

# ROSE: Robust Caches for Amazon Product Search

Chen Luo, Vihan Lakshman, Anshumali Shrivastava, Tianyu Cao, Sreyashi Nag, Rahul Goutam, Hanqing Lu, Yiwei Song, Bing Yin  
Amazon Search  
Palo Alto, California, USA

## ABSTRACT

Product search engines like Amazon Search often use caches to improve the customer user experience; caches can improve both the system’s latency as well as search quality. However, as search traffic increases over time, the cache’s ever-growing size can diminish the overall system performance. Furthermore, typos, misspellings, and redundancy widely witnessed in real-world product search queries can cause unnecessary cache misses, reducing the cache’s utility. In this paper, we introduce **ROSE**, a **ROBuSt** cache, a system that is tolerant to misspellings and typos while retaining the look-up cost of traditional caches. The core component of **ROSE** is a randomized hashing schema that makes **ROSE** able to index and retrieve an arbitrarily large set of queries with constant memory and constant time. **ROSE** is also robust to any query intent, typos, and grammatical errors with theoretical guarantees. Extensive experiments on real-world datasets demonstrate the effectiveness and efficiency of **ROSE**. **ROSE** is deployed in the Amazon Search Engine and produced a significant improvement over the existing solutions across several key business metrics.

## CCS CONCEPTS

• **Information systems** → **Query log analysis**; *Query intent*; *Query reformulation*.

## KEYWORDS

Product Search, Robust Cache, Data Mining

## 1 INTRODUCTION

Online shopping has become an essential part of consumers’ daily lives in recent years and has seen a dramatic increase in demand during the ongoing COVID-19 global pandemic. As a critical component of an e-commerce website, the product search engine connects the customer intent with the product selections. Improving the product search engine’s performance is critical to a better shopping experience. Two key factors impact the search engine’s performance: (1) The response time to a customer request and (2) Providing high-quality results that match the customers’ intent.

User studies show that slow responses cause perceived interruptions to the shopping experience and even site abandonment.

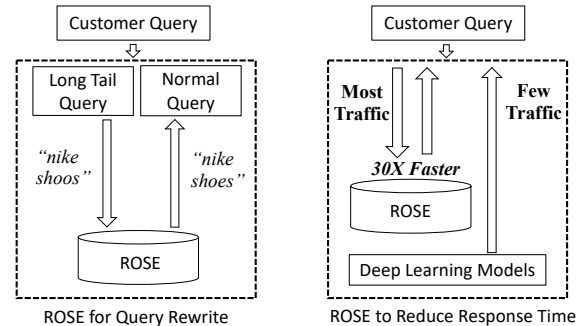
Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM 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).

ISIR-eCom 2022 @ WSDM-2022, Feb 25, 2022, Phoenix, AZ

© 2022 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00

<https://doi.org/10.1145/1122445.1122456>



**Figure 1: ROSE helps improve the search quality and system performance of product search engine. With ROSE, most of the search traffic is covered with single digit milliseconds latency. ROSE also improves search quality by mapping long-tailed queries to normal queries in near constant time.**

The response time is also a key factor for the product search engine’s throughput planning. Modern product search engines are usually composed of different expensive machine learning models [1, 8, 9, 13, 20, 28, 29, 31], such as relevance matching models [18], ranking models [2], and query annotation models [27]. Serving the entirety of search traffic through expensive deep learning models is prohibitive in real-world product search engines due to latency limitations, and cost considerations [14]. Thus, instead of serving all queries through these expensive deep learning models, a more practical solution is to serve frequent queries from a cache.

However, traditional caches suffer from the trade-off between the cache miss rate and the cache size. Having a small cache size will lead to a high cache miss rate. On the other hand, as product search engines scale, the set of frequently occurring queries becomes prohibitively large, and grows due to morphological variants of queries with the same intent. For instance, “Nike shoes”, “Nike shoe”, and “Nike’s shoe” may all be cached queries due to their frequency. These queries all share the same intent, and they artificially inflate the cache size and diminish performance. Therefore, designing a robust cache that is invariant to typos and morphological differences is critical for scaling real-world search services since it enables increasing the cache hit rate without correspondingly increasing the latency and memory footprint.

Moreover, one key issue that hurts the quality of search results is the presence of low-performing queries, which are queries for which the search engine fails to return high-quality results. Analyses show that most of these failure cases are due to typographical errors [24]. These low-performing queries are usually lexically or semantically similar to some frequently searched, well-performing queries that produce satisfactory results. Thus, if we could map these low-performing queries to a frequently searched query with

the same intent via a robust caching mechanism, we would be able to improve search quality. Furthermore, this query mapping process would also reduce latency since product search engines typically cache these frequently issued queries and their corresponding behavioral information for faster serving, as in Fig 1.

