# Efficient Sketching and Nearest Neighbor Search Algorithms for Sparse Vector Sets

SEBASTIAN BRUCH, Northeastern University, United States of America
FRANCO MARIA NARDINI, ISTI-CNR, Italy
COSIMO RULLI, ISTI-CNR, Italy
ROSSANO VENTURINI, University of Pisa, Italy

Sparse embeddings of data form an attractive class due to their inherent interpretability: Every dimension is tied to a term in some vocabulary, making it easy to visually decipher the latent space. Sparsity, however, poses unique challenges for Approximate Nearest Neighbor Search (Anns) which finds, from a collection of vectors, the $k$ vectors closest to a query. To encourage research on this underexplored topic, sparse Anns featured prominently in a BigANN Challenge at NeurIPS 2023, where approximate algorithms were evaluated on large benchmark datasets by throughput and accuracy. In this work, we introduce a set of novel data structures and algorithmic methods, a combination of which leads to an elegant, effective, and highly efficient solution to sparse Anns. Our contributions range from a theoretically-grounded sketching algorithm for sparse vectors to reduce their effective dimensionality while preserving inner product-induced ranks; a geometric organization of the inverted index; and the blending of local and global information to improve the efficiency and efficacy of Anns. Empirically, our final algorithm, dubbed Seismic, reaches sub-millisecond per-query latency with high accuracy on a large-scale benchmark dataset using a single CPU.

Additional Key Words and Phrases: maximum inner product search, sketching, sparse vectors

## 1 INTRODUCTION

Embeddings have gained increasing popularity since the introduction of Large Language Models (Llms) [40]. These models learn a high-dimensional vector[1] representation of their input such that the distance between any two input variables is approximately preserved in the target vector space.

As a concrete example, consider Neural Information Retrieval (Nir) where short text passages are embedded into a vector space, wherein the similarity between two vectors reflects the semantic similarity between the passages they represent. Representing text that way enables a more effective matching of queries to passages by semantic similarity [6]. That in turn has dramatically changed applications such as web search, recommendation, question answering, and more.

While embedding vectors are typically dense (i.e., each coordinate is almost surely nonzero), there are neural models that produce *sparse* vectors (i.e., in a matrix of embeddings, columns and rows are mostly zero). There are a few reasons why sparse embeddings are enticing, the most important of which is their inherent interpretability.

Sparse embeddings are interpretable because one may ground each dimension in a "term" from some fixed "vocabulary," so that by inspecting the nonzero coordinates of an embedding, one can intuitively understand what latent features are extracted from the input and represented in the embedding space. In Nir, for instance, sparse embedding models [24, 25, 27, 35, 43] tie each embedding dimension to a word in the source language. A coordinate that is nonzero in an embedding indicates that the corresponding word is semantically relevant to the input.

---

[1]We use "vector" and "point" interchangeably throughout this work.

---

Authors' addresses: Sebastian Bruch, Northeastern University, Boston, MA, United States of America, s.bruch@northeastern.edu; Franco Maria Nardini, ISTI-CNR, Pisa, Italy, francomaria.nardini@isti.cnr.it; Cosimo Rulli, ISTI-CNR, Pisa, Italy, cosimo.rulli@isti.cnr.it; Rossano Venturini, University of Pisa, Pisa, Italy, rossano.venturini@unipi.it.

## 1.1 Sparse Maximum Inner Product Search

Embeddings—dense or sparse—allow us to disentangle data representation from "search": the task of finding $k$ dataset items with the largest similarity to an arbitrary query. That is because, if data of any modality (e.g., text, image, video, or audio) are projected into the same vector space, search is simply the problem of *Nearest Neighbor Search (Nns)*, defined formally below:

DEFINITION 1 (NEAREST NEIGHBOR SEARCH). *Given a set $\mathcal{S} \subset \mathbb{R}^d$ with distance function $\delta : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ and query point $q \in \mathbb{R}^d$, Nns finds the $k$-subset of $\mathcal{S}$ consisting of $k$ closest points to $q$:*

$$\underset{u \in \mathcal{S}}{\overset{(k)}{\arg\min}}\, \delta(q, u).$$

For efficiency reasons, however, the problem is often relaxed to *Approximate Nearest Neighbor Search (Anns)*, which may erroneously include farther away points in the returned set [5]. Its accuracy is quantified as follows, with $|\cdot|$ denoting the size of its argument:

DEFINITION 2. *Suppose $S \subset \mathcal{S}$ is the set of $k$ nearest neighbors of query $q$ in $\mathcal{S}$. The accuracy of an Anns algorithm that returns a set of $k$ points $S'$, denoted accuracy@$k$, is $|S \cap S'|/k$.*

In this document, we are interested in *inner product* as a measure of *similarity* between vectors. We emphasize that algorithms for Anns by inner product can be trivially applied to Anns by cosine similarity, thereby covering the vast majority of use cases. Furthermore, we focus squarely on an instance of the problem over a set of *sparse* vectors, resulting in the following problem:

DEFINITION 3 (SPARSE MAXIMUM INNER PRODUCT SEARCH (SPARSE MIPS)). *Consider a set $\mathcal{S} \subset \mathbb{R}^d$ of sparse vectors and query point $q$. Sparse Mips is an instance of Nns that finds the subset of $k$ points in $\mathcal{S}$ most similar to $q$ by inner product, denoted by $\langle \cdot, \cdot \rangle$:*

$$\underset{u \in \mathcal{S}}{\overset{(k)}{\arg\max}}\, \langle q, u \rangle.$$

## 1.2 Challenges in Sparse MIPS

While Dense and Sparse MIPS are conceptually the same, sparsity poses unique challenges. First, any computation involving sparse vectors is incompatible with widely-adopted modern platforms (e.g., Graphics Processing Units and Tensor Processing Units). That is because such hardware platforms are designed to process data that is represented as contiguous arrays of floating point numbers—the standard way to store and represent dense vectors. While it is possible to represent sparse vectors in that format, that would be impractical, wasteful, and ultimately ineffective in high dimensions—often in the order of tens to hundreds of thousands of dimensions.

Second is a dearth of efficient representations. Contrast that to dense vectors for which a plethora of algorithms exist to reduce the amount of space taken up by a vector while (a) approximately preserving distances and (b) enabling fast distance computation [28, 29, 31–33, 69, 70].

The two challenges above—computational and storage inefficiencies—make Sparse MIPS far more resource-intensive than the same operation over dense vector sets. That is further exacerbated as datasets grow larger in the number of vectors or the number of nonzero coordinates per vector.

To counter some of these issues, a number of efforts have explored specialized algorithms for Sparse MIPS [7, 8, 26, 49, 51]. In fact, to encourage the community to design practical algorithms, the 2023 BigANN Challenge [65] at the Neural Information Processing Systems (NeurIPS) conference dedicated a track to this problem.

### 1.3   Seismic: Our Proposed Algorithm

Inspired by BigANN, we presented in [9] a novel Sparse Mips algorithm that we called Seismic (**S**pilled Clust**e**ring of **I**nverted Lists with **S**ummaries for **M**aximum **I**nner Produ**c**t Search). The algorithm, which we describe in detail in Section 4, admits effective and efficient Sparse Mips.

Seismic is realized by the amalgamation of a set of novel data structures and algorithmic contributions, each of which may be of independent interest. One such contribution is an elegant sketching algorithm to reduce the effective dimensionality of sparse vectors into a lower effective dimensionality such that inner products are approximately preserved with high probability. This sketching algorithm, which we call Set $\alpha$-Mass Subvector Sketch (Set $\alpha$-MSS), is critical to the success of Seismic. We describe this in more detail and present an extensive analysis in Section 3.

Another novel contribution is a data structure. Rather, Seismic *enhances* an existing data structure called the *inverted index* (i.e., a mapping from dimensions to an *inverted list* containing vectors with a nonzero coordinate in that dimension). The enhancement is through the organization of inverted lists into geometrically-cohesive blocks—an approach previously explored for other tasks [1, 13]. Each block is then summarized by a representative point with the idea that the summaries help identify the subset of blocks that can be dismissed without evaluation during search. This logic turns out to be a critical factor in enabling highly efficient search.

Our evaluation of Seismic against state-of-the-art baselines in [9], which included the top (open-source) submissions to the BigANN Challenge, empirically demonstrated that average per-query latency was in the *microsecond* territory on various sparse embeddings of MsMarco v1 [59]. Impressively, Seismic outperformed the winning solutions of the BigANN Challenge by a factor of at least 3.4× at 95% accuracy@10 on Splade—a sparse embedding model—and 12× on Efficient Splade—another sparse embedding model—with the margin widening substantially as accuracy increases. Other baselines were consistently *one to two orders of magnitude* slower than Seismic.

In a follow-up work [10], we further improved Seismic's query processing logic. We did so by incorporating another data structure that dovetails with the inverted index from before. The structure is the $\kappa$-Nearest Neighbor graph ($\kappa$-nn graph), a directed graph where each data point is connected to its $\kappa$ nearest neighbors.

When the inverted index returns an approximate set of nearest neighbors, we use the $\kappa$-nn graph to expand that set and form a larger pool of possible nearest neighbors. We do so by taking each point in the returned set, collecting the points that are one hop away from it in the $\kappa$-nn graph, and inserting the new points into the larger set. The expanded set is then re-evaluated and the closest $k$ points are returned as the (refined) approximate nearest neighbors.

The reason this procedure improves the accuracy of Sparse Mips is that it blends the "global" structure of the embedding space (captured by the inverted index) with the "local" structure (captured by the graph). This complementarity is not new in and of itself; it has guided prior research [34, 44]. Our experiments reported in [10], confirm the intuition above: given a fixed memory budget, providing the $\kappa$-nn graph to Seismic improves latency by up to a factor of 2× for the same level of accuracy. This was especially advantageous if a near perfect accuracy is desired.

Finally, in another empirical study [11], we strengthened our claims by applying Seismic to a large-scale dataset of about 140 million sparse embedding vectors. As we reported in that work, Seismic scales well to such massive datasets: it indexes that collection in a fraction of the time required by other algorithms; and, it produces 10 approximate nearest neighbors per query point in a mere 3 milliseconds with 98% accuracy.

### 1.4 Contributions

This work consolidates and extends our contributions from [9–11] and presents a comprehensive evaluation of the final algorithm on benchmark datasets. Concretely, we make the following contributions in this work:

- We present the evolution of Set $\alpha$-MSS, our sketching algorithm for sparse vectors, along with a comprehensive theoretical analysis. This analysis, which explains the choices we made in the design and implementation of the algorithm, was not included in our previous publications.
- We present a unified algorithm from the pieces scattered in separate publications, and tie it back to our theoretical analysis.
- We extend our empirical evaluation of Seismic to include new embedding models and datasets; compare against other baseline Sparse Mips algorithms; and, examine the sensitivity of Seismic to its hyperparameters and study the algorithm's generalizability.

### 1.5 Outline

The remainder of this paper is organized as follows. Section 2 reviews related work to position our work in the literature. In Section 3, we give a detailed account of our sketching algorithm with theoretical analysis. The main algorithm, Seismic, is then described in Section 4. We evaluate Seismic in Section 5 and conclude our work in Section 6.

## 2 RELATED WORK

This section reviews notable related research. We first summarize the thread of work on sparse embedding vectors by highlighting their use cases. We then venture into the field of text retrieval and review the algorithms presented in that literature. The reason for the latter topic is that we borrow the inverted index data structure from that literature—although as we explained earlier, we enhance it in novel ways and adapt it to Sparse Mips.

### 2.1 Sparse Embedding Vectors

Sparse embeddings were investigated in the context of text retrieval long before Llms [72]. But the rise of Llms supercharged this research and led to a flurry of activity on the topic [2, 15–17, 24, 26, 39, 43, 73]. See Figure 1(a) for an example sparse embedding of text used in text retrieval.

What the cited references claim and show empirically is that Sparse Mips over their proposed sparse embeddings (e.g., Splade [24]) can reach state-of-the-art performance for the task of text retrieval on a suite of benchmark datasets like MsMarco v1 [59] and Beir [66]. All that while also lending interpretability to the embeddings of text documents and queries, as argued earlier.

Unsurprisingly, due to their highly interpretable nature, sparse embeddings also play a pivotal role in the literature on interpreting Llms. That literature is peppered with algorithms that "probe" an Llm by way of learning a sparse representation of its layers.

