



# Sparkly: A Simple yet Surprisingly Strong TF/IDF Blocker for Entity Matching

Derek Paulsen  
University of Wisconsin-Madison  
Informatica Inc.  
dpaulsen2@wisc.edu

Yash Govind  
Apple Inc.  
yash\_govind@apple.com

AnHai Doan  
University of Wisconsin-Madison  
Informatica Inc.  
anhai@cs.wisc.edu

## ABSTRACT

Blocking is a major task in entity matching. Numerous blocking solutions have been developed, but as far as we can tell, blocking using the well-known tf/idf measure has received virtually no attention. Yet, when we experimented with tf/idf blocking using Lucene, we found it did quite well. So in this paper we examine tf/idf blocking in depth. We develop Sparkly, which uses Lucene to perform top-k tf/idf blocking in a distributed share-nothing fashion on a Spark cluster. We develop techniques to identify good attributes and tokenizers that can be used to block on, making Sparkly completely automatic. We perform extensive experiments showing that Sparkly outperforms 8 state-of-the-art blockers. Finally, we provide an in-depth analysis of Sparkly’s performance, regarding both recall/output size and runtime. Our findings suggest that (a) tf/idf blocking needs more attention, (b) Sparkly forms a strong baseline that future blocking work should compare against, and (c) future blocking work should seriously consider top-k blocking, which helps improve recall, and a distributed share-nothing architecture, which helps improve scalability, predictability, and extensibility.

### PVLDB Reference Format:

Derek Paulsen, Yash Govind, and AnHai Doan. Sparkly: A Simple yet Surprisingly Strong TF/IDF Blocker for Entity Matching. PVLDB, 16(6): 1507 - 1519, 2023.  
doi:10.14778/3583140.3583163

### PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/anhaidgroup/sparkly>.

## 1 INTRODUCTION

Entity matching (EM) finds data instances that refer to the same real-world entity. Most EM solutions proceed in two steps: blocking and matching. Given two tables  $A$  and  $B$  to match, the *blocking step* uses heuristics to quickly remove tuple pairs  $(a \in A, b \in B)$  judged unlikely to match. The *matching step* then applies a matcher to the remaining tuple pairs to predict match/no-match.

Both the blocking and matching steps have received significant attention (e.g., [1, 5–7, 12, 14, 26, 27, 30, 32]). In this paper we focus on the blocking step. Over the past 30 years, numerous blocking solutions have been developed. The goal is to maximize *recall* (the

fraction of true matches that survive blocking) while minimizing *the output size* and *the runtime* (see Section 2). In the past few years, as a part of the Magellan project at UW-Madison, which develops a comprehensive open-source EM platform [19], we have implemented many of the proposed blocker types, develop new blocker types [38], and applied them to many real-world EM tasks in domain sciences and industry [16]. While doing this, we found that a relatively simple blocking solution that uses the tf/idf similarity measure, as implemented in the open-source Apache Lucene library, seems to work quite well.

This is rather surprising because as far as we can tell, tf/idf based blocking has received virtually no attention. For example, the book “Data Matching” [6] and several recent EM surveys [14, 32] do not discuss any tf/idf solutions for blocking, and we are not aware of any recent work proposing tf/idf solutions. Yet we found tf/idf blocking highly promising in many informal experiments.

As a result, in this paper we perform an in-depth examination of tf/idf blocking. We begin by developing a solution called Sparkly Manual, which takes as input two tables  $A$  and  $B$  with the same schema, and outputs tuple pairs  $(a \in A, b \in B)$  judged likely to match. There are two key ideas underlying Sparkly Manual. First, *it performs top-k blocking*. For each tuple  $t$  of the larger table, say  $B$ , it finds the top  $k$  tuples in  $A$  with the highest tf/idf scores ( $k$  is pre-specified), then pairs these tuples with  $t$  and outputs the pairs.

Second, Sparkly Manual *performs the above top-k computations in a distributed shared-nothing fashion*, using Lucene on a Spark cluster (hence the name Sparkly, which stands for Spark + Lucene + Python). Specifically, it uses Lucene to build an inverted index  $I$  for table  $A$  on the driver node of the Spark cluster, ships the index  $I$  to all worker nodes, distributes the tuples of table  $B$  to the worker nodes, then uses Lucene to perform top-k computations for the tuples at the worker nodes. Thus, the worker nodes operate in parallel and share no dependencies. Each node processes a subset of tuples in  $B$ .

We compare Sparkly Manual with 8 state-of-the-art (SOTA) blockers on 15 datasets that have been extensively used in recent EM work [25, 29, 38]. *Surprisingly, Sparkly Manual outperforms all of the above blockers*. It achieves higher or comparable recall at a much smaller output size, and the performance gap is quite significant in several cases (see the experiment section).

While appealing, Sparkly Manual has a limitation. It requires the user to manually identify the attributes to be blocked on, e.g., product title, or name and phone. Then it computes the tf/idf score between any two tuples  $a \in A, b \in B$  using only these attributes, after 3-gram tokenization.

It can be difficult for users to identify good blocking attributes. So we develop Sparkly Auto, which automatically identifies a set of

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing [info@vldb.org](mailto:info@vldb.org). Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.  
Proceedings of the VLDB Endowment, Vol. 16, No. 6 ISSN 2150-8097.  
doi:10.14778/3583140.3583163

good blocking attributes, together with an appropriate tokenizer for each attribute. The key observation underlying Sparkly Auto is that *a good blocking attribute helps to discriminate between matches and non-matches*. We propose techniques to quantify discriminativeness, then to effectively search a large space for the optimal combination of attributes and tokenizers that maximizes this quantity.

We show that *Sparkly Auto achieves comparable or higher recall than Sparkly Manual, yet runs much faster*. In particular, Sparkly Auto can block large datasets at reasonable time and cost, e.g., blocking tables of 10M tuples under 100 minutes on an AWS cluster of 10 commodity nodes, costing only \$12.5, and blocking tables of 26M tuples under 130 minutes on an AWS cluster of 30 nodes, costing \$67.5. This suggests that Sparkly Auto can already be practical for many real-world EM problems.

We conclude by discussing questions that arise in light of Sparkly’s strong performance. In summary, the contributions and takeaways of this paper are as follows:

- We develop Sparkly, a tf/idf blocker that uses Lucene to perform top-k blocking on a Spark cluster. We develop techniques to automatically identify good attributes and tokenizers to block on.
- Extensive experiments show that Sparkly outperforms 8 state-of-the-art blockers. This is rather surprising because tf/idf blocking has received virtually no attention in the past 30 years. *The takeaway here is that tf/idf blocking needs more attention, and that Sparkly forms a strong baseline that future blocking work should compare against.*
- We provide an in-depth analysis of Sparkly’s performance. *The takeaway here is that future blocking work should seriously consider top-k blocking, which helps improve recall, and a distributed share-nothing architecture, which helps improve scalability, predictability, and extensibility.*
- Based on the above analysis, we identify a number of promising research directions for blocking.

For more information on Sparkly, see [34], which provides the code, all experiment datasets (except Hospital, which is private), and a longer technical report.

## 2 BLOCKING FOR ENTITY MATCHING

**EM, Blocking, Matching:** Many EM variations exist [6, 12]. A common EM variation [38], which we consider in this paper, is as follows: given two tables  $A$  and  $B$  with the same schema, find all tuple pairs  $(a \in A, b \in B)$  that refer to the same real-world entity. We call these pairs *matches*.

Considering all pairs in  $A \times B$  takes too long for large tables. So EM is typically performed in two steps: *blocking* and *matching* [6, 12]. *The blocking step* uses heuristics to quickly remove a large number of pairs judged unlikely to match. *The matching step* applies a rule- or ML-based matcher to each remaining pair, to predict match or non-match. Figure 1 illustrates these steps. Here blocking keeps only those pairs that share the same state. In this paper we focus on the blocking step.

**Existing Blocker Types, Threshold vs. Top-k Blocking:** Numerous blocking solutions have been developed (see [5, 27, 32] for

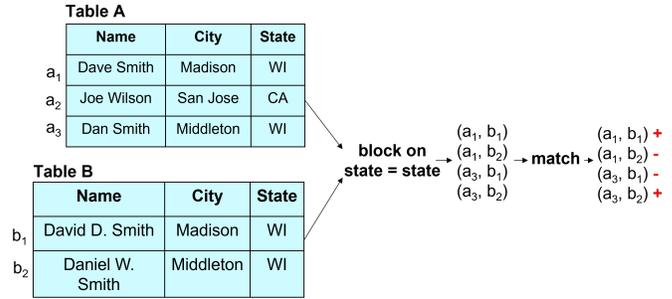


Figure 1: The blocking and matching steps of EM.

surveys). They fall roughly into five types: *sort*, *hash*, *similarity-based*, *rule-based*, and *composite*. *Sorted neighborhood* computes for each tuple a key, sorts tuples based on keys, then outputs a pair of tuples if their keys are within a pre-defined distance. *Hash-based methods* compute for each tuple a hash value (a.k.a. key), groups all tuples sharing the same hash value into a *block*, then outputs a pair of tuples if they belong to the same block. Examples of such methods include attribute equivalence, phonetic blocking, suffix array, etc. [5, 27, 32]. Most existing blocking methods are hash-based.

*Similarity-based methods* output only those tuple pairs where the similarity score between the tuples exceeds a pre-specified threshold, or one tuple is within the kNN (k-nearest) neighborhood of the other tuple [29].

We refer to the above two options as *threshold blocking* and *top-k blocking*, respectively. As we will see later, Sparkly uses top-k blocking, which we show to be critical to achieve high recall.

Similarity scores that have been considered for blocking include syntactic scores such as Jaccard, cosine, edit distance [12], and semantic scores such as those computed using word embedding/deep learning (DL) techniques [38].

*Rule-based methods* employ multiple blocking rules, where each rule can employ multiple predicates (e.g., if the Jaccard score of the titles is below 0.6 and the years are not equivalent, then the two papers do not match) [15]. Given a set of rules, the blocker figures out the best way to create a workflow and execute it using indexes [15]. Finally, *composite methods* generalizes rule-based blocking and can combine multiple blocking methods in a complex *pre-specified* workflow. Examples include canopy blocking [12] and the union of a DL method with a rule-based method in [38].

**Recent Research Directions:** In recent years researchers have pursued several directions regarding the above five blocker types [27, 32]. They have examined how to scale blocking methods (e.g., using Hadoop/Spark) [9] and how to apply DL (e.g., to develop novel hash-based [13] and similarity-based blockers [38]).

As discussed earlier, hash-based methods generate *blocks of tuples*, then output pairs whose tuples belong to the same block. A novel recent direction, called *meta-blocking*, examines how to manage these blocks (e.g., remove/prune blocks) [32]. A related direction is *token blocking*, in which each block contains all tuples that share a particular token. These blocks can be managed using meta blocking. Another interesting direction, called *schema-agnostic*, drops the assumption that Tables  $A$  and  $B$  share the same schema. The well-known JedAI EM platform implements many meta-blocking,

token-blocking, and schema-agnostic techniques [31, 33]. Other important directions include learning blockers [15], using the feedback from the matcher to improve the blocker, and explaining blockers [5, 27, 32].

