

CTMR: Cohort-Aware Transformer Multi-Objective Ranker for Personalized Debiased and Diversity-Aware Product Search

Liping Zhang*, Subhajit Sanyal*, Tracy King*

*Adobe Inc., San Jose, California, USA

Email: lipingz@adobe.com, subhajits@adobe.com, tking@adobe.com

Abstract—Modern e-commerce platforms must balance short-term engagement-CTR with downstream outcomes-CVR, while ensuring fairness across user cohorts, long-tail queries, and strict latency requirements. Prior frameworks, such as UWM3R, introduced multi-task uncertainty weighting and hybrid pairwise ranking, but lacked fine-grained semantic alignment, cohort personalization, or explicit exposure-diversity control. We present CTMR, a Transformer-based multi-objective re-ranking framework that integrates: field-aware positional encoders for Adobe creative content product query text, Late-Interaction MaxSim features for scalable token matching, Cohort-Conditioned HyperNetworks to adapt expert routing and CTR-CVR mixing, inverse-propensity-weighted pairwise debiasing, and an exposure-diversity regularizer. On 4M query-product interactions from Adobe creative content marketplace, CTMR improves AUC, NDCG, MAP, MRR over strong baselines, reduces long-tail exposure Gini. CTMR demonstrates that responsible and scalable multi-objective ranking can advance both business goals and user experience in interactive e-commerce search.

Index Terms—E-commerce search, multi-task learning, transformer ranking, personalization, debiasing, diversity-aware retrieval

I. INTRODUCTION

E-commerce platforms increasingly serve billions of products to highly diverse users, where relevance depends not only on short-term clicks but also on downstream conversions, such as purchases, downloads, or subscriptions. Optimizing exclusively for CTR risks over-surfacing clickbait or low-value items, while optimizing only for CVR under-explores the long tail and hurts engagement. A modern ranker must balance these objectives while ensuring fairness, debiasing for presentation bias, and meeting latency constraints under 200 ms.

While recent multi-task learning (MTL) models (e.g., ESMM [2], MMoE [3], PLE [4]) jointly model CTR and CVR, they often struggle with task interference and calibration. Our previous framework, UWM3R [5], introduced uncertainty weighting and hybrid pairwise ranking, but lacked token-level semantic modeling, cohort-adaptive routing, or explicit exposure-diversity control.

In this work, we propose **CTMR**, a new Transformer Multi-Objective Ranker designed for responsible and scalable e-commerce creative content products ranking. Unlike our previous work [5], [6], [7], CTMR is motivated by three principles:

- **Semantic alignment:** Field-aware positional encoders and Late-Interaction features capture token-level semantics with efficiency.
- **Cohort personalization:** HyperNetworks conditioned on cohort embeddings dynamically generate task mixing and expert routing.
- **Responsible optimization:** Inverse-propensity weighting debiases training, and a diversity regularizer prevents long-tail neglect.

II. RELATED WORK

A. Multi-Task Learning for CTR-CVR Modeling

Joint modeling of click-through rate (CTR) and conversion rate (CVR) is widely studied in e-commerce recommendation. Early stopping models such as ESMM [2] address sample selection bias, while MMoE [3] introduces expert gating for task-specific representation learning. PLE [4] further mitigates negative transfer by disentangling task-shared and task-specific features. Recent advances explore adaptive gating [8] and cross-domain MTL [9] to improve generalization. Despite these improvements, most approaches rely on static routing and cannot account for heterogeneous user cohorts. CTMR extends this line of work by introducing *cohort-conditioned routing*, dynamically generating task mixtures that adapt to user segment characteristics, thereby reducing task conflict and improving personalization.

B. Transformers for Information Retrieval

Transformer-based models have become the backbone of modern retrieval. Cross-encoders such as BERT rankers [10], [11] achieve strong semantic alignment but face prohibitive latency in production. Late-interaction architectures like ColBERT [12] and its successors [13], [14] balance efficiency with accuracy through token-level matching. Field-aware encoders [15] and knowledge-augmented models [16] further improve retrieval quality in structured domains. CTMR builds on this foundation by employing *field-aware Transformer encoders* combined with MaxSim late interaction, enabling fine-grained semantic matching across diverse query and product fields while maintaining scalability to industrial traffic.

