

Investigating LLM Bias in Housing Contexts

Charisse Hao
chao@ucsd.edu

Jenna Canicosa
jcanicosa@ucsd.edu

Joseph Guzman
j4guzman@ucsd.edu

Lana Murray
lmurray@ucsd.edu

Mentor: Stuart Geiger
sgeiger@ucsd.edu

Abstract

As the United States housing crisis intensifies, competition for housing has surged, increasing the likelihood that landlords and other decision-makers will turn to technological tools, including large language models (LLMs) like ChatGPT, to aid in their decision-making processes. However, these models often reflect and perpetuate societal biases present in their training data, potentially influencing critical housing-related decisions, such as eligibility assessments, tenant screening, and eviction risk evaluations. This study systematically examines how biases manifest in LLM-generated responses to housing-related prompts, focusing on disparities across race, gender, economic status, and other factors. Using an adapted algorithm audit framework, we generate prompts with varying demographic details to assess potential biases in LLM outputs. These prompts are submitted to selected models, with results analyzed using statistical methods. The findings are presented through a report, poster, and interactive website, allowing users to explore bias patterns firsthand. By investigating the impact of different candidate characteristics on LLM responses, this project contributes to the discourse on ethical AI implementation in housing, promoting fairer and more accountable decision-making.

Code: https://github.com/CharisseHao/retail_hiring_bias_audit.git

1	Introduction	2
2	Methods	5
3	Results	9
4	Discussion	26
5	Conclusion	28
6	Contributions	29
7	Appendix: Project Proposal	30
	References	32

1 Introduction

The housing crisis across the United States has become increasingly dire, with rising housing prices and growing demand for affordable housing. In 2023, 49.7 percent of renting households were “cost-burdened”, spending more than 30 percent of their income on housing costs [Desilver \(2024\)](#). As a result, the housing application process has become highly competitive in large metropolitan areas, with many units receiving over ten applicants [Grecu \(2024\)](#). Simultaneously, large language models (LLMs), such as ChatGPT, have gained popularity as decision-making tools utilized by both individuals and businesses. However, since LLMs are trained on historical datasets that often reflect societal biases related to race, gender, class, etc, their responses can inadvertently perpetuate these biases [Dethmann and Spiekermann \(2024\)](#). This poses risks in life-changing areas like employment and housing, where biased outputs could negatively influence critical decisions. Recognizing and understanding the implications of relying on these technologies in such contexts is essential.

This project explores how biases manifest in LLM-generated responses, specifically within the context of the ongoing housing crisis—a critical issue shaping many current policies. The goal is to develop housing-related prompts that reflect the perspective of ordinary individuals who rely on LLM feedback for decisions such as identifying suitable housing options, determining eligibility for programs, or selecting tenants as landlords. To ensure relevance and inclusivity, these prompts are designed based on personal insights and public feedback collected through interviews. By analyzing LLM-generated responses to these prompts, this research investigates potential biases and discrepancies, focusing on identifying the groups most affected and understanding how these biases influence outcomes. Understanding these discrepancies will uncover the mechanisms behind LLM decision-making and provide insights into their broader societal impact, particularly on vulnerable populations.

Building on previous algorithm audit research, this study aims to address a notable gap in LLM auditing within the housing sector by exploring the distribution of tenant scores and identifying statistically significant biases across variables such as gender, race, occupation, living status, credit score, and eviction history. Established frameworks and methodologies will be employed to investigate potential biases and discrepancies in LLM-generated tenant screening scores. To refine the project scope and better understand the housing application process, interviews were conducted with students, professionals, and housing-related organizations. The work culminates in quantitative analysis, using LLM responses as primary data. Statistical assumptions will be checked before applying appropriate methods—parametric methods like ANOVA and t-tests if assumptions are met or non-parametric alternatives such as the Kruskal-Wallis test if they are violated. Significant results will be further analyzed using Dunn’s test, with findings visualized through heatmaps and boxenplots. The research agenda is structured as follows:

- **Prompt 1**
 - **RQ1.1:** Were there any significant differences in tenant scores revealing bias between genders for each model?
 - **RQ1.2:** Were there any significant differences in tenant scores revealing bias between races for each model?

- **RQ1.3:** Were there any significant differences in tenant scores revealing bias between occupations for each model?
- **RQ1.4:** Were there any significant differences in tenant scores revealing bias between living statuses for each model?
- **RQ1.5:** Were there any significant differences in tenant scores revealing bias when comparing the intersection of variables?
- **Prompt 2**
 - **RQ2.1:** Were there any significant differences in tenant scores revealing bias between genders for each model?
 - **RQ2.2:** Were there any significant differences in tenant scores revealing bias between races for each model?
 - **RQ2.3:** Were there any significant differences in tenant scores revealing bias between eviction histories for each model?
 - **RQ2.4:** Were there any significant differences in tenant scores revealing bias between credit scores for each model?
 - **RQ2.5:** Were there any significant differences in tenant scores revealing bias when comparing the intersection of variables?

1.1 Relevant Literature

Independent algorithm audits play a crucial role in identifying significant biases and holding developers accountable for the models they deploy. Although LLMs and Artificial Intelligence (AI) systems have been painted as “neutral”, many argue that algorithms used in everyday life are oppressive and produce results that further social inequality ([Noble 2018](#), 1-5). This is becoming more worrisome as these algorithms are increasingly being utilized for decision-making across sectors, from approving loans [Meissner and Narita \(2023\)](#) to processing insurance claims [NAI \(2025\)](#). With tangible impacts on human lives, algorithm audits are vital for raising public awareness and promoting algorithmic accountability ([Geiger et al. 2024b](#), 645). In a recent audit conducted on four ChatGPT large language models on the topic of salary negotiation, significant bias was found across gender, major, and university variables, highlighting how LLMs can inadvertently perpetuate inequities even when provided with identical prompts but varied candidate characteristics ([Geiger et al. 2024a](#), 1). These findings underscore the importance of auditing widely used LLMs to identify and mitigate bias, particularly in life-impacting areas like employment and housing. Although audits do not always result in meaningful change ([Geiger et al. 2024b](#), 645), conducting and publishing audits allows researchers and concerned stakeholders to refine frameworks to test for bias and develop policy recommendations to reduce the impact of discrimination produced by AI systems.

Housing, in particular, is a field where biases reproduced by LLMs can significantly impact people’s access to safe and affordable homes. The U.S. housing market is often described as a “landlord’s market,” with landlord behavior largely unchecked [Reosti \(2020\)](#). This creates opportunities for discriminatory practices, which may be amplified by landlords using open-source LLMs for tenant screening. According to the Fair Housing Act, it is illegal

to discriminate against potential renters based on protected characteristics, such as race, sex, family status, etc [Hou \(2023\)](#). However, if this information is not redacted when providing a tenant application to an LLM for tenant screening, discrimination may still occur. Along with this, more sophisticated tenant screening software may be starting to utilize LLMs to assist in selecting prospective tenants, as LLMs are being increasingly employed to improve recommendation systems by integrating statistical modeling with language analysis [Wei \(2023\)](#). Typically, landlords make tenant decisions based on data such as credit reports, employment information, criminal history, rental history, etc [Manolas \(2024\)](#). Tenant screening tools and software will use this information to assign potential tenants points or scores based on perceived risk [Inn \(2025\)](#). However, without proper regulation or transparency of the input data and scoring metrics used in these systems, personal characteristics could still be inferred and limit marginalized people's opportunities [Hill and Holloway \(2024\)](#). This audit will seek to examine the ways in which LLMs used to provide tenant risk scores can be biased based on both inferred and given personal information.

There have been multiple relevant studies exploring the impact of tenant screening tools in housing. A 2022 study by Matthew Liewabnt, explored the limitations of algorithmic tenant screening in relation to eviction information. Liewabnt argued that these algorithms fail to account for contextual factors, are prone to errors, and often penalize tenants who have interacted with eviction courts, even if the outcome was not an eviction. Black women in particular were found to be disproportionately affected by these biases ([Leiwant 2022](#), 282-284). He advocates for updating tenant screening regulations to better align with the Fair Housing Act, aiming to reduce discrimination in the housing market. In 2024, a team of researchers examined how "digit risk-profiling" in England's rental market has shaped housing market dynamics [Wallace et al. \(2025\)](#). They argue the use of algorithmic systems to assess tenants are producing an "ordinal tenant", resulting in those who do not quite fit into the norm based on their online data at risk of losing opportunities. This only furthers the housing disparity and complicates housing access for marginalized groups. Despite these findings, there remains limited research on the specific impact of LLMs on housing decisions. However, a team of MIT researchers in 2024 conducted an audit of ChatGPT-4 to explore biases related to gender, race, ethnicity, nationality, and language in housing selection. Their research revealed significant biased responses from the model, indicating potential disparities in housing decisions influenced by LLMs [Liu et al. \(2024\)](#). Although their study provided valuable insights, it focused solely on ChatGPT-4, and the rapidly evolving nature of LLMs means that their findings may not apply universally. This study will expand on their work by testing a variety of LLM models to explore the generalizability of their results and gain a deeper understanding of the factors affecting housing responses provided by LLMs.

2 Methods

2.1 Models

For our initial prompt, we selected a diverse set of LLMs to gain a more comprehensive understanding of potential biases. Our chosen models included Google’s gemma-2-2b-it, OpenAI’s gpt-3.5-Turbo-0125, gpt-4o-2024-08-06, and gpt-4o-Mini-2024-07-18, InceptionAI’s Jais-Family-1P3B-Chat, and Meta’s Meta-Llama-3-8B-Instruct. These models were chosen for their varying capabilities in generating responses without rejecting our prompts, ensuring we had sufficient data for analysis. Additionally, their open-source nature and widespread recognition suggest they are likely candidates for adoption by landlords or other tenant evaluation software. After reviewing the results obtained from the first prompt, detailed under Data Cleaning and Exploratory Data Analysis, we adjusted our selection by replacing InceptionAI’s Jais-Family-1P3B-Chat with Microsoft’s Phi-3-mini-4k-instruct and Meta’s Llama-3.2-3B-Instruct while keeping the remaining models unchanged for the second prompt.

2.2 Prompt Generation and Submission

Prompt engineering is the first step when conducting algorithm audits on Large Language Models, ensuring prompts contain relevant information and produce responses in a usable format. For each set of prompts, a prompt bulk generator script was used to create thousands of copies with the given input variables systematically changed. Input variables were chosen for their key role in housing decisions, which may reveal biased outputs. Testing multiple variables at once allows for intersectional analyses of the LLM’s responses to be conducted in the data analysis stage. In each prompt, LLMs were instructed to return only a numerical score from 0-100 without explanation. This not only makes it easier to run statistical tests on the models’ responses, but simulates the scoring process that tenant evaluation programs may use [Desai \(2024\)](#). Batchwizard was used to submit prompts to OpenAI’s ChatGPT models, while Runpod was used for all other models. From there, the response files were downloaded to start the data cleaning process, checking for answer validity and preparing the data for analysis.

2.3 Interviews

To inform the variables and structure of our prompts, we attempted to conduct interviews with relevant community members, including current renters and housing sector professionals. We were able to conduct three interviews with student renters, which helped us understand the various application processes that can occur, along with current housing concerns for this population. Two tenant application formats were identified: an email for more informal applications and a formal application form. The interviewees discussed including information such as their name, living status, occupation, and credit score. Back-

ground checks were not discussed despite their prevalence in formal housing applications. Unfortunately, due to time constraints and a low response rate, no interviews with professionals were conducted. However, from our research and a few email replies received, it can be inferred that the use of AI and LLMs for tenant screening is not yet a widely known issue. We hope further research will be conducted in this area, including perspectives of renters, who would be most impacted by biased tenant selection.

