




# Pattern Masking for Dictionary Matching: Theory and Practice

Panagiotis Charalampopoulos<sup>1</sup> · Huiping Chen<sup>2</sup> · Peter Christen<sup>3</sup> ·  
Grigorios Loukides<sup>4</sup> · Nadia Pisanti<sup>5</sup> · Solon P. Pissis<sup>6,7</sup>  ·  
Jakub Radoszewski<sup>8</sup>

Received: 1 December 2021 / Accepted: 20 January 2024  
© The Author(s) 2024

## Abstract

Data masking is a common technique for sanitizing sensitive data maintained in database systems which is becoming increasingly important in various application areas, such as in record linkage of personal data. This work formalizes the Pattern Masking for Dictionary Matching (PMDM) problem: given a dictionary  $\mathcal{D}$  of  $d$  strings, each of length  $\ell$ , a query string  $q$  of length  $\ell$ , and a positive integer  $z$ , we are asked to compute a smallest set  $K \subseteq \{1, \dots, \ell\}$ , so that if  $q[i]$  is replaced by a wildcard for all  $i \in K$ , then  $q$  matches at least  $z$  strings from  $\mathcal{D}$ . Solving PMDM allows providing data utility guarantees as opposed to existing approaches. We first show, through a reduction from the well-known  $k$ -Clique problem, that a decision version of the PMDM problem is NP-complete, even for binary strings. We thus approach the problem from a more practical perspective. We show a combinatorial  $\mathcal{O}((d\ell)^{|K|/3} + d\ell)$ -time and  $\mathcal{O}(d\ell)$ -space algorithm for PMDM for  $|K| = \mathcal{O}(1)$ . In fact, we show that we cannot hope for a faster combinatorial algorithm, unless the combinatorial  $k$ -Clique hypothesis fails (Abboud et al. in SIAM J Comput 47:2527–2555, 2018; Lincoln et al., in: 29th ACM-SIAM Symposium on Discrete Algorithms (SODA), 2018). Our combinatorial algorithm, executed with small  $|K|$ , is the backbone of a greedy heuristic that we propose. Our experiments on real-world and synthetic datasets show that our heuristic finds nearly-optimal solutions in practice and is also very efficient. We also generalize this algorithm for the problem of masking multiple query strings simultaneously so that every string has at least  $z$  matches in  $\mathcal{D}$ . PMDM can be viewed as a generalization of the decision version of the dictionary matching with mismatches problem: by querying a PMDM data structure with string  $q$  and  $z = 1$ , one obtains the minimal number of mismatches of  $q$  with any string from  $\mathcal{D}$ . The query time

---

✉ Huiping Chen  
h.chen.13@bham.ac.uk

✉ Solon P. Pissis  
solon.pissis@cwi.nl

Extended author information available on the last page of the article

or space of all known data structures for the *more restricted* problem of dictionary matching with at most  $k$  mismatches incurs some exponential factor with respect to  $k$ . A simple exact algorithm for PMDM runs in time  $\mathcal{O}(2^\ell d)$ . We present a data structure for PMDM that answers queries over  $\mathcal{D}$  in time  $\mathcal{O}(2^{\ell/2}(2^{\ell/2} + \tau)\ell)$  and requires space  $\mathcal{O}(2^\ell d^2/\tau^2 + 2^{\ell/2}d)$ , for any parameter  $\tau \in [1, d]$ . We complement our results by showing a two-way polynomial-time reduction between PMDM and the Minimum Union problem [Chlamtáč et al., ACM-SIAM Symposium on Discrete Algorithms (SODA) 2017]. This gives a polynomial-time  $\mathcal{O}(d^{1/4+\epsilon})$ -approximation algorithm for PMDM, which is tight under a plausible complexity conjecture. This is an extended version of a paper that was presented at International Symposium on Algorithms and Computation (ISAAC) 2021.

**Keywords** String algorithms · Dictionary matching · Wildcards · Record linkage · Query term dropping

## 1 Introduction

Let us start with a true incident to illustrate the essence of the computational problem formalized in this work. In the Netherlands, water companies bill the non-drinking and drinking water separately. The 6th author of this paper had direct debit for the former but not for the latter. When he tried to set up the direct debit for the latter, he received the following masked message by the company:

```
Is this you?
Initial: S.   Name: P*****s
E-mail address: s*****13@g***l.com
Bank account number: NL10RABO*****11.
```

The rationale of the data masking is: the client should be able to identify themselves to help the companies *link* the client's profiles, without inferring the identity of any other client via a *linking attack* [33, 67], so that clients' privacy is preserved.<sup>1</sup> Thus, the masked version of the data is required to conceal as few symbols as possible, so that the client can recognize their data, but also to correspond to a sufficient number of other clients, so that it is hard for a successful linking attack to be performed.

This requirement can be formalized as the Pattern Masking for Dictionary Matching (PMDM) problem: Given a dictionary  $\mathcal{D}$  of  $d$  strings, each of length  $\ell$ , a query string  $q$  of length  $\ell$ , and a positive integer  $z$ , PMDM asks to compute a smallest set  $K \subseteq \{1, \dots, \ell\}$ , so that if  $q[i]$ , for all  $i \in K$ , is replaced by a wildcard,  $q$  matches at least  $z$  strings from  $\mathcal{D}$ . The PMDM problem applies data masking, a common operation to sanitize personal data maintained in database systems [27, 68, 69]. In particular, PMDM lies at the heart of record linkage of databases containing personal data [24, 48, 50, 63, 72, 74], which is the main application we consider in this work.

*Record linkage* is the task of identifying records that refer to the same entities across databases, in situations where no entity identifiers are available in these databases [22, 39, 58]. This task is of high importance in various application domains featuring

<sup>1</sup> In linking attacks, the adversary is a data recipient who uses the released data of an individual together with publicly available data, or with background knowledge, to infer the individual's identity.

personal data, ranging from the health sector and social science research, to national statistics and crime and fraud detection [24, 45]. In a typical setting, the task is to link two databases that contain names or other attributes, known collectively as quasi-identifiers (QIDs) [73]. The similarity between each pair of records (a record from one of the databases and a record from the other) is calculated with respect to their values in QIDs, and then all compared record pairs are classified into matches (the pair is assumed to refer to the same person), non-matches (the two records in the pair are assumed to refer to different people), and potential matches (no decision about whether the pair is a match or non-match can be made) [22, 39].

