

# A SOFT AND FAST PATTERN MATCHER FOR BILLION-SCALE CORPUS SEARCHES

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Researchers and practitioners in natural language processing and computational linguistics frequently observe and analyze the real language usage in large-scale corpora. For that purpose, they often employ off-the-shelf pattern-matching tools, such as `grep`, and keyword-in-context concordancers, which is widely used in corpus linguistics for gathering examples. Nonetheless, these existing techniques rely on surface-level string matching, and thus they suffer from the major limitation of not being able to handle orthographic variations and paraphrasing—notable and common phenomena in any natural language. In addition, existing continuous approaches such as dense vector search tend to be overly coarse, often retrieving texts that are unrelated but share similar topics. Given these challenges, we propose a novel algorithm that achieves *soft* (or semantic) yet efficient pattern matching by relaxing a surface-level matching with word embeddings. Our algorithm is highly scalable with respect to the size of the corpus text utilizing inverted indexes. We have prepared an efficient implementation, and we provide an accessible web tool. Our experiments demonstrate that the proposed method (i) can execute searches on billion-scale corpora in less than a second, which is comparable in speed to surface-level string matching and dense vector search; (ii) can extract harmful instances that semantically match queries from a large set of English and Japanese Wikipedia articles; and (iii) can be effectively applied to corpus-linguistic analyses of Latin, a language with highly diverse inflections.

## 1 INTRODUCTION

Recent advances in natural language processing (NLP) and corpus linguistics are largely driven by the availability of massive text corpora (Gao et al., 2020; Biderman et al., 2023b; Raffel et al., 2019; Van Der Zwaan et al., 2017; McGillivray et al., 2020). The field of NLP, driven by the grand goal of building intelligent chatbots and machine translation systems, has achieved remarkable breakthroughs over the past decade, largely thanks to self-supervised representation learning from vast corpora (Mikolov et al., 2013; Pennington et al., 2014; Devlin et al., 2019; Radford et al., 2019). Notable causal language models (LMs), capable of passing high-stakes exams (Katz et al., 2024; OpenAI, 2024; Jung et al., 2023), serving as general-purpose problem solvers (Brown et al., 2020), and engaging in realistic conversations with beloved characters (De Freitas, 2023; Chen et al., 2024), owe their foundational abilities to next-token prediction learned from large-scale data (Brown et al., 2020; Dubey et al., 2024; Riviere et al., 2024). In corpus linguistics (McEnery & Hardie, 2011), computational linguistics (Jurafsky & Martin, 2024), and digital humanities (Jensen, 2014)—disciplines aiming to uncover the scientific and computational principles of human language—such vast linguistic resources, consisting of the language use itself, have become more indispensable than ever in this era of large-scale corpora (Van Der Zwaan et al., 2017; McGillivray et al., 2020).

Given this context, the demand for efficient pattern matchers that enable rapid searches across massive corpora is higher than ever (Liu et al., 2024; Smadja, 1993). For example, when harmful information, misinformation, or memorized privacy information is generated by a large LM, it is necessary to identify the corresponding training instances where this information originated (Guo et al., 2022; Wang et al., 2024c; Ippolito et al., 2023; Biderman et al., 2023a). Another example is when researchers wish to determine which of two linguistic phenomena of interest is more frequent; conducting exhaustive searches across the largest available corpora is essential (Biber, 2015) for this purpose. Notably, many linguistic phenomena—word frequency as a simple example—follow

Table 1: *Hard* pattern matching (e.g., `grep`) and *soft* pattern matching (our method). Hard matching is based on *surface-level* comparison and does not match if, e.g., a synonym is used; Soft matching is based on *semantic* comparison and robust to such variations thanks to word embeddings.

Pattern	Theorem 1	John was born in
<i>Hard</i> matching (e.g., <code>grep</code> )	... thanks to <b>Theorem 1</b> ...	... <b>John was born in</b> 1345 ...
<i>Soft</i> matching (ours)	... thanks to <b>Theorem 1</b> ... ... <b>Theorem 3</b> holds because ... ... By <b>Lemma 5</b> , we may assume ... ... <b>Equation 1</b> describes ...	... <b>John was born in</b> 1345 ... ... <b>Edward had died in</b> May 1910 ... ... <b>Robert was born in</b> England ... ... the Emperor <b>Henry, died in</b> 1125 ...

a power-law distribution (Zipf, 1951; Heap, 1980; Kobayashi & Tanaka-Ishii, 2018). Consequently, only with vast corpora can we observe and analyze the various rare events residing in the long tail, which constitute the majority of linguistic phenomena.

One of the challenges when working with large corpora is that pattern-matching tools based on string matching (Hakak et al., 2019), such as `grep` (Bambenek & Klus, 2009) and `ripgrep` (Gallant, 2024), or linguist-oriented tools often referred to as KWIC (Anthony, 2013; Culy & Lyding, 2010; Schweinberger, 2024), primarily rely on surface-level exact matching as their core strategy. Because natural languages are characterized by their flexibility and richness in how humans can express similar concepts in different ways (McKeown, 1979; Witteveen & Andrews, 2019; Ganitkevitch et al., 2013), strictly exact string matching may not meet users’ demands. For example, on top of the query word in standard spelling, it is often desirable to catch non-standard spellings as well, such as *how r u* instead of *how are you*, which are widespread particularly on the web and in texting (Schulz, 2018). Additionally, it is also desirable to catch different inflected word forms such as *sing*, *sang*, *sung*, *sings*, and *singing*, which differ only in their morphological features and share the same lemma (base form) (Don, 2014; Embick, 2015). Rule-based exact matching is particularly hard when the target language is morphologically complex and exhibits irregular inflectional patterns.

To resolve the mismatch between the symbolic nature of existing pattern matchers and the diverse orthographic and morphological variations inherent in natural language, we have developed a *soft* (or semantic) yet efficient pattern matcher (Section 3). The core strategy is based on the simple idea of *softening* the matching process from binary  $\{0, 1\}$  values to continuous values, using word embeddings. By adopting inverted indexes and several other techniques, our tool can enumerate all softly matching instances in billion-scale corpora in less than a second. Table 1 shows specific examples of its operation. Our proposed method is capable not only of enumerating exact matches but also of flexibly listing semantically similar instances, even when their surface forms differ. For example, given the query “*Theorem 1*”, the method can retrieve instances such as “*Lemma 5*”. This characteristic substantially enhances key tasks in both NLP and corpus linguistics. For instance, it improves the filtering of harmful texts, and facilitates more efficient example retrieval, particularly for languages with complex morphological features (Section 4).

The contributions of our study are summarized as follows.

- We developed a *soft* (or semantic) pattern-matching algorithm—a relaxation of the exact string matching—using word embeddings (Section 3). By leveraging inverted indexing, we achieved high scalability (Sections 2.1, 2.3 and 3.3). For instance, the running time for searching over an English corpus with 3.4B words was less than 0.1 seconds without GPUs.
- We developed an easy-to-use web demo to facilitate interaction with the soft matching algorithm.<sup>1</sup> This demo is particularly beneficial for NLP researchers and developers who want to analyze LMs in a data-driven manner, as well as for language learners and humanities researchers who are conducting language analysis on large corpora from the perspectives of corpus linguistics and digital humanities (Section 4.1)

<sup>1</sup> Anonymized URL: <http://54.238.189.64/>. Also see Footnote 4.

- For quantitative evaluation, we verified that our method achieves complete enumeration in less than a second on billion-scale corpora, which are commonly used in training large LMs (Section 4.2). This performance is comparable in speed to dense vector search; moreover, our method enables the retrieval of examples closely aligned with specific queries, rather than just broad topical similarities.
- We conducted qualitative evaluations in specific scenarios to verify the usefulness of the proposed method in both the fields of NLP and corpus linguistics. In experiments assuming the training of LMs on large-scale corpora, we provided search examples for identifying and removing harmful instances contained within the corpus (Section 4.2). In experiments assuming the analysis of classical Western languages by linguists, we chose Latin—a language with highly complex inflections—as an example. We demonstrated that a corpus search with our tool can flexibly extract morphologically and semantically similar usage examples from the corpus (Section 4.3).

## 2 RELATED WORK

The procedures for finding sentences (or lines, documents) that match a given pattern (query) are prevalent across nearly all areas of computer science and data science, making it difficult to enumerate all related research. In this section, we review particularly relevant work in NLP and computational linguistics, which are the primary focus of this study, as well as in the closely related areas of string matching. Beyond the fields discussed here, we believe many other domains, such as information retrieval, time series analysis, and semantic web, also have connections to this research.

### 2.1 NATURAL LANGUAGE PROCESSING AND CORPUS SEARCH

The significant progress in NLP over the past decade is largely attributed to deep-learning-based self-supervised representation learning (Mikolov et al., 2013; Pennington et al., 2014; Bojanowski et al., 2017; Devlin et al., 2019; Radford et al., 2019; Brown et al., 2020; Dubey et al., 2024). **The vast raw corpora** without label annotations serve as the source of such models’ capabilities (Gao et al., 2020; Biderman et al., 2023b; Raffel et al., 2019), and, therefore, these corpora are continually referenced and analyzed for model improvement and evaluation (Biderman et al., 2023b). For instance, as LMs can memorize and elicit facts written in training corpora, they are reported to generate text containing privacy-related information (Huang et al., 2022; Li et al., 2023; Lukas et al., 2023) or content that could be used for terrorism or violence (Gehman et al., 2020; Schick et al., 2021; Kumar et al., 2023). In this context, Ippolito et al. (2023) work on filtering out verbatim memorization, requiring searches across large-scale training corpora.<sup>2</sup> ***N*-gram substring pattern matchers**, as representative tools for corpus search in NLP, have been particularly well-developed even before the advent of deep learning (Sekine, 2008; Sekine & Dalwani, 2010). In terms of data structures, inverted indexes (Sekine & Dalwani, 2010; Rogozinski & Kuc, 2016), and suffix arrays including their variants (Sekine & Dalwani, 2010; Yamamoto & Church, 2001; Liu et al., 2024; Burrows et al., 1994; Ferragina & Manzini, 2005; Langmead et al., 2009) are typically employed. Our algorithm is based on inverted indexes due to its algorithmic requirements, which we discuss in Section 3. Here, regardless of which data structure is used, it is important to note that existing *n*-gram matching techniques assume exact surface-level matching, making it extremely challenging to search and enumerate all relevant sentences while handling the complex characteristics of natural language, such as paraphrasing, orthographic variation, and inflection. **Dense vector search** (Khattab & Zaharia, 2020; Karpukhin et al., 2020; Izacard et al., 2022; Wang et al., 2024a) has recently gained widespread popularity as the foundational technology for retrieval-augmented generation (RAG) (Lewis et al., 2020; Khandelwal et al., 2020; Guu et al., 2020; Izacard et al., 2024). Dense vector search is highly compatible with approximate nearest neighbor search (Malkov & Yashunin, 2020; Jégou et al., 2011), and it offers a significant advantage in reducing the issue of hallucination (Ayala & Bechard, 2024; Gao et al., 2023) in scenarios requiring factual knowledge. However, dense vector search is a relatively “coarse” method, primarily used to retrieve documents that are topically similar (e.g., everything related to the U.S. presidential election), making it less suitable for the cases that require more specific *n*-gram level queries.