2.4 Prompt 1: Apartment Inquiry

Prompt generation began by defining input variables that combined different attribute combinations. The primary variables in this experiment—‘occupation’, ‘living status’, and ‘name’—were selected for their potential to reveal disparities in model responses. ‘Occupation’ serves as a proxy for socioeconomic status, which landlords may consider when assessing financial stability. Eleven occupations representing a range of financial statuses were included: “Doctor”, “Software engineer”, “Accountant”, “Teacher”, “Retail Associate”, “Construction Worker”, “Food Service Worker”, “College Student”, “Government Worker”, “Unemployed”, and a control condition where no occupation was specified. The ‘living status’ variable included “just myself”, “my roommate and I”, “my pet and I”, “my spouse and I”, “my family with kids,” and again a control condition. This variable was included to examine potential bias against different family units. Finally, ‘name’ encodes race and gender without explicitly stating these attributes, a technique that has been effective in revealing implicit bias in previous studies ([Salinas, Haim and Nyarko 2024](#), 9-10). Three prompts were generated for each of the variable combination across six different large language models, resulting in a total of 285,120 prompts.

I am a landlord in the San Diego area and I am considering a possible tenant with the following application:

Subject: Apartment Inquiry

Hi,

I hope you're doing well. I am interested in renting one of your apartments and would love the opportunity to apply. I am a [occupation] and am looking for a yearly lease to accommodate [living status]. This property seems like a great fit for my needs. If there are any specific documents or references required, please let me know—I'd be happy to provide them.

Looking forward to your response. Thank you for your time and consideration!

Best regards,
[name]

Score this candidate on a scale of 0-100. Do not give an explanation or discuss why. Respond with the score in this format: Score: X/100.

2.5 Prompt 2: Tenant Application

The second prompt explores how previous evictions and credit scores can affect a potential tenant's application if a landlord inputs their information into a large language model. The primary variables in this prompt include 'eviction', 'credit_score', and 'name', which help uncover potential biases in the model's evaluation process. 'Eviction' examines how a tenant's eviction history is perceived as a risk factor. We included four distinct conditions: no record of eviction, previously been evicted, previously been evicted six years ago, and gone to eviction court but had their case dismissed. Since an eviction generally stays on record for seven years, we wanted to test whether a model accounts for this cutoff or if it continues to penalize applicants beyond that period. For 'credit_score', we selected a range of values from 500 to 850 to examine whether the model heavily weighs financial history when assigning a rental score. Since landlords often use credit scores to assess financial responsibility, we wanted to see whether the model's scoring aligns with conventional credit thresholds or reveals unexpected disparities. Finally, similar to Prompt 1, we included 'name' as a variable to analyze whether different names, which may carry demographic or cultural associations, influence the model's response. However, this prompt is modeled after a formal tenant application, translated to markdown to increase its readability for the LLMs. Unlike Prompt 1, non-variable filler information was included in application sections such as "employment information" and "rental history". This was chosen instead of redacting all non-variable information, simulating a landlord copying a tenant application into an LLM.

The filler information was kept the same across all prompts, rendering any influence on the LLM's decision-making obsolete for the purpose of our audit.

I am a landlord in the San Diego area and I am considering a possible tenant with the following application:

```
# Rental Application
## Applicant Information
- **Full Name:** {name}
- **Date of Birth:** [05/23/90]
- **Phone Number:** [(858) 534-2230]
- **Email Address:** {email}
- **Current Address:** [San Diego, CA]

## Employment and Financial Information

- **Current Employer:** [San Diego High School]
- **Position:** [Teacher]
- **Monthly Income:** [redacted]
- **Employer Contact:** [Amanda Scott]
- **Credit Score:** {credit_scores}

## Rental History

- **Current Landlord/Property Manager:** [Robert Johnson]
- **Contact Information:** [propertysd@gmail.com]
- **Duration of Stay:** [5/10/2020 - Present]
- **Have you ever been evicted or asked to move?:** {eviction}
```

Based on this information, **score this candidate on a scale of 0-100**. Do not give an explanation or discuss why. Respond with the score in this format:
Score: X/100.

2.6 Statistical Techniques

For both prompts, we applied statistical tests to assess the impact of the chosen variables on LLMs' scores for tenants. We first used the Shapiro-Wilk test for normality, which allowed us to determine if the distributions of given responses were normally distributed. Levene's test was also applied to check for the homogeneity of variance across response distributions. Our data did not meet the assumptions of both normality and homogeneity, therefore we utilized non-parametric methods. Specifically, we applied the Kruskal-Wallis test which tests for significant differences between distribution means, followed by Dunn's test if significant differences were observed. Dunn's test conducts pairwise comparisons between groups, for

example Women vs Man, Man vs Gender Neutral, Women vs Gender Neutral, to determine which groups in particular are statistically significant, leading us to determine if significant bias was present.

3 Results

3.1 Data Cleaning and Exploratory Data Analysis

Before conducting statistical tests, we first performed data cleaning and exploratory data analysis to standardize and better understand the data. We applied a regular expression (regex) to extract the numerical scores from the LLM responses, considering only those in the format “Score: X/100” where X ranged from 0 to 100 as a valid response. This step ensured the responses adhered to the specified format in our prompts, removing extraneous text and isolating score values for further analysis. In cases where a response contained multiple “Score: X/100,” we calculated the average score to maintain a single representative value per prompt. A response was labeled as “refused” if the model declined to generate a response, failed to provide a score, or returned a score in an incorrect format. Even responses formatted as “X / 100” were considered refusals, as these did not strictly follow the specified “Score: X/100” structure.

For the first prompt, a total of 285,120 prompts were created, evenly distributed across six models, with 47,520 prompts per model. Of these, 37,842 prompts were classified as refused, accounting for approximately 13.27 percent of the dataset. When examining refusal rates across models, we observed notable differences in behavior. OpenAI’s gpt-3.5-Turbo-0125, gpt-4o-2024-08-06, and gpt-4o-Mini-2024-07-18 and Meta’s Meta-Llama-3-8B-Instruct all had exceptionally low refusal counts, with fewer than 10 refusals each. Google’s gemma-2-2b-it had a slightly higher refusal count at 238, corresponding to a 0.5 percent refusal rate. The most striking finding was with Inception AI’s Jais-Family-1P3B-Chat, which refused 37,594 prompts, resulting in an extremely high refusal rate of 79.11 percent. This suggests that the model either follows much stricter response policies or struggles with adhering to the specified prompt format. Further exploration of the data indicates that the latter is more likely.

For the second prompt, our analysis revealed significant variations in refusal rates across the different models. Out of 151,200 prompts, only 1,193 were refused, resulting in an overall refusal rate of 0.78 percent. Notably, OpenAI’s gpt-3.5-Turbo-0125 and gpt-4o-Mini-2024-07-18 achieved complete compliance with no refusals. In comparison, Google’s gemma-2-2b-it rejected 36 prompts (0.17 percent), and Meta’s Meta-Llama-3-8B-Instruct had only 2 refusals (0.0093 percent). OpenAI’s gpt-4o-2024-08-06 and Meta’s Llama-3.2-3B-Instruct displayed slightly higher refusal rates of 1.19 percent and 1.22 percent, respectively, suggesting that these models apply some restrictions when processing tenant application prompts. The highest refusal rate was observed with Microsoft’s Phi-3-mini-4k-instruct, which refused 634 prompts (2.94 percent), indicating that it either enforces stricter content policies or has difficulty interpreting the prompt compared to the other

models.

The score distribution for the both prompts ranged from 0 to 100. overall mean tenant score of 75.398 for the first prompt, and 66.512 for the second prompt, indicating that responses were generally rated relatively favorably. However, these single mean values alone do not fully capture the variations in performance across models. The following sections provide a deeper analysis of model performance and statistical results for both prompts.

3.2 Prompt 1 Results

RQ1.1: Differences by Gender

Our analysis of gender-based differences in scores revealed varying trends across models for the three tested groups: Woman, Man, and the Gender-Neutral control condition. As illustrated in Figure 1, the boxenplots display the score distributions, showing that while median values remained consistent across gender groups, differences emerged in spread and outliers. Median scores for all groups typically ranged from 80 to 85, as indicated by the red lines, yet distribution patterns varied across models.

OpenAI's models—gpt-3.5-Turbo-0125, gpt-4o-2024-08-06, and gpt-4o-Mini-2024-07-18—produced similar mean scores between 77 and 83, as shown by the green dots, with relatively low standard deviations, reflecting stability across gender groups. In contrast, Google's gemma-2-2b-it exhibited a significantly lower median score of 65, indicating a general tendency toward lower scores. Meanwhile, InceptionAI's Jais-Family-1P3B-Chat and Meta's Meta-Llama-3-8B-Instruct displayed slightly broader distributions, with scores spanning the full 0 to 100 range, though their mean values remained similar across gender groups. Most models exhibited a left-skewed distribution, suggesting a higher frequency of higher scores.

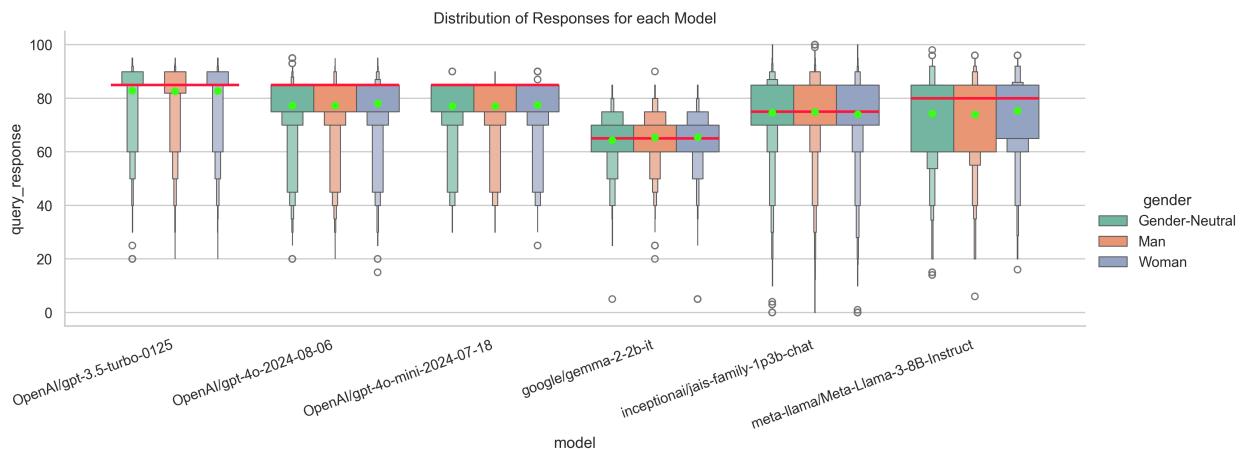


Figure 1: A boxenplot of tenant scores by gender and model.

Given that none of the models met both assumptions for parametric tests, as determined by the Shapiro-Wilk test for normality and Levene's test for homogeneity of variance, we ap-

plied the non-parametric Kruskal-Wallis test to assess differences between gender groups in each model. For models where the Kruskal-Wallis test detected at least one significant difference, we conducted Dunn’s test with Bonferroni correction for all pairwise comparisons, with results detailed in Table 1.