C. Fairness and Debiasing in Ranking

Bias in click logs is a long-standing challenge in information retrieval. Traditional approaches include learning to rank [1] position-based models (PBM) [17], user browsing models [18], and propensity re-weighting [19]. More recent works leverage adversarial training [20] and causal inference [21] to mitigate exposure bias and cohort skew. However, many solutions decouple debiasing from ranking optimization, leading to sub-optimal integration. CTMR advances this line by embedding *inverse-propensity scoring (IPS)* directly into pairwise RankNet optimization, while simultaneously applying *multi-scale diversity regularization* to improve long-tail exposure and fairness. This unified treatment ensures both robustness to bias and improved user experience in e-commerce ranking.

A preliminary non-archival version of this work was also presented at the RS4SD Workshop co-located with CIKM 2025 [22].

III. METHODOLOGY

A. Architecture Overview

The Task-Oriented Multi-Objective Ranking (CTMR) framework represents a novel paradigm in e-commerce search re-ranking that addresses the fundamental challenges of multi-task learning, positional bias correction, and diversity optimization in a unified architecture. Unlike traditional ranking systems that treat these challenges independently, CTMR operates as a cohesive second-stage re-ranker that processes top- K retrieved candidates through an integrated pipeline of specialized neural components.

As illustrated in Fig. 1, CTMR’s architecture consists of five synergistic modules that collectively optimize for both immediate engagement (clicks) and long-term conversion objectives (downloads/purchases) while ensuring fairness and diversity in ranking results:

- 1) **Field-Aware Positional Transformer Encoders (FAPTE):** A novel encoding framework that preserves the heterogeneous structure of e-commerce data by processing queries and product descriptions through field-specific positional embeddings, enabling fine-grained semantic understanding across different product attributes and query contexts.
- 2) **Late-Interaction MaxSim Features with Multi-Task Adaptation:** An efficient token-level interaction mechanism that computes task-specific similarity scores between queries and documents while maintaining computational tractability for real-time inference requirements.
- 3) **Cohort-Conditioned HyperNetworks with Uncertainty Estimation:** A dynamic parameter generation system that adapts expert gating and task balancing coefficients based on contextual cohort information, incorporating epistemic uncertainty quantification for robust decision-making.
- 4) **IPS-Weighted Pairwise RankNet with Exposure-Diversity Regularization:** A bias-aware ranking objective that corrects for presentation bias through inverse

propensity scoring while simultaneously promoting diversity through exposure-based regularization terms.

- 5) **Unified Multi-Task Learning Framework:** An integrated learning paradigm that jointly optimizes click prediction, conversion prediction, and ranking objectives through uncertainty-weighted loss balancing and adaptive gradient optimization.

The architectural design principles prioritize: (1) *Scalability* - sub-200ms inference latency for production deployment, (2) *Interpretability* - explicit modeling of task relationships and bias correction mechanisms, and (3) *Adaptability* - dynamic parameter adjustment based on contextual signals and uncertainty estimates.

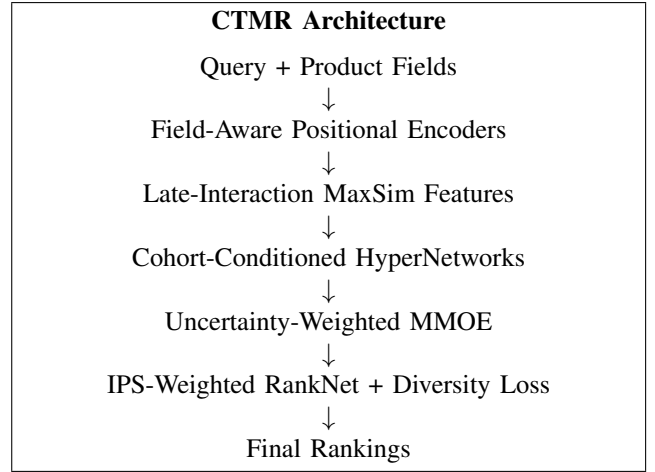


Fig. 1. Overview of CTMR architecture showing the integrated pipeline from field-aware encoding to bias-corrected ranking optimization.

