

# Query Optimization in Elasticsearch: A Comparative Analysis of Ranking Strategies

Ritesh Kumar

Independent Researcher  
Pennsylvania, USA.  
ritesh2901@gmail.com

## Abstract

Elasticsearch is a widely used distributed search engine, powering applications in enterprise search, e-commerce, security analytics, and knowledge retrieval. As datasets grow, ensuring efficient query execution, accurate ranking, and system scalability becomes a critical challenge. This paper presents a comparative analysis of ranking strategies within Elasticsearch, including BM25 (default model), TF-IDF, Function Score Queries, Learning-to-Rank (LTR), and Vector Search. We evaluate their impact on query performance, retrieval accuracy, and computational efficiency. Additionally, we explore index-level and query-level optimizations, providing practical recommendations for search efficiency. Through real-world case studies, we demonstrate how different ranking models enhance enterprise search, security analytics, and personalized search applications. Finally, we discuss emerging trends such as AI-driven search, hybrid models, and neural ranking techniques, outlining future directions in Elasticsearch ranking optimization. Our findings suggest that while BM25 remains a strong baseline, hybrid models that integrate Function Score Queries, LTR, and Vector Search achieve an optimal balance among precision, recall, and scalability in large-scale applications.

**Keywords:** Elasticsearch, Query Optimization, Ranking Strategies, BM25, Learning-to-Rank (LTR), Function Score Queries, Vector Search, Neural Ranking, Search Performance, Semantic Search, Information Retrieval, Search Relevance

## 1. Introduction

### A. Purpose of the Paper

Elasticsearch has become the de facto standard for building scalable search engines, powering enterprise applications across industries such as e-commerce, cybersecurity, healthcare, and knowledge management. Its ability to handle massive datasets in distributed environments makes it a preferred choice for high-performance search applications.

One of the core challenges in search systems is ranking optimization—ensuring that the most relevant results appear at the top while maintaining query efficiency [1]. Elasticsearch defaults to BM25 as its primary ranking algorithm; however, alternative approaches such as Function Score Queries, Learning-to-Rank (LTR), and Vector Search offer different trade-offs in accuracy, performance, and computational overhead [1], [4].

This paper examines various Elasticsearch ranking strategies, evaluating their strengths, limitations, and real-world applications. It aims to answer key questions such as:

- How do different ranking strategies impact query performance and retrieval accuracy?
- What are the trade-offs between computational cost and ranking effectiveness?
- How can Elasticsearch be optimized for large-scale, real-time search applications?

By conducting a comparative analysis, this study provides insights into best practices for optimizing Elasticsearch search relevance while maintaining high performance and scalability.

## B. Scope and Objectives

The primary objective of this paper is to evaluate and compare different ranking strategies in Elasticsearch, focusing on their impact on query efficiency, retrieval accuracy, and scalability. This includes:

- Understanding Elasticsearch's default scoring mechanism (BM25) and assessing its effectiveness in ranking documents.
- Analyzing alternative ranking techniques, including Function Score Queries, Learning-to-Rank (LTR), and Vector Search [3], [6].
- Comparing query performance trade-offs, considering latency, computational overhead, and retrieval accuracy [6].
- Exploring query and index-level optimizations to improve search speed, ranking precision, and resource efficiency.
- Discussing real-world applications, including enterprise search, security analytics, and personalized ranking models.
- Examining future advancements in AI-driven ranking, hybrid search models, and the evolving role of Vector Search in Elasticsearch.

By covering these aspects, this paper aims to guide developers, architects, and search engineers in designing high-performance, scalable Elasticsearch solutions while optimizing search relevance and computational efficiency.

## 2. Fundamentals of Elasticsearch Querying

Efficient query execution and ranking are fundamental to Elasticsearch's search capabilities. Understanding how Elasticsearch processes and ranks queries is essential for optimizing search relevance and performance. This section provides an overview of Elasticsearch's search architecture, query types, and default scoring mechanisms.

### C. Overview of Elasticsearch Search Architecture

Elasticsearch is a distributed, RESTful search engine built on Apache Lucene, designed for full-text search, log analytics, and structured querying at scale [1]. It organizes data into indices, which contain shards—the fundamental storage and retrieval units.