**Evaluating Blockers:** Most existing works evaluate blockers in three aspects: *recall*, *output size*, and *runtime*. Let  $G \subseteq A \times B$  be the set of (unknown) gold matches, and  $C \subseteq A \times B$  be the set of tuple pairs output by a blocker  $Q$ . Then the *recall* of  $Q$  is  $|C \cap G|/|G|$ , the fraction of gold matches in the output of  $Q$ . The *output size* is  $|C|$ , and the *runtime* is measured from when the blocker receives the two tables  $A$  and  $B$  until when it outputs  $C$ .

Other aspects considered important, especially in industry, include the ease of tuning, the ability to block on arbitrarily large tables (e.g., those with billions of tuples) without crashing, extensibility (e.g., with more blocking methods/rules), the ability to estimate the total blocking time, explainability, and the ability to run the blocker easily in a variety of environments (e.g., a single laptop, a Spark cluster, a Kubernetes cluster), among others.

In this paper, we will evaluate blockers using the above three popular aspects: recall, output size, and runtime. We will briefly discuss Sparkly regarding some additional aspects, but deferring a thorough evaluation of these aspects to future work.

### 3 THE SPARKLY SOLUTION

We now describe the tf/idf measure used in keyword search (KWS), the open-source KWS library Lucene, then Sparkly, which uses Lucene to perform blocking for EM.

#### 3.1 The TF/IDF Family of Scoring Functions

TF/IDF is a well-known family of scoring functions for ranking documents in KWS [24]. To explain, consider a set of documents  $\mathcal{D} = \{D_1, \dots, D_N\}$ , where each document  $D_i$  is a string (e.g., article, email). Given a user query  $Q$ , which is also a string, we want to find documents in  $\mathcal{D}$  that are most relevant to  $Q$ . To do so, we compute a score  $s(D, Q)$  for each document  $D$ , then return the documents ranked in decreasing score.

A well-known scoring function [12], TFIDF-cosine, is as follows. First we tokenize each document  $D$  into a bag of *tokens*, also called *terms*. Next, we convert document  $D$  into a vector  $V_D$  of *weights*, one weight per term, where the weight for term  $t$  is  $V_D(t) = tf(t, D) \cdot idf(t)$ . Here  $tf(t, D)$  is the *frequency of term  $t$*  in document  $D$ , i.e., the number of times it occurs in  $D$ . The quantity  $idf(t)$  is the *inverse document frequency of term  $t$* , defined as  $\log(N/df(t))$ , where  $N$  is the number of documents in  $\mathcal{D}$ , and  $df(t)$  is the number of documents that contain term  $t$ .

We tokenize and convert query  $Q$  into a vector of weights  $V_Q$  in a similar fashion. Finally, we compute score  $s(D, Q)$  to be the cosine of the angle between the two vectors  $V_D$  and  $V_Q$ :

$$s(D, Q) = \left[ \sum_t V_D(t) \cdot V_Q(t) \right] / \left[ \sqrt{\sum_t V_D(t)^2} \cdot \sqrt{\sum_t V_Q(t)^2} \right], \quad (1)$$

where  $t$  ranges over all terms in  $D$  and  $Q$ . This definition captures the intuition that if a term  $t$  of query  $Q$  occurs often in a document  $D$ , then  $D$  is likely to be relevant to  $Q$  and score  $s(D, Q)$  should be high. This is reflected in the use of the term frequency  $tf(t, D)$ . A higher

$tf(t, D)$  leads to a higher weight for  $t$  in  $V_D$ , and consequently a higher  $s(D, Q)$ . But this should not be true if term  $t$  also occurs in many other documents. In such cases term  $t$  should be discounted, i.e., its weight in  $V_D$  should be low, and this is accomplished by multiplying the term frequency  $tf(t, D)$  with the inverse document frequency  $idf(t)$ .

Over the years, many tf/idf scoring functions have been proposed. Among them, the following function, called *Okapi BM25*, has become most popular, and is the default scoring function used by Lucene [36]:

$$s(D, Q) = \sum_{t \in Q} \frac{tf(t, D) \cdot (k_1 + 1)}{tf(t, D) + k_1 \cdot (1 - b + b \cdot \frac{|D|}{avgdl})} \cdot idf(t), \quad (2)$$

where  $idf(t) = \log(\frac{N-df(t)+0.5}{df(t)+0.5} + 1)$ , and  $k_1$  and  $b$  are free parameters, often set as  $k_1 \in [1.2, 2.0]$  and  $b = 0.75$ . The tech report explains the intuition behind BM25 (see also [36]), which has been shown to work quite well for KWS [18, 24].

#### 3.2 The Lucene KWS Library

Many open-source software for KWS have been developed. Among them Apache Lucene has become most popular [18]. The latest releases of Lucene, since 2015, have used state-of-the-art techniques in KWS to be both accurate and fast [18].

Specifically, Lucene uses BM25 as the default scoring function, ensuring highly accurate KWS results. It has also been extensively optimized, to be very fast for top-k querying, i.e., *given a query  $Q$  and a set of documents  $\mathcal{D}$ , find the top  $k$  documents in  $\mathcal{D}$  that have the highest BM25 score with  $Q$* , for a pre-specified  $k$  (typically up to a few hundreds). To do this, naively we can use an inverted index to find all documents in  $\mathcal{D}$  that share at least one term with  $Q$ , compute BM25 scores for all of them, then sort and return the top  $k$  documents. This however would be very slow, because the set of documents sharing at least one term with  $Q$  is often very large.

To solve this problem, Lucene uses a recently developed KWS technique called *block-max WAND* [3, 10, 11]. This technique allows Lucene at query time to perform a *branch-and-bound search* to find the top  $k$ . This way, Lucene can avoid examining a huge number of documents, and can generally find the top-k documents very fast, as we will see in the experiment section.

Lucene has become the library of choice for a wide variety of KWS applications. Two other popular open-source KWS systems, Solr and Elasticsearch, build on Lucene. *As a library, Lucene provides two key API functions: indexing and querying*. Solr (started in 2004) and Elasticsearch (started in 2010) use these API functions, but provide extensive support for indexing and querying a large number of documents on a cluster of machines.

#### 3.3 The Sparkly Solution

We now describe Sparkly, which takes as input two tables  $A$  and  $B$  with the same schema, and outputs a table  $C$  consisting of tuple pairs  $(a \in A, b \in B)$  judged likely to match.

To do so, Sparkly uses two key ideas. First, *it performs top-k blocking*. Specifically, it builds an inverted index  $I$  for the smaller table, say table  $A$ . Then for each tuple  $b$  in table  $B$ , it probes  $I$  to find the top  $k$  tuples in  $A$  with the highest tf/idf scores (where  $k$  is pre-specified), then pairs these tuples with  $b$  and outputs the pairs.

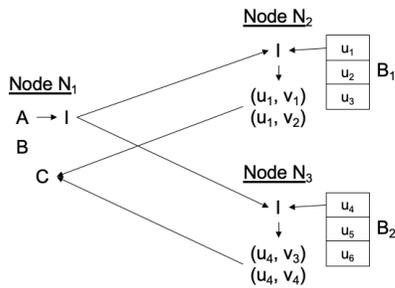


Figure 2: Sparkly's execution on a 3-node cluster.

Second, Sparkly executes the above steps in a distributed share-nothing fashion, using Lucene on a Spark cluster. We now describe the execution in detail, using the 3-node Spark cluster in Figure 2.

Build the inverted index  $I$  of table  $A$ : Suppose that tables  $A$  and  $B$  reside on the primary node  $N_1$  (see Figure 2), and that  $A$  is the smaller table, i.e., having fewer tuples than  $B$ . Sparkly chops table  $A$  horizontally into multiple chunks, each containing multiple tuples, starts multiple threads on the entire Spark cluster, sends each chunk to a thread, which calls Lucene's indexing procedure to create an inverted index for that chunk. Sparkly then combines these inverted indexes into a single inverted index  $I$  for table  $A$ , and writes  $I$  to the local disk of node  $N_1$ .

Ship index  $I$  and tuples of table  $B$  to the secondary nodes: Sparkly then ships index  $I$  to the local disks of the secondary nodes  $N_2$  and  $N_3$  (see Figure 2). Next, it chops table  $B$  (on primary node  $N_1$ ) into chunks, each containing multiple tuples (currently set to 500), sends each chunk to a secondary node and assigns to a thread on that node. Figure 2 shows that a chunk  $B_1$  of table  $B$  consisting of tuples  $u_1, u_2, u_3$  is sent to a thread on node  $N_2$ , and that another chunk  $B_2$  consisting of tuples  $u_4, u_5, u_6$  is sent to a thread on node  $N_3$ .

Find top- $k$  tuples in table  $A$  for each tuple of table  $B$ : Each thread now goes through the tuples in the assigned chunk. For each tuple, it probes index  $I$  to find the top  $k$  tuples in table  $A$  with the highest tf/idf scores, pair these tuples with the probing tuple, then sends the pairs back to the primary node  $N_1$ . (The thread only sends back the IDs, not the full tuples.)

Consider again the thread for chunk  $B_1$  with tuples  $u_1, u_2, u_3$  (under "Node  $N_2$ " in Figure 2). Suppose  $k = 2$ . This thread first processes tuple  $u_1$ : it probes index  $I$  to find the top 2 tuples in table  $A$  with the highest tf/idf scores with  $u_1$ . Suppose these are tuples  $v_1, v_2$ . Then the thread creates the pairs  $(u_1, v_1), (u_1, v_2)$  and send them back to node  $N_1$  (see the figure). Next, the thread processes tuple  $u_2$ , then tuple  $u_3$ . Similarly, Figure 2 shows how a thread on node  $N_3$  processes chunk  $B_2 = \{u_4, u_5, u_6\}$ .

When Sparkly has processed all chunks of table  $B$ , and all pairs sent back from the secondary nodes have been collected into a table  $C$  on primary node  $N_1$ , Sparkly terminates, returning  $C$  as the blocking output.

The tf/idf scoring function: All that is left is to describe the scoring function used by Sparkly. First, we ask that the user manually identify a set of attributes to block on. Typically these are "identity"

attributes, such as name, phone, address, product title, brand, etc. Next, for each tuple (in table  $A$  or  $B$ ), we concatenate the values of these attributes into a single string  $s$ , lowercase all characters in  $s$ , tokenize  $s$  into a bag of 3-gram tokens, and remove all non-alphanumeric tokens.

Let  $B_t$  be the bag of 3-grams for a tuple  $t$ . When indexing table  $A$ , for each tuple  $t \in A$ , we index only  $B_t$ , not the entire tuple  $t$ . Finally, when querying, we compute the BM25 score between two tuples  $u, v$  to be the BM25 score between  $B_u$  and  $B_v$ .

**Discussion:** We now discuss the rationales behind the main design decisions of Sparkly.

Use top- $k$  instead of thresholding: This is the most important decision that we made. We use top- $k$  instead of thresholding for the following reasons. First, it is often easier to select a value for  $k$  than a threshold  $\alpha$ . Given a value for  $k$ , we know precisely how big the blocking output will be, and the rule of thumb is to select  $k$  that produces the largest blocking output that the matching step can handle, because the larger the blocking output, the higher the recall. On the other hand, we often do not have any guidances on how to select a good threshold  $\alpha$ .

