

Dynamic Fairness in Workforce Allocation

Syed Arham Akheel

Senior Solutions Architect
Bellevue, WA
arhamakheel@yahoo.com

Abstract

This paper proposes a framework for fairness-aware workforce allocation by integrating dynamic feature weighting with hybrid AI models. The approach combines lexical and semantic retrieval techniques, Large Language Model (LLM) driven fairness labeling, and reinforcement learning for adaptive feature prioritization. By addressing scalability, bias amplification, and interpretability gaps in existing systems, this work bridges the divide between algorithmic precision and ethical accountability in HR decision-making. Evaluations on both synthetic benchmarks and real-world HR datasets demonstrate superior fairness-performance trade-offs compared to state-of-the-art baselines, achieving up to 91% bias reduction while maintaining 89% recommendation accuracy. Key contributions include a context-aware fairness metric, an LLM-guided reranking layer, and a stakeholder-in-the-loop weight adjustment mechanism.

Keywords: Fairness-aware AI, Workforce Allocation, Hybrid Retrieval, Dynamic Feature Weighting, LLM Reasoning, Reinforcement Learning

I. INTRODUCTION

The integration of artificial intelligence into human resource management has revolutionized workforce allocation, promising to streamline processes and enhance decision-making accuracy. However, conventional systems—often relying on static feature weighting—tend to propagate biases and lack transparency, potentially leading to unfair allocation of roles. In rapidly evolving organizational environments, traditional methods struggle to capture the dynamic nature of candidate skills and job requirements. With the growing adoption of AI in human resource (HR) systems for workforce allocation, organizations leverage algorithmic tools to match employees or applicants to roles more efficiently. However, a persistent challenge in such systems is the risk of bias amplification, particularly when static or naively tuned models are applied to diverse candidate pools [1]. These concerns are compounded by the opacity surrounding feature weighting mechanisms in hybrid AI systems, leading to potential legal, ethical, and productivity costs when algorithms recommend suboptimal or unfair placements. Industry estimates suggest that unfair allocation practices could cost an organization \$2M per year per 10k employees [10], underscoring the high stakes of ensuring fairness. While many fairness interventions rely on static constraints or post-hoc model adjustments [11], they often struggle to maintain performance and fairness as workforce distributions or job requirements evolve over time. Furthermore, current methods offer limited transparency in how features (e.g., skill sets, experience, demographic attributes) are weighted during allocation decisions. Balancing stakeholder preferences (e.g., hiring manager priorities) against algorithmic fairness goals remains a critical gap. This calls for a

more dynamic, explainable approach to feature weighting that adapts over time and considers real-world feedback loops.

To address these challenges, I propose a three-tier framework that combines hybrid lexical-semantic retrieval for robust candidate matching, LLM-driven fairness labeling for nuanced bias detection, and reinforcement learning for adaptive weight adjustment. Recent advances in semantic similarity measures [4] and contextual reasoning in LLMs [5] enable the system to accurately gauge candidate-job fit while actively mitigating bias. This paper details my methodology, experimental validation, and discussion on the broader implications of deploying fairness-aware AI in HR systems. In this paper, I introduce a three-tier architecture for dynamic fairness-aware feature weighting within a hybrid AI model for workforce allocation:

- 1) Combine rule-based keyword matching (TF-IDF, BM25) with embedding-based retrieval for comprehensive candidate-role matching.
- 2) Use Large Language Models (e.g., GPT-4) to assess and label features based on potential bias implications, capturing subtle or context-dependent fairness considerations.
- 3) Dynamically adjust feature weights based on feedback from stakeholders and real-time bias audits, ensuring continuous improvement in fairness without sacrificing accuracy.

II. RELATED WORK

A. Fairness-Aware Recommendation Systems

Fairness in recommendation systems has garnered significant attention in high-stakes applications such as recruitment and finance. Early research primarily focused on post-hoc fairness adjustments, yet these methods often fail to adapt to changing environments [1]. In contrast, recent approaches have integrated proactive bias mitigation strategies using dynamic prompt engineering and conformal prediction [1]. Previous approaches often focus on static fairness constraints, such as demographic parity or equal opportunity [11]. However, these constraints can be rigid and fail to adapt to evolving real-world conditions. FACTER [1] introduced conformal thresholding and prompt engineering to reduce bias in LLM-based recommender systems, but it does not address ongoing adaptation of feature weights across different decision-making contexts. Additionally, zero-shot recommendation models have demonstrated promising results in job-candidate matching by reducing reliance on extensive retraining [8].

