

# UniDex: Rethinking Search Inverted Indexing with Unified Semantic Modeling

Zan Li<sup>†</sup>, Jiahui Chen<sup>†</sup>, Yuan Chai<sup>\*†</sup>, Xiaoze Jiang<sup>\*†</sup>, Xiaohua Qi<sup>†</sup>, Zhiheng Qin, Runbin Zhou, Shun Zuo, Guangchao Hao, Kefeng Wang, Jingshan Lv, Yupeng Huang, Xiao Liang, Han Li  
Kuaishou Technology, Beijing, China  
{lizan07,chenjiahui11,chaiyuan,jiangxiaoze,qixiaohua03}@kuaishou.com

## Abstract

Inverted indexing has traditionally been a cornerstone of modern search systems, leveraging exact term matches to determine relevance between queries and documents. However, this term-based approach often emphasizes surface-level token overlap, limiting the system’s generalization capabilities and retrieval effectiveness. To address these challenges, we propose UniDex, a novel model-based method that employs unified semantic modeling to revolutionize inverted indexing. UniDex replaces complex manual designs with a streamlined architecture, enhancing semantic generalization while reducing maintenance overhead. Our approach involves two key components: UniTouch, which maps queries and documents into semantic IDs for improved retrieval, and UniRank, which employs semantic matching to rank results effectively. Through large-scale industrial datasets and real-world online traffic assessments, we demonstrate that UniDex significantly improves retrieval capabilities, marking a paradigm shift from term-based to model-based indexing. Our deployment within Kuaishou’s short-video search systems further validates UniDex’s practical effectiveness, serving hundreds of millions of active users efficiently.

## CCS Concepts

• Information systems → Novelty in information retrieval.

## Keywords

Information Retrieval; Sparse Retrieval; Representation learning

## ACM Reference Format:

Zan Li<sup>†</sup>, Jiahui Chen<sup>†</sup>, Yuan Chai<sup>\*†</sup>, Xiaoze Jiang<sup>\*†</sup>, Xiaohua Qi<sup>†</sup>, Zhiheng Qin, Runbin Zhou, Shun Zuo, Guangchao Hao, Kefeng Wang, Jingshan Lv, Yupeng Huang, Xiao Liang, Han Li. 2025. UniDex: Rethinking Search Inverted Indexing with Unified Semantic Modeling. In *The Preprint Version of UniDex, Beijing, China*. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/xxxxx>

<sup>\*</sup>Corresponding author. <sup>†</sup>Equal contribution.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

UniDex, Beijing, China

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN xxxxx/2025/09  
<https://doi.org/10.1145/xxxxx>

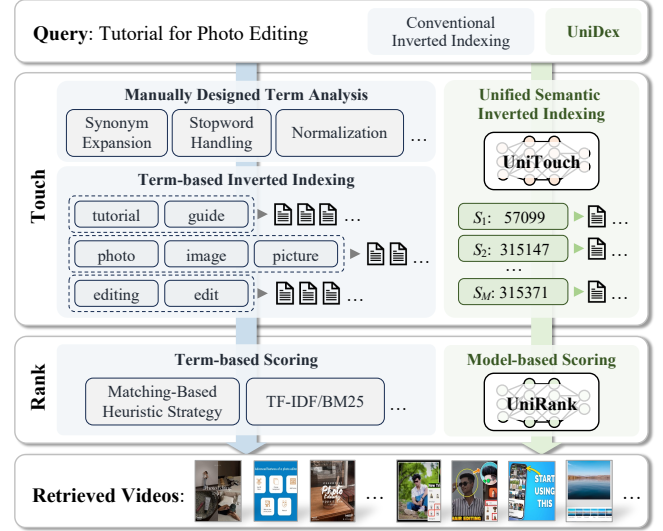


Figure 1: Comparison of Conventional Inverted Indexing and UniDex. UniDex substitutes the various manually designed term-based Touch and Rank components with the unified semantic modeling approaches, UniTouch and UniRank.

## 1 Introduction

Inverted indexing has long served as a foundational component of modern search systems due to its efficiency in large-scale information retrieval [1, 9, 41]. The conventional paradigm tokenizes documents into discrete terms and builds an index mapping each term to its associated documents [9, 10]. Relevance between queries and documents is then primarily determined through exact or approximate term matches [31, 35]. While effective in terms of efficiency, this approach inherently emphasizes surface-level lexical overlap, thereby limiting generalization and hindering retrieval effectiveness [6]. Breaking away from this term-centric paradigm to enhance semantic retrieval remains an open challenge.

As illustrated in the blue blocks of Figure 1, industrial-scale inverted indexing systems are typically organized into two stages: the **Touch** stage and the **Rank** stage. The Touch module performs term-level retrieval to construct a candidate set, while the Rank module scores these candidates with respect to the query. Over the years, a variety of manually engineered heuristics have been introduced at both stages to improve recall and ranking quality [1, 3, 41]. For example, synonym expansion, stopword handling, and term normalization are employed during Touch to enlarge the candidate pool, while heuristic strategies such as BM25 [41] or TF-IDF [35],

often combined with handcrafted term-importance features, are used in Rank to refine relevance estimation [29, 35, 41]. Despite their widespread adoption [6], these heuristics are costly to maintain, brittle across domains, and fundamentally constrained by the limitations of term-level modeling.

With the rise of deep semantic models [16, 18, 38], many recent studies attempt to integrate neural representations into inverted indexing [9, 10, 19]. However, most of these works remain tied to the term-based paradigm: semantic models are often used to generate richer sets of query or document terms, but the retrieval process itself still relies on lexical overlap. Consequently, such methods [1, 19] inherit the weaknesses of traditional indexing, including reliance on handcrafted matching strategies and limited semantic generalization.

Motivated by these limitations, we propose **UniDex**, a unified model-based framework that rethinks inverted indexing from a semantic perspective. As shown in the green blocks of Figure 1, UniDex replaces manually designed term-level components in both Touch and Rank with two semantic modeling modules: **UniTouch** and **UniRank**. Specifically, UniTouch maps queries and documents into discrete semantic IDs (SIDs) through a dual-tower encoder with an integrated quantization module. To better capture the inherent ambiguity of queries and the compositional semantics of documents, UniTouch represents each input with multiple learnable tokens, each token corresponding to a potential semantic aspect. For efficient retrieval, we design a Max–Max matching strategy that aligns query tokens with document tokens in a manner consistent with the inverted indexing lookup logic. This design ensures that a document can be retrieved if it matches at least one semantic aspect of the query, while preserving the scalability of inverted indexing. As a result, UniTouch can recall documents whose semantics are close to the query even when they do not share overlapping terms, thus significantly improving generalization over traditional term-based inverted indexes. On top of this retrieval stage, UniRank employs another semantic model with token-level interactions to precisely rank the retrieved candidates, going beyond conventional heuristic term-matching approaches. Together, UniTouch and UniRank form UniDex: a unified retrieval–ranking pipeline that combines the scalability of inverted indexing with the expressive power of semantic modeling. Unlike traditional pipelines, UniDex eliminates dependence on handcrafted rules, generalizes more effectively across diverse queries and domains, and substantially reduces engineering overhead, making it a practical solution for real-world search platforms.

We comprehensively evaluate UniDex on large-scale industrial datasets as well as real-world online traffic. The main contributions of this work are summarized as follows:

(1) We present the first paradigm shift of inverted indexing from a term-based to a model-based framework. This reformulation opens new directions for the evolution of retrieval systems and carries significant implications for both research and industry.

(2) We propose **UniDex**, a unified semantic modeling framework that replaces manually engineered term-matching and multi-path relevance computations with semantic ID–based indexing and model-driven ranking. This design substantially reduces maintenance overhead while providing stronger semantic generalization.