Second, we observe that real-world data is often so noisy that the similarity scores of many matching tuple pairs can be quite low (see Section 5). This makes it very difficult to set threshold  $\alpha$ . A high threshold kills off many matches, producing low recall. A low threshold often blows up the blocking output size in unpredictable ways. In contrast, in such cases we observe that the matching tuples are often still within the top- $k$  "distance" of each other, making top- $k$  retrieval still effective, as we show in Section 4.

Finally, Lucene and many other KWS systems are highly optimized runtime-wise for top- $k$  search, but not for threshold search.

Do top- $k$  on just one side instead of both sides: Currently we do top- $k$  only from table  $B$  into table  $A$ . Another option is to do top- $k$  on both sides: from  $B$  into  $A$  and from  $A$  into  $B$ , then return the union of the two outputs. We experimented with this option but found that it can significantly increase runtime yet improve recall only minimally. It also complicates coding (e.g., we have to write code to remove duplicate pairs from the outputs of both sides).

Do top- $k$  from the larger table: We index the smaller table, say table  $A$ , then do top- $k$  probing from the larger table  $B$  because indexing the smaller table takes less time and produces a smaller index  $I$ . Shipping this smaller index  $I$  to the secondary Spark nodes takes less time. Finally, probing from the larger table rather than the smaller one tends to produce higher recall, given the same  $k$  value.

Ship the index and tuples of table  $B$  to the secondary nodes: This is the second most important decision that we made. The challenge here is to find an efficient way to do distributed top- $k$  probing on a Spark cluster. Toward this goal, recall that we create the inverted index  $I$  for table  $A$  on the primary node  $N_1$ . Table  $B$  also resides on  $N_1$ . So the simplest solution is to do all top- $k$  probeings there, using only the cores of  $N_1$ . However,  $N_1$  has a limited number of cores (e.g., 16, 32), so it can run only a limited number of threads, severely limiting how much top- $k$  probing we can do in parallel.

The next solution is to send the tuples of  $B$  to the secondary Spark nodes, then do top- $k$  probing from the secondary nodes into the index  $I$  on primary node  $N_1$ . This way, the secondary nodes can

run a much larger number of threads. Unfortunately, when these threads contact primary node  $N_1$  to do top-k probing, they would need to rely on the threads running on the cores of  $N_1$  to do the actual probing into index  $I$ . So once again, the limited number of threads on  $N_1$  becomes the bottleneck for scaling.

As a result, we decided to ship the index  $I$  and the tuples of  $B$  to the secondary nodes. Each secondary node then runs multiple threads, each doing top-k probing using the copy of  $I$  on that node. So we can do as many top-k probings in parallel as the number of threads on the secondary nodes. This produces a share-nothing parallel solution that is highly modular and can scale horizontally as we add more secondary Spark nodes.

Partitioning very large tables  $A$  and  $B$ : A major concern is whether shipping index  $I$  would take too long, because it can be very large. This turned out not to be the case. For example, in our experiments, indexing a table of 10M tuples produces indexes of size 1.3-2GB, and shipping these takes 21-32 seconds (see Section 4).

Still, one may ask what if the tables have 500M or 5B tuples? Would the indexes become too big to fit on the disks of Spark nodes? *Our solution is to break table  $A$  (the smaller table, to be indexed) into partitions of say 50M tuples, then process the partitions sequentially.* For example, if table  $A$  has 100M tuples, then we break  $A$  into partitions  $A_1$  and  $A_2$  each having 50M tuples, then run two blocking tasks:  $A_1$  vs.  $B$  and  $A_2$  vs.  $B$ . Finally, we combine the top-k results produced by these tasks. *This guarantees that Sparkly never has to build and ship indexes for more than 50M tuples.*

Use Lucene instead of ElasticSearch or Solr: We use Lucene because it provides highly effective procedures to index a table and do top-k probing, which are exactly what we need. ElasticSearch (ES) and Solr build on top of Lucene and provide a lot more capabilities that we do not need (e.g., sharding) yet can cause complications. For example, when we first built Sparkly, we used ES and observed two problems. First, it took much longer (and more pain) to install Sparkly, because we had to install ES as a part of the process. Second, Sparkly has less applicability, because we could not run it in certain environments, e.g., on a Kubernetes cluster, because we cannot ensure data locality (i.e., when Spark performs a top-k query, the query will go to an ES instance installed on the same node). So we switched to Lucene, which addresses the above problems.

### 3.4 Selecting Attributes and Tokenizers

So far we ask an expert user to *manually* select a set of attributes. Then we *concatenate* the values of these attributes into a string, tokenize it *using a default (3-gram) tokenizer*, then index and search on the tokenized string. We call this solution Sparkly Manual.

Sparkly Manual works well, but can suffer from three problems. First, it can be difficult even for expert users to select good blocking attributes. Second, concatenating the attributes is problematic because the importance of a token depends on which attribute it appears in. Finally, using a single tokenizer is also problematic because different attributes may best benefit from different tokenizers. To address these problems, we will automatically select blocking attributes and associated tokenizers, as elaborated below.

**Problem Definition:** First we formally define this selection problem. Let the attributes of tables  $A$  and  $B$  be  $F = \{f_1, \dots, f_n\}$ .

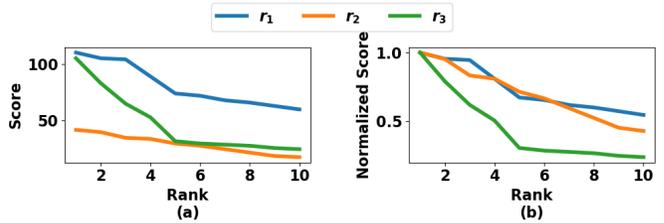


Figure 3: Illustrating the discriminativeness of configs.

Let  $T = \{t_1, \dots, t_m\}$  be a set of tokenizers (e.g., 3-gram, word-level). We define a *configuration*  $L$  (or *config*  $L$  for short) as a set of (attribute, tokenizer) pairs  $L = \{(f_{i1}, t_{i1}), \dots, (f_{ip}, t_{ip})\}$ , where  $f_{ij} \in F, t_{ij} \in T, j = 1 \dots p$ . Thus, in a config  $L$  different attributes can use different tokenizers.

Let  $\mathcal{L}$  be the set of all configs. Our goal is to find the config  $L \in \mathcal{L}$  that maximizes the recall. To define recall, we begin by defining the similarity score. Given two tuples  $b \in B, a \in A$ , we define their similarity score with respect to config  $L$  as the sum of the BM25 scores of the individual attributes in the config (tokenized using the assigned tokenizers). Formally, we have  $s(b, a, L) = \sum_{j=1}^p s_j[t_{ij}(b.f_{ij}), t_{ij}(a.f_{ij})]$ . Here  $t_{ij}(b.f_{ij})$  applies the tokenizer  $t_{ij}$  to the value of attribute  $f_{ij}$  of tuple  $b$ , producing a bag of tokens, and  $t_{ij}(a.f_{ij})$  produces another bag of tokens. Then  $s_j[t_{ij}(b.f_{ij}), t_{ij}(a.f_{ij})]$  computes the BM25 score between these two bags of tokens.

Next, we define the blocking output for when the above similarity score  $s(b, a, L)$  is used. For any tuple  $b \in B$ , let  $Q(b, A, k, L)$  be the list of top-k tuples from  $A$  that has the highest tf/idf scores, as defined by  $s(b, a, L)$ , with  $b$ . Let  $C(b, A, k, L)$  be the set of all pairs  $(b, v)$  where  $v \in Q(b, A, k, L)$ . Then the blocking output can be defined as  $C(B, A, k, L) = \cup_{b \in B} C(b, A, k, L)$ .

Finally, let  $recall(C(B, A, k, L))$  be the fraction of true matches in  $C(B, A, k, L)$ . Then our problem is to find a config  $L \in \mathcal{L}$  that maximizes  $recall(C(B, A, k, L))$  for a given  $k$ .

The above problem raises two challenges: how to estimate the recall of a config and how to find the config with the highest recall in a large space of configs. We now address these challenges.

**Estimating the Recall of a Config:** Given a config  $L$ , it is not possible to estimate its recall  $recall(C(B, A, k, L))$  because we do not know the true matches. To address this problem, we make the key observation that *it is possible to estimate the discriminative power of a config  $L$ , which captures its ability to tell apart the matches from the non-matches.* We can then search for the config with the maximal discriminativeness, on the heuristic assumption that this config is likely to achieve high recall.

*Example 3.1.* To motivate, consider a tuple  $b \in B$  and three singleton configs  $L_1, L_2, L_3$  involving attributes  $f_1, f_2, f_3$ , respectively. Let  $r_1, r_2, r_3$  be the top-k lists for  $b$ , produced by querying the inverted index  $I$  of table  $A$  using the above 3 configs, respectively. Figure 3.a shows the top-k lists  $r_1, r_2, r_3$ . Note that each top-k list contains tuple IDs in  $A$ , already sorted in decreasing BM25 scores. For each list, the figure shows the scores, plotted against the ranks of where they appear in the list.

Figure 3.a suggests that for the above tuple  $b \in B$ ,  $r_3$  is quite “discriminative”, because it “slopes down” steeply (i.e., the top few tuples of  $r_1$  have very high scores while the rest of the tuples have much lower scores). In fact, the curve  $r_3$  appears more discriminative than  $r_1$  and  $r_2$ , which do not “slope down” as much.

It may appear that we can measure this discriminativeness as the area under the curve (AUC): smaller AUC means higher discriminativeness. In Figure 3.a, this is indeed true for  $r_3$  and  $r_1$ :  $AUC(r_3) < AUC(r_1)$  and  $r_3$  is more discriminative than  $r_1$ . But it is not true for  $r_3$  and  $r_2$ , because  $AUC(r_2) < AUC(r_3)$ , yet  $r_2$  is not more discriminative than  $r_3$ . The problem is that the BM25 scores of the curves (generated by using different configs) are not comparable, and hence the AUCs are also not comparable. To address this, we normalize the BM25 scores of each curve to be between  $[0,1]$  (by dividing the original scores in each curve by the maximum score). Figure 3.b shows the normalized curves. Now it is indeed the case that smaller AUC means higher discriminativeness.

Thus, we can define the discriminativeness of a config  $L$  for a table  $B$  (given a table  $A$  and an inverted index  $I$ ) as the average discriminativeness of config  $L$  for each tuple in  $B$ :  $meanAUC(B, L, k) = \frac{1}{|B|} \sum_{b \in B} AUC(b, L, k)$ .

In turn, we can define the discriminativeness of config  $L$  for a tuple  $b$  in  $B$  as the normalized AUC. Let  $r(b, L, k) = ((v_1, s_1), \dots, (v_{k'}, s_{k'}))$  be the top- $k$  tuple list retrieved from index  $I$ , for record  $b \in B$ , scored according to config  $L$ , sorted in decreasing order of score  $s_1, \dots, s_{k'}$  ( $k' \leq k$  because only tuples with positive score can be in the list). Then we can compute the area under the curve as  $AUC(b, L, k) = \frac{1}{k' \cdot s_1} \sum_{i=1}^{k'-1} s_{i+1} + \frac{s_i - s_{i+1}}{2}$ . The tech report [34] explains how we arrive at this formula.