To solve these challenges, we propose **ROSE** to cache well-performing or frequent queries to improve the response time and search quality of the product search engine. The core component of **ROSE** is a randomized hashing structure that indexes the query set while preserving the lexical or semantic information. Specifically, our paper includes the following contributions:

- **Operational System:** We introduce **ROSE**, a comprehensive end-to-end solution for caching queries for product search. **ROSE** can index and perform look-ups on web-scale data in constant time and constant memory and is faster than other alternatives by orders of magnitude.
- **Technical Novelty:** We invented a system that combines multiple powerful randomized algorithmic techniques, including locality sensitive hashing, reservoir sampling, and count-based  $k$ -selection, in a novel way that together allow us to scale up **ROSE** to massive query sets while maintaining constant-time retrieval.
- **Real-World Impact:** We deployed **ROSE** in the Amazon product search engine, showing improvements in system performance and business metrics when compared to the existing solution.

## 2 BACKGROUND AND RELATED WORK

We provide some background knowledge and formally define the problem of indexing and retrieval for robust caches.

### 2.1 Robust Caches

Caching previously seen data is a central component of many latency-critical applications such as search engines and databases. A *cache-hit* happens when an incoming query matches one of the cached queries. Many well-established methods such as hash tables and Bloom filters [16] focus on detecting exact key matches. However, these methods are not robust to typos, lexical variants, and semantic variants. For instance, the queries “Nike sheos”, “Nike shoes”, “Nike shoe”, and “Nike’s shoe” are variants of the same query “nike shoes”. These variants will cause a cache miss under traditional exact-match caching methods. Historically, making a caching method robust required expensive algorithms to compute string similarity distances between a query key and the cached keys. In this paper, we propose **ROSE**, a robust caching framework that handles variants of input keys with near-constant time and near-constant memory.

### 2.2 Locality-Sensitive Hashing

**2.2.1 LSH Definition.** We use locality-sensitive hashing (LSH), a randomized hashing method, as the core algorithm in our work. LSH is a family of hash functions such that a function uniformly sampled from the family has the property that, under the hash mapping, similar points have a high probability of having the same hash value [32]. Taking  $\mathcal{H}$  to be a family of hash functions mapping  $\mathbb{R}^d$  to a discrete set  $\{0, \dots, U - 1\}$  we can formally the notion of locality-sensitive hashing.

**Definition 2.1. LSH Family** A family  $\mathcal{H}$  is called  $(R, cR, p, q)$ -sensitive if any two points  $x, y \in \mathbb{R}^d$  and  $h$  chosen uniformly from  $\mathcal{H}$  satisfies the following property:

- if  $\text{sim}(x, y) \geq R$ , then  $\Pr_{h \in \mathcal{H}}[h(x) = h(y)] \geq p$ .
- if  $\text{sim}(x, y) \leq cR$ , then  $\Pr_{h \in \mathcal{H}}[h(x) = h(y)] \leq q$ .

A collision occurs when the hash values for two points are equal:  $h(x) = h(y)$ . The collision probability is proportional to some monotonic function of similarity between the two points, i.e.,  $\Pr[h(x) = h(y)] \propto f(\text{sim}(x, y))$ , where  $\text{sim}(x, y)$  is the similarity under consideration and  $f$  is a monotonically increasing function. In other words, similar items are more likely to collide with each other under an LSH mapping [23].

**2.2.2 Minwise Hashing.** Minwise hashing (MinHash) is the LSH for Jaccard similarity [4]. The minwise hashing family of functions applies a random permutation  $\pi$  on a given set  $S$ , and stores only the minimum value after the permutation mapping. Given two sets,  $S_1$  and  $S_2$ , the probability of these sets having the same MinHash value is the Jaccard similarity between the two sets:  $\Pr[\min \pi(S_1) = \min \pi(S_2)] = \frac{|S_1 \cap S_2|}{|S_1 \cup S_2|}$ .

In this work, we also leverage recent algorithmic advances in computing weighted MinHashes [7]. Weighted MinHashing allows us to take the importance of items in set into account. For instance, with product-type preserving hashing, we would like to assign a higher weight to the tokens corresponding to the product type. Formally, as in [6], we define the weighted Jaccard similarity between two non-negative  $n$ -dimensional real vectors  $x$  and  $y$  as  $J(x, y) = \frac{\sum_i \min(x_i, y_i)}{\sum_i \max(x_i, y_i)}$ . As in the unweighted case, the probability of two weighted sets having the same weighted MinHash value is precisely the weighted Jaccard similarity.

**2.2.3 Densified One Permutation Hashing.** Another critical ingredient in allowing us to build performant robust caches is the recent progress in efficiently computing minhashes. In particular, we utilize the densified one permutation hashing (DOPH) algorithm introduced by Shrivastava [22]. DOPH is ideal for ultra-high dimensional and sparse datasets and, in an improvement over prior minhashing techniques, can generate multiple hashes in one pass over the data. Ultimately, DOPH allows us to compute  $k$  minhashes of a sparse vector with  $d$  nonzero entries in time  $O(d + k)$  as opposed to the  $O(dk)$  guarantee offered by classical algorithms.

