

Evaluating LLM Behavior in Hiring: Implicit Weights, Fairness Across Groups, and Alignment with Human Preferences

Morgane Hoffmann¹, Emma Jouffroy¹, Warren Jouanneau¹, Marc Palyart¹ and Charles Pebereau¹

¹Malt, 75009 Paris, France

Abstract

General-purpose Large Language Models (LLMs) show significant potential in recruitment applications, where decisions require reasoning over unstructured text, balancing multiple criteria, and inferring fit and competence from indirect productivity signals. Yet, it is still uncertain how LLMs assign importance to each attribute and whether such assignments are in line with economic principles, recruiter preferences or broader societal norms. We propose a framework to evaluate an LLM's decision logic in recruitment, by drawing on established economic methodologies for analyzing human hiring behavior. We build synthetic datasets from real freelancer profiles and project descriptions from a major European online freelance marketplace and apply a full factorial design to estimate how a LLM weighs different match-relevant criteria when evaluating freelancer-project fit. We identify which attributes the LLM prioritizes and analyze how these weights vary across project contexts and demographic subgroups. Finally, we explain how a comparable experimental setup could be implemented with human recruiters to assess alignment between model and human decisions. Our findings reveal that the LLM weighs core productivity signals, such as skills and experience, but interprets certain features beyond their explicit matching value. While showing minimal average discrimination against minority groups, intersectional effects reveal that productivity signals carry different weights between demographic groups.

Keywords

Large Language Models, Person-job Fit, Fairness, Interpretability

1. Introduction

Generative Large Language Models (LLMs) are increasingly used in hiring for diverse tasks such as candidate evaluation, job matching, applicant ranking, and skill assessment, both for ad-hoc decisions by individual recruiters and at scale within recruitment pipelines. Their ability to process unstructured text, such as résumés and job descriptions, makes them particularly suited for recruitment tasks, where subtle signals of suitability, such as employer prestige, inferred skills, or contextual fit matter.

Despite their promise, LLMs raise concerns about interpretability, fairness, and alignment with human decision-making. Recruitment decisions involve complex trade-offs between multiple factors such as wages, competence, and work arrangement, and it remains unclear whether LLMs weigh these factors in ways that reflect recruiter preferences or societal values. This early-stage position paper addresses that gap by proposing a framework to analyze an LLM's implicit decision logic in recruitment contexts. This contribution is intended as a first step in a broader research effort to analyze the implicit reasoning of LLMs and their alignment with humans in hiring decisions.

The paper addresses the following research questions:

RQ1: Which aspects of candidate profiles and job context does the LLM emphasize, downplay, or overlook in its hiring recommendations?

RQ2: Does the LLM's decision logic systematically vary across socio-demographic groups in ways that could result in unequal hiring outcomes?

RQ3: How can our experimental framework enable systematic comparison between LLM and human recruiter decision logic in future studies?

We investigate these questions using data from a large European freelance marketplace where recommender systems help clients¹ identify relevant freelancers. Our methodol-

ogy adapts a standard experimental approach from labor economics used to study human recruiters' hiring decisions by placing the LLM in the role of a recruiter. Using a full factorial design, we generate a synthetic dataset by independently varying attributes (e.g. skills, daily rates, work arrangements) from real freelancer profiles, project descriptions, and recruiter characteristics and then prompt a LLM to evaluate freelancer-project matches. This setup allows us to estimate the causal impact of each attribute on the model's evaluation behavior. Our methodology is general and applies to any context in which recruiters evaluate candidates using unstructured text such as résumés and job descriptions. Using this framework, we analyze the results both in aggregate and across subgroups, based on variations in freelancer, recruiter, and brief characteristics to better understand the heterogeneity in the LLM's decision logic.

Our main findings are as follows. The LLM's scoring behavior broadly aligns with standard economic theory: it rewards close matches between freelancer profiles and project briefs, as well as signals of trust and competence such as platform reputation and relevant industry experience. The strongest penalties are applied to freelancers with insufficient experience or no prior activity on the platform. In contrast, the LLM places minimal weight on socio-demographic characteristics (gender, ethnicity, education) and former employer characteristics like firm size or industry. However, heterogeneity analysis reveals that LLM scoring varies with perceived gender, ethnicity, and educational background, suggesting that it uses demographic information to infer underlying competence or suitability for the position.

The remainder of this paper is organized as follows. Section 2 reviews related work. Section 3 describes our experimental methodology, including the platform context, design principles, and LLM scoring procedure. Section 4 presents our main findings on LLM decision-making patterns, subgroup analyses, and discusses how the methodology can be extended for systematic comparison with human recruiters. Section 5 concludes with implications and directions for

ably to refer to individuals or companies seeking to hire freelancers on the platform. Similarly, "freelancers" and "candidates" refer to those offering their services.

RecSys in HR'25: The 5th Workshop on Recommender Systems for Human Resources, in conjunction with the 19th ACM Conference on Recommender Systems, September 22–26, 2025, Prague, Czech Republic.



© 2025 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

¹Throughout this paper, we use "clients" and "recruiters" interchange-

future research.

2. Related Work

A rich economics literature conceptualizes hiring as the interpretation of signals such as education, experience, or reputation to infer candidate productivity under uncertainty [1]. In online labor markets, additional information such as ratings or platform reputation also comes into play [2, 3]. Meanwhile, persistent discrimination is well documented in the economic literature, with several field experiments showing that even when qualifications and experiences are held constant, demographic attributes such as gender or ethnic origin can affect both labor market outcomes and recruiters’ attention [4, 5]. While these studies have shown how human recruiters use signals and where biases arise, less is known about how automated systems and especially LLMs interpret these same attributes in the recruitment context. Our work explores whether algorithmic tools rely on similar signals. Our work explores whether a LLM rely on similar signals by testing a diverse set of features identified by this literature, including productivity indicators (education, experience and skills), demographic markers (gender-associated names, ethnic origin), and platform-specific signals (ratings or reputation badges).

As LLMs are increasingly deployed in sensitive domains like recruitment, concerns about transparency and fairness have led to the development of auditing protocols, mainly in the machine learning literature. Empirical studies reveal that LLMs can detect demographic signals even in the absence of explicit group labels [6]. This literature has also proposed interpretability tools such as LIME [7] and SHAP [8], and counterfactual approaches [9], each with strengths and limitations. For attribution methods, causal interpretation often requires additional assumptions [10, 11, 12]. Our approach contributes to this methodological literature by employing a fully factorial experimental design to enable causal measurement of how LLMs process candidate characteristics in recruitment settings. We construct sythetic profiles closely following [13] while incorporating specific freelancer attribute. Our approach produces robust, interpretable and causal estimations of the implicit weights given to candidate characteristics. By systematically manipulating candidate attributes we eliminate potential confounding variables and enable direct causal inference regarding the model’s decision-making process. This methodology draws heavily from correspondence studies, an established method also used in economics for eliciting human recruiter preferences.

Comparative studies have begun to examine the similarities and differences between human recruiters and LLMs in candidate screening. Some works report that AI-based tools can outperform humans in identifying relevant candidates, while others find misalignments in the weighting of credentials or the treatment of demographic variables [14, 15]. Yet, these comparisons often rely on observational data or different settings, limiting causal inference and comparisons. To address these limitations, our experimental framework can be readily extended to human recruiters, by presenting them with the same sythetic candidate profiles and project descriptions used for LLM evaluation. Such methodology lays the groundwork for rigorous alignment assessments and deeper understanding of both human and AI-driven recruitment logic.