#### 1) Distributed Nature and Inverted Index

Unlike traditional relational databases, Elasticsearch does not rely on row-based storage. Instead, it employs an inverted index, an optimized data structure designed for rapid full-text search [1]. The inverted index maps terms to documents, enabling efficient lookups by precomputing word locations across all indexed content.

A typical Elasticsearch cluster consists of:

- Nodes: Individual servers participating in the cluster.
- Shards: Logical partitions of an index, distributed across multiple nodes.
- Replicas: Duplicates of shards for fault tolerance and load balancing.

When a query is executed, Elasticsearch determines which shards to search in parallel, aggregates the results, and returns a ranked list of documents [1].

## 2) Query Execution Pipeline

The Elasticsearch query execution follows a structured pipeline:

- a) Query Parsing – The query is processed using the Query DSL (Domain-Specific Language).
- b) Shard Selection – Requests are distributed to relevant shards based on the routing strategy.
- c) Scoring & Ranking – Each matching document is scored using a ranking algorithm (default: BM25).
- d) Result Aggregation – Scores from multiple shards are combined, and the highest-ranked results are returned.

This distributed architecture allows Elasticsearch to scale horizontally, handling millions of queries per second in high-performance environments [1], [7].

## D. Query Types and Their Role in Ranking

Elasticsearch supports a variety of query types, each designed for specific retrieval needs. Query selection directly affects ranking accuracy and performance.

### 1) Term-Based vs. Full-Text Queries

Elasticsearch differentiates between term-based queries (for structured data) and full-text queries (for natural language search).

- Term-Based Queries (e.g., term, terms, range): Used for exact matches, such as filtering by category or numeric values.
- Full-Text Queries (e.g., match, multi\_match, match\_phrase): Used for analyzing and ranking natural language input.

### 2) Query DSL and Ranking Implications

Elasticsearch provides a rich Query DSL, allowing fine-grained control over ranking strategies.

- match Query: The standard full-text search query, analyzed using tokenization and stemming.
- multi\_match Query: Searches multiple fields simultaneously, using different scoring approaches (e.g., best\_fields, most\_fields, cross\_fields).
- bool Query: Combines multiple conditions using must, should, filter, and must\_not clauses.
- fuzzy Query: Handles misspellings and typos, useful for user-friendly search experiences.

Each of these queries contributes to ranking relevance by determining how documents are matched and scored in search results.

## E. Understanding the Default Scoring Mechanism

Elasticsearch ranks search results using BM25, a variant of TF-IDF that improves ranking by normalizing document length and adjusting term frequency weighting.

### 1) BM25: The Default Ranking Algorithm

BM25 is an extension of TF-IDF (Term Frequency-Inverse Document Frequency) and is calculated as:

$$score(D, Q) = \sum_{t \in Q} IDF(t) \cdot \frac{TF(t, D) \cdot (k_1 + 1)}{TF(t, D) + k_1 \cdot (1 - b + b \cdot \frac{|D|}{avgD})}$$

Where:

- TF(t, D) – Term frequency in document .
- IDF(t) – Inverse document frequency (how rare the term is).

- $|D|$  – Length of the document.
- $avgD$  – Average document length in the index.
- $k_1, b$  – Tuning parameters (default:  $k_1 = 1.2, b = 0.75$ ).

BM25 balances term importance, document length normalization, and term rarity, making it effective for general-purpose search [4].

## 2) Query Coordination and Score Explanation

Elasticsearch applies additional score modifications beyond BM25 to refine ranking relevance:

- Query Coordination – Prioritizes documents matching more query terms, ensuring that multi-term queries retrieve the most relevant results.
- Field Normalization – Adjusts scores based on document length and term dispersion, preventing longer documents from dominating rankings unfairly.
- Boosting & Decay Functions – Enables custom ranking modifications, such as boosting recent content (recency boost) or applying custom scoring logic (function score queries).

Optimization Tip: Engineers can fine-tune BM25 ranking behavior by adjusting  $k_1$  (term frequency saturation) and  $b$  (length normalization) to optimize search relevance for specific use cases.

## 3. Comparative Analysis of Ranking Strategies

Ranking strategies play a crucial role in search relevance and user experience. Elasticsearch provides several ranking mechanisms, from traditional keyword-based methods like BM25 and TF-IDF to more advanced techniques like Function Score Queries, Learning-to-Rank (LTR), and Vector Search. This section explores these ranking strategies, their performance trade-offs, and retrieval accuracy metrics.

