# Enhancing E-commerce Search: Reinforcement Learning for Improved Language Model Fine-tuning

The landscape of e-commerce search is rapidly evolving, with advanced language models increasingly taking center stage in understanding user queries and retrieving relevant products. Recent research demonstrates that reinforcement learning (RL) techniques can significantly enhance language model fine-tuning for e-commerce search, addressing longstanding challenges like the semantic gap between brief user queries and product descriptions, handling long-tail queries, and improving product coverage. This comprehensive report examines how reinforcement learning approaches are transforming language model fine-tuning to deliver more accurate, diverse, and personalized search results in e-commerce platforms.

## The Evolution of E-commerce Search Systems

Traditional e-commerce search relies on structured search approaches that often struggle with understanding natural language queries and user intent. Most modern search applications, ad platforms, and recommender systems share a similar multitier information retrieval architecture with a candidate selection (retrieval) phase and a candidate ordering (ranking) phase. The retrieval phase reduces the possible candidates from millions to hundreds, while the ranking phase fine-tunes the ordering of candidates presented to customers[1]. While effective, this approach faces significant limitations when the best candidate items aren't included in the initial retrieval set.

E-commerce platforms particularly struggle with the complexity of generating detailed item titles from brief queries, the presence of noise in item titles with weak language ordering, issues with long-tail queries, and result interpretability[2]. These challenges create a substantial opportunity for reinforcement learning to enhance both the retrieval and ranking components of search systems. Unlike traditional methods that use static retrieval models, RL approaches can dynamically adapt to user behavior and feedback, creating a more responsive and accurate search experience.

## Understanding Reinforcement Learning in Language Model Fine-tuning

Reinforcement learning represents a paradigm shift in how language models are trained and optimized for specific tasks. Unlike supervised learning approaches that require extensive labeled datasets, RL enables models to learn through interactions with an environment, guided by reward signals that indicate desirable outcomes. When applied to language models, the model itself becomes the policy, and its action is the generation of the next token. A reward

model provides feedback on generated outputs, encouraging those that satisfy user preferences and penalizing those that don't[3].

## Reinforcement Learning from Human Feedback (RLHF)

The RLHF training process involves several key phases. First, a language model undergoes pretraining on a large corpus of text data from the internet. This establishes the foundational understanding of language patterns, syntax, and semantics. Next, a reward model is developed through training on human preference data, essentially teaching the model what constitutes a "good" response according to human evaluators. Finally, the model is fine-tuned using reinforcement learning techniques like Proximal Policy Optimization (PPO), which help maintain a balance between exploration of new possibilities and exploitation of known effective strategies[4].

The reward model plays a crucial role in this process, assigning numerical values to gauge the quality of responses. These values guide the reinforcement learning algorithm in adjusting the language model's weights to maximize future rewards. Techniques like Kullback-Leibler (KL) divergence are employed to measure differences between probability distributions, ensuring the model doesn't deviate too far from its original behavior while still improving based on feedback[4].

## Reinforcement Fine-Tuning (ReFT/RFT)

More recently, OpenAI and other organizations have introduced Reinforcement Fine-Tuning (ReFT or RFT), a specialized approach to using reinforcement learning for model customization. ReFT enables developers to customize models using dozens to thousands of high-quality tasks and grade responses with provided reference answers. This technique reinforces how the model reasons through similar problems and improves its accuracy on specific tasks in that domain[5].

What makes ReFT particularly valuable is its efficiency. According to OpenAI, ReFT can effectively fine-tune a model with just a few dozen examples, making it ideal for domains where data is scarce and costly, such as the medical sector[6]. This represents a significant advancement over traditional fine-tuning methods that typically require much larger datasets.

## Query Reformulation Through Reinforcement Learning

One of the most promising applications of reinforcement learning in e-commerce search is query reformulation (QR), which addresses the lexical gap between user queries and product descriptions to enhance search performance[7]. QR is particularly valuable in e-commerce, where users often provide brief, ambiguous queries that don't directly match product descriptions or categories.

## RLQR: Reinforcement Learning for Query Reformulations

A notable approach in this area is RLQR (Reinforcement Learning for Query Reformulations), which generates high-quality diverse reformulations aimed at maximizing product coverage. This technique has demonstrated a 28.6% increase in product coverage compared to standard generative models, outperforming state-of-the-art benchmarks by a significant margin[8] [9].

Unlike traditional generative approaches that may produce reformulations with low lexical diversity, RLQR explicitly optimizes for coverage, ensuring a wider variety of relevant products can be retrieved. This addresses a key limitation of conventional methods, where generated reformulations often fail to retrieve a large set of relevant products[8].