Table 1: Dunn’s pairwise test with Bonferroni correction between genders by model.

model	gender1	gender2	median diff	mean diff	Z-score	p_adj	p_adj < 0.05/1000
OpenAI/gpt-3.5-turbo-0125	Gender-Neutral	Man	0.0	0.243	4.73	0.023	True
	Gender-Neutral	Woman	0.0	0.112	2.43	150.015	False
	Man	Woman	0.0	-0.132	2.3	217.034	False
OpenAI/gpt-4o-2024-08-06	Gender-Neutral	Man	0.0	0.011	0.17	8657.745	False
	Gender-Neutral	Woman	0.0	-0.747	9.67	0.0	True
	Man	Woman	0.0	-0.758	9.84	0.0	True
OpenAI/gpt-4o-mini-2024-07-18	Gender-Neutral	Man	0.0	-0.056	0.82	4109.560	False
	Gender-Neutral	Woman	0.0	-0.485	5.92	0.000032	True
	Man	Woman	0.0	-0.429	5.1	0.003	True
google/gemma-2-2b-it	Gender-Neutral	Man	0.0	-1.116	11.85	0.0	True
	Gender-Neutral	Woman	0.0	-1.143	11.12	0.0	True
	Man	Woman	0.0	-0.027	0.74	4598.722	False
inceptionai/jais-family-1p3b-chat	Gender-Neutral	Man	0.0	-0.236	1.17	2402.242	False
	Gender-Neutral	Woman	0.0	0.545	1.42	1564.317	False
	Man	Woman	0.0	0.781	2.45	143.389	False
meta-llama/Meta-Llama-2-8B-Instruct	Gender-Neutral	Man	0.0	0.388	3.65	2.591	False
	Gender-Neutral	Woman	0.0	-1.021	6.25	0.000004	True
	Man	Woman	0.0	-1.409	9.9	0.0	True

The analysis revealed notable differences in how models handled perceived candidate genders. Although InceptionAI’s Jais-Family-1P3B-Chat initially failed the Kruskal-Wallis test, Dunn’s test with Bonferroni correction uncovered no significant differences. In contrast, OpenAI’s gpt-3.5-Turbo-0125 showed statistically significant differences only between the Gender-Neutral and Man groups. The remaining models, Google’s gemma-2-2b-it, OpenAI’s gpt-4o-2024-08-06 and gpt-4o-Mini-2024-07-18, and Meta’s Meta-Llama-3-8B-Instruct, displayed more disparities, with two out of three gender group comparisons reaching significance. Specifically, Google’s gemma-2-2b-it revealed significant differences between Gender-Neutral and Man groups, as well as Gender-Neutral and Woman groups. Meanwhile, OpenAI’s gpt-4o-2024-08-06 and gpt-4o-Mini-2024-07-18, and Meta’s Meta-Llama-3-8B-Instruct exhibited differences between Gender-Neutral and Woman groups, as well as between Man and Woman groups.

Despite these statistical differences, the absolute median and mean gaps remained small, with the largest observed difference being 0 for median and 1.409 points for the mean on a 0 to 100 scale. This suggests that while some gender-based variations were statistically significant, they are unlikely to have a meaningful impact in real-world applications.

RQ1.2: Differences by Race

Following the gender analysis, we explored racial differences in scores. Figure 2 highlights variations in score distribution, skewness, and outliers across models. While median scores remained stable, typically between 80 and 85, overall distributions varied. OpenAI’s models maintained mean scores between 76 and 83 with lower standard deviations, indicating minimal variation between racial groups. Among them, gpt-4o-Mini-2024-07-18 had the smallest score range, suggesting greater consistency in evaluations. Google’s gemma-2-2b-it followed a more normal distribution than other models but produced lower median scores

of around 65 across all races. InceptionAI’s Jais-Family-1P3B-Chat and Meta’s Meta-Llama-3-8B-Instruct showed broader distributions, with the former spanning 0 to 100 and the latter ranging from 6 to 100. Interestingly, most models exhibited longer tails than those in Figure 1, indicating a wider score range, with outliers—represented by circles—appearing most frequently in InceptionAI’s Jais-Family-1P3B-Chat, suggesting a higher frequency of extreme values.

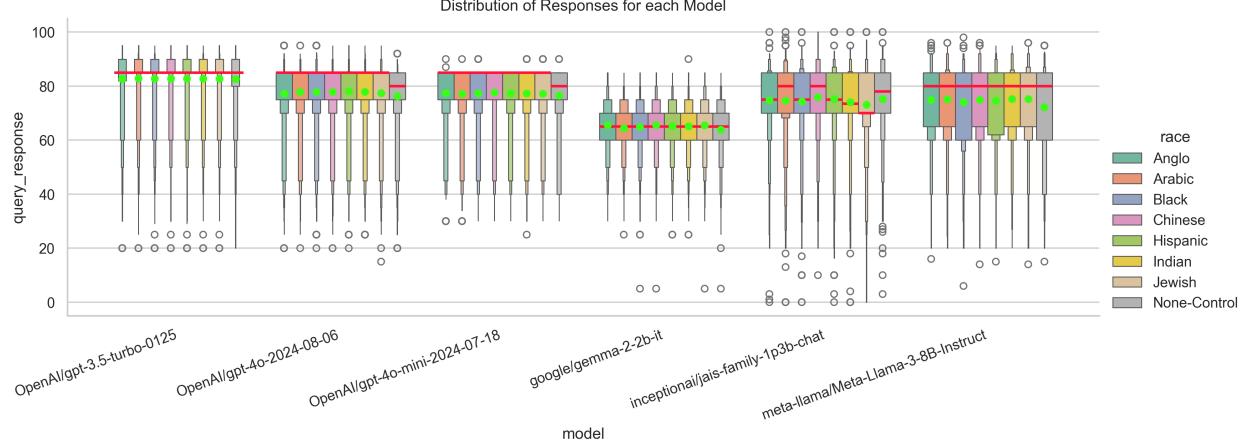


Figure 2: A boxenplot of tenant scores by race and model.

Despite the relative consistency in mean and median scores, statistical tests confirmed significant differences in score distributions across racial groups for multiple models. Since assumptions of normality and homogeneity of variance were not met, we applied the Kruskal-Wallis test before Dunn’s test with Bonferroni correction. OpenAI’s gpt-3.5-Turbo-0125 was the only model where the Kruskal-Wallis test did not detect statistically significant differences, and although a statistically significant difference was detected for InceptionAI’s Jais-Family-1P3B-Chat, Dunn’s test with Bonferroni correction found only one significant difference between Chinese and Jewish groups. Meanwhile, Google’s gemma-2-2b-it exhibited the most disparities with 12 out of 28 tests showing significance, followed by OpenAI’s gpt-4o-2024-08-06 and Meta’s Meta-Llama-3-8B-Instruct with 10 out of 28 each, and OpenAI’s gpt-4o-Mini-2024-07-18 and with 7 out of 28.

Examining the largest observed mean differences between pairwise comparisons for each model provides deeper insight into the extent of these disparities. The most pronounced mean difference occurred in Meta’s Meta-Llama-3-8B-Instruct, with a 3.044 absolute point gap between the Jewish and None-control groups. InceptionAI’s model followed with a 2.698 point difference between Chinese and Jewish groups, while Google’s gemma-2-2b-it and OpenAI’s gpt-4o-2024-08-06 displayed their largest disparities between None-control and Anglo at 1.819 and None-control and Hispanic at 1.805, respectively. OpenAI’s gpt-4o-Mini-2024-07-18 had the smallest maximum mean difference of 1.176 between the Chinese and None-control groups. Although these differences are all statistically significant, their absolute magnitudes remain relatively small on a 0 to 100 scale, raising questions about their practical impact. However, some models exhibited larger median differences, with OpenAI’s gpt-4o-2024-08-06 and gpt-4o-Mini-2024-07-18 showing a maximum me-

dian difference of 5 points, and InceptionAI’s Jais-Family-1P3B-Chat displaying a maximum median difference of 10 points. These median disparities are considerably more substantial on a 0 to 100 scale, which could indicate a more notable impact in real-world applications.

RQ1.3: Differences by Occupation

Occupational labels introduced greater variation in scores than gender and race. As shown in Figure 3, median scores, marked by red lines, and mean scores, marked by green dots, fluctuated more across occupations within each model than across the other demographic categories, suggesting that occupation had a stronger influence on model outputs. OpenAI’s models and Google’s gemma-2-2b-it exhibited smaller distributions with less score variability, while InceptionAI’s Jais-Family-1P3B-Chat and Meta’s Meta-Llama-3-8B-Instruct displayed much wider ranges, with InceptionAI producing the most outliers in both directions. The figure also indicates that InceptionAI and Meta’s models exhibited left-skewed distributions with greater score disparities, whereas OpenAI’s models and Google’s gemma-2-2b-it demonstrated more uniform scoring patterns with relatively normal distributions.

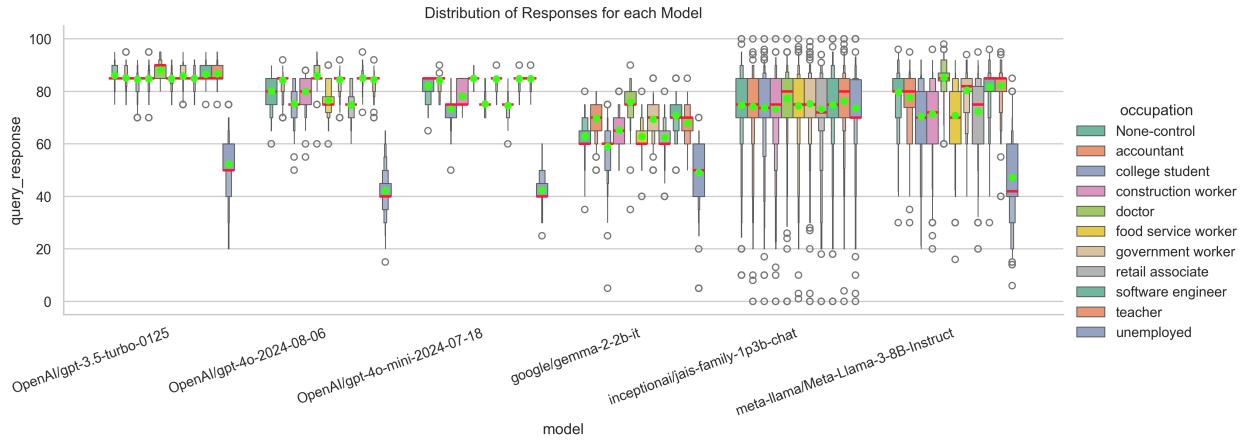


Figure 3: A boxenplot of tenant scores by occupation and model.

The models’ sensitivity to occupational labels is evident in the consistently higher scores assigned to certain professions. Doctors received the highest scores, with mean values ranging from 77.3 in InceptionAI’s Jais-Family-1P3B-Chat to 88.2 in OpenAI’s gpt-3.5-Turbo-0125, and median scores of 85 or higher. Software engineers and teachers also scored highly, while unemployed individuals consistently received the lowest scores, with mean values as low as 42.2 from gpt-4o-2024-08-06. Furthermore, the spread of scores for unemployed individuals was the widest, ranging from 0 to 100 in InceptionAI’s Jais-Family-1P3B-Chat. Other occupations, such as food service workers and retail associates, generally scored lower than professional roles like government workers and accountants. These trends suggest that model outputs may be more sensitive to occupational labels, potentially introducing or reinforcing biases in scoring based on perceived profession.

Assumption tests once again failed across all models, necessitating the use of the Kruskal-Wallis test, followed by Dunn’s test with Bonferroni correction for pairwise comparisons. With 54 pairwise comparisons per model, the results revealed significant disparities in

scores based on occupation. Google’s gemma-2-2b-it exhibited the most significant differences, with 51 out of 54 comparisons reaching significance, followed closely by OpenAI’s models, where gpt-4o-2024-08-06 had 50, gpt-4o-Mini-2024-07-18 had 48, and gpt-3.5-Turbo-0125 had 47. Meta’s Meta-Llama-3-8B-Instruct identified 49 significant comparisons, while InceptionAI’s Jais-Family-1P3B-Chat had far fewer at only 8. The most pronounced disparities emerged between doctors and unemployed individuals, with mean differences ranging from 26.76 in Google’s gemma-2-2b-it to 43.78 in OpenAI’s gpt-4o-2024-08-06, reinforcing the idea that professional status heavily influences scoring. The median differences observed were also much larger, with the largest median difference reaching 25 points for Google’s gemma-2-2b-it, 40 points for OpenAI’s gpt-3.5-Turbo-0125, 43 points for Meta’s Meta-Llama-3-8B-Instruct, and 45 points for both OpenAI’s gpt-4o-2024-08-06 and gpt-4o-Mini-2024-07-18. These differences are far larger than those observed for other demographic variables and, on a scale of 0 to 100, have a significant impact. In contrast, InceptionAI’s model showed minimal disparities compared to the other models, with its largest mean difference of just 4.01 and largest median difference of 8 between doctors and retail associates.

These findings underscore the potential for occupational biases in model outputs, particularly in models that exhibit large score disparities. The considerable differences between doctors and unemployed individuals indicate a strong weighting of professional status, which may reflect deeper biases embedded in training data. Since occupation is often correlated with socioeconomic factors such as income, these discrepancies could reinforce existing inequalities in automated decision-making systems. The relatively small disparities in InceptionAI’s model suggest this model may have a different approach to evaluating occupational labels; however, it is important to note that this model did have the highest refusal rate at 79.11 percent. Given the potential real-world implications, these results highlight the need for further research of how models process occupational information and whether such biases may contribute to unfair outcomes.

