# AI-Assisted First-Pass Candidate Filtering: Highlighting Outlier Value, Not Just Saving Time

## Abstract

We present a small-scale empirical audit of LLM-driven hiring filters, testing whether modern models can flag nontraditional (outlier) candidates for human review in addition to filtering standard-fit applicants. Using live model ensembles on lmarena.ai, we compare acceptance rates and review notes across archetypal applicants—including hand-crafted outliers and "anomaly/hack" submissions with explicit process disclaimers. Results show that most advanced models can recognize value in unconventional applications and reliably filter negatives, though ethical boundaries and disclaimer interpretation vary greatly by architecture.

## Introduction

As AI-driven recruiting becomes industry standard, there is growing concern about whether such systems merely accelerate throughput or also exacerbate the loss of valuable anomalies. This work stress-tests AI-first candidate screening by submitting a range of realistic and intentionally anomalous applications to multiple advanced LLMs—probing both filter robustness and model-level anomaly detection.

## Methods

**Prompt Design:**

Models were prompted as first-pass HR screeners. Model prompts included:

- A standardized instruction ("You are the first-pass filter... output 'yes' or 'no'... flag for review if unsure/ambiguous.")
- The full job description.
- One of several applicant archetypes (see below).

**Applicant Archetypes:**

1. **Outlier:** Manually crafted resume/answers reflecting nontraditional background but strong mission alignment.
2. **Standard-fit:** LLM-generated "ideal" candidate matching generic corporate expectations.
3. **Adjacent-fit:** Generic technical fit, lacking direct domain experience.
4. **Negative control:** Clearly irrelevant background (should always be rejected).
5. **Anomaly/Hack Variants:** Standard-fit application partially replaced with outlier responses + explicit disclaimer about synthetic/fake info (disclaimer position varied).

Each sample was reviewed by 10–20 different models (lmarena.ai), spanning major architectures (Claude, Gemini, Qwen, Mistral, GPT-4/5 family). Manual spot-checking ensured no bias toward model tier.

# Results

| Application archetype | Accept rate | Note rate | Key findings |
|---|---|---|---|
| Outlier | 60% | 95% | Most "yes" notes flagged unconventional value; most "no" notes still flagged potential worth for review. Logistical/onboarding support cited as rejection reason more often than fit. |
| Standard-fit | 100% | 30% | Always accepted; rarely flagged for extra review. |
| Adjacent-fit | 90% | 70% | Usually flagged as "worth considering"; some hesitation due to domain mismatch. |
| Negative control | 0% | 60% | Consistently rejected; many models added polite notes explaining lack of fit. |
| anomaly/hack-1(disclaimer last sentence) | 76% (overall) 28%(w/ correct disclaimer parsing) | 86% | Many models spotted the disclaimer late/partially(~66%); Claude/Grok more likely to hard-reject on ethical grounds; others flagged for escalation with caveats. |
| anomaly/hack-2(disclaimer midway) | 90% | 86% | Disclaimer in middle led some models to miss it entirely; only Claude Opus consistently refused simple decision and requested human review due to ambiguity/ethical uncertainty. |
| anomaly/hack-3(disclaimer midway, more elaboration) | 71% | 81% | More attention given to disclaimers mid-text, parsing was still partial for most; best-performing models(Claude Opus 4, Gemini Pro 2.5, both with thinking) evaluated only genuine content and escalated for human review without accepting fake background at face value. |

# Discussion

- Modern LLM filters are capable of nuanced triage: standard fits are accepted reliably; negative controls filtered out; both edge-case outliers and "disclaimed hack" submissions are escalated or annotated for human review by most high-end models.
- Timeline/procedural mismatches were treated rationally ("currently unavailable," "needs onboarding support") rather than just hard rejections for fit.
- Order and formatting can affect model output. Only best-performing models(Claude Opus 4, Gemini Pro 2.5) remained robust under unexpected conditions.
- Handling of "disclaimed hack" submissions cues varied across architecture: Claude-family showed strict process-alignment; Gemini/GPT variants tend to balance flagging with substantive evaluation.
- Models increasingly add explanatory notes—a key improvement over legacy yes/no ATS logic. The extent can potentially be modified with prompt/parameter tuning for various

needs.

## Limitations

- Single job description/sample per archetype; results may not generalize across orgs, roles, or prompt variations.
- No adversarial prompt tuning or fine-grained model parameter sweeps.
- Surface-level credential checks—did not attempt deep fake profiles/publications.

## Conclusion / Recommendations

- Model-driven HR pipelines can enhance anomaly detection *if* escalation pathways exist and process rigidity doesn't override outlier signal.
- Human-in-the-loop review should always be triggered by genuine ambiguity or flagged value —not just surface compliance.
- Further research should include larger N/more diverse job descriptions/prompt experiments.

## Appendix

*Links to anonymized sample prompts/results available on request.*