In practice, computing  $meanAUC(B, L, k)$  for a config  $L$  is too expensive, as we have to query index  $I$  with all tuples in  $B$ . So we approximate it using  $meanAUC(B', L, k)$ , where  $B'$  is a random sample of 10K tuples of  $B$  (and we set  $k$  to 250).

**Searching for a Good Config:** Our goal now is to find the config  $L$  that maximizes  $meanAUC(B', L, k)$ . The number of configs can be huge (e.g., in the millions). So we adopt a greedy search approach. First, we score all singleton configs (each using a single attribute/tokenizer pair) and find the top 10 configs with the lowest  $meanAUC$  scores. Next, we combine these configs to create “composite” configs, where each config has up to 3 attributes. We do not consider configs of more than 3 attributes because in our experience these configs take much longer to run yet only minimally improve recall, if at all. Finally, we score all configs and return the one with the lowest  $meanAUC$  score.

Since we use at most 10 singleton configs to create more configs of size up to 3 attributes, the total number of configs to score is at most 175, making exhaustive scoring of all configs possible.

We further speed up the above search using a technique called *early pruning*. To illustrate, consider again the problem of scoring all singleton configs to find the top 10 configs. Scoring a config means querying the inverted index  $I$  with all tuples  $b \in B'$ . Even though  $B'$  is small (currently set to 10K), this still takes time. So we score the configs using a sample  $B''$  which is a small subset of  $B'$ , use a statistical test to remove all configs for which we can say with high confidence that they will not make it into the top 10, then

Table 1: Datasets for our experiments.

Type	Dataset	Table A	Table B	#Matches	#Attr
Structured	Amazon-Google <sub>1</sub>	1,363	3,226	1,300	4
	Walmart-Amazon <sub>1</sub>	2,554	22,074	1,154	6
	DBLP-Google <sub>1</sub>	2,616	64,263	5,347	4
	DBLP-ACM <sub>1</sub>	2,616	2,294	2,224	4
	Hospital <sub>1</sub>	1,786	1,786	3,949	7
Textual	Songs-Songs <sub>1</sub>	1,000,000	1,000,000	1,292,023	5
	Amazon-Google <sub>2</sub>	1363	3,226	1,300	2
	Walmart-Amazon <sub>2</sub>	2,554	22,074	1,154	2
Dirty	Abt-Buy	1,081	1,092	1,097	3
	Amazon-Google <sub>3</sub>	1,363	3,226	1,300	4
	Walmart-Amazon <sub>3</sub>	2,554	22,074	1,154	6
	DBLP-Google <sub>2</sub>	2,616	64,263	5,347	4
	DBLP-ACM <sub>2</sub>	2,616	2,294	2,224	4
	Hospital <sub>2</sub>	1,786	1,786	3,949	7
	Songs-Songs <sub>2</sub>	1,000,000	1,000,000	1,292,023	5

expand  $B''$  with more tuples, re-score the remaining configs, and so on. Specifically: (1) Initialize the subsample  $B'' = \emptyset$  and  $S$  to be the set of all configs from which we have to compute the top 10 configs. (2) Expand the subsample  $B''$  by adding to it a small random sample of  $h$  tuples from  $B' \setminus B''$ . (3) Compute the  $meanAUC$  for all configs in  $S$  using  $B''$  and finding the set  $\hat{R}$  of the top-10 configs. (4) For each config  $L \in S \setminus \hat{R}$ , use the Wilcoxon signed-rank test [39] to determine (with high confidence) if its  $meanAUC$  score is greater than those of the configs in  $\hat{R}$ . If yes, then  $L$  is unlikely to ever be in the top 10. Remove  $L$  from  $S$ . (5) If  $S = \hat{R}$  or  $B'' = B'$ , return  $\hat{R}$  as the top-10 configs, otherwise go back to Step 2.

We also use the above early pruning procedure to search the space of the larger configs. We defer further details to the tech report [34].

## 4 EMPIRICAL EVALUATION

**Datasets:** We use 15 datasets described in Table 1, which come from diverse domains and sizes, and have been extensively used in recent EM work [23, 25, 29, 38] (except Hospital, which is private). Structured datasets have short atomic attributes such as name, age, city. Textual datasets have only 2-3 attributes that are textual blobs (e.g., title, description). For dirty EM, we focus on one type of dirtiness, which is widespread in practice [25] mainly due to information extraction glitches, where attribute values are “moved” into other attributes. Textual and dirty datasets are derived from the corresponding structured datasets (e.g., the textual dataset Amazon-Google<sub>2</sub> is derived from the structured dataset Amazon-Google<sub>1</sub>).

Later we use 6 additional datasets for certain experiments, as discussed in Section 4.5 and Section 5.

**Methods:** We compare Sparkly to 8 state-of-the-art (SOTA) EM blockers.

*Autoencoder, Hybrid, Union(DL,RBB):* A recent work [38] shows that deep learning (DL) based blockers significantly outperform many other blockers. So we compare Sparkly to the two best DL blockers: Autoencoder and Hybrid [38]. The work [38] also shows that combining the best DL blocker and RBB, a SOTA industrial blocker, produces even better recall at a minimal increase of blocking output size. As a result, we also compare Sparkly with that blocker, henceforth called Union(DL,RBB).

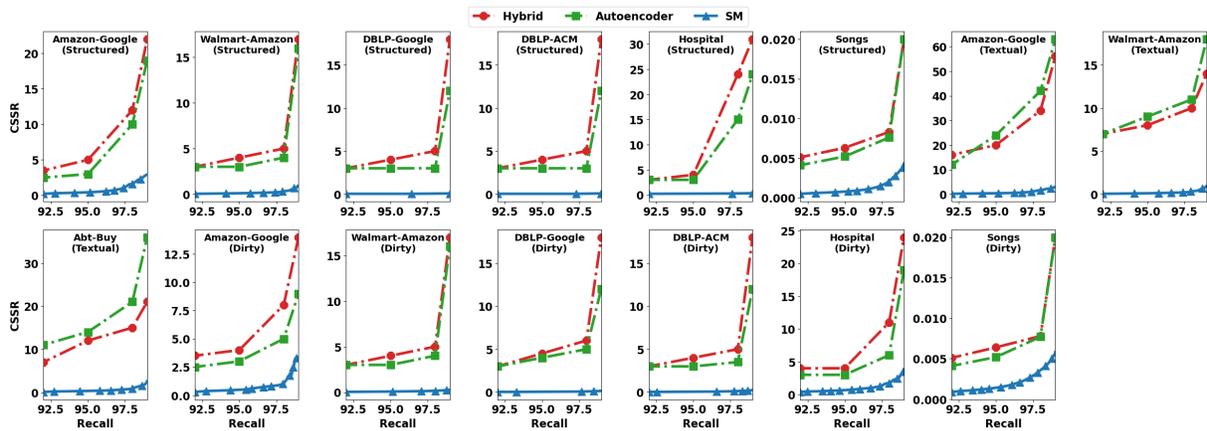


Figure 4: SM vs. the two best DL methods in terms of recall and blocking output size.

*PBW, DBW, JD*: Hash blockers have been very popular, and recent EM work has developed highly effective hash blockers, as captured in the pioneering JedAI open-source EM platform [31, 33]. These blockers hash each tuple to be matched into multiple blocks, one per each unique token in the tuple, then employ sophisticated methods to remove/clean blocks, among others. Based on personal communications with the JedAI authors, we compare Sparkly to 3 SOTA blockers in JedAI: PBW, DBW, and JD (we describe these methods in the technical report [34]).

*kNN-cosine, kNN-jaccard*: Finally, a recent work [29] shows that a kNN blocker outperforms many other blockers. This blocker finds all tuple pairs where a tuple is among the  $k$  nearest, i.e., most similar, neighbors of the other tuple, where the similarity measure is cosine over 5-gram tokenization. So we compare Sparkly with this blocker, denoted kNN-cosine. We also compare Sparkly to kNN-blockers where the similarity measure is cosine over 3-gram tokenization and Jaccard over 3-gram and 5-gram tokenization.

#### 4.1 Recall and Output Size

We begin by comparing Sparkly to existing methods in terms of recall and output size. To keep the comparison manageable, we first compare SM, the Sparkly version where the user manually selects the attributes to be blocked on (i.e., Sparkly Manual), with all SOTA blockers. Then we compare SM with SA, the Sparkly version that automatically selects blocking attributes (i.e., Sparkly Auto).

**Comparing SM to DL Methods:** Figure 4 compares SM with the two best DL blockers, Autoencoder and Hybrid [38]. The figure shows 15 plots, one per dataset. Consider the first plot, which is for the structured Amazon-Google dataset. Here the x-axis shows *the recall*  $R = |C \cap G|/|G|$ , where  $C$  is the blocking output and  $G$  is the set of all gold matches. The y-axis shows *the candidate set size ratio*  $CSSR = |C|/|A \times B|$ . Both axes show values *in percentage*. So a value of 85 on the x-axis means recall of 85%, and a value of 5 on the y-axis means CSSR of 5%. Like SM, the two DL blockers Autoencoder and Hybrid are also top-k. So we vary the value of  $k$  to generate the above plot. We generate the remaining 14 plots in a similar way. Note that the y-axes of the 15 plots vary significantly

in scale. This is necessary so that we can show the difference among the curves.

All 15 plots show that SM significantly outperforms the two DL blockers: for each recall value, SM achieves a much lower CSSR, and this gap widens dramatically as recall approaches 100%. For example, on the first plot, at recall of 98%, Sparkly-Man achieves CSSR of 2.5%, whereas the two DL blockers achieve CSSR of 10%. These gaps are bigger for textual datasets, suggesting that SM can better handle textual data than the DL blockers. The gaps are smaller but still quite significant on all dirty datasets.

The above two DL methods concatenate all attributes and then block on the concatenation. In the next experiment, we modified them to block on *the concatenation of only those attributes that SM blocks on*. Even in this case, SM still outperforms both DL methods on 14 datasets (sometimes by very large margins) and is comparable on 1 dataset (see the tech report for details).

**Comparing SM to Other Methods:** Next we compare SM with Union(DL,RBB), which combines the best DL blocker and RBB (a SOTA industrial blocker), and the three JedAI methods: PBW, DBW, and JD. It is very difficult to vary the parameters of these methods in such a way that generates meaningful recall-CSSR curves, because they do not have a top-k parameter that we can adjust. So we compare them with SM at  $k = 10, 20, 50$ , as shown in Table 2.

This table shows that SM is very *predictable*: it achieves high recall for all datasets (92.5-100% for  $k = 10$ , 96.4-100% for  $k = 20$ , 98.7-100% for  $k = 50$ ), and its output size is capped as  $k * |B|$ . In contrast, the remaining four methods are *unpredictable*. For example, PBW’s recall can be perfect (100%) but also can be as low as 74.5%, and its output size can be small but can also be as high as 4.2 billions for the structured dataset Songs. (We report no results for “S - D” because PBW was out of memory on this dataset, on a machine with more than 100G of RAM). Similarly, DBW’s recall can be as low as 84.7% and output size as high as 454.5M.