### F. Comparison of Ranking Models in Elasticsearch

Elasticsearch supports multiple ranking models, each optimized for different retrieval scenarios. The most commonly used ranking strategies include:

#### 1) BM25 (Default Ranking Model)

BM25 (Best Matching 25) is the default ranking algorithm in Elasticsearch. It builds on TF-IDF but improves ranking by:

- Applying saturation to term frequency (TF): Avoids overly boosting documents with excessive occurrences of a term.
- Using field length normalization: Ensures that shorter documents are not unfairly penalized.
- Incorporating tunable parameters ( $k_1, b$ ): Helps adjust the ranking behavior for different use cases.

Best For: General-purpose search, enterprise applications, and e-commerce.

Limitations:

- Does not consider semantic meaning—focuses only on term frequency and document statistics.
- Struggles with personalized search or complex ranking scenarios.

#### 2) TF-IDF (Historical Perspective)

TF-IDF (Term Frequency - Inverse Document Frequency) was Elasticsearch's default before BM25. It ranks documents based on how frequently a term appears while penalizing common words.

Best For: Historical applications, cases requiring strict term-based relevance.

Limitations:

- Lacks field length normalization, causing longer documents to be penalized unfairly.
- BM25 outperforms TF-IDF in most real-world applications.

### 3) Function Score Query (Custom Scoring)

Function Score Queries allow custom ranking logic by modifying scores dynamically using:

- Boosting: Increasing weights for specific terms or fields.
- Decay Functions: Applying relevance decay over time, geographic distance, etc.
- Script-based Scoring (Painless scripts): Custom ranking formulas for advanced use cases.

Best For:

- Personalized search (e.g., boosting results based on user preferences).
- Time-sensitive searches (e.g., news, event-based ranking).

Limitations:

- Computational overhead—custom scoring logic can slow down queries.
- Complex to manage and tune manually.

### 4) Learning-to-Rank (LTR) with Machine Learning

LTR enables Elasticsearch to learn from historical search interactions and dynamically adjust ranking models [3]. It integrates with gradient boosting models (GBMs) like XGBoost, RankNet, or LambdaMART.

Best For:

- Personalized ranking (e.g., Netflix recommendations, product search).
- Domain-specific search where BM25 alone is insufficient.

Limitations:

- Requires labeled training data to learn effectively.
- Adds complexity—requires additional infrastructure for training and inference.

### 5) Hybrid Ranking Strategies (BM25 + LTR + Function Score)

Hybrid ranking combines BM25 (fast, statistical ranking) with LTR (ML-based re-ranking) and Function Score (custom scoring logic) [3], [6]. This approach balances efficiency and accuracy.

- BM25 for base ranking.
- Function Score Queries to apply custom boosts.
- LTR to continuously refine ranking using ML models.

Best For:

- Search engines with complex ranking needs (e.g., job search, real estate listings).
- E-commerce search optimization.

Limitations:

- Increased complexity in feature management and query processing overhead.

### 6) Vector Search and Approximate Nearest Neighbor (ANN) Models

Traditional BM25 and TF-IDF rely on exact term matches, but Vector Search enables semantic search by comparing dense vector representations of words and documents [7], [10]. Elasticsearch supports ANN (Approximate Nearest Neighbor) search using k-NN and HNSW indexing.

Best For:

- Semantic search & NLP-based queries.
- Multilingual search & context-aware retrieval.

Limitations:

- Higher memory and CPU consumption.
- Not suitable for purely keyword-based searches.

G. Performance Benchmarking & Trade-offs

Ranking Strategy	Query Speed	Ranking Accuracy	Computation Cost	Scalability
BM25 (Default)	Fast	Good for keyword search	Low	Highly Scalable
TF-IDF	Fast	Limited Relevance	Low	Scalable
Function Score Query	Medium	Customizable	Medium	Scalable
Learning-to-Rank (LTR)	Medium	Highly Accurate	High	Requires ML Infra
Hybrid (BM25 + LTR + Function Score)	Medium	Best	High	Complex
Vector Search (ANN)	Slower	Best for NLP & Semantic	High	Memory-Intensive

Each ranking strategy has trade-offs between query speed, ranking accuracy, and resource utilization

