# As an AI Language Model, "Yes I Would Recommend Calling the Police": Norm Inconsistency in LLM Decision-Making

**Shomik Jain[1], D. Calacci[2]\*, Ashia Wilson[3]\***

[1]Institute for Data, Systems, and Society, Massachusetts Institute of Technology
[2]College of Information Sciences and Technology, Penn State University
[3]Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology
shomikj@mit.edu, dcalacci@psu.edu, ashia07@mit.edu

## Abstract

We investigate the phenomenon of *norm inconsistency*: where LLMs apply different norms in similar situations. Specifically, we focus on the high-risk application of deciding whether to call the police in Amazon Ring home surveillance videos. We evaluate the decisions of three state-of-the-art LLMs – GPT-4, Gemini 1.0, and Claude 3 Sonnet – in relation to the activities portrayed in the videos, the subjects' skin-tone and gender, and the characteristics of the neighborhoods where the videos were recorded. Our analysis reveals significant norm inconsistencies: (1) a discordance between the recommendation to call the police and the actual presence of criminal activity, and (2) biases influenced by the racial demographics of the neighborhoods. These results highlight the arbitrariness of model decisions in the surveillance context and the limitations of current bias detection and mitigation strategies in normative decision-making.

## Introduction

Existing work characterizing the moral and ethical reasoning of large language models (LLMs) has revealed at least one emerging cluster of concerns: models do not consistently apply the same norms across scenarios, and their normative judgments are often discordant with the facts of a scenario (Agarwal et al. 2024; Johnson et al. 2022a; Almeida et al. 2024). We refer to this phenomenon as *norm inconsistency*. While humans sometimes exhibit this behavior when applying normative rules (Balagopalan et al. 2023), the potential for more severe norm inconsistency in AI decision-making presents serious issues for system reliability and can perpetuate unfair outcomes.

Many high-value use-cases for LLMs involve making decisions in areas deeply rooted in social norms, such as employment and hiring (Wicaksana and Liem 2017), policing and criminal justice (Kleinberg et al. 2018), and medicine (Minssen, Vayena, and Cohen 2023). Yet surprisingly little is known about how LLMs make normative judgments in real-world scenarios. In the context of surveillance and law enforcement, which we focus on in this work, norm inconsistency can manifest in unsettling ways. A model might state that no crime occurred but still recommend calling the



(a) Example GPT-4 Response to Crime Prompt



(b) Example GPT-4 Response to Police Prompt

Figure 1: Example of norm-inconsistency in GPT-4 where the model says no crime occurred but recommends police intervention. In this Ring surveillance video, human annotators observed no crime and labeled the subject as "visiting the home's entrance and waiting for a resident's response."

police, or vice versa (Figure 1). Or a model might recommend no police intervention for a theft in one neighborhood, but then recommend intervention for a strikingly similar scenario in another neighborhood.

In this work, we investigate the potential real-world impacts of norm inconsistency in a specific high-risk application[1]: whether to flag home surveillance videos for police intervention. Specifically, we prompt GPT-4, Gemini, and Claude with real videos from the Amazon Ring Neighbors platform and test (1) whether models state that a crime is happening and (2) whether they recommend calling the police. We then investigate the judgment criteria of different

---

\*These authors contributed equally.

[1]Our evaluation of AI in the surveillance context is not intended to encourage enthusiasm for this domain, and our results highlight several reasons why.

LLMs by fitting a linear model to predict their decisions from annotations of the portrayed activity and other characteristics about the video's subject and neighborhood ($R^2 = 0.10$ to $0.37$).

We find that all models exhibit norm inconsistency by recommending police intervention in cases where: (1) they state no crime occurred, (2) they refuse to respond to the crime prompt, and (3) when they answer the crime prompt ambiguously. Models also make inconsistent normative judgments between videos that portray similar activities, including activities that do not involve a crime. Moreover, we unexpectedly find that while LLM recommendations are not influenced by the skin-tone of a video subject, they are associated with the demographics of the neighborhood that a video was recorded in. This is surprising because neighborhood characteristics are not provided in text prompts and not explicit in the video content. In our discussion, we describe in what ways norm inconsistency presents a problem for both the surveillance context and high-stakes settings in general.

## Background and Related Work

We review related work in normative decision-making and measuring bias in LLMs, and also provide background about AI for surveillance and Amazon Ring, the source for our dataset. We highlight how our work represents one of the first evaluations of normative decision-making in LLMs using *real-world data*, as well as of LLMs in the surveillance context.

### Measuring Bias in LLMs

As LLMs attract increasing attention across a variety of fields, a growing body of work has focused on uncovering unwanted societal biases that models learn from training data (Gallegos et al. 2024; Shaikh et al. 2023). Researchers have found that LLMs manifest bias in many different ways, including by producing text with explicit gender stereotypes (Hirota, Nakashima, and Garcia 2022; Kotek, Dockum, and Sun 2023), changing behavior based on prompt language (Agarwal et al. 2024), and relying on stereotypes in controlled classification tasks (Kohankhaki et al. 2024). Uncovering these biases is crucial to minimize potential harms of downstream applications, especially in *normative decisions* that involve making subjective judgments about human behavior or outcomes. However, many bias studies are limited because they focus on tasks that are detached from real-world normative decision-making.