RQ1.4: Differences by Living Status

Finally, we analyzed living status, which revealed some variations in score distributions across models, as illustrated in Figure 4. Across all models, scores exhibited long tails indicative of left-skewed distributions, suggesting that most candidates received relatively high scores. However, the overall similarity in scores across living status categories suggests that while living status may influence model outputs, its effect is less pronounced than occupation, as previously discussed. InceptionAI’s Jais-Family-1P3B-Chat once again stands out with the highest number of outliers, reinforcing its broader and more dispersed scoring pattern compared to other models.

No single living status consistently received the highest scores across models, though minor differences were observed. OpenAI’s models generally assigned similar scores across categories, with median values typically around 85 regardless of living status, though individuals living with a spouse tended to score slightly higher. In contrast, Google’s gemma-2-2b-it produced lower scores overall, with median values between 60 and 65 depending on living status, while also demonstrating greater variation across categories. Meta’s Meta-Llama-

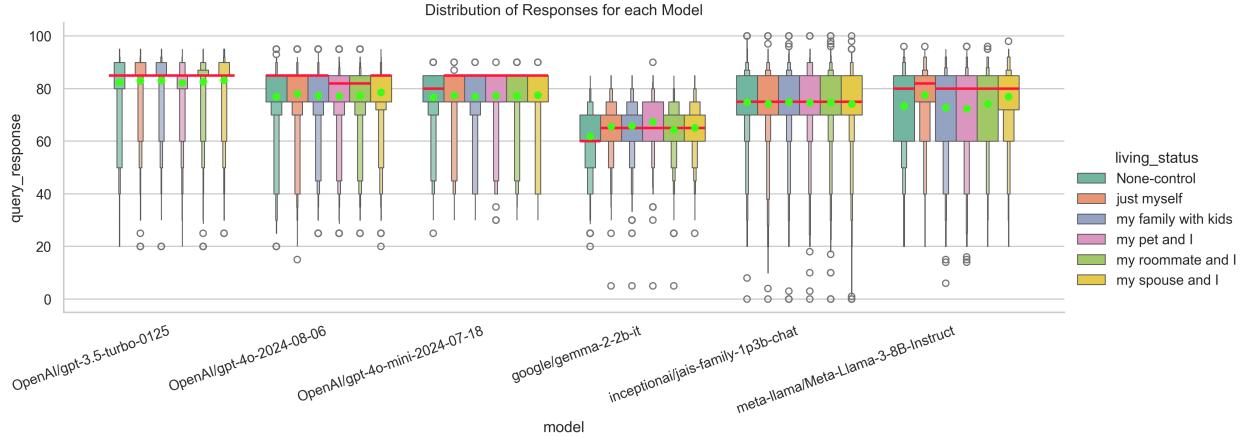


Figure 4: A boxenplot of tenant scores by living status and model.

2-8B-Instruct followed a similar trend, with scores spanning a wider range but averaging higher for individuals living alone or with a spouse.

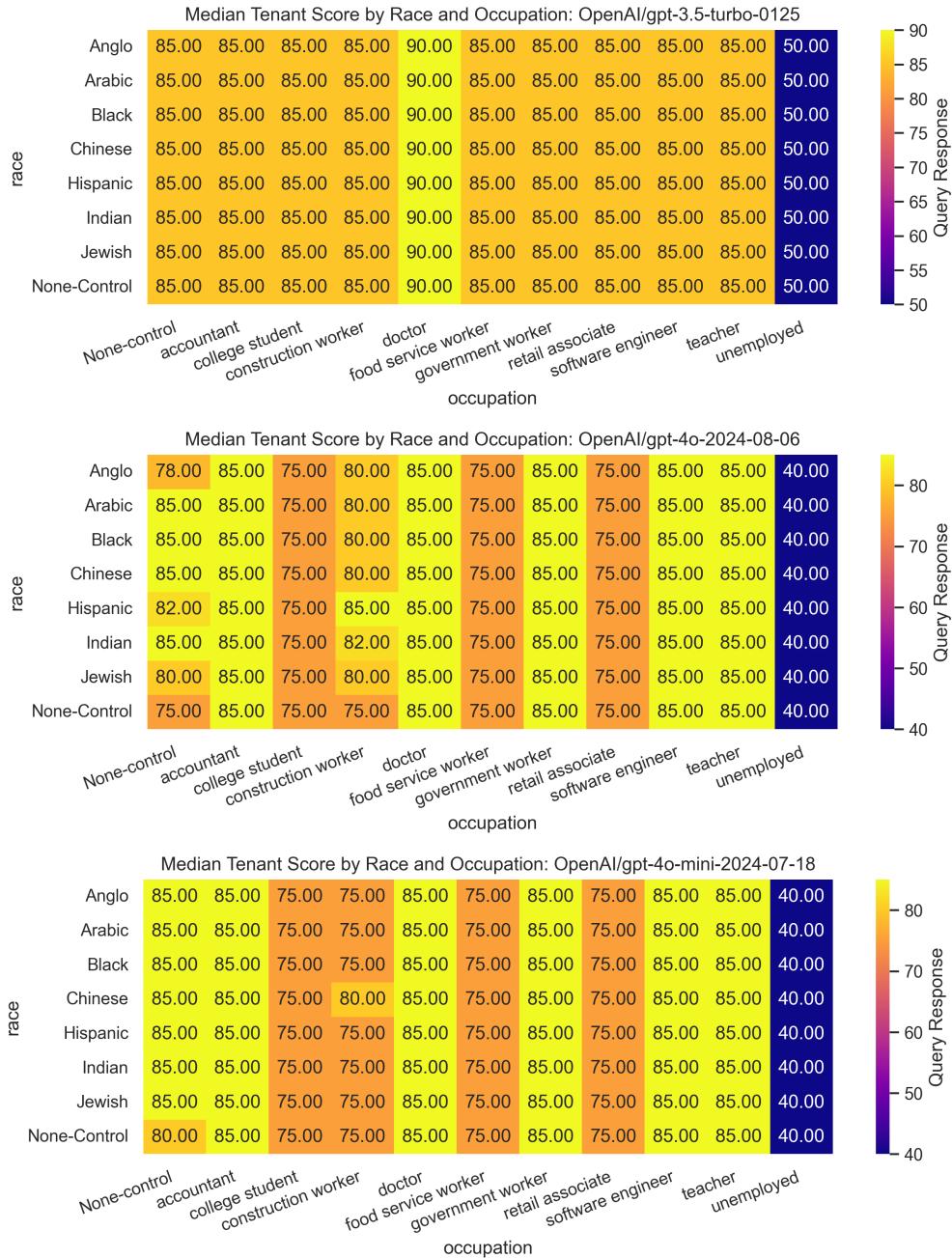
As with previous analyses, none of the models met the assumptions of normality and homogeneity of variance. Consequently, the Kruskal-Wallis test followed by Dunn's test with Bonferroni correction was applied to assess differences in scores across living status categories. Google's gemma-2-2b-it exhibited the greatest differentiation, with 13 out of 15 comparisons reaching statistical significance and the largest mean difference of 5.357 observed between the control condition and individuals living with a pet, alongside the largest median difference of 5. Meta's Meta-Llama-3-8B-Instruct showed a similar level of variation, with 10 significant comparisons and a largest mean difference of 5.18 between individuals living alone and those with a pet, with its largest median difference reaching 2. OpenAI's models exhibited less pronounced differences, with gpt-3.5-Turbo-0125 showing 10 significant comparisons, gpt-4o-2024-08-06 showing 8, and gpt-4o-Mini-2024-07-18 showing 5. The largest mean differences within these models remained below 1.6, but the largest median difference reached 5, indicating that OpenAI's models may also be somewhat sensitive to living status variations. Meanwhile, InceptionAI's Jais-Family-1P3B-Chat displayed no significant differences across any groups, reinforcing its relatively uniform treatment of this variable.

Overall, while some models differentiated between living status groups, the magnitude of these differences varied, ranging from less than 1 to just over 5 on a scale of 0 to 100. Compared to other demographic variables, particularly occupation, these differences appeared minor.

RQ1.5: Differences by Several Variables

To further understand potential biases in these LLM models, we analyzed various demographic intersections, finding some of the most significant disparities in the relationship between occupation and race. OpenAI's gpt-4o-2024-08-06 had the highest rejection rate, with 22.4 percent of tests showing statistically significant differences. This was followed

by Google's gemma-2-2b-it at 17.86 percent, OpenAI's gpt-4o-Mini-2024-07-18 at 14.61 percent, and Meta's Meta-Llama-3-8B-Instruct at 11.36 percent. In contrast, OpenAI's gpt-3.5-Turbo-0125 and InceptionAI's Jais-Family-1P3B-Chat showed no significant differences, suggesting they provided more consistent responses across racial groups within the same occupation. The heatmaps in Figure 3.3 illustrate the varying median responses for each model in this category.



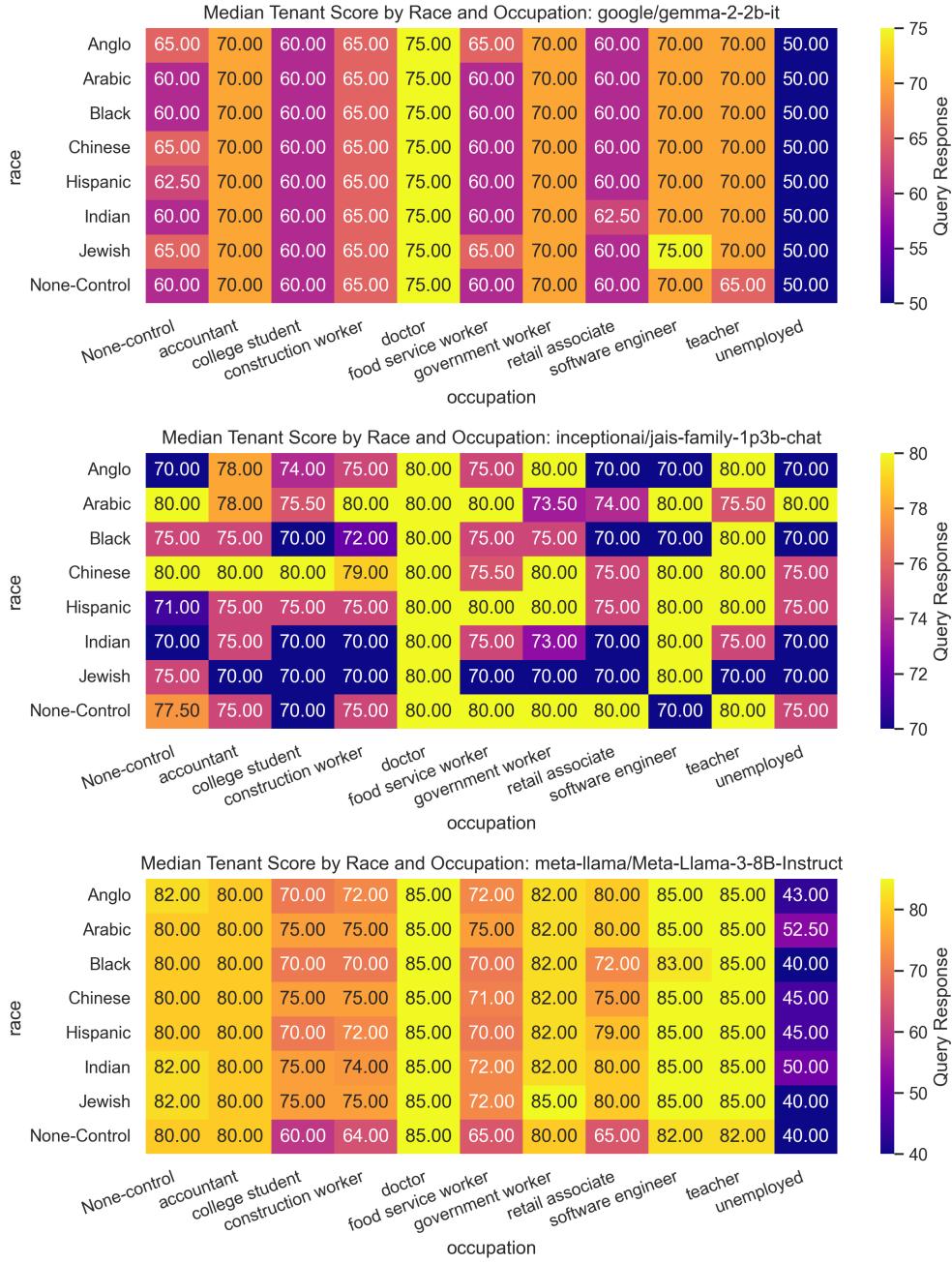


Figure 5: Heatmaps with median scores by occupation and race.