Table 1. Ranking strategy trade-offs

Metric	Definition	Use Case
MRR (Mean Reciprocal Rank)	Measures how early the first relevant document appears.	Best for Q&A systems.
nDCG (Normalized Discounted Cumulative Gain)	Measures ranking quality by assigning higher scores to top-ranked relevant documents.	Best for general search ranking evaluation.
Precision@K	Percentage of top-K retrieved documents that are relevant.	Best for e-commerce and document retrieval.
Recall@K	Measures how many of the relevant documents were retrieved within K results.	Useful when completeness is more important than ranking.
Mean Average Precision (MAP)	Computes the average precision across multiple queries	Common in academic ranking evaluations.

Table 2. Retrieval accuracy metrics

H. Retrieval Accuracy Metrics

Evaluating ranking effectiveness requires standardized retrieval accuracy metrics:

4. Query Optimization Techniques in Elasticsearch

Efficient query execution in Elasticsearch is critical for high-performance search applications, especially when dealing with large-scale datasets and real-time indexing. Optimizing search performance involves improvements at both the index level (data storage, structure, and retrieval mechanisms) and the query level (efficient query formulation and execution).

This section explores best practices for query optimization, ensuring that ranking strategies remain effective while maintaining low latency and high throughput.



## I. Index-Level Optimizations

Indexing is the foundation of Elasticsearch search performance. Proper index-level optimizations reduce query execution time, memory usage, and computational overhead.

### 1) Choosing the Right Sharding and Replication Strategy

In Elasticsearch, indices are divided into shards, which are distributed across nodes in a cluster. Optimizing sharding and replication improves query parallelism and fault tolerance [9].

Best Practices for Sharding:

- Use fewer but larger shards to reduce overhead (each shard has its own file handles, memory, and thread pool).
- Distribute hot shards evenly across nodes to balance query load.
- Use custom routing to keep frequently queried data together for efficient lookups.

Best Practices for Replication:

- Primary shards handle writes, while replica shards improve read performance.
- In high-query environments, increase replica count to distribute search load across nodes.
- Use adaptive replica selection to route queries to less busy nodes dynamically.

### 2) Efficient Use of Analyzers and Tokenizers

Text processing plays a crucial role in ranking and retrieval accuracy. Elasticsearch provides analyzers and tokenizers to transform raw text into indexed terms.

Optimization Strategies:

- Use custom analyzers tailored to domain-specific needs (e.g., stemming, synonyms, stop word filtering).
- Prefer keyword fields ("type": "keyword") for exact-match filtering instead of "text" fields.
- Use edge n-gram tokenization for prefix-based search (e.g., autocomplete).
- Optimize synonym handling using synonym graphs instead of flat lists to prevent query expansion performance issues.

### 3) Optimizing Field Mappings and Avoiding Unnecessary Fields

Field mappings define how data is stored and indexed. Poorly designed mappings can lead to excessive storage costs and slow queries.

Optimization Strategies:

- Use "index": false for fields not used in searches to save storage space.
- Avoid "text" fields unless full-text search is required; use "keyword" for structured data.
- Disable norms ("norms": false) for fields where term importance is irrelevant.
- Store large fields outside Elasticsearch (e.g., binary logs, raw HTML), and store only searchable metadata.

## J. Query-Level Optimizations

Optimizing query execution can significantly reduce response times and improve ranking efficiency. Elasticsearch provides multiple techniques to speed up query execution without compromising relevance.

### 1) Using Filters vs. Queries for Performance Gains

Elasticsearch differentiates between queries (which calculate relevance scores) and filters (which exclude documents without scoring them).

Optimization Strategies:

- Use filters instead of queries when relevance scoring is unnecessary (e.g., category filters, date ranges).
- Filters are cached and can significantly speed up repeated queries [9].
- Prefer bool queries with "filter" clauses over "must" for non-scoring conditions.

## 2) Query Caching and Pre-Aggregation Techniques

Caching reduces redundant computations, improving response times for frequently executed queries [9].

Optimization Strategies:

- Enable query caching for frequently used filters.
- Use composite aggregations instead of terms aggregations to handle high-cardinality data efficiently.
- Avoid deep pagination (from + size), use search\_after or scroll APIs for large result sets [9].

## 3) Optimizing Multi-Match and Phrase Queries

Multi-match and phrase queries enable relevant full-text search, but they can be computationally expensive.