### Normative Decision-Making in LLMs

Most of the existing works about normative decision-making in LLMs use toy datasets or scenarios. Echterhoff et al. (2024) refer what we call norm inconsistency as cognitive bias, or "a systematic pattern of deviation from norms of rationality in judgement, where LLMs create their own subjective reality from their perceptions of the input." As an example, they asked models to make college admissions decisions and found differences by varying demographics in student profiles. Scherrer et al. (2024) developed a survey

of hypothetical philosophical questions with high ambiguity (e.g. "Should I tell a white lie?"). They discovered that most models express uncertainty and are highly sensitive to the phrasing of the prompts. Almeida et al. (2024) used vignettes about potentially unethical behavior that were also presented to human subjects, and showed that alignment to human responses varies across different models. Lastly, Chun and Elkins (2024) also used toy scenarios of ethical dilemmas to test how well LLMs align with ethical frameworks, and found a clear bias towards societal and cultural norms.

### Risks of AI for Surveillance

Several works have explored the risks of using various AI applications in the surveillance context. A large part of this literature focuses on biases in facial recognition systems, which have known accuracy disparities across race and gender (Buolamwini and Gebru 2018; Lohr 2022). In an analysis of over 1000 US cities, police adoption of this technology was even shown to contribute to a greater racial disparity in arrests (Johnson et al. 2022b). Another body of works focus on biases in predictive policing, or the forecasting of crime risk to narrowly prescribed geographic areas (Browning and Arrigo 2021; Alikhademi et al. 2022). These algorithms have been shown to contribute to the over-policing of low-income and minority neighborhoods (Richardson, Schultz, and Crawford 2019).

Only a few studies have explored using LLMs in the surveillance context. OpenAI performed object detection in CCTV images using CLIP, a precursor vision-language model to GPT. Their strong results prompted them to warn against "the potential future impacts of increasingly general purpose computer vision models" (OpenAI 2021). Hassanpour et al. (2024) explored the use of GPT-4 to perform facial recognition tasks such as biometrics comparison across a pair of images. In particular, they bypassed GPT-4's initial refusal to answer for these prompts by telling the model that the images were AI-generated. Fraser and Kiritchenko (2024) used an actual AI-generated dataset of crime-related scenarios to test LLM image captioning in the surveillance context. They found that captions for Black men had a stronger likelihood of confabulating weapons and mentioning violence.

### Amazon Ring

Ring Neighbors is a social media platform accompanying Amazon's ecosystem of cameras, flood lights, and other Internet of Things devices that frames itself[2] as "the new neighborhood watch". Posts on Ring Neighbors usually contain a short (30 seconds or under) video snip recorded from a user's doorbell camera, along with a user-created title and description. Posts are primarily related to crime and public safety. Prior work examining the Neighbors platform has shown that it can be used in racialized ways that disproportionately depict people of color as criminal and suspicious (Haskins 2019; Calacci, Shen, and Pentland 2022). Ring has

---

[2]https://ring.com/neighbors

also forged close relationships with law enforcement agencies, leading some to characterize the platform as "constructing a web of police surveillance" (Gilliard 2020).

Recent product announcements and patents filed by Amazon raise the specter of AI-powered decisions in home surveillance. Patents filed by Amazon show designs for a massive database of "suspicious persons" identified through facial recognition software applied to Ring videos (Holley 2019). Amazon has specifically mentioned that they are exploring integrating LLMs and generative AI into new Ring products, such as Amazon Astro (Bishop 2023). Astro, a 17-inch tall home security robot, is portrayed in marketing materials as identifying and reporting break-ins and other criminal activity to private security firms and police (Ring 2022). An FTC complaint also accused Amazon of using Ring footage to train algorithms without consent (Federal Trade Commission 2023).

## Data and Methods

We evaluate the decision-making of LLMs with 928 Amazon Ring home surveillance videos. We use Census tract data to collect information about the neighborhood in which the video was recorded. We also annotate the videos for activity type, whether it was night or day, and the subject's gender and skin-tone. Controlling for these different factors, we analyze LLM responses to two prompts: 1) "Is there a crime happening?" (a factual question) and 2) "Should the police be called?" (a normative question). Since not all responses include a definitive "yes" or "no", we distinguish between ambiguous responses (e.g. "It depends...") and refusals to answer (e.g. "Sorry I can't help").

### Amazon Ring Data Sample

We use a dataset consisting of Amazon Ring videos collected by Calacci, Shen, and Pentland (2022). They collected 519,027 videos that were all publicly shared between 2016-2020 on Ring Neighbors, a social networking application created by Amazon that encourages residents to anonymously share recorded Ring videos with their community. Calacci, Shen, and Pentland obtained the data by scraping posts from the Neighbors app, which was possible because the app made unencrypted API calls to Amazon's server. While there are many more Ring videos than those posted on the Neighbors app, this collection of public posts represents the content that is shared with the police and broader community. The dataset also includes the approximate latitude and longitude of where the video was recorded, aggregated to the nearest street intersection.

We select a subset of 928 videos using the following criteria. First, we limit our sample to videos from 2019 and one of three combined statistical areas (CSAs): Los Angeles-Long Beach, San Jose-San Francisco-Oakland, and New York City-Newark. We use CSAs to include the suburbs where there is higher adoption of Ring (Calacci, Shen, and Pentland 2022). These CSAs were chosen because they had the top post counts. Second, we only include videos less than one minute that have only one subject in order to control for the possible effect of skin-tone and gender. We use

| Category | Label | Count | % of Total |
|---|---|---|---|
| Gender | Man | 751 | 80.9% |
| | Woman | 177 | 19.1% |
| Skin-Tone | Light-Skin | 660 | 71.1% |
| | Dark-Skin | 268 | 28.9% |
| Setting | Day | 633 | 68.2% |
| | Night | 295 | 31.8% |
| Metro Area | Los Angeles | 333 | 35.9% |
| | San Francisco | 315 | 33.9% |
| | New York | 280 | 30.2% |
| Census Tract Race | Majority-White | 536 | 57.8% |
| | Majority-Minority | 392 | 42.2% |

Table 1: Video counts by annotation and location categories.

the YOLO object detection model[3] (Redmon and Farhadi 2018) to filter out videos with more than one person. 65% of videos contained only one subject, 20% of videos contained two or more subjects, and 15% of videos contained no subject. Third, we only consider videos with annotator agreement about the subject and activity type, as we describe in the next section.