<sup>2</sup> Measuring influence of training corpus is also an important theme (Koh & Liang, 2017; Pruthi et al., 2020; Chen et al., 2021; Isonuma & Titov, 2024), but their application to vast corpora remains highly challenging.

## 2.2 CORPUS LINGUISTICS AND CORPUS SEARCH

Corpus linguistics is a subfield of linguistics that utilizes corpora to study language based on examples of ‘real-life’ language use (McEnery & Wilson, 2001). Digital humanities, a closely related field, is an interdisciplinary domain that aims to advance humanities through the application of data and information science. In digital humanities as well, quantitative analysis of texts using corpora is also actively pursued (Jensen, 2014). **Corpora** consist of collections of texts or spoken language data (Van Der Zwaan et al., 2017; McGillivray et al., 2020), offering real-world examples of how language is used in various contexts. In most settings, an ideal corpus is large in size to ensure the statistical reliability of certain linguistic phenomena and to cover a broader range of language use, including rare occurrences such as low-frequency words (Ha et al., 2009; Coole et al., 2020). **Search tools** (Hockey & Martin, 1987; TEI Consortium, 2023) are an indispensable element of corpus linguistics because they enable the identification of linguistic patterns, such as word frequencies, collocations, and syntactic structures. Traditional search methods commonly used in corpus linguistics include exact matching and its extension with regular expressions (Jurafsky & Martin, 2024). However, given the nature of natural languages—with their morphological complexity and unlimited paraphrases and synonyms with varying surface forms—rule-based search methods face extreme difficulties in exhaustively enumerating instances that closely match a given query.

## 2.3 STRING MATCHING

From an algorithmic point of view, our method is closely related to *string matching* algorithms (Hakak et al., 2019). In what follows, we briefly give an algorithmic comparison with some of them and elucidate their relationship to ours. **Offline string matching** algorithms process the corpus text to build an index that accelerates future searches, a technique widely applied in search engines and information retrieval systems such as Elasticsearch (Rogozinski & Kuc, 2016) and Apache Lucene (Grand et al., 2020). Inverted indexing is a well-known approach in this domain (Zobel & Moffat, 2006), and our method serves as a *soft* generalization of a string matching algorithm based on inverted indexing. **Online string matching** algorithms do not rely on preprocessing and handle the corpus text on the fly. Efficient algorithms include, e.g., the Knuth-Morris-Pratt, Karp-Rabin, and Boyer-Moore algorithms (Knuth et al., 1977; Karp & Rabin, 1987; Boyer & Moore, 1977). While widely used in tools such as `grep`, online string matching is also a basis of runtime verification (Bartocci et al., 2018), where system’s execution data is monitored. Although it is, in theory, possible to relax these algorithms with word embeddings, the number of *soft* word comparisons in such a relaxation would be linear in the size of the corpus text. In contrast, our algorithm requires only constant *soft* word comparisons in the size of the corpus text thanks to indexing. See Section 3.3 for the complexity analysis. **Approximate string matching** (or fuzzy matching) allows flexible matching, typically using edit distance to accommodate noisy or incomplete data (Navarro, 2001). Tools like `agrep` (Wu & Manber, 1992b;a) perform online matching that tolerates mismatches. Indexing-based algorithms are also proposed for approximate string matching (Boytsov, 2011). As a relaxation of exact string matching, approximate string matching is orthogonal to ours: approximate string matching focuses on relaxing surface-level comparison (e.g., for typos), while our approach focuses on semantic-level similarity based on word embeddings (e.g., for synonyms).

# 3 OUR ALGORITHM FOR SOFT PATTERN MATCHING

## 3.1 OVERVIEW

**Key aspect of design: Hard vs. Soft.** One key aspect we considered in designing the algorithm for soft pattern matching is *hard computation vs. soft computation*. *Hard* (or exact) pattern matching (such as `grep`) can process large texts blazingly fast. This is because hard matching requires only “*hard* computation”, e.g., bitwise comparison, which is highly efficient. *Soft* pattern matching, on the other hand, involves “*soft* computation”, e.g., cosine similarity of two vectors. Since such *soft* computation involves many floating-point arithmetic operations over high-dimensional vectors, it is typically much slower than *hard* computation. For instance, just taking the cosine similarity of a pair of word embeddings can be as expensive as thousands of simple bitwise operations.

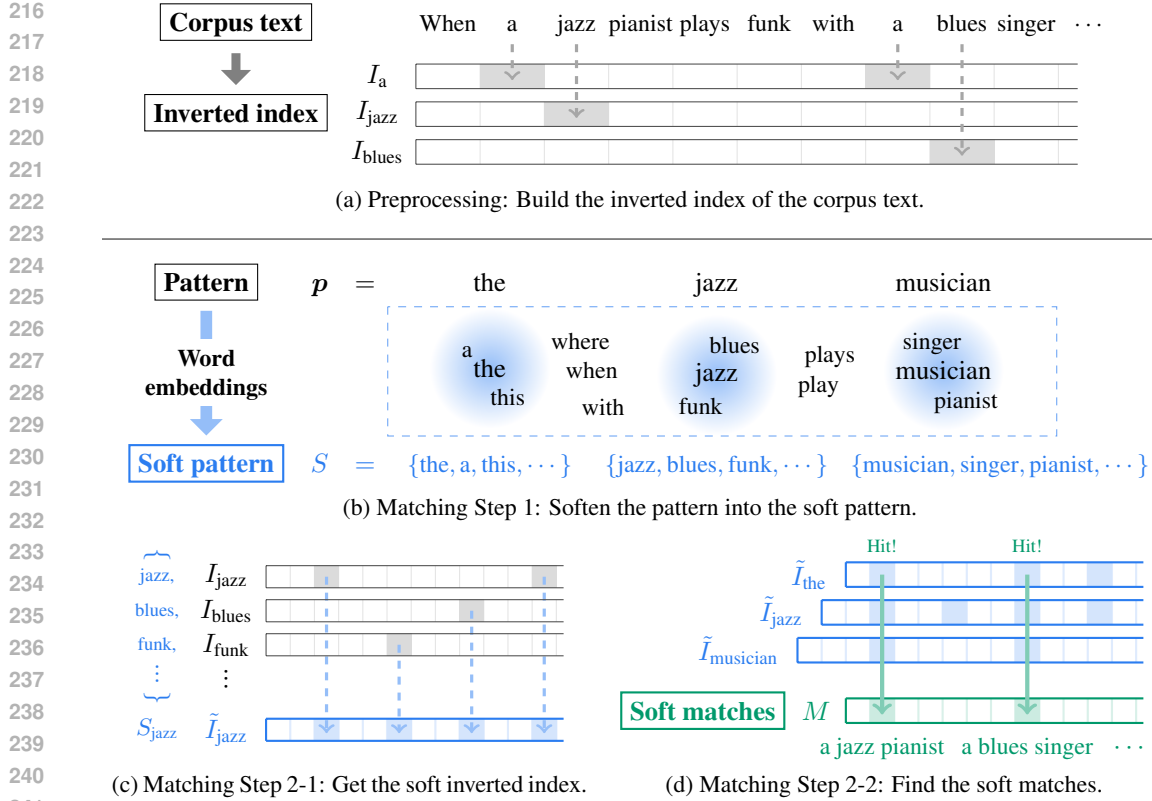


Figure 1: Illustration of our algorithm for soft pattern matching.

**Overview of our algorithm.** With this in mind, we designed the algorithm for soft pattern matching, which is precise and efficient. Our algorithm is illustrated in Figure 1. Our key idea is to run soft computation *only over the vocabulary*, not over the whole text as a naive algorithm would do. More specifically, our algorithm works in the following two steps: For each query, (Step 1) it first performs *soft* computation over the *vocabulary*, *softening* the pattern into the *soft pattern*  $S_1 S_2 \cdots S_n$  (Fig. 1b); (Step 2) then it performs *hard* computation that finds the *exact* positions in the corpus text that match the soft pattern (Figs. 1c and 1d). Here, the soft pattern  $S_1 S_2 \cdots S_n$  is the key data in play. It consists of the set of vocabulary words that *softly* match each word of the pattern in terms of word embeddings. The amount of soft computation (Step 1) is small, given that the vocabulary size is typically much smaller than the size of the corpus text. The amount of hard computation (Step 2) is also quite small, because it only handles filtered positions, not the whole corpus text.

### 3.2 DETAILS

**Problem formulation.** First, we formulate the *soft* pattern-matching problem by simply relaxing the classical *hard* (or exact) pattern-matching problem (Hakak et al., 2019) as follows:

**The soft pattern-matching problem.**

INPUT: The corpus text  $t = t_1 t_2 \cdots t_N \in \mathcal{V}^*$ , the pattern  $p = p_1 p_2 \cdots p_n \in \mathcal{V}^*$ , and the threshold  $\alpha \in (0, 1]$

OUTPUT: The set of soft match positions  $M$  with respect to the soft equivalence  $\approx_\alpha$ , i.e., the set  $M = \{i \in \{1, \dots, N\} \mid \forall k \in \{1, \dots, n\}. p_k \approx_\alpha t_{i+k-1}\}$

We define here the *soft equivalence*  $v \approx_\alpha v'$  between words as  $v \approx_\alpha v' \triangleq \cos(E(v), E(v')) \geq \alpha$ , to capture the similarity in terms of *word embeddings*  $E$ . Here, we let  $\mathcal{V} = \{v_1, v_2, \dots, v_L\}$  be the vocabulary and write  $E(v) \in \mathbb{R}^D$  for the word embedding of a word  $v \in \mathcal{V}$ . Also, we write  $\cos(e, e')$  for the cosine similarity  $\cos(e, e') \triangleq \frac{e \cdot e'}{\|e\| \|e'\|}$  of word embeddings  $e, e' \in \mathbb{R}^D$ .