(3) To the best of our knowledge, we achieve the first successful large-scale deployment of a model-based inverted indexing system in industry. UniDex has been integrated into Kuaishou’s short-video search platform, where extensive online A/B testing confirms its effectiveness and efficiency. It currently serves hundreds of millions of active users, demonstrating both the practicality and impact of our approach.

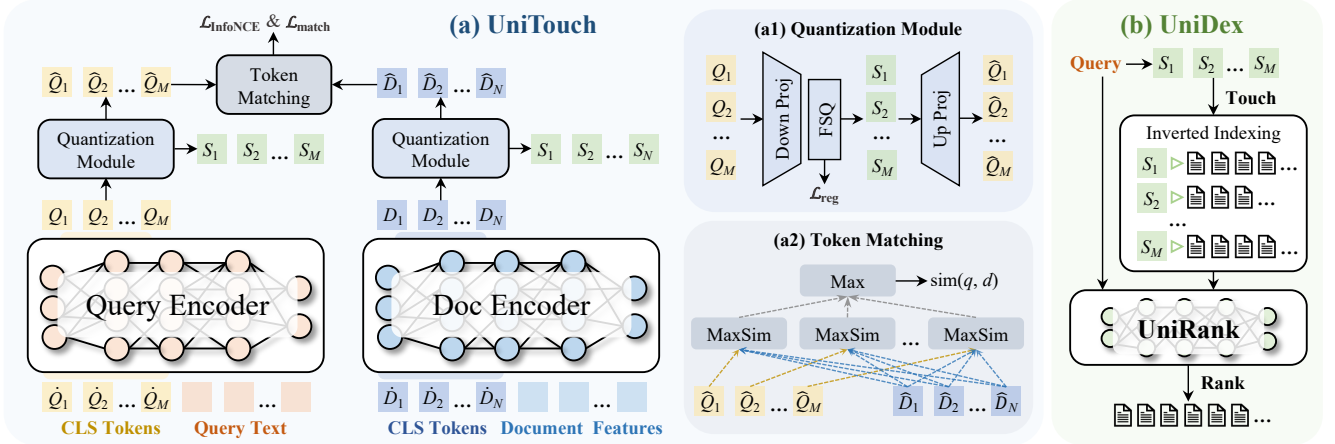
## 2 Related Works

### 2.1 Inverted Indexing

In the field of information retrieval, inverted indexing remains the most classical and widely adopted indexing structure. Its core principle is term-based indexing, which leverages the Bag-of-Words model to establish mappings between documents and terms, thereby enabling efficient matching and retrieval. Traditional approaches such as TF-IDF [35] and BM25 [41] compute query–document relevance by weighting term frequency and inverse document frequency. Despite their effectiveness, these methods inherently rely on exact lexical overlap, making them susceptible to vocabulary mismatch and unable to capture deeper semantic relationships. The rise of deep semantic modeling has driven extensive research into integrating inverted indexing with neural networks. For instance, DeepCT [6] leverages pre-trained language models to derive context-aware term weights, enabling finer-grained term weighting. However, it remains fundamentally limited by unresolved vocabulary mismatch. Generative approaches, such as doc2query [28] and doc2query-T5 [29], attempt to mitigate this limitation by predicting potential query terms to expand document representations, thus indirectly reinforcing salient terms; however, their training paradigm is inherently indirect. Another approach involves directly modeling the interactions between query or document tokens and the entire vocabulary, constructing an interaction matrix that is subsequently aggregated to obtain term-level importance scores. SparTerm [1] applies summation aggregation, while SPARTA [47] and EPIC [25] use max pooling, yet their representations are either insufficiently sparse without additional top- $k$  pruning or lack explicit sparsity-inducing regularization. The SPLADE family of models [3, 9, 10, 19] remedies these limitations by combining sparsity regularization with lexical expansion, thereby improving efficiency, interpretability, and retrieval effectiveness, though the paradigm remains largely constrained by token-level interactions and heuristic term-matching rules. This reliance limits generalization and hinders the ability to fully exploit model-based representations. In contrast, UniDex, as a unified semantic modeling approach, eliminates complex handcrafted rules and enhances generalization capabilities.

### 2.2 Other Retrieval Methods

Beyond inverted indexing, dense retrieval has emerged as a complementary paradigm that encodes queries and documents into a shared embedding space for semantic matching. This typically involves an offline document encoding and indexing phase, followed by an online query encoding and similarity search stage [7, 21, 33]. Early studies employed dual-encoder architectures built upon pre-trained language models (PLMs) such as BERT [7, 21], demonstrating strong performance on semantic retrieval benchmarks. Subsequent advances have sought to enhance dense retrievers through



**Figure 2: Overview of our framework. (a) The UniTouch model architecture, consisting of a Query Encoder, a Document Encoder, and Quantization Modules. (a1) The internal design of the Quantization Module. (a2) The token matching strategy in UniTouch, defined as the maximum entry of the token-level matching matrix. (b) The end-to-end retrieval pipeline of UniDex: a query is first discretized into semantic IDs via UniTouch, relevant documents are retrieved from the semantic inverted indexing, and the final candidate set is ranked by UniRank.**

task-specific pretraining (e.g., Condenser [14], coCondenser [12], PAIR [34]), improved negative sampling strategies [42, 44], and knowledge distillation from stronger teacher models [32, 33]. To capture fine-grained query-document interactions, token-level and multi-representation encoders have been proposed, including ColBERT [18], ColBERTv2 [36], ME-BERT [24], COIL [13], uniCOIL [22] and MVR[46]. These methods combine expressive embeddings with efficient late interaction mechanisms or sparse representations, thereby balancing retrieval effectiveness and computational efficiency. Efficiency-oriented research has further explored embedding compression, product quantization, and binary encoding to reduce memory footprint and latency [40, 43, 45]. More recently, large language models (LLMs) have been leveraged for dense retrieval owing to their strong semantic understanding, multi-task adaptability, and extended context modeling, enabling more sophisticated query-document interactions and effective long-document retrieval [4, 23]. Compared with our model-based inverted indexing framework, dense retrieval depends on computationally intensive neural similarity search, while our approach eliminates manual term-matching design and benefits from the inherent efficiency of inverted search. In addition, statistical retrieval methods are also used as complementary retrieval strategies [11, 15, 37]. Techniques such as item-to-item retrieval [15] and collaborative filtering [8, 37] exploit user-content associations to expand candidate sets. While effective in capturing behavioral relevance, these methods suffer from cold-start issues and the Matthew effect [11, 48], and are thus rarely deployed in isolation. Instead, they serve as auxiliary strategies when integrated into more robust retrieval frameworks. Our approach enhances inverted indexing by integrating semantic generalization capabilities, making it inherently more effective in cold-start scenarios than traditional statistical methods. This integration allows the system to better understand and relate to new and unseen queries. It elevates the potential for iterative improvements in inverted indexing.

### 3 Methodology

In this section, we first introduce the two components of UniDex: the UniTouch and UniRank models. Subsequently, we present the complete UniDex search pipeline.

#### 3.1 UniTouch

**3.1.1 Model Architecture.** The UniTouch model learns semantic representations of queries and documents, encoding them as discrete Semantic IDs (SIDs). As shown in Figure 2 (a), UniTouch consists of two major components: (i) a Query/Document Encoder, which captures semantic information, and (ii) Quantization Modules, which discretizes continuous embeddings into compact symbolic identifiers.

**Query Encoder.** Given a user query  $q$ , we first tokenize it into a sequence  $\{q_1, q_2, \dots, q_m\}$ . To capture the polysemy of queries, we append  $M$  learnable tokens  $\{\hat{Q}_1, \hat{Q}_2, \dots, \hat{Q}_M\}$  to the sequence. The encoder is instantiated as a  $L$ -layer BERT model. After the forward pass, we extract the embeddings of the  $M$  appended tokens, yielding  $\{Q_1, Q_2, \dots, Q_M\} \in \mathbb{R}^d$ , which collectively represent the semantic embedding of the query. Each token embedding captures a distinct semantic aspect of the query.