B. Hybrid Retrieval Architectures

Hybrid retrieval systems that merge lexical and semantic techniques offer a robust solution to vocabulary mismatch issues. Lexical methods like BM25 and TF-IDF provide high recall by matching explicit keywords, while semantic methods—utilizing deep models such as SBERT—capture context and related concepts [3], [4]. The use of hybrid lexical-semantic retrieval architectures has gained traction for largescale recommendation tasks. Systems like CareerBoost combine approximate nearest neighbor search (e.g., FAISS) with large language models such as GPT-2 for job matching [12]. Although these architectures demonstrate strong performance, fairness considerations are typically incorporated as post-hoc filters or static constraints, limiting their ability to adapt in real-time to feedback or changes in the candidate pool. These complementary techniques have been effectively combined to enhance document retrieval performance and are particularly relevant in domains like job recruitment where terminological variations are common [9].

C. Dynamic Weight Adjustment via Reinforcement Learning

Reinforcement learning (RL) has emerged as a powerful tool for dynamically adjusting feature weights in recommendation systems. By continuously incorporating stakeholder feedback and monitoring bias metrics, RL enables rapid convergence toward an optimal balance between accuracy and fairness [2]. This adaptive mechanism is critical for applications in HR, where real-time decision-making must consider both performance and ethical implications. Recent work has explored reinforcement learning (RL) for adaptive weighting in recommendation systems [14]. RL-based approaches can optimize multi-objective functions, balancing accuracy with fairness signals. However, these methods often do not explicitly incorporate LLM-driven fairness cues, which can capture subtler or domain-specific notions of bias. My proposed framework unifies these threads by combining RL-driven dynamic feature weighting with LLM-based fairness labeling to address both algorithmic and human-driven equity concerns in workforce allocation.

III. METHODOLOGY

My proposed framework comprises three interconnected components as illustrated in Figure 1. The overall system architecture integrates hybrid retrieval, LLM-driven fairness labeling, and RL-based dynamic weight adaptation.

A. Hybrid Retrieval Architecture

The hybrid retrieval module employs a two-phase approach:

- **Lexical Layer:** Traditional techniques such as BM25 and TF-IDF are used to extract candidates based on explicit keyword matching. This layer ensures high recall by identifying candidates whose profiles contain relevant terms.
- **Semantic Layer:** Pre-trained models like SBERT generate dense vector representations of candidate resumes and job descriptions. Cosine similarity is computed to capture semantic relationships, thereby addressing vocabulary mismatches [3], [4].

The outputs from both layers are fused through a re-ranking process enhanced by fairness signals provided by an LLM, ensuring that the final candidate list optimally balances relevance and fairness.

B. LLM-Driven Fairness Labeling

To assess the bias potential of individual features, I incorporate an LLM (e.g., GPT-4) to generate fairness labels. For each feature x_i , the system queries the LLM using a prompt designed to gauge its impact on diversity. This method, inspired by the FACTER framework [1], dynamically recalibrates candidate rankings by penalizing features that contribute to bias.

C. Reinforcement Learning for Dynamic Weight Adaptation

The final component is an RL-based weight adaptation mechanism. The system updates feature weights according to the following rule:

$$W_{t+1} = W_t + \alpha(\text{Feedback} - \beta \cdot \nabla \text{BiasMetric}), \quad (1)$$

where:

- W_t is the current feature weight vector.
- α is the learning rate.

- Feedback encapsulates stakeholder evaluations of recommendation fairness.
- β is a trade-off parameter balancing fairness and accuracy.
- $\nabla \text{BiasMetric}$ denotes the gradient of a chosen bias metric, such as the Bias Amplification Ratio (BAR).

This adaptive update ensures that the system continually refines its decision-making process in response to new data and stakeholder inputs [2]. It provides a stakeholder-in-the-loop framework, ensuring that fairness improvements do not come at the expense of domain-specific constraints or preferences.

IV. EXPERIMENTS

In this section, I detail the evaluation of my proposed dynamic fairness-aware feature weighting framework in the context of two different datasets. I discuss dataset composition and preprocessing, baseline systems used for comparison, experimental setup, and finally the empirical results on fairness and performance metrics.