JD produces much more reasonable output size across all datasets, but at the cost of lower recall 35.4-96.4%. Similar to JD, Union(DL,RBB) also produces reasonable output sizes (larger than those of JD), but varying recalls 83-99.9%.

Table 2: SM vs. the three JedAI methods and Union(DL,RBB) in terms of recall and blocking output size.

Dataset	PBW		DBW		JD		Union (DL,RBB)		Sparkly K=10		Sparkly K=20		Sparkly K=50	
	C	Recall	C	Recall	C	Recall	C	Recall	C	Recall	C	Recall	C	Recall
AG - S	24.5k	92.1	15.9k	89.2	5.9k	80.5	77.7k	98.8	33.3k	96.8	66.5k	97.8	165.9k	99.2
WA - S	1.5m	99.7	159.8k	93.8	88.3k	95.0	2.1m	98.9	220.7k	98.4	441.4k	99.0	1.1m	99.5
DG - S	430.5k	91.0	779.3k	99.6	53.1k	79.7	7.6m	99.6	641.1k	99.9	1.3m	100.0	3.2m	100.0
DA - S	8.1k	83.7	35.1k	99.9	2.3k	80.3	198.4k	99.9	22.9k	99.8	45.9k	100.0	114.7k	100.0
H - S	11.9k	100.0	4.0k	84.7	1.4k	35.4	209.8k	99.9	17.8k	100.0	35.4k	100.0	85.4k	100.0
S - S	4.2b	100.0	379.4m	99.8	2.5m	82.0	50m	98.7	10.0m	96.3	20.0m	97.9	50.0m	99.3
AG - T	24.5k	92.1	15.9k	89.2	5.9k	80.5	33.6k	85.0	33.3k	96.8	66.5k	97.8	165.9k	99.2
WA - T	1.5m	99.7	159.8k	93.8	88.3k	95.0	7.9m	83.0	220.7k	98.4	441.4k	99.0	1.1m	99.5
AB - T	4.7k	74.5	6.0k	88.6	1.2k	65.2	44.6k	95.7	10.9k	98.1	21.8k	98.9	54.5k	99.2
AG - D	38.8k	94.1	18.7k	91.3	6.4k	79.5	360.0k	99.3	33.3k	96.6	66.5k	98.2	166.0k	99.0
WA - D	1.1m	99.5	225.2k	97.4	88.1k	95.9	935.9k	97.9	220.7k	99.1	441.5k	99.7	1.1m	99.8
DG - D	4.0m	99.7	925.5k	98.8	180.5k	96.4	47.6m	99.8	642.2k	99.9	1.3m	100.0	3.2m	100.0
DA - D	12.5k	86.6	42.0k	97.2	4.7k	82.4	1.0m	99.8	22.9k	99.3	45.9k	99.8	114.7k	100.0
H - D	22.5k	100.0	31.2k	87.9	2.4k	56.1	136.8k	98.5	17.9k	94.0	35.6k	97.1	88.4k	98.7
S - D	—	—	454.5m	96.2	3.1m	68.3	50m	95.2	10.0m	92.5	20.0m	96.4	50.0m	98.8

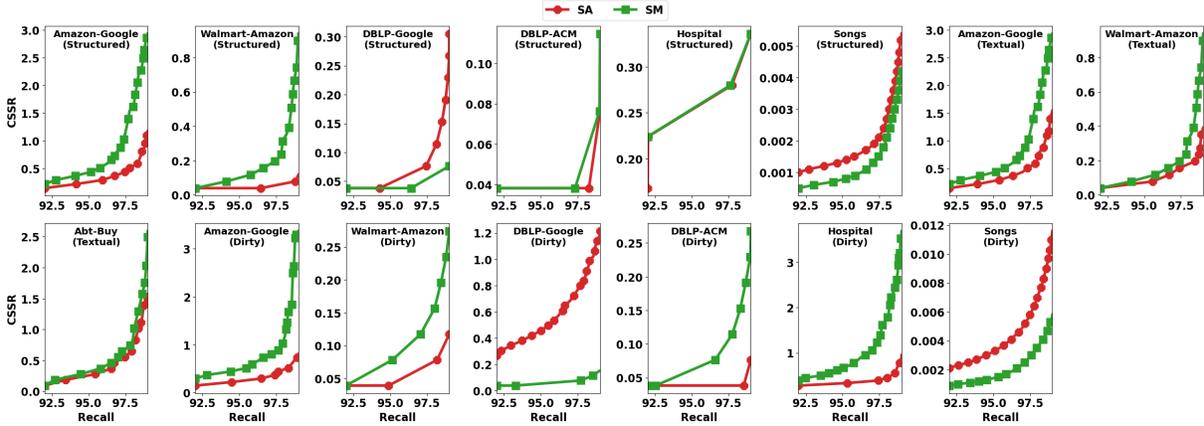


Figure 5: SM vs. SA in terms of recall and output size.

Finally, we compare SM with the kNN methods. The tech report [34] provides the details, including plots. Here we only summarize the findings. As mentioned earlier, a recent paper [29] finds that kNN-cosine using 5-gram tokenization outperforms many SOTA blockers. We find that on 10 datasets SM outperforms this method, sometimes by huge margins. On 1 dataset SM is comparable, and on the remaining 4 datasets SM is worse than kNN-cosine-5gram, but the performance gap is very small. We also find that kNN-cosine using 3grams is comparable to kNN-cosine using 5grams, and both outperform kNN using Jaccard (either 3grams or 5grams).

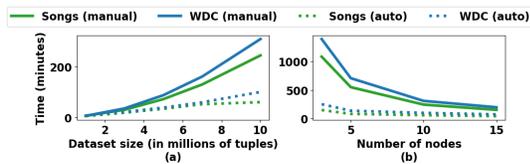
**Comparing SM to SA:** Figure 5 shows that SA outperforms SM on 10 datasets in terms of recall and output size, sometimes by a large margin. SA is worse than SM on the remaining 5 datasets, but the performance gap is very small. For example, at 98% recall, SA’s output size is at most 0.7% larger than that of SM. We discuss possible reasons for SA being worse than SM on several datasets in the tech report.

## 4.2 Runtime

We now examine the runtime of Sparkly. We ran all experiments on an AWS cluster of 10 nodes. Each node is an m5.4xlarge instance with 16 cores, 64G RAM, costing \$0.75/hour (as of July 2022).

Figure 6.a shows the runtime of SM (the solid lines) and SA (the dotted lines), as we vary the size of two datasets: Songs and WDC. Songs has 1M tuples. WDC is a large dataset of 26M tuples used by recent work [35] (we cannot use WDC for recall experiments because it does not have all gold matches). A point  $n$  on the x-axis reports the runtimes measured on a sample of  $n$  millions tuples randomly obtained from WDC, and on a sample of  $n$  millions tuples obtained by replicating Songs  $n$  times.

The above figure shows that both SM and SA scale (slightly superlinearly) as we increase the dataset size, and that SA is much faster than SM. This is because SA scores the attributes individually, instead of scoring their concatenation as SM. Concatenation produces long attributes, and performing top-k search on long attributes takes more time than on short attributes. Further, SA can



**Figure 6: Runtime for (a) varying dataset sizes, and (b) varying cluster sizes using datasets of 10M tuples each.**

use word tokenizers on some attributes, whereas SM only uses 3gram tokenizers. This produces fewer tokens, which often leads to faster top-k search.

It is noteworthy that SA can block datasets of size 10M under 100 minutes, incurring an AWS cost of only \$12.5.

Figure 6.b shows the runtime of Sparkly on the 10M WDC and 10M Songs datasets, as we vary the size of the AWS cluster. As the number of nodes goes from 3 to 15, runtime decreases significantly, as expected. As we increase the cluster size, eventually the overhead time (indexing, shipping, attribute/tokenizer selection, etc.) will dominate (compared to the top-k probing time).

### 4.3 Performance of Sparkly’s Components

We find the indexing time to be minimal. For example, on the same AWS cluster described above, indexing Songs at size 5M and 10M takes 76 and 115 seconds, respectively. The resulting index sizes are also reasonable. For Songs and WDC at size 1M, 5M, and 10M, the sizes are 137, 664, 1318 MB, and 214, 1034, and 2042 MB, respectively. Shipping these indexes to the Spark nodes takes minimal time. For Songs and WDC at 1M, 5M, and 10M, shipping the indexes takes 2.2, 7.2, 21 seconds, and 2.5, 11, 32 seconds, respectively.

Finally, recall that SA performs a search for a good set of attributes/tokenizers (to block on). On Songs and WDC 1M, 5M, and 10M, *without early pruning*, this search takes 4, 9.2, 15.6 mins and 4.6, 10.1, 17.2 mins, respectively. Early pruning cuts these times by up to 70%, to 1.2, 3.2, 6 mins, and 2, 5.3, 14 mins, respectively. The greedy method used by the searcher was quite effective. We performed exhaustive search on 11 datasets to find the optimal configs, and found that the greedy method found a config with a score within 0-0.8% of the optimal score on 10 datasets and within 10% on 1 dataset.

### 4.4 Sensitivity Analysis

We now vary the parameters of the major components to examine the sensitivity of Sparkly. We only summarize the findings here, deferring a detailed discussion to the TR.

**Blocking Attributes:** Recall that in SM, the user manually selects a set of attributes  $S$  to block on. We find that varying this set of attributes does impact the performance of SM, minimally by 0-1.5% CSSR in 11 datasets, and moderately by 2-5% CSSR in 4 datasets. This suggests that while manually selecting blocking attributes is a reasonable strategy, there is still room to improve, e.g., by automatically finding such attributes, as done in SA.

**Tokenizers:** Recall that SM uses a 3gram tokenizer. Next we examine replacing this tokenizer with a 2gram, 4gram, and word-based tokenizer, respectively. We find that changing the tokenizer can

significantly impact the performance (e.g., by up to 11.5% CSSR). Overall, the 3gram tokenizer (used by SM) is a good choice as it has reasonable performance on most datasets. The 2gram and 4gram tokenizers perform worst, with the 2gram tokenizer also incurring the longest runtime.

**BM25’s Parameters:** BM25 has two parameters:  $k_1$  (default value 1.2) and  $b$  (default value 0.75), to handle term saturation and document length. Varying  $k_1$  from 1 to 2 does not significantly change SM’s performance. This may be because in our setting blocking attributes do not have many terms, and they do not have high term frequency, so term saturation is not a major issue. Varying  $b$  from 0.5 to 1 changes SM’s performance more, by up to 2% CSSR. But we also find that  $b = 0.75$  provides a good default value for SM, as its curve is either the best curve or very close to the best curve on most datasets.

**Config Searcher’s Parameters:** SA uses a searcher to find a good blocking config. This searcher has four major parameters: (1) the size of  $B'$ , a sample of table  $B$  on which to score the configs (set to 10K), (2) the number of tuples returned in each querying  $k = 250$ , (3) the number of initial configs selected (set to 10), and (4) the max number of attributes considered in a config (set to 3).

Varying (1) from 5K to 15K changes SA minimally. Varying (2) from 200 to 300 again changes SA minimally (only up to 0.2% CSSR on 1 dataset). Similarly, varying (3) from 8 to 12 and varying (4) from 2 to 4 show minimal changes. In all cases, the default values for (1)-(4) provide a good curve, which is either the best or very near the best.