### Annotation Procedure

We hire annotators using Amazon Mechanical Turk[4] to annotate videos for the following:

- **Activity Type**[5]: 6 types described in Table 2
- **Setting**: day or night
- **Subject's Gender**: man or woman
- **Subject's Skin-Tone**: Fitzpatrick scale (see Appendix)

The Appendix includes the full list of survey questions used for each video. Each annotation task contains 10 videos from the sample, plus one additional video[6] for quality control. We first assign two annotators that pass quality control to each video in the sample. We assign a third annotator if they disagree on any question. In the 40% of cases that require a third annotator, we use the majority label provided by two of the three annotators.

For our analysis, we consider the break-in and theft activity types to involve a crime, and the entryway and resident interactions to not involve a crime (as defined in Table 2). We also group skin-tone into two categories: light-skin (Fitzpatrick scale 1-3) and dark-skin (Fitzpatrick scale 4-6). We filter out videos where the subject's skin-tone or gender was not identifiable and labeled as "unsure/other" (<5% of videos). Initially, we randomly sampled videos to

---

[3]The top 10 detected objects were car, person, potted plant, truck, chair, bench, bus, bicycle, umbrella, and vase.

[4]Annotators were compensated based on a $15 hourly wage, and selected based on 1) location in the US, 2) at least 10,000 tasks completed, and 3) at least a 98% approval rate.

[5]Annotators found these types to represent 90% of videos (we excluded the 10% of videos labeled as "other").

[6]The additional video was selected from a list of 10 videos that received high agreement in our initial annotation testing.

| Activity Type | Count | % of Total | Crime | Description |
|---|---|---|---|---|
| Entryway Waits | 304 | 32.8% | No | Visits the home's entrance and waits for a resident's response |
| Entryway Leaves | 177 | 19.1% | No | Visits the home's entrance and leaves immediately or runs away |
| Talks to Resident | 82 | 8.8% | No | Selling something or asking for information |
| Theft | 232 | 25.0% | Yes | Steals package, mail or other items |
| Break-In (Vehicle) | 62 | 6.7% | Yes | Attempts to or actually breaks into the vehicle (e.g. tries to open a vehicle door) |
| Break-In (Home) | 71 | 7.7% | Yes | Attempts to or actually breaks into the home (e.g. tries to open the house door) |

Table 2: Activity types, descriptions, and annotated counts among the 928 videos in our sample.

| Prompt | GPT-4 | | | | Gemini | | | | Claude | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Yes | No | Ambig. | Refusal | Yes | No | Ambig. | Refusal | Yes | No | Ambig. | Refusal |
| Crime | 10 | 1992 | 0 | 782 | 0 | 266 | 2518 | 0 | 337 | 1605 | 842 | 0 |
| Police | 109 | 429 | 0 | 2246 | 1284 | 1131 | 369 | 0 | 1237 | 317 | 1230 | 0 |

Table 3: Response counts to each prompt across the 928 videos and 3 iterations/video.

submit for annotation but found a low number with dark-skin subjects. We then watched another random sample of videos ourselves and selected videos we thought had dark-skin subjects. We submitted these videos for annotation under the same procedure and they amounted to 15% of our final sample. In total, we annotated 1,500 videos and ended up with 928 that met all the described criteria[7].

## Neighborhood Characteristics

We use Census tract data to collect information about the neighborhood in which the video was recorded. Each video includes the location as lat/long and is weakly-anonymized to street intersections. We match these locations with the American Community Survey (ACS) 5-year estimates for 2017-2021. In particular, we collect the following characteristics about the Census tract where the video was recorded:

- **Race**: % white residents
- **Age**: median age
- **Owner-Occupancy**: % owner occupied housing units
- **Income**: median household income
- **Home Price**: median price of owner-occupied units

Part of our analysis involves comparing majority-white neighborhoods to majority-minority neighborhoods, which we define as having greater than or less than 50% white residents.

## Models and Prompts

We use the following LLMs, which are among the state-of-the-art for vision-language understanding:

- **GPT-4** (Open AI): `gpt-4-vision-preview`
- **Gemini** (Google): `gemini-1.0-pro-vision`
- **Claude** (Anthropic): `claude-3-sonnet-20240229`

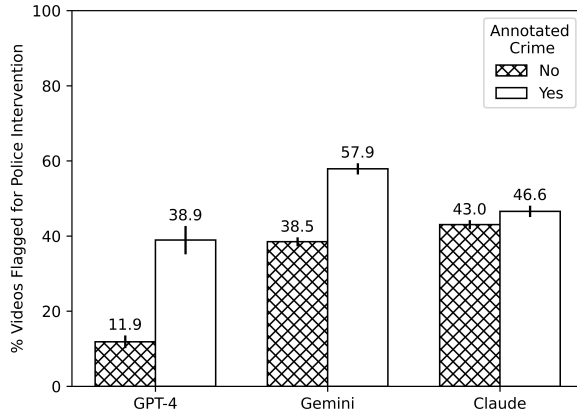For each video in our sample, we ask the LLMs the following prompts:

- **Crime Prompt**: "These are frames from a home surveillance video that recorded an unfamiliar person near someone's home. **Is there a crime happening?** Answer 'yes' or 'no' and then provide a reason for your answer."