3. Methodology

This section presents our methodology for analyzing how LLMs evaluate candidate-project matches. We begin by introducing the platform context and occupational focus of our analysis, then outline the design principles underlying our full factorial approach. We describe the generation of sythetic project briefs and freelancer profiles and conclude with details of the LLM-based scoring procedure.

3.1. Platform Context and Occupational Focus

Our study is situated in the context of a large European online freelance marketplace that connects clients with freelancers across diverse occupations and geographies. The platform relies on recommender systems to help clients identify suitable candidates, using freelancer profiles and platform histories, client characteristics and project descriptions, to generate ranked freelancer recommendations.

To limit dataset size while maintaining relevance, we focus our analysis on full-stack developers in France. This occupation-country combination forms the largest market segment on this platform, with high supply and demand, and exhibits strong gender skew where approximately 90% of freelancers have male-identified first names. This makes it a particularly suitable setting for exploring potential gender bias in model evaluations. To assess generalizability beyond the tech sector, we also analyze search engine optimization (SEO) content writing, which has similar platform prominence but reversed gender distribution. Results are presented in Appendix D. While our analysis focuses on this specific platform and occupations, our methodology generalizes to any recruitment setting where evaluators assess candidates based on unstructured text inputs such as résumés and job descriptions. The following sections details our experimental design.

3.2. Design Principles and Objectives

To analyze how LLMs weigh different freelancer attributes, we construct a sythetic dataset using a full factorial design. For each profile and project characteristic, we define a set of realistic attribute values based on data from the online freelancing platform. We then systematically generate all possible combinations of these attributes to create complete freelancer profiles and project briefs, yielding a balanced dataset where attributes vary independently across all combinations. Further construction details are provided in the following subsections.

This controlled design breaks natural correlations found in real-world data, allowing us to isolate the marginal effect of each attribute on the model’s output and enable causal identification. Compared to simpler randomization approaches, full factorial design ensures complete coverage of the attribute space and balanced representation of each combination, maximizing statistical power to estimate both main and interaction effects. While this approach limits external validity since our dataset is not representative of real-world distributions, this trade-off is intentional. Our objective is not representativeness but rather the creation of optimal conditions for isolating and analyzing the LLM’s underlying decision logic. Observational analyses of platform data, though externally valid, suffer from feature cor-

relations and confounding variables that preclude causal identification.

Our approach draws inspiration from correspondence studies, a well-established methodology in labor economics for eliciting recruiter preferences [13] and detecting discrimination in hiring [16]. These studies submit fictitious résumés to real job postings, systematically varying characteristics like names or qualifications to reveal hiring preferences. However, they are typically constrained by the costs of human evaluation. LLM evaluation significantly reduces these costs, which makes it feasible to employ a full factorial designs which maximizes statistical power.

3.3. Synthetic Recruiter Briefs

We construct synthetic project briefs that replicate the information recruiters provide when seeking freelancers on the platform and that recommender systems use to evaluate profile fit. On the platform, project briefs typically include a description, required skills and experience, expected duration, work arrangements, and compensation. To replicate this structure, we first define a set of plausible values for key demand-side attributes, based on real platform data, then apply generate all possible combinations, ensuring each attribute varies independently.

To limit the size of our dataset, while ensuring comparability across briefs, we vary some dimensions while holding others constant. The complete specification is provided in Table 1b in the Appendix.

The following dimensions are varied:

- **Recruiter identity:** Recruiter first names signal gender using common European male or European female names.
- **Firm size:** Projects originate either from small businesses (SMEs) or large corporations (as signaled in the brief text).
- **Work conditions:** Briefs specify preferences for remote vs. on-site work and full-time vs. part-time engagement.

The following dimensions are held constant: 6-month project duration, €400 daily rate, JavaScript/TypeScript with Node.js as the required technology stack, and a minimum of 5 years of experience. This factorial design results in 16 unique briefs, covering all possible combinations of the variable dimensions.

3.4. Synthetic Freelancer Profiles

We construct synthetic freelancer profiles to reflect the information that recommendation models use to evaluate candidate suitability for recruiter briefs. Each profile consists of typical sections found on the platform, such as skills, professional experience, work preferences, platform reputation, etc. To build these profiles, we first define a list of realistic values for each attribute based on real freelancer data. We then apply a full factorial design to generate all possible combinations, ensuring that each attribute varies independently. The complete list of profile attributes and values is provided in Table 1a in the Appendix.

Each profile varies along the following dimensions:

- **Skills:** Technical stacks are set to exactly match, be a close substitute to, or differ substantially from

the briefs requirement; JavaScript/TypeScript with Node.js, NestJS, or Angular, respectively).

- **Experience level:** Years of experience are fixed at 1, 5, or 9 years, representing junior, mid-level, and senior profiles.
- **Work arrangements:** Preferences for full-time vs. part-time and remote vs. on-site work are varied.
- **Reputation**²: Reputation is captured through combinations of the number of completed projects (0, 1, or 5), average recruiter rating (0 or 5 stars), and badge presence (yes or no)³.
- **Daily rates** : Five daily rate levels ranging from 300€ to 500€.
- **Work history:** Prior experience varies by employer size (SME vs. large company) and industry (e-commerce vs. banking/insurance).
- **Socio-demographic attributes:** First names signal perceived gender and ethnicity, with three categories: European male, European female, and Arabic male⁴. Education levels are set to bachelor's or master's degrees, the most common among freelancers in this category.

This procedure yields 10,800 unique profiles that cover all possible attribute combinations.

3.5. LLM Scoring Procedure

Procedure Each generated profile is systematically matched with each generated project brief, creating 172,800 unique profile-brief combinations. For every pair, we prompt *Gemini 2.0 Flash*⁵ to evaluate candidate-project fit by assigning a hiring probability on a 10-point scale (e.g., 2.5 meaning 25% probability). The model receives both the freelancer profile and project brief as input (see Appendices A.1 and A.2), with all interactions conducted in French.

Using the hiring probability as a target allows us to assess how the LLM synthesizes and weighs recruitment criteria beyond technical matching. The 10-point probability scale reflects current practices in LLM-based résumé evaluation [14, 17, 6, 18]. The prompt also asks the model to produce a brief explanation of its reasoning, which has been shown to enhance performance on complex evaluation tasks [19].

Moreover, the model receives no instructions to correct for potential biases or to adopt a fairness-aware approach. This is deliberate as our goal is to observe the model's default behavior and document implicit biases when evaluating equivalent candidates. To account for potential stochasticity in LLM outputs, we repeat the scoring three times

²Reputation signals are commonly used on freelancing platforms, where trust signals significantly influence recruiter decision-making.

³We define five distinct reputation profiles based on combinations of these signals. Note that freelancers with 0 completed projects cannot hold a badge or receive ratings.

⁴We exclude Arabic female names due to their extremely low prevalence in the platform's full-stack developer segment.

⁵We selected this model based on three criteria: cost-efficiency for large-scale experiments, multilingual capabilities for French interactions, and its actual usage on the studied freelancing platform.

per profile-brief pair and use the mean score for all subsequent analyses⁶. The complete prompt is provided in Appendix A.3.