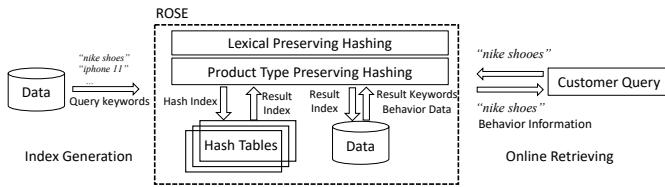
### 2.3 Problem Formulation

In this subsection, we formally define the indexing and retrieval tasks for a robust cache.

The indexing step is defined as follows: Given a set of queries to be cached  $\mathcal{D}$ , the indexing step constructs a cache structure  $C$  that preserves the similarity information of input queries such that we can map an unseen query  $q \notin \mathcal{D}$  to the most similar item in the cache during the retrieval phase.

Thus, given the cache structure  $C$ , we define the retrieval step as follows: Given an input query  $q$ , we are interested in efficiently computing  $X^* = \arg \max_{x \in \mathcal{D}} S(q, x)$ , where  $S(x, y)$  is the similarity between two items  $x$  and  $y$ .

Large-scale product search engines such as Amazon Search usually cache hundreds of millions of queries. Furthermore, the latency



**Figure 2: The overall framework of ROSE. ROSE contains two phases: (1) Cache Index Generation: generating the robust index using the input queries. (2) Online Retrieval: mapping the input query to one of the queries in the cache.**

threshold for cache is extremely low, milliseconds. Keeping these factors in mind, the indexing and retrieval steps for our robust cache are constrained by the following two conditions:

- *Constant Retrieval Time*: At run-time, the retrieval step must finish in constant time regardless of the cache size.
- *Constant Memory Usage*: The memory usage of the cache does not increase with the number of cached queries.

The similarity metric  $S(x, y)$  is a critical factor of robust caches. This paper shall focus on two intent similarity measures widely used in product search engines: (1) Lexical Preserving Caching and (2) Product Type Preserving Caching. We note in passing that one can fit any other intent similarity measures into our framework.

**Lexical Preserving Caching**: A lexical preserving cache maps a query to an existing set of queries while preserve the query’s lexical information. For example, suppose we have a query “red nike shoos”. This query has a typo so we assume that it is not cached by the search engine. In our cache, we may have the query “red nike shoes” that is lexically similar with “red nike shoes”. The problem of lexical preserving query mapping is to map “red nike shoos” to “red nike shoes” and perform the corresponding cache lookup.

**Product Type Preserving Caching**: Product type is the most important attribute of queries to a product search engine. [18]. For example, the product type for query “red nike shoes” is “SHOES.” Product type preserving caching maps the input query to a candidate query with the same product type.

### 3 ROSE: ROBUST CACHE VIA RANDOMIZED HASHING

In this section, we describe **ROSE**, a robust cache for queries via randomized hashing. **ROSE** contains two phases, Index Generation and Online Retrieval, as illustrated in Fig 2. We first introduce these two phases, followed by a theoretical analysis of **ROSE** in terms of both time and memory complexity.

#### 3.1 ROSE Index Generation

We design the index generation process of **ROSE** under two requirements. First, the cache needs to capture the query similarity, meaning that cache needs to take the similarity of queries into account when performing look-ups to be robust to typos and semantic variance. Secondly, due to the large-scale indexing space of real product search engines, the cache size needs to avoid scaling with the volume of queries.

To capture the textual similarity information, we use locality-sensitive hashing (LSH) [10] for the index generation phase. LSH

generates signatures for input data under a certain similarity measure. The signatures generated by LSH capture the similarity information between queries such that similar queries have a high probability of having the same hashing signature and thus colliding. Since LSH is a randomized procedure, we boost the probability of hashing similar queries together by maintaining  $L$  independent hash tables for our index. This work shall focus on two hashing strategies: lexical preserving hashing and product type preserving hashing. We will introduce the details of these two hash functions in Section 3.3 and Section 3.4, respectively.

However, under the locality-sensitive hashing framework, the size of the hash tables increases linearly with the volume of data [23] which leads to an explosion in memory footprint when working with web-scale data. To solve this problem, inspired by the work in [26], we use a reservoir sampling strategy to fix our cache’s memory usage and preserve the data’s similarity information.

The reservoir sampling algorithm [25] processes a stream of  $m$  numbers and generates  $R$  uniform samples by only using an array of size  $R$ , where  $R \ll m$ . Moreover, reservoir sampling only needs one pass over the data, and does not increase the computational complexity of the index generation process. We will provide a theoretical analysis of this sampling strategy as applied in our caching framework in Section 3.5.

Algorithm 1 summarizes the offline indexing algorithm of **ROSE**. In Algorithm 1,  $B$  denotes the number of buckets in the hash table,  $I_q$  denotes the index generated by the hash function. The function  $Rand(0, B)$  generates a random number between 0 and  $B$ .