**Algorithm 1:** Our soft pattern-matching algorithm.

---

```

270
271
272 Given : The corpus text  $\mathbf{t} = t_1 t_2 \dots t_N$  and its inverted index  $I_v \subseteq \{1, \dots, N\}$  over each vocabulary
273 word  $v \in \mathcal{V}$  such that  $i \in I_v$  if and only if  $v = t_i$ .
274 Input : The phrase pattern  $\mathbf{p} = p_1 p_2 \dots p_n$  and the threshold  $\alpha \in (0, 1]$ .
275 Output : The set of soft match positions  $M \subseteq \{1, \dots, N\}$ , such that  $i \in M$  holds if and only if
276  $p_k \approx_\alpha t_{i+k-1}$  holds over any  $k = 1, \dots, n$ .
277 // Step 1: Soften the pattern  $p_1 p_2 \dots p_n$  into the soft pattern  $S_1 S_2 \dots S_n$ 
278 1 for  $k \leftarrow 1 \dots n$  do
279   2 |  $S_k \leftarrow \emptyset$ 
280   3 | for  $v \in \mathcal{V}$  do
281     4 | | if  $v \approx_\alpha p_k$  then  $S_k \leftarrow S_k \cup \{v\}$ 
282 // Step 2-1: Get the soft inverted index  $\tilde{I}_k$  by relaxing  $I$  with  $S_k$ 
283 5 for  $k \leftarrow 1 \dots n$  do
284   6 |  $\tilde{I}_k \leftarrow \emptyset$ 
285   7 | for  $v \in S_k$  do  $\tilde{I}_k \leftarrow \tilde{I}_k \cup I_v$ 
286 // Step 2-2: Get the complete matching  $M$  by aggregating  $\tilde{I}_k$ 
287 8  $M \leftarrow \tilde{I}_1$ 
288 9 for  $k \leftarrow 2 \dots n$  do  $M \leftarrow M \cap \{i - (k - 1) \mid i \in \tilde{I}_k\}$ 
289 10 return  $M$ 

```

---

**Preprocessing.** Before processing queries, our algorithm preprocesses the whole corpus text to compute the *inverted index*  $I$  (Fig. 1a), following a standard technique in search engine indexing (Zobel & Moffat, 2006). The inverted index  $I$  maps each vocabulary word  $v \in \mathcal{V}$  to the set of positions  $I_v \subseteq \{1, \dots, N\}$  that represents the occurrences of the word  $v$  in the corpus text  $\mathbf{t}$ , that is, the set  $I_v = \{i \in \{1, \dots, N\} \mid v = t_i\}$ .

**Our algorithm.** Algorithm 1 and Figures 1b to 1d outline our algorithm for soft pattern matching. Our algorithm works in the following two steps:

1. For each query, it performs *soft* computation over the *vocabulary*, *softening* the pattern  $p_1 p_2 \dots p_n$  into the *soft pattern*  $S_1 S_2 \dots S_n$  (lines 1 to 4 in Algorithm 1, Figure 1b). Here,  $S_k$  is the set of vocabulary words that *softly* matches each word  $p_k$  of the pattern in terms of the *soft equivalence*  $\approx_\alpha$ , i.e.,  $S_k = \{v \in \mathcal{V} \mid v \approx_\alpha p_k\}$ .
2. Then, it performs *hard* computation that computes the *exact* set of positions  $M$  in the corpus text that matches the soft pattern  $S_1 S_2 \dots S_n$  by the following two sub-steps:
  - 2-1. Compute the *soft inverted index*  $\tilde{I}$ , which maps each pattern word  $p_k$  to the union  $\tilde{I}_k = \bigcup_{v \in S_k} I_v$  of the exact inverted index  $I$  over  $S_k$  (lines 5 to 7, Fig. 1c). The set  $\tilde{I}_k$  agrees with the set of positions that softly match the pattern word  $p_k$ , i.e.,  $\tilde{I}_k = \{i \in \{1, \dots, N\} \mid t_i \approx_\alpha p_k\}$ .
  - 2-2. Output the set of *soft matches*  $M$  by computing the shifting intersection  $M = \bigcap_{k=1}^n \{i - (k - 1) \mid i \in \tilde{I}_k\}$  over the soft inverted index  $\tilde{I}$  (lines 8 to 9, Fig. 1d). The resulting set precisely agrees with the expected set, i.e.,  $M = \{i \in \{1, \dots, N\} \mid \forall k \in \{1, \dots, n\}. p_k \approx_\alpha t_{i+k-1}\}$ .

We perform soft pattern matching using the index  $I$ . First, for each word  $p_k$  in  $\mathbf{p}$ , we construct the set  $\tilde{I}_k$  indicating the positions of the words in  $\mathbf{t}$  matching  $p_k$  (lines 1 to 4): we compare each word  $v \in \mathcal{V}$  in the vocabulary with  $p_k$  (line 4); we add  $I_v$  to  $\tilde{I}_k$  if we have  $v \approx_\alpha p_k$  (line 7). Then, we construct the result  $M$  by aggregating each  $\tilde{I}_k$  considering the position  $k$  of  $p_k$  in  $\mathbf{p}$ .

**Inverted index vs. Suffix array.** Our algorithm uses an inverted index rather than a suffix array (Ferragina & Manzini, 2005), which is crucial. A suffix array manages *exact sequences* of characters (or tokens) in the corpus text, for which *softening* patterns cannot be performed efficiently. In contrast, entries of an inverted index  $I_v$  just have the occurrences of *each* word type  $v$  and can be relaxed for soft matching simply by merging the sets, as illustrated in Fig. 1c.

### 3.3 FORMAL ANALYSIS

How efficient is our algorithm? To answer this with theoretical guarantees, we analyze the time and space complexity of the algorithm. Overall, the high-level observation is that the corpus text size  $N$  affects the time and space required for preprocessing and soft matching only *linearly*.

**Preprocessing.** The time complexity for constructing the inverted index  $I$  is  $\mathcal{O}(L \times N)$ , consisting of  $L \times N$  exact word comparisons, since each vocabulary word  $v \in \mathcal{V}$  is compared with each word  $t_i$  in the corpus text  $t$ . Notably, the constant factor here is typically very low, as only bitwise comparison is required. The total space required to store the inverted index  $I$  is only  $\mathcal{O}(N + L)$ , by simply storing the list of positions (of space  $\mathcal{O}(|I_v|)$ ) for each vocabulary word  $v \in \mathcal{V}$ . This is because each position  $i \in \{1, 2, \dots, N\}$  in the text  $t$  occurs exactly once in  $I$  and hence  $N = \sum_{v \in \mathcal{V}} |I_v|$ . The space complexity is also  $\mathcal{O}(N + L)$ , since no extra space is required.

**Soft matching.** Step 1 for softening the pattern takes  $\mathcal{O}(n \times L)$  in time, where the dominant part is  $n \times L$  soft word comparisons (taking the cosine similarity of word embeddings) between each pattern word  $p_k$  and each vocabulary word  $v \in \mathcal{V}$  (line 4). The space complexity is  $\mathcal{O}(\sum_{k=1}^n |S_k|)$ . Notably, the time and the space for Step 1 are independent of the corpus text size  $N$ . Step 2 for finding the set of soft matches takes  $\mathcal{O}(K)$  in time and space, where  $K$  is the total size of the soft inverted index  $K \triangleq \sum_{k=1}^n |\tilde{I}_k|$  (or roughly the total number of ‘candidate’ positions for matches). Remarkably, this is only linear to the corpus text size  $N$ .

## 4 EMPIRICAL EVALUATION

We address the following research questions in the experiments:

- RQ1:** Does `SoftMatcha` scale to a billion-scale corpus, which is the typical size of training corpus for large causal LMs? (Section 4.2)
- RQ2:** Does `SoftMatcha` perform as expected in typical scenarios in NLP and corpus linguistics? (Sections 4.2 and 4.3)

### 4.1 IMPLEMENTATION

We implemented the algorithm in Section 3 as a tool named `SoftMatcha`.<sup>3</sup> Our implementation adopts the following designs for efficiency. Firstly, we design our inverted index using a sparse matrix in a compressed sparse row (CSR) format to reduce memory consumption and enable efficient access to the index for each word, i.e.,  $I_v$ . The index  $I$  is represented by a one-dimensional array, with the positions of each word contiguously allocated in memory. Secondly, we compile some time-consuming operations to native code using NUMBA (Lam et al., 2015). Specifically, line 4 in Algorithm 1 calculates word embedding similarities over  $n \times L$  times, i.e., the time complexity is  $\mathcal{O}(n \times L \times D)$ , where  $D$  is the dimension of each word embedding. To speed up this calculation, we leverage the vectorized instruction set (SIMD) for parallel processing. In addition, finding the intersection of two large sets in line 9 in Algorithm 1 is computationally expensive. We employ just-in-time (JIT) compilation and parallel loops for efficient comparisons of elements. To ensure reproducibility, we included all the source code as supplementary material. Also, upon acceptance of this paper, we will release our source code on GitHub and an installable package on PyPI.

We provide a demo environment<sup>4</sup> for users to explore the capabilities of our method and experience how diverse and linguistically natural the search results are and how quickly results can be obtained. The corpora used in experiments in Sections 4.2 and 4.3 are available. To prevent excessive output, for English and Japanese, only subsets of the corpora have been incorporated. Furthermore, search results are presented in smaller batches—to improve usability—with additional results available upon request, rather than all at once.

<sup>3</sup> The word *Matcha* means powdered green tea in Japanese and is here a pun for ‘matcher’. *Matcha softo* in Japanese means soft-serve ice cream of green tea flavor.

<sup>4</sup> Anonymized URL: <http://54.238.189.64/>. The screenshot of the demo is presented in Fig. 3, Appendix C. Upon acceptance, we will release the demo on the web publicly in conjunction with the source code.

Table 2: Running time (sec) of indexing and search in the English and Japanese Wikipedia articles. JIT compilation only affects the search speed.

	En (3.4B words)		Ja (1.1B words)	
	Indexing	Search	Indexing	Search
Exact matching	685.8	0.005	242.5	0.022
<i>SoftMatcha</i>	685.8	0.098	242.5	0.055
— JIT compile	685.8	0.107	242.5	0.082
Dense vector search	1036.5	0.389	320.4	0.283

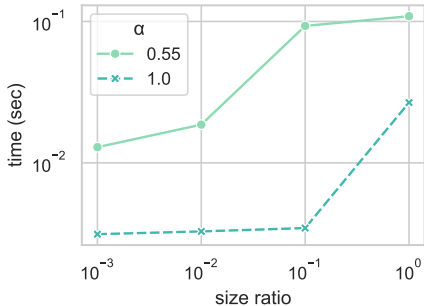


Figure 2: Running time on subsampled English Wikipedia corpora of varying sizes.

#### 4.2 CASE STUDY IN NATURAL LANGUAGE PROCESSING — BILLION-SCALE CORPUS SEARCH

In this section, we simulate scenarios where NLP researchers use *SoftMatcha* on billion-scale English and Japanese corpora. We focus on two types of evaluations: (i) quantitative evaluation — we assess the running time using large data comparable in size to those used for training standard large LMs; and (ii) qualitative evaluation — we examine the tool’s effectiveness in detecting toxic instances within large LM training corpora.

**Why English and Japanese?** English is the dominant language in modern NLP and accounts for the largest portion of training data for large LMs (Brown et al., 2020; Chowdhery et al., 2022; Touvron et al., 2023). In contrast, Japanese presents a typologically distinct language, with orthographic and structural features that differ significantly from English, such as: (i) character types (alphabetic vs. hiragana, katakana, and kanji); (ii) syntax (SVO vs. SOV, head-initial vs. head-final); and (iii) word segmentation (space-separated vs. no word separation) (Haspelmath et al., 2005). Furthermore, the substantial size of Japanese corpora<sup>5</sup> makes it suitable for testing the scalability of our approach.