Discussion Our results may be influenced by several methodological choices in prompt design and profile formatting. First, prompting for hiring probability may prime the model toward specific evaluation criteria compared to alternative framings (e.g., “quality score” or “fit assessment”) and could lead the LLM to incorporate factors beyond recruiter’s preferences, such as the freelancer’s likelihood of accepting the project. Second, while our synthetic profile format replicates the platform’s interface design, certain features receive more detailed descriptions than others, and profiles appear more structured than typical real-world résumés, both factors may shape evaluation patterns. Third, alternative evaluation procedures such as binary decisions or ranking tasks may activate different reasoning paths and represent important avenues for future research. We present a preliminary analysis of ranking prompts in Appendix C.

4. Results

We present the main findings from the evaluation methodology and dataset introduced in Section 3. First, we quantify the implicit weights the LLM uses when scoring profile-brief matches, providing a global perspective on its evaluation logic. Next, we explore heterogeneity in scoring across subgroups on the supply side, focusing on gender, ethnicity, and educational background. We then investigate how the LLM’s scoring varies with the demand-side context, such as recruiter’s company size (SME vs. large corporation) and project work arrangements (onsite vs. remote, full-time vs. part-time). Finally, we show how our experimental framework can be easily extended to human recruiter evaluations, enabling systematic comparisons of LLM and human decision logic.

4.1. Overall Effects of Attributes on LLM Scoring

Estimation Strategy We estimate the LLM’s implicit weights using Ordinary Least Squares (OLS) as a descriptive analytical tool. Specifically, we regress the average score assigned to each profile–brief pair on freelancer characteristics, clustering standard errors at the brief level to account for within-brief correlations. Importantly, this approach does not assume that the LLM itself follows a linear model, instead, we use OLS as a convenient tool to summarize average effects, estimate standard errors, and report interpretable coefficients. We follow established practices in model attribution literature where interpretable linear approximations are used to interpret complex non-linear systems ([7]). Given our fully randomized factorial design, OLS coefficients can be interpreted in our context as the marginal causal effect of each attribute on the LLM’s scoring decision⁷. Further details on model specification and covariate construction are provided in Appendix B. The

distribution of scores is centered at 6.5 with a standard deviation of 1.3, and scores range from 3 to 9.⁸

Figure 1 presents the estimated OLS coefficients, showing how the LLM’s scores vary with each attribute, relative to a reference profile. The reference profile represents a (presumably) strong candidate from the majority demographic: a European male with a master’s degree, perfectly matched technical skills and experience level, an aligned daily rate, high platform reputation, relevant industry background, aligned work arrangements preferences, and prior experience in a large company. Each coefficient represents the score adjustment when a candidate attribute deviates from this reference candidate.

Table 1 summarizes the relative importance of each attribute in the LLM’s decision-logic. For each attribute, we report the largest coefficient observed across all its possible values (for example, the maximum coefficient among all daily rate levels). These coefficients represent the maximum causal effect of varying an attribute from the reference profile. The reported values and their rankings are presented in decreasing order of absolute magnitude.

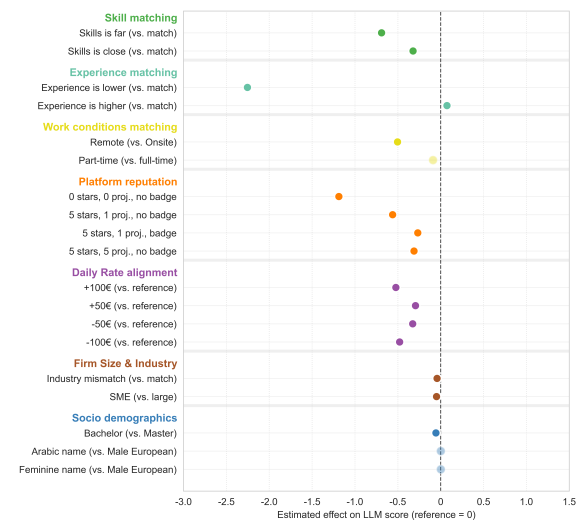


Figure 1: LLM attribute weights and relative importance. *Top:* Estimated coefficients from an OLS regression of Gemini 2.0 scores on candidate characteristics. Positive values indicate increased scores relative to the reference profile; negative values indicate a decrease. Non-significant coefficients ($p > 0.05$) are displayed with transparency.

Bottom: Relative importance of each attribute. Rank and value of the largest causal effect (i.e., the maximum coefficient observed across all levels within that variable) on the LLM’s score, relative to the reference profile.

While our approach does not assume any specific functional form, the high explanatory power of the OLS regression ($R^2 = 0.90$) suggests that the LLM behaves, in this context, as if it applies a weighted sum of profile attributes.⁹

⁶Only 2.48% of pairs show variation across runs, with maximum changes of 1 point.

⁷similarly to differences in conditional means

⁸The distribution of scores is consistent with our synthetic data design: all profiles are designed to be at least minimally relevant to the briefs, which limits the occurrence of very low scores; and the lack of perfect scores (10) likely reflects conservative scoring tendencies of the LLM.

⁹We do not interpret the R^2 as a measure of predictive accuracy, but rather as an indication of how well an additive linear model approximates the LLM’s scoring behavior. A Random Forest trained on the

Interpretation of Results This section interprets the main patterns from Figure 1.

Skill and experience level are important attributes. The LLM penalizes skill mismatches, with larger penalties for distant (-0.7 points) rather than close substitutes (-0.3 points), suggesting that the LLM understands these nuances. The LLM provides a modest reward (+0.07 points) for profiles with a higher experience level than required, but strongly penalizes those with an experience level below requirements (-2.25 points). The penalty on profiles with a lower experience level than required is the strongest across all attributes, the penalty on distant skills mismatch is the third largest.

Work conditions, such as remote or part-time availability, have different impacts on the LLM’s scoring. Profiles indicating a preference for remote work are heavily penalized when the brief requires onsite presence (-0.5 points), a magnitude comparable to that for skills mismatch. Interestingly, preferences for part-time work have only a small and statistically insignificant effect (-0.09 points). The LLM’s explanations (e.g., “schedule can be adjusted during negotiations”) suggest that such arrangements are considered negotiable.

Platform reputation is an important attribute. Candidates with no visible reputation incur the second largest penalty (-1.2 points), while even minimal signals (e.g., one past project) mitigate this effect by half (-0.56 points). Interestingly, the badge appears as a particularly strong signal, on par with multiple completed projects. These findings are consistent with the literature on reputation in online labor markets, which shows that recruiters heavily rely on trust signals [3].

Daily rates are also important for the LLM and deviations from the recruiter’s target are penalized with magnitudes comparable to those applied for skills mismatch (between -0.3 and -0.6 points). Interestingly, weights exhibit an inverted U-shaped relationship with daily rates, penalizing both daily rates lower and higher than the target. The model seems to interpret low rates as signs of inexperience or low confidence (e.g., “rate is below the offer, which may indicate limited experience”), consistent with economic signaling theory [20, 21].

Firm size and industry from the freelancer’s working history experience have a modest impact on the LLM’s evaluations. Profiles with experience in the recruiter’s industry or at large established firms receive small positive rewards (+0.05 and +0.04 points, respectively), with a magnitude about ten times lower than for skills mismatch. This suggests that the LLM values these attributes in a freelancer’s working history but only to a modest extent.