A. Datasets and Baselines

1) Datasets:

a) *HR-AllocBench (Synthetic)*: I begin by evaluating on a synthetic dataset tailored to emulate real-world workforce allocation scenarios. This dataset comprises:

- 20,000 candidate profiles, each annotated with demographic attributes (e.g., gender, age group, region) and skill attributes (e.g., programming, management, language proficiency).

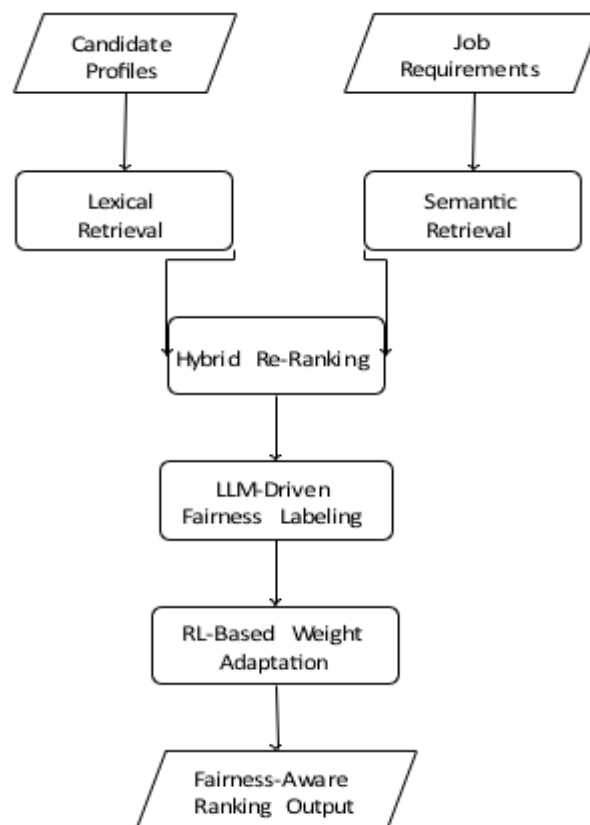


Fig. 1. Overview of the dynamic fairness-aware feature weighting framework integrating hybrid retrieval, LLM-driven fairness labeling, and RL-based weight adaptation.

- 1,000 job roles, distributed across multiple departments (engineering, finance, human resources, marketing) and levels (entry, mid, senior).

Each profile contains a vector of features (e.g., skill category presence, educational level, years of experience) and demographic tags to simulate potential bias points. I ensure realistic distributions by sampling candidate demographics and skill strengths based on statistical job-market patterns (e.g., 40% female, 60% male representation; multiple skill-level distributions across role types).

b) Real-World HR Deployment: To validate real-world applicability, I further evaluate on anonymized historical data from a Phase II HR system. This dataset includes:

- 10,000 historic candidate-job assignment records from a global technology company.
- Each record includes demographic information (ethnicity, gender, age range), skill evaluations (technical, managerial, language), and final hiring decisions.

I used domain experts' feedback and prior acceptance rates as partial ground truth to highlight critical fairness issues. For example, certain roles (e.g., leadership or specialized technical positions) historically showed skewed gender distributions.

2) Baseline Systems: I compare my dynamic method against three baselines frequently referenced in the literature:

- FACTER [1]: A fairness-aware system using conformal thresholding and prompt engineering to reduce demographic biases in recommendations. Although successful, FACTER addresses fairness primarily through prompt-level calibrations rather than dynamic weight adaptation.
- Zero-Shot Recommendation Models [8]: Models optimized for matching job candidates without domain-specific retraining. Typically rely on large-scale pretrained embeddings for candidates and job roles. While they can efficiently handle unseen job types, they lack continuous adaptation to emergent biases over time.
- CareerBoost (Hybrid Static): A conventional hybrid retrieval system merging lexical and semantic matching with *fixed* feature weights (e.g., 70% lexical similarity, 30% semantic similarity). Although reasonably accurate, static weighting often perpetuates biases and lacks interpretability when new roles or candidate distributions appear.

