimination": [
373             "Data bias",
374             "Algorithmic bias"
375         ],
376         "Human Incompetence": [
377             "Technical"
378         ]
379     },
380     "Sub-subclass": {
381         "Data bias": [
382             "Gender",
383             "Race",
384             "Sexual Orientation"
385         ],
386         "Algorithmic bias": [
387             "Other"
388         ]
389     },
390     "Area of AI Application": "content filtering",
391     "Online": "yes"
392 },
393 "14": {
394     "Country": "Worldwide",
395     "State": "",
396     "City": "",
397     "Continent": "Worldwide",
398     "Company": "Google",
399     "Company city": "Mountain View",
400     "Company state": "California",
401     "Affected population": [
402         "Online Population"
403     ],
404     "Number of people actually affected": "Unknown",
405     "Number of people potentially affected": "Unknown",
406     "Classes of irresponsible AI use": [
407         "Discrimination"
408     ],
409     "Subclasses": {
410         "Discrimination": [
411             "Data bias"
412         ]
413     },
414     "Sub-subclass": {
415         "Data bias": [
416             "Gender",
417             "Race",

```

```

418         "Sexual Orientation"
419     ]
420     },
421     "Area of AI Application": "Natural Language
Processing",
422     "Online": "yes"
423 },
424 "15": {
425     "Country": "Worldwide",
426     "State": "",
427     "City": "",
428     "Continent": "Worldwide",
429     "Company": "Amazon.com",
430     "Company city": "Seattle",
431     "Company state": "Washington",
432     "Affected population": [
433         "LGBTQ"
434     ],
435     "Number of people actually affected": "Unknown",
436     "Number of people potentially affected": "Unknown",
437     "Classes of irresponsible AI use": [
438         "Discrimination",
439         "Human Incompetence"
440     ],
441     "Subclasses": {
442         "Discrimination": [
443             "Data bias"
444         ],
445         "Human Incompetence": [
446             "Technical"
447         ]
448     },
449     "Sub-subclass": {
450         "Data bias": [
451             "Sexual Orientation"
452         ]
453     },
454     "Area of AI Application": "Content filtering",
455     "Online": "yes"
456 }
457 }

```

Listing 2.4. ToT framework with few-shot learning and CoT classification results of 17 incidents

E**Zero-Shot Learning Classification Results**

The zero-shot learning classification generally got the geolocation fields and affected population information without much issue.

```

1  [
2    {
3      "id": "0",
4      "Country": "Worldwide",
5      "State": "",
6      "City": "",
7      "Continent": "Worldwide",
8      "Company": "Multiple companies involved in AI development",
9      "Company city": "",
10     "Company state": "",
11     "Affected population": "Global population",
12     "Number of people actually affected": "Unknown",
13     "Number of people potentially affected": "Unknown",
14     "Classes of irresponsible AI use": "Environmental Impact",
15     "Subclasses": "",
16     "Sub-subclass": "",
17     "Area of AI Application": "General AI development and application",
18     "Online": "No",
19     "title": "AI's Environmental Footprint",
20     "description": "The environmental footprint of AI, particularly in training large models, is significant. According to a study by researchers at the University of Massachusetts, the energy used in training certain popular large AI models can produce about 626,000 pounds of carbon dioxide. This amount is equivalent to roughly 300 round-trip flights between New York and San Francisco, highlighting the substantial carbon footprint associated with advanced AI technologies. This data underscores the need for more sustainable practices in the field of AI to mitigate its impact on climate change.",
21     "date occurrence": "2023-07-18",
22     "date publicly known": "2023-07-18"
23   },
24   {
25     "id": "1",
26     "Country": "Worldwide",
27     "State": "",
28     "City": "",
29     "Continent": "Worldwide",
30     "Company": "YouTube",

```