**Setup. Corpora:** we utilized the LLM-jp corpus v2.0 (LLM-jp et al., 2024), which contains organized collections of Wikipedia articles in both English (3.4B words) and Japanese (1.1B words). **Word embeddings:** we used GloVe glove-wiki-gigaword-300 (Pennington et al., 2014) with a threshold  $\alpha = 0.55$  for the English word embeddings, and fastText facebook/fasttext-ja-vectors (Grave et al., 2018) with a threshold  $\alpha = 0.50$  for the Japanese word embeddings. **Baseline methods:** we compared the search results of *SoftMatcha* with those of exact matching, i.e., pattern matching with a standard inverted index, and dense vector search using intfloat/multilingual-e5-large model (Wang et al., 2024b). In dense vector search, we encoded each training instance in the corpus and built the graph-based index using hierarchical navigable small worlds (HNSW) (Malkov & Yashunin, 2020) FAISS implementation (Douze et al., 2024) for the approximate nearest neighbor search. **Computational environment:** we measured the running time on 152 core CPUs (Intel® Xeon® Platinum 8368 CPU @ 2.40GHz) and a 226 GiB main memory for the exact matching and *SoftMatcha*, and on 8 NVIDIA A100 GPUs for the dense vector search.

**Running time.** We measured the indexing time and search time in both the English and Japanese Wikipedia articles. Table 2 demonstrates that *SoftMatcha* is faster than dense vector search in both indexing and search time and takes the same indexing time as exact matching. Next, we investigated the relationship between the search speed and the corpus size. We constructed subsets of the English Wikipedia, whose sizes are  $\{10^{-3}, 10^{-2}, 10^{-1}\}$  times that of the original corpus of 3.4B tokens, by randomly sampling from the original corpus, and measured the search time with each subset. Figure 2 shows the results. We observed that the search time increased only sublinearly, not linearly, with respect to the increase in corpus size. We thus confirmed that our algorithm works effectively for billion-scale corpus searches.

<sup>5</sup> Japanese ranks fifth in the number of Wikipedia entries, while all the other languages in the top 10 use an alphabetic writing system (Wikipedia, 2024).



Table 3: Results of billion-scale corpus search in the English and Japanese Wikipedia articles.

	Exact matching	SoftMatcha	Dense vector search
Query: “homemade bombs” in English Wikipedia (3.4B words)			
# Hits	107	1,473	n/a (depends on the top- <i>k</i> )
Match examples	<i>homemade bombs</i>	<i>homemade bombs</i> <i>home-made grenades</i> <i>homemade missiles</i>	Article: Survival Under Atomic Attack Article: Mark 24 nuclear bomb Article: List of common misconceptions
Query: “手製爆弾 ( <i>homemade bombs</i> )” in Japanese Wikipedia (1.1B words)			
# Hits	27	42	n/a (depends on the top- <i>k</i> )
Match examples	手製爆弾 ( <i>homemade bombs</i> )	手製爆弾 ( <i>homemade bombs</i> ) 手製手榴弾 ( <i>home-made grenades</i> )	Article: 花火 ( <i>fireworks</i> ) Article: 手持ち花火 ( <i>consumer fireworks</i> )

**Search results.** Table 3 shows the results of the billion-scale corpus search. We queried “homemade bombs” in the English Wikipedia corpus and “手製爆弾 (*homemade bombs*)” in the Japanese Wikipedia corpus. In the dense vector search, we selected some retrieved examples from the top-10 search results. The table demonstrates that our SoftMatcha extends the exact matching. Note that the results of exact matching are a subset of the results of SoftMatcha. In Japanese, the dense vector search retrieved an article of “花火” (*fireworks*), which does not contain the contents related to the query. In contrast, SoftMatcha lists all contents that have the query pattern or its similar pattern. In addition, while dense vector search allows semantic similarity search, it cannot identify where the query pattern is located in a text. To summarize, a text that exactly contains the query pattern is not always retrieved in dense vector search, i.e., the recall is not always 100%, while exact matching and SoftMatcha can return all texts that contain the query pattern, and SoftMatcha also enables matching of semantically similar patterns.

#### 4.3 CASE STUDY IN CORPUS LINGUISTICS — RETRIEVING LATIN EXAMPLES

**Why Latin?** Latin is a morphologically complex fusional language where a lemma (dictionary form) may exhibit numerous different word forms depending on the morphological features (*e.g.*, voice, mood, tense, aspect, person, and number for verbals; gender, number, and case for nominals) that the word bears. For example, a transitive verb may conjugate to more than 100 different finite verb forms. This morphological complexity makes it harder to search through a corpus by exact matching. Furthermore, due to Latin’s philological significance and the vast body of accumulated literature, there has been a persistent demand for advanced search tools to facilitate corpus analysis (Bamman & Smith, 2012). For these reasons, Latin is a suitable touchstone to test the utility of our proposed soft pattern matcher, particularly for humanities researchers and language learners.

**Setup. Corpora:** we used two corpora: one from the Perseus Project (Crane, 2023) (5M tokens) and the Augustinian Sermon Parallelism (ASP) Dataset (Bothwell et al., 2023) (0.1M tokens). **Word embeddings:** we used the pre-trained fastText embeddings facebook/fasttext-la-vectors (Grave et al., 2018).

**Search results.** Table 4 shows the returned matches with the queries *factus est* (‘it/he is done’ or ‘it/he is made’) and *equus est* (‘it is a horse’). It is evident that SoftMatcha effectively links the queries to semantically similar words, as seen in *mortuus* ‘dead’ and *creatus* ‘created’ matched with the query *factus* ‘done, finished, made’, and *bos* ‘cow’, *currus* ‘chariot’, and *Minotaurus* ‘Minotaur’ with *equus* ‘horse’. Interestingly, the tool is also able to catch different word forms with different morphological features while sharing the same lemma, which are highlighted in pink. For example, although the Latin copula verb in the query *est* ‘is’ exhibits highly irregular conjugation patterns, the matches successfully include their inflected forms such as *sunt*, *esset*, *erat*, and *fuit*. Furthermore, the matches *mortuus* and *creatus* are not only semantically similar to the query word *factus* but also have morphological features in common (perfect, participle, masculine, nominative, and singular).

Table 4: Latin match examples by SoftMatcha with queries *factus est* ‘he/it was done/finished/made’ and *equus est* ‘it is a horse’. The top row of an interlinear gloss represents words, the middle row their morphological analysis, and the bottom row the free translation. Morphological matches are highlighted in pink, semantics matches in blue, and exact matches in green. The numbers in parentheses represent the cosine similarities between the query words and their corresponding matched words in the embedding space.

492	Query: <i>factus est</i>	Query: <i>equus est</i>
493	<i>fact-us</i> <i>est</i>	<i>equus</i> <i>est</i>
494	do.PASS.PF.PTCP-M.NOM.SG <b>be</b> .IND.PRS.3SG	horse.M.NOM.SG <b>be</b> .IND.PRS.3SG
495	<b>‘he/it is done/finished/made’</b>	<b>‘it is a horse’</b>
496		
497	Match: <i>facta sunt</i> (0.56, 0.56)	Match: <i>bos est</i> (0.48, 1.00)
498	<i>fact-a</i> <i>sunt</i>	<i>bos</i> <i>est</i>
499	do.PASS.PF.PTCP-N.NOM.PL <b>be</b> .IND.PRS.3PL	cow.F.NOM.SG <b>be</b> .IND.PRS.3SG
500	<b>‘they are done’</b> or ‘they are facts’	<b>‘it is a cow.’</b>
501	Match: <i>mortuus esset</i> (0.53, 0.58)	Match: <i>currus fuit</i> (0.44, 0.63)
502	<i>mortu-us</i> <i>esset</i>	<i>currus</i> <i>fuit</i>
503	<b>die</b> .ACT.PF.PTCP-M.NOM.SG <b>be</b> .SUB.IMPF.3SG	chariot.M.NOM.SG <b>be</b> .IND.PF.3SG
504	<b>‘he/it was dead’</b>	<b>‘it was a chariot’</b>
505	Match: <i>creatus erat</i> (0.65, 0.65)	Match: <i>Minotaurus esset</i> (0.49, 0.58)
506	<i>creat-us</i> <i>erat</i>	<i>Minotaurus</i> <i>esset</i>
507	<b>create</b> .PASS.PF.PTCP-M.NOM.SG <b>be</b> .IND.IMPF.3SG	<b>Minotaur</b> .M.NOM.SG <b>be</b> .SUB.IMPF.3SG
508	<b>‘he/it was created’</b>	<b>‘it would be the Minotaur’</b>
509		

## 5 CONCLUSION

In this paper, we propose a new pattern-matching algorithm that can flexibly handle the orthographic diversity of natural languages while also performing efficiently on large-scale corpora. Our algorithm, combining word embeddings and inverted indexing, achieves inference speed independent of corpus size. We have also developed and released a simple-to-use web demo for researchers and practitioners. For quantitative evaluation, we confirm that the developed tool can enumerate all search results within one second on billion-scale corpora, a typical scenario in training large LMs. For qualitative evaluation, we assess the tool in typical scenarios in NLP (detecting harmful examples from large-scale Japanese and English Wikipedia corpora) and computational linguistics (example retrieval from Latin, a language with highly diverse inflections), observing that our method can retrieve semantically similar examples that hard pattern matchers would miss.

## ETHICS STATEMENT

The corpora used for training large LMs, composed of web corpora, including Common Crawl (Common Crawl, 2024) and digitized book data such as (Zhu et al., 2015), can be said to encompass a substantial portion of humanity’s linguistic knowledge. Inevitably, these corpora include a significant amount of ethically problematic content. This includes personal information (Subramani et al., 2023), as well as information that could be “beneficial” for violence and terrorism (or in more conventional terms, harmful information) (Albalak et al., 2024). Moreover, as these corpora continue to grow exponentially and our languages possess numerous paraphrases, comprehensively identifying such harmful learning resources is far from a simple task. Our research aims to substantially simplify this crucial activity for human ethics: identifying harmful learning instances within vast linguistic resources Section 4.2. We hope that, by leveraging our tool, NLP researchers and developers will contribute to providing LLMs as a safe and reliable social infrastructure.

## REPRODUCIBILITY STATEMENT

Our source code is attached to the submission form. We will make the source code available on GitHub. For `SoftMatcha`, we will also release an installable package on PyPI upon acceptance. In all experiments, we used open datasets. Further details on the embeddings, corpora, computational environment, and protocols used in experiments are explained in Section 4.