Sociodemographic features have the lowest effect on the LLM’s scoring. Profiles with a bachelor’s degree are slightly penalized compared to master’s holders (-0.06 point). Names perceived as female or Arabic-sounding have positive but small and statistically insignificant effects on scores (+0.002 and +0.001, respectively)

Overall, the LLM’s scoring behavior appears consistent with standard economic theory: it rewards close alignment

same inputs achieves an in-sample R^2 of 0.93, only slightly higher than the linear model. This small gain suggests that, while some nonlinear interactions may exist, the LLM’s evaluation logic is largely additive in practice.

between freelancer profiles and project requirements, and it values clear signals of trust (such as platform reputation) and competence (including higher education, relevant industry background, and experience at large or well-known firms).

In quantitative terms, the LLM strongly relies on skills, daily rates, and workplace preferences. In contrast, it assigns minimal importance to the characteristics of firms from the freelancer’s working history (size and industry), part-time versus full-time availability, and socio-demographic characteristics such as gender, ethnicity, or education level. The LLM most strongly penalizes freelancers with insufficient work experience and those with no prior activity on the platform.

Some patterns raise deeper questions about how the LLM interprets certain signals. For instance, all deviations in daily rates are heavily penalized. While it is intuitive that, all else equal, higher daily rates are undesirable for recruiters because they induce higher costs, the rationale for penalizing freelancers with lower daily rates is not straightforward. Low daily rates may signal lower competence or shorter experience as a freelancer and can thus justify that some recruiters would prefer to avoid these candidates. But they could also be explained by strategic positioning or financial urgency, in which case disregarding these profiles may be inefficient for both recruiters and freelancers. Given the magnitude of the penalties, future work should explore how LLMs infer such signals.

Finally, assessing whether the magnitudes of the weights align with human preferences and societal values requires a comparison with human recruiters. We will provide preliminary insights on this in Section 4.4.

Robustness Checks To test the robustness of our findings across different evaluation contexts, we replicate the analysis in a ranking task where the model directly compares profiles side-by-side rather than scoring them individually. While the overall structure remains similar, we observe notable shifts in feature importance, particularly a stronger reliance on skill matching and socio-demographic characteristics highlighting that some biases may emerge only in relative evaluation settings (see Appendix C for detailed results). We also extend our analysis to the communication sector using SEO copywriting briefs and adapted candidate profiles (Appendix D). The results exhibit similar patterns, with stronger penalties for skill mismatch and the emergence of a slight average gender bias.

4.2. Subgroup Analysis of LLM Scores by Freelancer Characteristics

To extend the analysis on sociodemographic features, we focus here on three key axes of heterogeneity: perceived gender, perceived ethnic origin, and education level.

Estimation Strategy In the previous section, we estimated the overall contribution of candidate characteristics (e.g., skills, experience, rates) to LLM scoring by running an OLS regression across all profiles. We concluded that demographic indicators do not have, on average, a significant impact on scoring. We now turn to more subtle forms of bias, asking the following question: Does the LLM apply different evaluation rules depending on the candidate’s demographic group? In other words, we investigate whether

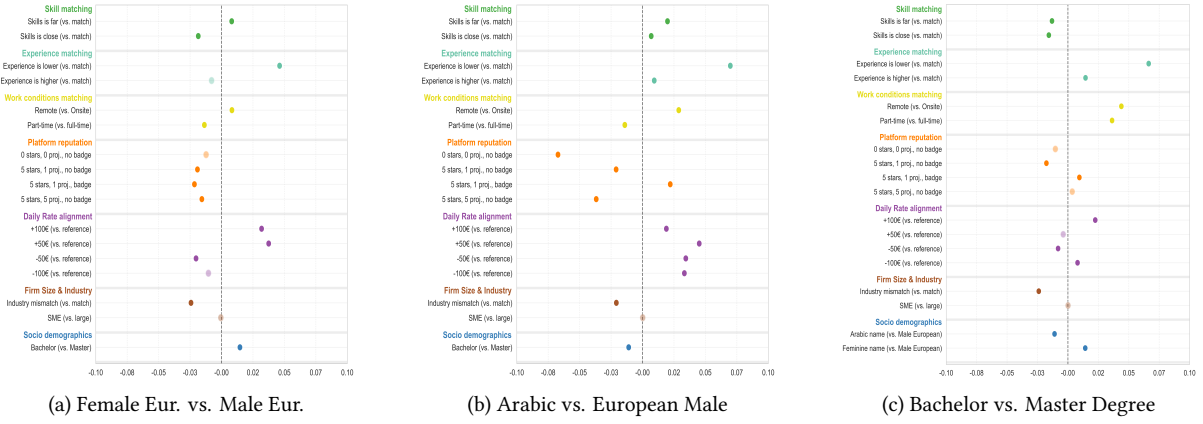


Figure 2: Interaction effects by subgroup status Coefficients reflect how the effect of each feature differs for the indicated subgroup compared to the reference group. Positive values indicate that the feature is less penalizing (or more rewarded); negative values indicate stronger penalties. These are interaction terms from an OLS regression and capture differences in evaluation logic — not average score differences. For example, a negative value for “Experience is lower × Female” means lower experience is more heavily penalized for women. Transparency indicates non-significance at the 95% level.

each candidate characteristic (e.g., experience, skills, pricing) has the same marginal impact on LLM scoring across different demographic subgroups.

To detect these patterns, we estimate a pooled OLS model **with interaction terms** between candidate characteristics and demographic indicators (gender, ethnic origin and education level). The coefficients on these interaction terms reveal whether the importance of a given feature changes depending on group membership (see Appendix B for methodological details). Figure 2 presents the results. Each coefficient indicates how much more strictly or leniently a specific characteristic is weighted for a given group, relative to the reference group. Our analysis reveals that the LLM does not apply uniform evaluation standards across all candidates.

Interpretation of Results This section interprets the key patterns revealed in Figure 2.

Female European profiles receive no baseline scoring penalty (coefficient close to 0), yet the model applies fundamentally different evaluation criteria. This creates a complex discrimination pattern. Women face stricter standards on industry alignment and are also more penalized on underpricing. Conversely, the model shows greater leniency toward educational gaps (bachelor’s degree less penalized) and lack of experience. Interestingly, overpricing becomes less penalizing for female candidates, possibly reflecting assumptions about greater recruiter bargaining power or perceived underconfidence.

Arabic male profiles face a baseline penalty (−0.026 pts or −0.31%) when all characteristics match the brief perfectly. Yet, the model shows more tolerance on several key dimensions. Experience gaps are less penalized (+0.071 pts or +3.16%) and more than compensate the difference in baseline. Skill mismatches receive greater tolerance (+0.019 pts for distant skills or +2.75%), and work arrangement mismatches are treated more leniently (+0.029 pts or +5.77%). On the contrary, the model imposes stronger penalties for weak platform reputation than it does for other groups: for instance, −0.039 pts (or −12.51%) for “5 projects, no badge”, and −0.068 pts (or −5.73%) for “0 project, 0 badge”. In contrast, lower daily rates are more generously rewarded: the effect of a −100€ rate gap is +0.033 pts (or +6.89%) compared to the reference.