```

31   "Company city": "San Bruno",
32   "Company state": "California",
33   "Affected population": "Online Child Population",
34   "Number of people actually affected": "Unknown",
35   "Number of people potentially affected": "Unknown",
36   "Classes of irresponsible AI use": "Human Incompetence,
Mental Health, Other",
37   "Subclasses": "Administrative",
38   "Sub-subclass": "",
39   "Area of AI Application": "content filtering",
40   "Online": "yes",
41   "title": "Googles YouTube Kids App Presents
Inappropriate Content",
42   "description": "YouTubes content filtering and
recommendation algorithms exposed children to disturbing
and inappropriate videos.",
43   "date occurrence": "2015-05-19",
44   "date publicly known": "2017-03-26"
45 },
46 {
47   "id": "2",
48   "Country": "United States",
49   "State": "New Jersey",
50   "City": "Robbinsville",
51   "Continent": "North America",
52   "Company": "Amazon",
53   "Company city": "Seattle",
54   "Company state": "Washington",
55   "Affected population": "Warehouse Workers",
56   "Number of people actually affected": "24",
57   "Number of people potentially affected": "54",
58   "Classes of irresponsible AI use": "Human Incompetence",
59   "Subclasses": "Technical",
60   "Sub-subclass": "",
61   "Area of AI Application": "Warehouse Automation",
62   "Online": "No",
63   "title": "Warehouse robot ruptures can of bear spray and
injures workers",
64   "description": "Twenty-four Amazon workers in New Jersey
were hospitalized after a robot punctured a can of bear
repellent spray in a warehouse.",
65   "date occurrence": "2018-12-05",
66   "date publicly known": "2018-12-06"
67 },
68 {
69   "id": "3",
70   "Country": "Indonesia",
71   "State": "",
72   "City": "Jakarta",
73   "Continent": "Asia",

```

```

74     "Company": "Boeing",
75     "Company city": "Chicago",
76     "Company state": "Illinois",
77     "Affected population": "Airline Passengers, Crew Members",
78     "Number of people actually affected": "189",
79     "Number of people potentially affected": "Unknown",
80     "Classes of irresponsible AI use": "Human Incompetence,
Other",
81     "Subclasses": "Technical",
82     "Sub-subclass": "",
83     "Area of AI Application": "Aviation safety systems",
84     "Online": "no",
85     "title": "Crashes with Maneuvering Characteristics
Augmentation System (MCAS)",
86     "description": "A Boeing 737 crashed into the sea,
killing 189 people, after faulty sensor data caused an
automated maneuvering system to repeatedly push the plane
's nose downward.",
87     "date occurrence": "2018-10-27",
88     "date publicly known": "2019-03-13"
89 },
90 {
91     "id": "4",
92     "Country": "USA",
93     "State": "Arizona",
94     "City": "Tempe",
95     "Continent": "North America",
96     "Company": "Uber",
97     "Company city": "San Francisco",
98     "Company state": "California",
99     "Affected population": "Pedestrians, Autonomous Vehicle
Test Subjects",
100     "Number of people actually affected": "1",
101     "Number of people potentially affected": "Unknown",
102     "Classes of irresponsible AI use": "Human Incompetence,
Other",
103     "Subclasses": "Technical",
104     "Sub-subclass": "",
105     "Area of AI Application": "Autonomous Vehicles",
106     "Online": "no",
107     "title": "Uber AV Killed Pedestrian in Arizona",
108     "description": "An Uber autonomous vehicle (AV) in
autonomous mode struck and killed a pedestrian in Tempe,
Arizona.",
109     "date occurrence": "2018-03-18",
110     "date publicly known": "2018-03-22"
111 },
112 {
113     "id": "5",

```

```

114     "Country": "United States",
115     "State": "",
116     "City": "",
117     "Continent": "North America",
118     "Company": "Not specified in the article",
119     "Company city": "Not specified in the article",
120     "Company state": "",
121     "Affected population": "Patients undergoing robotic
122     surgery",
123     "Number of people actually affected": "Unknown",
124     "Number of people potentially affected": "Unknown",
125     "Classes of irresponsible AI use": "Human Incompetence,
126     Technical Difficulties",
127     "Subclasses": "Technical, Malfunction, Design flaws",
128     "Sub-subclass": "",
129     "Area of AI Application": "Robotic surgery",
130     "Online": "No",
131     "title": "Collection of Robotic Surgery Malfunctions",
132     "description": "Study on database reports of robotic
133     surgery malfunctions (8,061), including those ending in
134     injury (1,391) and death (144), between 2000 and 2013.",
135     "date occurrence": "2015-07-13",
136     "date publicly known": "2015-07-20"
137 },
138 {
139     "id": "6",
140     "Country": "Worldwide",
141     "State": "",
142     "City": "",
143     "Continent": "Worldwide",
144     "Company": "Microsoft",
145     "Company city": "Redmond",
146     "Company state": "Washington",
147     "Affected population": "Twitter Users, Online Population
148     ",
149     "Number of people actually affected": "Unknown",
150     "Number of people potentially affected": "Unknown",
151     "Classes of irresponsible AI use": "Human Incompetence,
152     Disinformation, Other",
153     "Subclasses": "Technical, Textual",
154     "Sub-subclass": "",
155     "Area of AI Application": "Social Media Interaction",
156     "Online": "yes",
157     "title": "TayBot",
158     "description": "Microsoft's Tay, an artificially
159     intelligent chatbot, was released on March 23, 2016 and
160     removed within 24 hours due to multiple racist, sexist,
161     and anit-semitic tweets generated by the bot.",
162     "date occurrence": "2016-03-24",
163     "date publicly known": "2019-11-24"

```