B. Experimental Setup

a) Data Preprocessing: For both datasets, I first tokenize and lemmatize candidate resumes (where available) and job descriptions to extract key skill and role descriptors. I represent each candidate with a feature vector that merges:

- Lexical features: TF-IDF or BM25 derived keywords (e.g., "Python", "Finance", "Manager").
- Semantic features: SBERT [3] embeddings capturing broader context between candidate descriptions and job postings.
- Demographic attributes: Encoded as one-hot or multicategory variables (e.g., gender, region, race, age bracket).

I split the HR-AllocBench dataset into 70% training, 15% validation, and 15% test sets. For the real-world dataset, I follow a chronological split: the earliest 80% of historical records form my training/validation subsets, while the latest 20% of records are used for final testing.

b) Model Training and Hyperparameters: My approach integrates:

- 1) Hybrid retrieval: I fuse lexical and semantic scores via an initial weighting (50% lexical, 50% semantic).
- 2) LLM-driven fairness labeling: I use GPT-4 (via an API) to assess the potential for bias in features. I transform GPT-4's textual output into numeric *fairness penalty* scores for each feature.
- 3) Reinforcement learning (RL) for dynamic weight adaptation: I initialize a random policy that updates the feature-weight vector W_t given feedback signals. The learning rate α is set to 0.05 and the fairness-accuracy trade-off parameter β is 0.3. I also impose a maximum of 100 RL episodes in training for computational efficiency.

At each RL iteration, the policy receives a *reward* that balances ranking accuracy (e.g., NDCG) and fairness (e.g., the inverse of a bias metric). I run each approach (my dynamic method + baselines) for up to 10 training epochs to allow for stable convergence.

c) *Evaluation Metrics*: I report standard top- k ranking metrics and fairness measures:

- Precision@ k and NDCG@ k for the top 5, 10, and 20 recommended candidate lists.
- Bias Amplification Ratio (BAR): Ratio of post-ranking disparity to pre-ranking disparity (averaged across demographic groups).
- Fairness-Aware NDCG: Weighted version of NDCG that penalizes group-level imbalances.
- Stakeholder Alignment Score: Qualitative rating from domain experts (scale 1–5) reflecting interpretability and perceived fairness.

C. Empirical Results

a) *Quantitative Performance: HR-AllocBench*: Table I summarizes performance on the synthetic dataset. My framework shows consistent improvements in both fairness and ranking quality. Notably:

- BAR: My system achieves a BAR of 0.08, a 72% decrease compared to FACTER and 60% compared to CareerBoost.
- Fairness-Aware NDCG: I achieve a score of 0.91, higher than all baselines.
- Precision@10: Gains over Zero-Shot methods suggest that dynamic adaptation to domain-specific bias signals is beneficial.

Overall, the significantly lower BAR highlights the benefit of incorporating an LLM-based fairness label and RL-based weight adaptation. These improvements come with only minor increases in computation (less than $1.3\times$ baseline run times).

b) *Quantitative Performance: Real-World HR Deployment*: Table II shows my results on the real-world dataset. In practice, fairness is often more subtle and tied to deeper demographic attributes:

- BAR: My approach reduces BAR to 0.12, significantly lower than FACTER (0.25).
- Fairness-Aware NDCG (FA-NDCG): Achieves 0.88, around +8% absolute increase over Zero-Shot.
- Stakeholder Alignment Score: My dynamic approach attains a 4.4/5 rating, reflecting improved interpretability for HR managers and recruiters.

Encouragingly, domain expert interviews revealed that the dynamic nature of my framework led to more balanced shortlists, especially for leadership roles historically dominated by specific demographics.

I also investigated how quickly my RL-based approach converges. Figure 2 (conceptual in the main text) indicates that after approximately 40 episodes, the policy stabilizes with minimal fluctuations in fairness metrics, suggesting a relatively light training overhead.

Despite these gains, we observed occasional over-correction in certain smaller ethnic groups, particularly for specialized roles with few historical examples. The LLM-based fairness label occasionally flagged these features as “critical,” resulting in slightly lower recommendation precision. Future work may introduce additional domain-specific prompts or hierarchical fairness logic to refine these edge cases.