Unfortunately, potential matches happen quite often [8]. A common approach [63, 72] to deal with potential matches is to conduct a manual clerical review, where a domain expert looks at the attribute values in record pairs and then makes a manual match or non-match decision. At the same time, to comply with policies and legislation, one needs to prevent domain experts from inferring the identity of the people represented in the manually assessed record pairs [63]. The challenge is to achieve desired data protection/utility guarantees; i.e. enabling a domain expert to make good decisions without inferring peoples' identities.

To address this challenge, we can solve PMDM twice, for a potential match  $(q_1, q_2)$ . The first time we use as input the query string  $q_1$  and a reference dictionary (database)  $\mathcal{D}$  containing personal records from a sufficiently large population (typically, much larger than the databases to be linked). The second time, we use as input  $q_2$  instead of  $q_1$ . Since each masked  $q$  derived by solving PMDM matches at least  $z$  records in  $\mathcal{D}$ , the domain expert would need to distinguish between at least  $z$  individuals in  $\mathcal{D}$  to be able to infer the identity of the individual corresponding to the masked string. The underlying assumption is that  $\mathcal{D}$  contains one record per individual. Also, some wildcards from one masked string can be superimposed on another to ensure that the expert does not gain more knowledge from combining the two strings, and the resulting strings would still match at least  $z$  records in  $\mathcal{D}$ . Thus, by solving PMDM in this setting, we provide privacy guarantees alike  $z$ -map [70]; a variant of the well-studied  $z$ -anonymity [65] privacy model.<sup>2</sup> In  $z$ -map, each record of a dataset must match at least  $z$  records in a reference dataset, from which the dataset is derived. In our setting, we consider a pattern that is not necessarily contained in the reference dataset. Offering such privacy is desirable in real record linkage systems where databases containing personal data are being linked [24, 50, 74]. On the other hand, since each masked  $q$  contains the minimum number of wildcards, the domain expert is still able to use the masked  $q$  to meaningfully classify a record pair as a match or as a non-match.

Offering such utility is again desirable in record linkage systems [63]. Record linkage is an important application for our techniques, because no existing approach can provide privacy and utility guarantees when releasing linkage results to domain experts [49]. In particular, existing approaches [49, 50] recognize the need to offer privacy by preventing the domain expert from distinguishing between a small number of individuals, but they provide *no algorithm* for offering such privacy, let alone an algorithm offering utility guarantees as we do.

<sup>2</sup> The notation used for such privacy models is generally  $k$  instead of  $z$ , e.g.  $k$ -anonymity [66, 70].

A secondary application where PMDM is of importance is *query term dropping*, an information retrieval task that seeks to drop keywords (terms) from a query, so that the remaining keywords retrieve a sufficiently large number of documents. This task is performed by search engines, such as Google [7], and by e-commerce platforms such as e-Bay [51], to improve users' experience [35, 71] by making sufficiently many search results available to users. For example, e-Bay applies query term dropping, removing one term, in our test query:

```
Query: vacuum database cleaner
Query results: 0 results found for vacuum database cleaner
              42 results found for vacuum cleaner
```

We could perform query term dropping by solving PMDM in a setting where strings in a dictionary correspond to document terms and a query string corresponds to a user's query. Then, we provide the user with the masked query, after removing all wildcards, and with its matching strings from the dictionary. Two remarks are in order for this application. First, we consider a setting where the keyword order matters. This occurs, for example, when using *phrase search* in Google.<sup>3</sup> Second, since the dictionary may contain strings of different length, PMDM should be applied only to the dictionary strings that have the same length as the query string.

Query term dropping is a relevant application for our techniques, because existing techniques [71] do not minimize the number of dropped terms. Rather, they drop keywords randomly, which may unnecessarily shorten the query, or drop keywords based on custom rules, which is not sufficiently generic to deal with all queries. More generally, our techniques can be applied to drop terms from any top- $z$  database query [42] to ensure there are  $z$  results in the query answer.

**Related algorithmic work** Let us denote the wildcard symbol by  $\star$  and provide a brief overview of works related to PMDM, the main problem considered in this paper.

- *Partial Match*: Given a dictionary  $\mathcal{D}$  of  $d$  strings over an alphabet  $\Sigma = \{0, 1\}$ , each of length  $\ell$ , and a string  $q$  over  $\Sigma \sqcup \{\star\}$  of length  $\ell$ , the problem asks whether  $q$  matches any string from  $\mathcal{D}$ . This is a well-studied problem [13, 18, 44, 56, 59, 60, 64]. Patrascu [59] showed that any data structure for the Partial Match problem with cell-probe complexity  $t$  must use space  $2^{\Omega(\ell/t)}$ , assuming the word size is  $\Theta(d^{1-\epsilon}/t)$ , for any constant  $\epsilon > 0$ . The key difference to PMDM is that the wildcard positions in the query strings are fixed.
- *Dictionary Matching with  $k$ -errors*: A similar line of research to that of Partial Match has been conducted under the Hamming and edit distances, where, in this case,  $k$  is the maximum allowed distance between the query string and a dictionary string [10, 11, 14, 16, 26, 77]. The structure of Dictionary Matching with  $k$ -errors is very similar to Partial Match as each wildcard in the query string gives  $|\Sigma|$  possibilities for the corresponding symbol in the dictionary strings. On the other hand, in Partial Match the wildcard positions are fixed.