The tool that does the probing is known as a Sparse Autoencoder (Sae) [14, 61]. As illustrated in Figure 1(b), Saes learn to project a dense vector (obtained from an Llm's internal layer) to a sparse vector, then back to the same dense vector space. The learning is unsupervised: Saes learn to reconstruct their input. What is interesting, however, is the sparse embedding that is learned in the middle: It is argued that this projection "reverses" superposition [23] and polysemanticity [14], allowing us to "find concise, human-understandable explanations for what neural networks are doing internally" [3, 14, 30, 38, 57, 61]. In this context, Sparse Mips can be used to gather insight not from individual vectors, but a large collection of them.
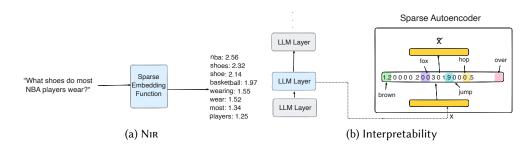
Fig. 1. Example usage of sparse embeddings. In (a), a text input is embedded into a sparse vector (shown as a set of nonzero coordinate-value pairs), where terms that are semantically similar to the input have nonzero values. Such embeddings enable accurate and interpretable query-document matching in text retrieval. In (b), a Sparse Autoencoder probes one layer of an Llm, translating its dense output into a sparse vector and back (such that $x \approx \tilde{x}$). The sparse vector reveals the features that the layer has learned, shedding light on what that layer processes.

## 2.2 Text Retrieval Algorithms

The information retrieval literature offers a wide array of algorithms tailored to retrieval on text collections [67]: given a text query, the task is to find the top $k$ text documents that are "relevant" to the query. Theses algorithms are often *exact* and scale easily to massive datasets. MaxScore [68] and WAND [4], and their intellectual descendants [21, 22, 52, 53], are examples that perform search over a "bag-of-words" representation of text, such as BM25 [63] or TF-IDF [64].

These algorithms operate on an inverted index, augmented with additional data to speed up query processing. One that features prominently is the maximum attainable partial inner product—an upper-bound. This enables the possibility of navigating the inverted lists, one document at a time, and deciding quickly if a document may belong to the result set. Effectively, such algorithms (safely) *prune* the parts of the index that cannot be in the top-$k$ set. That is why they are often referred to as *dynamic pruning* techniques.

Although efficient in practice, dynamic pruning methods are designed specifically for text collections. Importantly, they ground their performance on several pivotal assumptions: nonnegativity, higher sparsity rate for queries, and a Zipfian shape of the length of inverted lists. These assumptions are valid for TF-IDF or BM25, which is the reason why dynamic pruning works well and the worst-case time complexity of search is seldom encountered in practice.

These assumptions do not necessarily hold for sparse embeddings, however. Embedding vectors may be real-valued, with a sparsity rate that is closer to uniform across dimensions [7, 48]. Mackenzie *et al.* [49] find that sparse embeddings reduce the odds of pruning or early-termination in the document-at-a-time (DaaT) and Score-at-a-Time (SaaT) paradigms. That is not to say there have not been attempts to modify inverted index-based algorithms to make them compatible with sparse embeddings [45, 51, 55]. Unlike these *exact* methods, our approach is *approximate* by design, offering far greater flexibility to sacrifice accuracy for latency per operational requirements.

The most similar work to ours is [8]. The authors investigate if Anns algorithms for *dense* vectors port over to *sparse* vectors. They focus on *inverted file* (IVF) where vectors are partitioned into clusters during indexing, with only a fraction of clusters scanned during search. They show that IVF serves as an efficient solution for sparse Anns. Interestingly, the authors cast IVF as dynamic pruning and turn that insight into a novel organization of the inverted index for sparse Anns for general sparse vectors. Our index structure can be viewed as an extension of theirs.

Finally, we briefly describe another Anns algorithm over dense vectors: HNSW [50], an algorithm that constructs a graph where each data point is a node and two nodes are connected if they are deemed "similar." Similarity is based on Euclidean distance, but [58] shows inner product too can be used to form the graph (known as IP-HNSW). As we learn in the presentation of our empirical analysis, algorithms that adapt IP-HNSW [58] to sparse vectors work remarkably well.

## 3  SPARSE VECTOR SKETCHING

As we noted earlier, one of the challenges with sparse vectors that have a large number of nonzero coordinates is that exact Nns algorithms struggle to perform the task when the latency budget is tight. To make search efficient, we explore a method to reformulate the problem as Sparse Mips.

The cornerstone of our reformulation is a sketching algorithm that reduces the number of nonzero coordinates such that the relative order of two points induced by their inner product with a query is approximately preserved. This section presents our sketching algorithm, which we call Set $\alpha$-MSS, along with a thorough analysis.

### 3.1  Notation

Before we describe the algorithm, let us introduce our notation. We use lower-case letters (e.g., $u$) to denote a vector and write $u_i$ for its $i$-th coordinate. A sparse vector can be identified as a set of nonzero entries $\{(i, u_i) \mid u_i \neq 0\}$.

Let $\mathbb{B}_1^d = \{u \in [0, 1]^d \mid \|u\|_1 \leq 1\}$ be the nonnegative region of the $\ell_1$ unit ball. Suppose we have a set $\mathcal{S} \subset \mathbb{B}_1^d$ of nonnegative *sparse* vectors. We assume $\mathcal{S} \sim \mathcal{D}$ for some unknown data distribution $\mathcal{D}$. We use superscript to enumerate: $u^{(j)}$ is the $j$-th vector in $\mathcal{S}$.

Let us define a few concepts that we frequently refer to. The $\ell_p$ norm of a vector denoted by $\|\cdot\|_p$ is defined as $\|x\|_p = (\sum_i |x_i|^p)^{1/p}$. We call the $\ell_p$ norm of a vector its $\ell_p$ *mass*. Additionally, we use the notion of a *subvector* throughout this work: The subvector $\tilde{u}$ of sparse vector $u$ is a vector with the property that $\tilde{u}_i \neq 0 \implies u_i \neq 0 \land \tilde{u}_i = u_i$. Finally, a vector $u$ *restricted to a set of dimensions* $\mathcal{I}$ is a vector $u_\mathcal{I}$ such that $(u_\mathcal{I})_i = u_i$ if $i \in \mathcal{I}$ and 0 otherwise.

### 3.2  $\ell_1$ Threshold Sampling

A recent study [18] presented a sketching algorithm for sparse vectors that rests on the following principle: Coordinates that contribute more heavily to the $\ell_2$ norm squared weigh more heavily on the inner product between vectors. Based on that intuition, they report that if we were to drop the nonzero coordinates of a sparse vector with a probability proportional to their contribution to the $\ell_2$ norm squared, we can reduce the size of vectors while approximately maintaining their inner products. We refer to their sketching algorithm as $\ell_2$-Threshold Sampling (TS).

*3.2.1  The $\ell_1$-TS algorithm.* We build on that same intuition, but consider instead the $\ell_1$ norm. We make this choice because, on the datasets we use in our experiments, $\ell_1$-TS yields inner product estimates that are far more accurate than the $\ell_2$ variant. Our algorithm detailed in Algorithm 1 is similar to the $\ell_2$-TS procedure and enjoys similar properties. Given the $\ell_1$-TS sketches of two vectors $u$ and $v$, their inner product is estimated as follows:

$$\sum_{i \in K_u \cap K_v} \frac{u_i v_i}{\min(1, u_i \tau_u, v_i \tau_v)}, \tag{1}$$

where $K_u$, $K_v$, $\tau_u$, and $\tau_v$ are produced by Algorithm 1.

*3.2.2  Analysis of $\ell_1$-TS.*

---

**Algorithm 1:** $\ell_1$ Threshold Sampling

---

**Input**: $u \in \mathbb{B}_1^d$; target sketch size $d_\circ$; uniformly random hash function $h : \{1, \ldots, d\} \rightarrow [0, 1]$.
**Result:** Sketch $ts(u) = \{K_u, V_u, \tau_u\}$, where $K_u$ is a subset of dimensions from $\{1, \ldots, d\}$; $V_u$ contains $u_i$ for all $i \in K_u$; and, $\tau_u \in \mathbb{R}_+$.

1: $K_u \leftarrow \emptyset; V_u \leftarrow \emptyset$
2: **for** $i$ such that $u_i \neq 0$ **do**
3: $\quad \tau_i \leftarrow d_\circ \cdot \frac{u_i}{\|u\|_1}$.
4: $\quad$ **if** $h(i) \leq \tau_i$ **then**
5: $\quad\quad K_u \leftarrow K_u \cup \{i\}; V_u \leftarrow V_u \cup \{u_i\}$
6: **return** $ts(u) = \{K_u, V_u, \tau_u\}$ where $\tau_u = d_\circ/\|u\|_1$.

---

THEOREM 3.1. *For vectors $u, v \in \mathbb{B}_1^d$ and sketch size $d_\circ$, let $ts(u)$ and $ts(v)$ denote sketches produced by Algorithm 1. Let $W$ be the inner product estimated by Equation* (1). *We have that $\mathbb{E}[W] = \langle u, v \rangle$ and*

$$\text{Var}[W] \leq \frac{1}{d_\circ} \langle u, v \rangle \big( \|u\|_1 + \|v\|_1 \big).$$

PROOF. Let $\mathcal{I}$ denote the set of all indices $i$ for which $u_i \neq 0$ and $v_i \neq 0$. For any $i \in \mathcal{I}$, let $\mathbb{1}_i$ denote the indicator random variable for the event that $i \in K_u \cap K_v$. Note that, for $i \neq j$, $\mathbb{1}_i$ is independent from $\mathbb{1}_j$, because the hash values $h(i)$ and $h(j)$ are drawn independently. Now, $\mathbb{P}[\mathbb{1}_i = 1]$ over all hash functions is equal to:

$$p_i = \min \left( 1, \frac{d_\circ u_i}{\|u\|_1}, \frac{d_\circ v_i}{\|v\|_1} \right) = \min(1, \tau_u u_i, \tau_v v_i).$$

Considering the independence of $\mathbb{1}_i$'s, we can write $W$ as $\sum_{i \in \mathcal{I}} \mathbb{1}_i \cdot u_i v_i / p_i$. By linearity of expectation, we have that:

$$\mathbb{E}[W] = \sum_{i \in \mathcal{I}} p_i \frac{u_i v_i}{p_i} = \sum_{i \in \mathcal{I}} u_i v_i = \langle u, v \rangle.$$

As for the variance, we again take advantage of independence and write:

$$\begin{aligned}
\text{Var}[W] &\leq \frac{1}{d_\circ} \sum_{i \in \mathcal{I}, \, p_i \neq 1} \frac{u_i^2 v_i^2}{\min(u_i/\|u\|_1, v_i/\|v\|_1)} = \frac{1}{d_\circ} \sum u_i^2 v_i^2 \max(\|u\|_1/u_i, \|b\|_1/v_i) \\
&= \frac{1}{d_\circ} \sum \max(u_i v_i^2 \|u\|_1, u_i^2 v_i \|v\|) \leq \frac{1}{d_\circ} \sum u_i v_i^2 \|u\|_1 + u_i^2 v_i \|v\|_1 \\
&= \frac{1}{d_\circ} \big( \|u\|_1 \sum u_i v_i^2 + \|v\|_1 \sum u_i^2 v_i \big) \leq \frac{1}{d_\circ} \big( \|u\|_1 \langle u, v \rangle + \|v\|_1 \langle u, v \rangle \big) \\
&= \frac{1}{d_\circ} \langle u, v \rangle \big( \|u\|_1 + \|v\|_1 \big).
\end{aligned}$$

$\square$

By applying the one-sided Chebyshev's inequality, we can derive the following corollary.

COROLLARY 3.2. *For all values of $\epsilon, \delta \in (0, 1)$ and sparse vectors $u, v \in \mathbb{B}_1^d$, $\ell_1$-TS with sketch size $d_\circ$ returns inner product $W$ such that, with probability at least $1 - \delta$:*

$$W - \langle u, v \rangle \leq \sqrt{\frac{1 - \delta}{\delta d_\circ} \langle u, v \rangle \big( \|u\|_1 + \|v\|_1 \big)}.$$

### 3.3 $\alpha$-Mass Subvector Sketch

The problem with Algorithm 1 is that the inner product of two sketches is a *weighted* product of surviving entries. That weight is a function of *both* vectors. That implies, when one of the vectors is a dynamic query point $q$, the weights have to be calculated for every data point $u$ during search, making Sparse MIPS complex and time consuming.

In this section, we use the principle underlying $\ell_1$-TS, but make its stochastic nature deterministic and simplify the inner product estimate computation. Doing so has an important caveat: the inner product of two sketches is an *under-estimate* of the true inner product. But as we will see later, this modification has its own advantages.