---

#### Algorithm 1 ROSE Indexing Generation

---

```

1: procedure ROSE INDEX GENERATION
2:   Input: number of hash tables  $L$ , number of hashes  $K$ 
3:   Initialize the hash functions  $H$ , array of hash tables  $T$ .
4:   for each query  $q$  in  $D$  do
5:      $I_q = H(q, K, L)$ .
6:     if  $T[I_q]$  is full then
7:        $R = Rand(0, B)$ .
8:       if  $R < B$  then
9:          $T[I_q] = q$ 
10:      end if
11:      Continue
12:    end if
13:     $T[I_q] = q$ 
14:  end for
15:  return  $T$ 
16: end procedure

```

---

#### 3.2 ROSE Online Retrieval

Given a search query, we perform a robust cache lookup by first computing the LSH signature of this query and looking up the corresponding bucket in the hash tables. We then rank the similarity of the cached queries within the bucket to the new search and return the top result. However, under the standard LSH schema [16], we still have to calculate the pairwise similarities inside the bucket to retrieve the top result, which can be expensive, especially since product search engines typically maintain strict latency budgets.

To avoid this expensive pairwise similarity computation, we use the strategy of count-based  $k$ -selection inspired by [15]. Across the  $L$  different hash tables, we observe that the cached entries with the greatest number of collisions with the new query are more similar to the query text. This observation allows us to estimate the actual ranking in an unbiased manner. We count each data point’s frequency of occurrence in the aggregated reservoirs and rank all the data points based on the frequency. By using this strategy, the online retrieval process runs in constant time, as shown in Section 3.5.

### 3.3 Lexical Preserving Hashing

Our goal for lexical preserving hashing is to design a hash function that preserves the lexical similarity among input queries. To achieve this in product search, we use the Jaccard similarity to measure the similarity between two queries, defined as the ratio of character spans that two query keywords share, and use minhash [3] as the corresponding LSH scheme.

Given a query  $Q$  of  $n$  characters and  $m$  words, we slice these keywords into a set of subsequences consisting of character-level sequences and word unigrams, denoted by  $\mathcal{S}(Q) = \{c_i\}_{i=1}^n \cup \{c_i c_{i+1}\}_{i=1}^{n-1} \dots \cup \{w_i\}_{i=1}^m$ , where  $c_i$  and  $w_i$  denote the  $i$ -th character and word of the query, respectively. The length of the character subsequence is a hyper-parameter. We find that a subsequence length of 3 gave us the best results. We then use the recent advances in densified one permutation hashing (DOPH) [22] to compute the minhash signatures of  $\mathcal{S}(Q)$  efficiently.

### 3.4 Product Type Preserving Hashing

In a product search engine, understanding the product type information of a query is crucial to showing relevant results that match the customer’s intent and avoid, for instance, returning dishwasher accessories in response to a search for dishwashers. Thus, when performing a cache lookup, it is critical that we map the original query to one that preserves the original product type intent.

To preserve the product information, we add weights to product type tokens in the query. The product type tokens are extracted by a production NER model [30]. We use the same process as lexical preserving hashing to generate the token set  $\mathcal{S}(Q)$  for the input query. We then assign weights to the tokens in  $\mathcal{S}(Q)$  by the following strategy: If a token is not a product type token, we give a weight of 1.0. Otherwise, we assign weight  $W > 1$  to this token. Here,  $W$  is a hyperparameter in our algorithm. In our real-world experiments, we find  $W = 10$  gave us the best results. To generate the hash signatures of the weighted set  $\mathcal{S}(Q)$ , we leverage recent advances in efficiently computing weighted minhash signatures [7, 11, 21].

### 3.5 Theoretical Analysis

In this subsection, we analyze the complexity of our algorithm.

**Indexing Step Time Complexity:** In the proposed algorithm, the average time complexity of computing the hashes for one query is  $O(LT)$ , where  $L$  is the number of repetitions of LSH and  $T$  is the average number of tokens per query. The complexity of generating the entire robust cache structure is  $O(LNT)$  for a dataset with  $N$  queries. In practice,  $L$  and  $T$  are small constants much less than  $N$ , so we can consider asymptotic time complexity to be  $O(N)$ . This

linear time complexity of building the cache gives our method a significant scaling advantage to cache a massive amount of data.

**Retrieval Step Time Complexity:** The time complexity of ROSE’s retrieval step is  $O(LT \cdot BL)$ .  $O(LT)$  is the complexity of calculating the hash values for the incoming query.  $O(BL)$  is the time complexity of  $k$ -selection in the combined sets, where  $B$  is the bucket size. Therefore, the retrieval step’s overall time complexity is  $O(L^2BT)$ , independent of the cache size  $N$ . In practice,  $L$ ,  $B$  and  $T$  are small constants. As a result, cache retrieval’s time complexity is constant, which gives ROSE the decisive advantage for latency-critical services like product search.

**Memory Complexity:** The memory usage of ROSE is  $O(B \cdot N_B \cdot L)$ , where  $N_B$  is the number of buckets in one hash table.  $N_B$  is a hyperparameter and is a constant number independent of the cache size. We can see that the memory usage is not increasing with the size of the cache. This enables ROSE to achieve fast retrieval speeds on massive data with minimal memory costs, an ideal combination for industry-scale search engines.