### 4.5 Additional Experiments

We now examine how Sparkly performs on very large datasets, how it compares runtime-wise to DL methods, and whether DL methods can achieve higher accuracy, given larger datasets (to train on).

It is very difficult to find very large *public* datasets with *complete gold*, i.e., all true matches (without which we cannot compute the blocking recall). After an extensive search, we settle on three datasets: BC, MB, and WDC. BC (Big Citations) blocks two tables of 2.5M and 1.8M paper citations. MB (Music Brainz) blocks a table of 20M songs (against itself), and WDC blocks a table of 26M product descriptions [35]. BC and MB have complete gold, but WDC does not (see the tech report).

Table 3 shows the results. First, we deployed an AWS cluster of 30 m5.4xlarge nodes (16 cores, 64G RAM, \$0.75/hour, per node), then ran Sparkly on all three datasets (see the first three rows of the table). Each row lists the results of SM and SA separated by “/”. Column “Time” shows the total time in minutes, while the next three columns show the recall at  $k = 10, 25, 50$ . We cannot compute recall for WDC as it does not have the complete gold.

The first three rows show that *Sparkly scales to very large datasets, and that SA is much faster than SM*, taking only 130 and 168 mins to block WDC 26M and MB 20M, respectively, at a reasonable cost of less than \$67.5 on AWS. Sparkly achieves high recall on MB and BC at  $k = 50$ .

Skipping the 4th row of Table 3 (which we discuss later), we now consider the DL method Autoencoder. Unfortunately we had

**Table 3: Sparkly and DL methods on large datasets.**

Method	Dataset	Time	Recall @ 10	Recall @ 25	Recall @ 50
Sparkly	WDC 26M	603/130	-	-	-
	MB 20M	449/168	79/95	87/97	91/98
	BC 2.5M	44/11	99/79	100/89	100/94
	MB 10M	132/61	85/96	91/98	94/98
Autoencoder	WDC 10M	925	-	-	-
	MB 10M	691	30	35	40
	BC 2.5M	146	81	84	85
Hybrid	BC 2.5M	2719	73	76	78

tremendous difficulties scaling Autoencoder to large datasets. Autoencoder is a prototype code used in the paper [38], for datasets of up to 1M tuples. It runs on a single GPU and uses many Python libraries that are not well suited to large datasets (e.g., the SVD implementation of Sklearn). So when applied to large datasets, Autoencoder quickly exhausts memory and crashes, and there is no easy way to modify it to run in a distributed setting (where it can use a lot more GPU memory).

After extensive optimization efforts, we managed to apply Autoencoder to BC 2.5M, WDC 10M, and MB 10M, on a SOTA hardware available to us (32t/16c CPU with 64G RAM coupled with RTX 2080ti GPU with 11G RAM). Table 3 shows the results. While it is not entirely fair to compare the runtimes of Autoencoder and Sparkly, because they run on *different* hardware, it is still interesting to note that Autoencoder takes much more time than Sparkly, e.g., 691 vs 132/61 mins (for SM/SA) on MB 10M, and 146 vs 44/11 mins on BC 2.5M. Autoencoder spent most time in preprocessing and self-supervised training.

Hybrid is far more complex than Autoencoder, and we only managed to run it on BC 2.5M (it ran out of memory on WDC 5M and MB 5M). Even on BC, its runtime is already very high (2719 mins). This suggests that *existing prototype DL blockers do not scale to large datasets, requiring a lot more future work on this topic.*

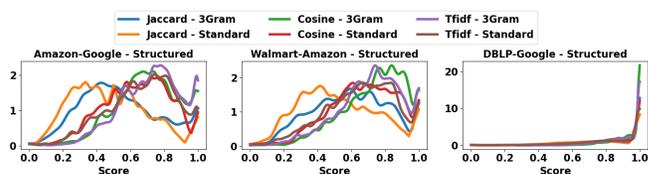
Both Autoencoder and Hybrid achieve far lower recall at  $k = 50$  than Sparkly (see the rows for BC 2.5M and MB 10M), suggesting that *these methods still cannot exploit larger datasets to achieve higher accuracy than Sparkly.* In the tech report we discuss how the remaining SOTA methods also do not scale to these large datasets.

## 5 DISCUSSION & FUTURE WORK

We now discuss questions that may arise in light of Sparkly’s strong blocking performance.

**Why Little Attention So Far?** TF/IDF has long been studied in the matching step of EM [8]. It is not clear why it has not been studied in the blocking step. A possible reason is that so far EM researchers have preferred to focus on hash-based blocking, which is conceptually simple and easy to scale [5, 27, 32]. Even when considering similarity-based blocking, researchers have preferred similarity measures that seem simpler and more amenable to scaling, e.g., edit distance, Jaccard, overlap, cosine [5, 27, 32]. TF/IDF appears difficult to scale. A straightforward application of inverted indexes is slow, and it was not obvious how to do much better.

Up until 2015, Lucene was also slow, making tf/idf blocking using Lucene impractical. Then it adopted the block-max WAND indexing technique and became much faster [18]. As this paper has



**Figure 7: Distributions of the scores of gold matches.**

shown, a combination of the “new” Lucene and Spark has now made tf/idf blocking very competitive, and suggests that going forward it should receive more attention.

**Top-k vs Thresholding:** We believe doing top-k search is critical. To see this, we compute the similarity scores of *all gold matches* on each dataset, for tf/idf, Jaccard, and cosine measures. Figure 7 shows the smoothed histograms of these scores, on the first 3 datasets (see the TR for similar results on the remaining 12 datasets). Only on the 3rd dataset would most of these scores be in a narrow range near 1.0 (e.g., between 0.8 and 1), making thresholding work well, i.e., achieving high recall. In the first 2 datasets, these scores are spread all over the range [0,1]. This suggests that thresholding cannot achieve high recall, unless we set the threshold very low, which blows up the blocking output size.

In contrast, Section 4 shows that doing top-k (as in Sparkly) achieves high recall without blowing up the output size. Thus, we believe that *future blocking solutions should seriously consider doing top-k, rather than thresholding.*

**Do We Need Both TF and IDF?** But top-k alone is not enough. kNN-cosine and kNN-jaccard also use top-k, yet underperform Sparkly, which uses tf/idf. So it seems we need tf/idf too. But do we need both tf and idf for blocking? To answer this question, we performed several experiments. We summarize the findings here (deferring the details to the TR). When we remove idf from SM (by dropping “idf(t)” in the BM25 formula of Equation 2 in Section 3), we find that SM greatly outperforms SM-no-idf.

Similarly, when we remove idf from TFIDF-cosine, by dropping “idf(t)” from the formula  $V_D(t) = tf(t, D) \cdot idf(t)$  (see Section 3, recall that TFIDF-cosine is the well-known tf/idf measure that computes  $s(D, Q) = [\sum_t V_D(t) \cdot V_Q(t)] / [\sqrt{\sum_t V_D(t)^2} \cdot \sqrt{\sum_t V_Q(t)^2}]$ ), we find that TFIDF-cosine greatly outperforms TFIDF-cosine-no-idf on many datasets. This suggests that *idf is important for blocking.*

Interestingly, when we performed a similar experiment where we removed tf, we did not see a clear trend. TFIDF-cosine and TFIDF-cosine-no-tf perform largely the same. SM is minimally better than SM-no-tf on some datasets, minimally worse on some others, and about the same on the remaining datasets.

So tf seems to have minimal effect on the 15 datasets. We believe this is because the attributes to block on (e.g., product title, person name) tend to be short, where few tokens repeat multiple times. To test this hypothesis, we consider Companies, a highly textual dataset where each tuple is a long document describing a company, and we need to block using the entire tuple (this dataset was used in the paper [25] to evaluate DL methods for matching). Indeed on this dataset where many tokens repeat multiple times, SM greatly outperforms SM-no-tf, e.g., achieving 62% vs 33% recall at  $k = 50$ . Similarly, TFIDF-cosine greatly outperforms TFIDF-cosine-no-tf.

This result suggests that *tf's effect on short blocking attributes is minimal, but can be significant on long textual blocking attributes.*

**Which Scoring Function Works Best?** The above results suggest using both *tf* and *idf*. But many scoring functions (which combine *tf* and *idf*) exist. TFIDF-cosine and BM25 are just two functions among them. So which scoring function works best? To explore, we examine four scoring functions: TFIDF-cosine, TFIDF-jacc, SM, and SM+. TFIDF-cosine is the cosine similarity function using *tf/idf*, described earlier, while TFIDF-jacc is the function  $fms^{apx}$  described in [4], which can be viewed as the Jaccard similarity function using *idf*. SM uses BM25, and SM+ is an extension of SM that we will describe shortly.

Our experiments shows that TFIDF-jacc is somewhat worse than TFIDF-cosine (see the TR). But surprisingly, TFIDF-cosine is better than SM on many datasets. We found this is because TFIDF-cosine's scoring function incorporates the *tf* and *idf* of each term from both the query  $Q$  and the document  $D$  sides (see Equation 1 in Section 3), but the BM25 scoring function used by SM does not. It incorporates only *tf* and *idf* of each term from the document  $D$ 's side, not from the query  $Q$ 's side. In other words, *TFIDF-cosine treats Q and D uniformly, whereas the BM25 function of SM does not.* This makes sense in keyword search, where typically  $Q$  has few terms, each occurring only once and all terms in  $Q$  are important. But these are not true in EM, where  $Q$  is a tuple in Table  $B$  (and is often as long as  $D$ , which is a tuple in Table  $A$ ). Here, it makes more sense to treat  $Q$  and  $D$  uniformly, like TFIDF-cosine.

Using the above observation, we modify BM25 to incorporate *tf* and *idf* from  $Q$ 's side (see the TR), producing the SM+ solution. We find that SM+ performs very well, being either the best solution or very close to the best solution on all datasets. We modified SA similarly to produce the SA+ solution. We find that SA+ outperforms SA: better or equal to SA in 12 datasets, and worse in 3 datasets. In general, SA+ is either the best solution or very close to the best solution in 14 datasets (see the TR).

Another point worth noting is that BM25 (or the modified BM25 used in SM+ and SA+) enables very fast incremental index update (when tuples get added or removed), whereas TFIDF-cosine does not. All these results suggest that *using tf/idf weighting and the BM25 scoring function, but modifying it to incorporate tf and idf from the query side, is a promising future direction to explore.*

**Blocking Numeric Datasets:** Next we examine how well Sparkly blocks numeric datasets. It is very difficult to find numeric public datasets with complete gold. After an extensive search, we settle on two datasets AW and RE. Each dataset matches the columns of all tables within a data lake. As such, each dataset consists of a single table  $X$  where each tuple describes a column (listing column name, table name, the average length of the column's values, the average/min/max of the column's values if it is numeric, etc.). The goal then is to match  $X$  with itself. AW (AdventureWorks) and RE (Real Estate) have 799 and 451 tuples, respectively. See the TR for more details.

Here we found that SM is better than SA. On AW and RE, SM achieves 81% and 94% recall at  $k = 50$  compared to 79% and 67% for SA. This is because SA was confused by numeric attributes and so picked up some numeric attributes to block on. SM is comparable

to kNN-jaccard and kNN-cosine, and is much better than the two DL methods and JedAI (see the TR).