The PMDM problem is a generalization of the decision version of the Dictionary Matching with  $k$ -errors problem (under Hamming distance): by querying a data

<sup>3</sup> An indicative example of a query in Google is “free blue tv” which yielded 7050 results, as blue tv is a well-known app, whereas the query “blue free tv” yielded only 5 results.

structure for PMDM with string  $q$  and  $z = 1$ , one obtains the minimum number of mismatches of  $q$  with any string from  $\mathcal{D}$ , which suffices to answer the decision version of the Dictionary Matching with  $k$ -errors problem. The query time or space of all known data structures for Dictionary Matching with  $k$ -mismatches incurs some exponential factor with respect to  $k$ . In [25], Cohen-Addad et al. showed that, in the pointer machine model, for the reporting version of the problem, one cannot avoid exponential dependency on  $k$  either in the space or in the query time. In the word-RAM model, Rubinfeld showed that, conditional on the Strong Exponential Time Hypothesis [15], any data structure that can be constructed in time polynomial in the total size  $|\mathcal{D}|$  of the strings in the dictionary cannot answer queries in time strongly sublinear in  $|\mathcal{D}|$ .

We next provide a brief overview of other algorithmic works related to PMDM.

- *Dictionary Matching with  $k$ -wildcards*: Given a dictionary  $\mathcal{D}$  of total size  $N$  over an alphabet  $\Sigma$  and a query string  $q$  of length  $\ell$  over  $\Sigma \sqcup \{\star\}$  with up to  $k$  wildcards, the problem asks for the set of matches of  $q$  in  $\mathcal{D}$ . This is essentially a parameterized variant of the Partial Match problem. The seminal paper of Cole et al. [26] proposed a data structure occupying  $\mathcal{O}(N \log^k N)$  space allowing for  $\mathcal{O}(\ell + 2^k \log \log N + |\text{output}|)$ -time querying. This data structure is based on recursively computing a heavy-light decomposition of the suffix tree and copying the subtrees hanging off light children. Generalizations and slight improvements have been proposed in [12, 34, 54]. In [12] the authors also proposed an alternative data structure that instead of a  $\log^k N$  factor in the space complexity has a multiplicative  $|\Sigma|^{k^2}$  factor. Nearly-linear-sized data structures that essentially try all different combinations of letters in the place of wildcards and hence incur a  $|\Sigma|^k$  factor in the query time have been proposed in [12, 52]. On the lower bound side, Afshani and Nielsen [2] showed that, in the pointer machine model, essentially any data structure for the problem in scope must have exponential dependency on  $k$  in either the space or the query time, explaining the barriers hit by the existing approaches.
- *Enumerating Motifs with  $k$ -wildcards*: Given an input string  $s$  of length  $n$  over an alphabet  $\Sigma$  and positive integers  $k$  and  $z$ , this problem asks to enumerate all motifs over  $\Sigma \sqcup \{\star\}$  with up to  $k$  wildcards that occur at least  $z$  times in  $s$ . As the size of the output is exponential in  $k$ , the enumeration problem has such a lower bound. Several approaches exist for efficient motif enumeration, all aimed at reducing the impact of the output's size: efficient indexing to minimize the output delay [6, 37]; exploiting a hierarchy of wildcards positions according to the number of occurrences [9]; defining a subset of motifs of fixed-parameter tractable size (in  $k$  or  $z$ ) that can generate all the others [61, 62], or defining maximality notions meaning a subset of the motifs that implicitly include all the others [31, 36].

**Our Contributions** We consider the word-RAM model of computations with  $w$ -bit machine words, where  $w = \Omega(\log(d\ell))$ , for stating our results. We make the following contributions:

1. (Sect. 3) A reduction from the  $k$ -Clique problem to a decision version of the PMDM problem, which implies that PMDM is NP-hard, even for strings over a binary alphabet. The reduction also implies conditional hardness of the PMDM problem. We also present a generalized reduction from the  $(c, k)$ -Hyperclique problem [55].
2. (Sect. 4) A combinatorial  $\mathcal{O}((d\ell)^{k/3} + d\ell)$ -time and  $\mathcal{O}(d\ell)$ -space algorithm for PMDM if  $k = |K| = \mathcal{O}(1)$ , which is optimal if the combinatorial  $k$ -Clique hypothesis is true.
3. (Sect. 5) We consider a generalized version of PMDM, referred to as MPMDM: we are given a collection  $\mathcal{M}$  of  $m$  query strings (instead of one query string) and we are asked to compute a smallest set  $K$  so that, for every  $q$  from  $\mathcal{M}$ , if  $q[i]$ , for all  $i \in K$ , is replaced by a wildcard, then  $q$  matches at least  $z$  strings from dictionary  $\mathcal{D}$ . We show an  $\mathcal{O}((d\ell)^{k/3} z^{m-1} + d\ell)$ -time algorithm for MPMDM, for  $k = |K| = \mathcal{O}(1)$  and  $m = \mathcal{O}(1)$ .
4. (Sect. 6) A data structure for PMDM that answers queries over  $\mathcal{D}$  in  $\mathcal{O}(2^{\ell/2}(2^{\ell/2} + \tau)\ell)$  time and requires space  $\mathcal{O}(2^\ell d^2/\tau^2 + 2^{\ell/2}d)$ , for any parameter  $\tau \in [1, d]$ .
5. (Sect. 7) A polynomial-time  $\mathcal{O}(d^{1/4+\epsilon})$ -approximation algorithm for PMDM, which we show to be tight under a plausible complexity conjecture.
6. (Sect. 8) A greedy heuristic based on the  $\mathcal{O}((d\ell)^{k/3} + d\ell)$ -time algorithm.
7. (Sect. 9) An extensive experimental evaluation on real-world and synthetic data demonstrating that our heuristic finds nearly-optimal solutions in practice and is also very efficient. In particular, our heuristic finds optimal or nearly-optimal solutions for PMDM on a dataset with six million records in less than 3 s.

We conclude this paper with a few open questions in Sect. 10.

This paper is an extended version of a paper that was presented at ISAAC 2021 [17]. Compared to [17], Sects. 8 and 9 are new, while Sects. 1, 5 and 6 contain additional details that were omitted from [17] due to space constraints.