**Error Analysis:** Due to the randomized nature of LSH, we note that it is possible to map the original query to an unrelated bucket with some small, but nonzero probability. However, we can dramatically reduce this error probability by maintaining  $L$  independent hash tables. In the following analysis, we choose the lexical preserving hashing for our study. All the theoretical analysis can be generalized to product type preserving hashing.

Given two query keywords  $X$  and  $Y$ , let  $C$  denote the sum of independent indicator random variables representing the number of buckets in which  $X$  and  $Y$  collide. We note that we can apply the following Chernoff bounds [17] :

$$\begin{aligned} Pr[C \leq (1 - \delta)Lp_{(x,y)}] &\leq \exp\left(\frac{-Lp_{(x,y)}\delta^2}{2}\right) \\ Pr[C \geq (1 + \delta)Lp_{(x,y)}] &\leq \exp\left(\frac{-Lp_{(x,y)}\delta^2}{3}\right), \end{aligned} \quad (1)$$

where  $L$  denotes the number of repetitions (hash tables) and  $p_{(x,y)}$  represents the collision probability between queries  $X$  and  $Y$ . These bounds show that the random variable  $C$  falls off from its expected value exponentially quickly with the number of repetitions.

To illustrate the behavior of these bounds, consider two queries “candi” and “corn” against the target query “candy.” Assuming that the pair (“candi”, “candy”) represents the same query while the tuple (“corn”, “candy”) is unrelated, we note that the former pair has a Jaccard similarity of 2/3 while the latter has a value of 2/7. We can use the first Chernoff bound above to upper bound the probability that the related pair (“candi”, “candy”) will not map to the same bucket in the hash table. Similarly, we can use the second bound to bound the probability that (“corn”, “candy”) will map to the same bucket. In both cases, we are upper bounding the probability of a failure case. Assuming that we set  $L = 10$ , we choose  $\delta$  accordingly and find that

$$\begin{aligned} Pr[C^{(\text{candi}, \text{candy})} \leq L/2] &\leq 1.6 \cdot 10^{-6} \\ Pr[C^{(\text{corn}, \text{candy})} \geq L/2] &\geq 9 \cdot 10^{-4} \end{aligned}$$

This practical illustration demonstrates the effectiveness of our proposed method. We will provide empirical study in Section 4.

## 4 OFFLINE EXPERIMENTS

**Dataset:** We sampled approximately 60 million well-performing queries from Amazon search logs as our cache’s target set. Following the same evaluation strategy in [19], our evaluation dataset samples queries from three buckets: a) **NQ**: Normal Queries, which are those in the top tercile of frequency, b) **HQ**: Hard queries sampled from the middle tercile of queries by frequency, and c) **LTQ**: Long-tail queries in the bottom tercile of frequency.

We randomly selected these queries from the search logs over one month. Each of these three sets contains 1000 queries. We obtained the re-mapped results for the queries from various query caching strategies and used a group of highly trained human judges to assign a binary relevance grade (relevant or irrelevant) to each returned query with respect to the original query’s intent. This relevant grade is used for calculating the performance metrics of different methods.

**Experimental Design:** We designed the experiments to answer two critical questions: a) **Robustness**: How accurate is **ROSE**’s retrieval process? b) **Efficiency**: How efficient is **ROSE**’s indexing and retrieval process? To answer these two questions, we test the following methods:

- **R-LP**: This method is our proposed method, **ROSE**, with lexical preserving hashing. The number of hash tables is  $L = 36$  and the number of hashes is  $K = 3$ .
- **R-PT**: This method is our proposed method, **ROSE**, with product type preserving hashing. All the other hyperparameters are the same as **ROSE-LP**.
- **EC** [5]: This is the exact-match cache implemented as a standard hash map. In the retrieval phase, the Exact-only cache returns the exact match candidates.
- **BF**: This is a cache structure designed by replacing **ROSE**’s retrieval algorithm with brute force search. We use edit distance as our similarity measure, computed via a dynamic programming algorithm<sup>1</sup>.
- **FC**: This is a designed cache structure for embedding vectors by replacing **ROSE**’s indexing and retrieval algorithm with FAISS [12]. We obtain the embeddings for each input query using a semantic product embedding model [18]. We choose the hyperparameters suggested by [12].

We adopted three commonly used metrics for offline evaluation: Precision, Recall, and  $F_1$  Measure. To compute these metrics, we utilize human judgments of relevance. We also analyze the speed of different methods for the indexing generation time and the online retrieval time.