#### 3.3.1 The sketching algorithm.

DEFINITION 4 ($\alpha$-MASS SUBVECTOR SKETCH). *Consider a vector $u \in \mathbb{B}_1^d$ and a constant $\alpha \in (0, 1]$. Take a permutation $\pi$ such that $|u_{\pi_i}| \geq |u_{\pi_{i+1}}|$ and denote by $j$ the smallest integer such that:*

$$\sum_{i=1}^{j} |u_{\pi_i}| \geq \alpha \|u\|_1.$$

*We call $\tilde{u}$ made up of $\{(\pi_i, u_{\pi_i})\}_{i=1}^{j}$, the $\alpha$-Mass Subvector Sketch ($\alpha$-MSS) of $u$. Furthermore, the inner product estimate of two sketches $\tilde{u}$ and $\tilde{v}$ is defined simply as $\langle \tilde{u}, \tilde{v} \rangle = \sum \tilde{u}_i \tilde{v}_i$.*

It is easy to see that the sketch of Definition 4 leads to *biased* inner product estimates. However, given an $\alpha$-MSS of $u$ and a $\beta$-MSS of $v$, we can present the following bound on the error of the estimated inner product.

THEOREM 3.3. *Consider an $\alpha$-MSS $\tilde{u}$ of $u$ and $\beta$-MSS $\tilde{v}$ of $v$, where $u, v \in \mathbb{B}_1^d$. Defining $\mathcal{I} = \{i \text{ s.t. } u_i \neq 0 \land v_i \neq 0\}$, the following bound holds:*

$$\langle u, v \rangle - \langle \tilde{u}, \tilde{v} \rangle \leq (1 - \alpha)\|u\|_1 \cdot (1 - \beta)\|v\|_1 + \|u_{\mathcal{I}}\|_1 \cdot (1 - \beta)\|v\|_1 + \|v_{\mathcal{I}}\|_1 \cdot (1 - \alpha)\|u\|_1,$$

*where $u_{\mathcal{I}}$ is the vector $u$ restricted to $\mathcal{I}$.*

PROOF. Let $\mathbb{1}_p$ denote the indicator function which is 1 if the predicate $p$ holds.

$$\langle u, v \rangle - \langle \tilde{u}, \tilde{v} \rangle = \sum_{i \in \mathcal{I}} u_i v_i - \tilde{u}_i \tilde{v}_i = \sum_{i \in \mathcal{I}} \mathbb{1}_{\tilde{u}_i=0 \lor \tilde{v}_i=0} \cdot u_i v_i$$

$$= \sum_{i \in \mathcal{I}} \mathbb{1}_{\tilde{u}_i=0 \land \tilde{v}_i=0} \cdot u_i v_i + \mathbb{1}_{\tilde{u}_i=0 \land \tilde{v}_i\neq0} \cdot u_i v_i + \mathbb{1}_{\tilde{u}_i\neq0 \land \tilde{v}_i=0} \cdot u_i v_i.$$

Note that $\|\tilde{u}\|_1 \geq \alpha\|u\|_1$, $\|\tilde{v}\|_1 \geq \beta\|v\|_1$, and $\sum u_i v_i \leq \|u\|_1\|v\|_1$. As such the first term is bounded above by $\sum_{i \in \mathcal{I}} \mathbb{1}_{\tilde{u}_i=0} u_i \cdot \sum_{i \in \mathcal{I}} \mathbb{1}_{\tilde{v}_i=0} v_i$, which in turn is bounded by $(1 - \alpha)\|u\|_1 \cdot (1 - \beta)\|v\|_1$. The last two terms follow a similar reasoning. That yields the desired bound. □

#### 3.3.2 Empirical evidence.
Let us demonstrate the result above on the MSMARCO v1 dataset [59]—a benchmark text collection to be described in more detail later in this work—embedded into a sparse vector space with SPLADE [25].[2] We take every query and data vector, sort its entries by coordinate in descending order, and measure the fraction of the $\ell_1$ mass preserved by considering a given number of top coordinates.

For query points, the top 10 entries yield 0.75-mass subvectors. For data points, the top 50 (about 30% of nonzero entries), give 0.75-mass subvectors. We illustrate our measurements in Figure 2(a). Next, consider the inner product estimate using $\alpha$-MSS between query and data points. This is

---

[2]The cocondenser-ensembledistill checkpoint was obtained from https://huggingface.co/naver/splade-cocondenser-ensembledistil.
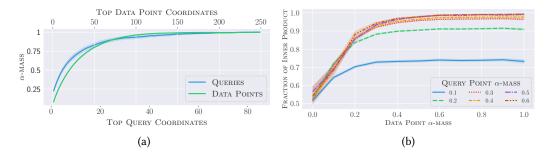
Fig. 2. Empirical demonstration of $\alpha$-MSS on the SPLADE embeddings of the MsMarco v1 dataset: (a) Fraction of $\ell_1$ mass preserved by keeping only the largest nonzero entries. This shows the number of coordinates that must be preserved (horizontal axis) in order to construct an $\alpha$-mass subvector (vertical axis). (b) Fraction of true inner product (with 95% confidence intervals) preserved by $\alpha$-MSS.

shown in Figure 2(b). We observe that, a 0.2-MSS sketch of queries and documents preserves an average of 85% of the true inner product.

## 3.4 Set $\alpha$-Mass Subvector Sketch

The sketching algorithms presented thus far apply to a vector independently of other vectors. As such, given a collection of vectors $\mathcal{S}$ and a constant $\alpha$, the sketch of every vector $u \in \mathcal{S}$ has a mass that is at least $\alpha\|u\|_1$. While that may be desirable for datasets where every vector has an equal likelihood of maximizing the inner product with an arbitrary query $q$, it becomes wasteful for datasets where that likelihood across vectors is non-uniform.

Let us elaborate by considering an extreme case. Suppose $u, v \in \mathbb{B}_1^d$ are such that $u = tv$ for some $t \in [0, 1)$. It is clear that, $\langle q, u \rangle < \langle q, v \rangle$ for any $q \in \mathbb{B}_1^d$. If $u$ and $v$ belong to the set $\mathcal{S}$, and noting that our ultimate goal is to find the solution to Sparse Mips over $\mathcal{S}$ for $q$, then it is not necessary for the sketches of $u$ and $v$ to preserve the same fraction $\alpha$ of their $\ell_1$ mass. Instead, it makes more sense for $\alpha$ to *adapt* to the likelihood of $u$ and $v$ being the solution to Sparse Mips: if $\tilde{v}$ is the $\alpha_v$-MSS of $v$ and $\tilde{u}$ is the $\alpha_u$-MSS of $u$, $\alpha_u < \alpha_v$ is a reasonable choice.

While the extreme example above rarely occurs, the property that some vectors are more likely to be the solution of Sparse Mips than others appears to hold in many practical datasets and query distributions. For instance, consider the benchmark MsMarco v1 dataset and its queries: nearest neighbors of each query generally have a larger $\ell_1$ norm along a query's active dimensions (i.e., dimensions corresponding to nonzero coordinates) than father neighbors. To demonstrate this, we sample 1,000 queries from the dataset, find their nearest neighbor as well as the $k$-th nearest neighbor for some large $k > 1$. We then compute and plot in Figure 3 the ratio of the $\ell_1$ norm of the two neighbors. The plot clearly shows this phenomenon.

In this section, we wish to sketch the entire data matrix (i.e., the set $\mathcal{S}$) at once such that each vector is sketched in the context of other vectors in an adaptive manner. In particular, on the basis of the empirical observation above, we wish the amount of mass preserved for each vector to depend on the magnitude of its $\ell_1$ mass relative to other vectors'. The remainder of this section describes our sketching algorithm, which we call Set $\alpha$-Mass Subvector Sketch.

*3.4.1 Set $\alpha$-Mass Subvector Sketch.* We define the following notion first which is then used in the description of the algorithm. The algorithm is schematically described in Figure 4.
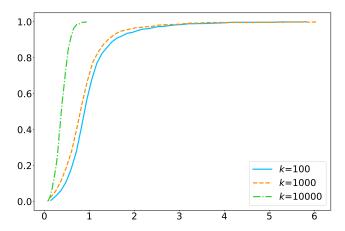
Fig. 3. We sample 1,000 queries from MsMarco v1, find their nearest neighbor (say, $u$) and the $k$-th nearest neighbor (say, $v$) for some large $k > 1$. We then compute the ratio of their $\ell_1$ norms, $\|v_{\mathcal{I}}\|_1 / \|u_{\mathcal{I}}\|_1$, where $\mathcal{I}$ is the set of active dimensions in the query. The figure shows the cumulative distribution (vertical axis) of this ratio (horizontal axis). As $k$ increases, so that we include farther points, the ratio becomes smaller, showing that points that are closer to the query have a larger $\ell_1$ norm along the active dimensions of the query.
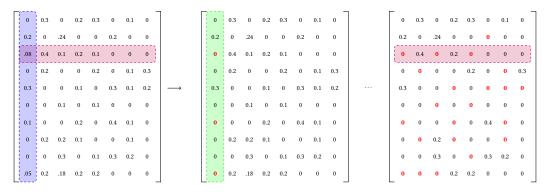


Fig. 4. Application of Set $\alpha$-MSS with $\alpha = 0.4$ to a toy set of sparse vectors $\mathcal{S}$ rendered as a matrix (leftmost): vectors become rows (purple box is a vector); dimensions columns (blue box is a dimension). In column $i$, the largest $\alpha$ portion of the nonzero entries are kept while the rest are reset. The center matrix shows that process for the first dimension. Rightmost is the final matrix $\tilde{\mathcal{S}}$ where each vector is sketched in the context of other vectors.

DEFINITION 5. *Let $L_i = \{j \text{ s.t. } u_i^{(j)} \neq 0 \; \forall \; u^{(j)} \in \mathcal{S}\}$ be the list of vectors in the set $\mathcal{S}$ whose $i$-th coordinate is nonzero. We say that $\mathcal{S}$ is $\rho$-dense ($\rho \in [0, 1]$) in the $i$-th dimension with $\rho = |L_i|/|\mathcal{S}|$.*

DEFINITION 6 (SET $\alpha$-MASS SUBVECTOR SKETCH). *Fix $\alpha \in (0, 1]$. Represent a set of sparse vectors $\mathcal{S} \subset \mathbb{B}_1^d$ in matrix form with vectors as rows and dimensions as columns. For each column $i$ with density $\rho_i$, keep the $\lambda_i = \lceil \alpha \rho_i |\mathcal{S}| \rceil$ largest values, and set all other values to $0$. The $j$-th row in the resulting matrix $\tilde{\mathcal{S}}$ is the Set $\alpha$-MSS of the $j$-th vector.*

### 3.4.2 Analysis of Set $\alpha$-MSS.

In this section, we show that Set $\alpha$-MSS has the property that its parameter $\alpha$ adapts to vectors based on their $\ell_1$ mass. To do that, we first establish that the Set $\alpha$-MSS of a fixed vector $u$ in a randomly sampled set $\mathcal{S}$ is its $\alpha$-MSS in expectation.

THEOREM 3.4. *In expectation over sets* $\mathcal{S} \sim \mathcal{D}$, *where dimensions are independent or weakly-dependent, the Set* $\alpha$-MSS $\tilde{u}$ *of* $u$ *is the* $\alpha$-MSS *of* $u$: $\mathbb{E}_{\mathcal{S} \sim \mathcal{D}}[\|\tilde{u}\|_1] = \alpha\|u\|_1$.