## 2 Definitions and Notation

**Strings** An *alphabet*  $\Sigma$  is a finite nonempty set whose elements are called *letters*. We assume throughout an integer alphabet  $\Sigma = [1, |\Sigma|]$ . Let  $x = x[1] \cdots x[n]$  be a *string* of length  $|x| = n$  over  $\Sigma$ . For two indices  $1 \leq i \leq j \leq n$ ,  $x[i..j] = x[i] \cdots x[j]$  is the *substring* of  $x$  that starts at position  $i$  and ends at position  $j$  of  $x$ . By  $\varepsilon$  we denote the *empty string* of length 0. A *prefix* of  $x$  is a substring of  $x$  of the form  $x[1..j]$ , and a *suffix* of  $x$  is a substring of  $x$  of the form  $x[i..n]$ . A *dictionary* is a collection of strings. We also consider alphabet  $\Sigma_\star = \Sigma \sqcup \{\star\}$ , where  $\star$  is a *wildcard* letter that is not in  $\Sigma$  and *matches* all letters from  $\Sigma_\star$ . Then, given a string  $x$  over  $\Sigma_\star$  and a string  $y$  over  $\Sigma$  with  $|x| = |y|$ , we say that  $x$  *matches*  $y$  if and only if  $x[i] = y[i]$  or  $x[i] = \star$ , for all  $1 \leq i \leq |x|$ . Given a string  $x$  of length  $n$  and a set  $S \subseteq \{1, \dots, n\}$ , we denote by  $x_S = x \otimes S$  the string obtained by first setting  $x_S = x$  and then  $x_S[i] = \star$ , for all  $i \in S$ . We then say that  $x$  is *masked* by  $S$ .

The main problem considered in this paper is the following.

PATTERN MASKING FOR DICTIONARY MATCHING (PMDM)

**Input:** A dictionary  $\mathcal{D}$  of  $d$  strings, each of length  $\ell$ , a string  $q$  of length  $\ell$ , and a positive integer  $z$ .

**Output:** A smallest set  $K \subseteq \{1, \dots, \ell\}$  such that  $q_K = q \otimes K$  matches at least  $z$  strings from  $\mathcal{D}$ .

We refer to the problem of computing only the size  $k$  of a smallest set  $K$  as PMDM-SIZE. We also consider the data structure variant of the PMDM problem in which  $\mathcal{D}$  is given for preprocessing, and  $q, z$  queries are to be answered on-line. Throughout, we assume that  $k \geq 1$  as the case  $k = 0$  corresponds to the well-studied dictionary matching problem for which there exists a classic optimal solution [3]. We further assume  $z \leq d$ ; otherwise the PMDM has trivially no solution. In what follows, we use  $N$  to denote  $d\ell$ .

**Tries.** Let  $\mathcal{M}$  be a finite set containing  $m > 0$  strings over  $\Sigma$ . The *trie* of  $\mathcal{M}$ , denoted by  $\mathcal{R}(\mathcal{M})$ , contains a node for every distinct prefix of a string in  $\mathcal{M}$ ; the root node is  $\varepsilon$ ; the set of leaf nodes is  $\mathcal{M}$ ; and edges are of the form  $(u, \alpha, u\alpha)$ , where  $u$  and  $u\alpha$  are nodes and  $\alpha \in \Sigma$  is the label. The *compacted trie* of  $\mathcal{M}$ , denoted by  $\mathcal{T}(\mathcal{M})$ , contains the root, the branching nodes, and the leaf nodes of  $\mathcal{R}(\mathcal{M})$ . Each maximal branchless path segment from  $\mathcal{R}(\mathcal{M})$  is replaced by a single edge, and a fragment of a string  $M \in \mathcal{M}$  is used to represent the label of this edge in  $\mathcal{O}(1)$  space. The size of  $\mathcal{T}(\mathcal{M})$  is thus  $\mathcal{O}(m)$ . The most well-known example of a compacted trie is the suffix tree of a string: the compacted trie of all the suffixes of the string [75]. To access the children of a trie node by the first letter of their edge label in  $\mathcal{O}(1)$  time we use perfect hashing [32]. In this case, the claimed complexities hold *with high probability* (w.h.p., for short), that is, with probability at least  $1 - N^{-c}$  (recall that  $N = d\ell$ ), where  $c > 0$  is a constant fixed at construction time. Assuming that the children of every trie node are sorted by the first letters of their edge labels, randomization can be avoided at the expense of a  $\log |\Sigma|$  factor incurred by binary searching for the appropriate child.

### 3 NP-hardness and Conditional Hardness of PMDM-SIZE

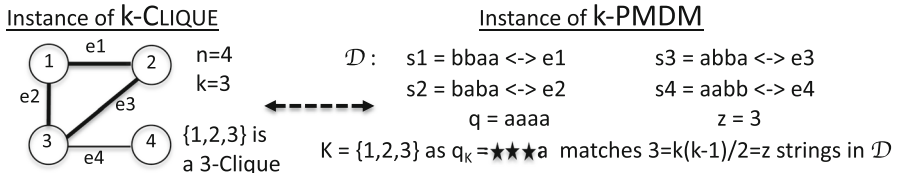
We show that the following decision version of PMDM-SIZE is NP-complete.

*k*-PMDM

**Input:** A dictionary  $\mathcal{D}$  of  $d$  strings, each of length  $\ell$ , a string  $q$  of length  $\ell$ , and positive integers  $z \leq d$  and  $k \leq \ell$ .

**Output:** Is there a set  $K \subseteq \{1, \dots, \ell\}$  of size  $k$ , such that  $q_K = q \otimes K$  matches at least  $z$  strings from  $\mathcal{D}$ ?

Our reduction is from the well-known NP-complete *k*-CLIQUE problem [46]: Given an undirected graph  $G$  on  $n$  nodes and a positive integer  $k$ , decide whether  $G$  contains a clique of size  $k$  (a *clique* is a subset of the nodes of  $G$  that are pairwise adjacent).