**Overall Performance:** The results of all methods under the five metrics are presented in Table 1. Compared with other methods, **ROSE** performs the best on all three datasets. Specifically, **ROSE** offers a relative performance gain of 1.2% in Recall and 2.0% in  $F_1$  over the best baselines averaged across the three datasets. In particular, we find that the improvements of **ROSE-PT** on Normal queries and Hard Queries are more significant than Long-tail queries. On the other hand, **ROSE-LP** performs better on long-tail queries compared to **ROSE-PT**. Additionally, **ROSE** not only achieves a superior quality over these competing methods, but does

<sup>1</sup><https://www.geeksforgeeks.org/edit-distance-dp-5/>

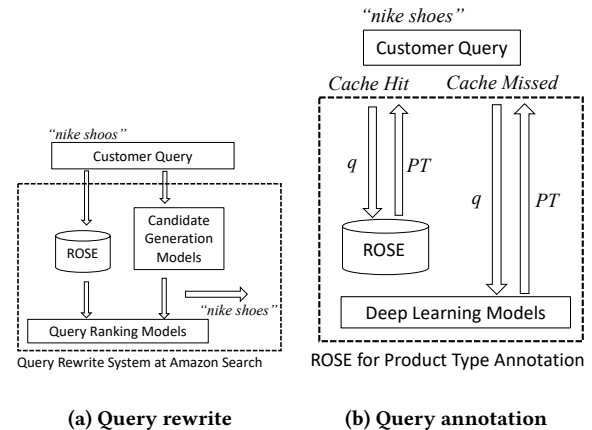


Figure 3: System overview for ROSE in Amazon Search

so with much better efficiency. **ROSE** is significantly faster in terms of index generation time and online retrieval time. In particular, **ROSE-LP** completes the index generation process in 65 minutes while **ROSE-PT** requires 75 minutes. Compared with other caches such as BF-Cache and FAISS-Cache, **ROSE** has a decisive speed advantage. **ROSE** can finish the online retrieval process in around 2ms, while FAISS-cache needs 120ms and BF-Cache requires 65 minutes. In summary, **ROSE** shows strong retrieval performance with extremely low latency and minimal cost, which makes it a compelling solution for latency-critical services such as product search engines.

## 5 SYSTEM DEPLOYMENT IN AMAZON

### 5.1 ROSE for Query Rewrite

We deployed **ROSE** within the Amazon.com product search engine to rewrite problematic user queries, such as those with typos, to alternative queries that provide a better user experience. We refer to this system as **ROSE-QR**. Leveraging lexical-preserving hashing, **ROSE-QR** maps an incoming query to one of the existing cached queries that have high-quality results according to lexical similarity.

We ran an online A/B experiment on the Amazon search engine to test **ROSE-QR**’s impact on the user experience. In the online experiment, users in the treatment group saw expanded search results from the alternative queries generated by **ROSE-QR**. Professional human judges measured the quality of the top search results shown in each arm of the experiment. We tracked the reduction of recall failures when the search engine does not return enough results for the user queries. We also measured business metrics such as revenue and purchased units. Our system did a better job in providing more relevant results, as measured by human evaluators, and significantly improved several business metrics as shown in Table 2.

### 5.2 ROSE for Product Type Annotation

The intended product type, such as *shoes* in the query "red nike shoes", is the most critical information in a user query. Identifying the correct product type from the query helps the search engine

Data set	Metrics	ROSE-LP	ROSE-PT	Exact-Cache	BF-Cache	FAISS-Cache
Normal Queries	Precision	0.88±0.03	0.96±0.01	<b>1.00±0.00</b>	0.90±0.02	0.96±0.08
	Recall	0.81±0.02	<b>0.90±0.04</b>	0.50±0.04	0.88±0.02	0.89±0.09
	F1-Measure	0.84±0.05	<b>0.93±0.08</b>	0.70±0.04	0.89±0.09	0.92±0.03
Hard Queries	Precision	0.78±0.01	0.90±0.03	<b>1.00±0.00</b>	0.80±0.03	0.89±0.07
	Recall	0.80±0.09	<b>0.86±0.05</b>	0.52±0.05	0.79±0.08	0.85±0.07
	F1-Measure	0.79±0.06	<b>0.88±0.09</b>	0.39±0.06	0.79±0.08	0.87±0.07
Long-tail Queries	Precision	0.77±0.03	0.73±0.06	<b>1.00±0.00</b>	0.76±0.04	0.75±0.02
	Recall	<b>0.79±0.04</b>	0.76±0.03	0.12±0.03	0.75±0.03	0.78±0.02
	F1-Measure	0.74±0.05	0.74±0.05	0.21±0.03	0.75±0.04	<b>0.76±0.05</b>
Index Generation Time		65m±5m	75m±5m	10m±1m	0m±0m	120m±10m
Cache Retrieval Time		1.8ms±0.3 ms	2.1ms±0.2ms	0.1ms±0.1ms	65m±3m	120ms±20ms

**Table 1: Offline Experiment Results. ROSE shows strong performance across both retrieval quality and system performance.**

retrieve the correct products and display a search result page layout customized for each product type.

We implemented **ROSE** to cache the intended product type of 5-10 million frequent queries. For an incoming tail query, **ROSE** maps the query to a few cached queries and uses the retrieved cached product types as the prediction for the tail query’s product type. To evaluate the impact on user experience, we used our **ROSE** product type prediction model to filter out irrelevant search results with the wrong product types. We deployed this system in the Amazon.com product search engine and measured the search defect rate with and without product type recognition. We define the product type defect rate as the number of products in the top 16 results with the wrong product type. From Table 2, we observe that, by using **ROSE**, the defect rate decreased by 1.7%, a significant improvement to the user experience.