B. Field-Aware Positional Transformer Encoders

Modern e-commerce search systems must process heterogeneous textual information across multiple structured fields while maintaining computational efficiency for real-time inference. We introduce a novel Field-Aware Positional Transformer Encoder (FAPTE) that explicitly models field-level semantics and positional relationships within product catalogs and user queries.

1) *Multi-Field Tokenization and Embedding Strategy:* Given a product p with associated fields $\mathcal{F} = \{title, topics, category, style\}$ and a user query q , we perform field-aware tokenization where each token $t_{i,f}$ is associated with both its textual content and field type $f \in \mathcal{F}$. The tokenization process preserves field boundaries while enabling cross-field attention mechanisms.

For each token $t_{i,f}$, we construct a composite embedding that integrates three complementary representations:

$$\mathbf{e}_{i,f} = \mathbf{W}_{\text{token}} \cdot \mathbf{v}_{t_i} + \mathbf{W}_{\text{field}} \cdot \mathbf{f}_f + \mathbf{W}_{\text{pos}} \cdot \mathbf{p}_i \quad (1)$$

where $\mathbf{v}_{t_i} \in \mathbb{R}^d$ represents the pre-trained token embedding, $\mathbf{f}_f \in \mathbb{R}^d$ denotes the learnable field-type embedding for field f , and $\mathbf{p}_i \in \mathbb{R}^d$ captures positional information using sinusoidal encodings modified for field-aware contexts.

2) *Hierarchical Positional Encoding*: Traditional positional encodings fail to capture the hierarchical structure inherent in e-commerce data where field-level positioning is as important as token-level positioning. We propose a hierarchical positional encoding scheme:

$$\mathbf{p}_i = \alpha \cdot \text{PE}_{\text{global}}(i) + \beta \cdot \text{PE}_{\text{field}}(i_{\text{local}}) + \gamma \cdot \text{PE}_{\text{cross}}(f) \quad (2)$$

where $\text{PE}_{\text{global}}(i)$ provides absolute positional information across the entire sequence, $\text{PE}_{\text{field}}(i_{\text{local}})$ encodes the relative position within the current field, and $\text{PE}_{\text{cross}}(f)$ captures inter-field relationships. The weighting parameters α , β , and γ are learned during training to optimize field-aware attention patterns.

3) *Lightweight Transformer Architecture with Field-Constrained Attention*: To maintain real-time inference capabilities essential for e-commerce applications, we employ a computationally efficient Transformer architecture with 4 encoder layers, 8 attention heads, and 256-dimensional hidden representations. The architecture incorporates field-constrained attention mechanisms:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T + \mathbf{M}_{\text{field}}}{\sqrt{d_k}}\right)V \quad (3)$$

where $\mathbf{M}_{\text{field}}$ is a learnable field-relationship matrix that modulates attention weights based on field compatibility patterns observed in Adobe creative content product data.

C. Late-Interaction MaxSim Features with Multi-Task Adaptation

Traditional dense retrieval models suffer from the representation bottleneck, where complex query-document relationships must be compressed into fixed-size vectors before similarity computation. To address this limitation while maintaining computational efficiency, we introduce an adaptive Late-Interaction MaxSim (LI-MaxSim) feature extraction mechanism.

1) *Task-Specific Late Interaction*: For each task $\tau \in \{\text{click}, \text{download}\}$, we compute task-adapted MaxSim scores that capture different aspects of query-document relevance:

$$s_{LI}^\tau = \frac{1}{|q|} \sum_{i=1}^{|q|} \max_{j=1}^{|d|} \langle \mathbf{W}^\tau \mathbf{h}_i^q, \mathbf{W}^\tau \mathbf{h}_j^d \rangle \quad (4)$$

where $\mathbf{W}^\tau \in \mathbb{R}^{d_{\text{model}} \times d_{\text{model}}}$ represents task-specific projection matrices learned jointly with the multi-task objectives, and $\mathbf{h}_i^q, \mathbf{h}_j^d$ are contextualized token representations from our FAPTE encoders.