PROOF. Each $u_i$ of a vector $u$ is included in the Set $\alpha$-MSS of $u$ if its rank among the $i$-th dimension of all vectors in $\mathcal{S}$ is less than $\lambda_i$: $r(u_i) \le \lambda_i$. The probability of that event is equal to $\lambda_i/\rho_i|\mathcal{S}| = \alpha$. We have that:

$$\mathbb{E}_{\mathcal{S} \sim \mathcal{D}}[\|\tilde{u}\|_1] = \sum_{i:\, u_i \neq 0} u_i \, \mathbb{P}[\mathbb{1}_{r(u_i) \le \lambda_i}] = \sum_{i:\, u_i \neq 0} u_i \alpha = \alpha\|u\|_1.$$

□

The result above is in expectation. We are now interested in the deviation of $\|\tilde{u}\|_1$ from $\alpha\|u\|_1$. We present the following result to that end:

THEOREM 3.5. *Given the conditions of Theorem 3.4,* $\epsilon \in (0, 1)$, *and a fixed* $\mathcal{S}$, *we have for each* $u \in \mathcal{S}$ *that:*

$$\mathbb{P}[\|\tilde{u}\|_1 - \alpha\|u\|_1 \ge \epsilon] \le \exp\left(-\frac{\epsilon^2}{\alpha(1-\alpha)\|u\|_2^2 + \frac{2}{3}\|u\|_1\epsilon}\right).$$

PROOF. Given the variance of $\|\tilde{u}\|_1$, and noting that $\|\tilde{u}\|_1 \in (0, \|u\|_1]$ is a bounded random variable, we can apply the one-sided Bernstein concentration inequality to derive the bound in the theorem. All that is left then is to derive the variance, which we do as follows. Letting $r(u_i)$ take the same definition as in the proof of Theorem 3.4:

$$\text{Var}[\|\tilde{u}\|_1] = \sum_{i:\, u_i \neq 0} \text{Var}[u_i \mathbb{1}_{r(u_i) \le \lambda_i}] = \sum u_i^2 \text{Var}[\mathbb{1}_{r(u_i) \le \lambda_i}] = \|u\|_2^2(\alpha - \alpha^2).$$

That completes the proof.

□

It should be noted that, due to the independence assumption, the above results hold even if we restricted the space to an arbitrary subspace. In particular, we may restrict the space to the subspace spanned by the basis vectors corresponding to the nonzero coordinates of a query. Rewriting Theorem 3.4 gives that $\mathbb{E}_{\mathcal{S} \sim \mathcal{D}}[\|\tilde{u}_{\mathcal{I}}\|_1] = \alpha\|u_{\mathcal{I}}\|_1$, where $\mathcal{I}$ is the set of active dimensions in a query.

Let us examine what this bound implies. Suppose that $\alpha \ll 1$, so that $\alpha\|u\|_1 \ll \epsilon$. In the MSMARCO v1 dataset, for instance, the maximum $\ell_1$ norm is 0.7, so that choosing $\alpha = 0.1$ gives sketches whose $\ell_1$ norm is less than 0.07. If $\|u\|_2 \ll 1$ (this value is less than $1e-3$ in MSMARCO v1) then the first term in the denominator of the exponent vanishes, so that we are left with $\exp(-3\epsilon/2\|u\|_1)$.

Now, fix $\epsilon$ and take vectors $u$ and $v$ such that $\|u\|_1 \ll \|v\|_1$. Then the upper-bound on the deviation is larger for $v$ than it is for $u$, implying that, with high probability, the Set $\alpha$-MSS of $v$ results is a $\alpha_v$-MSS with $\alpha_v > \alpha$, while the Set $\alpha$-MSS of $u$ results in an $\alpha_u$-MSS of $u$ with $\alpha_u \approx \alpha$.

That gives us the effect we sought to induce. To see that, suppose $\|u\|_1 \ll \|v\|_1$ implies with high probability that $\langle q, u \rangle < \langle q, v \rangle$—we have verified this earlier. Because $\alpha_u < \alpha_v$, more coordinates of $v$ are preserved in its sketch with high probability than coordinates of $u$. As such, according to Theorem 3.3, the estimate of $\langle q, v \rangle$ is more accurate than $\langle q, u \rangle$, with high probability.

## 4 SEISMIC: OUR SPARSE MIPS ALGORITHM

In this section, we describe our Sparse MIPS algorithm dubbed SEISMIC. The algorithm relies on three data structures: an inverted index, a forward index, and (optionally) a $\kappa$-Nearest Neighbor graph ($\kappa$-NN graph). In a nutshell, we use the inverted index (i.e., a mapping from dimensions to an inverted list: the id of data points whose coordinate in that dimension is nonzero) to identify a set of approximate nearest neighbors; the $\kappa$-NN graph (i.e., a graph where each node represents a data point, with edges that connect it to its $\kappa$ nearest neighbors) to refine and expand that approximate
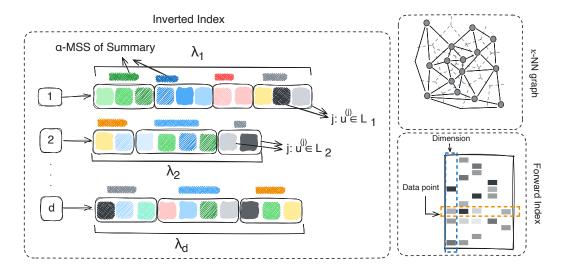
Fig. 5. The design of SEISMIC. We apply Set $\alpha$-MSS to the vector collection and build an inverted index. Each inverted list is then independently partitioned into geometrically-cohesive blocks, where a block is a set of data point identifiers along with an $\alpha$-MSS of a "summary" vector. The inner product of a query with the summary approximates the inner product attainable with the data points in that block. The forward index stores the complete vectors (i.e., all nonzero dimensions and coordinates). The $\kappa$-NN graph, which is stored as a table, records a mapping from a data point identifier to its $\kappa$ closest neighbors (by inner product).

set; and the forward index (i.e., a mapping from data point id to its full vector representation) for exact inner product computation. Figure 5 gives an overview of the overall design.

SEISMIC is novel in the following ways. First, by applying Set $\alpha$-MSS to a vector collection, it significantly reduces the size of the collection. Second, it partitions inverted lists into geometrically-cohesive blocks to reduce the number of data points for which it must compute exact inner product. Finally, it attaches a *summary* to each block, whose inner product with a query approximates the inner product of the query with data points referenced in the block. We expand on these in detail in the following sections.

## 4.1 Inverted Index

*4.1.1 Sketching.* SEISMIC heavily relies on the Set $\alpha$-MSS algorithm—with a fixed $\alpha$, our first hyperparameter—as introduced in Section 3.4. The sketch was designed so that a small subset of a sparse vector's nonzero coordinates approximates its inner product with a query point with arbitrary accuracy, while preserving rankings of data points with high probability.

*4.1.2 Blocking of Inverted Lists.* SEISMIC also introduces a novel blocking strategy on inverted lists. It partitions the $i$-th inverted list, $L_i$, into $\beta|L_i|$ small blocks—our second hyperparameter. The rationale behind a blocked organization of an inverted list is to group together data points that are *similar*, so as to facilitate a *dynamic pruning* strategy, to be described shortly.

We defer the determination of similarity to a clustering algorithm. In other words, the data points referenced in the $i$-th inverted list are given as input to a clustering algorithm, which subsequently partitions them into $\beta|L_i|$ clusters. Each cluster is then turned into one block, consisting of the id of data points that belong to the same cluster. Conceptually, each block is "atomic" in the following sense: if the dynamic pruning algorithm decides we must visit a block, *all* the data points referenced in that block are fully evaluated (i.e., their inner product with a query is computed exactly).

We note that any geometrical clustering algorithm may be readily used. We use a shallow variant [12] of K-Means as follows. Given an inverted list $L_i$, we uniformly-randomly sample $\beta|L_i|$ vectors from it, $\{\mu^{(k)}\}_{k=1}^{\beta|L_i|}$, and use their corresponding data point as cluster representatives. For each $j \in L_i$, we find $k^* = \arg\max_k \langle u^{(j)}, \mu^{(k)} \rangle$, and assign $j$ to the $k^*$-th cluster.

*4.1.3 Per-block Summary Vectors.* The query processing algorithm, to be presented later, requires an efficient way to determine if a block should be fully evaluated. To that end, Seismic leverages the concept of a *summary* vector: a $d$-dimensional sparse vector that "represents" the data points referenced in a block. The summary vectors are stored in the inverted index, one per block, and are meant to serve as an efficient way to compute an approximation of the inner product between a query and the data points referenced in each block.

One realization of this idea is to upper-bound the full inner product attainable by data points referenced in a block. In other words, the $i$-th coordinate of the summary vector of a block would contain the maximum $u_i$ among the data points referenced in that block.

More precisely, our summary function $\phi : 2^S \to [0, 1]^d$ takes a block $B$ from the universe of all blocks, and produces a vector whose $i$-th coordinate is simply:

$$\phi(B)_i = \max_{j \in B} u_i^{(j)}. \tag{2}$$

This summary is *conservative*: its inner product with the query is no less than the inner product between the query and any of the data points referenced in that block: $\langle q, \phi(B) \rangle \geq \langle q, u^{(j)} \rangle$ for all $j \in B$ and an arbitrary query $q$.

The caveat, however, is that the number of nonzero entries in summary vectors grows quickly with the block size. That is the source of two potential issues: 1) the space required to store summaries increases; and 2) as inner product computation takes time proportional to the number of nonzero entries, the time to evaluate a block could become unacceptably high.

We address that caveat by sketching the summary vectors using $\gamma$-MSS—$\gamma$ being our third hyperparameter—with the understanding that sketching takes away the conservatism property of the summary vectors. We further reduce the size of summaries by applying scalar quantization. With the goal of reserving a single byte for each value, we subtract the minimum value $m$ from each summary entry, and divide the resulting range into 256 sub-intervals of equal size. A value in the summary is replaced with the index of the sub-interval it maps to. To reconstruct a value approximately, we multiply the id of its sub-interval by the size of the sub-intervals, then add $m$.

## 4.2 Forward Index

As noted earlier, Seismic blends together three data structures. The first is an inverted index that tells us which data points to examine. To make it practical, we applied several layers of approximations that allow us to gain efficiency with a possible loss in accuracy. A forward index (built from the original vector collection) helps correct those errors and offers a way to compute the exact inner products between a query and the data points identified for full evaluation.

We must note that, data points referenced in the same block are not necessarily stored consecutively in the forward index. This is simply infeasible because the same data point may belong to different inverted lists and, thus, to different blocks. Because of this layout, computing the inner products may incur many cache misses, which are detrimental to query latency. In our implementation, we extensively use prefetching instructions to mitigate this effect.

## 4.3 $\kappa$-nn Graph

Our third data structure provides information that allows Seismic to *refine* the candidate pool during query processing—we will describe that logic shortly. In particular, it constructs a graph

---

**Algorithm 2:** Indexing with SEISMIC.

---

**Input:** Collection $\mathcal{S} \in \mathbb{B}_1^d$ of sparse vectors; $\alpha$: parameter of Set $\alpha$-MSS to sketch $\mathcal{S}$;
$\beta \in (0, 1)$: factor determining the number of blocks per inverted list; $\gamma$: parameter of $\gamma$-MSS
for the summaries; $\kappa$: outgoing degree in the $\kappa$-NN graph.

**Result:** SEISMIC index.

1: $\mathcal{N}(u) \leftarrow \arg\max_{v \in \mathcal{S}}^{(\kappa)} \langle u, v \rangle$ for all $u \in \mathcal{S}$                    ▷ The $\kappa$-nn graph.

2: $\tilde{\mathcal{S}} \leftarrow$ Set $\alpha$-MSS of $\mathcal{S}$

3: **for** $i \in \{1, \ldots, d\}$ **do**

4:     $L_i \leftarrow \{j \mid u_i^{(j)} \neq 0, \ u^{(j)} \in \tilde{\mathcal{S}}\}$                    ▷ The inverted index

5:     CLUSTER $L_i$ into $\beta|L_i|$ partitions, $\{B_{i,j}\}_{j=1}^{\beta|L_i|}$                    ▷ The blocks

6:     **for** $1 \leq j \leq \beta$ **do**

7:         $S_{i,j} \leftarrow \gamma$-MSS of $\phi(B_{i,j})$                    ▷ Sketched summaries using Equation (2)

8: **return** $\mathcal{N}, \{L_i\}_{i=1}^d, \{B_{i,j}, S_{i,j}\} \ \forall_{1 \leq i \leq d, \ 1 \leq j \leq \beta|L_i|}$

---

from a collection of data points, wherein each node represents a data point with an outgoing edge to $\kappa$ other data points whose inner product with that node is maximal.

The resulting structure can be stored as a look-up table consisting of pairs of one data point id and a list of its $\kappa$ closest neighbors—our fourth hyperparameter. The $\kappa$-NN graph allows us to identify the nearest neighbors of a data point quickly. Formally, let us denote by $\mathcal{N}(u)$ the set of $\kappa$ closest data points to data point $u \in \mathcal{S}$:

$$\mathcal{N}(u) = \arg\max_{v \in \mathcal{S}}^{(\kappa)} \langle u, v \rangle. \tag{3}$$

Setting $\kappa = 0$ effectively removes the $\kappa$-NN graph from indexing and query processing altogether.