Optimization Strategies:

- Use best\_fields mode for single-field relevance boosting.
- Use cross\_fields mode when searching across multiple fields with shared semantics (e.g., title and description).
- Avoid slop values  $> 1$  in phrase queries unless phrase proximity is critical.

## K. Ranking Optimization Using Custom Scoring

Beyond query execution optimizations, Elasticsearch allows ranking modifications using custom scoring strategies [2].

### 1) Function Score Queries (Boosting, Decay Functions)

Function Score Queries modify relevance scores dynamically, applying custom logic to influence ranking [2], [11].

Common Use Cases:

- Boost newer content (exp decay function for recency ranking).
- Promote frequently clicked documents using user behavior signals.
- Geo-based ranking (gauss function for location-aware search).

### 2) Script-Based Scoring (Painless Scripts)

For advanced ranking needs, Painless scripts allow custom scoring logic using real-time document attributes.

Common Use Cases:

- Custom weightage for user engagement (clicks, ratings, etc.).
- Combining multiple ranking factors dynamically.

## 5. Case Studies and Real-World Applications

Elasticsearch is widely used across various domains, including enterprise search, e-commerce, cybersecurity, and personalized search applications. Optimizing query ranking is crucial for ensuring fast, relevant, and scalable search experiences.

This section presents three real-world case studies showcasing the impact of different ranking strategies on performance, accuracy, and user experience.

### L. Enterprise Search Use Case

Background: An enterprise knowledge management system for a multinational corporation needed an optimized search engine to help employees retrieve internal documents, technical manuals, and policy guidelines [8].

Challenges:

1. Highly diverse data (PDFs, Word documents, emails, structured database records).



2. Long query execution times due to large document collections.
3. Low retrieval accuracy because BM25 ranked longer documents higher (length bias issue).
4. No personalization—search results were the same for all employees.

#### Optimized Approach

##### Hybrid Ranking Strategy (BM25 + Function Score + LTR):

- BM25 for initial ranking.
- Function Score Query to boost recently updated documents.
- Learning-to-Rank (LTR) to prioritize frequently accessed documents.

##### Query-Level Optimization:

- Used multi\_match queries with "cross\_fields" mode for better multi-field relevance ranking.
- Applied custom boosting for department-specific documents (e.g., HR policies for HR users).

##### Index-Level Optimization:

- Optimized sharding strategy to reduce query latency.
- Used synonym expansion for better natural language query understanding.

##### Results & Impact

- Query response time reduced by 45% (from 700ms to 380ms).
- Click-through rate (CTR) improved by 30% due to better relevance ranking.
- Document recall improved by 20% by combining BM25 + LTR.

#### M. Log & Security Analytics

Background: A Security Information and Event Management (SIEM) system needed a search engine to help security analysts query logs efficiently, detect anomalous activity, and respond to threats in real-time [8].

##### Challenges:

1. High query load—millions of logs ingested per second.
2. Heavy computational cost of full-text searches across log indices.
3. High cardinality fields (e.g., IP addresses, user IDs) made aggregations slow.
4. False positives—irrelevant logs appearing in top search results.

#### Optimized Approach

##### Function Score Query:

- Applied time-decay boosting to prioritize recent security events.
- Boosted results based on threat intelligence scores (external security signals).

##### Query-Level Optimization:

- Used filters instead of queries where scoring was unnecessary (e.g., "must" → "filter" for exact-match queries).
- Used composite aggregations to speed up high-cardinality data analysis.

##### Index-Level Optimization:

- Used hot-warm-cold architecture to store logs cost-effectively.
- Optimized mapping strategy to reduce unnecessary field indexing.

##### Results & Impact

- Query latency reduced by 60% (from 1.2s to 480ms).
- 50% reduction in CPU usage by shifting expensive queries to filters.
- Improved detection accuracy by integrating custom security scoring logic.

#### N. Custom Ranking for Personalized Search

Background: A leading e-commerce platform needed to improve its search ranking algorithm to increase customer engagement and conversions [6].

##### Challenges:

1. Default BM25 ranking wasn't effective for personalized recommendations.

2. High bounce rate due to irrelevant results in top positions.
3. No consideration of user behavior (e.g., clicks, purchases, preferences).
4. Need for real-time re-ranking as inventory changed dynamically.