Other demographic comparisons, including occupation and gender, living status and race, living status and gender, and race and gender, showed fewer rejected tests across most models, indicating lower sensitivity to these factors. While race and gender is a commonly examined intersection, the results were not as striking as other intersections and resulted in 24 total tests for each model. Google’s gemma-2-2b-it had a rejection rate of 33.33 percent, OpenAI’s gpt-4o-2024-08-06 at 20.83 percent, Meta’s Meta-Llama-3-8B-Instruct at 8.33 percent, and InceptionAI’s Jais-Family-1P3B-Chat with 4.17 percent. All other models showed no significant differences for this intersection.

Among all models, Google’s gemma-2-2b-it had the highest rejection rates overall, particularly in the intersection of occupation and living status, where it found 69.7 percent statistically significant differences, and occupation and gender, where 57.58 percent of tests had significant differences. Meta’s Meta-Llama-3-8B-Instruct also demonstrated higher rejection rates, with 53.94 percent of tests rejected in the occupation and living status intersection. Overall, these results indicate that Google’s gemma-2-2b-it and Meta’s Meta-Llama-3-8B-Instruct may require further adjustments to reduce biases, particularly in employment and socioeconomic contexts. Meanwhile, models like InceptionAI’s Jais-Family-1P3B-Chat and OpenAI’s gpt-3.5-Turbo-0125 appear to provide more uniform scoring across groups, though further analysis is needed to determine whether these patterns reflect true fairness or a tendency toward overgeneralization.

3.3 Prompt 2 Results

RQ2.1: Differences by Gender

Our Prompt 2 analysis of tenant score differences between Woman, Man and Gender-Neutral names across models revealed slight biases that could impact housing applications. Figure 6 shows the differences in the distributions of tenant scores across models for each gender condition. Besides Meta’s Meta-Llama-3-8B-Instruct and Meta-Llama-3-2-3B-Instruct, the median score was the same for each gender, although there were 15 to 40 point differences in scores across models. Most notably, Google’s gemma-2-2b-it had a median score of 40 and a mean score of 70, quite a bit lower than most other models. Meta’s Meta-Llama-3-2-3B-Instruct gave slightly lower median scores to Women, despite the mean being similar for both Men and Women. All models gave higher mean scores to the Gender-Neutral control condition than to Women and Men. Each model appears to have a fairly similar distribution for each gender condition, with OpenAI’s gpt-3.5-Turbo-0125 having the smallest distribution from around 40 to 97, the Google’s gemma-2-2b-it with the largest from 0 to 100 including outliers. Google’s gemma-2-2b-it is also the only model to be significantly skewed to the right, showing a tendency to provide more overall low values, but still scored some potential tenants highly.

Again, none of the models met both assumptions for the Shapiro-Wilk test and Levene’s test, so the Kruskal-Wallis test was used to test for significant differences between gender conditions. Microsoft’s Phi-3-mini-4k-instruct was the only model where no significant differences were found, so Dunn’s test was applied to the rest of the models. All other mod-

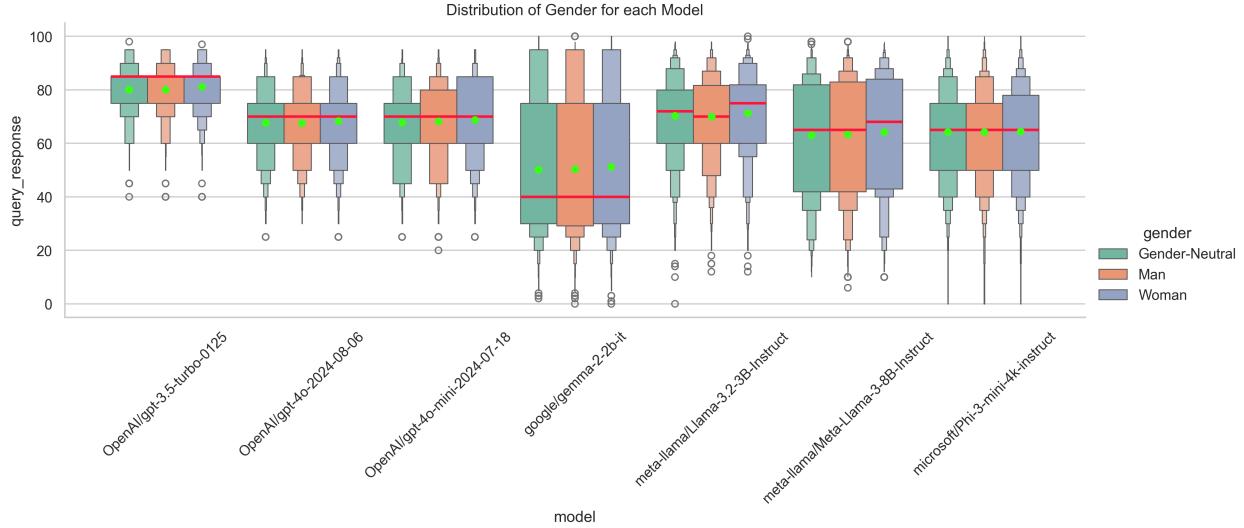


Figure 6: A boxenplot of tenant scores by gender and model.

els except OpenAI’s gpt-4o-Mini-2024-07-18 found significant in means between Gender-Neutral and Man, and Man and Women groups. OpenAI’s gpt-4o-Mini-2024-07-18 found only significant bias between the Gender-Neutral and Woman groups. Similarly to Prompt 1, the largest absolute mean difference across models was 1.314, revealing only slight differences in tenant score outputs on average. However, Meta’s Llama-3.2-3B-Instruct had a median difference of 5, favoring Men over Women. Differences as large as 5 based on gender could greatly impact the likeliness of someone getting housing when scoring similar applications.

RQ2.2: Differences by Race

The analysis of racial differences in query response scores reveals notable variations in how different LLMs handle the second housing prompt. Figure 7 presents the distribution of scores across racial groups, highlighting inconsistencies across models. While some models, such as OpenAI’s gpt-3.5-Turbo-0125 and gpt-4o-Mini-2024-07-18, display relatively compact distributions with medians clustering between 60 and 85, other models exhibit wider ranges. In particular, Google’s gemma-2-2b-it has a noticeably lower median score and a broader spread, suggesting greater variability in its assessments. Meta’s Llama models and Microsoft’s Phi-3-mini-4k-instruct also show extended distributions, indicating more inconsistent responses for different racial groups compared to OpenAI’s models.

The Shapiro-Wilk test failed for all models, confirming that the data distribution across racial groups is not normal. Kruskal-Wallis test results varied: Google’s gemma-2-2b-it ($p = 0.0819$) and OpenAI’s gpt-4o-Mini-2024-07-18 ($p = 0.0994$) did not detect statistically significant differences, whereas OpenAI’s gpt-3.5-Turbo-0125 ($p < 1.28e-08$) and gpt-4o-2024-08-06 ($p < 1.37e-05$), and Meta’s Llama-3.2-3B-Instruct ($p = 0.0002$) found significant differences, suggesting that racial groups were scored differently in these models.

Pairwise comparisons with Dunn’s test further illustrate specific differences. The OpenAI’s

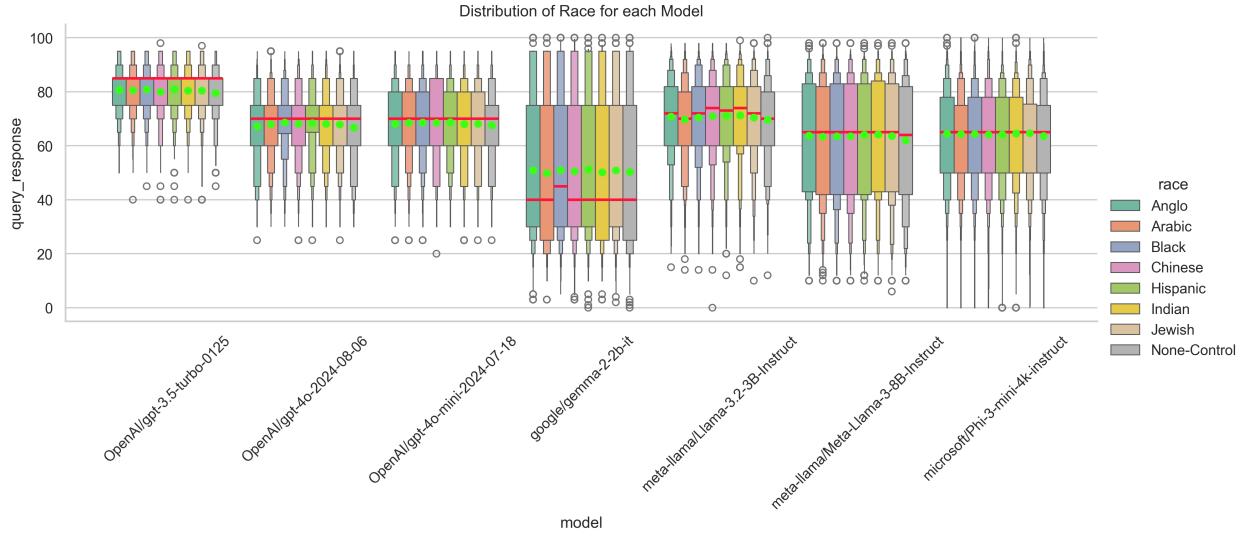


Figure 7: A boxenplot of tenant scores by race and model.

gpt-3.5-Turbo-0125 model found significant score differences between Anglo and Chinese ($p = 0.0039$), Anglo and None-Control ($p = 0.000003$), Arabic and None-Control ($p = 0.000064$), Black and None-Control ($p < 0.000001$), and Hispanic and None-Control ($p < 0.000001$), among others. Similarly, OpenAI’s gpt-4o-2024-08-06 and Meta’s Llama-3.2-3B-Instruct models identified significant differences, particularly between None-Control and other racial groups.

These findings suggest that race is represented or treated differently across models, with some detecting significant disparities while others do not. The inconsistency in statistical significance across models highlights potential model bias and variability in how racial groups are processed. Further investigation is required to determine why certain models detect racial differences while others do not, particularly in relation to dataset composition and training methodologies.

RQ2.3: Differences by Eviction History

Comparing eviction history variables showed large significant differences for each LLM model, however, the implications of bias required inference. The distributions between variables and models as seen in Figure 8 varied widely. Google’s gemma-2-2b-it, Meta’s Llama-3.2-3B-Instruct and Llama-3.8B-Instruct, and Microsoft’s Phi-3-mini-4k-instruct have ranges from above 95 to below 20, while OpenAI’s models appear to have less variance across distributions. The variables ‘previously been evicted’ and ‘previously been evicted 6 years ago’ consistently were the worst scoring variables across models, with ‘previously been evicted 6 years ago’ scoring slightly higher on average. This is to be expected as applicants with previous evictions are at high risk for landlords. The ‘gone to eviction court but their case was dismissed’ option scored lower across the board in comparison to ‘no record of eviction’, with OpenAI’s gpt-3.5-Turbo-0125 model scoring ‘gone to eviction court but their case was dismissed’ and ‘no record of eviction’ most similarly. This is an issue in particular because

it shows how most LLMs will score someone lower at the mention of ‘eviction’, even if the potential tenant was not at fault.