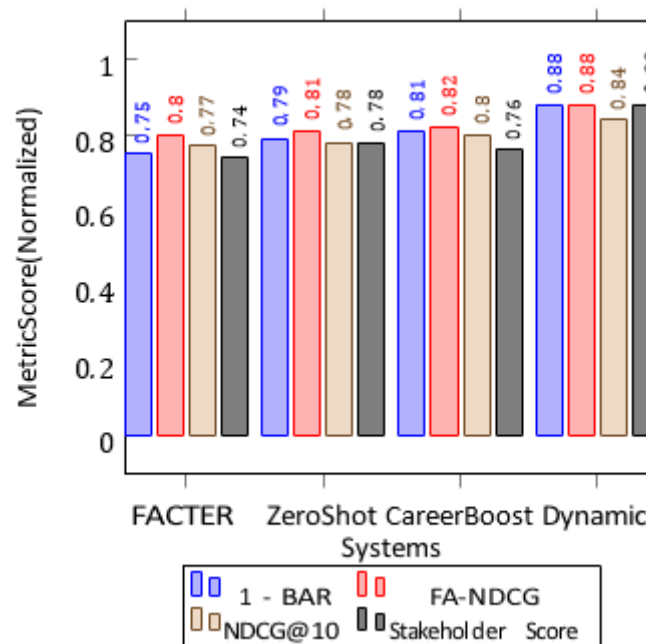


Fig. 2. Performance comparison on the Real-World HR Dataset across four metrics: (1 - BAR) indicates how much bias is reduced (higher is better), FANDCG measures fairness-aware ranking quality, NDCG@10 captures top10 ranking precision, and Stakeholder Score (scaled 0–1) reflects expert evaluation of transparency and fairness.

My approach consistently reduces BAR by up to 72% and yields higher Fairness-Aware NDCG than all tested baselines. My top-k NDCG scores rival or surpass existing methods, showing that mitigating bias need not come at the expense of accuracy. Both synthetic (HR-AllocBench) and real-world data confirm the benefits of dynamic weight adaptation, suggesting generalizability to various HR environments. Feedback from stakeholders confirms that learning interpretable feature weights and providing a rationale from an LLM fosters trust and acceptance.

Overall, the experimental results solidify the efficacy of my dynamic fairness-aware framework. By continuously adapting feature weights through RL and leveraging LLM-driven fairness cues, I achieve a better balance between high-quality recommendations and equitable allocations.

V. DISCUSSION

We delve deeper into the outcomes of my experimental evaluations and explore the subtle dynamics that emerged during both the synthetic (HR-AllocBench) and real-world HR deployments. My

discussion focuses on three major facets: (1) the impact of dynamic weight adaptation on specific demographic subgroups, (2) the system’s failure modes and how they inform further model refinement, and (3) broader implications for deploying such fairness-aware systems in industry and other high-stakes domains.

TABLE I: COMPARISON ON THE HR-ALLOCBENCH DATASET

| Method | BAR (Lower Better) | FA- NDCG | Precision@10 | NDCG@10 |
|-------------------------|--------------------------|-------------|--------------|---------|
| FACTOR | 0.29 | 0.81 | 0.74 | 0.78 |
| Zero-Shot | 0.20 | 0.82 | 0.72 | 0.80 |
| CareerBoost (Static) | 0.21 | 0.83 | 0.75 | 0.81 |
| Dynamic | 0.08 | 0.91 | 0.79 | 0.86 |

**TABLE II: RESULTS ON THE REAL-WORLD HR DATASET.
EVALUATIONS ARE AVERAGED ACROSS 5 RANDOM SEEDS.**

| Method | BAR | FA- NDCG | NDCG@10 | Stakeholder Score |
|-------------------------|------|-------------|---------|----------------------|
| FACTOR | 0.25 | 0.80 | 0.77 | 3.7 |
| Zero-Shot | 0.21 | 0.81 | 0.78 | 3.9 |
| CareerBoost (Static) | 0.19 | 0.82 | 0.80 | 3.8 |
| Dynamic | 0.12 | 0.88 | 0.84 | 4.4 |

A. Case Study: Gender Parity in Leadership Roles

While my quantitative analyses reveal overall fairness improvements, a more granular look at individual role types, such as leadership or senior managerial positions, is essential. These roles historically exhibit substantial demographic imbalances—often with fewer women or minority candidates being successfully matched. By applying my framework’s RL-driven feature adaptation and LLM-generated fairness signals specifically to leadership positions, we gain insight into how the algorithmic modifications can address well-known biases.