### 4.4 Recap of Indexing

We summarize the discussion above in Algorithm 2. When indexing a collection $\mathcal{S} \subset \mathbb{B}_1^d$, we first construct a $\kappa$-NN graph (Line 1). Next, we apply Set $\alpha$-MSS to the collection (Line 2). For every dimension $i \in \{1, \ldots, d\}$, we form its inverted list, recording only the data point identifiers (Line 4). Next, we apply clustering to each inverted list separately to derive at most $\beta|L_i|$ blocks (Line 5). Once data points are assigned to the blocks, we then build the $\gamma$-MSS of each block's summary (Line 7).

### 4.5 Query Processing

Algorithm 3 shows the query processing logic in SEISMIC. We take an $\alpha_q$-MSS $\tilde{q}$ the query $q$ (Line 1). We then utilize our dynamic pruning strategy (Lines 4–7) that allows the algorithm to skip blocks in the inverted lists corresponding to nonzero coordinates in $\tilde{q}$.

SEISMIC adopts a coordinate-at-a-time traversal (Line 3) of the inverted index. For each nonzero coordinate $i$ of $\tilde{q}$, the algorithm sorts the blocks in $L_i$ by the inner product of their summary vectors with $q$ (referred to as "summary score") in descending order. Note that, inner products of $q$ with summaries in an inverted list can be computed efficiently with sparse matrix-vector multiplication. Sorting the blocks is also cheap as $\beta|L_i|$ is typically small.

The algorithm then visits the blocks in the sorted order. The data points referenced in a block are fully evaluated if the block's summary score is greater than a fraction (denoted by heap_factor $\in$ (0, 1], the second hyperparameter) of the minimum inner product in the Min-HEAP. That means the forward index retrieves the complete data point referenced in the selected block and computes

---

**Algorithm 3:** Query processing with Seismic.

---

**Input:** $q$: query; $k$: number of points to return; $\alpha_q$: parameter to $\alpha_q$-MSS to sketch the query; heap_factor: a correction factor to rescale the summary inner product; $\mathcal{N}$, $L_i$'s and $(B_{i,j}, S_{i,j})$'s produced by Algorithm 2.

**Result:** A Heap with the top-$k$ data points.

1:   $\tilde{q} \leftarrow \alpha_q$-MSS of $q$
2:   Heap$\leftarrow \emptyset$
3:   **for all** $i$ *s.t.* $\tilde{q}_i \neq 0$ **do**
4:      $r_j \leftarrow \langle q, S_{i,j} \rangle \quad \forall_{1 \leq j \leq \beta|L_i|}$
5:      **for all** $B_{i,j} \in L_i$ sorted in decreasing order by $r_j$ **do**
6:        **if** Heap.len() $= k$ and $r_j < \frac{\text{Heap.min()}}{\text{heap\_factor}}$ **then**
7:          **continue**                                ▷ Skip the block
8:        **for** $d \in B_{i,j}$ **do**
9:          $p = \langle q, \text{ForwardIndex}[d] \rangle$      ▷ Fully evaluate the data point with id $d$
10:         **if** Heap.len() $< k$ or $p >$ Heap.min() **then**
11:          Heap.insert($(p, d)$)
12:         **if** Heap.len() $= k + 1$ **then**
13:          Heap.pop_min()
14:   $S \leftarrow$ the ids of the vectors in the Heap
15:   **for** $j \in S$ **do**
16:      **for** $v \in \mathcal{N}(u^{(j)})$ **do**
17:        $p = \langle q, \text{ForwardIndex}[v] \rangle$     ▷ Fully evaluate points in the expanded set $S_+$
18:        **if** Heap.len() $< k$ or $p >$ Heap.min() **then**
19:          Heap.insert($(p, v)$)
20:        **if** Heap.len() $= k + 1$ **then**
21:          Heap.pop_min()
22:   **return** Heap

---

their inner products with $q$. A data point whose inner product is greater than the minimum score in the Min-Heap is inserted into the heap.

Once the above concludes, we take the set of data points in the heap, denoted by $S$. As Lines 14–21 show, we then form the *expanded set* $S_+ = \bigcup_{u \in S} (\{u\} \cup \mathcal{N}(u))$, compute scores for data points in $S_+$, and return the top-$k$ subset. As noted previously, $\kappa = 0$ disables the expansion of $S$.

## 5 EXPERIMENTS

We now evaluate Seismic experimentally. Specifically, we are interested in investigating the performance of Seismic along the following axes: accuracy, latency, space usage, indexing time, and scalability.

We begin by introducing our setup, including datasets and embedding models. We then evaluate Seismic with $\kappa = 0$, so that the $\kappa$-nn graph is disabled. Next, we enable the $\kappa$-nn graph and study its impact on the evaluation axes above. We conclude our empirical evaluation by scaling up the collection of vectors by an order of magnitude.

### 5.1 Setup

*5.1.1 Datasets.* We experiment on three publicly-available datasets. These include:

- MsMarco v1 [59]: a collection of 8.8M text passages in English with a set of 6,980 queries;

- Natural Questions (NQ): a Beir [66] collection of 2.68M questions and 7,842 queries; and,
- MsMarco v2: a collection of 138 million passages with 3,903 queries.

We use the last dataset in our scalability study only.

*5.1.2  Sparse embedding models.* We evaluate all Sparse Mips algorithms on the embeddings of the above datasets using four embedding models:

- Splade [25]. Each nonzero entry is the importance weight of a term in the BERT [20] Word-Piece [71] vocabulary consisting of 30,000 terms. We use the `cocondenser-ensembledistil`[3] version of Splade. The number of nonzero entries in data points (query points) is, on average, 119 (43) for MsMarco v1; 153 (51) for NQ; and, 127 (44) for MsMarco v2.
- Efficient Splade [35]. Similar to Splade, but whose embeddings typically have a larger number of nonzero coordinates in data points and a smaller number in query points. For example, there are 181 (5.9) nonzero entries in MsMarco v1 data (query) points. We use the `efficient-splade-V-large`[4]. We refer to this model as E-Splade.
- Splade-v3 [36]. An improved variant of Splade that incorporates a modified objective and distillation strategy. The number of nonzero entries in data (query) points is, on average, 168 (24) for MsMarco v1.
- uniCoil-T5 [39, 42]. Expands a text passge with relevant terms generated by DocT5Query [60]. There are, on average, 68 (6) nonzero entries in MsMarco v1 data (query) points.

We note that, we generate the Splade and E-Splade embeddings using the original code published on GitHub.[5]

*5.1.3  Baselines.* We compare Seismic with six state-of-the-art Sparse Mips algorithms. Two of these are the winning solutions of the "Sparse Track" at the 2023 BigANN Challenge[6] at NeurIPS:

- GrassRMA: A graph-based method for dense vectors adapted to sparse vectors that appears in the BigANN challenge as "sHnsw."[7]
- PyAnn: Another graph-based ANN solution.[8]

The other baselines are as follows:

- Ioqp [46]: Impact-sorted query processor written in Rust. We choose Ioqp because it is known to outperform JASS [41], another widely-adopted impact-sorted query processor.
- SparseIvf [8]: An inverted index where lists are partitioned into blocks through clustering. At query time, after finding the $N$ closest clusters to the query, a coordinate-at-a-time algorithm traverses the inverted lists. The solution is approximate because only the data points that belong to top $N$ clusters are considered.
- Pisa [54]: An inverted index-based library based on `ds2i` [62] that uses highly-optimized blocked variants of WAND. Pisa is *exact* as it traverses inverted lists in a rank-safe manner.
- kANNolo [19]: a new implementation of HNSW [50] written in Rust for approximate nearest neighbors search, targeting both dense and sparse domains.[9] We choose this solution because, as Delfino *et al.* show, it outperforms both GrassRMA and PyAnn in terms of search time.

---

[3]https://huggingface.co/naver/splade-cocondenser-ensembledistil
[4]Checkpoints at https://huggingface.co/naver/efficient-splade-V-large-doc and https://huggingface.co/naver/efficient-splade-V-large-query.
[5]https://github.com/naver/splade
[6]https://big-ann-benchmarks.com/neurips23.html
[7]C++ code is publicly available at https://github.com/Leslie-Chung/GrassRMA.
[8]C++ code is publicly available at https://github.com/veaaaab/pyanns.
[9]Code obtained from https://github.com/TusKANNy/kannolo

We also considered the method by Lassance *et al.* [37]. Their approach statically prunes either inverted lists (by keeping $p$-quantile of elements), data points (by keeping a fixed number of top entries), or all coordinates whose value is below a threshold. While simple, their method is only able to speed up query processing by 2–4× over Pisa on E-Splade embeddings of MsMarco v1. We found it to be ineffective on Splade and generally far slower than GrassRMA and PyAnn. As such we do not include it in our discussion.

*5.1.4 Evaluation Metrics.* We evaluate all methods using three metrics:

- Latency ($\mu$sec.). The time elapsed, in *microseconds*, from the moment a query vector is presented to the index to the moment it returns the approximate set of $k$ nearest data vectors running in single thread mode. Latency does not include embedding time.
- Accuracy. The percentage of true nearest neighbors recalled in the returned set, per Definition 2. By measuring the recall of an approximate set given the exact top-$k$ set, we study the impact of the different levers in an algorithm on its overall accuracy as a retrieval engine.
- Index size (MiB). The space the index occupies in memory.

*5.1.5 Reproducibility and Hardware Details.* We implemented Seismic in Rust.[10] We compile Seismic by using the version 1.77 of Rust and use the highest level of optimization made available by the compiler. We conduct experiments on a server equipped with one Intel i9-9900K CPU with a clock rate of 3.60 GHz and 64 GiB of RAM. The CPU has 8 physical cores and 8 hyper-threaded ones. We query the index using a single thread.

## 5.2 Seismic without the $\kappa$-nn Graph

We now present our experimental results for a configuration of Seismic where $\kappa = 0$, effectively removing the $\kappa$-nn graph from indexing and search subroutines.

*5.2.1 Hyperparameter tuning.* We build Ioqp and Pisa indexes using Anserini[11] and apply recursive graph bisection [47]. For Ioqp, we vary the *fraction* (of the total collection) hyperparameter in [0.1, 1] with step size of 0.05. For SparseIvf, we sketch data points using Sinnamon$_{\text{Weak}}$ and a sketch size of 1,024, and build $4\sqrt{N}$ clusters, where $N$ is the number of data points in the collection. For GrassRMA and PyAnn, we build different indexes by running all possible combinations of $ef_c \in \{1000, 2000\}$ and $M \in \{16, 32, 64, 128, 256\}$. During search we test $ef_s \in [5, 100]$ with step size 5, then [100, 400] with step 10, [100, 1000] with step 100, and finally [1000, 5000] with step 500. We apply early stopping when accuracy saturates.

Our grid search for Seismic on MsMarco v1 is over: $\alpha \in \{0.05, 0.1, 0.15, 0.2, 0.3, 0.4, 0.5\}$; $\beta \in \{0.1, 0.2, 0.3\}$ with step $x$, and $\gamma \in [0.1, 0.5]$ with step size 0.1. Best results on Splade are achieved when $\alpha = 0.1$, $\beta = 0.3$, and $\gamma = 0.4$. Best results with E-Splade with $\alpha = 0.1$, $\beta = 0.2$, and $\gamma = 0.4$. Best results for uniCoil-T5 with $\alpha = 0.2$, $\beta = 0.2$, and $\gamma = 0.4$. The grid search for Seismic on NQ is over: $\alpha \in \{0.05, 0.1, 0.15, 0.2, \}$; $\beta \in \{0.1, 0.2, 0.3\}$, and $\gamma \in \{0.4, 0.5, 0.6\}$. Best results are achieved when $\alpha = 0.1$, $\beta = 0.3$, and $\gamma = 0.5$. Seismic uses 8-bit scalar quantization for summaries, and matrix multiplication to efficiently compute the product of a query vector with all quantized summaries of an inverted list.

*5.2.2 Accuracy-Latency Trade-off.* Table 1 details Anns performance in terms of average per-query latency for Splade, E-Splade, and uniCoil-T5 on MsMarco v1, and Splade on NQ. We frame the results as the trade-off between effectiveness and efficiency. In other words, we report mean per-query latency at a given accuracy level.