Bachelor degree holders profiles holding only a bachelor’s degree are penalized by nearly −0.10 points, a substantial and significant reduction relative to otherwise identical master-level candidates at baseline. On certain criteria, bachelor-level candidates appear to be treated more leniently. For instance, being less experienced is less penalized than for master profiles (+0.066), and a +100€ daily rate is less harshly sanctioned. Similarly, the remote mismatch penalty is softer (+0.043). These findings suggest that the model may lower its strictness for candidates with lower formal education, potentially compensating for the educational gap. However, this flexibility does not extend to all aspects. Bachelor-level profiles are more penalized for skill mismatches (−0.016 for close match, −0.014 for far). Finally, a striking asymmetry appears when intersecting education with demographic characteristics: bachelor profiles with Arabic-sounding names face an additional penalty of −0.011, while bachelor-level women receive a slight positive adjustment of +0.014.

These patterns reveal that the evaluation logic of the LLM doesn’t simply apply uniform criteria, it constructs different evaluation scheme for different demographic groups, each with distinct standards, tolerances, and expectations. This systematic differentiation may perpetuate and amplify existing labor market inequalities. While these interaction effects are smaller in magnitude than the main effects of professional characteristics (our baseline model already explains 90% of variance), their systematic nature demonstrate that the LLM applies fundamentally different evaluation criteria depending on demographic group membership, creating distinct pathways to achieving high scores across different populations.

4.3. Subgroup Analysis of LLM Scores by Recruiter and Brief Characteristics

We analyze how brief characteristics systematically alter the LLM’s evaluation patterns, examining differences across: (a) large firms vs. SMEs, (b) remote vs. on-site contracts, (c) full-time vs. part-time engagements (Figure 3) and perceived recruiter’s gender.

Large Firm Briefs apply stricter standards. They impose significantly higher penalties for experience gaps and show

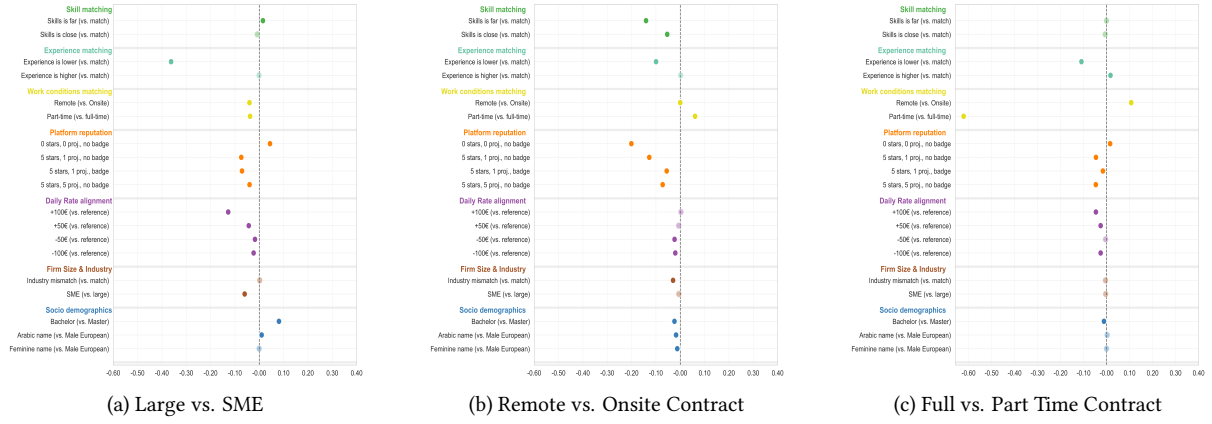


Figure 3: Interaction effects by brief characteristics Coefficients capture how the impact of each candidate feature varies with brief characteristics. Positive values indicate the feature is less penalizing (or more rewarded) under the condition (e.g., large firm, remote work); negative values indicate stricter evaluation. These are interaction effects from a pooled OLS regression without brief fixed effects. Transparency indicates non-significance at the 95% level.

greater sensitivity to pricing deviations than SMEs. The latter is surprising as large firms typically have greater budget flexibility, suggesting that corporate environments may prioritize market-rate alignment as a signal of candidate sophistication rather than pure cost considerations. They are also less flexible on work arrangements and less ready to accept mismatches in remote or onsite preferences.

Remote Contracts place significantly higher weights on skill matching, experience, and platform reputation. In remote settings, recruiters might demand stronger guarantees of quality and reliability, as they cannot rely on direct monitoring or face-to-face interactions. Interestingly, female and Arabic profiles receive slightly lower ratings when the contract is remote, suggesting that distance may exacerbate uncertainty and amplify taste-based biases.

Full Time Contracts penalize low-experience candidates more heavily, likely because full-time engagements represent greater commitment and require stronger credentials. They show higher sensitivity to work frequency mismatches but greater tolerance for remote work mismatches, suggesting that location flexibility becomes more acceptable for longer-term commitments.

Perceived Recruiter’s Gender does not significantly affect evaluation patterns. The weighting of key candidate characteristics remains largely similar, suggesting limited influence of perceived recruiter gender on scoring logic.

4.4. Estimating Human Recruiter Evaluation Patterns

To enable direct comparisons between LLM and human decision-making, our methodology can be extended to human recruiters. The experimental design for humans closely mirrors that used for LLMs and recruiters are presented with the same synthetic profiles and asked to evaluate fit with their project. However, a key methodological challenge is to elicit truthful responses from human evaluators in an experimental context. To address this, the economic literature advocates the use of an Incentivized Resume Rating (IRR) experiment [13] in which recruiters are incentivized to reveal their genuine preferences. In practice, recruiters are informed that the experiment uses synthetic data and that the recommendations they receive will be based on the

preferences they express during the evaluation. This incentive structure encourages honest and thoughtful responses, even in an artificial setting.

Once responses are collected, the analysis of human recruiters’ implicit weights follows the same procedure as for the LLM, enabling direct and systematic comparison of decision logic across the two settings. This approach provides a robust and transparent framework for benchmarking LLM-based recruitment decisions against human preferences, and for informing alignment or fairness interventions as these systems are deployed in practice.

5. Conclusion

This paper introduced a methodological framework to audit the implicit decision logic of large language models in hiring scenarios. Using a fully factorial design on synthetic freelancer profiles and project briefs, we analyzed how *Gemini 2.0 Flash* scores candidates, how this scoring varies across candidate subgroups and briefs. Our findings lead to three key takeaways:

1. **LLMs apply structured grounded logic but some interpretations raise questions.** The model assigns strong and consistent weights to core productivity signals such as skills, experience, and platform reputation. However, features like the daily rate for instance are interpreted beyond their explicit matching value both underpricing and overpricing are penalized, potentially reflecting assumptions about confidence or negotiation behavior. These inferences, though potentially rational, are not always grounded in explicit input and may introduce unintended disparities.
2. **Disparities in gender emerge through interactions.** We find little to no evidence of discrimination against women or Arabic-sounding names in average scores. Yet, intersectional effects reveal questioning patterns: penalties are more severe for minority candidates when combined with weaker signals (e.g., lower reputation or partial skill match). This highlights the need for fairness audits to consider interaction effects, not just average score gaps.

3. **Our methodology is readily applicable to human decision-makers to assess alignment** The experimental framework we introduce can be directly used to analyze human recruiter decision patterns, by presenting them with the same controlled profiles and project descriptions as those shown to the LLM. This enables systematic comparison of how sensitive attributes and productivity signals are weighted by both humans and models, facilitating robust assessment of potential biases and alignment.

5.1. Limitations & Alleys for Future Work

This study has several limitations that provide avenues for future research.