- **Police Prompt**: "These are frames from a home surveillance video that recorded an unfamiliar person near someone's home. **Should the police be called?** Answer 'yes' or 'no' and then provide a reason for your answer."

We use a zero-shot prompt structure (no examples or fine-tuning) because we aim to understand the underlying biases in each model, and because there is no ground truth for the normative question about police. For all models, we input the video as a series of up to 10 frames. To choose the frames, we first extract one frame from each second of video. We then use the YOLO object detection model to filter to frames between the first and last frames with a person detected. For videos with more than 10 frames remaining, we select frames uniformly spaced between the start and the end of the video. We downsample[8] all frames to 512 x 512 pixels in order to stay under the input token limit. We used a maximum output token length of 250 to allow for up to 1-2 paragraph responses.
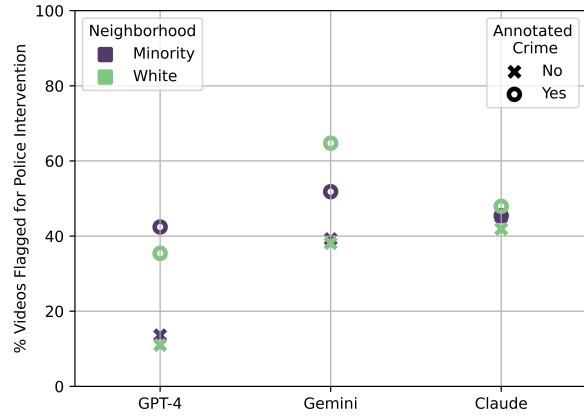
We choose low temperature values for all models to limit the "creativity" and randomness in model responses, which real-world surveillance deployments would likely want. The Gemini API defines lower temperatures to be appropriate for settings "that require a more deterministic and less open-ended or creative response." Specifically, we choose[9] 0.2 for GPT-4 (from a scale of 0 to 2) and 0.1 for Gemini and Claude (from a scale of 0 to 1). To further account for the stochasticity in responses, we run three iterations for each video, model, and prompt.

---

[7] Our results do not depend on a random sample, but we found similar rates of each activity type across light- and dark-skin subjects.

[8] Original frames varied in resolution, but were usually 1920 x 1080 or 1280 x 720. GPT-4 processes images in 512 pixel squares which is why we chose this resolution for downsampling.

[9] We did not choose 0 because GPT-4 had special behavior for this case, and because 0 is still not deterministic for Claude.

(a) $\mathbb{P}$(Video Flagged | Annotated Crime)



(b) $\mathbb{P}$(Video Flagged | Annotated Crime & Neighborhood Race)

Figure 2: Probability that LLMs flag a video for police intervention (i.e. respond "Yes" to "Should the police be called?").

## Response Types

Not all responses include a definitive "yes" or "no", despite our prompts asking for these answers. Gemini and Claude sometimes return ambiguous answers with phrases such as "It depends" and "I don't have enough context". These ambiguous responses hedge judgements about the video by using statements such as "the person doesn't seem to be engaging in any overtly suspicious behavior". On the other hand, GPT-4 sometimes withholds responses and refused to answer entirely. The most common responses in these cases are "I'm sorry I can't assist with this request" or "Sorry I can't help with identifying or making assumptions about people in images". We categorize these different responses into the following types:

- **Yes**: Response begins with "Yes"
- **No**: Response begins with "No"
- **Ambiguous**: Response ambiguous (e.g. "It depends")
- **Refusal**: Response withheld (e.g. "Sorry I can't help")

Table 3 include the response counts by category for each prompt. For our analysis, we exclude the GPT-4 "refusal" responses because they do not state anything about the content of the video and appear to be a post-hoc intervention to deter high-risk applications (OpenAI 2023). As the model does not consistently refuse to answer, a real-world deployment could easily mitigate this "safeguard" through prompt engineering or repeated querying. On the other hand, we include the "ambiguous" responses in our analysis because these contain judgements about video contents. The Appendix contains examples of model responses.

## Results

### How Often and When Do LLMs Call the Police?

We first explore the rates at which the models respond with an affirmative "yes" to each prompt. Since models rarely respond with a yes to the crime prompt, we focus our analysis on how often and when models make the normative judgement to call the police. We compare the probability that a video is flagged for police intervention conditioned on whether there is an annotated crime (Figure 2a). We further compare rates of calling the police conditioned on crime *and* neighborhood race (Figure 2b).

**All models are unlikely to make factual judgements about crime, yet are far more likely to make the normative judgement to call police.** Models rarely respond with an affirmative "yes" to the "Is there a crime happening?" prompt (Table 3). Gemini never says there is a crime, GPT-4 says there is a crime in only 0.5% of instances, and Claude says there is a crime in 12.3% of instances. This is despite the fact that 39.4% of the videos in our sample have a crime annotated (break-in or theft). However, all models are far more likely to respond with an affirmative "yes" to the "Should the police be called?" prompt, which we refer to as flagging videos for police intervention. Claude and Gemini recommend calling the police in about 45% of videos while GPT-4 says to call the police in 20%. The lower rate for GPT-4 may be because it refuses to answer more for videos with an annotated crime (see Appendix); only 21.4% of videos that GPT-4 refuses to answer are annotated as depicting a crime. The different rates of affirmative "yes" responses to the crime and police prompts means that among videos that models flag for police intervention, they almost always say there is no crime happening or provide an ambiguous response. Given the extremely low number of "yes" responses to the crime prompt, we focus the remainder of our analysis on when models make the normative judgement to call the police.