---

[10]Our code is publicly available at https://github.com/TusKANNy/seismic.
[11]https://github.com/castorini/anserini

Table 1. Mean latency ($\mu$sec/query) at different accuracy@10 cutoffs with speedup (in parenthesis) gained by SEISMIC over others.

| | SPLADE on MsMarco v1 | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Accuracy (%) | 90 | | 91 | | 92 | | 93 | | 94 | | 95 | |
| IOQP | 17,423 | (96.8×) | 17,423 | (96.8×) | 18,808 | (90.0×) | 21,910 | (104.8×) | 24,382 | (101.6×) | 31,843 | (98.3×) |
| SPARSEIVF | 4,169 | (23.2×) | 4,984 | (27.7×) | 6,442 | (30.8×) | 7,176 | (34.3×) | 8,516 | (35.5×) | 10,254 | (31.6×) |
| GRASSRMA | 807 | (4.5×) | 867 | (4.8×) | 956 | (4.6×) | 1,060 | (5.1×) | 1,168 | (4.9×) | 1,271 | (3.9×) |
| PYANN | 489 | (2.7×) | 539 | (3.0×) | 603 | (2.9×) | 654 | (3.1×) | 845 | (3.5×) | 1,016 | (3.1×) |
| KANNOLO | 537 | (3.0×) | 586 | (3.3×) | 629 | (3.0×) | 673 | (3.2×) | 743 | (3.1×) | 836 | (2.6×) |
| **SEISMIC** ($\kappa = 0$) | 180 | – | 180 | – | 209 | – | 209 | – | 240 | – | 324 | – |
| | E-SPLADE on MsMarco v1 | | | | | | | | | | | |
| IOQP | 7,857 | (45.9×) | 8,382 | (46.8×) | 8,892 | (49.7×) | 9,858 | (51.3×) | 10,591 | (35.9×) | 11,536 | (37.8×) |
| SPARSEIVF | 4,643 | (27.2×) | 5,058 | (28.3×) | 5,869 | (32.8×) | 6,599 | (34.4×) | 7,555 | (25.6×) | 8,962 | (29.4×) |
| GRASSRMA | 2,074 | (12.1×) | 2,658 | (14.8×) | 2,876 | (16.1×) | 3,490 | (18.2×) | 4,431 | (15.0×) | 5,141 | (16.9×) |
| PYANN | 1,685 | (9.9×) | 1,702 | (9.5×) | 2,045 | (11.4×) | 2,409 | (12.5×) | 3,119 | (10.6×) | 4,522 | (14.8×) |
| KANNOLO | 1,310 | (7.7×) | 1,432 | (8.0×) | 1,649 | (9.2×) | 1,864 | (9.7×) | 2,369 | (8.0×) | 2,864 | (9.4×) |
| **SEISMIC** ($\kappa = 0$) | 171 | – | 179 | – | 179 | – | 192 | – | 295 | – | 305 | – |
| | UNICOIL-T5 on MsMarco v1 | | | | | | | | | | | |
| IOQP | 22,278 | (210.2×) | 25,060 | (217.9×) | 26,541 | (215.8×) | 30,410 | (230.4×) | 33,327 | (173.6×) | 34,061 | (168.6×) |
| SPARSEIVF | 6,375 | (60.1×) | 7,072 | (61.5×) | 8,192 | (66.6×) | 9,207 | (69.8×) | 10,306 | (53.7×) | 12,308 | (60.9×) |
| GRASSRMA | 1,318 | (12.4×) | 1,434 | (12.5×) | 1,812 | (14.7×) | 2,004 | (15.2×) | 2,168 | (11.3×) | 2,668 | (13.2×) |
| PYANN | 1,133 | (10.7×) | 1,456 | (12.7×) | 1,741 | (14.2×) | 1,755 | (13.3×) | 2,061 | (10.7×) | 2,973 | (14.7×) |
| KANNOLO | 1,085 | (10.2×) | 1,188 | (10.3×) | 1,301 | (10.6×) | 1,468 | (11.1×) | 1,691 | (8.8×) | 1,961 | (9.7×) |
| **SEISMIC** ($\kappa = 0$) | 106 | – | 115 | – | 123 | – | 132 | – | 192 | – | 202 | – |
| | SPLADE on NQ | | | | | | | | | | | |
| IOQP | 8,313 | (47.2×) | 8,854 | (43.0×) | 9,334 | (38.3×) | 11,049 | (45.3×) | 11,996 | (42.2×) | 14,180 | (42.8×) |
| SPARSEIVF | 3,862 | (21.9×) | 4,309 | (20.9×) | 4,679 | (19.2×) | 5,464 | (22.4×) | 6,113 | (21.5×) | 6,675 | (20.2×) |
| GRASSRMA | 1,000 | (5.7×) | 1,138 | (5.5×) | 1,208 | (5.0×) | 1,413 | (5.8×) | 1,549 | (5.5×) | 2,091 | (6.3×) |
| PYANN | 610 | (3.5×) | 668 | (3.2×) | 748 | (3.1×) | 870 | (3.6×) | 1,224 | (4.3×) | 1,245 | (3.8×) |
| KANNOLO | 471 | (2.7×) | 521 | (2.5×) | 567 | (2.3×) | 606 | (2.5×) | 702 | (2.5×) | 805 | (2.4×) |
| **SEISMIC** ($\kappa = 0$) | 176 | – | 206 | – | 244 | – | 244 | – | 284 | – | 331 | – |

The results on these datasets show SEISMIC's remarkable relative efficiency, reaching a latency that is often one to two orders of magnitude smaller. Overall, SEISMIC consistently outperforms all baselines at all accuracy levels, including GRASSRMA and PYANN, which in turn perform better than other inverted index-based baselines—confirming the findings of the BigANN Challenge.

We make a few additional observations. IOQP appears to be the slowest method across datasets. This is not surprising considering IOQP is not designed with the distributional properties of sparse embeddings in mind. SPARSEIVF generally improves over IOQP, but SEISMIC speeds up query processing further. In fact, the minimum speedup over IOQP (SPARSEIVF) on MsMarco v1 is 84.6× (22.3×) on SPLADE, 30.7× (20.9×) on E-SPLADE, and 181× (54.8×) on UNICOIL-T5.

SEISMIC consistently outperforms GRASSRMA and PYANN by a substantial margin, ranging from 2.6× (SPLADE on MsMarco v1) to 21.6× (E-SPLADE on MsMarco v1) depending on the level of accuracy. In fact, as accuracy increases, the latency gap between SEISMIC and the two graph-based methods widens. This gap is much larger when query vectors are sparser, such as with E-SPLADE embeddings. That is because, when queries are highly sparse, inner products between queries and data points become smaller, reducing the efficacy of a greedy graph traversal. As one statistic, PYANN over E-SPLADE embeddings of MsMarco v1 visits roughly 40,000 data points to reach 97% accuracy, whereas SEISMIC evaluates just 2,198 data points.

Table 2. Index size and build time.

| SPLADE on MsMARCO v1 | | |
| --- | --- | --- |
| Model | Index size (MiB) | Index build time (min.) |
| IOQP | 2,195 | - |
| SPARSEIVF | 8,830 | 44 |
| GRASSRMA | 10,489 | 267 |
| PYANN | 5,262 | 137 |
| kANNOLO | 7,369 | 147 |
| **SEISMIC** ($\kappa = 0$) | 6,598 | 2 |

We highlight that PISA is the slowest (albeit, *exact*) solution. On MsMARCO v1, PISA processes queries in about 100,325 microseconds on SPLADE embeddings. On E-SPLADE and UNICOIL-T5, its average latency is 7,947 and 9,214 microseconds, respectively. We note that its high latency on SPLADE is largely due to the much larger number of nonzero entries in queries.

Finally, we must note that we have conducted a comparison of SEISMIC with another recent algorithm named BMP [56] for completeness. While BMP performs similarly to SEISMIC in terms of the index size and accuracy, it has a considerably larger per-query latency. As one data point, setting the memory budget to 8GB and performing search over the SPLADE embeddings of MsMARCO v1, BMP executes each query in 3.2 milliseconds on average while SEISMIC takes 0.7 milliseconds per query to reach the same level of accuracy.

*5.2.3 Space and Build Time.* Table 2 records the time it takes to index the entire MsMARCO v1 collection embedded with SPLADE with different methods, and the size of the resulting index. We left out the build time for IOQP because its index construction has many external dependencies (such as Anserini and graph bisection) that makes giving an accurate estimate difficult. We perform this experiment on a machine with two Intel Xeon Silver 4314 CPUs clocked at 2.40GHz, with 32 physical cores plus 32 hyper-threaded ones and 512 GiB of RAM. We build the indexes by using multi-threading parallelism with 64 cores.

Trends for other datasets are similar to those reported in Table 2. Notably, indexes produced by approximate methods are larger. That makes sense: using more auxiliary statistics helps narrow the search space dynamically and quickly. Among the highly efficient methods, the size of SEISMIC's index is mild, especially compared with GRASSRMA. Importantly, SEISMIC builds its index in a fraction of the time it takes PYANN or GRASSRMA to index the collection.

## 5.3 SEISMIC with the $\kappa$-NN Graph

Now that we have compared the performance of SEISMIC against all baseline algorithms, we move to present experimental results for a configuration of SEISMIC *with* the $\kappa$-NN graph (i.e., $\kappa > 0$). As we will see, including the $\kappa$-NN graph leads to better accuracy and faster search, but at the expense of increased index construction time and size.

We note that, because in the preceding section we concluded that SEISMIC is far faster and more accurate than baseline ANNS algorithms, we limit our comparison in this section to different flavors of SEISMIC: one without the $\kappa$-NN graph, and another with.

Before we present our results, let us elaborate on how we construct the $\kappa$-NN graph. Constructing the (exact) graph is expensive due to its quadratic time complexity. As such, we relax the construction to an approximate (but almost-exact) graph: we connect each data point to its approximate set of

| Embeddings | Budget | Accuracy (%) | 90 | 91 | 92 | 93 | 94 | 95 | 96 | 97 | 98 | 99 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SPLADE | 1.5× | SEISMIC ($\kappa = 0$) | 183 (1.5×) | 207 (1.5×) | 236 (1.5×) | 313 (2.0×) | 348 (1.8×) | 517 (2.3×) | | – | – | – |
| | | SEISMIC ($\kappa > 0$) | 126 | 136 | 154 | 154 | 196 | 226 | 333 | 539 | – | – |
| | 2× | SEISMIC ($\kappa = 0$) | 174 (1.4×) | 174 (1.3×) | 209 (1.6×) | 209 (1.5×) | 249 (1.6×) | 255 (1.6×) | 301 (1.8×) | 349 (1.7×) | 480 (2.0×) | 715 (1.8×) |
| | | SEISMIC ($\kappa > 0$) | 121 | 134 | 134 | 143 | 156 | 156 | 168 | 203 | 245 | 387 |
| SPLADE-v3 | 1.5× | SEISMIC ($\kappa = 0$) | 298 (1.1×) | 298 (1.1×) | 298 (1.1×) | 298 (1.1×) | 298 (1.1×) | 438 (1.6×) | 438 (1.6×) | | – | – |
| | | SEISMIC ($\kappa > 0$) | 281 | 281 | 281 | 281 | 281 | 281 | 281 | 281 | 424 | – |
| | 2× | SEISMIC ($\kappa = 0$) | 152 (0.6×) | 152 (0.6×) | 176 (0.7×) | 176 (0.7×) | 201 (0.8×) | 222 (0.9×) | 254 (1.0×) | 288 (1.1×) | 328 (1.2×) | 567 (1.5×) |
| | | SEISMIC ($\kappa > 0$) | 253 | 253 | 253 | 253 | 253 | 253 | 253 | 253 | 266 | 371 |

Table 3. Mean latency ($\mu$sec/query) at different accuracy@10 cutoffs with speedup (in parenthesis) gained by adding the $\kappa$-NN graph on MsMarco v1. The "Budget" column indicates the memory budget as a multiple of the size of the forward index.

top-$\kappa$ points. To find the $\kappa$ approximate set, we use SEISMIC (without the $\kappa$-NN graph). Storing the graph takes $(\lfloor \log_2(N - 1) \rfloor + 1)N\kappa$ bits, where $N$ is the size of the collection.

*5.3.1 Hyperparameter tuning.* When building the two SEISMIC indexes, we first fix a memory budget as a multiple of the size of the forward index. We then sweep the hyperparameters as follows to find the best configuration that results in an index no larger than the budget: $\alpha \in \{0.05, 0.1, 0.15, 0.2, 0.3, 0.4, 0.5\}$; $\beta \in \{0.1, 0.2, 0.3\}$, and $\gamma \in [0.1, 0.6]$ with step size 0.1 and when the $\kappa$-NN graph is enabled, $\kappa \in \{10, 20, 30, 40, 50\}$. We set $\alpha_q \in [0.1, 0.9]$ with step size 0.1 and heap_factor $\in \{0.7, 0.8, 0.9, 1.0\}$, and report the best configuration.