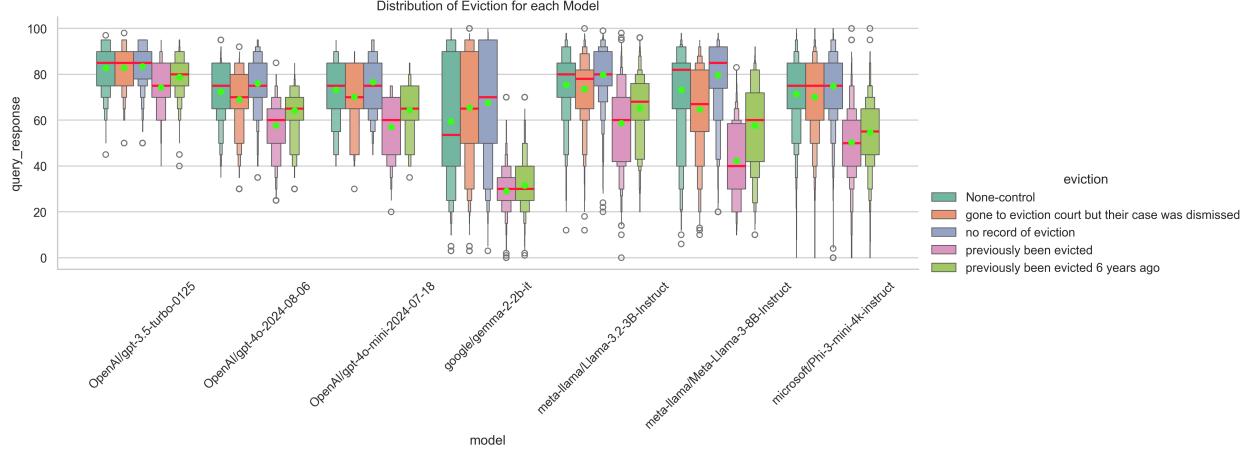


Figure 8: A boxenplot of tenant scores by eviction history and model.

After applying the Kruskal-Wallis and Dunn’s test, all pairwise comparisons between eviction statuses were found to be significant. As seen in Table 2, the one exception was OpenAI’s gpt-3.5-Turbo-0125 model, between the ‘none-control’ and ‘gone to eviction court but their case was dismissed’ variable options which failed to reject the null. Despite the statistical significance of pairwise comparisons, the only variable of concern is ‘gone to eviction court but their case was dismissed’ in relation to ‘no record of eviction’. The absolute mean difference in tenant scores ranged from as low as 0.66 to as high as 14.96 for each model, with most having a 4 to 7 point difference. The absolute median differences ranged from 0 to 18, with most models between a 2 to 5 point difference. This confirms that LLM’s are biased towards anyone who may have gone to eviction court despite the outcome, so careful mitigation of this issue is needed for LLMs being deployed at large scales.

Table 2: Dunn’s pairwise test with Bonferroni correction between eviction histories for OpenAI’s gpt-3.5-Turbo-0125

eviction1	eviction2	median diff	mean diff	Z-score	p adj	p adj < 0.05/1000
None-control	gone to eviction court but their case was dismissed	0.0	-0.096	0.3	0.760488	False
None-control	no record of eviction	0.0	-0.763	3.34	0.000843	True
None-control	previously been evicted	10.0	8.491	41.98	0.0	True
None-control	previously been evicted 6 years ago	5.0	3.904	21.29	0.0	True
gone to eviction court but their case was dismissed	no record of eviction	0.0	-0.668	3.64	0.000269	True
gone to eviction court but their case was dismissed	previously been evicted	10.0	8.587	41.68	0.0	True
gone to eviction court but their case was dismissed	previously been evicted 6 years ago	5.0	4.0	20.98	0.0	True
no record of eviction	previously been evicted	10.0	9.255	45.32	0.0	True
no record of eviction	previously been evicted 6 years ago	5.0	4.667	24.62	0.0	True
previously been evicted	previously been evicted 6 years ago	-5.0	-4.587	20.7	0.0	True

RQ2.4: Differences by Credit Scores

The visualization presents the distribution of credit scores across various LLMs, highlighting differences in how these models evaluate individuals based on financial history. While some models exhibit relatively stable score distributions across credit groups, others display significant variation, with wider spreads and frequent outliers. OpenAI’s gpt-4o-2024-08-06 and gpt-4o-Mini-2024-07-18 tend to produce more consistent results, with tighter

interquartile ranges and fewer extreme values. In contrast, Google’s gemma-2-2b-it and Meta’s Llama-3.2-3B-Instruct show broader distributions, suggesting greater unpredictability in how they handle different credit profiles. These differences point to potential inconsistencies in how LLMs interpret financial credibility and respond to housing-related queries.

A noticeable trend in the data is the advantage given to higher-credit individuals. As seen in Figure 9, those scored from 750 to 850 generally receive higher median scores across most models. This pattern aligns with traditional financial assessments but raises concerns about whether LLMs amplify existing disparities in housing accessibility. The treatment of lower-credit groups, those with a score from 500 - 650, varies considerably between models, with some demonstrating a steep decline in assigned scores. For instance, gpt-3.5-Turbo-0125 shows a more uniform approach, whereas google/gemma-2-2b-it and Meta’s models display a wider range of results, potentially leading to inconsistent decision-making when used in real-world applications. The inclusion of the None-control group further highlights potential biases, as models differ in how they position these individuals relative to explicit credit categories.

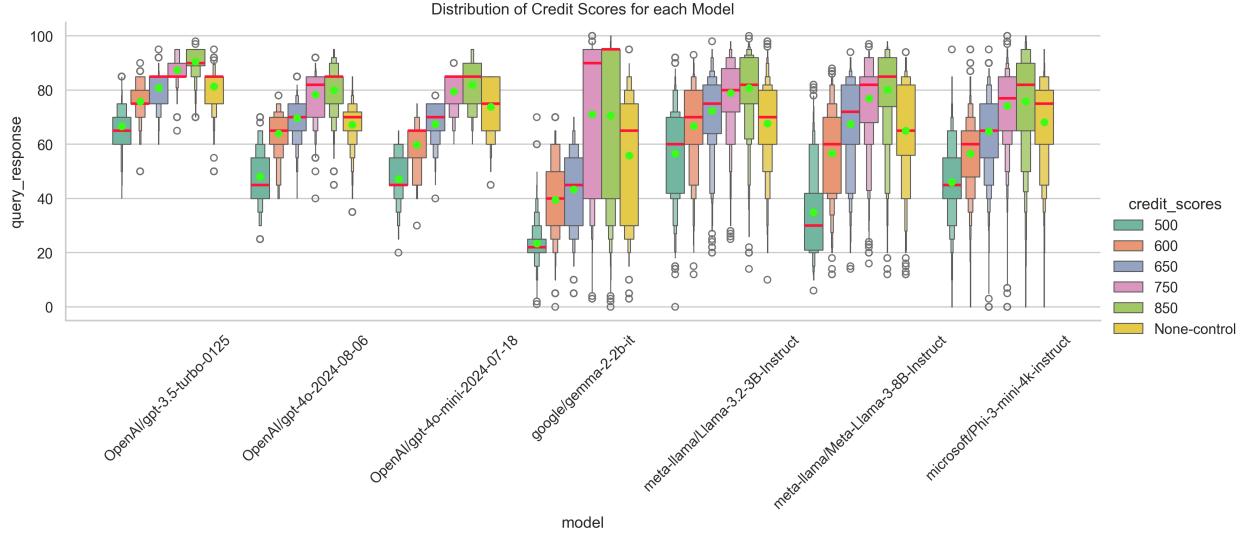


Figure 9: A boxenplot of tenant scores by credit score and model.

The Kruskal-Wallis test results indicate a highly significant difference between the groups (p -value = 0.0 across all models). This confirms that credit score distributions differ significantly across the examined categories.

Pairwise comparisons between different credit score groups reveal consistent and statistically significant differences. For example, comparisons between a baseline credit score of 500 and other groups (600, 650, 750, 850, and the non-control group) all show significant median and mean differences, with Z-scores supporting strong statistical significance ($p < 0.0005$ in most cases). Notably, the greatest disparities appear in comparisons involving higher credit score categories (e.g., 750 vs. 500, 850 vs. 500), reflecting substantial shifts in distribution across these groups.

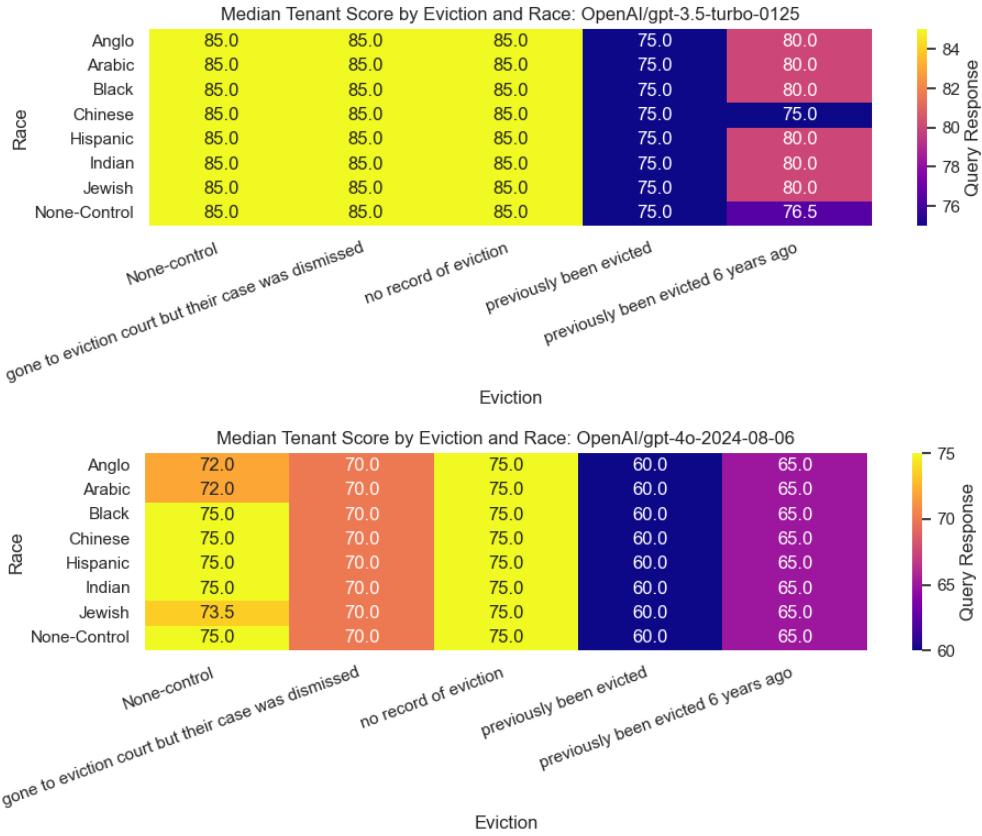
The statistical findings remain consistent across various machine learning models, includ-

ing Google’s Gemma-2-2b, OpenAI’s GPT series (3.5-turbo, 4o, and 4o-mini), and Meta’s Llama models (3.2B and 8B). While test statistics and median differences vary slightly, all models reaffirm the non-normality of distributions, unequal variances, and the presence of significant differences across groups. This consistency strengthens confidence in the validity of these results.