## Learning to Retrieve by Trying (LeReT)

Another innovative framework is Learning to Retrieve by Trying (LeReT), which explores search queries and uses preference-based optimization to improve their quality. LeReT can improve absolute retrieval accuracy by up to 29% and downstream generator evaluations by 17%[10]. The framework observes that language models can learn to search for relevant facts by trying different queries and learning to prioritize those that successfully produce relevant results.

LeReT's simplicity and flexibility allow it to be applied to arbitrary off-the-shelf retrievers, making it a promising technique for improving general language model pipelines. This approach is particularly valuable for complex or indirect topics where language models often struggle with posing the right search queries[10].

## Cluster-Based Language Models with Reinforcement Learning

Query clustering represents another powerful approach to enhancing e-commerce search through reinforcement learning. By grouping similar queries together, models can better adapt to the specific characteristics and intents behind different query types, ultimately delivering more accurate and personalized results.

## Cluster Language Model for E-commerce Retrieval

A novel method proposed for improving product search in e-commerce utilizes a cluster language model. This approach aims to address the limitations of bi-encoder architecture while maintaining a minimal additional training burden. The method first fine-tunes a pre-trained language model on query-product pairs using a bi-encoder approach, forming a baseline model. It then clusters training queries into k clusters and refines the baseline model for each query cluster using a novel labeling and refinement strategy[11].

The key insight behind this approach is that while the baseline model effectively captures general semantic relationships, it may not be sensitive to the specific characteristics of different query clusters. By refining the model for each cluster, it better captures the nuances of various query types and provides more accurate and personalized retrieval results[11].

## Generative Retrieval with Preference Optimization

Building on the cluster-based approach, generative retrieval with preference optimization addresses several challenges in e-commerce search, including the complexity of generating detailed item titles from brief queries and issues with long-tail queries. This framework transforms the task of generating titles from queries into generating multi-span identifiers from queries, which simplifies the generation process[2].

By employing multi-span identifiers to represent raw item titles, the framework aligns with human preferences using click data and employs a constrained search method to identify key spans for retrieving the final item. This enhances result interpretability while achieving competitive performance on real-world datasets[2].

## Learning-to-Rank-and-Retrieve with Reinforcement Learning

Traditional learning-to-rank (LTR) approaches have been effective in improving the ranking of top-k candidates by learning from customer feedback or actions. However, these benefits apply only to candidates selected during the retrieval phase. If the best candidate isn't in the candidate set, even the most sophisticated ranking model won't help customers find what they want[1].

### The LTR&R Approach

Learning-to-rank-and-retrieve (LTR&R) extends the learning-to-rank approach to include retrieval, making both components dynamic and leveraging customer feedback. This approach uses contextual multiarmed bandits, a form of reinforcement learning that optimizes the trade-off between exploring new retrieval strategies and exploiting known ones to minimize "regret"[1].

By applying reinforcement learning to both retrieval and ranking phases, e-commerce platforms can continuously improve their search systems based on user interactions. This helps address the candidate selection problem, ensuring that the most relevant products are included in the initial retrieval set before ranking occurs.

## Implementation Frameworks and Approaches

Implementing reinforcement learning for e-commerce search requires careful consideration of efficiency, scalability, and integration with existing systems. Several frameworks and approaches have emerged to address these challenges.

### Hybrid Pipelines Combining Offline and Online Learning

A promising approach is a hybrid pipeline that balances efficiency and effectiveness by combining offline knowledge distillation with online reinforcement learning. This pipeline creates a lightweight but efficient student model through offline knowledge distillation, then refines query rewriting dynamically using real-time feedback through online reinforcement learning[7].

A key innovation in this approach is using language models as simulated human feedback, enabling scalable reward signals and cost-effective evaluation without manual annotations. This addresses one of the major challenges in reinforcement learning: obtaining high-quality feedback at scale[7].

## When to Choose RL Over Traditional Fine-tuning

Research suggests three sufficient conditions for choosing reinforcement fine-tuning over supervised fine-tuning:

1. When you don't have labeled data, but can verify the correctness of the output

2. When you have some labeled data, but not much (less than 100 labeled examples)

3. When task performance improves significantly with chain-of-thought reasoning at inference time[12]

These conditions highlight reinforcement learning's particular strengths in scenarios with limited training data, which is common in specialized e-commerce domains or for new product categories with limited historical data.

## Performance Improvements and Real-world Impact

The application of reinforcement learning to language model fine-tuning for e-commerce search has demonstrated significant performance improvements across various metrics and real-world implementations.