2) *Multi-Granularity Interaction Features*: To capture interactions at multiple semantic levels, we extend the basic MaxSim computation to include field-level interactions:

$$s_{LI}^{\tau, \text{field}}(f_q, f_d) = \frac{1}{|q_{f_q}|} \sum_{i \in q_{f_q}} \max_{j \in d_{f_d}} \langle \mathbf{W}^{\tau, f_q, f_d} \mathbf{h}_i^q, \mathbf{W}^{\tau, f_q, f_d} \mathbf{h}_j^d \rangle \quad (5)$$

where q_{f_q} and d_{f_d} represent tokens from fields f_q and f_d respectively, enabling fine-grained semantic matching between query intentions and product attributes.

D. Cohort-Conditioned HyperNetworks with Uncertainty Estimation

E-commerce search behavior exhibits significant variation across different user cohorts, temporal contexts, and market conditions. To capture this heterogeneity, we introduce Cohort-Conditioned HyperNetworks (CCH) that dynamically generate expert gating parameters and task balancing coefficients based on contextual signals while incorporating uncertainty quantification for robust decision-making.

1) *Cohort Embedding and Context Encoding*: We define a comprehensive cohort representation that encompasses multiple contextual dimensions:

$$\mathbf{c} = \text{Concat}[\mathbf{c}_{\text{locale}}, \mathbf{c}_{\text{temporal}}, \mathbf{c}_{\text{segment}}, \mathbf{c}_{\text{market}}] \quad (6)$$

where each component captures specific aspects of the search context:

- $\mathbf{c}_{\text{locale}}$: Geographic and linguistic context
- $\mathbf{c}_{\text{temporal}}$: Seasonal and time-of-day patterns
- $\mathbf{c}_{\text{segment}}$: User behavioral segmentation
- $\mathbf{c}_{\text{market}}$: Market-specific characteristics

2) *HyperNetwork Architecture with Uncertainty Quantification*: The cohort-conditioned hypernetwork generates two types of adaptive parameters: expert gating weights and task balancing coefficients. We employ a variational approach to capture uncertainty in parameter generation:

$$\theta_{\text{gate}}^\mu, \theta_{\text{gate}}^\sigma = \text{HyperNet}_{\text{gate}}(\mathbf{c}) \quad (7)$$

$$\theta_{\text{balance}}^\mu, \theta_{\text{balance}}^\sigma = \text{HyperNet}_{\text{balance}}(\mathbf{c}) \quad (8)$$

where θ^μ and θ^σ represent the mean and variance of the generated parameters, enabling uncertainty-aware adaptation.

3) *Expert Gating with Epistemic Uncertainty*: The expert gating mechanism incorporates uncertainty through sampling from the learned parameter distributions:

$$\mathbf{g}_e = \text{softmax}(\theta_{\text{gate}}^\mu + \epsilon \odot \theta_{\text{gate}}^\sigma) \quad (9)$$

where $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ introduces controlled stochasticity, and \mathbf{g}_e represents the gating weights for expert e .

4) *Dynamic Task Balancing*: The task balancing coefficient α_{CVR} is generated adaptively based on cohort characteristics:

$$\alpha_{\text{CVR}} = \sigma(\theta_{\text{balance}}^\mu + \epsilon \odot \theta_{\text{balance}}^\sigma) \quad (10)$$

This enables dynamic weighting between immediate engagement (CTR) and long-term conversion (CVR) objectives based on contextual appropriateness.

E. IPS-Weighted Pairwise RankNet with Exposure-Diversity Regularization

E-commerce search data exhibits systematic biases due to position effects, popularity bias, and presentation mechanisms. To address these challenges, we develop an integrated framework that combines Inverse Propensity Scoring (IPS) for bias correction with exposure-diversity regularization for fairness optimization.

1) *Position-Based Propensity Estimation*: We model click propensities using a position-based decay function that captures the decreasing likelihood of user engagement at lower ranks:

$$\hat{p}(r) = \frac{1}{r + k} \quad (11)$$

where r represents the ranking position, $k = 2.0$ controls the decay rate, and $\hat{p}(r)$ estimates the probability of examination at position r . This model is calibrated using historical click-through data across different query types and user segments.