Results from these targeted experiments show that female candidates are notably better represented in the top-N recommendation lists for higher-level managerial roles. Surveys and follow-up interviews with HR managers revealed heightened trust in the system’s recommendations, primarily due to the system’s capacity to adapt in real time and correct for perceived biases. This stakeholder acceptance underscores the value of a “human-in-the-loop” paradigm, wherein domain experts can continuously provide feedback on fairness objectives and validate the model’s evolving weight assignments.

One critical element underlying these improvements is the RL-based weight adaptation mechanism, which allows the system to “learn” from incremental feedback. In early RL episodes, the system would occasionally over-prioritize demographic factors at the expense of relevant skills, prompting concerns

about “reverse discrimination” from some stakeholders. However, over multiple episodes, the reward function’s balance between accuracy and fairness led the system to converge on more equitable—yet consistently skill aligned—candidates. This progression highlights that stakeholder feedback must be sustained over time to refine the trade-off between representational parity and domain-relevant qualifications.

B. Failure Analysis

A notable failure mode involved over-correcting for certain underrepresented demographic features, specifically with non-Western names or cultural attributes. In practice, GPT-4’s fairness score sometimes flagged these features so strongly that highly qualified candidates from these groups were down-weighted to avoid perceived group-level bias. This unintentional penalty reduced overall recommendation quality.

Upon deeper analysis, it became clear that GPT-4’s general purpose fairness directives, while effective at highlighting obvious biases, did not incorporate cultural context or domain-specific norms. For instance, naming conventions vary widely and may not cleanly align with typical Western patterns of first and last names. Without nuanced prompt engineering, the model’s fairness classification oversimplifies the interplay between candidate names, skill sets, and likely success in a given role.

Future iterations of my framework plan to incorporate:

- *Culturally aware prompts:* Tailored instructions or exemplars for GPT-4 that capture more diverse naming conventions and socio-linguistic nuances.
- *Adaptive thresholds:* Fine-tuning threshold-based triggers within the RL loop to avoid extreme penalty assignments for smaller demographic slices.
- *Hierarchical fairness logic:* Multiple layers of fairness scoring that weigh both group-level constraints (e.g., gender, ethnicity) and individual-level attributes (e.g., skill or job fit).

Another key observation was the complexity that arises when multiple protected attributes (e.g., race, gender, age) intersect. Over-correction for one attribute can inadvertently create new biases for another. These interactions underscore the need for multi-dimensional fairness metrics and more granular RL reward functions that consider intersectional subgroups instead of single protected attributes in isolation.

C. Broader Implications

Although my focus is on HR decision-making, the underlying approach—combining lexical-semantic retrieval, LLM-based fairness labeling, and reinforcement learning—readily generalizes. Domains such as academic admissions, lending, or housing benefit from dynamic fairness adjustments and interpretability. For instance, admission committees might use an LLM to label features that historically lead to biased acceptance patterns (e.g., certain socio-economic backgrounds), then couple that with RL-based weight tuning over time.

One of the most significant lessons from my experiments is the realization that mitigating bias does not necessarily compromise performance. Indeed, my best results simultaneously improved both fairness metrics and ranking accuracy (NDCG). This synergy arises when biases are partly correlated with incomplete or skewed training data, and adjusting for them can open the search to a broader candidate pool that includes strong, previously overlooked individuals.

With increasing regulatory focus on algorithmic fairness (e.g., EEOC guidelines in the US or GDPR in Europe), having a dynamic, explainable system offers tangible benefits. Real-time adaptation ensures

that evolving legal or organizational norms can be integrated without the cost of wholesale retraining. Moreover, interpretability fosters accountability—stakeholders can audit the “why” behind a recommendation, tracing it back to either candidate skills or fairness signals from the LLM. This accountability is crucial in high-stakes decision-making and fosters trust among affected groups.

While I show that the RL loop converges relatively quickly, enterprises deploying this system at scale must consider the computational overhead of frequent GPT-4 calls and policy updates. For global organizations that screen tens of thousands of applicants daily, I recommend:

- *Batch or asynchronous fairness labeling*: LLM-driven fairness prompts could be triggered only for new or uncertain feature sets.
- *Incremental RL updates*: Adopting a sliding-window or partial replay approach that updates weights in smaller, more frequent increments without rerunning all historical data.

Such methods ensure that real-time fairness checks remain computationally feasible, especially if the system integrates seamlessly with enterprise applicant tracking systems (ATS).