SM however can still be improved. For example, a separate work on schema matching for data lakes (under preparation, from where we obtained the above two datasets) extended SM to incorporate rules such as "if the average values of two numeric columns are too far apart, e.g., one value is greater than 10 times the other value, then do not include the columns as a pair in the blocking output". This solution achieves recall 99% and 97% at  $k = 50$  for AW and RE. It turns out that we can naturally incorporate many such rules into the querying function of Lucene.

Thus, the results suggest that *Sparkly still performs better or comparable to SOTA methods on numeric data. Further, it can significantly be improved with rules exploiting the properties of numeric attributes, and many such rules can naturally be incorporated into Lucene.* Future work should explore this direction, and also improve SA to avoid picking up numeric attributes to block on.

**Scalability, Predictability, Extensibility, and Updates:** We believe Sparkly scales due to four reasons. First, *it decomposes blocking into executing a large number of independent tasks* (each querying the inverted index  $I$  using a chunk of tuples in table  $B$ ). Second, it executes these tasks on a Spark cluster in *a parallel share-nothing fashion*, by shipping the index  $I$  to the Spark worker nodes, then execute the tasks there. Third, it can execute each task fast, by *using the block-max WAND indexing technique* of Lucene. Finally, when the table  $A$  is too big (e.g., 200M tuples), it breaks  $A$  into smaller partitions (e.g., of 50M tuples each) and blocks each partition against table  $B$  (this minimizes the problem of having the indexes and other intermediate data structures grow uncontrolled as the table sizes increase, eventually crashing the cluster). As such, *Sparkly can maximally utilize the entire Spark cluster, and scale horizontally by adding more nodes.*

The above architecture is also *predictable*. Recall that we chop table  $B$  into chunks of tuples, and execute these chunks (i.e., use them to query index  $I$ ) on the worker nodes. After executing a few hundred chunks, it is possible to use the execution times of these chunks to estimate with high accuracy how much longer Sparkly will run.

The above architecture is also *extensible*, in that we can add other kinds of blocking. For example, consider hash blocking. We can create a hash  $H$  of table  $A$  and ship it to the worker nodes. Then given a tuple  $b \in B$ , we can consult the inverted index  $I$  to obtain a top- $k$  result, consult the hash  $H$  to obtain a result, union or intersect the two results, then send the output back to the driver node. In general, we can ship all kinds of indexes to each worker node, then do processing for each tuple  $b \in B$  using these indexes.

Thus, we believe that *future blocking solutions should seriously consider an architecture similar to Sparkly, which can provide significant benefits in scaling, predictability, and extensibility.*

Finally, Sparkly can naturally handle updates. An appealing property of BM25 is that it enables very fast incremental index updates, when we add or remove tuples from the indexed table. In contrast, TFIDF-cosine requires the entire index to be rebuilt from scratch. Fast index update can provide moderate to significant benefits in many cases, e.g., when matching two tables or matching a tuple in real-time into a table, as we discuss in the TR.

**Limitations of Sparkly:** We now analyze cases where Sparkly may achieve low recall. Consider a match  $g = (u \in B, v \in A)$ . There are three possible reasons why tuple  $v$  may not make it into the top- $k$  list of  $u$  (thus excluding  $g$  from the blocking output). First,  $v$  may have a low tf/idf score with  $u$ . This can happen due to numeric data (as we saw earlier), dirty data, missing values, synonyms, and natural variations (e.g., “Robert Smith” vs. “Ben Smith”). Second,  $v$  may have a high tf/idf score with  $u$ , but there are more than  $k$  true matches for  $u$ , so  $v$  is excluded.

Finally, many other non-matching tuples in  $A$  may have high tf/idf scores with  $u$ , crowding out  $v$ . For example, if we block just on the person name, then the top- $k$  list for a particular “David Smith” may contain tuples of many other “David Smith”-s, because David Smith is a very common name. As another example, the non-match (“iPhone 12 mint condition 32G *white* with case”, “iPhone 12 mint condition 32G *black* with case”) may have a high tf/idf score, because tf/idf fails to locate the colors and realize that if the colors do not match then the tuples do not match.

To address the above limitations, possible solutions include ways to clean and standardize the data, using a dynamic  $k$  value (e.g., if all tf/idf scores in the current top- $k$  list is high, then increase  $k$  then query again), and performing information extraction to isolate important attributes such as color.

**Future Directions:** Based on the above discussion, we propose the following research directions. First, we should study tf/idf blocking in more depth, improving Sparkly and similar tf/idf systems.

Second, we should develop much bigger benchmarks for blocking. The current datasets are too small (or some are large but have no gold matches, so we cannot evaluate blocking recall), and do not contain enough variety.

Third, we should develop effective methods to clean and standardize datasets, as this can significantly improve blocking recall and minimize the need to use sophisticated but costly blocking methods. Recent large language models (LLMs, a.k.a. foundation models, such as GPT-3, T5) may be promising here, as they can help clean and summarize the tuples.

Fourth, we should study how to improve existing blocking solutions and develop new ones, using Sparkly as a benchmark.

Finally, it is likely that an ideal blocking solution will have to use multiple blocking techniques, including tf/idf and others, to maximize recall. It must also scale, i.e., block tables of hundreds of millions of tuples in a reasonable time on a reasonable hardware at a reasonable cost. Toward this goal, the scalable share-nothing architecture of Sparkly provides a promising starting point.

## 6 ADDITIONAL RELATED WORK

EM has been a long-standing challenge in data management [1, 6, 7, 12, 14, 26, 30]. There has been multiple academic efforts on building scalable EM systems such as JedAI [31, 33], Magellan [19], and CloudMatcher [17].

Over the past decades, numerous blocking solutions have been developed. See [5, 27, 32] for surveys, and see Section 2 for a discussion of the main blocker categories. However, tf/idf blocking has received virtually no attention, as far as we can tell. The closest work that we have found is the recent work [28], which performs token blocking, i.e., hashing each tuple to multiple blocks, each corresponding to a token in the tuple. This work removes blocks

that correspond to tokens of low tf/idf values. The work [4] develops a scoring function that can be viewed as the Jaccard similarity function using IDF. We evaluated this function in Section 5.

TF/IDF has long been used in IR and Web search [24]. Lucene was released in 1999. For a long time it was somewhat slow and inaccurate, and was largely ignored by researchers [18]. In 2015, however, Lucene adopted cutting-edge techniques such as BM25 and block-max WAND. It is now viewed as quite accurate and fast, and has attracted attention from IR researchers [18]. TF/IDF has long been used in the matching step of EM [8].

The work [40] has studied top- $k$  search, but only for similarity measures such as Jaccard, cosine, dice, and overlap, for string matching. As far as we can tell, top- $k$  tf/idf search has been studied intensively by IR researchers (resulting in the block-max WAND technique), but not by database researchers. The work [41] develops AutoBlock, which was shown by [38] to underperform the DL methods Autoencoder and Hybrid, which underperform Sparkly. The work [22] also addresses blocking. But it maximizes recall while keeping precision (i.e., the fraction of pairs in the blocking output that are correct matches) above a threshold. We consider a fundamentally different problem of maximizing recall for any given  $k$  (i.e., any given blocking output size).

The share-nothing architecture of Sparkly is reminiscent of share-nothing architectures for parallel processing of relational data [37], and our Spark-based probing method for blocking is reminiscent of distributed/parallel joins for relational data [20, 21]. But here we consider the novel context of blocking for EM. Finally, the work [2] describes an industrial blocking solution at Amazon, which uses meta blocking to manage token-centric blocks and uses sophisticated techniques to scale.

## 7 CONCLUSIONS

Despite decades of research, tf/idf blocking has received very little attention. Yet anecdotal evidence suggests that it can do very well. As a result, in this paper we have performed an in-depth examination of tf/idf blocking.

We developed Sparkly, a novel solution that performs top- $k$  tf/idf blocking, using Lucene and Spark in a distributed share-nothing architecture. We developed techniques to select good attribute/tokenizer pairs to block on, making Sparkly completely automatic. Extensive experiments show that Sparkly outperforms 8 state-of-the-art blocking solutions and scales to large datasets.

Overall, our work suggests that tf/idf blocking should receive more attention, that future blocking work should consider Sparkly as a baseline, and that the distributed share-nothing architecture of Sparkly provides a promising starting point to build blocking solutions that are scalable, predictable, and extensible.

## ACKNOWLEDGMENTS

We are grateful to the reviewers whose insightful comments greatly improve this paper. We thank George Papadakis and Themis Palpanas for assistance with running experiments with JedAI and feedback on an earlier draft, and Saravanan Thirumuruganathan for assistance with the DeepBlocker software.