**Document Encoder.** The document encoder is designed analogously. For a document  $d$  tokenized as  $\{d_1, d_2, \dots, d_n\}$ , we append  $N$  learnable tokens  $\{\hat{D}_1, \hat{D}_2, \dots, \hat{D}_N\}$ . The forward pass produces  $N$  embeddings  $\{D_1, D_2, \dots, D_N\} \in \mathbb{R}^d$ , with each encoding a distinct semantic component of the document. Since documents typically contain richer and more diverse semantics than queries, we set  $N > M$ . During training, the query and document encoders share model parameters, ensuring semantic alignment across domains.

**Quantization Module.** To convert continuous embeddings into discrete semantic IDs, we introduce a quantization module. The detailed network of quantization module is shown in Figure 2 (a1). Taking a query embedding  $Q_i \in \mathbb{R}^d$  as an example, we first project

it into a lower-dimensional space:

$$Q'_i = \text{DownProj}(Q_i) \in \mathbb{R}^{d_q}, \quad d_q \ll d. \quad (1)$$

Next, we apply the *Finite Scalar Quantization (FSQ)* [27] algorithm, which discretizes each dimension of  $Q'_i$  into one of  $K$  bins:

$$S_i = \text{FSQ}(Q'_i) \in \{0, 1, \dots, K-1\}^{d_q}, \quad (2)$$

$$\text{FSQ}(Q'_i) \triangleq \text{Round}[(K-1)\sigma(Q'_i)], \quad (3)$$

where  $\sigma(\cdot)$  is the sigmoid function. Thus, each  $S_i$  corresponds to one of the possible discrete codes  $K^{d_q}$ , which we interpret as the semantic IDs of the query. Finally, the discrete code  $S_i$  is mapped back into the  $d$ -dimensional embedding space via an up-projection:

$$\hat{Q}_i = \text{UpProj}(S_i) \in \mathbb{R}^d. \quad (4)$$

This reconstructed embedding  $\hat{Q}_i$  allows the model to integrate discrete IDs into downstream tasks while maintaining compatibility with continuous representation spaces. During optimization, we adopt the Element-wise Gradient Scaling (EWGS) [20] strategy to propagate gradients through the non-differentiable quantization step. This approach mitigates training instability introduced by Straight-Through Estimator (STE) [2], leading to faster convergence and more stable learning dynamics in the quantization module. Empirically, we observe that setting  $K = 2$  (i.e., binary quantization) and  $d_q = 19$  yields favorable quantization performance. Under this configuration, each semantic ID corresponds to an integer within the range  $[0, 2^{19}]$ , balancing representation capacity and model efficiency.

**Token Matching.** To measure the semantic relevance between queries and documents, we adopt a “Max-Max” matching strategy, as illustrated in Figure 2 (a2). Following the term-match essence of inverted indexing, a document should be retrieved if it matches at least one semantic aspect of the query. For instance, consider the query “apple”, which may refer to the fruit, the technology company, or even the record label. A document discussing global technology trends may contain sections on Google, Microsoft, and Apple. Even if only the “Apple Inc.” aspect overlaps with the query, the document should still be retrieved.

Formally, the similarity between the query and the document is then defined as the maximum entry of the token-level matching matrix (other matching methods detailed in Section 4.3):

$$\text{sim}(q, d) = \max_{i \in [M]} \max_{j \in [N]} s(\hat{Q}_i, \hat{D}_j) = \max_{i \in [M]} \max_{j \in [N]} \frac{\hat{Q}_i \cdot \hat{D}_j}{\|\hat{Q}_i\| \cdot \|\hat{D}_j\|}. \quad (5)$$

**3.1.2 Learning Objectives.** To optimize the UniTouch model, we introduce multiple training objectives that jointly enhance both semantic representation learning and the alignment between SIDs of query and documents. The training data are derived from real-world search logs. Each training instance consists of a triplet  $\{q, \mathcal{D}, \mathcal{Y}\}$ , where  $q$  denotes the user query,  $\mathcal{D} = \{d_1, d_2, \dots, d_n\}$  represents the set of candidate documents (including those displayed to the user as well as unexposed documents sampled as negatives), and  $\mathcal{Y} = \{l_1, l_2, \dots, l_n\}$  contains the associated relevance labels, which reflect the semantic relevance between queries and documents.

**Contrastive Learning.** We first utilize a list-wise contrastive learning strategy to refine semantic representations. For a given query

$q$ , we compute the InfoNCE loss over its candidate document set. Specifically, for each document  $d \in \mathcal{D}$ , its loss is formulated as:

$$\mathcal{L}_{\text{InfoNCE}} = -\frac{1}{|\mathcal{D}|} \sum_{d_i \in \mathcal{D}} \log \frac{\exp(\text{sim}(q, d_i)/\tau)}{\sum_{d_j \in \{d_i\} \cup \mathcal{N}(d_i)} \exp(\text{sim}(q, d_j)/\tau)}, \quad (6)$$

where  $\tau$  is a temperature parameter, and  $\mathcal{N}(d_i) = \{d_k \in \mathcal{D} | l_i > l_k\}$  denotes the set of negatives. In our design, negatives include documents with lower relevance labels as well as documents from other queries within the same batch. Compared with point-wise learning, contrastive optimization captures finer-grained semantic distinctions and provides stronger supervision for representation learning.

**Matching Loss.** To enable UniTouch to be more effectively applied in inverted indexing structures, we further introduce a matching loss that reinforces consistency between query SIDs and the SIDs of highly relevant documents. The intuition is that if a document is highly relevant to a query, their discretized semantic identifiers should overlap in at least part of the code space. During training, this objective is only applied to documents with the highest relevance grade:

$$\mathcal{L}_{\text{match}} = \frac{1}{|\mathcal{D}'|} \sum_{d_i \in \mathcal{D}'} (1 - \text{sim}(q, d_i)), \quad (7)$$

where  $\mathcal{D}' = \{d_i \in \mathcal{D} | l_i = \max(\mathcal{Y})\}$  denote the set of documents with highest relevance label.

**Quantization Regularization.** Finally, to improve quantization stability when  $K = 2$ , we introduce a regularization loss that encourages embeddings to stay away from the decision boundary (i.e., 0.5) in each dimension:

$$\mathcal{L}_{\text{reg}} = \frac{1}{M} \sum_{i \in [M]} \|\sigma(Q'_i) - 0.5\| - 0.5\|^2. \quad (8)$$

This objective prevents embeddings from oscillating around the quantization thresholds, thereby stabilizing training and ensuring robust semantic discretization.

**3.1.3 Semantic Inverted Indexing.** After training the UniTouch model, we deploy it into the online search platform. The key idea is to replace traditional term-based inverted indexing with a *semantic inverted indexing* constructed from the model’s discrete SIDs.

Specifically, we first perform offline inference on the document encoder to compute SIDs for all candidate items, and then build an inverted index where each entry corresponds to a particular SID. At query time, the query encoder generates SIDs for the incoming user query in real time. An important feature of our design is that the  $M$  query SIDs represent multiple potential semantic interpretations of the query. Consequently, the retrieval process aggregates results from all  $M$  query SIDs. Formally, the final candidate set is obtained as the union of documents retrieved under each query SID, ensuring broad semantic coverage and robust recall performance.