## Quantifiable Improvements in Retrieval Metrics

Multiple studies have shown substantial improvements in key e-commerce search metrics through reinforcement learning approaches:

- RLQR demonstrated a 28.6% increase in product coverage compared to standard generative models[8] [9]

- LeReT improved absolute retrieval accuracy by up to 29% and downstream generator evaluations by 17%[10]

- Generative retrieval with preference optimization achieved competitive performance on real-world datasets, with online A/B tests demonstrating superiority in improving conversion gains[2]

These improvements translate directly to enhanced user experiences, with customers able to find more relevant products more quickly. The diversity of retrieved products also increases, exposing users to a wider range of potentially interesting items.

## Efficiency and Scalability Benefits

Beyond accuracy improvements, reinforcement learning approaches also offer efficiency advantages. ReFT can effectively fine-tune models with just a few dozen examples, making it viable for domains where large labeled datasets are unavailable[6]. This efficiency is crucial for e-commerce platforms that need to frequently update their models to accommodate new products, changing trends, and seasonal variations.

The MiniELM model proposed in one framework balances performance and efficiency through offline knowledge distillation and online reinforcement learning, addressing the high inference

latency and cost challenges that often come with using large language models in online settings[7].

## Challenges and Future Directions

Despite the promising results, implementing reinforcement learning for e-commerce search presents several challenges that need to be addressed.

### Current Limitations

Fixed reward models may suffer from inaccurate off-distribution assessments since policy optimization continuously shifts language models' data distribution. Repeatedly collecting new preference data from the latest models can alleviate this issue but makes the resulting system more complicated and difficult to optimize[13].

Additionally, the computational requirements for reinforcement learning can be substantial, potentially limiting its application in resource-constrained environments. Integration with existing e-commerce systems also presents challenges, requiring careful planning and potentially significant infrastructure updates.

### Emerging Techniques and Future Potential

Reward Learning on Policy (RLP) represents an emerging approach to address some current limitations. This unsupervised framework refines a reward model using policy samples to keep it on-distribution, introducing robust representations of policy samples and using synthetic preference generation to simulate high-quality preference data[13].

The integration of reinforcement learning with other approaches like Retrieval-Augmented Generation (RAG) also shows promise. While RAG enhances language models by enabling them to access up-to-date information from internal knowledge bases without retraining, combining it with reinforcement learning could create even more powerful and adaptive search systems[14].

## Conclusion

Reinforcement learning has emerged as a transformative approach to improving language model fine-tuning for e-commerce search. By enabling dynamic adaptation based on user feedback and interactions, RL addresses many limitations of traditional methods, particularly in scenarios with limited training data or complex user intents.

The various approaches discussed—from query reformulation using RLQR and LeReT to cluster-based models and hybrid pipelines—demonstrate the versatility and effectiveness of reinforcement learning across different aspects of e-commerce search. The significant improvements in retrieval accuracy, product coverage, and conversion rates highlight the real-world impact these techniques can have.

As reinforcement learning techniques continue to evolve and become more efficient, their integration into e-commerce search systems will likely accelerate, driving further improvements in user experience and business outcomes. For e-commerce platforms looking to enhance their

search capabilities, investment in reinforcement learning represents a promising direction with substantial potential returns in customer satisfaction and engagement.

❋

1. https://www.amazon.science/blog/from-structured-search-to-learning-to-rank-and-retrieve

2. https://arxiv.org/html/2407.19829v1

3. https://www.width.ai/post/reinforcement-learning-from-human-feedback

4. https://www.labellerr.com/blog/reinforcement-learning-from-human-feedback/

5. https://openai.com/form/rft-research-program/

6. https://www.datacamp.com/blog/reinforcement-fine-tuning

7. https://arxiv.org/html/2501.18056v1

8. https://www.amazon.science/publications/enhancing-e-commerce-product-search-through-reinforcement-learning-powered-query-reformulation

9. https://www.semanticscholar.org/paper/Enhancing-E-commerce-Product-Search-through-Query-Agrawal-Merugu/9cd7b13442bc58cc414a72a630065ba3e2fb2e4e

10. https://arxiv.org/html/2410.23214v2

11. https://aclanthology.org/2024.ecnlp-1.15.pdf

12. https://predibase.com/blog/how-reinforcement-learning-beats-supervised-fine-tuning-when-data-is-scarce

13. https://arxiv.org/html/2403.19279v1

14. https://www.k2view.com/blog/retrieval-augmented-generation-vs-fine-tuning/