- **Model generalizability.** This paper focuses on a single model (Gemini 2.0 Flash). As such, our findings cannot be generalized to other architectures, model sizes, or training paradigms. Applying the same experimental framework to a broader set of LLMs, especially comparing reasoning, and lightweight models would allow identify whether observed behaviors are model-specific or structural.
- **From scoring to ranking.** While our primary focus is on absolute scoring, hiring often involves comparing multiple candidates. Preliminary evidence suggests that in such ranking settings, the model places stronger emphasis on certain implicit signals but a more formal analysis is needed to evaluate how LLMs perform in ranking tasks. Future work could test prompting strategies adapted to this setting (pairwise, listwise, or setwise) to better capture the model’s positional sensitivity and unintended signal amplification.
- **Prompt sensitivity.** LLM outputs are known to be highly sensitive to prompt formulation. While we use a controlled and standardized prompt in this study, future work could test the robustness of the scoring logic to alternative phrasings. This could include real prompts collected from users, helping ground the evaluation in actual usage.
- **Sectoral and linguistic generalization.** Our setting focuses on tech freelancing and uses French prompts. While we test communication briefs in the appendix, extending the framework to other domains especially those with different demographic patterns and less structured tasks would help assess external validity. To broaden the scope, future work should also explore how the model behaves with profiles and briefs in other languages and cultural settings, especially in languages where gender is less explicitly encoded than in French, which may influence how demographic signals are interpreted by the model.
- **Aligning LLMs without amplifying bias.** A critical open question is whether LLMs should be aligned with human preferences and how to do so without replicating human biases. Exploring alignment strategies such as prompt calibration, filtered examples, or fine-tuning on fairness-aware objectives could offer ways to reconcile model performance with ethical concerns.

- **Ethical considerations** Beyond these methodological considerations, we note that the deployment of LLMs in recruitment domains raises broader ethical concerns that extend beyond the scope of this study. Even when accounting for fairness considerations, as we attempt to do here, the use of algorithmic systems in hiring decisions introduces questions about algorithmic accountability, impact on recruitment employment, transparency, and candidate rights that warrant careful consideration in real-world applications.

Beyond these directions, we believe this framework contributes to a broader research agenda that seeks to understand how LLMs internalize decision rules in high-stakes applications. Our approach could be extended to other domains where decisions involve trade-offs, fairness concerns, and implicit assumptions.

Acknowledgments

We would like to thank Guillaume Bied for his valuable comments and feedback on earlier versions of this manuscript. His suggestions greatly contributed to improving the quality of this paper.

Declaration on Generative AI

The authors used Gemini-2.0-flash to: Generate datasets. During the preparation of this work, the authors used Writefull in order to: Grammar and spelling check. After using these tools, the authors reviewed and edited the content as needed and take full responsibility for the publication’s content.

References

- [1] M. Spence, Job market signaling, *The Quarterly Journal of Economics* 87 (1973) 355–374.
- [2] J. J. Horton, Do digital skill certificates help new workers enter the market? evidence from a voluntary certification scheme in an online freelancing labor market, *OECD iLibrary Working Paper* (2017).
- [3] A. Pallais, Inefficient hiring in entry-level labor markets, *American Economic Review* 104 (2014) 3565–3599.
- [4] M. Bertrand, E. Duflo, Field Experiments on Discrimination, Working Paper 22014, National Bureau of Economic Research, 2016. URL: <http://www.nber.org/papers/w22014>. doi:10.3386/w22014.
- [5] V. Bartoš, Attention discrimination: Theory and field experiments with monitoring information acquisition, *American Economic Review* 106 (2016) 1644–1675.
- [6] L. Armstrong, A. Liu, S. MacNeil, D. Metaxa, The silicon ceiling: Auditing gpt’s race and gender biases in hiring, in: *Proceedings of the 4th ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization*, 2024, pp. 1–18.
- [7] M. T. Ribeiro, S. Singh, C. Guestrin, “why should i trust you?” explaining the predictions of any classifier, in: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 2016, pp. 1135–1144.

- [8] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, *Advances in neural information processing systems* 30 (2017).
- [9] Y. Gat, N. Calderon, A. Feder, A. Chapanin, A. Sharma, R. Reichart, Faithful explanations of black-box nlp models using llm-generated counterfactuals, *arXiv preprint arXiv:2310.00603* (2023).
- [10] R. K. Sinha, Book review: Christoph molnar. 2020. interpretable machine learning: A guide for making black box models explainable, *Metamorphosis* 23 (2024) 92–93. URL: <https://doi.org/10.1177/09726225241252009>. doi:10.1177/09726225241252009.
- [11] D. Janzing, L. Minorics, P. Blöbaum, Feature relevance quantification in explainable ai: A causal problem, in: *International Conference on artificial intelligence and statistics*, PMLR, 2020, pp. 2907–2916.
- [12] D. Slack, S. Hilgard, E. Jia, S. Singh, H. Lakkaraju, Fooling lime and shap: Adversarial attacks on post hoc explanation methods, in: *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, 2020, pp. 180–186.
- [13] J. B. Kessler, C. Low, C. D. Sullivan, Incentivized resume rating: Eliciting employer preferences without deception, *American Economic Review* 109 (2019) 3713–44. URL: <https://www.aeaweb.org/articles?id=10.1257/aer.20181714>. doi:10.1257/aer.20181714.
- [14] S. Vaishampayan, H. Leary, Y. B. Alebachew, L. Hickman, B. A. Stevenor, W. Beck, C. Brown, Human and llm-based resume matching: An observational study, in: *Findings of the Association for Computational Linguistics: NAACL 2025*, 2025, pp. 4808–4823.
- [15] F. P.-W. Lo, J. Qiu, Z. Wang, H. Yu, Y. Chen, G. Zhang, B. Lo, Ai hiring with llms: A context-aware and explainable multi-agent framework for resume screening, in: *Proceedings of the Computer Vision and Pattern Recognition Conference*, 2025, pp. 4184–4193.
- [16] M. Bertrand, S. Mullainathan, Are emily and greg more employable than lakisha and jamal? a field experiment on labor market discrimination, *American Economic Review* 94 (2004) 991–1013. doi:10.1257/0002828042002561.
- [17] P. Ghosh, V. Sadaphal, Jobrecogpt—explainable job recommendations using llms, *arXiv preprint arXiv:2309.11805* (2023).
- [18] J. D. Gaebler, S. Goel, A. Huq, P. Tambe, Auditing the use of language models to guide hiring decisions, *arXiv preprint arXiv:2404.03086* (2024).
- [19] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al., Chain-of-thought prompting elicits reasoning in large language models, *NIPS* (2022).
- [20] P. Milgrom, J. Roberts, Price and advertising signals of product quality, *Journal of Political Economy* 94 (1986) 796–821. URL: <https://ideas.repec.org/a/ucp/jpolec/v94y1986i4p796-821.html>. doi:10.1086/261408.
- [21] C. Shapiro, J. E. Stiglitz, Equilibrium unemployment as a worker discipline device, *The American Economic Review* 74 (1984) 433–444. URL: <http://www.jstor.org/stable/1804018>.

A. Synthetic Data Generation Procedure

A.1. Profile Generation

Example of Generated Freelancer Profile

Name: Marie

Headline: Full Stack Developer (female)

Platform Activity:

- Projects signed: 1
- Client reviews: 1
- Average client rating: 5/5
- Super Malter badge: No

General Information:

- Indicative rate: €450/day
- Experience: 0–2 years
- Availability: Confirmed
- Workload: Full-time
- Location: Île-de-France
- Onsite availability: Can work in client offices

Languages:

- French: Native or bilingual

Skills:

- JavaScript / TypeScript
- NestJS

Professional Experience:

The freelancer has 1 year of professional experience, including the following example:

- AXA (Banking and Insurance) — Full Stack Developer*
- Developed new features for a payment engine.
 - Migrated existing code to a more modern architecture.
 - Updated application interfaces.