*5.3.2 Results.* Table 3 compares the two configurations of SEISMIC on MsMarco v1, vectorized using two of the embedding models. In these experiments, we consider two memory budgets expressed as multiples of the size of the collection (4GB for SPLADE and 5.6GB for SPLADE-v3). We choose the best index configuration that respects the budget and report the fastest configuration reaching the accuracy cutoffs from 90% to 99%. We also report the speedup in parenthesis gained by adding the $\kappa$-NN graph.

It is clear that adding the $\kappa$-NN graph to SEISMIC leads to faster search across accuracy cutoffs. On SPLADE, this speedup is 2.2× with a budget of 1.5×, and is up to 1.7× with a 2.0× memory budget. The same trend can be observed on SPLADE-v3. Note that, because SPLADE-v3 embeddings are less sparse, the memory impact of storing the $\kappa$-NN graph is smaller.

Interestingly, the "exact" algorithm PISA [54] caps at 99% accuracy due to quantization, which enables significant memory savings. At 99% cutoff, SEISMIC (with $\kappa$-NN graph) takes 409$\mu s$, which is *two orders of magnitude* faster than PISA's 95,818$\mu s$ per query.

## 5.4 Scalability of SEISMIC

In this section, we apply KANNOLO and SEISMIC to MsMarco v2, a massive dataset of 138 million embedding vectors, and compare their performance. As before, let us explain our protocol for tuning the hyperparameters of the algorithm before presenting experimental results.

*5.4.1 Hyperparameter tuning.* As in the preceding section, we allocate memory budgets for the index as multiples of the dataset size, which is about 66GB with 16-bit floats for values and 16-bit integers for components. We only consider hyperparameters that result in indexes that are up to 1.5× the dataset size in one scenario, and 2× in another.

KANNOLO has the following hyperparameters: $M$ (number of neighbors per node), and $ef_c$ (number of nodes scanned to build a node's neighborhood). We sweep $M \in \{16, 32, 64\}$ and $ef_c = 500$. For SEISMIC, we set: $\alpha \in 0.1, 0.2, 0.3$ $\beta = 0.4$, and $\gamma = 0.4$. We let $\kappa$ take values in $\{10, 20\}$.

Given an index, we process queries using the following hyperparameters and evaluate each configuration separately. For SEISMIC, we vary $\alpha_q$ in $[0.1, 0.9]$, with step 0.1, heap_factor in $\{0.6, .., 1.0\}$. For KANNOLO, we vary $ef_s$ in $[10, 50]$ with step 5, $[50, 100]$ with step 10, $[100, 1500]$
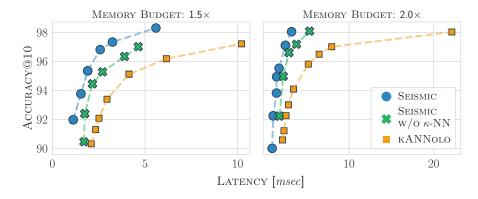
Fig. 6. Comparison of κANNolo and Seismic (with and without κ-nn graph) by accuracy@10 as a function of query latency. We allow hyperparameters that result in an index whose size is at most 1.5× (left) or 2× (right) the size of the dataset. Note that, unlike previous figures, latency is measured in *milliseconds*.

with step 100. For both methods, we consider the shortest time it takes the algorithm to reach top-10 accuracy of $\{90\%, 91\%, \dots, 99\%\}$.

*5.4.2 Hardware details.* Due to a much larger collection size, we use a different machine for experiments in this section. We run experiments on a NUMA server with 1TiB of RAM and four Intel Xeon Gold 6252N CPUs (2.30 GHz), with a total of 192 cores (96 physical and 96 hyper-threaded). Using the numactl[12] tool, we enforce that the execution of the retrieval algorithms is done on a single CPU and its local memory. In this way, we avoid performance degradation due to non-uniform memory accesses across different CPUs. We note that, each CPU node has enough local memory (256GiB) to store entire indexes.

*5.4.3 Accuracy and latency.* Our first experiments, visualized in Figure 6, compare κANNolo and Seismic by accuracy as a function of query latency for top-10 Anns. Seismic remains Pareto-efficient relative to κANNolo in both memory budget settings.

Interestingly, when we allow the algorithms to use more memory (i.e., 2× the dataset size), the gap between Seismic and κANNolo widens as accuracy increases. That is thanks to the κ-nn graph, which gives an advantage in high accuracy scenarios (>= 97%), as validated by a comparison of Seismic with and without the κ-nn graph. As an example, in the 1.5× scenario, Seismic with (without) κ-nn graph is 3.0× (2.3×) faster than κANNolo at 97% accuracy cutoff. The speedups increase to 6.9× (4.4×) in the 2× scenario at 98% accuracy cutoff.

Observe that the κ-nn graph has a notable impact on memory: Storing it costs $(\lfloor \log_2(N-1) \rfloor + 1)N\kappa$ bits, with $N$ denoting the size of the collection, which translates to 27 bits per document in our setup. This analysis suggests that reducing the memory impact of the κ-nn graph, for example, using delta encoding to store the ids of the documents, may help Seismic become even faster at a given memory budget.

*5.4.4 Index construction time and size.* Seismic builds the MsMarco v1 index much more quickly than other methods as previously shown in Table 2. This advantage holds strong on MsMarco v2, as we show in Table 4: Seismic (without the κ-nn graph) builds its index in about 30 minutes, while κANNolo in the sparse domain takes about 12.5 hours. The reported configuration refers to the 1.5× memory budget scenario.

---

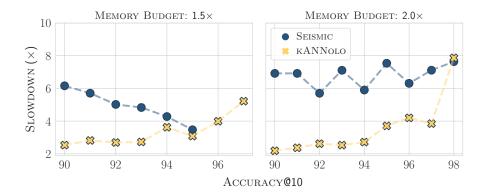[12]https://linux.die.net/man/8/numactl

Fig. 7. Scaling laws of Seismic and sparse kANNolo. For each accuracy cutoff, we measure the ratio between the latency of a method on MsMarco v2 and on MsMarco v1.

Enabling the $\kappa$-nn graph in Seismic substantially increases the indexing time. That is to be expected. Despite the construction of the $\kappa$-nn graph being approximated with Seismic itself, this process involves the daunting task of searching through approximately 140 million points. Moreover, as the collection size grows, not only would one have more queries to search, but there would be a larger number of documents to search over. As such, the inference benefits of the $\kappa$-nn graph shown in Figure 6 come at a high cost at indexing time.

### 5.5  Scaling Laws

In Figure 7, we present the search time scaling laws for Seismic and kANNolo. For each accuracy threshold on the horizontal axis, we plot on the vertical axis the ratio of the search time on MsMarco v2 to search time on MsMarco v1. In doing so, we report the amount of slowdown between the two datasets.

Considering the fact that MsMarco v2 is approximately 15× larger than MsMarco v1, from Figure 7 we learn that both methods scale effectively with dataset size. Concerning Seismic, all configurations show a slowdown of less than 8×, except for accuracy thresholds of 97 and 98. Regarding kANNolo, in the 1.5× memory budget scenario, slowdowns are reduced, ranging from 2.5× to 5.2×. In the 2× memory budget, the slowdown is small for mid accuracy cuts (90-94) and then, as with Seismic, it increases, reaching a peak of 7.9× at 98% accuracy.

## 6  CONCLUDING REMARKS

We presented Seismic, a novel approximate algorithm that facilitates effective and efficient retrieval over sparse embedding vectors. We showed empirically its remarkable efficiency on a number of embeddings of publicly-available datasets. Seismic outperforms existing methods, including the winning algorithms at the BigANN Challenge in NeurIPS 2023 that use similar-sized (or larger)

Table 4. Index size and build time for winning configurations.

| Model | Configuration | Index size (GB) | Build time (hours) |
|---|---|---|---|
| kANNolo | M=32 | 104.4 | 10.1 |
| Seismic (with $\kappa$-nn graph) | $\lambda = 4e4, \kappa = 20$ | 98.4 | 16.6 |
| Seismic (without $\kappa$-nn graph) | $\lambda = 6e4$ | 98.5 | 0.5 |

indexes. Incredibly, one configuration of Seismic performs top-10 Anns over nearly 140 million vectors with almost exact accuracy within just 3 milliseconds.

The efficiency of Seismic is in large part due to a sketching algorithm we designed specifically for sparse vectors. As we explained the evolution of the sketch, dubbed Set $\alpha$-MSS, and presented theoretical results, the sketch can dramatically reduce the number of nonzero coordinates of a nonnegative sparse vector while preserving relative pairwise distances with arbitrarily high probability. That enabled us to reduce the size our novel inverted index, which in turn leads to fewer data points being evaluated.

# REFERENCES

[1] Ismail Sengor Altingovde, Engin Demir, Fazli Can, and Özgür Ulusoy. 2008. Incremental cluster-based retrieval using compressed cluster-skipping inverted files. *ACM Transactions on Information Systems* 26, 3, Article 15 (June 2008), 36 pages.

[2] 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:2010.00768 [cs.IR]

[3] Dan Braun, Jordan Taylor, Nicholas Goldowsky-Dill, and Lee Sharkey. 2024. Identifying Functionally Important Features with End-to-End Sparse Dictionary Learning. arXiv:2405.12241 [cs.LG]

[4] Andrei Z. Broder, David Carmel, Michael Herscovici, Aya Soffer, and Jason Zien. 2003. Efficient Query Evaluation Using a Two-Level Retrieval Process. In *Proceedings of the 12th International Conference on Information and Knowledge Management* (New Orleans, LA, USA). 426–434.

[5] Sebastian Bruch. 2024. *Foundations of Vector Retrieval.* Springer Nature Switzerland.

[6] Sebastian Bruch, Claudio Lucchese, and Franco Maria Nardini. 2023. Efficient and Effective Tree-based and Neural Learning to Rank. *Foundations and Trends® in Information Retrieval* 17, 1 (2023), 1–123.

[7] Sebastian Bruch, Franco Maria Nardini, Amir Ingber, and Edo Liberty. 2023. An Approximate Algorithm for Maximum Inner Product Search over Streaming Sparse Vectors. *ACM Transactions on Information Systems* 42, 2, Article 42 (November 2023), 43 pages.

[8] Sebastian Bruch, Franco Maria Nardini, Amir Ingber, and Edo Liberty. 2024. Bridging Dense and Sparse Maximum Inner Product Search. *ACM Transactions on Information Systems* 42, 6, Article 151 (Aug. 2024), 38 pages.

[9] 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* (Washington DC, USA). 152–162.

[10] Sebastian Bruch, Franco Maria Nardini, Cosimo Rulli, and Rossano Venturini. 2024. Pairing Clustered Inverted Indexes with $\kappa$-NN Graphs for Fast Approximate Retrieval over Learned Sparse Representations. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management* (Boise, ID, USA). 3642–3646.

[11] Sebastian Bruch, Franco Maria Nardini, Cosimo Rulli, Rossano Venturini, and Leonardo Venuta. 2025. Investigating the Scalability of Approximate Sparse Retrieval Algorithms to Massive Datasets. In *Advances in Information Retrieval*. 437–445.

[12] Flavio Chierichetti, Alessandro Panconesi, Prabhakar Raghavan, Mauro Sozio, Alessandro Tiberi, and Eli Upfal. 2007. Finding near Neighbors through Cluster Pruning. In *Proceedings of the Twenty-Sixth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems* (Beijing, China). 103–112.

[13] Enes Recep Cinar and Ismail Sengor Altingovde. 2023. Exploiting Cluster-Skipping Inverted Index for Semantic Place Retrieval. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Taipei, Taiwan). 1981–1985.

[14] Hoagy Cunningham, Aidan Ewart, Logan Riggs, Robert Huben, and Lee Sharkey. 2023. Sparse Autoencoders Find Highly Interpretable Features in Language Models. arXiv:2309.08600 [cs.LG]

[15] Zhuyun Dai and Jamie Callan. 2019. Context-Aware Sentence/Passage Term Importance Estimation For First Stage Retrieval. arXiv:1910.10687 [cs.IR]

[16] Zhuyun Dai and Jamie Callan. 2020. Context-Aware Document Term Weighting for Ad-Hoc Search. In *Proceedings of The Web Conference* (Taipei, Taiwan). 1897–1907.

[17] Zhuyun Dai and Jamie Callan. 2020. Context-Aware Term Weighting For First Stage Passage Retrieval. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval* (Virtual Event, China). 1533–1536.