## 3.2 UniRank

Once UniTouch produces a large pool of retrieved candidates, traditional term-based ranking methods become less effective, as many of the retrieved items may not share explicit lexical overlap with the query. To address this challenge, we propose UniRank, a semantic relevance ranking model that operates on top of the

semantic inverted indexing built by UniTouch. UniRank replaces multiple handcrafted lexical matching heuristics with a unified neural ranking framework, thereby simplifying the indexing pipeline and improving generalization.

**3.2.1 Model Architecture.** The architecture of UniRank follows a dual-tower design similar to UniTouch. Both the query and document encoders append a set of learnable tokens after their respective input sequences to capture fine-grained semantic information. Unlike UniTouch, which emphasizes broad coverage for retrieval, UniRank focuses on precise semantic alignment for ranking. Inspired by ColBERT [18], UniRank computes token-level interactions between query and document embeddings to achieve high-resolution semantic matching. Formally, given  $M$  query embeddings  $\{Q_1, \dots, Q_M\}$  and  $N$  document embeddings  $\{D_1, \dots, D_N\}$ , the matching score is defined as:

$$\text{sim}(q, d) = \sum_{i=1}^M \max_{j \in [N]} \frac{Q_i \cdot D_j}{\|Q_i\| \cdot \|D_j\|}. \quad (9)$$

To train the UniRank model, we use a list-wise contrastive objective as described in Equation 6 to optimize the ranking capability, and distill the fine-grained ranking scores into the model through mean squared error (MSE) loss to enhance its relevance discrimination ability.

**3.2.2 Deployment.** During deployment, the document encoder pre-computes embeddings for all candidate items in an offline stage, and these embeddings are stored in memory. At query time, the query encoder generates embeddings in real time. The final ranking score is computed according to Equation 9, and documents are sorted accordingly to produce the ranked results.

### 3.3 UniDex

The overall search pipeline, termed UniDex, integrates UniTouch and UniRank into a unified two-stage framework, as shown in Figure 2 (b). In the first stage, UniTouch leverages semantic inverted indexing to efficiently retrieve a broad set of candidate documents. By discretizing query and document embeddings into semantic IDs, UniTouch enables scalable retrieval that captures multiple semantic interpretations of a query. This ensures high recall while maintaining efficiency comparable to traditional term-based inverted indexing. In the second stage, UniRank takes the retrieved candidate pool and performs fine-grained semantic ranking. While UniTouch focuses on recall-oriented retrieval, UniRank emphasizes precision by computing detailed token-level interactions between queries and documents. This step filters out semantically weaker candidates and surfaces the most relevant results at the top of the list.

Compared to conventional inverted indexing, which relies on exact term matching, synonym expansion, and manually designed matching heuristics, UniDex offers a more generalizable and efficient solution. By replacing term-level indices with SIDs, UniDex eliminates the dependence on synonym rewriting and alleviates the computational overhead of handcrafted strategies. As a result, the system achieves robust generalization, scalable efficiency, and improved retrieval quality.

## 4 Experiments

In this section, we present the implementation details, and both offline and online experimental analyses related to UniDex.

### 4.1 Implementation Details

**Datasets.** To rigorously evaluate UniDex, we construct large-scale training and testing datasets derived from real-world search logs of the Kuaishou App, ensuring diversity and representativeness across user behaviors. Specifically, the training dataset comprises approximately 120 million user search sessions, which cover over 3 billion query-video pairs and span multiple stages of the online search pipeline, including recall, pre-ranking, and ranking. Each session contains up to 30 candidate videos, which are uniformly sampled across different stages and divided into multiple tiers to better capture varying degrees of relevance. Hard negative samples are sampled from the earlier stages of the search pipeline, while high-quality positive samples are constructed by integrating scores from the online fine-grained ranking model with explicit user feedback. Each query-video session is represented with rich features, including textual features of both queries and videos, video consumption data, and user feedback signals. For offline evaluation, we select 10 million videos from the video library to form the large-scale candidate pool and extract 50,000 user search sessions from the online system to construct the test set. This design allows us to rigorously assess retrieval and ranking performance under realistic and scalable conditions.

**Details of the UniDex.** The UniTouch encoder is implemented using a 24-layer BERT [7], initialized with internally pre-trained weights to leverage domain-specific knowledge. The model employs a hidden dimension of 1024, and training is conducted with a batch size of 32 using the Adam optimizer. We adopt an initial learning rate of  $2 \times 10^{-5}$ , with a linear warm-up for the first 2,000 steps, followed by a cosine decay schedule. Maximum sequence lengths are set to 32 for queries and 256 for videos, balancing computational efficiency with sufficient contextual coverage. Semantic information is encoded through 19-dimensional 2-bit quantized vectors, with 3 SIDs generated per query and 8 SIDs per video, enabling fine-grained semantic alignment between queries and video candidates. For the UniRank module, the settings of all training-related experimental hyperparameters are consistent with those adopted in UniTouch, and UniRank encodes the semantic information of both the query and the video into four 128-dimensional dense vectors each.

**Evaluation Metrics.** Following Chen et al. [5], to quantitatively assess retrieval and ranking performance, we adopt Recall@300 and Mean Reciprocal Rank (MRR@10) as our primary evaluation metrics. Recall@300 evaluates the ability of the model to retrieve relevant candidates within the top 300 results, while MRR@10 measures the ranking quality by considering the position of the first relevant item, truncated at the top 10 ranks. Formally, for Recall@300, let  $t_i$  denote the number of true positive instances among the top 300 retrieved results for the  $i$ -th query,  $y_i$  denote the total number of positive instances associated with that query, and  $n$  denote the total number of queries. Recall@300 can be written as:  $\text{Recall@300} = \frac{1}{n} \sum_{i=1}^n \frac{t_i}{y_i}$ . For MRR@10, let  $t_i$  denote the rank position of the first relevant item within the top 10 retrieved results

**Table 1: Performance comparison of different models on the large-scale video search dataset.**

Model			Recall@300(%)		MRR@10(%)	
	Touch Module	Rank Module	RS	CS	RS	CS
Sparse Retrievals	Inverted Indexing	BM25 [41]	49.56	46.10	22.21	18.94
		DeepCT [6]	52.05	48.60	23.58	20.42
		SPLADE [9]	54.61	50.74	24.21	22.91
		SPLADE-Max [10]	56.56	51.18	25.04	23.27
Kuaishou Baseline	Touch-Base	Rank-Base	55.33	51.12	27.50	24.92
	UniTouch-24L	Rank-Base	66.06	62.45	31.66	26.10
	Touch-Base	UniRank	56.24	51.73	29.67	25.89
UniDex (Ours)	UniTouch-6L	UniRank	65.21	61.20	32.29	27.13
	UniTouch-12L		68.56	63.02	33.24	28.11
	UniTouch-24L		<b>70.74</b>	<b>65.80</b>	<b>34.06</b>	<b>28.42</b>
Dense Retrievals (For Refer)	DPR [17]		69.57	64.38	34.08	28.11
	ANCE [39]		70.02	65.73	34.56	28.62
	ColBERT [18]		70.98	66.16	34.72	29.10
	TriSampler [42]		<b>73.09</b>	<b>67.75</b>	<b>35.27</b>	<b>29.96</b>

for the  $i$ -th query (with  $t_i = \infty$  if no relevant item appears in the top 10). MRR@10 can be denoted as:  $\text{MRR@10} = \frac{1}{n} \sum_{i=1}^n \frac{1}{t_i}$ . We conduct evaluations on two distinct subsets: the *ranking subset* (RS) and the *click subset* (CS). The RS test set comprises videos recommended to users by the search system and is used to evaluate alignment with system-level preferences. In contrast, the CS test set treats user-clicked items as positives, capturing immediate user interest and reflecting real-world interaction signals. This dual evaluation strategy enables a comprehensive analysis of both systemic relevance and user-centric engagement.