**Fig. 1** An example of the reduction from  $k$ -CLIQUE to  $k$ -PMDM. The solution for both is  $\{1, 2, 3\}$  as shown. Note that, for  $k = 4$ , the instance of 4-PMDM would need  $z = 6$  matches; neither this many matches can be found in  $\mathcal{D}$  nor a 4-clique can be found in the graph

**Theorem 3.1** Any instance of the  $k$ -CLIQUE problem for a graph with  $n$  nodes and  $m$  edges can be reduced in  $\mathcal{O}(nm)$  time to a  $k$ -PMDM instance with  $\ell = n$ ,  $d = m$  and  $\Sigma = \{a, b\}$ .

**Proof** Let  $G = (V, E)$  be an undirected graph on  $n = |V|$  nodes numbered 1 through  $n$ , in which we are looking for a clique of size  $k$ . We reduce  $k$ -CLIQUE to  $k$ -PMDM as follows. Consider the alphabet  $\{a, b\}$ . Set  $q = a^n$ , and for every edge  $(u, v) \in E$  such that  $u < v$ , add string  $a^{u-1}ba^{v-u-1}ba^{n-v}$  to  $\mathcal{D}$ . Set  $z = k(k - 1)/2$ . Then  $G$  contains a clique of size  $k$ , if and only if  $k$ -PMDM returns a positive answer. This can be seen by the fact that cliques of size  $k$  in  $G$  are in one-to-one correspondence with subsets  $K \subseteq \{1, \dots, n\}$  of size  $k$  for which  $q_K$  matches  $z$  strings from  $\mathcal{D}$ : the elements of  $K$  correspond to the nodes of a clique and the  $z$  strings correspond to its edges.  $k$ -PMDM is clearly in NP and the result follows.  $\square$

An example of the reduction from  $k$ -CLIQUE to  $k$ -PMDM is shown in Fig. 1.

**Corollary 3.2**  $k$ -PMDM is NP-complete for strings over a binary alphabet.

Any algorithm solving PMDM-SIZE can be trivially applied to solve  $k$ -PMDM.

**Corollary 3.3** PMDM-SIZE is NP-hard for strings over a binary alphabet.

**Remark 3.4** Given an undirected graph  $G$ , an *independent set* is a subset of nodes of  $G$  such that no two distinct nodes of the subset are adjacent. Let us note that the problem of computing a maximum clique in a graph  $G$ , which is equivalent to that of computing the maximum independent set in the complement of  $G$ , cannot be  $n^{1-\epsilon}$ -approximated in polynomial time, for any  $\epsilon > 0$ , unless  $P = NP$  [38, 78]. In Sect. 7, we show a polynomial-time  $\mathcal{O}(d^{1/4+\epsilon})$ -approximation algorithm for PMDM. We remark that this algorithm and Theorem 3.1 do not contradict the inapproximability results for the maximum clique problem, since our reduction from  $k$ -CLIQUE to  $k$ -PMDM cannot be adapted to a reduction from maximum clique to PMDM-SIZE.

Theorem 3.1 shows that solving  $k$ -PMDM efficiently even for strings over a binary alphabet would imply a breakthrough for the  $k$ -CLIQUE problem for which it is known that, in general, no fixed-parameter tractable algorithm with respect to parameter  $k$  exists unless the Exponential Time Hypothesis (ETH) fails [19, 43]. That is,  $k$ -CLIQUE has no  $f(k)n^{o(k)}$  time algorithm, and is thus W[1]-complete (again, under the ETH hypothesis). On the upper bound side,  $k$ -CLIQUE can be trivially solved in



$\mathcal{O}(n^k)$  time (enumerating all subsets of nodes of size  $k$ ), and this can be improved to  $\mathcal{O}(n^{\omega k/3})$  time for  $k$  divisible by 3 using square matrices multiplication ( $\omega$  is the exponent of square matrix multiplication). However, for general  $k \geq 3$  and any constant  $\epsilon > 0$ , the  $k$ -CLIQUE hypothesis states that there is no  $\mathcal{O}(n^{(\omega/3-\epsilon)k})$ -time algorithm and no combinatorial  $\mathcal{O}(n^{(1-\epsilon)k})$ -time algorithm [1, 55, 76]. Thus, conditional on the  $k$ -CLIQUE hypothesis, and since  $d\ell = nm = \mathcal{O}(n^3)$  and  $\ell = n$  (Theorem 3.1), we cannot hope to devise a combinatorial algorithm for  $k$ -PMDM with runtime  $\mathcal{O}((d\ell)^{(1-\epsilon)k/3})$  or  $\mathcal{O}(\ell^{(1-\epsilon)k})$  for any constant  $\epsilon > 0$ . In Sect. 4, we show a combinatorial  $\mathcal{O}(d\ell + \min\{(d\ell)^{k/3}, \ell^k\})$ -time algorithm, for constant  $k \geq 3$ , for the optimization version of  $k$ -PMDM (seeking to maximize the matches), which can then be trivially applied to solve  $k$ -PMDM in the same time complexity, thus matching the above conditional lower bound. Additionally, under the  $k$ -CLIQUE hypothesis, even with the aid of algebraic techniques, one cannot hope for an algorithm for  $k$ -PMDM with runtime  $\mathcal{O}((d\ell)^{(\omega/9-\epsilon)k})$  or  $\mathcal{O}(\ell^{(\omega/3-\epsilon)k})$ , for any constant  $\epsilon > 0$ .

In fact, as we show next, by reducing from the  $(c, k)$ -HYPERCLIQUE problem, which is not known to benefit from fast matrix multiplication [55], we obtain stronger conditional lower bounds for some values of  $d$  and  $\ell$ .

A hypergraph  $H$  is a pair  $(V, E)$ , where  $V$  is the set of nodes of  $H$  and  $E$  is a set of non-empty subsets of  $V$ , called hyperedges. The  $(c, k)$ -HYPERCLIQUE problem is defined as follows: Given a hypergraph  $H = (V, E)$  such that all of its hyperedges have size  $c$ , does there exist a set  $S$  of  $k > c$  nodes in  $V$  so that every subset of  $c$  nodes from  $S$  is a hyperedge? We call set  $S$  a  $(c, k)$ -hyperclique in  $H$ . We will reduce  $(c, k)$ -HYPERCLIQUE to  $k$ -PMDM in time  $\mathcal{O}(|V| \cdot |E|)$ .