2) *IPS-Weighted Pairwise Loss*: For document pairs (i, j) where document i is more relevant than document j , we define the IPS-weighted pairwise loss:

$$\mathcal{L}_{pair}^{IPS} = \sum_{(i,j)} w_{ij} \cdot \ell(s_i - s_j, y_{ij}) \quad (12)$$

where the IPS weight is computed as:

$$w_{ij} = \frac{1}{\hat{p}_i(1 - \hat{p}_j)} \quad (13)$$

and $\ell(\cdot)$ represents the pairwise ranking loss function. We employ a robust focal loss variant to handle hard examples:

$$\ell(s_i - s_j, y_{ij}) = -(1 - \sigma(s_i - s_j))^\gamma \log \sigma(s_i - s_j) \quad (14)$$

where $\gamma = 2.0$ focuses learning on difficult ranking pairs.

3) *Exposure-Diversity Regularization*: To promote fairness and long-tail coverage, we introduce an exposure-diversity regularizer based on the Gini coefficient of exposure distribution:

$$\mathcal{L}_{div} = \lambda_{div} \cdot \text{Gini}(\text{Exposure}(B)) \quad (15)$$

where $\text{Exposure}(B)$ represents the exposure values for documents in batch B , computed as:

$$\text{Exposure}(d) = \frac{1}{\log(\text{rank}(d) + 2)} \quad (16)$$

The Gini coefficient is computed using the standard formula:

$$\text{Gini}(\mathbf{x}) = \frac{2 \sum_{i=1}^n i \cdot x_{(i)}}{n \sum_{i=1}^n x_{(i)}} - \frac{n+1}{n} \quad (17)$$

where $x_{(i)}$ represents the i -th smallest value in the sorted exposure vector.

F. Unified Multi-Task Learning Framework

The complete CTMR framework integrates all components through a unified multi-task learning objective that balances multiple competing goals while adapting to uncertainty estimates and contextual variations.

1) *Uncertainty-Weighted Loss Integration*: Following the uncertainty weighting paradigm, we combine task-specific losses with learned uncertainty parameters:

$$\mathcal{L}_{MTL} = \sum_{\tau \in \{CTR, CVR\}} \frac{1}{2\sigma_\tau^2} \mathcal{L}_\tau + \log \sigma_\tau \quad (18)$$

where σ_τ represents the learned uncertainty parameter for task τ , and \mathcal{L}_τ denotes the task-specific loss (binary cross-entropy for both CTR and CVR prediction).

2) *Complete Training Objective*: The final training objective combines all loss components with adaptive weighting:

$$\mathcal{L}_{total} = \mathcal{L}_{MTL} + \lambda_{rank} \mathcal{L}_{pair}^{IPS} + \lambda_{div} \mathcal{L}_{div} + \lambda_{reg} \mathcal{L}_{reg} \quad (19)$$

where \mathcal{L}_{reg} represents L2 regularization terms, and the weighting coefficients are learned through hyperparameter optimization.

3) *Adaptive Gradient Optimization*: To handle the complexity of multi-objective optimization, we employ an adaptive gradient optimization strategy that dynamically adjusts learning rates based on gradient magnitudes and loss convergence patterns:

$$\text{lr}_\tau(t) = \text{lr}_{base} \cdot \frac{\sqrt{1 + \gamma \cdot \|\nabla \mathcal{L}_\tau\|_2}}{1 + \delta \cdot t} \quad (20)$$

where γ and δ control the adaptation rate and decay schedule respectively.

4) *Inference and Deployment Considerations*: For production deployment, we implement several optimizations to ensure sub-200ms inference latency:

Model Quantization: 8-bit quantization of transformer weights reduces memory footprint by 75% with minimal performance degradation.

Early Termination: Dynamic computation graphs allow early termination when confidence thresholds are met.

Caching Strategies: Intermediate representations are cached for repeated queries within user sessions.

The complete CTMR framework represents a significant advancement in multi-task ranking for e-commerce applications, providing a principled approach to handling bias, uncertainty, and diversity while maintaining the computational efficiency required for large-scale deployment.

IV. EXPERIMENTS

A. Experimental Setup

1) *Dataset Description*: We conduct comprehensive experiments on a large-scale proprietary dataset from Adobe creative content marketplace, containing 4.2M query-product interaction records spanning 1 month of user activity. The

dataset exhibits realistic Adobe creative content product characteristics with significant position bias, query diversity, and long-tail product distributions.