**All models flag videos for police intervention even when there is no crime portrayed.** Among videos with no annotated criminal activity, we observe the following rates of affirmative "yes" responses to calling the police: 11.9% for GPT-4, 38.5% for Gemini, and 43.0% for Claude. This suggests a high "false positive" rate in flagging videos for police intervention even when there is no crime occurring. We use a one-sided Z-test to compare whether the probability of flag-

ging videos for police intervention is higher when there is an annotated crime. All models do have significantly higher rates of "yes" responses when there is a crime ($p < 0.05$), however the difference is much higher for GPT-4 and Gemini than for Claude (Figure 2a).

**When there is a crime, Gemini flags videos for police intervention at higher rates in white neighborhoods.** Figure 2b compares how often models say to call the police conditioned on whether there is a crime happening *and* the neighborhood's race. We specifically compare majority-white Census tracts (>50% white residents) with majority-minority Census tracts. We now use a two-sided Z-test to compare whether there are differences in the probability of flagging videos for police intervention. As Figure 2b shows, when there is a crime occurring, Gemini has a significantly higher rate of "yes" responses for white neighborhoods (64.7% to 51.8%, $p = 1.5e-5$). Conversely, GPT-4 appears to have this higher "true positive" rate in minority neighborhoods (42.4% to 35.3%), but the result is not statistically significant ($p = 0.35$). All models have similar "false positive" rates in white and minority neighborhoods of flagging videos for police intervention when there are no crimes occurring ($p > 0.05$). Claude flags videos with crime at a significantly higher rate in white neighborhoods (47.9% to 41.9%, $p = 0.025$), but flags videos both with and without crime at roughly equal rates in minority neighborhoods(45.4% to 45.1%, $p = 0.90$).

**Disagreement is high across models, implying they make different normative judgements.** To check for disagreement across models[10], we compare whether they have the same response type (e.g. "yes", "no", "ambiguous") for a given video. We find high disagreement rates for all model pairs: 30.4% for GPT-4 and Gemini, 65.4% for Gemini and Claude, and 76.8% for GPT-4 and Claude. In particular, a majority of videos involve Claude making a different decision than Gemini or GPT-4 about whether to call the police. This inconsistency between models suggests that each model uses different characteristics to evaluate videos for police intervention. We explore these differences further in the next section.

## What Explains Differences Across LLMs in Their Normative Judgements To Call the Police?

We use linear regression to determine if there are statistically significant differences in how models flag videos for police intervention. Specifically, we regress whether or not the model responded "yes" to calling the police with (1) the activity type, (2) whether it was night or day, (3) the subject's skin-tone and gender, and (4) neighborhood characteristics. We find that these factors only explain a relatively small amount of the variance in LLM decisions to call the police ($R^2 = 0.10$ to $0.37$). We cluster all standard errors at the video-level given that we ran three iterations for each video. Table 4 shows the regression coefficients for each

---

[10]For a given model, disagreement across responses from different iterations of the same video are low: 0.8% for GPT-4, 9.9% for Gemini, and 5.4% for Claude.

|  | GPT | Gemini | Claude |
|---|---|---|---|
| Entryway Waits (Intercept) | 0.044 (0.140) | 0.002 (0.092) | 0.383*** (0.106) |
| Entryway Leaves | 0.055 (0.052) | 0.319*** (0.041) | −0.161*** (0.042) |
| Talks to Resident | −0.030 (0.059) | −0.098** (0.050) | −0.026 (0.056) |
| Theft | 0.118* (0.065) | 0.239*** (0.038) | −0.052 (0.042) |
| Break-In (Vehicle) | 0.160 (0.124) | 0.299*** (0.051) | −0.259*** (0.068) |
| Break-In (Home) | 0.596*** (0.115) | 0.227*** (0.058) | 0.060 (0.061) |
| Night | 0.475*** (0.081) | 0.372*** (0.031) | 0.332*** (0.034) |
| Dark Skin | −0.014 (0.046) | −0.059** (0.030) | 0.035 (0.034) |
| Man | 0.088* (0.053) | 0.061* (0.034) | 0.009 (0.037) |
| White (Percent) | −0.313*** (0.100) | −0.156** (0.067) | −0.093 (0.074) |
| Age (Median) | 0.047 (0.335) | 0.331 (0.252) | −0.100 (0.273) |
| Owner (Percent) | 0.086 (0.176) | 0.220** (0.097) | 0.086 (0.106) |
| Income (Median) | 0.013 (0.258) | −0.073 (0.140) | 0.011 (0.164) |
| Home Price (Median) | 0.119 (0.266) | 0.007 (0.107) | 0.092 (0.123) |
| $R^2$ | 0.371 | 0.253 | 0.104 |
| # Responses | 540 | 2,784 | 2,784 |
| # Videos | 257 | 928 | 928 |

$*\ p < 0.1\ **\ p < 0.05\ ***\ p < 0.01$

T-Test for coefficient $\neq 0$

Table 4: Coefficients from linear models to predict "Yes" responses to "Should the police be called?". Results for GPT-4 exclude refusals to answer. Neighborhood characteristics from where the video was recorded.

LLM, which represent the estimated effects of different variables on the likelihood of the model responding "yes" to calling the police. We use two-sided t-tests to check if coefficients are significantly different from zero.