## 4.2 Main Results

To assess the effectiveness of UniDex, we carry out comprehensive experiments on the large-scale video search dataset and compare its performance with existing sparse retrieval methods, dense retrieval models, and the inverted indexing baseline currently used by Kuaishou production (Online Benchmark). All experiments are conducted under consistent training and testing conditions to ensure a fair comparison.

**Compared Methods.** We consider the following three categories of models: (1) *Sparse Retrieval*. The compared sparse retrieval models include BM25 [41] and its neural extensions. DeepCT [6] improves lexical matching by estimating the contextual importance of individual terms. More recent advances, including SPLADE [9] and SPLADE-Max [10], leverage sparsity-inducing regularization combined with lexical expansion to achieve a competitive trade-off between retrieval accuracy and efficiency. (2) *Online Benchmark*. To evaluate the experimental results against the baseline in a real production environment, we compare UniDex with the inverted indexing retrieval framework currently deployed in Kuaishou’s search system. This framework consists of two modules, referred to as *Touch-Base* and *Rank-Base* in our experiments. (3) *Dense Retrieval*.

Finally, we introduce existing dense retrieval methods. It is important to note that dense retrieval fundamentally differs from inverted indexing, yet both serve as crucial recall mechanisms in search systems. We provide a comparison of complex dense retrieval methods as a reference for evaluation. DPR [17] uses a dual-encoder framework to map queries and passages into low-dimensional dense vectors, enabling efficient retrieval via dot-product similarity. ANCE [39] introduces approximate nearest neighbor contrastive learning to address the training bottleneck in dense retrieval. ColBERT [18] enhances query and document representations with a token-level late-interaction mechanism, enabling multi-vector retrieval. TriSampler [42] further refines hard negative sampling strategies.

**Advantages of Unified Semantic Modeling.** On the one hand, UniDex achieves significant growth compared to sparse models. As demonstrated in the *Sparse Retrievals* block of Table 1, even our lightweight model (UniTouch-6L and UniRank) outperforms the strongest sparse methods, achieving an 8.65% improvement in Recall@300 and a 7.25% improvement in MRR@10 on the RS dataset. On the other hand, at its peak performance, UniDex can match the effectiveness of dense retrieval methods. As shown in the *Dense Retrievals* block of Table 1, our top-performing UniDex model (UniTouch-24L and UniRank) lags behind the leading model by 2.35% in Recall@300 and 1.21% in MRR@10 on the RS dataset. This outcome is consistent with our expectations. Dense retrieval and inverted indexing recall represent two distinct retrieval strategies. Dense retrieval is more complex and demands greater computational resources. Additionally, it is more sensitive to increases in retrieval scale, such as the expansion of candidate sets, compared to inverted indexing recall. The above analysis demonstrates that through unified semantic modeling, UniDex significantly enhances the recall capability of inverted indexing retrieval, offering a new approach for this retrieval method.



**Table 2: Comparison of various token-match mechanisms.**

Method	Recall@300(%)		MRR@10(%)	
	RS	CS	RS	CS
UniDex-Max-Sum	34.58	32.19	20.10	17.67
UniDex-Max-Mean	42.41	39.60	23.73	19.83
UniDex-Max-Max	<b>70.74</b>	<b>65.80</b>	<b>34.06</b>	<b>28.42</b>

**Table 3: Results of the influence of learning objectives, where ML represents the Matching Loss, and QR represents the Quantization Regularization.**

Method	Recall@300(%)		MRR@10(%)	
	RS	CS	RS	CS
UniDex	<b>70.74</b>	<b>65.80</b>	<b>34.06</b>	<b>28.42</b>
UniDex w/o ML	67.88	63.05	33.59	27.96
UniDex w/o QR	70.62	65.54	34.01	28.25

**Effectiveness of the Comparison with Online Benchmark.**

We further compare UniDex with Kuaishou’s strongest online product baseline, displayed in the *Kuaishou Baseline* block of Table 1. Let’s take Recall@300 on RS as an example for analysis. Compared to *Touch-Base* and *Rank-Base*, our comprehensive model shows a 15.41% improvement. This significant enhancement in metrics is due to our unified semantic modeling, which improves the semantic understanding capabilities of inverted indexing retrieval and provides a robust alternative to traditional inverted indexing methods. For a more detailed analysis, applying UniTouch alone without modifying the rank module results in a 10.73% increase. However, there’s still a 4.68% gap compared to UniDex. At the same time, replacing only the *Rank-Base* with UniRank yields a gain of 0.91%. This suggests that the UniTouch module is more critical, as it can access more semantically relevant results for the rank module. Furthermore, the unified semantic modeling of UniTouch and UniRank can collaborate to achieve optimal results.

**Expanding the Iteration Space of Inverted Indexing.** As shown in the *UniDex* block of Table 1, we explore various configurations of the UniTouch module (maintain UniRank at a consistent 24 layers) and observe that as the number of model layers increases, the metrics consistently show improvement (we set UniTouch-24L and UniRank as our online experiment setting). This indicates that further iterations of UniDex could yield even greater business benefits. It creates opportunities for continuous iterations of inverted indexing, thereby raising the business potential. Meanwhile, the unified modeling approach of UniDex removes the necessity for extensive manual term expansion and the design of manual term matching relevance methods. This streamlines the iteration process and significantly reduces maintenance costs.

**4.3 Ablation Study**

We perform ablation studies to assess the impact of several key configurations in UniDex on final performance, with particular

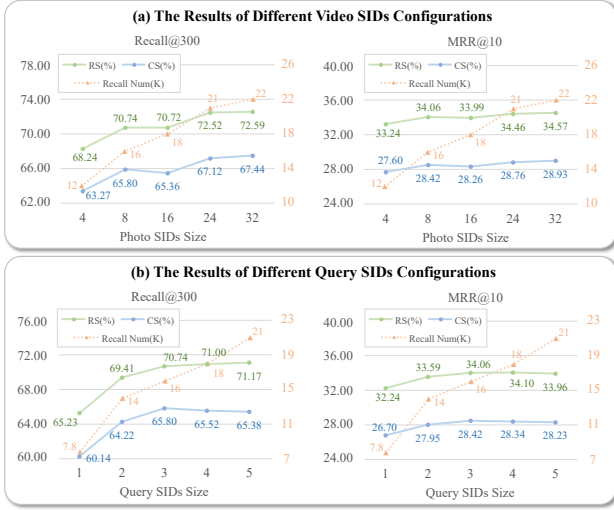
**Table 4: Analysis results of the FSQ codebook space on the experimental effects, where FSQ vectors are 2-bit quantized vectors with dimension  $d_q$ , and the number of SIDs on the query and video sides are 3 and 8, respectively. “\*\*” indicates the setting we adopted.**

Method	$d_q$	Recall@300(%)		MRR@10(%)		Recall Num.
		RS	CS	RS	CS	
UniDex	16	<b>71.60</b>	<b>66.23</b>	<b>34.18</b>	<b>28.47</b>	26K
	18	70.84	65.91	34.03	28.39	19K
	19*	70.74	65.80	34.06	28.42	16K
	20	69.81	64.78	33.75	28.08	13.5K
	22	68.76	63.38	33.48	27.71	10K
	24	66.16	61.92	32.58	27.11	7.6K

emphasis on the token-matching mechanism, training strategy, dimensionality of the FSQ codebook space, and the number of SIDs produced by the query and video.