```

155 },
156 {
157   "id": "9",
158   "Country": "United States",
159   "State": "New York",
160   "City": "New York City",
161   "Continent": "North America",
162   "Company": "New York City Department of Education",
163   "Company city": "New York City",
164   "Company state": "New York",
165   "Affected population": "Students, Teachers",
166   "Number of people actually affected": "Unknown",
167   "Number of people potentially affected": "Unknown",
168   "Classes of irresponsible AI use": "Human Incompetence,
Other",
169   "Subclasses": "Administrative",
170   "Sub-subclass": "",
171   "Area of AI Application": "Educational assessment and
teacher evaluation",
172   "Online": "yes",
173   "title": "NY City School Teacher Evaluation Algorithm
Contested",
174   "description": "An algorithm used to rate the
effectiveness of school teachers in New York has resulted
in thousands of disputes of its results.",
175   "date occurrence": "2012-02-25",
176   "date publicly known": "2013-10-18"
177 },
178 {
179   "id": "10",
180   "Country": "United States",
181   "State": "",
182   "City": "San Diego",
183   "Continent": "North America",
184   "Company": "Starbucks",
185   "Company city": "Seattle",
186   "Company state": "Washington",
187   "Affected population": "Low-income single mothers, Retail
workers, Baristas",
188   "Number of people actually affected": "Unknown",
189   "Number of people potentially affected": "130,000 (
Starbucks baristas nationwide)",
190   "Classes of irresponsible AI use": "Human Incompetence,
Mental Health, Other",
191   "Subclasses": "Technical",
192   "Sub-subclass": "",
193   "Area of AI Application": "Workforce management, Employee
scheduling",
194   "Online": "No",

```

```

195     "title": "Kronos Scheduling Algorithm Allegedly Caused
196     Financial Issues for Starbucks Employees",
197     "description": "Kronos scheduling algorithm and its
198     use by Starbucks managers allegedly negatively impacted
199     financial and scheduling stability for Starbucks
200     employees, which disadvantaged wage workers.",
201     "date occurrence": "2014-08-14",
202     "date publicly known": "2015-06-02"
203   },
204   {
205     "id": "11",
206     "Country": "United States",
207     "State": "Florida",
208     "City": "Fort Lauderdale",
209     "Continent": "North America",
210     "Company": "Northpointe",
211     "Company city": "Traverse City",
212     "Company state": "Michigan",
213     "Affected population": "African American",
214     "Number of people actually affected": "Unknown",
215     "Number of people potentially affected": "Unknown",
216     "Classes of irresponsible AI use": "Discrimination",
217     "Subclasses": "Data bias, Algorithmic bias",
218     "Sub-subclass": "Race, Feedback loop",
219     "Area of AI Application": "Criminal justice risk
220     assessment",
221     "Online": "yes",
222     "title": "Northpointe Risk Models",
223     "description": "An algorithm developed by Northpointe and
224     used in the penal system is two times more likely to
225     incorrectly label a black person as a high-risk re-
226     offender and is two times more likely to incorrectly
227     label a white person as low-risk for reoffense according
228     to a ProPublica review.",
229     "date occurrence": "2016-05-23",
230     "date publicly known": "2016-05-22"
231   },
232   {
233     "id": "13",
234     "Country": "Worldwide",
235     "State": "",
236     "City": "",
237     "Continent": "Worldwide",
238     "Company": "Google's Alphabet",
239     "Company city": "Mountain View",
240     "Company state": "California",
241     "Affected population": "Online Population",
242     "Number of people actually affected": "Unknown",
243     "Number of people potentially affected": "Unknown",

```