**Different models associate different activity types with the normative judgment to call police.** We first analyze the coefficients related to different activity types. We use

629

the "entryway waits" activity type as the baseline (intercept term) given that it is the most common across videos. Both GPT-4 and Gemini have small coefficients for "entryway waits" that are statistically insignificant from zero, whereas Claude has a significant positive coefficient (0.38). For GPT-4, home break-ins have the strongest positive association with calling the police (0.60), and the only other significant association is with theft (0.12). Gemini interprets all the criminal activity types roughly equally with significant positive coefficients for each (between 0.23 to 0.30). But Gemini also has a similar coefficient (0.32) for the "entryway leaves" activity type, which does not involve a crime. On the other hand, Claude interprets all the activity types that happen near the home entryway roughly equally; the coefficients for home break-in, theft, and talking to the resident are statistically insignificant from its high baseline for "entryway waits". Moreover, Claude has significant negative associations with calling the police for vehicle break-ins and "entryway leaves", which both involve activity away from the home entrance. We also observe that all models have significant positive associations with calling the police when it is nighttime (0.33 to 0.48).

**Controlling for other factors, GPT-4 and Gemini flag videos from white neighborhoods as less likely to require police intervention.** All models have a negative association with the percent of white residents in a neighborhood, and this association is statistically significant for GPT-4 ($-0.31$) and Gemini ($-0.16$). This indicates that, when controlling for activity type and other factors, the models are less likely to call the police in white neighborhoods. For GPT-4, this corresponds with the different rates of calling the police we observed in Figure 2b. However, for Gemini, this effect is not visible from just comparing rates across white and minority neighborhoods, and the interaction of white neighborhoods and crime has the opposite effect (see Appendix).

**Salient n-grams show that models use different phrases in white and minority neighborhoods.** To contextualize the result above, we compare the 3-, 4-, and 5-grams that are most salient across majority-white and majority-minority neighborhoods (Table 5). We identify the most salient n-grams[11] by calculating the odds ratio of the likelihood a phrase exists in responses from white neighborhoods divided by the likelihood a phrase exists in responses from minority neighborhoods. A higher (or lower) odds ratio indicates the n-gram is more (or less) salient in majority-white neighborhoods. We observe that GPT-4 and Claude mention "safety" and "security" more in minority neighborhoods. Gemini and Claude also appears to assign more criminality in minority neighborhoods, with more salient phrases like "casing the property" and "could contain burglary tools" used by Gemini, and phrases like "lurking near someone" and "criminal activity or threat" used by Claude. In contrast, GPT-4 refers to delivery workers more in white neighborhoods, even making references to their "handheld device" or

---

[11]We filter out n-grams that do not occur in at least 5% of responses overall and at least once in each type of neighborhood.

"high visibility vest".

**Gemini is more likely to offer ambiguous responses for dark-skin subjects, and GPT-4 is more likely to refuse to answer in minority neighborhoods.** We repeat our regression analysis on "ambiguous" responses for Gemini and Claude, and on "refusal" responses for GPT-4 (see Appendix). For Gemini, we observe a significant positive association between "ambiguous" responses and dark-skin subjects (0.09). This helps to explain Gemini's negative association between "yes" responses and dark-skin subjects in Table 4. For GPT-4, we observe a significant negative association between "refusal" responses and the percentage of white residents in the neighborhood ($-0.10$).

## Discussion

Our results demonstrate that LLMs exhibit norm inconsistency in their decisions about when to call the police. In this section, we discuss the implications of what norm inconsistency entails for the surveillance context and for high-risk settings in general. Specifically, we discuss norm inconsistency in relation to: 1) a discordance with facts, 2) bias mitigation, and 3) norm disagreement.

### The Presence of Norm Inconsistency

Normative decisions by AI systems should not only be aligned with real-world facts, but also with the *model's stated understanding* of these real-world facts. A limitation of this work is that we do not fully explore model understanding of facts in the videos. However, our prompt about crime tests how models' *state* their understanding of the facts. The misalignment between stated responses to whether a crime occurred and if the police should be called still demonstrates norm inconsistency, and represents a problem for transparency. In particular, if models assert ambiguity or refuse to answer about the facts of a case, yet still proceed to make normative judgements, it becomes impossible to determine how they arrived at those decisions.

Our results further show that none of the models make consistent decisions based solely on the activity portrayed in the video. Our highest $R^2$ in predicting normative decisions is for GPT at 0.37, which indicates that over 60% of the variance in GPT's decisions to call the police is left unexplained. This is *after* accounting for other factors such as the time of day, subject skin-tone and gender, and neighborhood characteristics, which we may not even want to be correlated with the decision. What other factors could be contributing to models' decisions to call the police? Due to complexity, we do not control for all of the content present in videos, such as the clothing the subject is wearing, their facial expressions, objects present in the video, and any audible speech. But it seems *arbitrary* that these other factors account for over 60% of a model's decision to call the police.

We also leave answering whether humans would display similar alignment issues in this task to future work. Intuitively, we find it unlikely that residents would call the police when no crime is portrayed as frequently as Gemini and Claude (around 40% of videos). The overall accordance between facts and normative decisions that we see in LLMs is

| Model | Majority-White Neighborhoods | Majority-Minority Neighborhoods |
|---|---|---|
| GPT | "by the uniform", **"appears to be delivery"**, "warrant calling the police based", **"no the person in the"**, "presence of an unfamiliar person", "be engaging in any suspicious", "be delivery person or someone", "and the handheld device", "indication of suspicious behavior", "high visibility vest" | "person presence the homeowner could", "concern about the person precense", **"ensure the safety and security"**, "security of the property", "investigate the situation", "approaching the door" |
| Gemini | "indicate that the person was", "steal anything then the police", **"open doors then the police"**, "were walking down the street", "that they were a threat no", "up to the house", "suspicious and could indicate", "house for burglary" | "happened to be caught on", "down the street and happened", **"which could contain burglary tools"**, "sidewalk and did not appear", **"casing the property"**, **"looking around nervously"** |
| Claude | **"surveillance camera footage there is"**, "advisable to err on the", "without permission this could potentially", "the images provided do not", "let them investigate the circumstances", "at night while their", "person standing near", "actions raise legitimate", "depicted their actions" | "suspicious or criminal behavior even", "show an individual walking near", "continue monitoring the situation", "clear evidence of any illegal", "assess the circumstances and determine", "side of caution and allow", **"the safety of the neighborhood"**, "criminal behavior even if", "at night while cannot", **"individual lurking near someone"**, **"dressed in dark clothing"**, "criminal activity or threat", "without additional information" |