**Effect of Token-Matching Mechanisms.** The token-matching mechanism of UniDex is crucial for adapting to inverted indexing, as detailed in Section 3.1. We evaluate the impact of various token-matching objectives between query-video pairs on performance, including Max-Sum ( $\text{sim}(q, d) = \sum_i^M \max_{j \in [N]} s(\hat{Q}_i, \hat{D}_j)$ ), Max-Mean ( $\text{sim}(q, d) = \frac{1}{M} \sum_i^M \max_{j \in [N]} s(\hat{Q}_i, \hat{D}_j)$ ), and Max-Max ( $\text{sim}(q, d) = \max_{i \in [M]} \max_{j \in [N]} s(\hat{Q}_i, \hat{D}_j)$ ). The results in Table 2 evident that the Max-Max interaction significantly outperforms the other methods. This is mainly because Max-Max better preserves consistency between the training and retrieval phases of the UniTouch, and aligns well with the inverted indexing paradigm. In inverted indexing, retrieval occurs only when at least one SID from either the query side or the video side matches. In contrast, Max-Sum and Max-Mean require that all SIDs on the query side reach a certain level of similarity with SIDs on the video side to enable retrieval, which deviates from the mechanism of inverted indexing.

**Influence of Learning Objectives.** Given that contrastive loss is the main optimization objective in the retrieval domain and has been extensively explored in prior research [26, 30], our focus here is on examining the effects of *matching loss* and *quantization regularization*. Matching loss (ML) ensures that similar queries and videos receive similar SIDs. When disabled ML, shown in Table 3, it (UniDex w/o ML) leads to a 2.86% decrease in Recall@300 on the RS. The sequence matching of SIDs for semantically similar content significantly affects performance. This is especially crucial in inverted indexing, where SIDs should contain as many identical values as possible for similar semantics. Then, removing quantization regularization (QR) from UniDex (UniDex w/o QR) results in a reduction of a fraction of a percent across the metrics. Although the reduction in metrics is minor, it holds significant importance in engineering practice. Variations in model training precision and the numerical precision of online inference optimizations can lead to uncontrollable changes in the decimal places of float values in the results. During quantization, errors may be incorrectly mapped to 0 or 1 (since we utilize binary quantization, detailed in Section



**Figure 3: The results of configurations with different SIDs quantities on RS and CS test set. (a) Comparison results of different SIDs quantities for Video, and (b) Comparison results of different SIDs quantities for Query.**

3.1), negatively impacting online applications. Thus, we design this regularization to keep the model away from the decision boundary of 0.5, mitigating uncontrollable changes caused by device precision issues.

**Analysis on FSQ Codebook Space.** The FSQ’s quantized codebook space significantly affects both the effectiveness of semantic models and the quality of inverted indexing. To clearly explain these effects, we conduct a series of controlled experiments, as detailed below. To ensure that our experiments are comparable, we keep the number of SIDs constant for both queries and videos in all tests. Specifically, we create six experimental groups with FSQ codebook spaces ranging from  $2^{16}$  to  $2^{24}$  (i.e.,  $d_q$  from 16 to 24). As reflected in Table 4, there is a clear trend: the Recall@300 metric gradually decreases as the codebook space grows. This is primarily due to the rapid expansion of the codebook’s quantization dimensions, which causes the semantic space to become more dispersed. Such dispersion, in turn, reduces the recall rate during the retrieval phase to some extent. Conversely, enhancing recall is associated with an increase in the number of retrieval results. For instance, when  $d_q = 16$ , the recall results peak at 26K. This growth in recall results significantly adds to the system’s overall burden, leading to higher latency and increased consumption of computational resources. To achieve a better balance between retrieval rate and system efficiency, we set  $d_q = 19$ . At this dimension, the decrease in Recall@30 and MRR@10 compared to  $d_q = 16$  is minimal, while the number of retrieved items is significantly reduced by 10K.

**Impact of SIDs Count on Query/Video.** In the UniDex framework, the number of SIDs associated with queries and videos significantly impacts both model performance and online retrieval efficiency. On the query side, the number of SIDs directly impacts the scope of recall. Since queries usually consist of short, semantically clear text, we empirically investigate configurations with 1

**Table 5: The results of UniDex in online A/B test compared to the production baseline. We consider the user’s satisfaction (Sat.) and resource costs (Cost).**

UniDex	Sat.	CTR ↑	VPD ↑	LPC ↑	MRS ↑
		+0.185%	+0.287%	+0.352%	+0.346%
Cost	Core ↓	Memory ↓	Latency ↓		
	-20550	-37TB	-25%		

to 5 SIDs. On the contrary, on the video side, the text is generally longer and contains rich semantic information, prompting us to explore settings ranging from 4 to 32. As shown in Figure 3 (a), we start by examining the impact of Video SIDs size on the results. Clearly, as the size of SIDs increases, both Recall@300 and MRR@10 improve. Additionally, it can be observed that increasing the size from 8 to 24 (a threefold increase) results in an approximate 2% improvement in Recall@300 and about a 0.4% improvement in MRR. However, this change also results in an increase of 5K in the number of retrieved items, which can significantly raise inference resource consumption in the ranking module. Therefore, we select 8 as our final setting. Subsequently, Figure 3 (b) demonstrates how the query-side SIDs size affects the outcomes. We notice that when the size surpasses 3, further improvements in metrics are minimal. This is likely because queries are typically short, and their semantic information can be effectively captured by just 3 SIDs.

#### 4.4 Online Testing

We further evaluate UniDex in real-world settings by deploying it within Kuaishou’s short-video search system. A 5-day online A/B test was conducted, with both the experimental and control groups randomly assigned 10% of the actual search traffic.

**Improving User Satisfaction with Search Results.** We focus on four online metrics: page click-through rates (CTR), video playback duration (VPD), long play count (LPC), and the mean relevance score of the top 4 results (MRS). As shown in *Sat.* block of Table 5, UniDex outperforms the advanced production baseline in terms of CTR, VPD, and LPC. This indicates that users are clicking on more videos and spending more time watching relevant content, significantly improving their experience. Additionally, MRS serves as a key monitoring metric, reflecting changes in the fine-grained relevance scores of the top four results shown to users after implementing our model. A MRS improvement of +0.346% suggests that UniDex can deliver more relevant videos compared to traditional term-based inverted indexing systems. This enhancement not only improves the system’s overall relevance representation capability but also positively influences the search system’s long-term evolution.

**Minimizing System Resource Expenditure.** UniDex has replaced the traditional term-based inverted indexing used online. In the conventional pipeline, manual term matching and relevance computation consume significant resources. With the full deployment of UniDex, a substantial amount of inference resources has been saved. As shown in the *Cost* block of Table 5, it saves 20550



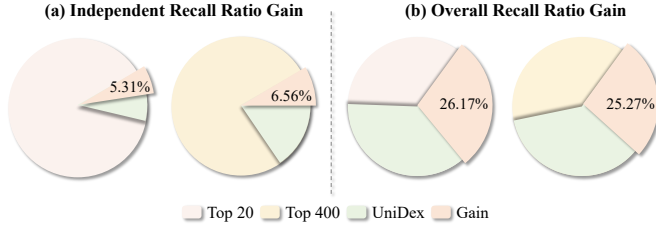


Figure 4: The variations in recall ratio following the implementation of our comprehensive model, UniDex.

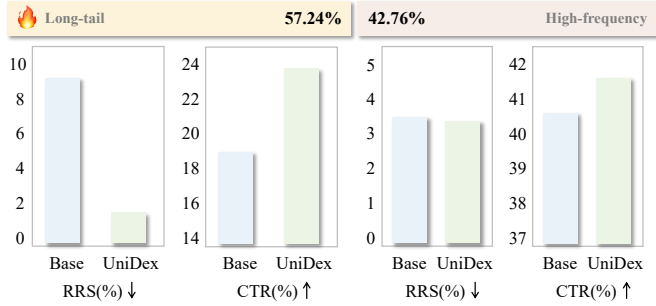


Figure 5: Data analysis of Long-tail and High-frequency Queries based on improvements in Independent Recall Ratio for the Top 20. The top percentage reflects the contribution levels of queries at various frequencies. RRS denotes the Relevant Result Scarcity rate, while CTR refers to the Click-through Rate of the results.

cores and 37TB of storage resources, resulting in a 25% improvement in system response time. The reduction in computational and storage demands leads to lower operational costs, allowing for more efficient allocation of resources. This efficiency enables the system to handle higher volumes of concurrent user requests, effectively increasing its scalability and robustness. Furthermore, the 25% improvement in system response time enhances the user experience by providing faster access to relevant information, which can boost user engagement and satisfaction.