**Theorem 3.5** *Any instance of the  $(c, k)$ -HYPERCLIQUE problem for a hypergraph with  $n$  nodes and  $m$  hyperedges each of size  $c$  can be reduced in  $\mathcal{O}(nm)$  time to a  $k$ -PMDM instance with  $\ell = n$ ,  $d = m$  and  $\Sigma = \{a, b\}$ .*

**Proof** We reduce  $(c, k)$ -HYPERCLIQUE to  $k$ -PMDM as follows. Consider the alphabet  $\{a, b\}$ . Set  $q = a^n$ , and for every hyperedge  $e_i \in E$ , add the binary string  $x_i$  to  $\mathcal{D}$  such that  $x_i[j] = b$  if and only if  $j \in e_i$ . Set  $z = \binom{k}{c}$ . Then  $H$  contains a  $(c, k)$ -hyperclique if and only if  $k$ -PMDM returns a positive answer. This can be seen by the fact that  $(c, k)$ -hypercliques  $H$  are in one-to-one correspondence with subsets  $K \subseteq \{1, \dots, n\}$  of size  $k$  for which  $q_K$  matches  $z$  strings from  $\mathcal{D}$ : the elements of  $K$  correspond to the nodes of a  $(c, k)$ -hyperclique and the  $z$  strings correspond to its hyperedges.  $\square$

The  $(c, k)$ -HYPERCLIQUE hypothesis states that there is no  $\mathcal{O}(n^{(1-\epsilon)k})$ -time algorithm, for any  $k > c > 2$  and  $\epsilon > 0$ , that solves the  $(c, k)$ -HYPERCLIQUE problem. For a discussion on the plausibility of this hypothesis and for more context, we refer the reader to [55, Sect. 7]. Theorem 3.5 shows that solving  $k$ -PMDM efficiently even for strings over a binary alphabet would imply a breakthrough for the  $(c, k)$ -HYPERCLIQUE problem. In particular, assuming that the  $(3, k)$ -HYPERCLIQUE hypothesis is true, due to Theorem 3.5, and since  $d\ell = nm = \mathcal{O}(n^4)$ , we cannot hope to devise an algorithm for  $k$ -PMDM requiring time  $\mathcal{O}((d\ell)^{(1-\epsilon)k/4})$  or  $\mathcal{O}(\ell^{(1-\epsilon)k})$ , for any  $k > 3$  and  $\epsilon > 0$ .

## 4 Exact Algorithms for a Bounded Number $k$ of Wildcards

We consider the following problem, which we solve by exact algorithms. These algorithms will form the backbone of our effective and efficient heuristic for the PMDM problem (see Sect. 8).

**HEAVIEST  $k$ -PMDM**

**Input:** A dictionary  $\mathcal{D}$  of  $d$  strings, each of length  $\ell$ , a string  $q$  of length  $\ell$ , and a positive integer  $k \leq \ell$ .

**Output:** A set  $K \subseteq \{1, \dots, \ell\}$  of size  $k$  such that  $q_K = q \otimes K$  matches the maximum number of strings in  $\mathcal{D}$ .

We will show the following result, which we will employ to solve the PMDM problem.

**Theorem 4.1** HEAVIEST  $k$ -PMDM for  $k = \mathcal{O}(1)$  can be solved in  $\mathcal{O}(N + \min\{N^{k/3}, \ell^k\})$  time, where  $N = d\ell$ .

Recall that a *hypergraph*  $H$  is a pair  $(V, E)$ , where  $V$  is the set of nodes of  $H$  and  $E$  is a set of non-empty subsets of  $V$ , called *hyperedges*—in order to simplify terminology we will simply call them edges. Hypergraphs are a generalization of graphs in the sense that an edge can connect more than two nodes. Recall that the size of an edge is the number of nodes it contains. The *rank* of  $H$ , denoted by  $r(H)$ , is the maximum size of an edge of  $H$ .

We refer to a hypergraph  $H[K] = (K, \{e : e \in E, e \subseteq K\})$ , where  $K$  is a subset of  $V$ , as a  $|K|$ -*section*.  $H[K]$  is the hypergraph induced by  $H$  on the nodes of  $K$ , and it contains all edges of  $H$  whose elements are all in  $K$ . A hypergraph is *weighted* when each of its edges is associated with a weight. We define *the weight* of a weighted hypergraph as the sum of the weights of all of its edges. In what follows, we also refer to weights of nodes for conceptual clarity; this is equivalent to having a singleton edge of equal weight consisting of that node.

We define the following auxiliary problem on hypergraphs (see also [21]).

**HEAVIEST  $k$ -SECTION**

**Input:** A weighted hypergraph  $H = (V, E)$ , with  $E$  given as a list, and an integer  $k > 0$ .

**Output:** A subset  $K$  of size  $k$  of  $V$  such that  $H[K]$  has maximum weight.

When  $k = \mathcal{O}(1)$ , we preprocess the edges of  $H$  as follows in order to have  $\mathcal{O}(1)$ -time access to any queried edge. We represent each edge as a string, whose letters correspond to its elements in increasing order. Then, we sort all such strings lexicographically using radix sort in  $\mathcal{O}(|E|)$  time and construct a trie over them. An edge can then be accessed in  $\mathcal{O}(k \log k) = \mathcal{O}(1)$  time by a forward search starting from the root node of the trie.

A polynomial-time  $\mathcal{O}(n^{0.697831+\epsilon})$ -approximation for HEAVIEST  $k$ -SECTION, for any  $\epsilon > 0$ , for the case when all hyperedges of  $H$  have size at most 3 was shown in [21] (see also [5]).

**Instance of HEAVIEST  $k$ -PMDM**

$\mathcal{D}$ : s1 = abcda q = aaaaa  
 s2 = aadba k = 3  
 s3 = acaba d = 6  
 s4 = adaca  $\ell = 5$   
 s5 = bbaac  
 s6 = acdaa

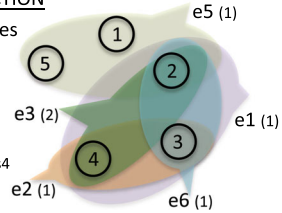
$K = \{2,3,4\}$  of size  $k=3$  has maximum number of matches (5) as  $q_K = a\star\star\star a$  matches s1, s2, s3, s4, and s6.