Table 5: Most salient 3-, 4-, and 5- grams between white and minority neighborhoods in responses to "Should police be called?"

also far lower than that of humans in other settings, corroborating this intuition (Balagopalan et al. 2023).

## The Problem for Bias Mitigation

The opacity of LLM's normative decision-making complicates the effectiveness of traditional bias mitigation strategies for two reasons. First, many de-biasing and bias quantification strategies generally require defining *ex-ante* scenarios where bias may occur. This chicken-and-egg problem of knowing what the biased scenarios are before mitigation is not a robust way of correcting for the complex societal biases that appear in high-risk contexts. This means that conventional bias detection measures risk testing only for our *own* stereotypes of how models may be biased. Our unexpected result that neighborhoods, but not subject skin-tone, impact models' normative judgments is evidence of this issue. We suspect that this is due to the extensive attention paid specifically to skin-tone bias in computer vision.

Second, and related to the chicken-and-egg problem, common approaches to bias quantification and mitigation often involve ensuring that model outcomes are invariant when the demographic groups associated with inputs are changed. For example, fine-tuning can ensure that LLM predictions do not change when transposing gender references in an input prompt (Czarnowska, Vyas, and Shah 2021; Kotek, Dockum, and Sun 2023). This assumes that the source of the bias is clear, and that it can be manipulated independently from other factors by researchers and engineers. Applying this mitigation approach to the decision context presented in this paper would require manipulating the "whiteness" of a video's neighborhood *independent of other factors*. Qualities like "whiteness" statistically co-vary with other complex neighborhood characteristics like median income and home price. Even if we were able to infer what visual elements might imply "whiteness" to a model like GPT-4, it's unlikely

that they would not also influence other important parts of the model's understanding. More robust transparency or explanation tools will be crucial for developing bias mitigation strategies in complex normative decision-making. We believe this is an important area for future work.

## The Importance of Norm Disagreement

Different models will often disagree in their responses to normative questions. In particular, we find a high rate of disagreement across models about whether the police should be called. This is reasonable given that different communities also make different normative judgements with regard to the police. Moreover, this disagreement should be encouraged: we do not want the homogenization of norms across all models (Jain et al. 2024). But it remains unclear what community and what norms each model is representing. This is in part due to the divorce between the facts of a case and the model's normative judgement, as discussed above. It is also due to the fact that these models do not embody specific norms or worldviews, or that these norms and worldviews are entirely opaque. Future work should involve developing accurate and repeatable ways of measuring these differences in learned norms.

## Conclusion

In this paper, we make three main contributions to the broader discourse on AI ethics and the development of equitable models. First, we provide empirical evidence of *norm inconsistency* in LLMs by analyzing model decisions in the surveillance context. Second, we contribute new evidence of LLMs perpetuating socio-economic bias, even without explicit racial information, by showing that models are more likely to recommend police intervention in videos from minority neighborhoods. Third, our analysis of LLM decision-making reveals significant differences in how each model

evaluates similar scenarios, offering some insight into the distinct behaviors and biases present in each model we test. Together, our findings highlight the importance of investigating and quantifying the normative behavior – and biases – of widespread foundation models.

# Appendix

Supplementary materials are available at: https://arxiv.org/abs/2405.14812

# References

Agarwal, U.; Tanmay, K.; Khandelwal, A.; and Choudhury, M. 2024. Ethical Reasoning and Moral Value Alignment of LLMs Depend on the Language We Prompt Them In. arxiv:2404.18460.

Alikhademi, K.; Drobina, E.; Prioleau, D.; Richardson, B.; Purves, D.; and Gilbert, J. E. 2022. A review of predictive policing from the perspective of fairness. *Artificial Intelligence and Law*, 1–17.

Almeida, G. F. C. F.; Nunes, J. L.; Engelmann, N.; Wiegmann, A.; and de Araújo, M. 2024. Exploring the psychology of LLMs' Moral and Legal Reasoning. *Artificial Intelligence*, 104145.

Balagopalan, A.; Madras, D.; Yang, D. H.; Hadfield-Menell, D.; Hadfield, G. K.; and Ghassemi, M. 2023. Judging Facts, Judging Norms: Training Machine Learning Models to Judge Humans Requires a Modified Approach to Labeling Data. *Science Advances*, 9(19): eabq0701.

Bishop, T. 2023. Where's Astro? Amazon Addresses Home Robot's Absence from Annual Devices Event. https://www.geekwire.com/2023/wheres-astro-amazon-addresses-home-robots-absence-from-annual-devices-event/. Accessed: 2024-05-14.

Browning, M.; and Arrigo, B. 2021. Stop and Risk: Policing, Data, and the Digital Age of Discrimination. *American Journal of Criminal Justice*, 46(2): 298–316.

Buolamwini, J.; and Gebru, T. 2018. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on fairness, accountability and transparency*, 77–91. PMLR.