#### 4.5 Further Analysis

To better understand the advantages of UniDex over traditional term-based inverted indexing, we conduct a more in-depth analysis. **Enhancing Retrieval Effectiveness.** The recall proportion during the fine-grained ranking phase reflects the retrieval effectiveness of the model. We examine two types of recall ratios: the overall recall ratio, which includes any video retrieved through inverted indexing, and the independent recall ratio, which counts a video only if it is retrieved exclusively through this method. As illustrated in Figure 4, UniDex demonstrates enhancements in both recall ratios. Notably, there is an increase of over 25% in the overall recall ratio. This significant improvement can be attributed to UniDex’s ability to greatly enhance the generalization capability of inverted indexing, allowing it to retrieve more videos that overlap with semantic retrieval. In contrast, earlier term-based inverted indexing techniques were limited in capturing deeper search semantics. Furthermore, when



Figure 6: The visualization of the UniDex retrieval process.

we expand the number of ranked videos from the top 20 to the top 400, UniDex’s independent recall rate increased by an additional 1.25%. This indicates that integrating unified semantic modeling with inverted indexing provides a powerful complementary enhancement to the existing search pipeline, enabling the retrieval of more videos that were previously unretrievable.

**Advancing Long-tail Semantic Characterization.** We conduct a detailed analysis of the sources of improvements in independent recall ratios, categorizing them by queries of varying frequencies. As shown in Figure 5, long-tail queries contribute 14.48% (57.24% vs 42.76%) more than high-frequency queries. We also evaluate the Relevant Result Scarcity (RRS) rate and Click-through Rate (CTR) separately for both long-tail and high-frequency queries. A lower RRS score indicates better relevance of the retrieval results, while a higher CTR score suggests greater user satisfaction. Our findings reveal that UniDex significantly lowers the RRS and boosts the CTR for long-tail queries, highlighting its superior semantic representation capabilities compared to term-based inverted indexing. While UniDex demonstrates improvements over the baseline in representing high-frequency queries, these advancements are less pronounced compared to the significant gains observed for long-tail queries. This discrepancy may be due to the fact that term-based inverted indexing methods excel at capturing the semantics of high-frequency queries. However, long-tail queries demand a deeper textual understanding, which renders simple term-matching methods inadequate. By introducing unified semantic modeling, our UniDex model enhances the semantic representation capabilities of inverted indexing, resulting in notable improvements for long-tail queries. This enhanced semantic capability for long-tail queries addresses more nuanced search needs, thereby playing a crucial role for the platform.

**Balancing Generalization and Matching.** To provide a clearer understanding of how UniDex operates, we perform a case analysis, with the results presented in Figure 6. When the user inputs the query "Labubu 5th Generation", UniDex first activates the UniTouch module, converting the query into a series of SIDs (e.g., "61011", "61267", "62829"). These SIDs are then fed into a pre-stored inverted

indexing that links SIDs to corresponding video lists, resulting in a set of candidate videos. Finally, the UniRank module is employed to rank the videos within this candidate set. In this case, most videos linked to SID-61011 are not retrieved by the traditional inverted indexing because the term “Labubu” does not appear in the video content. However, UniDex can utilize unified semantic modeling to understand the various linguistic expressions of the entity “Labubu”, allowing it to retrieve videos with text descriptions that include the Chinese name for Labubu. The enhancement of semantic generalization raises the upper limit of inverted index retrieval, creating new opportunities for business innovation and iteration. Meanwhile, we also observe that SID-62829 employs a term omission strategy similar to term-based inverted indexing. Most of the videos associated with it are related to Labubu but lack crucial information such as “5th”. This capability becomes particularly crucial when the candidate set contains only a few items related to the query, as it allows for the retrieval of as many vaguely relevant videos as possible. Consequently, UniDex retains the term-matching advantages of term-based inverted indexing while enhancing generalization.

## 5 Conclusion

The paper presents **UniDex**, a novel model-based approach that fundamentally transforms inverted indexing from a traditional term-based paradigm to a unified semantic modeling framework. By effectively replacing complex manual term-matching strategies with a streamlined semantic model, UniDex significantly reduces maintenance overhead and enhances retrieval capabilities through improved semantic generalization. Our successful deployment of UniDex in Kuaishou’s short-video search systems, supported by extensive online A/B testing, demonstrates its practical effectiveness and scalability, serving millions of users efficiently. UniDex not only addresses the limitations of existing approaches but also opens new avenues for future innovations in search technology.