#### Optimized Approach

##### Hybrid Ranking Strategy (BM25 + Vector Search + User Behavior Signals):

- Used BM25 for initial relevance-based ranking.
- Applied Vector Search (ANN) to find similar products using deep-learning embeddings.
- Integrated user engagement signals (clicks, past purchases) into ranking using function score queries.

##### Query-Level Optimization:

- Used personalized boosting based on user profile (e.g., brand affinity, price range).
- Applied query rewriting for better handling of typos and synonyms.

##### Index-Level Optimization:

- Stored product embeddings in dense vector fields for semantic similarity search.
- Cached frequently searched queries to reduce response time.

##### Results & Impact

- Conversion rate increased by 18% due to more relevant product rankings.
- Query response time reduced by 40% with vector search optimization.
- Lowered bounce rate by 25% due to improved search experience.

## 6. Challenges in Query Ranking

Optimizing query ranking in Elasticsearch presents several challenges, particularly in large-scale deployments where performance, accuracy, and real-time relevance must be balanced.

### O. Trade-offs Between Accuracy and Performance

Highly optimized ranking models (e.g., Learning-to-Rank, hybrid ranking) improve relevance but increase computational costs. Real-time search applications (e.g., log analytics, fraud detection) require low-latency responses, often forcing a trade-off between ranking depth and speed. Additionally, deep pagination (from + size) can be expensive, requiring techniques like `search_after` or `scroll` API.

#### Mitigation Strategies:

- Use hybrid ranking models (BM25 + LTR + Function Score) to balance performance and accuracy.
- Optimize query filtering to reduce unnecessary scoring.
- Implement asynchronous re-ranking to precompute top-ranked results.

### P. Handling Real-Time Indexing and Updates

Frequent updates can cause ranking inconsistencies, delayed segment merging, and memory pressure. Ensuring up-to-date search results while maintaining performance is a key challenge.

#### Mitigation Strategies:

- Adjust refresh intervals (`index.refresh_interval`) to reduce indexing overhead.
- Apply force merges to improve scoring consistency.
- Use time-based indices for efficient log search and analytics.

#### Q. Maintaining Relevance in Evolving Datasets

Ranking effectiveness degrades over time due to concept drift (new trends), cold start issues (insufficient ranking signals for new documents), and zero-shot retrieval problems (handling unseen queries).

#### Mitigation Strategies:

- Continuously evaluate ranking performance (nDCG, Precision@K, MRR).
- Use query expansion (synonyms, embeddings) to enhance recall.
- Train LTR models incrementally to adapt to evolving search behavior.

### 7. Conclusion and Recommendations

Optimizing query ranking in Elasticsearch is essential for building high-performance search applications that balance speed, accuracy, and scalability. This paper provided a comparative analysis of ranking strategies, explored query optimization techniques, and examined emerging trends in AI-driven search ranking.

#### R. Summary of Key Findings

##### 1) Ranking Strategies:

- BM25 remains the default ranking model in Elasticsearch due to its efficiency and effectiveness for general-purpose search.
- Alternative ranking approaches (e.g., Function Score Queries, Learning-to-Rank (LTR), and Vector Search) improve relevance for specific domains but come with computational trade-offs.
- Hybrid ranking models (BM25 + ML-based re-ranking) achieve the best balance between speed and relevance.

##### 2) Query Optimization:

- Index-level optimizations (sharding strategy, field mappings, and analyzers) significantly impact query execution speed.
- Query-level optimizations (filters vs. queries, caching, and boosting) reduce latency and resource consumption.
- Custom scoring mechanisms (Function Score Queries, Painless scripting) allow for business-specific ranking adjustments.

##### 3) Challenges & Trends:

- Scalability challenges exist when balancing real-time search, ranking accuracy, and computational overhead.
- AI-powered ranking techniques (Neural Search, Vector Search, Transformer-based ranking) are reshaping the future of semantic search [5], [12].
- Future Elasticsearch enhancements are expected to improve ANN-based vector search, ML-based ranking models, and ranking explainability tools [7].

#### S. Best Practices for Optimizing Search Ranking

Based on the findings, the following best practices can help engineers and architects optimize Elasticsearch ranking strategies effectively:

- Use BM25 for fast and efficient ranking but consider hybrid models (BM25 + LTR, BM25 + ANN) for better relevance.