## REFERENCES

- Alon Albalak, Yanai Elazar, Sang Michael Xie, Shayne Longpre, Nathan Lambert, Xinyi Wang, Niklas Muennighoff, Bairu Hou, Liangming Pan, Haewon Jeong, Colin Raffel, Shiyu Chang, Tatsunori Hashimoto, and William Yang Wang. A Survey on Data Selection for Language Models. *arXiv [cs.CL]*, 26 February 2024. URL <http://arxiv.org/abs/2402.16827>.
- Laurence Anthony. A critical look at software tools in corpus linguistics. *Linguistic Research*, 30(2):141–161, 2013.
- Orlando Ayala and Patrice Bechard. Reducing hallucination in structured outputs via retrieval-augmented generation. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 6: Industry Track)*, pp. 228–238, Stroudsburg, PA, USA, 2024. Association for Computational Linguistics. URL <https://aclanthology.org/2024.naacl-industry.19.pdf>.
- John Bambenek and Agnieszka Klus. *Grep Pocket Reference: A Quick Pocket Reference for a Utility Every UNIX User Needs*. O’Reilly Media, 2009.
- David Bamman and David Smith. Extracting two thousand years of latin from a million book library. *J. Comput. Cult. Herit.*, 5(1), April 2012. ISSN 1556-4673. doi: 10.1145/2160165.2160167. URL <https://doi.org/10.1145/2160165.2160167>.
- Ezio Bartocci, Yliès Falcone, Adrian Francalanza, and Giles Reger. Introduction to runtime verification. In *Lectures on Runtime Verification*, volume 10457 of *Lecture Notes in Computer Science*, pp. 1–33. Springer, 2018.
- Douglas Biber. Corpus-Based and Corpus-Driven Analyses of Language Variation and Use. In *The Oxford Handbook of Linguistic Analysis*. Oxford Academic, 2nd edition, 2015. doi: 10.1093/oxfordhb/9780199677078.013.0008.
- Stella Biderman, USVSN Sai Prashanth, Lintang Sutawika, Hailey Schoelkopf, Quentin Anthony, Shivanshu Purohit, and Edward Raff. Emergent and Predictable Memorization in Large Language Models, 2023a. URL <https://arxiv.org/abs/2304.11158>.
- Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, Usvsn Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. Pythia: A suite for analyzing large language models across training and scaling. *arXiv [cs.CL]*, 3 April 2023b. URL <http://arxiv.org/abs/2304.01373>.

- 594 Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching word vectors  
595 with subword information. *Trans. Assoc. Comput. Linguist.*, 5:135–146, December 2017. URL  
596 <https://aclanthology.org/Q17-1010.pdf>.  
597
- 598 Stephen Bothwell, Justin DeBenedetto, Theresa Crnkovich, Hildegund Müller, and David Chiang.  
599 Introducing Rhetorical Parallelism Detection: A New Task with Datasets, Metrics, and Base-  
600 lines. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference*  
601 *on Empirical Methods in Natural Language Processing*, pp. 5007–5039, Singapore, December  
602 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.305. URL  
603 <https://aclanthology.org/2023.emnlp-main.305>.
- 604 Robert S. Boyer and J Strother Moore. A Fast String Searching Algorithm. *Commun. ACM*, 20  
605 (10):762–772, 1977. doi: 10.1145/359842.359859. URL [https://doi.org/10.1145/](https://doi.org/10.1145/359842.359859)  
606 [359842.359859](https://doi.org/10.1145/359842.359859).
- 607 Leonid Boytsov. Indexing methods for approximate dictionary searching: Comparative analysis.  
608 *ACM J. Exp. Algorithmics*, 16(1), 2011.  
609
- 610 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-  
611 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-  
612 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,  
613 Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Ma-  
614 teusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam Mc-  
615 Candlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language Models are  
616 Few-Shot Learners. *Advances in Neural Information Processing Systems*, 33:1877–1901,  
617 2020. URL [https://proceedings.neurips.cc/paper\\_files/paper/2020/](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf)  
618 [file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf).
- 619 Michael Burrows, D J Wheeler D I G I T A L, Robert W. Taylor, David J. Wheeler, and David  
620 Wheeler. A block-sorting lossless data compression algorithm. 1994. URL [https://api.](https://api.semanticscholar.org/CorpusID:2167441)  
621 [semanticscholar.org/CorpusID:2167441](https://api.semanticscholar.org/CorpusID:2167441).
- 622 Julius Caesar. *Commentarii de bello gallico. With explanatory notes, lexicon, maps, indexes, etc.*  
623 Eldredge & Brother, Philadelphia, 1882.  
624
- 625 Yi-Pei Chen, Noriki Nishida, Hideki Nakayama, and Yuji Matsumoto. Recent Trends in Person-  
626 alized Dialogue Generation: A Review of Datasets, Methodologies, and Evaluations. In *Pro-*  
627 *ceedings of the 2024 Joint International Conference on Computational Linguistics, Language*  
628 *Resources and Evaluation (LREC-COLING 2024)*, pp. 13650–13665, 2024. URL [https://](https://aclanthology.org/2024.lrec-main.1192.pdf)  
629 [aclanthology.org/2024.lrec-main.1192.pdf](https://aclanthology.org/2024.lrec-main.1192.pdf).
- 630 Yuanyuan Chen, Boyang Li, Han Yu, Pengcheng Wu, and Chunyan Miao. HyDRA: Hypergra-  
631 dient Data Relevance Analysis for Interpreting Deep Neural Networks. *Proc. Conf. AAAI Ar-*  
632 *tif. Intell.*, 35(8):7081–7089, 18 May 2021. URL [https://cdn.aaai.org/ojs/16871/](https://cdn.aaai.org/ojs/16871/16871-13-20365-1-2-20210518.pdf)  
633 [16871-13-20365-1-2-20210518.pdf](https://cdn.aaai.org/ojs/16871/16871-13-20365-1-2-20210518.pdf).
- 634 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam  
635 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh,  
636 Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam  
637 Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James  
638 Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Lev-  
639 skaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin  
640 Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret  
641 Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick,  
642 Andrew M Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica  
643 Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Bren-  
644 nan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas  
645 Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. PaLM: Scaling language modeling with Pathways.  
646 *arXiv [cs.CL]*, 5 April 2022. URL <http://arxiv.org/abs/2204.02311>.
- 647 Common Crawl. Common Crawl - Open Repository of Web Crawl Data, 2024. URL [https://](https://commoncrawl.org/)  
[commoncrawl.org/](https://commoncrawl.org/).

- 648 Matthew Coole, Paul Rayson, and John Mariani. LexiDB: Patterns & Methods for Corpus Linguistic  
649 Database Management. In *Proceedings of the Twelfth Language Resources and Evaluation*  
650 *Conference*, pp. 3128–3135, 2020. URL [https://aclanthology.org/2020.lrec-1.](https://aclanthology.org/2020.lrec-1.383.pdf)  
651 [383.pdf](https://aclanthology.org/2020.lrec-1.383.pdf).
- 652 Gregory Crane. The Perseus Digital Library and the future of libraries. *International Journal on*  
653 *Digital Libraries*, 2023. doi: 10.1007/s00799-022-00333-2.
- 654 Chris Culy and Verena Lyding. Double Tree: An Advanced KWIC Visualization for Expert Users.  
655 In *2010 14th International Conference Information Visualisation*, pp. 98–103, 2010. doi: 10.  
656 1109/IV.2010.24.
- 657 Daniel De Freitas. Announcing our Series A and our new AI model, C1.2. [https://blog.](https://blog.character.ai/character-ai/)  
658 [character.ai/character-ai/](https://blog.character.ai/character-ai/), 23 March 2023. Accessed: 2024-10-1.
- 659 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of  
660 Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Con-*  
661 *ference of the North American Chapter of the Association for Computational Linguistics*, pp.  
662 4171–4186, Stroudsburg, PA, USA, 2019. Association for Computational Linguistics. URL  
663 <https://aclanthology.org/N19-1423.pdf>.
- 664 Jan Don. *Morphological theory and the morphology of English*. Edinburgh University Press, 2014.  
665 doi: 10.1515/9780748645145.
- 666 Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-  
667 Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The Faiss library, 2024.  
668 URL <https://arxiv.org/abs/2401.08281>.
- 669 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
670 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony  
671 Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark,  
672 Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere,  
673 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris  
674 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong,  
675 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny  
676 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino,  
677 Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael  
678 Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Ander-  
679 son, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah  
680 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan  
681 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-  
682 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy  
683 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak,  
684 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-  
685 wala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini,  
686 Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yearry, Laurens van der  
687 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo,  
688 Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Man-  
689 nat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova,  
690 Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal,  
691 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur  
692 Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhar-  
693 gava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong,  
694 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic,  
695 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sum-  
696 baly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa,  
697 Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang,  
698 Sharath Rparathy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende,  
699 Soumya Batra, Spencer Whitman, Sten Sootla, Stéphane Collot, Suchin Gururangan, Sydey  
700 Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom,  
701 Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta,

702 Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-  
703 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang,  
704 Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur,  
705 Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre  
706 Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha  
707 Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand,  
708 Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein,  
709 Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples,  
710 Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco,  
711 Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman,  
712 Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto  
713 De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Bran-  
714 don Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina  
715 Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai,  
716 Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li,  
717 Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana  
718 Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil,  
719 Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Ar-  
720 caute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco  
721 Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella  
722 Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory  
723 Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang,  
724 Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Gold-  
725 man, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman,  
726 James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer  
727 Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe  
728 Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie  
729 Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun  
730 Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal  
731 Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva,  
732 Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian  
733 Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson,  
734 Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Ke-  
735 neally, Michael L Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel  
736 Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mo-  
737 hammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navy-  
738 ata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong,  
739 Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli,  
740 Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux,  
741 Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao,  
742 Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li,  
743 Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott,  
744 Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Sa-  
745 tadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-  
746 say, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang  
747 Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen  
748 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho,  
749 Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser,  
750 Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Tim-  
751 othy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan,  
752 Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu  
753 Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Con-  
754 stable, Xiaocheng Tang, Xiaofang Wang, Xiaoqian Wu, Xiaolan Wang, Xide Xia, Xilun Wu,  
755 Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi,  
Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef  
Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The Llama 3 herd of models. *arXiv  
[cs.AI]*, 31 July 2024. URL <http://arxiv.org/abs/2407.21783>.