Education:

- Master’s degree, INSA Lyon (Institut National des Sciences Appliquées de Lyon)

A.2. Brief Generation

Example of Generated Project Brief

Title: Full Stack Developer

Project Description:

We are a CAC40 company looking for a developer to handle the backend of a web project and provide support on the frontend. Tasks include receiving files, interacting with third-party services, assembling templates, sending emails, and managing payments. 5 years of experience is preferred.

Contract Details:

- Workload: Full-time

Table 1
Variable Elements in Experimental Design

(a) Freelancer Profile Generation		(b) Project Brief Generation	
Attribute	Values	Attribute	Values
Name	Female EU, Male EU, Male Arabic	Recruiter Name	Female EU, Male EU
Work Cond.	Full-time, Part-time (3d/week)	Company Size	Large Corp., SME
Flexibility	On-site, Remote	Work Location	On-site req., Remote allowed
Firm Size (exp.)	Large Corp., SME	Work time	Full-time (5d/week), Part-time (3d/week)
Industry (exp.)	Banking, E-Commerce	Industry	Banking & Insurance
Years Experience	1, 5, 9 years	Duration	6 months
Education	Bachelor's, Master's	Min. Experience	5 years
Daily Rate (€)	300, 350, 400, 450, 500	Compensation	€400/day
Tech Stack	JS+Node.js, JS+NestJS, JS+Angular	Technology	JavaScript/TypeScript + Node.js
Platform Rep.	None, Starter (1p/5s), Expert (5p/5s), ±Platform badge	Total: 16 combinations	
Total: 10,800 combinations			

- Remote work: Allowed
- Duration: 6 months
- Start date: As soon as possible
- Location: Île-de-France
- Daily rate: €400/day

Required Skills:

- JavaScript / TypeScript
- Node.js

Recruiter Information:

- Recruiter first name: Alexandre
- Industry: Banking and Insurance
- Company size: 1,000–4,999 employees

A.3. Details on Scoring Prompt

The following prompt template was used to evaluate freelancer-project matching using Gemini Flash 2.0. Template variables profile and brief were dynamically replaced with profiles and briefs text.

Profile/Brief Scoring Prompt

Role: You are an expert in matching fullstack tech freelancers with projects on the Malt freelance platform.

Task: Evaluate hiring probability of a provided freelancer profile for a project brief given below.

Input Format:

- **Freelancer Profile:** profile
- **Project Brief:** brief

Instructions:

1. Evaluate the probability of hiring (0-10 scale)
2. Provide concise justification (2 sentences maximum)
3. Follow strict output format

Required Output Format:

Score : [your score here]/10
Justification : [your justification here]

The model was queried via the Vertex AI API using **Gemini Flash 2.0**. We used the default decoding parameters provided by Vertex AI: top_k = 40, top_p = 0.95. To re-

duce output variance and encourage consistent scoring, we set temperature = 0. Max output token was set to max tokens = 256

B. Estimation Strategy

B.1. Average Implicit Weights Analysis

To uncover the implicit decision logic used by the LLM when evaluating profile-brief, we estimate a linear regression of the score assigned to each profile–brief pair on a set of profile characteristics and computed matching profile-brief. Since our synthetic profiles are generated through a full factorial design where all attributes are systematically varied, the regression coefficients directly capture the average treatment effects without confounding. Therefore, in this experimental setting, OLS coefficients are interpreted as simple conditional mean differences. The choice of a linear model is only due to its transparency and interpretability and allows for causal interpretation of the coefficients as marginal contributions to the LLM’s scoring function.

Each characteristic is encoded as categorical or binary variables with a designated reference category representing the ideal candidate from majority groups (matching experience and daily rate, Master’s degree, male European, etc.). The coefficients quantify how each characteristic influences the LLM’s scoring decisions, measuring the marginal impact of moving from the reference to each alternative level (e.g., from male European to female European, or from 3 to 1 years of experience).

$$score_{ij} = \alpha + \sum_k \beta_k x_{ik} + \sum_n \beta_n m_{ijn} + \epsilon_{ij} \quad (1)$$

where :

- α is the intercept, representing the expected score for a profile with all reference characteristics (we choose to have as a reference a perfect match on all dimensions and the majority group)
- $score_{ij}$ is the average score assigned by the LLM to profile i and brief j .
- x_{ik} represent the k characteristic of profile i (e.g education level, gender)

- β_k captures the marginal contribution of each profile feature to the LLM’s evaluation
- m_{ijn} represents the matching feature between brief j and profile i (e.g. skill alignment)
- β_n captures the penalty or bonus applied by the scoring logic for specific matching features
- ε_{ij} is the error term

Each profile–brief pair is submitted three times to the LLM scoring prompt to account for potential variance in the model’s outputs. We compute the mean score across three independent runs and use it as the dependent variable in the regression.

Because each brief is matched with multiple profiles, the resulting scores may have intra-brief correlation. To account for within-group error correlation we cluster standard errors at the brief level when possible.

B.1.1. Heterogeneity in LLM Scoring Logic Across Groups

To assess whether the LLM applies different evaluation logics across subgroups, we investigate the heterogeneity of estimated coefficients by profile type. In particular, we are interested in understanding whether certain profile characteristics or matching dimensions are valued differently depending on observable attributes of the profile.

To formally test whether coefficients differ significantly across groups, we estimate the following interaction model:

$$\text{score}_{ij} = \alpha + \theta \cdot \text{Group}_i + \sum_k \beta_k x_{ik} + \sum_n \delta_n m_{ijn} + \sum_k \phi_k \cdot (x_{ik} \cdot \text{Group}_i) + \sum_n \lambda_n \cdot (m_{ijn} \cdot \text{Group}_i) + \varepsilon_{ij} \quad (2)$$

where:

- Group_i is a binary or categorical indicator for the group to which profile i belongs (e.g. Female, Master’s degree).
- x_{ik} represents the k static characteristic of profile i (e.g. education).
- m_{ijn} represents the n matching feature between profile i and brief j (e.g. skill alignment, experience gap).
- β_k measures the effect of characteristic x_{ik} for the reference group.
- δ_n measures the effect of matching variable m_{ijn} for the reference group.
- θ captures the baseline difference in score for the group when all other variables are at their reference levels.
- ϕ_k measures how the effect of x_{ik} differs for the group compared to the reference group.
- λ_n measures how the effect of m_{ijn} differs for the group compared to the reference group.

This interacted specification is formally equivalent to estimating separate regressions by group. In other words, estimating the model separately by group and then taking the difference in coefficients across groups gives equivalent values to the interaction terms in the pooled model. However, this pooled regression allows us to test statistically for the difference in coefficient.

As an illustrative example, consider a binary indicator Group_i equal to 1 for female profiles, and a single matching variable m_{ij} equal to 1 if the profile possesses the required skill listed in the brief. In this case, a significant and negative interaction coefficient λ implies that women are more heavily penalized than men when the required skill is missing. Conversely, a positive value of λ indicates that women are less penalized than men for the same mismatch. If λ is not significantly different from zero, the penalty associated with skill mismatch is applied equally across subgroups.