```

234   "Classes of irresponsible AI use": "Discrimination, Human
      Incompetence",
235   "Subclasses": "Data bias, Algorithmic bias, Technical",
236   "Sub-subclass": "Race, Sexual Orientation, Other,
      Feedback loop, Other",
237   "Area of AI Application": "content filtering, online
      discussion moderation",
238   "Online": "yes",
239   "title": "High-Toxicity Assessed on Text Involving Women
      and Minority Groups",
240   "description": "Google's Perspective API, which assigns a
      toxicity score to online text, seems to award higher
      toxicity scores to content involving non-white, male,
      Christian, heterosexual phrases.",
241   "date occurrence": "2017-02-27",
242   "date publicly known": "2021-02-09"
243 },
244 {
245   "id": "14",
246   "Country": "Worldwide",
247   "State": "",
248   "City": "",
249   "Continent": "Worldwide",
250   "Company": "Google",
251   "Company city": "Mountain View",
252   "Company state": "California",
253   "Affected population": "Ethnic and religious minorities,
      LGBTQ community",
254   "Number of people actually affected": "Unknown",
255   "Number of people potentially affected": "Unknown",
256   "Classes of irresponsible AI use": "Discrimination",
257   "Subclasses": "Data bias",
258   "Sub-subclass": "Race, Sexual Orientation",
259   "Area of AI Application": "Sentiment analysis, natural
      language processing",
260   "Online": "yes",
261   "title": "Biased Sentiment Analysis",
262   "description": "Google Cloud's Natural Language API
      provided racist, homophobic, and antisemitic sentiment
      analyses.",
263   "date occurrence": "2017-10-26",
264   "date publicly known": "2017-10-25"
265 },
266 {
267   "id": "167",
268   "Country": "United States",
269   "State": "California",
270   "City": "Stanford",
271   "Continent": "North America",
272   "Company": "Stanford Graduate School of Business",

```

```

273     "Company city": "Stanford",
274     "Company state": "California",
275     "Affected population": "LGBTQ",
276     "Number of people actually affected": "Unknown",
277     "Number of people potentially affected": "Unknown",
278     "Classes of irresponsible AI use": "Discrimination,
Pseudoscience, Mental Health",
279     "Subclasses": "Data bias, Facial",
280     "Sub-subclass": "Sexual Orientation",
281     "Area of AI Application": "Facial recognition analysis",
282     "Online": "yes",
283     "title": "Researchers' Homosexual-Men Detection Model
Denounced as a Threat to LGBTQ Peoples Safety and
Privacy",
284     "description": "Researchers at Stanford Graduate School
of Business developed a model that determined, on a
binary scale, whether someone was homosexual using only
his facial image, which advocacy groups such as GLAAD and
the Human Rights Campaign denounced as flawed science
and threatening to LGBTQ folks.",
285     "date occurrence": "2017-09-07",
286     "date publicly known": "2017-10-08"
287 },
288 {
289     "id": "382",
290     "Country": "United Kingdom",
291     "State": "",
292     "City": "London",
293     "Continent": "Europe",
294     "Company": "Meta",
295     "Company city": "Menlo Park",
296     "Company state": "California",
297     "Affected population": "Online Female Population",
298     "Number of people actually affected": "1",
299     "Number of people potentially affected": "Unknown",
300     "Classes of irresponsible AI use": "Mental Health,
Disinformation",
301     "Subclasses": "Textual",
302     "Sub-subclass": "",
303     "Area of AI Application": "content filtering",
304     "Online": "yes",
305     "title": "Instagram's Exposure of Harmful Content
Contributed to Teenage Girls Suicide",
306     "description": "Instagram was ruled by a judge to have
contributed to the death of a teenage girl in the UK
allegedly through its exposure and recommendation of
suicide, self-harm, and depressive content.",
307     "date occurrence": "2017-11-21",
308     "date publicly known": "2022-09-30"
309 },

```