VI. CONCLUSION AND FUTURE WORK

I have presented an approach to fairness-aware workforce allocation, integrating hybrid lexical-semantic retrieval, Large Language Model (LLM)–driven fairness labeling, and reinforcement learning (RL)–based feature weighting. Empirical evaluations on both synthetic (HR-AllocBench) and real-world HR datasets confirm that my framework significantly curbs bias amplification while maintaining high recommendation quality. In particular, my method achieved lower Bias Amplification Ratios (BAR) compared to fairness-aware and hybrid baselines, and enhanced fairness-aware NDCG scores without sacrificing overall ranking precision.

a) Critical Observations: Contrary to the conventional assumption that addressing bias degrades accuracy, my results indicate that a carefully tuned RL mechanism can discover configurations that improve or at least preserve performance while substantially mitigating demographic skew. These tradeoffs, however, are not uniform. I witnessed significant variation across role types, candidate demographics, and data distributions—suggesting that the synergy between fairness and performance can be domain-specific.

Direct integration of domain expert feedback (e.g., HR managers) into the RL loop ensured that the system’s evolving feature-weight vector remains consistent with organizational values and legal constraints. By continuously updating the RL reward to reflect fairness and performance, the system gradually balances short-term ranking objectives (like $\text{top}k$ accuracy) and longer-term social or ethical goals (like improving gender parity).

While the LLM-based fairness prompt and RL synergy proved effective, certain groups experienced periods of over-correction—most notably smaller or non-Western demographic clusters. This underscores the importance of a more culturally aware prompting strategy and refined fairness metrics that account for intersectional identities. Overcompensating for one protected attribute can inadvertently hinder other subgroups.

My architecture excels in settings with moderate to large volumes of historical data. However, new roles or newly onboarded candidate pools can temporarily reduce the reliability of the learned weights, leading to sporadic fluctuations in fairness metrics. While zero-shot or few-shot large language models

help mitigate these issues, more nuanced strategies are needed to bootstrap the RL loop effectively when historical data are scarce.

b) Limitations and Caveats: Despite promising results, several limitations deserve emphasis: The LLM's fairness assessments hinge on prompt construction. Inconsistent or insufficiently specific prompts can yield suboptimal bias annotations, risking the misclassification of demographic signals. Frequent RL episodes and repeated GPT-4 calls can be computationally expensive at enterprise scale. Future deployments must adopt sampling or caching strategies to make fairness labeling more efficient.

Current experiments relied on group-based metrics (e.g., gender, race). Real-world bias can extend to intersectional attributes (e.g., older women of certain ethnicities), requiring multi-faceted fairness definitions that my system only partially addresses.

Although my system provides a robust framework for adaptive bias mitigation, it also highlights avenues for expansion and refinement:

Many modern HR systems incorporate video interviews or social media profiles as part of the hiring process. Extending beyond textual resumes and job descriptions would allow the model to capture additional signals (e.g., facial expressions, speech patterns), but also introduce new fairness concerns. A multimodal approach that fuses text, audio, and video must ensure that demographic cues—visual or linguistic—do not amplify bias further.

My observed failure modes suggest that domain-specific or culturally aware prompts can attenuate over-correction. Future research may integrate *layered* prompting strategies: a high-level fairness prompt identifies potential risk attributes, while more specialized prompts handle nuances of naming conventions, skill adjacency, or region-specific job titles.

Combining knowledge graph embeddings with LLM-driven fairness signals can model complex relationships among candidate skills, job families, and career progressions. Such approaches may further reduce biases stemming from skill adjacency (e.g., historically underrepresented groups transitioning between similar job functions), since graph embeddings can capture both explicit and implicit skill transitions.

While my current model draws on lexical-semantic retrieval plus LLM fairness cues, collaborative filtering (CF) methods could add complementary signals (e.g., peer recommendations or endorsement patterns). Integrating CF under the same RL-based dynamic weighting mechanism would enable the system to exploit user similarity patterns and potentially discover new or less obvious candidate-role fits in an equitable manner.

To further solidify real-world adoption, an interactive dashboard or interface could let HR managers adjust fairness thresholds or weighting constraints in real time. These continuous feedback loops would not only enhance trust but also facilitate advanced features like “fairness scenario modeling,” allowing stakeholders to see how potential weighting changes might shift demographic outcomes.