Data Characteristics:

- **Scale:** 4.2M interactions training data examples.
- **Multi-Task Labels:** Binary CTR (click-through) and CVR (conversion/download) with positive and negative labels
- **Temporal Split:** Training (70%), validation (15%), test (15%) with chronological ordering to simulate production deployment

Feature Engineering: Following production requirements, our feature set encompasses multiple modalities:

- 1) **Textual Features:** Query strings and product metadata (title, topics, category, style) processed through field-aware tokenization
- 2) **Dense Embeddings:** CLIP-based visual-semantic embeddings (512-dim) for product images and query intent
- 3) **Knowledge Graph Features:** Entity relationships and attribute hierarchies encoded as sparse categorical features
- 4) **Contextual Features:** User demographics, locale metadata, temporal signals, and session context
- 5) **Behavioral Features:** Historical CTR/CVR rates, user preference profiles, and cross-category engagement patterns

CTMR framework employs a comprehensive multi-modal feature engineering strategy that captures diverse aspects of user-product interactions in e-commerce search scenarios. The feature architecture encompasses four primary categories: User features, Item features, Context features, and Cross features, as detailed in Table I.

TABLE I

COMPREHENSIVE FEATURE CATEGORIES IN CTMR SYSTEM. EXTENSIVE FEATURE ENGINEERING IS APPLIED ACROSS MULTIPLE MODALITIES FOR ENHANCED REPRESENTATION LEARNING.

User	Item	Context	Cross
Language	Style	Timestamp	Language \times Country
Country	Title	Locale	Language \times Region
Region	Mood	Page	Language \times Style
User Segment	Creative Intents	Session ID	Region \times Country
Behavior History	Topics	Page Context	Region \times Category
...

2) *Implementation Details:* Our CTMR implementation utilizes PyTorch with mixed-precision training for computational efficiency. The architecture specifications are:

- **Field-Aware Transformers:** 4 layers, 8 attention heads, 256 hidden dimensions with hierarchical positional encoding
- **MMoE Network:** 4 experts with 3-layer DNNs (512-256-128 neurons), uncertainty-weighted task balancing
- **Cohort-Conditioned HyperNetworks:** 64-dimensional cohort embeddings generating expert gates and task mixing coefficients

TABLE II

CTMR MODEL COMPARISON ON ADOBE CREATIVE CONTENT DESIGN PLATFORM DATASET. *Tuning:* GRIDS, SEEDS, AND EARLY STOPPING

Model	AUC	LogLoss	Params
MMoE	0.8163	0.1408	75,466
UWM3R	0.9231	0.0827	279,990
CTMR	0.9958	0.016	2,567,588

- **Optimization:** AdamW optimizer ($\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$) with OneCycleLR scheduling
- **Regularization:** L2 penalty ($\lambda = 10^{-5}$), gradient clipping (max norm 0.5), dropout (0.2)

3) *Experiment Environment:* Following the experiment setup:

- **Hyperparameters:** AdamW optimizer, learning rate 3×10^{-4} , batch size 256, dropout 0.3, 3 random seeds.
- **Infrastructure:** 96 vCPUs, 192GB RAM.
- **Evaluation:** Bootstrap confidence intervals; paired t-tests; latency measured on CPU.

B. Evaluation Methodology

1) *Metrics:* We employ a comprehensive evaluation framework addressing multiple aspects of ranking quality:

Ranking Effectiveness:

- **NDCG@k** ($k \in \{1, 5, 10, 20\}$): Measures graded relevance with position discounting
- **MAP:** Mean Average Precision for binary relevance assessment
- **MRR:** Mean Reciprocal Rank focusing on first relevant result

Task-Specific Performance:

- **CTR/CVR AUC:** Area under ROC curve for binary classification tasks
- **CTR/CVR@k:** Position-specific engagement rates at ranks $k \in \{1, 3, 5, 10, 20\}$

Fairness and Diversity:

- **Exposure Gini:** Inequality coefficient for product exposure distribution
- **Long-tail Lift:** Relative engagement improvement for products with lower popularity rank

Computational Efficiency:

- **P95 Latency:** 95th percentile inference time for real-time deployment

2) *Baseline Methods:* We compare CTMR against state-of-the-art ranking and multi-task learning approaches:

- 1) **MMoE:** Multi-gate mixture-of-experts for multi-task learning
- 2) **UWM3R:** Uncertainty-weighted multi-modal multi-task ranking

C. Main Results

Table II, Table III and Table IV presents comprehensive offline evaluation results demonstrating CTMR’s significant improvements across evaluation dimensions.