**Instance of HEAVIEST  $k$ -SECTION**

Hypergraph  $H$  with  $\ell = 5$  nodes  
 a weight-2 edge (e3) and  
 4 weight-1 edges

- e1 = {2,3,4} from s1
- e2 = {3,4} from s2
- e3 = e4 = {2,4} from s3 and s4
- e5 = {1,2,5} from s5
- e6 = {2,3} from s6

$K = \{2,3,4\}$  of size  $k=3$  is s.t.  $H \times K$  has maximum weight 5



**Fig. 2** An example of the reduction from HEAVIEST  $k$ -PMDM to HEAVIEST  $k$ -SECTION. The solutions are at the bottom. Each edge has its weight in brackets and the total weight is  $d = 6$

Two remarks are in place. First, we can focus on edges of size up to  $k$  as larger edges cannot, by definition, exist in any  $k$ -section. Second, HEAVIEST  $k$ -SECTION is a generalization of the problem of deciding whether a  $(c, k)$ -hyperclique (i.e. a set of  $k$  nodes whose subsets of size  $c$  are all in  $E$ ) exists in a graph, which in turn is a generalization of  $k$ -CLIQUE. Unlike  $k$ -CLIQUE, the  $(c, k)$ -hyperclique problem is not known to benefit from fast matrix multiplication in general; see [55] for a discussion on its hardness.

**Lemma 4.2** HEAVIEST  $k$ -PMDM can be reduced to HEAVIEST  $k$ -SECTION for a hypergraph with  $\ell$  nodes and  $d$  edges in  $\mathcal{O}(N)$  time, where  $N = d\ell$ .

**Proof** We first compute the set  $M_s$  of positions of mismatches of  $q$  with each string  $s \in \mathcal{D}$ . We ignore strings from  $\mathcal{D}$  that match  $q$  exactly, as they will match  $q$  after changing any set of letters of  $q$  to wildcards. This requires  $\mathcal{O}(d\ell) = \mathcal{O}(N)$  time in total.

Let us consider an empty hypergraph (i.e. with no edges)  $H$  on  $\ell$  nodes, numbered 1 through  $\ell$ . Then, for each string  $s \in \mathcal{D}$ , we add  $M_s$  to the edge-set of  $H$  if  $|M_s| \leq k$ ; if this edge already exists, we simply increment its weight by 1.

We set the parameter  $k$  of HEAVIEST  $k$ -SECTION to the parameter  $k$  of HEAVIEST  $k$ -PMDM. We now observe that for  $K \subseteq V$  with  $|K| = k$ , the weight of  $H[K]$  is equal to the number of strings that would match  $q$  after replacing with wildcards the  $k$  letters of  $q$  at the positions corresponding to elements of  $K$ . The statement follows.  $\square$

An example of the reduction in Lemma 4.2 is shown in Fig. 2.

The next lemma gives a straightforward solution to HEAVIEST  $k$ -SECTION. It is analogous to algorithm SMALL- $\ell$ , presented in Sect. 6, but without the optimization in computing sums of weights over subsets. It implies a linear-time algorithm for HEAVIEST 1-SECTION.

**Lemma 4.3** HEAVIEST  $k$ -SECTION, for any constant  $k$ , can be solved in  $\mathcal{O}(|V|^k + |E|)$  time and  $\mathcal{O}(|V| + |E|)$  space.

**Proof** For every subset  $K \subseteq V$  of size at most  $k$ , we sum the weights of all edges corresponding to its subsets. There are  $\binom{|V|}{k} = \mathcal{O}(|V|^k)$  choices for  $|K|$ , each having  $2^k - 1$  non-empty subsets: for every subset, we can access the corresponding edge (if it exists) in  $\mathcal{O}(1)$  time.  $\square$

We next show that for the cases  $k = 2$  and  $k = 3$ , there exist more efficient solutions. In particular, we provide a linear-time algorithm for HEAVIEST 2-SECTION.

**Lemma 4.4** HEAVIEST 2-SECTION can be solved in  $\mathcal{O}(|V| + |E|)$  time.

**Proof** Let  $K$  be a set of nodes of size 2 such that  $H[K]$  has maximum weight. We decompose the problem in two cases. For each of the cases, we give an algorithm that considers several 2-sections such that the heaviest of them has weight equal to that of  $H[K]$ .

*Case 1* There is an edge  $e = K$  in  $E$ . For each edge  $e \in E$  of size 2, i.e. edge in the classic sense, we compute the sum of its weight and the weights of the nodes that it is incident to. This step requires  $\mathcal{O}(|E|)$  time.

*Case 2* There is no edge equal to  $K$  in  $E$ . We compute  $H[\{v_1, v_2\}]$ , where  $v_1, v_2$  are the two nodes with maximum weight, i.e. max and second-max. This step takes  $\mathcal{O}(|V|)$  time.

In the end, we return the heaviest 2-section among those returned by the algorithms for the two cases, breaking ties arbitrarily.  $\square$

We next show that for  $k = 3$  the result of Lemma 4.3 can be improved when  $|E| = o(|V|^2)$ .

**Lemma 4.5** HEAVIEST 3-SECTION can be solved in time  $\mathcal{O}(|V| \cdot |E|)$  using  $\mathcal{O}(|V| + |E|)$  space.

**Proof** Let  $K$  be a set of nodes of size 3 such that  $H[K]$  has maximum weight. We decompose the problem into the following three cases.

*Case 1* There is an edge  $e = K$  in  $E$ . We go through each edge  $e \in E$  of size 3 and compute the weight of  $H[e]$  in  $\mathcal{O}(1)$  time. This takes  $\mathcal{O}(|E|)$  time in total. Let the edge yielding the maximum weight be  $e_{\max}$ .