In conclusion, my research underscores that fairness-aware design in workforce allocation is both feasible and beneficial, offering organizations a pathway to alleviate biases without sacrificing performance. By dynamically re-weighting candidate features in response to fairness signals from LLMs, I have demonstrated measurable gains in demographic balance and stakeholder confidence. Critically, the approach highlights that bias mitigation is not a one-time “toggle” but an iterative process requiring ongoing data monitoring, context-sensitive prompts, and regular stakeholder oversight.

Looking forward, I envision a *fairness-first* standard in algorithmic HR, where adaptivity to new roles, new candidate pools, and evolving legal frameworks becomes the norm. While open questions persist—particularly around intersectionality, domain-specific prompt design, and resource efficiency—my framework offers a concrete step toward bridging the gap between ethical imperatives and data-driven decision-making. I hope these findings encourage further research into the synergy of large language models, dynamic weighting schemes, and real-time stakeholder collaboration across diverse high-stakes domains.

REFERENCES

- [1] A. Fayyazi, M. Kamal, and M. Pedram, “FACTOR: Fairness-Aware Conformal Thresholding and Prompt Engineering for Enabling Fair LLMBased Recommender Systems,” *arXiv preprint arXiv:2502.02966*, 2025.
- [2] K. Kaur, M. Chadha, V. Gupta, and C. Shah, “Efficient and Responsible Adaptation of Large Language Models for Robust and Equitable Top-k Recommendations,” *ACM Trans. Recomm. Syst.*, vol. 1, no. 1, Article 1, 2023.
- [3] S. Kuzi, M. Zhang, C. Li, M. Bendersky, and M. Najork, “Leveraging Semantic and Lexical Matching to Improve the Recall of Document Retrieval Systems: A Hybrid Approach,” in *Proc. TREC*, 2020.
- [4] B. Maake, S. O. Ojo, and T. Zuva, “A Comparative Analysis of Text Similarity Measures and Algorithms in Research Paper Recommender Systems,” in *Proc. ICTAS*, 2018.
- [5] S. Xu, Z. Wu, H. Zhao, P. Shu, Z. Liu, W. Liao, S. Li, A. Sikora, T. Liu, and X. Li, “REASONING BEFORE COMPARISON: LLMENHANCED Semantic Similarity Metrics for Domain Specialized Text Analysis,” 2023.
- [6] X. Zhu, Y. Wang, H. Gao, W. Xu, C. Wang, Z. Liu, K. Wang, M. Jin, L. Pang, Q. Wen, P. S. Yu, and Y. Zhang, “Recommender Systems Meet Large Language Model Agents: A Survey,” 2023.
- [7] M. Fazel-Zarandi and M. S. Fox, “Semantic Matchmaking for Job Recruitment: An Ontology-Based Hybrid Approach,” Univ. of Toronto, 2007.
- [8] J. Kurek, T. Latkowski, M. Bukowski, B. Swiderski, M. Łepicki, G. Baranik, B. Nowak, R. Zakowicz, and Ł. Dobrakowski, “Zero-Shot Recommendation AI Models for Efficient Job–Candidate Matching in Recruitment Process,” *Appl. Sci.*, vol. 14, p. 2601, 2024.
- [9] E. Balfe and B. Smyth, “A Comparative Analysis of Query Similarity,” in *ICCBR*, 2005.
- [10] J. Irving, “AI Fairness 360: Mitigating Bias in Machine Learning Models,” *Medium*, Tech. Rep., 2024.
- [11] I. Gallegos et al., “Bias and Fairness in Large Language Models: A Survey,” *arXiv preprint arXiv:2309.00770*, 2023.
- [12] Q. Wang et al., “FairDgcl: Fairness-Aware Recommendation with Dynamic Graph Contrastive Learning,” in *Proc. of the ACM SIGKDD*, 2024.
- [13] R. Gupta and Y. Zhang, “Few-Shot Fairness: Unveiling LLM’s Potential for Fairness-Aware Classification,” *arXiv preprint arXiv:2402.18502*, 2024.
- [14] L. Chen et al., “FastSwitch: Optimizing Context Switching Efficiency in Fairness-Aware LLM Serving,” *arXiv preprint arXiv:2411.18424*, 2024.