- 756 David Embick. *The morpheme: A theoretical introduction*. De Gruyter Mouton, Berlin, München,  
757 Boston, 2015. doi: 10.1515/9781501502569.
- 758
- 759 Paolo Ferragina and Giovanni Manzini. Indexing compressed text. *J. ACM*, 52(4):552–581, July  
760 2005. ISSN 0004-5411. doi: 10.1145/1082036.1082039. URL [https://doi.org/10.](https://doi.org/10.1145/1082036.1082039)  
761 [1145/1082036.1082039](https://doi.org/10.1145/1082036.1082039).
- 762
- 763 Andrew Gallant. ripgrep: ripgrep recursively searches directories for a regex pattern while respect-  
764 ing your gitignore, 2024. URL <https://github.com/BurntSushi/ripgrep>.
- 765
- 766 Juri Ganitkevitch, Benjamin Van Durme, and Chris Callison-Burch. PPDB: The Paraphrase  
767 Database. In *Proceedings of the 2013 Conference of the North American Chapter of the As-*  
768 *sociation for Computational Linguistics: Human Language Technologies*, pp. 758–764, 2013.  
769 URL <https://aclanthology.org/N13-1092.pdf>.
- 770
- 771 Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason  
772 Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The pile:  
773 An 800GB dataset of diverse text for language modeling. *arXiv [cs.CL]*, 31 December 2020.  
774 URL <http://arxiv.org/abs/2101.00027>.
- 775
- 776 Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng  
777 Wang, and Haofen Wang. Retrieval-augmented generation for large language models: A survey.  
778 *arXiv [cs.CL]*, 18 December 2023. URL <http://arxiv.org/abs/2312.10997>.
- 779
- 780 Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A Smith. Real-  
781 ToxicityPrompts: Evaluating neural toxic degeneration in language models. In *Findings of*  
782 *the Association for Computational Linguistics: EMNLP 2020*, pp. 3356–3369, Stroudsburg,  
783 PA, USA, November 2020. Association for Computational Linguistics. URL [https://](https://aclanthology.org/2020.findings-emnlp.301.pdf)  
784 [aclanthology.org/2020.findings-emnlp.301.pdf](https://aclanthology.org/2020.findings-emnlp.301.pdf).
- 785
- 786 Adrien Grand, Robert Muir, Jim Ferenczi, and Jimmy Lin. From MAXSCORE to Block-Max Wand:  
787 The Story of How Lucene Significantly Improved Query Evaluation Performance. In *ECIR (2)*,  
788 volume 12036 of *Lecture Notes in Computer Science*, pp. 20–27. Springer, 2020.
- 789
- 790 Edouard Grave, Piotr Bojanowski, Prakhar Gupta, Armand Joulin, and Tomas Mikolov. Learning  
791 Word Vectors for 157 Languages, 2018. URL <https://arxiv.org/abs/1802.06893>.
- 792
- 793 Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. A Survey on Automated Fact-Checking.  
794 *Transactions of the Association for Computational Linguistics*, 10:178–206, 2022. doi: 10.1162/  
795 [tacl.a\\_00454](https://doi.org/10.1162/tacl.a_00454). URL <https://aclanthology.org/2022.tacl-1.11>.
- 796
- 797 Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. REALM: Retrieval-  
798 Augmented Language Model Pre-Training. In *Proceedings of the 37th International Conference*  
799 *on Machine Learning, ICML’20*. JMLR.org, 2020.
- 800
- 801 Le Quan Ha, Philip Hanna, Ji Ming, and F. J. Smith. Extending Zipf’s law to n-grams for large  
802 corpora. *Artificial Intelligence Review*, 2009. doi: 10.1007/s10462-009-9135-4.
- 803
- 804 Saqib Hakak, Amirrudin Kamsin, Palaiahnakote Shivakumara, Gulshan Amin Gilkar, Wazir Zada  
805 Khan, and Muhammad Imran. Exact String Matching Algorithms: Survey, Issues, and Future  
806 Research Directions. *IEEE Access*, 7:69614–69637, 2019. doi: 10.1109/ACCESS.2019.2914071.  
807 URL <https://doi.org/10.1109/ACCESS.2019.2914071>.
- 808
- 809 Martin Haspelmath, Martin S Dryer, David Gil, and Bernard Comrie. *The World Atlas of Language*  
810 *Structures*. Oxford University Press, London, England, 21 July 2005.
- 811
- 812 H. S. Heap. Information retrieval: Computational and theoretical aspects. *Libr. Q.*, 50(1):153–154,  
813 January 1980. URL <http://dx.doi.org/10.1086/629887>.
- 814
- 815 Susan Hockey and Jeremy Martin. The Oxford Concordance Program version 2. *Literary and*  
816 *Linguistic Computing*, 2(2):125–131, January 1987. doi: 10.1093/lc/2.2.125.

- 810 Jie Huang, Hanyin Shao, and Kevin Chen-Chuan Chang. Are Large Pre-Trained Language Models  
811 Leaking Your Personal Information? In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 2038–2047, December 2022. URL <https://aclanthology.org/2022.findings-emnlp.148.pdf>.
- 814 Daphne Ippolito, Florian Tramer, Milad Nasr, Chiyuan Zhang, Matthew Jagielski, Katherine Lee,  
815 Christopher Choquette Choo, and Nicholas Carlini. Preventing generation of verbatim memoriza-  
816 tion in language models gives a false sense of privacy. In C. Maria Keet, Hung-Yi Lee, and Sina  
817 Zarrieß (eds.), *Proceedings of the 16th International Natural Language Generation Conference*,  
818 pp. 28–53, Prague, Czechia, September 2023. Association for Computational Linguistics. doi: 10.  
819 18653/v1/2023.inlg-main.3. URL <https://aclanthology.org/2023.inlg-main.3>.
- 820 Masaru Isonuma and Ivan Titov. Unlearning Traces the Influential Training Data of Language Mod-  
821 els. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*  
822 *(Volume 1: Long Papers)*, pp. 6312–6325, Stroudsburg, PA, USA, 2024. Association for Compu-  
823 tational Linguistics. URL <https://aclanthology.org/2024.acl-long.343.pdf>.
- 824 Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand  
825 Joulin, and Edouard Grave. Unsupervised Dense Information Retrieval with Contrastive Learn-  
826 ing, 2022. URL <https://arxiv.org/abs/2112.09118>.
- 827  
828 Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane  
829 Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. Atlas: Few-shot Learning  
830 with Retrieval Augmented Language Models. *J. Mach. Learn. Res.*, 24(1), March 2024. ISSN  
831 1532-4435.
- 832  
833 Hervé Jégou, Matthijs Douze, and Cordelia Schmid. Product Quantization for Nearest Neighbor  
834 Search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(1):117–128, 2011.
- 835  
836 Kim Ebensgaard Jensen. Linguistics in digital humanities: (computational) corpus linguistics.  
837 *MedieKultur: Journal of Media and Communication Research*, 30(57), 2014. doi: 10.7146/  
838 mediekultur.v30i57.15968.
- 839  
840 Leonard B Jung, Jonas A Gudera, Tim L T Wiegand, Simeon Allmendinger, Konstantinos Dimi-  
841 triadis, and Inga K Koerte. ChatGPT Passes German State Examination in Medicine With Pic-  
842 ture Questions Omitted. *Dtsch. Arztebl. Int.*, 120(21):373–374, 30 May 2023. URL <http://dx.doi.org/10.3238/arztebl.m2023.0113>.
- 843  
844 Daniel Jurafsky and James H. Martin. *Speech and Language Processing: An Introduction to Natural*  
845 *Language Processing, Computational Linguistics, and Speech Recognition with Language Mod-*  
846 *els*. 3rd edition, 2024. URL <https://web.stanford.edu/~jurafsky/slp3/>. Online  
847 manuscript released August 20, 2024.
- 848  
849 Richard M. Karp and Michael O. Rabin. Efficient Randomized Pattern-Matching Algorithms. *IBM*  
850 *J. Res. Dev.*, 31(2):249–260, 1987.
- 851  
852 Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi  
853 Chen, and Wen-tau Yih. Dense Passage Retrieval for Open-Domain Question Answering. In  
854 *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*  
855 *(EMNLP)*, pp. 6769–6781, Online, November 2020. Association for Computational Linguistics.  
856 doi: 10.18653/v1/2020.emnlp-main.550. URL <https://www.aclweb.org/anthology/2020.emnlp-main.550>.
- 857  
858 Daniel Martin Katz, Michael James Bommarito, Shang Gao, and Pablo Arredondo. GPT-4 Passes  
859 the Bar Exam. *Philos. Trans. A Math. Phys. Eng. Sci.*, 382(2270):20230254, 15 April 2024. URL  
860 <http://dx.doi.org/10.1098/rsta.2023.0254>.
- 861  
862 Urvasi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. General-  
863 ization through Memorization: Nearest Neighbor Language Models. In *International Confer-*  
*ence on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HklBjCEKvH>.