[18] Majid Daliri, Juliana Freire, Christopher Musco, Aécio Santos, and Haoxiang Zhang. 2023. Sampling Methods for Inner Product Sketching. arXiv:2309.16157 [cs.DB]

[19] Leonardo Delfino, Domenico Erriquez, Silvio Martinico, Franco Maria Nardini, Cosimo Rulli, and Rossano Venturini. 2025. kANNolo: Sweet and Smooth Approximate k-Nearest Neighbors Search. In *Advances in Information Retrieval: 47th European Conference on Information Retrieval, ECIR 2025, Lucca, Italy, April 6–10, 2025, Proceedings, Part IV* (Lucca, Italy). 400–406.

[20] 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.

[21] Constantinos Dimopoulos, Sergey Nepomnyachiy, and Torsten Suel. 2013. Optimizing Top-k Document Retrieval Strategies for Block-Max Indexes. In *Proceedings of the Sixth ACM International Conference on Web Search and Data Mining* (Rome, Italy). 113–122.

[22] Shuai Ding and Torsten Suel. 2011. Faster Top-k Document Retrieval Using Block-Max Indexes. In *Proceedings of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Beijing, China). 993–1002.

[23] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. 2022. Toy Models of Superposition. arXiv:2209.10652 [cs.LG] https://arxiv.org/abs/2209.10652

[24] Thibault Formal, Carlos Lassance, Benjamin Piwowarski, and Stéphane Clinchant. 2021. SPLADE v2: Sparse Lexical and Expansion Model for Information Retrieval. arXiv:2109.10086 [cs.IR]

[25] Thibault Formal, Carlos Lassance, Benjamin Piwowarski, and Stéphane Clinchant. 2022. From Distillation to Hard Negative Sampling: Making Sparse Neural IR Models More Effective. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Madrid, Spain). 2353–2359.

[26] Thibault Formal, Carlos Lassance, Benjamin Piwowarski, and Stéphane Clinchant. 2023. Towards Effective and Efficient Sparse Neural Information Retrieval. *ACM Transactions on Information Systems* (December 2023).

[27] 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* (Virtual Event, Canada). 2288–2292.

[28] Tiezheng Ge, Kaiming He, Qifa Ke, and Jian Sun. 2014. Optimized Product Quantization. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 36, 4 (2014), 744–755.

[29] Ruiqi Guo, Philip Sun, Erik Lindgren, Quan Geng, David Simcha, Felix Chern, and Sanjiv Kumar. 2020. Accelerating large-scale inference with anisotropic vector quantization. In *Proceedings of the 37th International Conference on Machine Learning*. Article 364, 10 pages.

[30] Zhengfu He, Wentao Shu, Xuyang Ge, Lingjie Chen, Junxuan Wang, Yunhua Zhou, Frances Liu, Qipeng Guo, Xuanjing Huang, Zuxuan Wu, Yu-Gang Jiang, and Xipeng Qiu. 2024. Llama Scope: Extracting Millions of Features from Llama-3.1-8B with Sparse Autoencoders. arXiv:2410.20526 [cs.LG]

[31] Herve Jégou, Matthijs Douze, and Cordelia Schmid. 2011. Product Quantization for Nearest Neighbor Search. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 33, 1 (2011), 117–128.

[32] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2021. Billion-Scale Similarity Search with GPUs. *IEEE Transactions on Big Data* 7, 3 (2021), 535–547.

[33] Yannis Kalantidis and Yannis Avrithis. 2014. Locally Optimized Product Quantization for Approximate Nearest Neighbor Search. In *2014 IEEE Conference on Computer Vision and Pattern Recognition*. 2329–2336.

[34] Hrishikesh Kulkarni, Sean MacAvaney, Nazli Goharian, and Ophir Frieder. 2023. Lexically-Accelerated Dense Retrieval. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Taipei, Taiwan). 152–162.

[35] Carlos Lassance and Stéphane Clinchant. 2022. An Efficiency Study for SPLADE Models. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Madrid, Spain). 2220–2226.

[36] Carlos Lassance, Hervé Déjean, Thibault Formal, and Stéphane Clinchant. 2024. SPLADE-v3: New baselines for SPLADE. arXiv:2403.06789 [cs.IR]

[37] Carlos Lassance, Simon Lupart, Hervé Déjean, Stéphane Clinchant, and Nicola Tonellotto. 2023. A Static Pruning Study on Sparse Neural Retrievers. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Taipei, Taiwan). 1771–1775.

[38] Tom Lieberum, Senthooran Rajamanoharan, Arthur Conmy, Lewis Smith, Nicolas Sonnerat, Vikrant Varma, Janos Kramar, Anca Dragan, Rohin Shah, and Neel Nanda. 2024. Gemma Scope: Open Sparse Autoencoders Everywhere All At Once on Gemma 2. In *Proceedings of the 7th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*. 278–300.

[39] Jimmy Lin and Xueguang Ma. 2021. A Few Brief Notes on DeepImpact, COIL, and a Conceptual Framework for Information Retrieval Techniques. arXiv:2106.14807 [cs.IR]

[40] Jimmy Lin, Rodrigo Frassetto Nogueira, and Andrew Yates. 2021. *Pretrained Transformers for Text Ranking: BERT and Beyond.* Morgan & Claypool Publishers.

[41] Jimmy Lin and Andrew Trotman. 2015. Anytime Ranking for Impact-Ordered Indexes. In *Proceedings of the 2015 International Conference on The Theory of Information Retrieval* (Northampton, Massachusetts, USA). 301–304.

[42] Xueguang Ma, Ronak Pradeep, Rodrigo Nogueira, and Jimmy Lin. 2022. Document expansion baselines and learned sparse lexical representations for ms marco v1 and v2. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 3187–3197.

[43] Sean MacAvaney, Franco Maria Nardini, Raffaele Perego, Nicola Tonellotto, Nazli Goharian, and Ophir Frieder. 2020. Expansion via Prediction of Importance with Contextualization. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval* (Virtual Event, China). 1573–1576.

[44] Sean MacAvaney, Nicola Tonellotto, and Craig Macdonald. 2022. Adaptive Re-Ranking with a Corpus Graph. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management* (Atlanta, GA, USA). 1491–1500.

[45] Joel Mackenzie, Antonio Mallia, Alistair Moffat, and Matthias Petri. 2022. Accelerating Learned Sparse Indexes Via Term Impact Decomposition. In *Findings of the Association for Computational Linguistics: EMNLP 2022*. 2830–2842.

[46] J. Mackenzie, M. Petri, and L. Gallagher. 2022. IOQP: A simple Impact-Ordered Query Processor written in Rust. In *Proc. DESIRES*. 22–34.

[47] J. Mackenzie, M. Petri, and A. Moffat. 2021. Faster Index Reordering with Bipartite Graph Partitioning. In *Proc. SIGIR*. 1910–1914.

[48] Joel Mackenzie, Andrew Trotman, and Jimmy Lin. 2021. Wacky Weights in Learned Sparse Representations and the Revenge of Score-at-a-Time Query Evaluation. arXiv:2110.11540 [cs.IR]

[49] Joel Mackenzie, Andrew Trotman, and Jimmy Lin. 2023. Efficient Document-at-a-Time and Score-at-a-Time Query Evaluation for Learned Sparse Representations. *ACM Transactions on Information Systems* 41, 4 (2023).

[50] Yu A. Malkov and D. A. Yashunin. 2020. Efficient and Robust Approximate Nearest Neighbor Search Using Hierarchical Navigable Small World Graphs. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 42, 4 (4 2020), 824–836.

[51] Antonio Mallia, Joel Mackenzie, Torsten Suel, and Nicola Tonellotto. 2022. Faster Learned Sparse Retrieval with Guided Traversal. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Madrid, Spain). 1901–1905.

[52] Antonio Mallia, Giuseppe Ottaviano, Elia Porciani, Nicola Tonellotto, and Rossano Venturini. 2017. Faster BlockMax WAND with Variable-Sized Blocks. In *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Shinjuku, Tokyo, Japan). 625–634.

[53] Antonio Mallia and Elia Porciani. 2019. Faster BlockMax WAND with Longer Skipping. In *Advances in Information Retrieval*. 771–778.

[54] Antonio Mallia, Michal Siedlaczek, Joel Mackenzie, and Torsten Suel. 2019. PISA: Performant Indexes and Search for Academia. In *Proceedings of the Open-Source IR Replicability Challenge co-located with 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, OSIRRC@SIGIR 2019, Paris, France.* 50–56.

[55] Antonio Mallia, Torsten Suel, and Nicola Tonellotto. 2024. Faster Learned Sparse Retrieval with Block-Max Pruning. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Washington DC, USA). 2411–2415.

[56] Antonio Mallia, Torsten Suel, and Nicola Tonellotto. 2024. Faster Learned Sparse Retrieval with Block-Max Pruning. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Washington DC, USA). 2411–2415.

[57] Samuel Marks, Can Rager, Eric J Michaud, Yonatan Belinkov, David Bau, and Aaron Mueller. 2025. Sparse Feature Circuits: Discovering and Editing Interpretable Causal Graphs in Language Models. In *The 13th International Conference on Learning Representations*.

[58]  Stanislav Morozov and Artem Babenko. 2018. Non-metric Similarity Graphs for Maximum Inner Product Search. In *Advances in Neural Information Processing Systems*.

[59]  Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. MS MARCO: A Human Generated MAchine Reading COmprehension Dataset.

[60]  Rodrigo Nogueira, Jimmy Lin, and AI Epistemic. 2019. From doc2query to docTTTTTquery. *Online preprint* 6 (2019), 2.

[61]  Charles O'Neill and Thang Bui. 2024. Sparse Autoencoders Enable Scalable and Reliable Circuit Identification in Language Models. arXiv:2405.12522 [cs.CL]

[62]  Giuseppe Ottaviano and Rossano Venturini. 2014. Partitioned Elias-Fano Indexes. In *Proceedings of the 37th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. 273–282.

[63]  Stephen E. Robertson, Steve Walker, Susan Jones, Micheline Hancock-Beaulieu, and Mike Gatford. 1994. Okapi at TREC-3.. In *TREC (NIST Special Publication, Vol. 500-225)*, Donna K. Harman (Ed.). National Institute of Standards and Technology (NIST), 109–126.

[64]  Gerard Salton and Christopher Buckley. 1988. Term-weighting approaches in automatic text retrieval. *Information Processing and Management* 24, 5 (1988), 513–523.

[65]  Harsha Vardhan Simhadri, Martin Aumüller, Amir Ingber, Matthijs Douze, George Williams, Magdalen Dobson Manohar, Dmitry Baranchuk, Edo Liberty, Frank Liu, Ben Landrum, Mazin Karjikar, Laxman Dhulipala, Meng Chen, Yue Chen, Rui Ma, Kai Zhang, Yuzheng Cai, Jiayang Shi, Yizhuo Chen, Weiguo Zheng, Zihao Wan, Jie Yin, and Ben Huang. 2024. Results of the Big ANN: NeurIPS'23 competition. arXiv:2409.17424 [cs.IR]

[66]  Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. 2021. BEIR: A Heterogeneous Benchmark for Zero-shot Evaluation of Information Retrieval Models. In *35th Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.

[67]  Nicola Tonellotto, Craig Macdonald, and Iadh Ounis. 2018. Efficient Query Processing for Scalable Web Search. *Foundations and Trends in Information Retrieval* 12, 4–5 (December 2018), 319–500.

[68]  Howard Turtle and James Flood. 1995. Query Evaluation: Strategies and Optimizations. *Information Processing and Management* 31, 6 (November 1995), 831–850.

[69]  David P. Woodruff. 2014. Sketching as a Tool for Numerical Linear Algebra. *Foundations and Trends in Theoretical Computer Science* 10, 1–2 (Oct 2014), 1–157.

[70]  Xiang Wu, Ruiqi Guo, Ananda Theertha Suresh, Sanjiv Kumar, Daniel N Holtmann-Rice, David Simcha, and Felix Yu. 2017. Multiscale Quantization for Fast Similarity Search. In *Advances in Neural Information Processing Systems*, Vol. 30.

[71]  Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Łukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2016. Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation.

[72]  Hamed Zamani, Mostafa Dehghani, W. Bruce Croft, Erik Learned-Miller, and Jaap Kamps. 2018. From Neural Re-Ranking to Neural Ranking: Learning a Sparse Representation for Inverted Indexing. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management* (Torino, Italy). 497–506.

[73]  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.