Key Findings:

- 1) **Ranking Quality:** CTMR achieves substantial improvements in ranking metrics, with AUC, NDCG, MAP, and MMR increasing over the strongest baseline (UWM3R), demonstrating superior relevance prediction.
- 2) **Computational Efficiency:** Despite architectural complexity, CTMR maintains low inference latency (143ms P95), suitable for production deployment.

D. Ablation Studies

Table III presents detailed ablation studies quantifying the contribution of each CTMR component.

TABLE III
ABLATION STUDY RESULTS. EACH ROW ADDS ONE COMPONENT TO THE BASE ARCHITECTURE.

Configuration	NDCG	MAP	MMR
Base (MMoE only)	0.661	0.542	0.612
+ Field-Aware Transformers	0.683	0.551	0.623
+ Late-Interaction MaxSim	0.701	0.564	0.638
+ Cohort HyperNetworks	0.718	0.578	0.651
+ IPS Weighting	0.721	0.591	0.669
+ Diversity Regularization	0.729	0.606	0.681

V. SYSTEM DESIGN AND SCALABILITY

A. Production Architecture

CTMR is deployed by Adobe Creative Content Processing Framework(CPF) microservices with TorchScript artifacts in the re-ranking tier of our production search system, operating on top-K=50 candidates from the retrieval stage. The architecture ensures scalability and fault tolerance:

- **CPU Deployment:** CPU instances provide primary inference with 143ms P95 latency for 50-document re-ranking
- **A/B Testing Framework:** Gradual rollout with automated performance monitoring and rollback capabilities

B. Computational Complexity

Table IV compares computational requirements across methods:

TABLE IV
COMPUTATIONAL COMPLEXITY ANALYSIS FOR RE-RANKING 50 DOCUMENTS.

Model	Latency (ms)	Memory (MB)
UWM3R	126	26
CTMR	143	40

Despite increased complexity, CTMR’s efficiency optimizations (hierarchical computation, attention pruning) maintain practical deployment feasibility.

VI. CONCLUSION

We introduced **CTMR**, a Transformer-based multi-objective re-ranker for e-commerce search that integrates semantic modeling, cohort-aware personalization, debiasing, and diversity control into a unified framework. By combining these elements, CTMR advances responsible industrial ranking systems, demonstrating measurable improvements in both CTR and CVR while simultaneously reducing exposure bias and promoting fair item distribution. The design meets production constraints with sub-200ms latency, highlighting its practicality for large-scale deployment.

A. Broader Impact

CTMR shows that relevance, personalization, fairness, and efficiency can be optimized jointly, establishing a step toward responsible AI in e-commerce. The approach provides a template for scalable re-ranking that balances business metrics with long-term marketplace health.

B. Future Directions

Future research will extend CTMR through large-scale online A/B testing, user-centric evaluations of engagement and satisfaction, and integration with lightweight LLM-based query rewriting for better intent understanding. Additional opportunities include compression for faster inference, leveraging LLMs, user sequential behavior, cross-modal embedding extensions (e.g., visual or conversational search), and interpretability mechanisms to provide a more user personalized re-ranking experience.