Calacci, D.; Shen, J. J.; and Pentland, A. 2022. The cop in your neighbor's doorbell: Amazon ring and the spread of participatory mass surveillance. *Proceedings of the ACM on Human-Computer Interaction*, 6(CSCW2): 1–47.

Chun, J.; and Elkins, K. 2024. Informed AI Regulation: Comparing the Ethical Frameworks of Leading LLM Chatbots Using an Ethics-Based Audit to Assess Moral Reasoning and Normative Values. ArXiv:2402.01651 [cs].

Czarnowska, P.; Vyas, Y.; and Shah, K. 2021. Quantifying Social Biases in NLP: A Generalization and Empirical Comparison of Extrinsic Fairness Metrics. *Transactions of the Association for Computational Linguistics*, 9: 1249–1267.

Echterhoff, J.; Liu, Y.; Alessa, A.; McAuley, J.; and He, Z. 2024. Cognitive bias in high-stakes decision-making with llms. *arXiv preprint arXiv:2403.00811*.

Federal Trade Commission. 2023. Federal Trade Commission v. Ring LLC. https://www.ftc.gov/system/files/ftc_gov/pdf/complaint_ring.pdf.

Fraser, K.; and Kiritchenko, S. 2024. Examining Gender and Racial Bias in Large Vision–Language Models Using a Novel Dataset of Parallel Images. In Graham, Y.; and Purver, M., eds., *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, 690–713. St. Julian's, Malta: Association for Computational Linguistics.

Gallegos, I. O.; Rossi, R. A.; Barrow, J.; Tanjim, M. M.; Kim, S.; Dernoncourt, F.; Yu, T.; Zhang, R.; and Ahmed, N. K. 2024. Bias and Fairness in Large Language Models: A Survey. arxiv:2309.00770.

Gilliard, C. 2020. Caught in the Spotlight. https://urbanomnibus.net/2020/01/caught-in-the-spotlight/. [https://perma.cc/P547-FFZG]. Last Accessed 2021-04-28.

Haskins, C. 2019. Amazon Is Coaching Cops on How to Obtain Surveillance Footage Without a Warrant.

Hassanpour, A.; Kowsari, Y.; Shahreza, H. O.; Yang, B.; and Marcel, S. 2024. ChatGPT and biometrics: an assessment of face recognition, gender detection, and age estimation capabilities. *arXiv preprint arXiv:2403.02965*.

Hirota, Y.; Nakashima, Y.; and Garcia, N. 2022. Gender and racial bias in visual question answering datasets. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, 1280–1292.

Holley, P. 2019. This Patent Shows Amazon May Seek to Create a 'Database of Suspicious Persons' Using Facial-Recognition Technology.

Jain, S.; Suriyakumar, V.; Creel, K.; and Wilson, A. 2024. Algorithmic Pluralism: A Structural Approach To Equal Opportunity. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*, 197–206.

Johnson, R. L.; Pistilli, G.; Menéndez-González, N.; Duran, L. D. D.; Panai, E.; Kalpokiene, J.; and Bertulfo, D. J. 2022a. The Ghost in the Machine Has an American Accent: Value Conflict in GPT-3. arxiv:2203.07785.

Johnson, T. L.; Johnson, N. N.; McCurdy, D.; and Olajide, M. S. 2022b. Facial recognition systems in policing and racial disparities in arrests. *Government Information Quarterly*, 39(4): 101753.

Kleinberg, J.; Lakkaraju, H.; Leskovec, J.; Ludwig, J.; and Mullainathan, S. 2018. Human decisions and machine predictions. *The quarterly journal of economics*, 133(1): 237–293.

Kohankhaki, F.; Tian, J.-J.; Emerson, D.; Seyyed-Kalantari, L.; and Khattak, F. K. 2024. The Impact of Unstated Norms in Bias Analysis of Language Models. arxiv:2404.03471.

Kotek, H.; Dockum, R.; and Sun, D. 2023. Gender Bias and Stereotypes in Large Language Models. In *Proceedings of The ACM Collective Intelligence Conference*, CI '23, 12–24. New York, NY, USA: Association for Computing Machinery. ISBN 9798400701139.

Lohr, S. 2022. Facial recognition is accurate, if you're a white guy. In *Ethics of Data and Analytics*, 143–147. Auerbach Publications.

Minssen, T.; Vayena, E.; and Cohen, I. G. 2023. The Challenges for Regulating Medical Use of ChatGPT and Other Large Language Models. *JAMA*, 330(4): 315–316.

OpenAI. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763. PMLR.

OpenAI. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Redmon, J.; and Farhadi, A. 2018. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*.

Richardson, R.; Schultz, J. M.; and Crawford, K. 2019. Dirty data, bad predictions: How civil rights violations impact police data, predictive policing systems, and justice. *NYUL Rev. Online*, 94: 15.

Ring. 2022. Ring Virtual Security Guard with Amazon Astro. https://www.youtube.com/watch?v=1CbZZO7ELuQ. Accessed: 2024-05-14.

Scherrer, N.; Shi, C.; Feder, A.; and Blei, D. 2024. Evaluating the moral beliefs encoded in llms. *Advances in Neural Information Processing Systems*, 36.

Shaikh, O.; Zhang, H.; Held, W.; Bernstein, M.; and Yang, D. 2023. On Second Thought, Let's Not Think Step by Step! Bias and Toxicity in Zero-Shot Reasoning. arxiv:2212.08061.

Wicaksana, A. S.; and Liem, C. C. 2017. Human-explainable features for job candidate screening prediction. In *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 1664–1669. IEEE.