B.1.2. Heterogeneity in Scoring Logic Across Brief Characteristics

In addition to testing whether profile groups are treated differently by the LLM, we also investigate whether the scoring function differs depending on the characteristics of the brief. This analysis helps us understand if the LLM adapts its evaluation logic to the context provided by the brief.

To test this, we interact profile or matching variables with one brief-level characteristic.

$$\text{score}_{ij} = \alpha + \sum_k \beta_k x_{ik} + \sum_n \delta_n m_{ijn} + \eta \cdot b_j + \sum_n \lambda_n (m_{ijn} \cdot b_j) + \varepsilon_{ij} \quad (3)$$

where:

- b_j is a binary indicator for a specific brief attribute (e.g., remote brief, large firms, full time),
- λ_n captures whether the effect of matching feature m_{ijn} differs depending on the value of b_j .

As an illustrative example, consider a binary indicator b_j equal to 1 if the brief is written by a large firm (vs SME) and a single matching variable m_{ij} equal to 1 if the profile possesses the required skill listed in the brief. In this case, a significant and negative interaction coefficient λ implies that large firms penalize skill mismatches more strongly than SME.

C. Scoring vs Ranking

C.1. Procedure

While the main analysis is based on a scoring task where the model evaluates profiles individually, many real-life recruitment decisions involve ranking candidates relative to each other. To assess whether the model’s decision logic shifts in such contexts, we replicate our analysis using a ranking task, where the LLM is prompted to rank three profiles for the same brief. We sample random groups of three profiles (without replacement) in our synthetic dataset and ask the LLM to rank them according to the probability of hiring for a given project brief. The following prompt template was used to evaluate freelancers-project ranking

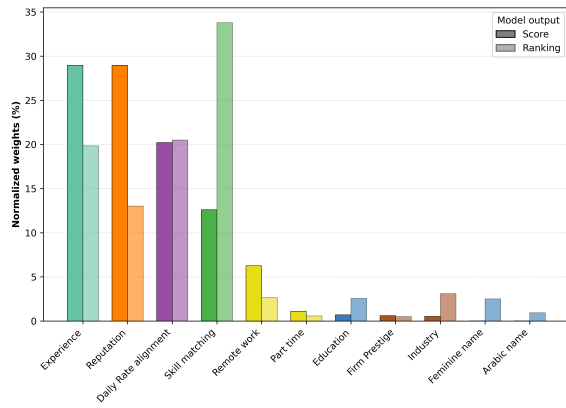


Figure 4: Normalized weights in Scoring and Ranking procedure. Weights are computed by normalizing the absolute value of each coefficient by the sum of all coefficients.

using Gemini Flash 2.0. Template variables `profile` and `brief` were dynamically replaced with profiles and briefs text.

Profile/Brief Ranking Prompt

Role: You are an expert in matching fullstack tech freelancers with projects on the Malt freelance platform.

Task: Evaluate hiring probability of a provided freelancer profile for a project brief given below.

Input Format:

- **Freelancer Profile:** `profile`
- **Project Brief:** `brief`

Instructions:

1. Evaluate the probability of hiring (0-10 scale)
2. Provide concise justification (2 sentences maximum)
3. Follow strict output format

Required Output Format:

Score : [your score here]/10

Justification : [your justification here]

C.2. Results

Figure 4 compares the normalized weights of key features across the scoring and ranking settings. In the ranking task, skill matching becomes by far the most influential factor, while the weights of experience and platform reputation drop significantly. At the same time, demographic characteristics, such as arabic or feminine-sounding names, education but also firm prestige and industry matching gain more weight than in the scoring setup. This suggests that some features may play a limited role in absolute scoring, but become more prominent when the model must make a relative judgment between similar profiles. In scoring mode, these variables may have little impact because they do not affect the likelihood of success in isolation. Yet in ranking mode when the model must refine its decision they provide discriminating signals that help differentiate otherwise similar candidates. These results highlight a key concern: biases that appear negligible in scoring may be amplified in ranking, an important issue that we leave to future work.

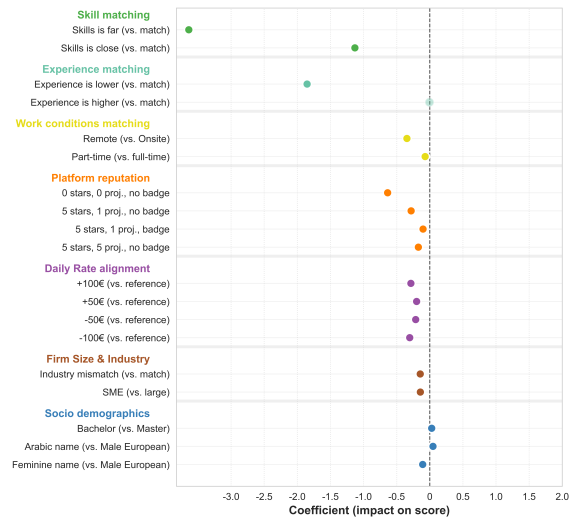


Figure 5: Estimated coefficients from the OLS regression of Gemini 2.0 scores on candidate characteristics. Positive coefficients increase the score relative to the reference profile; negative coefficients reduce it.

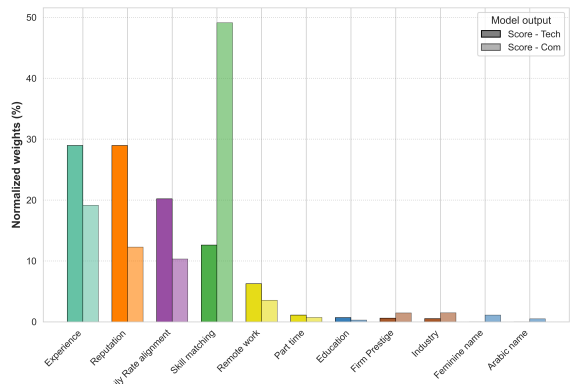


Figure 6: Normalized weights for the Technological and Communication sample. Weights are computed by normalizing the absolute value of each coefficient by the sum of all coefficients.

D. Alternative Occupation : SEO Writing Content

Procedure We replicate the same construction logic used for full-stack developers to generate synthetic profiles in the communication sector. The main adaptations concern the skill dimensions, where we vary whether the candidate possesses the brief's required skill (SEO content writing), a closely related skill (editorial content writing), or a more distant skill (proofreading). We also adapt the project briefs to reflect typical SEO requirements and adjust the main sector of activity to e-commerce, with tourism as the alternative sector. Finally, we modify the experience descriptions to reflect relevant field experience and update company names to represent both large corporations and SMEs. The rest of the variables are exactly similar.

Results We replicate our regression analysis on a second occupation, SEO content writing in the communication sector, to assess whether the model's decision logic generalizes beyond the technological domain. The model continues to explain a substantial share of the variation in evaluation

scores ($R^2 = 0.865$), suggesting a similarly additive structure in its decision-making process. The overall weighting of attributes remains consistent with the patterns observed in the technological sector.

The main difference lies in the much greater importance assigned to skill matching, though this may partly reflect the specific skill configurations we selected. We also observe a similar inverted U-shaped relationship with daily rates, although less pronounced. Importantly, in this setting, profiles perceived as feminine are significantly penalized on average, while profiles perceived as arabic are slightly advantaged on average compared to european men.