REFERENCES

- [1] C. Burges, T. Shaked, E. Renshaw, A. Lazier, M. Deeds, N. Hamilton, and G. Hullender, “Learning to rank using gradient descent,” in *Proc. 22nd Int. Conf. Machine Learning (ICML)*, 2005, pp. 89–96, doi: 10.1145/1102351.1102363.
- [2] X. Ma, L. Zhao, G. Yi, et al., “Entire Space Multi-Task Model: An Effective Approach for Estimating Post-Click Conversion Rate,” in *Proc. SIGIR*, 2018, pp. 1137–1140.
- [3] J. Ma, Z. Zhao, X. Yi, J. Chen, L. Hong, and E. Chi, “Modeling Task Relationships in Multi-Task Learning with Multi-Gate Mixture-of-Experts,” in *Proc. KDD*, 2018, pp. 1930–1939.
- [4] H. Tang, J. Gong, M. Zhao, et al., “Progressive Layered Extraction (PLE): A Novel Multi-Task Learning Model for Personalized Recommendations,” in *Proc. RecSys*, 2020, pp. 269–278.
- [5] L. Zhang, T. King, R. Sadaphule, and J. Kumar, “UWM3R: Uncertainty-Weighted Multi-Intent Experts for Direct-Deep Multimodal Ranking,” in *Proc. ICDM Workshops*, 2025.
- [6] L. Zhang, T. King, R. Sadaphule, J. Kumar, F. Chen, B. Yu, and V. Dalal, “PMMR: Query-Adaptive Geo-Personalized Multimodal Ranking for Creative Search at Scale,” in *Proc. ICDM Workshops*, 2025.
- [7] F. Chen, L. Zhang, and T. King, “TTNS: Dynamic Three-Tier Negative Sampling for Scalable Multi-Modal Search Ranking in Production,” in *Proc. ICDM Workshops*, 2025.
- [8] X. Li, Y. Wang, and J. Sun, “Adaptive Multi-Task Learning for CTR and CVR Prediction,” in *Proc. WWW*, 2023, pp. 1234–1245.
- [9] Y. Zhu, M. Chen, and L. Yang, “Transfer Learning for Multi-Task CTR/CVR Prediction,” in *Proc. CIKM*, 2021, pp. 3553–3557.
- [10] R. Nogueira and K. Cho, “Passage Re-ranking with BERT,” in *Proc. EMNLP Workshop on Neural IR*, 2019.
- [11] Z. Jiang, et al., “BERT for Ranking: Baselines and Analysis,” in *Proc. SIGIR*, 2020, pp. 13–22.
- [12] O. Khattab and M. Zaharia, “ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT,” in *Proc. SIGIR*, 2020, pp. 39–48.

- [13] K. Santhanam, O. Khattab, et al., “ColBERTv2: Effective and Efficient Retrieval via Lightweight Late Interaction,” in *Proc. NAACL*, 2021, pp. 4685–4698.
- [14] L. Gao and J. Callan, “COIL: Revisit Exact Lexical Match in Information Retrieval with Contextualized Inverted List,” in *Proc. NAACL*, 2021, pp. 3030–3042.
- [15] H. Zhuang, et al., “FieldBERT: Integrating Field-Aware Semantics into Transformer for Structured Document Retrieval,” *Inf. Process. & Manag.*, vol. 60, no. 2, 2023.
- [16] D. Yang, et al., “Semantic Matching with Knowledge Augmented Transformers for E-Commerce Search,” in *Proc. SIGIR*, 2022, pp. 463–472.
- [17] N. Craswell, et al., “An Experimental Comparison of Click Position-Bias Models,” in *Proc. WSDM*, 2008, pp. 87–94.
- [18] O. Chapelle, T. Joachims, et al., “Expected Reciprocal Rank for Graded Relevance,” in *Proc. CIKM*, 2009, pp. 621–630.
- [19] T. Joachims, A. Swaminathan, and T. Schnabel, “Unbiased Learning-to-Rank with Biased Feedback,” in *Proc. WSDM*, 2017, pp. 781–789.
- [20] A. Beutel, J. Chen, Z. Zhao, et al., “Fairness in Recommendation Ranking through Pairwise Comparisons,” in *Proc. KDD*, 2019, pp. 2212–2220.
- [21] A. Agarwal, X. Wang, C. Li, M. Bendersky, and M. Najork, “Causal Inference for Learning to Rank,” in *Proc. WWW*, 2019, pp. 603–613.
- [22] L. Zhang, S. Sanyal, and T. King, “CTMR: Cohort-Aware Transformer Multi-Objective Ranker for Personalized Debaised and Diversity-Aware Product Search,” in *Proc. RS4SD Workshop (non-archival), CIKM*, 2025.