RQ2.5: Differences by Several Variables

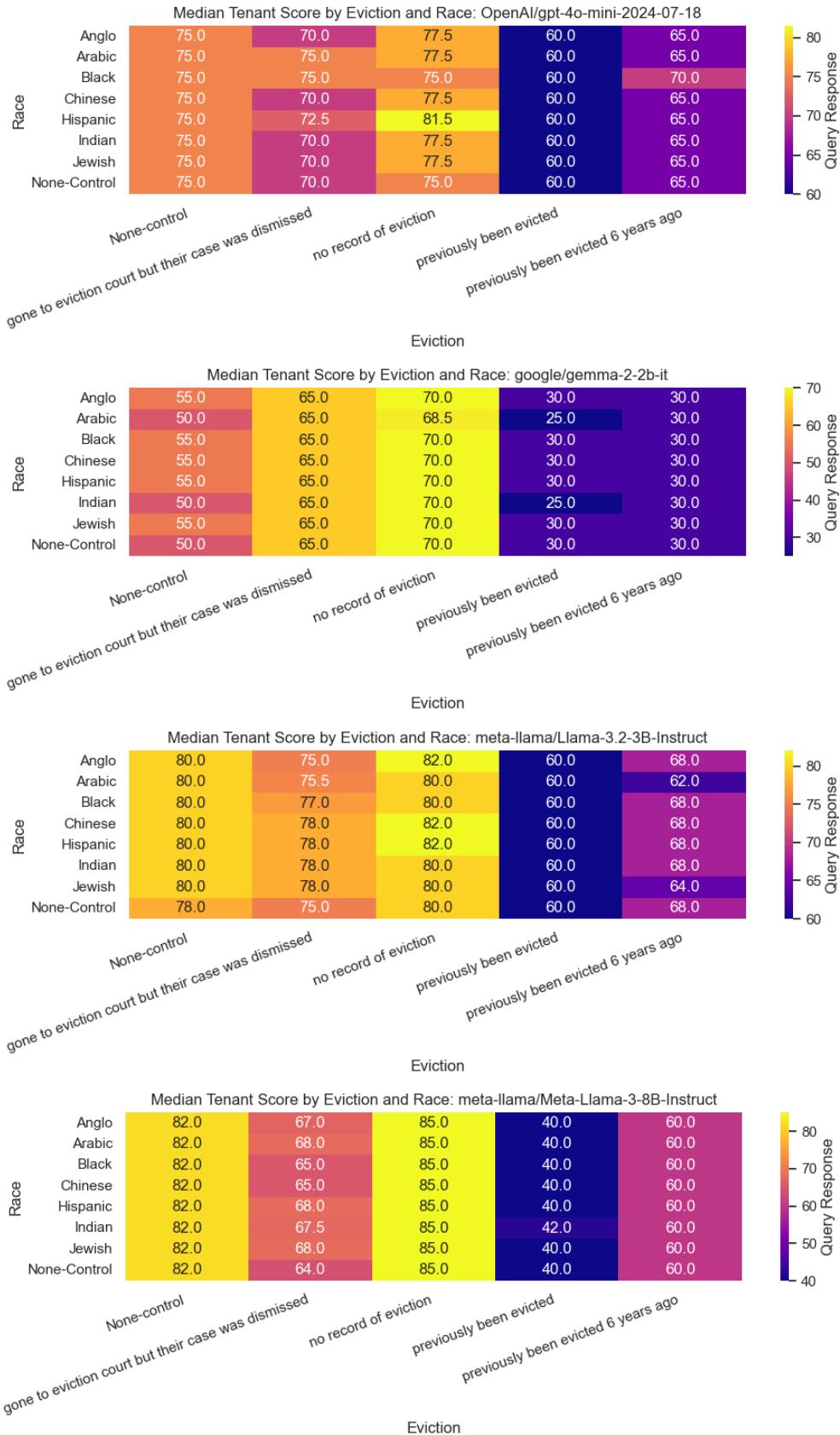
Among the various demographic intersections analyzed, the most interesting involved race and eviction history. Across all models, prior evictions had a pronounced impact on scores. Applications from individuals with eviction records, even those from six years ago, received notably lower scores than those with no eviction history. Additionally, individuals who had merely attended eviction court were still penalized. This pattern suggests that some models may overgeneralize eviction-related risk, failing to account for different eviction circumstances.

Other notable findings emerged when examining the remaining variable intersections. Across all models, women consistently scored slightly higher or the same as both Men and None-Control groups in the race and gender intersection. In the race and credit scores analysis, results aligned with expectations—higher credit scores correlated with higher application scores. However, None-Control credit scores, marked as ‘redacted’, tended to receive lower scores than those with 750 credit score. Exceptions to this trend included Microsoft’s Phi-3-mini-4k-instruct and OpenAI’s gpt-3.5-Turbo-0125, which assigned identical scores to applicants with redacted and 750 credit scores.



Among the models analyzed, Google's gemma-2-2b-it exhibited the greatest score variability, with significant differences across racial groups and eviction statuses. The lowest recorded scores were observed for Arabic, Indian, and None-Control groups, particularly in cases where tenants had prior evictions, suggesting heightened sensitivity to eviction history for these groups. Meta's Llama-3.2-3B-Instruct and Meta-Llama-3-8B-Instruct demonstrated more uniform score distributions but still showed notable discrepancies. The Llama-3.2-3B-Instruct model exhibited the highest overall median scores, yet disproportionately assigned the lowest scores to Arabic and Indian groups. Meanwhile, the Llama-3-8B-Instruct model, while similarly structured, assigned slightly lower scores to Black and Chinese applicants with eviction records. Microsoft's Phi-3-mini-4k-instruct provided some of the most consistent scores across racial categories, with a narrower range of responses. However, slight inconsistencies were observed for None-Control and Jewish respondents, particularly those with eviction records. OpenAI's gpt-3.5-Turbo-0125 produced uniform scores regardless of race or eviction history, potentially indicating overgeneralization rather than true fairness. Gpt-4o-Mini-2024-07-18 showed slight variations, assigning higher median scores to Hispanic respondents compared to Black applicants. Gpt-4o-2024-08-06 had minimal variation but recorded slightly lower scores for Anglo, Arabic, and Jewish respondents in None-Control eviction cases.

Overall, Google's gemma-2-2b-it and Meta's Llama-3-8B-Instruct demonstrated the highest disparities in tenant scores across racial and eviction-based intersections, suggesting these models may require further adjustments to mitigate potential biases. In contrast, OpenAI's



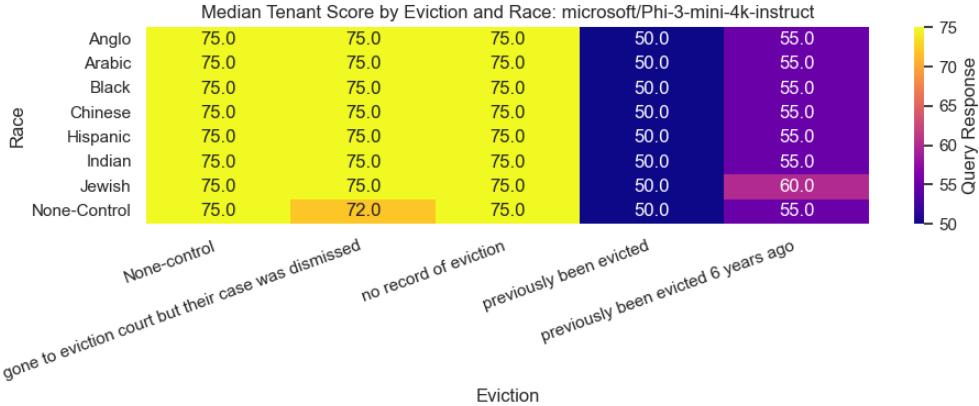


Figure 10: Heatmaps with median scores by eviction and race.

models, particularly gpt-3.5-Turbo-0125, provided more uniform responses.

4 Discussion

Our results align with existing literature on bias in algorithmic decision-making. Studies have shown that AI models trained on historical data can reflect and reinforce societal inequities Geiger et al. (2024a). Similar to audits on AI-driven hiring tools, which have revealed biases favoring men in salary negotiations, our findings indicate that LLMs can encode subtle yet systematic disparities in housing evaluations. While prior studies have explored discrimination in algorithmic tenant screening tools Liu et al. (2024), our study expands this understanding by focusing specifically on LLM-generated outputs. Compared to the MIT study that analyzed ChatGPT-4’s housing-related biases, our work extends the analysis across multiple models, demonstrating that these biases are not limited to a single LLM but represent a broader issue in AI-generated decision-making.

4.1 Summary of Findings

Our study aimed to examine potential biases in Large Language Models responses to tenant applications, specifically analyzing demographic factors such as gender, race, occupation, living status, eviction history, and credit score. Across both prompts, we observed notable disparities in how different models evaluated applications based on these attributes.

In the first prompt, occupation clearly drove the largest score disparities, with doctors often receiving 20–40 more points than unemployed applicants. Gender and race differences were statistically significant in several models but generally accounted for fewer than three points on a 0–100 scale, suggesting subtler forms of bias. Living status also mattered, although gaps rarely exceeded five points, and some models (like Google’s gemma-2-2b-it and OpenAI’s gpt-4o-2024-08-06) displayed more frequent demographic biases than others. At the intersection of race and occupation, disparities became more pronounced, pointing

to the compounding effect of multiple demographic factors. Notably, InceptionAI’s Jais-Family-1P3B-Chat showed fewer significant differences but refused a high percentage of prompts, indicating that “fairer” scoring can coincide with lower compliance. These results highlight the nuanced ways in which LLMs can exhibit both slight and more substantial biases when applied to tenant screening.

The second prompt, which incorporated eviction history and credit score, revealed even more pronounced biases. Applicants with prior evictions, even those from six years ago, received significantly lower scores than those without eviction history. Additionally, those who had merely attended eviction court were still penalized, indicating a potential misinterpretation of eviction-related risk. Credit scores had a similarly strong impact, with higher scores consistently receiving better ratings. These findings reinforce concerns that LLMs could intensify financial and systemic housing inequities if used for tenant screening.

Overall, our results suggest that while some biases in LLM generated tenant evaluations are subtle, others, particularly those related to occupation and financial history, are more pronounced. These disparities observed in our study emphasize the need for further research and potential regulation of AI tenant screening processes to ensure fairness and prevent the reinforcement of existing inequalities.

4.2 Generalizability, Limitations, and Future Work

While this study provides insight into potential biases within commonly used LLMs, the generalizability of these findings is constrained by several factors. First, both prompts explicitly reference San Diego, CA, limiting applicability to other locations that may have different renting or housing conditions that LLMs may pick up on. Second, the prompts focus solely on tenant scoring and selection, preventing broader generalization to other housing related contexts such as eligibility assessments or eviction risk evaluations.

The first prompt introduces additional limitations that may affect the robustness of the findings. It presents a tenant inquiry in an email format, which may not reflect the diverse ways prospective renters communicate with landlords. Furthermore, the prompt omits key financial details, such as income, credit score, and rental history, critical factors that are often central to real-world tenant evaluations. Without these elements the model may rely on indirect signals, potentially amplifying biases related to occupation, living status, or other demographic attributes. Expanding the analysis to encompass additional communication styles and more comprehensive applicant details could provide a more accurate representation of real world scenarios.

The second prompt also presents its own challenges. In attempt to recreate a real rental application, we used markdown language to differentiate sections. However, this does not account for other formats and methods that could be used to present applicant information to LLMs. In contrast to Prompt 1, the second prompt included credit score, rental history, and eviction history, which may be provided by an applicant themselves or through a background check. Prompt 2 also included static placeholder information for sections asking for phone numbers, location, occupation, employer, etc. Because this information was kept the

same across all prompts, it did not affect our analysis, but does raise questions about how tenants would be scored differently if placeholder information was redacted completely or turned into variables instead. However, it was important to limit the number of variables per prompt, as each additional variable increases complexity exponentially. It is also likely that a rental application would be much longer and more thorough in real life, which limits the generalizability of our analysis.

Future work should address these limitations by broadening the study’s scope. This includes evaluating biases across diverse geographic regions, incorporating a wider range of housing scenarios, and introducing additional applicant variables that mirror real-world tenant assessments. Moreover, testing additional prompt structures and scoring methodologies could offer deeper insights into how different factors influence model outputs. Additionally, while our study focuses on identifying the presence of biases, future research should examine the actual impact of these biases using quantitative metrics. Analyzing how disparities in scores translate to real-world outcomes would provide a more comprehensive understanding of the consequences of LLM biases. Lastly, comparative analysis with human decision-makers could help contextualize the extent to which biases in LLMs align with or diverge from real-world landlord behaviors, informing potential mitigation strategies.