## References

- [1] Yang Bai, Xiaoguang Li, Gang Wang, Chaoliang Zhang, Lifeng Shang, Jun Xu, Zhaowei Wang, Fangshan Wang, and Qun Liu. 2020. SparTerm: Learning Term-based Sparse Representation for Fast Text Retrieval. *arXiv preprint arXiv:2010.00768* (2020).
- [2] Yoshua Bengio, Nicholas Léonard, and Aaron Courville. 2013. Estimating or Propagating Gradients Through Stochastic Neurons for Conditional Computation. *arXiv preprint arXiv:1308.3432* (2013).
- [3] Sebastian Bruch, Franco Maria Nardini, Cosimo Rulli, and Rossano Venturini. 2024. Efficient Inverted Indexes for Approximate Retrieval over Learned Sparse Representations. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 152–162.
- [4] Haonan Chen, Zhicheng Dou, Kelong Mao, Jiongnan Liu, and Ziliang Zhao. 2024. Generalizing Conversational Dense Retrieval via LLM-Cognition Data Augmentation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 2700–2718.
- [5] Jiahui Chen, Xiaozhe Jiang, Zhibo Wang, Quanzhi Zhu, Junyao Zhao, Feng Hu, Kang Pan, Ao Xie, Maohua Pei, Zhiheng Qin, et al. 2025. UniSearch: Rethinking Search System with a Unified Generative Architecture. *arXiv preprint arXiv:2509.06887* (2025).
- [6] Zhuyun Dai and Jamie Callan. 2019. Deeper Text Understanding for IR with Contextual Neural Language Modeling. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 985–988.
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. 4171–4186.
- [8] Fethi Fkih. 2023. Enhancing item-based collaborative filtering by users’ similarities injection and low-quality data handling. *Data & Knowledge Engineering* 144 (2023), 102126.
- [9] Thibault Formal, Benjamin Piwowarski, and Stéphane Clinchant. 2021. SPLADE: Sparse Lexical and Expansion Model for First Stage Ranking. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2288–2292.
- [10] Thibault Formal, Benjamin Piwowarski, and Stéphane Clinchant. 2021. SPLADE: Sparse Lexical and Expansion Model for First Stage Ranking. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2288–2292.
- [11] Chongming Gao, Kexin Huang, Jiawei Chen, Yuan Zhang, Biao Li, Peng Jiang, Shiqi Wang, Zhong Zhang, and Xiangnan He. 2023. Alleviating Matthew Effect of Offline Reinforcement Learning in Interactive Recommendation. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 238–248.
- [12] Luyu Gao and Jamie Callan. 2022. Unsupervised Corpus Aware Language Model Pre-training for Dense Passage Retrieval. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 2843–2853.
- [13] Luyu Gao, Zhuyun Dai, and Jamie Callan. 2021. COIL: Revisit Exact Lexical Match in Information Retrieval with Contextualized Inverted List. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 3030–3042.
- [14] Luyu Gao, Zhuyun Dai, Tongfei Chen, Zhen Fan, Benjamin Van Durme, and Jamie Callan. 2021. Complement Lexical Retrieval Model with Semantic Residual Embeddings. In *European Conference on Information Retrieval*. Springer, 146–160.
- [15] Yue He, Yancheng Dong, Peng Cui, Yuhang Jiao, Xiaowei Wang, Ji Liu, and Philip S Yu. 2021. Purify and Generate: Learning Faithful Item-to-Item Graph from Noisy User-Item Interaction Behaviors. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 3002–3010.
- [16] Xiaozhe Jiang, Yaobo Liang, Weizhu Chen, and Nan Duan. 2022. XLM-K: Improving Cross-Lingual Language Model Pre-training with Multilingual Knowledge. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 36. 10840–10848.
- [17] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 6769–6781.
- [18] Omar Khattab and Matei Zaharia. 2020. ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 39–48.
- [19] Carlos Lassance, Hervé Déjean, Thibault Formal, and Stéphane Clinchant. 2024. SPLADE-v3: New baselines for SPLADE. *arXiv preprint arXiv:2403.06789* (2024).
- [20] Junghyup Lee, Dohyung Kim, and Bumsub Ham. 2021. Network Quantization with Element-wise Gradient Scaling. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Vol. 1. 6444–6453.
- [21] Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. 2019. Latent Retrieval for Weakly Supervised Open Domain Question Answering. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 6086–6096.
- [22] Jimmy Lin and Xueguang Ma. 2021. A Few Brief Notes on DeepImpact, COLL, and a Conceptual Framework for Information Retrieval Techniques. *arXiv preprint arXiv:2106.14807* (2021).
- [23] Zheng Liu, Chaofan Li, Shitao Xiao, Yingxia Shao, and Defu Lian. 2024. Llama2Vec: Unsupervised Adaptation of Large Language Models for Dense Retrieval. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 3490–3500.
- [24] Yi Luan, Jacob Eisenstein, Kristina Toutanova, and Michael Collins. 2021. Sparse, Dense, and Attentional Representations for Text Retrieval Open Access. *Transactions of the Association for Computational Linguistics* 9 (2021), 329–345.
- [25] Sean MacAvaney, Franco Maria Nardini, Raffaele Perego, Nicola Tonellotto, Nazli Goharian, and Ophir Frieder. 2020. Expansion via Prediction of Importance with Contextualization. *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval* (Jul 2020).
- [26] Alessandro Magnani, Feng Liu, Suthesh Chaidaroon, Sachin Yadav, Praveen Reddy Suram, Ajit Puthenpussery, Sijie Chen, Min Xie, Anirudh Kashi, Tony Lee, et al. 2022. Semantic Retrieval at Walmart. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 3495–3503.
- [27] Fabian Mentzer, David Minnen, Eirikur Agustsson, and Michael Tschannen. 2024. Finite Scalar Quantization: VQ-VAE Made Simple. *The Twelfth International Conference on Learning Representations*.
- [28] Rodrigo Nogueira, Jimmy Lin, and AI Epistemic. 2019. From doc2query to docTTTTTquery. *Online preprint* 6 (2019).
- [29] Rodrigo Nogueira, Wei Yang, Jimmy Lin, and Kyunghyun Cho. 2019. Document Expansion by Query Prediction. *arXiv preprint arXiv:1904.08375* (2019).
- [30] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation Learning with Contrastive Predictive Coding. *arXiv preprint arXiv:1807.03748* (2018).

- [31] Liang Pang, Yanyan Lan, Jiafeng Guo, Jun Xu, Jingfang Xu, and Xueqi Cheng. 2017. DeepRank: A New Deep Architecture for Relevance Ranking in Information Retrieval. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*. 257–266.
- [32] Xiaohua Qi, Renda Li, Long Peng, Qiang Ling, Jun Yu, Ziyi Chen, Peng Chang, Mei Han, and Jing Xiao. 2025. Data-free Knowledge Distillation with Diffusion Models. *arXiv preprint arXiv:2504.00870* (2025).
- [33] Yingqi Qu, Yuchen Ding, Jing Liu, Kai Liu, Ruiyang Ren, Wayne Xin Zhao, Daxiang Dong, Hua Wu, and Haifeng Wang. 2021. RocketQA: An Optimized Training Approach to Dense Passage Retrieval for Open-Domain Question Answering. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 5835–5847.
- [34] Ruiyang Ren, Shangwen Lv, Yingqi Qu, Jing Liu, Wayne Xin Zhao, QiaoQiao She, Hua Wu, Haifeng Wang, and Ji-Rong Wen. 2021. PAIR: Leveraging Passage-Centric Similarity Relation for Improving Dense Passage Retrieval. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. 2173–2183.
- [35] Thomas Roelleke and Jun Wang. 2008. TF-IDF uncovered: a study of theories and probabilities. In *Proceedings of the 31st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. 435–442.
- [36] Keshav Santhanam, Omar Khattab, Jon Saad-Falcon, Christopher Potts, and Matei Zaharia. 2022. ColBERTv2: Effective and Efficient Retrieval via Lightweight Late Interaction. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 3715–3734.
- [37] Teng Shi, Jun Xu, Xiao Zhang, Xiaoxue Zang, Kai Zheng, Yang Song, and Han Li. 2025. Retrieval Augmented Generation with Collaborative Filtering for Personalized Text Generation. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1294–1304.
- [38] Zhibo Wang, Xiaozhe Jiang, Zhiheng Qin, and Enyun Yu. 2025. Personalized Query Auto-Completion for Long and Short-Term Interests with Adaptive Detoxification Generation. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*. 5018–5028.
- [39] Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul Bennett, Junaid Ahmed, and Arnold Overwijk. 2021. Approximate Nearest Neighbor Negative Contrastive Learning for Dense Text Retrieval. In *International Conference on Learning Representations*.
- [40] Ikuya Yamada, Akari Asai, and Hannaneh Hajishirzi. 2021. Efficient Passage Retrieval with Hashing for Open-domain Question Answering. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*. 979–986.
- [41] Peilin Yang, Hui Fang, and Jimmy Lin. 2017. Anserini: Enabling the Use of Lucene for Information Retrieval Research. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1253–1256.
- [42] Zhen Yang, Zhou Shao, Yuxiao Dong, and Jie Tang. 2024. TriSampler: A Better Negative Sampling Principle for Dense Retrieval. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 38. 9269–9277.
- [43] Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Jiafeng Guo, Min Zhang, and Shaoping Ma. 2021. Jointly Optimizing Query Encoder and Product Quantization to Improve Retrieval Performance. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 2487–2496.
- [44] Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Jiafeng Guo, Min Zhang, and Shaoping Ma. 2021. Optimizing Dense Retrieval Model Training with Hard Negatives. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1503–1512.
- [45] Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Jiafeng Guo, Min Zhang, and Shaoping Ma. 2022. Learning Discrete Representations via Constrained Clustering for Effective and Efficient Dense Retrieval. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*. 1328–1336.
- [46] Shunyu Zhang, Yaobo Liang, Ming Gong, Daxin Jiang, and Nan Duan. 2022. Multi-View Document Representation Learning for Open-Domain Dense Retrieval. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 5990–6000.
- [47] Tiancheng Zhao, Xiaopeng Lu, and Kyusong Lee. 2021. SPARTA: Efficient Open-Domain Question Answering via Sparse Transformer Matching Retrieval. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 565–575.
- [48] Xu Zhao, Yi Ren, Ying Du, Shenzheng Zhang, and Nian Wang. 2022. Improving Item Cold-start Recommendation via Model-agnostic Conditional Variational Autoencoder. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 2595–2600.