```

310 {
311   "id": "451",
312   "Country": "United Kingdom",
313   "State": "",
314   "City": "London",
315   "Continent": "Europe",
316   "Company": "Stability AI",
317   "Company city": "London",
318   "Company state": "",
319   "Affected population": "Artists, Photographers, Content
320   Creators",
321   "Number of people actually affected": "Unknown",
322   "Number of people potentially affected": "
323   Thousands possibly millions",
324   "Classes of irresponsible AI use": "Copyright Violation",
325   "Subclasses": "",
326   "Sub-subclass": "",
327   "Area of AI Application": "AI art tool, Image generation
328   ",
329   "Online": "yes",
330   "title": "Stable Diffusion's Training Data Contained
331   Copyrighted Images",
332   "description": "Stability AI reportedly scraped
333   copyrighted images by Getty Images to be used as training
334   data for Stable Diffusion model.",
335   "date occurrence": "2022-10-16",
336   "date publicly known": "2023-01-16"
337 },
338 {
339   "id": "505",
340   "Country": "Belgium",
341   "State": "",
342   "City": "",
343   "Continent": "Europe",
344   "Company": "Chai Research",
345   "Company city": "Silicon Valley",
346   "Company state": "California",
347   "Affected population": "Chatbot Users",
348   "Number of people actually affected": "1",
349   "Number of people potentially affected": "Unknown",
350   "Classes of irresponsible AI use": "Mental Health, Human
351   Incompetence, Other",
352   "Subclasses": "Technical",
353   "Sub-subclass": "",
354   "Area of AI Application": "Conversational AI",
355   "Online": "yes",
356   "title": "Man Reportedly Committed Suicide Following
357   Conversation with Chai Chatbot",
358   "description": "A Belgian man reportedly committed
359   suicide following a conversation with Eliza, a language

```

```

351     model developed by Chai that encouraged the man to commit
352     suicide to improve the health of the planet.",
353     "date occurrence": "2023-03-27",
354     "date publicly known": "2023-03-27"
355 },
356 {
357     "id": "39",
358     "Country": "Worldwide",
359     "State": "",
360     "City": "",
361     "Continent": "Worldwide",
362     "Company": "Multiple companies including Adobe,
363     University of Washington, Lyrebird, and others mentioned
364     in the context of developing AI technology",
365     "Company city": "Multiple locations",
366     "Company state": "",
367     "Affected population": "General Public",
368     "Number of people actually affected": "Unknown",
369     "Number of people potentially affected": "Worldwide
370     population with internet access",
371     "Classes of irresponsible AI use": "Disinformation,
372     Mental Health, Other",
373     "Subclasses": "Textual, Image, Video, Audio",
374     "Sub-subclass": "",
375     "Area of AI Application": "Media and information
376     dissemination, including video conferencing, virtual
377     reality, and potentially any form of digital
378     communication",
379     "Online": "yes",
380     "title": "Deepfake Obama Introduction of Deepfakes",
381     "description": "University of Washington researchers made
382     a deepfake of Obama, followed by Jordan Peele",
383     "date occurrence": "2017-07-01",
384     "date publicly known": "2017-07-18"
385 }
386 ]

```

Listing 2.5. Zero-Shot Learning Classification Results of 17 incidents

F**Table comparison of classification results****Table 1.** Comparison of Classification Results for Incident 13 (Part 1)

Method	ID	Country	State	City	Continent	Company
Manual	13	United States	-	-	North America	Google LLC
Zero-Shot	13	Worldwide	-	-	Worldwide	Google's Alphabet
Few-Shot CoT	13	Worldwide	-	-	Worldwide	Google
ToT Few-Shot and CoT	13	United States	-	-	North America	Google's Alphabet

Table 2. Comparison of Classification Results for Incident 13 (Part 2)

Co. City	Co. State	Affected Pop.	Classes	Subclasses
Mountain View	California	Minorities, women, etc.	Discrimination	Data bias
Mountain View	California	Online Population	Discrimination, Human Incompetence	Data bias, Algorithmic bias, Technical
Mountain View	California	Online Users	Discrimination, Human Incompetence	Data bias, Algorithmic bias, Technical
Mountain View	California	Online Users	Discrimination, Human Incompetence, Disinformation	Data bias, Algorithmic bias, Technical, Textual

Table 3. Comparison of Classification Results for Incident 4 (Part 1)

Method	ID	Country	State	City	Continent	Company
Manual	4	United States	Arizona	Tempe	North America	Uber
Zero-Shot	4	USA	Arizona	Tempe	North America	Uber
Few-Shot CoT	4	USA	Arizona	Tempe	North America	Uber
ToT Few-Shot	4	USA	Arizona	Tempe	North America	Uber

Table 4. Comparison of Classification Results for Incident 4 (Part 2)

Co. City	Co. State	Affected Pop.	Classes	Subclasses	AI Application	Online
San Francisco	California	Pedestrians	Human in-competence	Technical	Autonomous driving	No
San Francisco	California	Pedestrians, AV Test Subjects	Human In-competence, Other	Technical	Autonomous Vehicles	No
San Francisco	California	Pedestrians, AV Users, Public	Human In-competence, Other	Technical	Autonomous Vehicles	No
San Francisco	California	Pedestrians, AV Test Subjects	Human In-competence, Other	Technical	Autonomous Vehicles	No