ROSE-QR Metric	Gain	Filter Method	Defects Rate
Revenue	+0.42%	No Filtering	11.1%
Purchases	+0.30%	<b>ROSE-PT</b>	9.4%
Click-Through Rate	+7.26%		

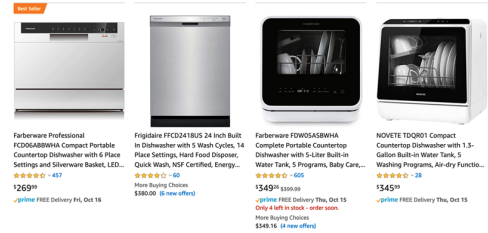
**Table 2: Production Impact of ROSE in Amazon Search**

## 6 CONCLUSION

With the dramatic growth in e-commerce adoption during the ongoing COVID-19 global pandemic, there is a pressing need to build scalable search systems that can gracefully handle this heightened consumer demand while preserving or even improving the search quality. Towards this end, in this paper we present **ROSE** for product search. **ROSE** is a robust cache that maps an online query to cached queries by preserving the query intent (lexically or semantically). Our proposed system is highly scalable and can deal with hundreds of millions of candidates in constant time and constant memory. We provide both a theoretical analysis of **ROSE** as well as an extensive offline evaluation. We deployed **ROSE** in the Amazon.com search engine and witnessed significant improvement over the existing solutions in terms of system performance and business metrics. In the future, we hope to apply the core ideas of **ROSE** towards improving the performance of additional information retrieval systems that may benefit from the flexibility and scalability of robust caches.



**(a) Top 4 results without product type restriction.**



**(b) Top 4 results with product type restriction.**

**Figure 4: The top-4 results for the query "dishwasher".**

## REFERENCES

- [1] Aman Ahuja, Nikhil Rao, Sumeet Katariya, Karthik Subbian, and Chandan K Reddy. 2020. Language-Agnostic Representation Learning for Product Search on E-Commerce Platforms. In *Proceedings of the 13th International Conference on Web Search and Data Mining*, 7–15.
- [2] Keping Bi, Choon Hui Teo, Yesh Dattatreya, Vijai Mohan, and W Bruce Croft. 2019. A Study of Context Dependencies in Multi-page Product Search. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. 2333–2336.
- [3] Andrei Broder. 2005. Algorithms for duplicate documents. URL: <http://www.cs.princeton.edu/courses/archive/spr05/cos598E/bib/Princeton.pdf> (: 01.02. 2015) (2005).
- [4] Andrei Z. Broder and Michael Mitzenmacher. 2001. Completeness and robustness properties of min-wise independent permutations. *Random Struct. Algorithms* 18, 1 (2001), 18–30.
- [5] Randal E Bryant, O'Hallaron David Richard, and O'Hallaron David Richard. 2003. *Computer systems: a programmer's perspective*. Vol. 2. Prentice Hall Upper Saddle River.
- [6] Flavio Chierichetti, Ravi Kumar, Sandeep Pandey, and Sergei Vassilvitskii. 2010. Finding the jaccard median. In *Proceedings of the twenty-first annual ACM-SIAM symposium on Discrete Algorithms*. SIAM, 293–311.
- [7] Tobias Christiani. 2020. DartMinHash: Fast Sketching for Weighted Sets. *arXiv preprint arXiv:2005.11547* (2020).
- [8] Yulong Gu, Zhuoye Ding, Shuaiqiang Wang, and Dawei Yin. 2020. Hierarchical User Profiling for E-commerce Recommender Systems. In *Proceedings of the 13th International Conference on Web Search and Data Mining*. 223–231.