- 864 Omar Khattab and Matei Zaharia. ColBERT: Efficient and Effective Passage Search via Contextu-  
865 alized Late Interaction over BERT. In *Proceedings of the 43rd International ACM SIGIR Confer-*  
866 *ence on Research and Development in Information Retrieval*, New York, NY, USA, 25 July 2020.  
867 ACM. URL <https://dl.acm.org/doi/10.1145/3397271.3401075>.
- 868 Donald E. Knuth, James H. Morris Jr., and Vaughan R. Pratt. Fast pattern matching in strings. *SIAM*  
869 *J. Comput.*, 6(2):323–350, 1977.
- 870 Tatsuru Kobayashi and Kumiko Tanaka-Ishii. Taylor’ s law for Human Linguistic Sequences. In  
871 *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume*  
872 *1: Long Papers)*, pp. 1138–1148, Stroudsburg, PA, USA, 2018. Association for Computational  
873 Linguistics. URL <https://aclanthology.org/P18-1105.pdf>.
- 874 Pang Wei Koh and Percy Liang. Understanding Black-box Predictions via Influence Functions.  
875 *arXiv [stat.ML]*, 14 March 2017. URL <http://arxiv.org/abs/1703.04730>.
- 876 Sachin Kumar, Vidhisha Balachandran, Lucille Njoo, Antonios Anastasopoulos, and Yulia Tsvetkov.  
877 Language Generation Models Can Cause Harm: So What Can We Do About It? An Actionable  
878 Survey. In *Proceedings of the 17th Conference of the European Chapter of the Association for*  
879 *Computational Linguistics*, pp. 3299–3321, Stroudsburg, PA, USA, 2023. Association for Com-  
880 putational Linguistics. URL [https://aclanthology.org/2023.eacl-main.241.](https://aclanthology.org/2023.eacl-main.241.pdf)  
881 pdf.
- 882 Siu Kwan Lam, Antoine Pitrou, and Stanley Seibert. Numba: a LLVM-based Python JIT compiler.  
883 In *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC, LLVM ’15*,  
884 New York, NY, USA, 2015. Association for Computing Machinery. ISBN 9781450340052. doi:  
885 10.1145/2833157.2833162. URL <https://doi.org/10.1145/2833157.2833162>.
- 886 B Langmead, C Trapnell, M Pop, and S Salzberg. Ultrafast and memory-efficient alignment of  
887 short DNA sequences to the human genome. *Genome Biol.*, 10:R25, 2009. doi: 10.1186/  
888 gb-2009-10-3-r25.
- 889 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,  
890 Heinrich Küttler, Mike Lewis, Wen-Tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe  
891 Kiela. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *arXiv [cs.CL]*,  
892 22 May 2020. URL <http://arxiv.org/abs/2005.11401>.
- 893 Haoran Li, Dadi Guo, Wei Fan, Mingshi Xu, Jie Huang, Fanpu Meng, and Yangqiu Song. Multi-  
894 step Jailbreaking Privacy Attacks on ChatGPT. In *Findings of the Association for Comput-*  
895 *ational Linguistics: EMNLP 2023*, pp. 4138–4153, Stroudsburg, PA, USA, December 2023.  
896 Association for Computational Linguistics. URL [https://aclanthology.org/2023.](https://aclanthology.org/2023.findings-emnlp.272.pdf)  
897 findings-emnlp.272.pdf.
- 898 Jiacheng Liu, Sewon Min, Luke Zettlemoyer, Yejin Choi, and Hannaneh Hajishirzi. Infini-gram:  
899 Scaling Unbounded n-gram Language Models to a Trillion Tokens. *arXiv [cs.CL]*, 30 January  
900 2024. URL <http://arxiv.org/abs/2401.17377>.
- 901 LLM-jp, :, Akiko Aizawa, Eiji Aramaki, Bowen Chen, Fei Cheng, Hiroyuki Deguchi, Rintaro  
902 Enomoto, Kazuki Fujii, Kensuke Fukumoto, Takuya Fukushima, Namgi Han, Yuto Harada,  
903 Chikara Hashimoto, Tatsuya Hiraoka, Shohei Hisada, Sosuke Hosokawa, Lu Jie, Keisuke Ka-  
904 mata, Teruhito Kanazawa, Hiroki Kanezashi, Hiroshi Kataoka, Satoru Katsumata, Daisuke Kawa-  
905 hara, Seiya Kawano, Atsushi Keyaki, Keisuke Kiryu, Hirokazu Kiyomaru, Takashi Kodama,  
906 Takahiro Kubo, Yohei Kuga, Ryoma Kumon, Shuhei Kurita, Sadao Kurohashi, Conglong Li, Taiki  
907 Maekawa, Hiroshi Matsuda, Yusuke Miyao, Kentaro Mizuki, Sakae Mizuki, Yugo Murawaki,  
908 Ryo Nakamura, Taishi Nakamura, Kouta Nakayama, Tomoka Nakazato, Takuro Niitsuma, Jiro  
909 Nishitoba, Yusuke Oda, Hayato Ogawa, Takumi Okamoto, Naoaki Okazaki, Yohei Oseki, Shin-  
910 taro Ozaki, Koki Ryu, Rafal Rzepka, Keisuke Sakaguchi, Shota Sasaki, Satoshi Sekine, Kohei  
911 Suda, Saku Sugawara, Issa Sugiura, Hiroaki Sugiyama, Hisami Suzuki, Jun Suzuki, Toyotaro  
912 Suzumura, Kensuke Tachibana, Yu Takagi, Kyosuke Takami, Koichi Takeda, Masashi Takeshita,  
913 Masahiro Tanaka, Kenjiro Taura, Arseniy Tolmachev, Nobuhiro Ueda, Zhen Wan, Shuntaro Yada,  
914 Sakiko Yahata, Yuya Yamamoto, Yusuke Yamauchi, Hitomi Yanaka, Rio Yokota, and Koichiro  
915 Yoshino. LLM-jp: A Cross-organizational Project for the Research and Development of Fully  
916 Open Japanese LLMs, 2024. URL <https://arxiv.org/abs/2407.03963>.

- 918 Nils Lukas, Ahmed Salem, Robert Sim, Shruti Tople, Lukas Wutschitz, and Santiago Zanella-  
919 Béguelin. Analyzing Leakage of Personally Identifiable Information in Language Models. In  
920 *2023 IEEE Symposium on Security and Privacy (SP)*. IEEE, May 2023. URL <http://dx.doi.org/10.1109/sp46215.2023.10179300>.  
921
- 922 Yu A. Malkov and D. A. Yashunin. Efficient and Robust Approximate Nearest Neighbor Search  
923 Using Hierarchical Navigable Small World Graphs. *IEEE Trans. Pattern Anal. Mach. Intell.*, 42  
924 (4):824–836, April 2020. ISSN 0162-8828. doi: 10.1109/TPAMI.2018.2889473. URL <https://doi.org/10.1109/TPAMI.2018.2889473>.  
925
- 926 Tony McEnery and Andrew Hardie. *Corpus Linguistics: Method, Theory and Practice*. Cambridge  
927 University Press, 2011.  
928
- 929 Tony McEnery and Andrew Wilson. *Corpus Linguistics*. Edinburgh University Press, Edinburgh,  
930 2001. doi: 10.1515/9781474470865.  
931
- 932 Barbara McGillivray, Thierry Poibeau, and Pablo Ruiz. Digital Humanities and Natural Language  
933 Processing: “Je t’aime... Moi non plus”. *Digital Humanities Quarterly*, 14(2), 2020. doi: 10.  
934 17863/CAM.55816.  
935
- 936 Kathleen R McKeown. Paraphrasing Using Given and New Information in a Question-Answer  
937 System. In *Proceedings of the 17th annual meeting on Association for Computational Linguistics*  
938 -, pp. 67–72, Morristown, NJ, USA, 1979. Association for Computational Linguistics. URL  
939 <https://aclanthology.org/P79-1016.pdf>.
- 940 Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient Estimation of Word Rep-  
941 resentations in Vector Space. *arXiv [cs.CL]*, 16 January 2013. URL <http://arxiv.org/abs/1301.3781>.  
942
- 943 Gonzalo Navarro. A Guided Tour to Approximate String Matching. *ACM Comput. Surv.*, 33(1):  
944 31–88, 2001.  
945
- 946 OpenAI. Introducing OpenAI o1. [https://openai.com/index/  
947 introducing-openai-o1-preview/](https://openai.com/index/introducing-openai-o1-preview/), 12 September 2024. Accessed: 2024-10-1.  
948
- 949 Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word  
950 representation. In Alessandro Moschitti, Bo Pang, and Walter Daelemans (eds.), *Proceedings  
951 of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp.  
952 1532–1543, Doha, Qatar, October 2014. Association for Computational Linguistics. doi: 10.  
953 3115/v1/D14-1162. URL <https://aclanthology.org/D14-1162>.
- 954 Garima Pruthi, Frederick Liu, Satyen Kale, and Mukund Sundararajan. Estimating Training Data  
955 Influence by Tracing Gradient Descent. *Advances in Neural Information Processing Systems*,  
956 33:19920–19930, 2020. URL [https://proceedings.neurips.cc/paper\\_files/  
957 paper/2020/file/e6385d39ec9394f2f3a354d9d2b88eec-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/e6385d39ec9394f2f3a354d9d2b88eec-Paper.pdf).
- 958 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language  
959 Models are Unsupervised Multitask Learners. 2019.  
960
- 961 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi  
962 Zhou, Wei Li, and Peter J Liu. Exploring the Limits of Transfer Learning with a Unified Text-to-  
963 Text Transformer. *arXiv [cs.LG]*, 23 October 2019. URL [http://arxiv.org/abs/1910.  
964 10683](http://arxiv.org/abs/1910.10683).
- 965 Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard  
966 Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya  
967 Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy  
968 Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin, Nikola Momchev, Matt  
969 Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Olivier Bachem, Alanna  
970 Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic,  
971 Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Paterson, Ben  
Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Chris

- 972 Welty, Christopher A Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijayku-  
973 mar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica  
974 Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn  
975 Cameron, Gus Martins, Hadi Hashemi, Hanna Klimeczak-Plucińska, Harleen Batra, Harsh Dhand,  
976 Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Peng  
977 Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost van Amersfoort,  
978 Josh Gordon, Josh Lipschultz, Josh Newlan, Ju-Yeong Ji, Kareem Mohamed, Kartikeya Badola,  
979 Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene,  
980 Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago, Lilly Mc-  
981 Nealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid,  
982 Manvinder Singh, Mark Iversen, Martin Görner, Mat Velloso, Mateo Wirth, Matt Davidow,  
983 Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moyni-  
984 han, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao,  
985 Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil  
986 Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Cullit-  
987 ton, Pradeep Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Ardeshir Rokni,  
988 Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah Perrin, Sébastien M R  
989 Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan Chan,  
990 Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky, Tulsee Doshi, Vihan Jain,  
991 Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei Wei, Wenming Ye,  
992 Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan Wei, Victor Cotruta,  
993 Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli Collins, Joelle Barral,  
994 Zoubin Ghahramani, Raia Hadsell, D Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol  
995 Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya,  
996 Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi, and Alek  
997 Andreev. Gemma 2: Improving Open Language Models at a Practical Size. *arXiv [cs.CL]*, 31 July  
998 2024. URL <http://arxiv.org/abs/2408.00118>.
- 999 Marek Rogozinski and Rafal Kuc. *Elasticsearch Server, 3rd edition*. Packt Publishing, 2016.
- 1000 Timo Schick, Sahana Udupa, and Hinrich Schütze. Self-Diagnosis and Self-Debiasing: A Proposal  
1001 for Reducing Corpus-Based Bias in NLP. *Trans. Assoc. Comput. Linguist.*, 9:1408–1424, 17 Dec-  
1002 cember 2021. URL <https://aclanthology.org/2021.tacl-1.84.pdf>.
- 1003 Sarah Schulz. *The Taming of the Shrew: Non-Standard Text Processing in the Digital Humanities*.  
1004 PhD thesis, University of Stuttgart, 2018.
- 1005 Martin Schweinberger. Concordancing with R. <https://ladal.edu.au/kwics.html>,  
1006 2024.
- 1007  
1008 Satoshi Sekine. A Linguistic Knowledge Discovery Tool: Very Large Ngram Database Search with  
1009 Arbitrary Wildcards. In *Coling 2008: Companion volume: Demonstrations*, pp. 181–184, 2008.  
1010 URL <https://aclanthology.org/C08-3010.pdf>.
- 1011  
1012 Satoshi Sekine and Kapil Dalwani. Ngram Search Engine with Patterns Combining Token, POS,  
1013 Chunk and NE Information. In *Proceedings of the Seventh International Conference on Lan-  
1014 guage Resources and Evaluation (LREC’10)*, 2010. URL [http://www.lrec-conf.org/  
1015 proceedings/lrec2010/pdf/158\\_Paper.pdf](http://www.lrec-conf.org/proceedings/lrec2010/pdf/158_Paper.pdf).
- 1016  
1017 Frank Smadja. Retrieving Collocations from Text: Xtract, 1993. URL [https://  
1018 aclanthology.org/J93-1007.pdf](https://aclanthology.org/J93-1007.pdf).
- 1019  
1020 Nishant Subramani, Sasha Luccioni, Jesse Dodge, and Margaret Mitchell. Detecting Personal In-  
1021 formation in Training Corpora: an Analysis. In *Proceedings of the 3rd Workshop on Trustwor-  
1022 thy Natural Language Processing (TrustNLP 2023)*, pp. 208–220, Stroudsburg, PA, USA, 2023.  
1023 Association for Computational Linguistics. URL [https://aclanthology.org/2023.  
1024 trustnlp-1.18.pdf](https://aclanthology.org/2023.trustnlp-1.18.pdf).
- 1025 TEI Consortium. TEI P5: Guidelines for electronic text encoding and interchange, 2023. URL  
<https://www.tei-c.org/Guidelines/P5>.