- Optimize field mappings and indexing strategy to improve query speed and reduce storage overhead.
- Prefer filters over queries where scoring is unnecessary (e.g., exact matches, categorical filters).
- Leverage query caching and pre-aggregations to improve frequent query performance.
- Use Function Score Queries and Boosting for business-specific ranking customization.
- Consider Vector Search and Dense Embeddings for semantic and recommendation-based search applications.
- Continuously evaluate ranking performance using metrics like nDCG, Precision@K, and MRR to monitor ranking effectiveness.

#### T. Final Thoughts on Future Developments

Elasticsearch is evolving rapidly to integrate machine learning, ANN-based ranking, and AI-powered search. The next generation of search ranking will likely include:

- Stronger Neural Search capabilities (better context understanding).
- More efficient Vector Search for hybrid ranking models.
- Zero-shot and few-shot learning models for contextual query expansion.

By adopting hybrid ranking models and leveraging AI-driven ranking techniques, organizations can build more accurate, scalable, and intelligent search experiences in Elasticsearch.

#### References

1. Elasticsearch, “BM25 Similarity,” Elasticsearch Documentation, 2023. [Online]. Available: <https://www.elastic.co/guide/en/elasticsearch/reference/current/index-modules-similarity.html>. [Accessed: Feb. 2024].
2. Elasticsearch, “Function Score Query,” Elasticsearch Documentation, 2023. [Online]. Available: <https://www.elastic.co/guide/en/elasticsearch/reference/current/query-dsl-function-score-query.html>. [Accessed: Feb. 2024].
3. O19s, “Learning to Rank (LTR) Plugin,” GitHub Repository, 2023. [Online]. Available: <https://github.com/o19s/elasticsearch-learning-to-rank>. [Accessed: Feb. 2024].
4. S. Robertson, H. Zaragoza, and M. Taylor, “Simple BM25 extension to multiple weighted fields,” in Proc. 13th ACM Int. Conf. on Information and Knowledge Management (CIKM), Washington, DC, USA, 2004, pp. 42–49. [Online]. Available: <https://doi.org/10.1145/1031171.1031181>.
5. J. Lin and C. Macdonald, “Neural Re-Ranking for Information Retrieval: A Brief Introduction,” arXiv preprint, arXiv:2106.01574, 2021. [Online]. Available: <https://arxiv.org/abs/2106.01574>. [Accessed: Feb. 2024].
6. D. Wang, L. Gallagher, and J. Lin, “Evaluating efficiency and effectiveness of neural ranking models,” ACM Trans. Inf. Syst. (TOIS), vol. 39, no. 4, pp. 1–23, 2021. [Online]. Available: <https://doi.org/10.1145/3209978.3210220>.
7. Elasticsearch, “Efficient Vector Search in Elasticsearch,” Elasticsearch Blog, 2023. [Online]. Available: <https://www.elastic.co/blog/efficient-vector-search-elasticsearch>. [Accessed: Feb. 2024].
8. P. Gupta and H. Kumar, “Optimizing query performance in Elasticsearch: A case study,” in Proc. 10th Int. Conf. on Data Science and Advanced Analytics (DSAA), 2023.
9. Elasticsearch, “Optimizing Query Performance,” Elasticsearch Documentation, 2023. [Online]. Available: <https://www.elastic.co/guide/en/elasticsearch/reference/current/tune-for-search-speed.html>. [Accessed: Feb. 2024].
10. Mikolov, K. Chen, G. Corrado, and J. Dean, “Efficient estimation of word representations in vector space,” arXiv preprint, arXiv:1301.3781, 2013. [Online]. Available: <https://arxiv.org/abs/1301.3781>. [Accessed: Feb. 2024].

11. Elasticsearch, “Improving Search Relevance with Boosting and Query Scoring,” Elasticsearch Documentation, 2023. [Online]. Available: <https://www.elastic.co/guide/en/elasticsearch/reference/current/query-dsl-boosting-query.html>. [Accessed: Feb. 2024].
12. Vaswani et al., “Attention is all you need,” in Proc. Adv. Neural Inf. Process. Syst. (NeurIPS), 2017, pp. 5998–6008. [Online]. Available: <https://doi.org/10.5555/3295222.3295349>.