- [9] Christian Hansen, Rishabh Mehrotra, Casper Hansen, Brian Brost, Lucas Maystre, and Mounia Lalmas. 2021. Shifting consumption towards diverse content on music streaming platforms. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*. 238–246.
- [10] Piotr Indyk and Rajeev Motwani. 1998. Approximate Nearest Neighbors: Towards Removing the Curse of Dimensionality. In *Proceedings of the Thirtieth Annual ACM Symposium on Theory of Computing* (Dallas, Texas, USA) (STOC '98). Association for Computing Machinery, New York, NY, USA, 604–613. <https://doi.org/10.1145/276698.276876>
- [11] Sergey Ioffe. 2010. Improved consistent sampling, weighted minhash and 11 sketching. In *2010 IEEE International Conference on Data Mining*. IEEE, 246–255.
- [12] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2017. Billion-scale similarity search with GPUs. *arXiv preprint arXiv:1702.08734* (2017).
- [13] Ting Liang, Guanxiong Zeng, Qiwei Zhong, Jianfeng Chi, Jinghua Feng, Xiang Ao, and Jiayu Tang. 2021. Credit Risk and Limits Forecasting in E-Commerce Consumer Lending Service via Multi-view-aware Mixture-of-experts Nets. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*. 229–237.
- [14] Heran Lin, Pengcheng Xiong, Danqing Zhang, Fan Yang, Ryoichi Kato, Mukul Kumar, William Headden, and Bing Yin. [n.d.]. Light Feed-Forward Networks for Shard Selection in Large-scale Product Search. ([n. d.]).
- [15] Chen Luo. 2020. *Some Rare LSH Gems for Large-scale Machine Learning*. Ph.D. Dissertation. Rice University.
- [16] Chen Luo and Anshumali Shrivastava. 2018. Arrays of (locality-sensitive) count estimators (ace): High-speed anomaly detection via cache lookups. (2018).
- [17] Michael Mitzenmacher and Eli Upfal. 2005. *Probability and Computing: Randomized Algorithms and Probabilistic Analysis*. Cambridge University Press, USA.
- [18] Priyanka Nigam, Yiwei Song, Vijai Mohan, Vihan Lakshman, Weitian Ding, Ankit Shingavi, Choon Hui Teo, Hao Gu, and Bing Yin. 2019. Semantic product search. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2876–2885.
- [19] Xichuan Niu, Bofang Li, Chenliang Li, Rong Xiao, Haochuan Sun, Hongbo Deng, and Zhenzhong Chen. 2020. A Dual Heterogeneous Graph Attention Network to Improve Long-Tail Performance for Shop Search in E-Commerce. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 3405–3415.
- [20] Jiayu Shi, Huaxiu Yao, Xian Wu, Tong Li, Zedong Lin, Tengfei Wang, and Bin-qiang Zhao. 2021. Relation-aware Meta-learning for E-commerce Market Segment Demand Prediction with Limited Records. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*. 220–228.
- [21] Anshumali Shrivastava. 2016. Simple and Efficient Weighted Minwise Hashing.. In *NIPS*. 1498–1506.
- [22] Anshumali Shrivastava. 2017. Optimal Densification for Fast and Accurate Minwise Hashing. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6–11 August 2017 (Proceedings of Machine Learning Research, Vol. 70)*, Doina Precup and Yee Whye Teh (Eds.). PMLR, 3154–3163. <http://proceedings.mlr.press/v70/shrivastava17a.html>
- [23] Anshumali Shrivastava and Ping Li. 2015. Improved Asymmetric Locality Sensitive Hashing (ALSH) for Maximum Inner Product Search (MIPS). In *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence, UAI 2015, July 12–16, 2015, Amsterdam, The Netherlands*, Marina Meila and Tom Heskes (Eds.). AUAI Press, 812–821. <http://auai.org/uai2015/proceedings/papers/96.pdf>
- [24] Zehong Tan, Canran Xu, Mengjie Jiang, Hua Yang, and Xiaoyuan Wu. 2017. Query rewrite for null and low search results in eCommerce. In *eCOM@ SIGIR*.
- [25] Jeffrey S Vitter. 1985. Random sampling with a reservoir. *ACM Transactions on Mathematical Software (TOMS)* 11, 1 (1985), 37–57.
- [26] Yiqiu Wang, Anshumali Shrivastava, Jonathan Wang, and Junghee Ryu. 2018. Randomized algorithms accelerated over cpu-gpu for ultra-high dimensional similarity search. In *Proceedings of the 2018 International Conference on Management of Data*. 889–903.
- [27] Rong Xiao, Jianhui Ji, Baoliang Cui, Haihong Tang, Wenwu Ou, Yanghua Xiao, Jiwei Tan, and Xuan Ju. 2019. Weakly Supervised Co-Training of Query Rewriting and Semantic Matching for e-Commerce. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*. 402–410.
- [28] Da Xu, Chuanwei Ruan, Evren Korpeoglu, Sushant Kumar, and Kannan Achan. 2020. Product knowledge graph embedding for e-commerce. In *Proceedings of the 13th international conference on web search and data mining*. 672–680.
- [29] Da Xu, Chuanwei Ruan, Evren Korpeoglu, Sushant Kumar, and Kannan Achan. 2021. Theoretical Understandings of Product Embedding for E-commerce Machine Learning. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*. 256–264.
- [30] Danqing Zhang, Zheng Li, Tianyu Cao, Chen Luo, Tony Wu, Hanqing Lu, Yiwei Song, Bing Yin, Tuo Zhao, and Qiang Yang. 2021. QUEACO: Borrowing Treasures from Weakly-labeled Behavior Data for Query Attribute Value Extraction. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 4362–4372.
- [31] Junhao Zhang, Weidi Xu, Jianhui Ji, Xi Chen, Hongbo Deng, and Keping Yang. 2021. Modeling Across-Context Attention For Long-Tail Query Classification in E-commerce. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*. 58–66.
- [32] Kang Zhao, Hongtao Lu, and Jincheng Mei. 2014. Locality Preserving Hashing. In *Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence, July 27–31, 2014, Québec City, Québec, Canada*, Carla E. Brodley and Peter Stone (Eds.). AAAI Press, 2874–2881. <http://www.aaai.org/ocs/index.php/AAAI/AAAI14/paper/view/8357>