- 1026 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
1027 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher,  
1028 Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy  
1029 Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn,  
1030 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel  
1031 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee,  
1032 Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra,  
1033 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi,  
1034 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh  
1035 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen  
1036 Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic,  
1037 Sergey Edunov, and Thomas Scialom. Llama 2: Open Foundation and Fine-Tuned Chat Models.  
1038 *arXiv [cs.CL]*, 18 July 2023. URL <http://arxiv.org/abs/2307.09288>.
- 1039 Janneke Van Der Zwaan, Wouter Smink, Anneke Sools, Gerben Westerhof, Bernard Veldkamp,  
1040 and Sytse Wiegersma. Flexible NLP Pipelines for Digital Humanities Research. In *4th Digital*  
1041 *Humanities Benelux Conference*, Utrecht, Netherlands, 2017.
- 1042 Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Ma-  
1043 jumder, and Furu Wei. Text Embeddings by Weakly-Supervised Contrastive Pre-training, 2024a.  
1044 URL <https://arxiv.org/abs/2212.03533>.
- 1045 Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. Mul-  
1046 tilingual e5 text embeddings: A technical report, 2024b. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2402.05672)  
1047 [2402.05672](https://arxiv.org/abs/2402.05672).
- 1048 Yuxia Wang, Revanth Gangi Reddy, Zain Muhammad Mujahid, Arnav Arora, Aleksandr Ruba-  
1049 shevskii, Jiahui Geng, Osama Mohammed Afzal, Liangming Pan, Nadav Borenstein, Aditya Pil-  
1050 lai, Isabelle Augenstein, Iryna Gurevych, and Preslav Nakov. Factcheck-bench: Fine-grained  
1051 evaluation benchmark for automatic fact-checkers, 2024c. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2311.09000)  
1052 [2311.09000](https://arxiv.org/abs/2311.09000).
- 1053 Wikipedia. Wikipedia:Multilingual statistics, 2024. URL [https://en.wikipedia.org/](https://en.wikipedia.org/wiki/Wikipedia:Multilingual_statistics)  
1054 [wiki/Wikipedia:Multilingual\\_statistics](https://en.wikipedia.org/wiki/Wikipedia:Multilingual_statistics).
- 1055 Sam Witteveen and Martin Andrews. Paraphrasing with Large Language Models. In *Pro-*  
1056 *ceedings of the 3rd Workshop on Neural Generation and Translation*, pp. 215–220, Strouds-  
1057 burg, PA, USA, November 2019. Association for Computational Linguistics. URL [https:](https://aclanthology.org/D19-5623.pdf)  
1058 [//aclanthology.org/D19-5623.pdf](https://aclanthology.org/D19-5623.pdf).
- 1059 Sun Wu and Udi Manber. Fast Text Searching Allowing Errors. *Commun. ACM*, 35(10):83–  
1060 91, 1992a. doi: 10.1145/135239.135244. URL [https://doi.org/10.1145/135239.](https://doi.org/10.1145/135239.135244)  
1061 [135244](https://doi.org/10.1145/135239.135244).
- 1062 Sun Wu and Udi Manber. Agrep – A Fast Approximate Pattern-Matching Tool. In *1992 Winter*  
1063 *USENIX Conference*, 1992b.
- 1064 Mikio Yamamoto and Kenneth W Church. Using Suffix Arrays to Compute Term Frequency and  
1065 Document Frequency for All Substrings in a Corpus. *Comput. Linguist. Assoc. Comput. Linguist.*,  
1066 27(1):1–30, March 2001. URL <https://aclanthology.org/J01-1001.pdf>.
- 1067 Yukun Zhu, Ryan Kiros, Richard Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Tor-  
1068 ralba, and Sanja Fidler. Aligning Books and Movies: Towards Story-like Visual Explan-  
1069 ations by Watching Movies and Reading Books. *arXiv [cs.CV]*, 22 June 2015. URL [http:](http://arxiv.org/abs/1506.06724)  
1070 [//arxiv.org/abs/1506.06724](http://arxiv.org/abs/1506.06724).
- 1071 George Kingsley Zipf. Human Behavior and the Principle of Least Effort: An Introduction to Human  
1072 Ecology. *Am. J. Psychol.*, 64(1):149, January 1951. URL [http://dx.doi.org/10.2307/](http://dx.doi.org/10.2307/1418618)  
1073 [1418618](http://dx.doi.org/10.2307/1418618).
- 1074 Justin Zobel and Alistair Moffat. Inverted Files for Text Search Engines. *ACM Comput. Surv.*, 38(2):  
1075 6, 2006. doi: 10.1145/1132956.1132959. URL [https://doi.org/10.1145/1132956.](https://doi.org/10.1145/1132956.1132959)  
1076 [1132959](https://doi.org/10.1145/1132956.1132959).

1080 Table 5: Additional examples of soft matches in Latin and their scores returned by `SoftMatcha`.  
 1081 The top row of an interlinear gloss represents words, the middle row their morphological analysis,  
 1082 and the bottom row the free translation. Morphological matches are highlighted in pink, semantic  
 1083 matches in blue, and exact matches in green. The scores in parentheses represent the cosine similar-  
 1084 ities between queries and their match.

1085		
1086	Query: <i>bellum gallicum</i> (book by Caesar (1882))	Query: <i>non potest</i>
1087	<i>bell-um gallic-um</i>	<i>non potest</i>
1088	war.N-NOM.SG gallic-N.NOM.SG	not can.IND.PRS.3SG
1089	‘the Gallic war’	‘he/she/it cannot’
1090	Match: <i>bellum etruscum</i> (1.00, 0.44)	Match: <i>non possit</i> (1.00, 0.59)
1091	<i>bell-um etrusc-um</i>	<i>non possit</i>
1092	war.N-NOM.SG Etruscan-N.NOM.SG	not can.SUB.PRS.3SG
1093	‘the Etruscan war’	‘he/she/it cannot’
1094	Match: <i>contra Caesarem</i> (0.49, 0.45)	Match: <i>nec posse</i> (0.53, 0.52)
1095	<i>contra Caesar-em</i>	<i>nec posse</i>
1096	against Caesar.M-ACC.SG	nor can.INF.PRS
1097	‘against Caesar’	‘not to be able’
1098	Query: <i>quo vadis</i> (Bible, John 13:36)	Query: <i>homo sapiens</i>
1099	<i>quo vadis</i>	<i>homo sapiens</i>
1100	where go.IND.PRS.2SG	human.M.NOM.SG wise.M.NOM.SG
1101	‘where do you go’	‘a wise human’
1102	Match: <i>ibi vadis</i> (0.42, 1.00)	Match: <i>homo honestus</i> (1.00, 0.39)
1103	<i>ibi vadis</i>	<i>homo honest-us</i>
1104	there go.IND.PRS.2SG	human.M.NOM.SG noble-M.NOM.SG
1105	‘you go there’	‘an honorable human’
1106	Match: <i>autem vocaris</i> (0.50, 0.48)	Match: <i>vir sincerus</i> (0.47, 0.40)
1107	<i>autem voc-aris</i>	<i>vir sincer-us</i>
1108	but summon-PASS.IND.PRS.2SG	man.M.NOM.SG pure-M.NOM.SG
1109	‘but you are summoned’	‘a pure man’
1110	Query: <i>quod erat demonstrandum</i> (Q.E.D.)	Query: <i>post meridiem</i> (p.m.)
1111	<i>quod erat demonstrand-um</i>	<i>after meridiem</i>
1112	which be.IND.IMPF.3SG show.GER-N.NOM.SG	after noon.M.ACC.SG
1113	‘which was to be shown’	‘after noon’
1114	Match: <i>haec erat forma</i> (0.46, 1.00, 0.41)	Match: <i>ante meridiem</i> (0.68, 1.00)
1115	<i>haec erat forma</i>	<i>ante meridiem</i>
1116	this.F.NOM.SG be.IND.IMPF.3SG form.F.NOM.SG	before noon.M.ACC.SG
1117	‘this was the form’	‘before noon’
1118	Match: <i>quod postea accipiamus</i> (1.00, 0.54, 0.47)	Match: <i>contra septentrionem</i> (0.51, 0.52)
1119	<i>quod postea accipiamus</i>	<i>contra septentrionem</i>
1120	which afterwards accept.ACT.SUB.PRS.1PL	against Ursa.Major-F.ACC.SG
1121	‘which we shall accept later’	‘against Ursa Major’ (i.e., ‘facing north’)
1122		
1123		
1124		
1125	A OTHER EXAMPLES OF SOFT MATCHES IN LATIN	
1126		
1127		
1128	Table 5 shows several additional examples of Latin soft matches by <code>SoftMatcha</code> .	
1129		
1130		
1131	B LIST OF GLOSSING ABBREVIATIONS	
1132		
1133	Table 6 lists the glossing abbreviations used in this paper.	

Table 6: List of the gloss abbreviations used in this paper.

Gloss	Meaning
1, 2, 3	1st, 2nd, 3rd person, respectively
ACC	accusative
ACT	active voice
F	feminine (gender)
FUT	future
GER	gerundive
IMPF	imperfective
IND	indicative mood
M	masculine (gender)
N	neuter (gender)
NOM	nominative case
PASS	passive voice
PF	perfective
PL	plural
PRS	present tense
PTCP	participle
SG	singular
SUB	subjunctive mood

The screenshot shows the SoftMatcha web interface. At the top, there is a search bar with the query "March 1, 2016" and a "Search" button. Below the search bar, the interface displays "Examples" with several buttons: "theorem 1", "march 1, 2016", "John was born in", "海鮮丼", "言語の普遍性", "フランスでは", "factus est", "equus est", and "bellum Gallicum". The "Results" section shows a search for "march 1, 2016" with 7738 hits and a search time of 0.317 sec. The results are displayed as text snippets with highlighted phrases and their corresponding scores. For example, the first result is "In 2015, Fey created and produced the television comedy Unbreakable Kimmy Schmidt with fellow 30 Rock-alumnus Robert Carlock . The series stars Ellie Kemper as the titular character who escapes from a doomsday cult and moves to New York . It also stars Fey 's former co-star Jane Krakowski, as well as Tituss Burgess ( who had previously appeared in four 30 Rock episodes ) and Carol Kane . Although it was originally produced for NBC, it was eventually sold to Netflix and immediately renewed for a second season . The show premiered on March 6 , 2015 to critical acclaim ." with a score of 1.00. Other results include "On July 16 , 2015 , the series was nominated for seven Primetime Emmy Awards, including Outstanding Comedy Series ." with a score of 0.91, and "In 2015, it was announced Fey would be the narrator for the Disney Nature film Monkey Kingdom , which was released in theaters on April 17 , 2015 ." with a score of 0.85.

Figure 3: The screenshot of our demo. Given the query “Match 1, 2016”, it returns lines including softly-matched patterns such as “July 16, 2015.”

## C WEB INTERFACE

Figure 3 shows a screenshot of our demo tool.