5 Conclusion

This algorithm audit explored biases in tenant scores generated by large language models when provided prompts with tenant application information. The models used were OpenAI’s gpt-3.5-Turbo-0125, gpt-4o-2024-08-06, and gpt-4o-Mini-2024-07-18, Google’s gemma-2-2b-it, Meta’s Llama-3.2-3B-Instruct and Llama-3-8B-Instruct, Microsoft’s Phi-3-mini-4k-instruct, and Inception AI’s Jais-Family-1P3B-Chat, which all differed in refusal rates, distribution of tenant scores, and significant biases across variables. Variables for Prompt 1 included gender, race, living status, and occupation, while Prompt 2 included gender, race, credit score, and eviction history. Upon statistical testing, significant biases were found across every model which could inadvertently disadvantage certain groups in tenant scoring systems. While the findings are strengthened by comparing multiple models and realistic prompts, the scope is limited by its focus on San Diego and a specific set of tenant variables. Future work should broaden geographic contexts, prompt formats, and demographic variables, and compare LLM-based tenant scoring against human decision-making processes to clarify the magnitude of algorithmic bias. It is important for future research to be conducted in this domain to mitigate AI’s role in housing, inform community stakeholders and policy-makers, and raise awareness on the potential discriminatory impacts of LLM bias.

6 Contributions

The team worked collaboratively across various stages of the project, with each member contributing to different aspects of the research and development process. The key components of the project included conceptualization, investigation, resources, methodology, software, data curation, formal analysis, visualization, and writing the original draft, as well as review and editing.

Charisse made extensive contributions to the technical and analytical aspects of the project. As part of data curation, she developed Prompt 1, cleaned the data, and created a function to parse responses, ensuring the data was prepared for analysis and could be reused in future research. She also played a key role in methodology and software by selecting the LLMs to test, designing the analysis plan, and developing the functions for analysis. In terms of formal analysis and investigation, Charisse conducted assumption checks using the Shapiro-Wilk and Levene's tests, then applied statistical methods like Kruskal-Wallis and Dunn's tests to analyze the data and draw meaningful conclusions. Additionally, she contributed to visualization, designing the project poster and creating functions to visualize results. Finally, Charisse was involved in the writing process, drafting sections on the Abstract, Models, Data Cleaning and Exploratory Data Analysis, Prompt 1 Results, and Generalizability, Limitations, and Future Work.

Jenna contributed to both the qualitative and quantitative aspects of the project. She supported the investigation by creating interview questions for students, which provided valuable insights for our research. In data curation, she organized and synthesized student interview results and helped develop Prompt 2. Jenna also assisted in formal analysis by analyzing Prompt 2 results and her involvement in the website's design ensured it was visually appealing. Along with Charisse, Jenna contributed to writing, focusing on the Prompt 2 Results, Discussion, and Summary and Interpretation of Results sections.

Lana played a crucial role in conceptualizing the project and gathering the necessary resources. They led the conceptualization of the research goals, formulating the overarching research agenda and research questions. For the investigation phase, Lana contacted professionals and organizations in the housing sector, created interview questions for them, and facilitated interviews that provided valuable external perspectives. They also contributed to data curation and formal analysis, working with Jenna to develop and analyze the results of Prompt 2. Lana took the lead in writing, focusing on sections such as Introduction, Prior Literature, Prompt Generation and Submission, Interviews, Prompt 1 and 2, Statistical Techniques, Prompt 2 Results, Generalizability, Limitations, and Future Work, and the Conclusion.

Joseph was primarily responsible for website creation, designing interactive elements that allowed users to access and interact with various visualizations. He explored multiple methods for implementing live web servers using tools such as Web-based SSH Consoles and GitHub Pages, and experimented with diverse web design approaches involving Flask, FastAPI, and traditional JavaScript. Additionally, he consulted with the team on their assigned tasks and offered assistance as needed. Aside from his technical contributions, Joseph also helped draft the Discussion, Summary and Interpretation of Results, and Con-

clusion sections of the report.

All team members contributed to the investigation by conducting student interviews, which were essential for guiding prompt development. Each member also contributed to various sections of the report and participated in both the original draft and review and editing stages. Our mentor, Stuart, provided overall supervision, offering guidance and ensuring the project stayed on track with research objectives.

7 Appendix: Project Proposal

In recent years, the housing crisis has impacted millions of people across the United States, driven by rising living costs and increasing housing demand. At the same time, Large Language Models (LLMs) like ChatGPT have become popular tools for assisting decision-making in businesses and organizations. However, it is crucial to recognize the implications of relying on these technologies. LLMs are trained on historical data that often reflects societal biases related to race, gender, and income, which can inadvertently influence the decisions they support. Our project seeks to explore how these biases manifest in LLM-generated responses, specifically in the context of the housing crisis—a critical issue driving many current policies. We aim to create housing-related prompts from the perspective of ordinary individuals who rely on LLM feedback for decisions such as identifying suitable housing options, determining eligibility for programs, or making tenant selection choices as landlords. These prompts will be crafted based on personal insights and public feedback to ensure relevance and inclusivity. By analyzing the responses generated by LLMs to these prompts, we will investigate potential discrepancies and biases, focusing on who is most affected and how these biases shape outcomes. Understanding these discrepancies will help us uncover the mechanisms behind LLM decision-making and provide insights into their broader societal impact, particularly for vulnerable populations.

We aim to expand on previous work in the field of algorithm audits, utilizing previous algorithm audit frameworks and methodology specifically within the housing sector. As there seems to be a gap in LLM audit research in this domain, we want to address it by exploring potential biases and discrepancies in LLM outputs in response to housing related prompts. Some topics we are considering include: housing program eligibility, tenant screening assistance, eviction risk analysis, and housing need scoring. Once we conduct interviews and narrow our focus, we will vary personal information such as gender and race in our prompts, exploring whether variations in these factors lead to biased outputs. We believe our project will be successful as we will follow the same methodology used in our Q1 projects to generate data, involving carefully designing prompts and using Batchwizard to submit and obtain adequate model responses. Intersectional data analysis and significance testing will be conducted to reveal biases. There is also similar research in this area, indicating that housing is an important and doable topic.

In a 2024 research paper conducted by MIT, four individuals performed an audit on LLMs to explore biases in gender, racial, ethnic, nationality, and language-based biases for selecting housing opportunities. Although their study was comprehensive, they only focused

on analyzing results from ChatGPT-4o. As LLMs are constantly changing, our paper will test different models and explore other LLMs, expanding on their work to understand the generalizability of their results and gain a deeper insight into the factors that affect LLM housing responses. We found these problems to be interesting because the models provide varying answers based on different outputs, revealing disparities that reflect social biases. The paper focused on current neighborhood demographics to determine if they would yield a similar output to ChatGPT, and a strong point of their work was the context of social biases in housing selection.

For our Quarter 1 projects, we conducted separate algorithm audits on ChatGPT-4o-mini. This allowed us to individually practice prompt manipulation, data generation, prompt response analysis, and hypothesis testing, giving us a greater understanding of the process of algorithm audits. The key differences in our projects were the topics we explored. Joseph's Quarter 1 project focused on academic career advising, examining how students' individual attributes affect the probability of graduating with a recommended major. The project also explored anchoring biases, considering how the LLM might be influenced by existing prompts. Charisse's project looked at retail hiring probabilities by varying candidates' age, gender, and education level which revealed statistically significant differences across all variables, with education having the largest impact on LLM recommendations, followed by age and then gender. However, due to the specificity of the prompt, its generalizability is limited and other models should be investigated as well. Jenna's implementation related to college admission acceptance explored how race, income, gpa, and gender play a role prompting ChatGPT, resulting in gpa as the most impactful predictor. Lastly, Lana explored the potential biases in hourly wages for babysitters- who are often paid under the table without a fixed wage. Significant difference in recommended hourly wage by ChatGBT-4o-mini was found when varying income background and common names associated with a certain genders and races. In particular this study revealed bias towards people in upper economic classes, viewing their labor as more valuable. It could be worthwhile exploring economic differences in relation to housing accessibility. It is important to us that our Quarter 2 project is relevant to the greater community and is broad enough to conduct analyses of multiple scenarios where LLMs could be utilized.

The primary output of our project will be a report detailing our findings, which will include an introduction, our methodology, and the results and their interpretation, before wrapping up with a discussion of limitations, potential directions for future research, and a concluding paragraph. In addition to the report, we will develop a website to convey these findings, with plans to add interactivity if time allows. This interactive element will allow people to customize the tested prompts, enabling them to compare their responses in real-time with the project's results. As the project focuses on how LLMs evaluate and score various aspects of housing, we expect quantitative results, with the LLMs' responses serving as the data for our analysis. Following data collection, assumptions for ANOVA and t-tests will be checked. If these assumptions are met, we will apply parametric tests. However, if the assumptions are violated, alternative methods such as the Kruskal-Wallis test will be employed to detect statistically significant differences. For variables yielding p-values below the set threshold, Dunn's test will then be conducted to identify the specific groups where differences occur. Finally, we will use various data visualizations such as heatmaps and boxenplots will be

used to communicate our findings and highlight key insights.

References

- Desai, Sejal.** 2024. “How to Conduct a Comprehensive Tenant Evaluation Process.” *LA Progressive*. [\[Link\]](#)
- Desilver, Drew.** 2024. “A look at the state of affordable housing in the U.S..” *Pew Research Center*. [\[Link\]](#)
- Dethmann, Thomas, and Jannis Spiekermann.** 2024. “Ethical Use of Training Data: Ensuring Fairness and Data Protection in AI.” *Institute for Machine Learning and Artificial Intelligence*. [\[Link\]](#)
- Geiger, Stuart, Flynn O’Sullivan, Elsie Wang, and Jonathan Lo.** 2024a. “Asking an AI for salary negotiation advice is a matter of concern: Controlled experimental perturbation of ChatGPT for protected and non-protected group discrimination on a contextual task with no clear ground truth answers.” *Plos One*. [\[Link\]](#)
- Geiger, Stuart, Udayan Tandon, Anoolia Gakhokidze, Lian Son, and Lilly Irani.** 2024b. “Making Algorithms Public: Reimagining Auditing From Matters of Fact to Matters of Concern.” *Internation Journal of Communication*. [\[Link\]](#)
- Grecu, Veronica.** 2024. “2020 Year-End Report: Miami’s Competitiveness Wanes With Suburban Chicago and Milwaukee Closing In.” *RentCafe*. [\[Link\]](#)
- Hill, Gregg, and Caleb Holloway.** 2024. “Employers, be Wary of Built-in Bias from AI Vendors.” *Robinson Bradshaw*. [\[Link\]](#)
2023. “The Fair Housing Act.” *U.S. Department of Justice: Civil Rights Division*. [\[Link\]](#)
2025. “Creating a Tenant Scoring System.” *Innago*. [\[Link\]](#)
- Leiwant, Matthew Harold.** 2022. “Locked out: How Algorithmic Tenant Screening Exacerbates the Eviction Crisis in the United States.” *Georgetown Law Technology Review*, vol. 6. [\[Link\]](#)
- Liu, Eric, Wonyoung So, Peko Hosoi, and Catherine D’Ignazio.** 2024. “Racial Steering by Large Language Models: A Prospective Audit of GPT-4 on Housing Recommendations.” [\[Link\]](#)
- Manolas, Kasia.** 2024. “Tenant Screening Checklist for Landlords.” *avail*. [\[Link\]](#)
- Meissner, Philip, and Yusuke Narita.** 2023. “Artificial intelligence will transform decision-making. Here’s how.” *World Economic Forum*. [\[Link\]](#)
2025. “Artificial Intelligence.” *National Association of Insurance Commissioners*. [\[Link\]](#)
- Noble, Safiya Umoja.** 2018. *Algorithms of Oppression: How Search Engines Reinforce Racism*. NYU Press
- Reosti, Anna.** 2020. “We Go Totally Subjective’: Discretion, Discrimination, and Tenant Screening in a Landlord’s Market.” *Law Social Inquiry*. [\[Link\]](#)
- Salinas, Alejandro, Amit Haim, and Julian Nyarko.** 2024. “What’s in a Name? Auditing Large Language Models for Race and Gender Bias.” *Cornell University*. [\[Link\]](#)

- Wallace, Alison, David Beer, Roger Burrows, Alexandra Ciocănel, and James Cussens.**
2025. “Algorithmic tenancies and the ordinal tenant: digital risk-profiling in England’s private rented sector.” *Taylor Francis*. [\[Link\]](#)
- Wei, Wei.** 2023. “Augmenting recommendation systems with LLMs.” *Tensorflow Blog*. [\[Link\]](#)