## REFERENCES

- [1] Nils Barlaug and Jon Atle Gulla. 2020. Neural networks for entity matching. *arXiv preprint arXiv:2010.11075* (2020).
- [2] Andrew Borthwick, Stephen Ash, Bin Pang, Shehzad Qureshi, and Timothy Jones. 2020. Scalable Blocking for Very Large Databases. In *ECML PKDD 2020 Workshops - Workshops of the European Conference on Machine Learning and Knowledge Discovery in Databases (ECML PKDD 2020): SoGood 2020, PDFL 2020, MLCS 2020, NFMCP 2020, DINA 2020, EDML 2020, XKDD 2020 and INRA 2020, Ghent, Belgium, September 14-18, 2020, Proceedings (Communications in Computer and Information Science)*, Irena Koprinska et al. (Eds.), Vol. 1323. Springer, 303–319. [https://doi.org/10.1007/978-3-030-65965-3\\_20](https://doi.org/10.1007/978-3-030-65965-3_20)
- [3] Andrei Z. Broder, Michael Herscovici, and Jason Zien. 2003. Efficient query evaluation using a two-level retrieval process. In *In Proc. of the 12th ACM Conf. on Information and Knowledge Management*.
- [4] Surajit Chaudhuri, Kris Ganjam, Venkatesh Ganti, and Rajeev Motwani. 2003. Robust and Efficient Fuzzy Match for Online Data Cleaning. In *Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data, San Diego, California, USA, June 9-12, 2003*, Alon Y. Halevy, Zachary G. Ives, and AnHai Doan (Eds.). ACM, 313–324. <https://doi.org/10.1145/872757.872796>
- [5] Peter Christen. 2011. A survey of indexing techniques for scalable record linkage and deduplication. *IEEE transactions on knowledge and data engineering* 24, 9 (2011), 1537–1555.
- [6] Peter Christen. 2012. *Data Matching - Concepts and Techniques for Record Linkage, Entity Resolution, and Duplicate Detection*. Springer. <https://doi.org/10.1007/978-3-642-31164-2>
- [7] Vassilis Christophides, Vasilis Efthymiou, Themis Palpanas, George Papadakis, and Kostas Stefanidis. 2020. An overview of end-to-end entity resolution for big data. *ACM Computing Surveys (CSUR)* 53, 6 (2020), 1–42.
- [8] William W. Cohen, Pradeep Ravikumar, and Stephen E. Fienberg. 2003. A Comparison of String Distance Metrics for Name-Matching Tasks. In *Proceedings of IJCAI-03 Workshop on Information Integration on the Web (IIWeb-03), August 9-10, 2003, Acapulco, Mexico*, Subbarao Kambhampati and Craig A. Knoblock (Eds.), 73–78. <http://www.isi.edu/info-agents/workshops/ijcai03/papers/Cohen-p.pdf>
- [9] Sanjib Das, Paul Suganthan G. C., AnHai Doan, Jeffrey F. Naughton, Ganesh Krishnan, Rohit Deep, Esteban Arcaute, Vijay Raghavendra, and Youngchoon Park. 2017. Falcon: Scaling Up Hands-Off Crowdsourced Entity Matching to Build Cloud Services. In *Proceedings of the 2017 ACM International Conference on Management of Data, SIGMOD Conference 2017, Chicago, IL, USA, May 14-19, 2017*, Semih Salihoglu, Wenchao Zhou, Rada Chirkova, Jun Yang, and Dan Suciu (Eds.). ACM, 1431–1446. <https://doi.org/10.1145/3035918.3035960>
- [10] Constantinos Dimopoulos, Sergey Nepomnyachiy, and Torsten Suel. 2013. Optimizing top-k document retrieval strategies for block-max indexes. In *Sixth ACM International Conference on Web Search and Data Mining, WSDM 2013, Rome, Italy, February 4-8, 2013*, Stefano Leonardi, Alessandro Panconesi, Paolo Ferragina, and Aristides Gionis (Eds.). ACM, 113–122. <https://doi.org/10.1145/2433396.2433412>
- [11] Shuai Ding and Torsten Suel. 2011. Faster top-k document retrieval using block-max indexes. In *Proceeding of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2011, Beijing, China, July 25-29, 2011*, Wei-Ying Ma, Jian-Yun Nie, Ricardo Baeza-Yates, Tat-Seng Chua, and W. Bruce Croft (Eds.). ACM, 993–1002. <https://doi.org/10.1145/2009916.2010048>
- [12] A. Doan, A. Halevy, and Z. Ives. 2012. *Principles of Data Integration*. Elsevier.
- [13] Muhammad Ebraheem, Saravanan Thirumuruganathan, Shafiq Joty, Mourad Ouzzani, and Nan Tang. 2018. Distributed representations of tuples for entity resolution. *PVLDB* 11, 11 (2018), 1454–1467.
- [14] Ahmed K. Elmagarmid, Panagiotis G. Ipeirotis, and Vassilios S. Verykios. 2007. Duplicate Record Detection: A Survey. *TKDE* 19, 1 (2007).
- [15] C. Gokhale, S. Das, A. Doan, J. F. Naughton, N. Rampalli, J. Shavlik, and X. Zhu. 2014. Corleone: hands-off crowdsourcing for entity matching. *SIGMOD*.
- [16] Yash Govind, Pradap Konda, Paul Suganthan G. C., Philip Martinkus, Palaniappan Nagarajan, Han Li, Aravind Soundararajan, Sidharth Mudgal, Jeffrey R. Ballard, Haojun Zhang, Adel Ardalan, Sanjib Das, Derek Paulsen, Amanpreet Singh Saini, Erik Paulson, Youngchoon Park, Marshall Carter, Mingju Sun, Glenn Moo Fung, and AnHai Doan. 2019. Entity Matching Meets Data Science: A Progress Report from the Magellan Project. In *Proceedings of the 2019 International Conference on Management of Data, SIGMOD Conference 2019, Amsterdam, The Netherlands, June 30 - July 5, 2019*, Peter A. Boncz, Stefan Manegold, Anastasia Ailamaki, Amol Deshpande, and Tim Kraska (Eds.). ACM, 389–403. <https://doi.org/10.1145/3299869.3314042>
- [17] Yash Govind, Erik Paulson, Palaniappan Nagarajan, Paul Suganthan G. C., AnHai Doan, Youngchoon Park, Glenn Fung, Devin Conathan, Marshall Carter, and Mingju Sun. 2018. CloudMatcher: A Hands-Off Cloud/Crowd Service for Entity Matching. *Proc. VLDB Endow.* 11, 12 (2018), 2042–2045. <https://doi.org/10.14778/3229863.3236255>
- [18] Adrien Grand, Robert Muir, Jim Ferenczi, and Jimmy Lin. 2020. From MAXSCORE to Block-Max Wand: The Story of How Lucene Significantly Improved Query Evaluation Performance. In *Advances in Information Retrieval - 42nd European Conference on IR Research, ECIR 2020, Lisbon, Portugal, April 14-17, 2020, Proceedings, Part II (Lecture Notes in Computer Science)*, Joemon M. Jose, Emine Yilmaz, João Magalhães, Pablo Castells, Nicola Ferro, Mário J. Silva, and Flávio Martins (Eds.), Vol. 12036. Springer, 20–27. [https://doi.org/10.1007/978-3-030-45442-5\\_3](https://doi.org/10.1007/978-3-030-45442-5_3)
- [19] Pradap Konda, Sanjib Das, Paul Suganthan GC, AnHai Doan, Adel Ardalan, Jeffrey R Ballard, Han Li, Fatemah Panahi, Haojun Zhang, Jeff Naughton, et al. 2016. Magellan: Toward building entity matching management systems. *PVLDB* 9, 13 (2016), 1581–1584.
- [20] Donald Kossmann. 2000. The State of the art in distributed query processing. *ACM Comput. Surv.* 32, 4 (2000), 422–469. <https://doi.org/10.1145/371578.371598>
- [21] Paraschos Kouttris, Semih Salihoglu, and Dan Suciu. 2018. Algorithmic Aspects of Parallel Data Processing. *Found. Trends Databases* 8, 4 (2018), 239–370. <https://doi.org/10.1561/19000000055>
- [22] Peng Li, Xiang Cheng, Xu Chu, Yeye He, and Surajit Chaudhuri. 2021. Auto-FuzzyJoin: Auto-Program Fuzzy Similarity Joins Without Labeled Examples. In *SIGMOD '21: International Conference on Management of Data, Virtual Event, China, June 20-25, 2021*, Guoliang Li, Zhanhui Li, Stratos Idreos, and Divesh Srivastava (Eds.). ACM, 1064–1076. <https://doi.org/10.1145/3448016.3452824>
- [23] Yuliang Li, Jinfeng Li, Yoshihiko Suhara, AnHai Doan, and Wang-Chiew Tan. 2020. Deep entity matching with pre-trained language models. *PVLDB* 14, 1 (2020), 50–60.
- [24] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. 2008. *Introduction to information retrieval*. Cambridge University Press. <https://doi.org/10.1017/CBO9780511809071>
- [25] Sidharth Mudgal, Han Li, Theodoros Rekatsinas, AnHai Doan, Youngchoon Park, Ganesh Krishnan, Rohit Deep, Esteban Arcaute, and Vijay Raghavendra. 2018. Deep learning for entity matching: A design space exploration. In *SIGMOD*.
- [26] Felix Naumann and Melanie Herschel. 2010. An introduction to duplicate detection. *Synthesis Lectures on Data Management* 2, 1 (2010), 1–87.
- [27] Kevin OHare, Anna Jurek-Loughrey, and Cassio P. de Campos. 2019. A review of unsupervised and semi-supervised blocking methods for record linkage. *Linking and Mining Heterogeneous and Multi-view Data* (2019), 79–105.
- [28] Kevin O'Hare, Anna Jurek-Loughrey, and Cassio P. de Campos. 2022. High-Value Token-Blocking: Efficient Blocking Method for Record Linkage. *ACM Trans. Knowl. Discov. Data* 16, 2 (2022), 24:1–24:17. <https://doi.org/10.1145/3450527>
- [29] G. Papadakis, M. Fischella, F. Schoger, G. Mandilaras, N. Augsten, and W. Nejdl. 2022. *Benchmarking Filtering Techniques for Entity Resolution*. Technical Report. arXiv:2022.12521v3.
- [30] George Papadakis, Ekaterini Ioannou, Emanouil Thanos, and Themis Palpanas. 2021. The Four Generations of Entity Resolution. *Synthesis Lectures on Data Management* 16, 2 (2021), 1–170.
- [31] George Papadakis, George Mandilaras, Luca Gagliardelli, Giovanni Simonini, Emanouil Thanos, George Giannakopoulos, Sonia Bergamaschi, Themis Palpanas, and Manolis Koubarakis. 2020. Three-dimensional Entity Resolution with JedAI. *Information Systems* 93 (2020), 101565.
- [32] George Papadakis, Dimitrios Skoutas, Emmanouil Thanos, and Themis Palpanas. 2020. Blocking and filtering techniques for entity resolution: A survey. *ACM Computing Surveys (CSUR)* 53, 2 (2020), 1–42.
- [33] George Papadakis, Leonidas Tsekouras, Emmanouil Thanos, Nikiforos Pittaras, Giovanni Simonini, Dimitrios Skoutas, Paul Isaris, George Giannakopoulos, Themis Palpanas, and Manolis Koubarakis. 2020. JedAI3: beyond batch, blocking-based Entity Resolution. In *EDBT*. 603–606.
- [34] D. Paulsen, Y. Govind, and A. Doan. 2022. Homepage of the Sparkly Blocking System. <https://github.com/anhaidgroup/sparkly>.
- [35] Anna Primpeli, Ralph Peeters, and Christian Bizer. 2019. The WDC Training Dataset and Gold Standard for Large-Scale Product Matching. In *Companion of The 2019 World Wide Web Conference, WWW 2019, San Francisco, CA, USA, May 13-17, 2019*, Sihem Amer-Yahia, Mohammad Mahdian, Ashish Goel, Geert-Jan Houben, Kristina Lerman, Julian J. McAuley, Ricardo Baeza-Yates, and Leila Zia (Eds.). ACM, 381–386. <https://doi.org/10.1145/3308560.3316609>
- [36] Rudi Seitz. 2022. *Understanding tfidf and BM25*. Technical Report. <https://kmlwlc.com/index.php/2020/03/20/understanding-tf-idf-and-bm-25>.
- [37] Michael Stonebraker. 1986. The Case for Shared Nothing. *IEEE Database Eng. Bull.* 9, 1 (1986), 4–9. <http://sites.computer.org/debull/86MAR-CD.pdf>
- [38] Saravanan Thirumuruganathan, Han Li, Nan Tang, Mourad Ouzzani, Yash Govind, Derek Paulsen, Glenn Fung, and AnHai Doan. 2021. Deep Learning for Blocking in Entity Matching: A Design Space Exploration. *Proc. VLDB Endow.* 14, 11 (2021), 2459–2472. <https://doi.org/10.14778/3476249.3476294>
- [39] Frank Wilcoxon. 1945. Individual Comparisons by Ranking Methods. *Biometrics Bulletin* 1, 6 (1945), 80–83. <http://www.jstor.org/stable/3001968>
- [40] C. Xiao, W. Wang, X. Lin, and H. Shang. 2009. Top-k set similarity joins. *ICDE*.
- [41] Wei Zhang, Hao Wei, Bunyamin Sisman, Xin Luna Dong, Christos Faloutsos, and Davd Page. 2020. AutoBlock: A hands-off blocking framework for entity matching. In *WSDM*. 744–752.