*Case 2* There is no edge of size larger than one in  $H[K]$ . We compute  $H[\{v_1, v_2, v_3\}]$ , where  $v_1, v_2, v_3$  are the three nodes with maximum weight, i.e. max, second-max and third-max. This step takes  $\mathcal{O}(|V|)$  time.

*Case 3* There is an edge of size 2 in  $H[K]$ . We can pick an edge  $e$  of size 2 from  $E$  in  $\mathcal{O}(|E|)$  ways and a node  $v$  from  $V$  in  $\mathcal{O}(|V|)$  ways. We compute the weight of  $H[(e \cup \{v\})]$  for all such pairs. Let the pair yielding maximum weight be  $(e', u')$ .

Finally, the maximum weight of  $H[K']$  for  $K' \in \{e_{\max}, \{v_1, v_2, v_3\}, e' \cup \{u'\}\}$  is equal to the weight of  $H[K]$ , breaking ties arbitrarily.  $\square$

We next address the remaining case of any arbitrarily large constant  $k \geq 4$ .

**Lemma 4.6** HEAVIEST  $k$ -SECTION for an arbitrarily large constant  $k \geq 4$  can be solved in time  $\mathcal{O}((|V| \cdot |E|)^{k/3})$  using  $\mathcal{O}(|V| + |E|)$  space.

**Proof** If  $|E| > |V|^2$ , then the simple algorithm of Lemma 4.3 solves the problem in time

$$\mathcal{O}(|V|^k + |E|) = \mathcal{O}\left(|V|^{k/3} \left(|V|^2\right)^{k/3} + |E|\right) = \mathcal{O}\left((|V| \cdot |E|)^{k/3}\right)$$

and linear space. We can thus henceforth assume that  $|E| \leq |V|^2$ .

Let  $K$  be a set of nodes of size at most  $k$  such that  $H[K]$  has maximum weight. If  $H[K]$  contains isolated nodes (i.e. nodes not contained in any edge), they can be safely deleted without altering the result. We can thus assume that  $H[K]$  does not contain isolated nodes, and that  $|V| \leq k|E|$  since otherwise the hypergraph  $H$  would contain isolated nodes.

We first consider the case that the rank  $r(H[K]) > 1$ , i.e. there is an edge of  $H[K]$  of size at least 2. We design a branching algorithm that constructs several candidate sets; the ones with maximum weight will have weight equal to that of  $H[K]$ . We will construct a set of nodes  $X$ , starting with  $X := \emptyset$ . For each set  $X$  that we process, let  $Z_X$  be the superset of  $X$  of size at most  $k$  such that  $H[Z_X]$  has maximum weight. We have the following two cases:

*Case 1* There is an edge  $e$  in  $H[Z_X]$  that contains at least two nodes from  $Z_X \setminus X$ . To account for this case, we select every possible such edge  $e$ , set  $X := X \cup e$ , and continue the branching algorithm.

*Case 2* Each edge in  $H[Z_X]$  contains at most one node from  $Z_X \setminus X$ . In this case we conclude the branching algorithm as follows. For every node  $v \in V \setminus X$  we compute its weight as the total weight of edges  $Y \cup \{v\} \in E$  for  $Y \subseteq X$  in  $\mathcal{O}(2^k) = \mathcal{O}(1)$  time. Finally, in  $\mathcal{O}(|V|k) = \mathcal{O}(|V|)$  time we select  $k - |X|$  nodes with largest weights and insert them into  $X$ . The total time complexity of this step is  $\mathcal{O}(|V|)$ . This case also works if  $|X| = k$  and then its time complexity is only  $\mathcal{O}(1)$ .

The correctness of this branching algorithm follows from an easy induction, showing that at every level of the branching tree there is a subset of  $K$ .

Let us now analyze the time complexity of this branching algorithm. Each branching in Case 1 takes  $\mathcal{O}(|E|)$  time and increases the size of  $|X|$  by at least 2. At every node of the branching tree we call the procedure of Case 2. It takes  $\mathcal{O}(|V|)$  time if  $|X| < k$ .

If the procedure of Case 2 is called in a non-leaf node of the branching tree, then its  $\mathcal{O}(|V|)$  running time is dominated by the  $\mathcal{O}(|E|)$  time that is required for further branching since we have assumed that  $|V| \leq k|E|$ . Hence, it suffices to bound (a) the total time complexity of calls to the algorithm for Case 2 in leaves that correspond to sets  $X$  such that  $|X| < k$  and (b) the total number of leaves that correspond to sets  $X$  such that  $|X| = k$ .

If  $k$  is even, (a) is bounded by  $\mathcal{O}(|E|^{(k-2)/2}|V|)$  and (b) is bounded by  $\mathcal{O}(|E|^{k/2})$ . Hence, (b) dominates (a) and we have

$$\mathcal{O}(|E|^{k/2}) = \mathcal{O}(|E|^{k/3}|E|^{k/6}) = \mathcal{O}(|E|^{k/3}|V|^{k/3}). \tag{1}$$

If  $k$  is odd, (a) is bounded by  $\mathcal{O}(|E|^{(k-1)/2}|V|)$  and (b) is bounded by  $\mathcal{O}(|E|^{(k-1)/2})$ , which is dominated by (a). By using (1) for  $k - 3$  we also have:

$$\begin{aligned} \mathcal{O}(|E|^{(k-1)/2} \cdot |V|) &= \mathcal{O}(|E|^{(k-3)/2} \cdot |E| \cdot |V|) \\ &= \mathcal{O}((|E| \cdot |V|)^{(k-3)/3} \cdot |E| \cdot |V|) = \mathcal{O}((|E| \cdot |V|)^{k/3}). \end{aligned}$$

We now consider the case that  $r(H[K]) = 1$ . We use the algorithm for Case 2 above that works in  $\mathcal{O}(|V|)$  time, which is  $\mathcal{O}(|V| \cdot |E|)$ .  $\square$

Lemmas 4.2–4.6 imply Theorem 4.1, which we iteratively employ to obtain the following result.

**Theorem 4.7** *PMDM can be solved in time  $\mathcal{O}(N + \min\{N^{k/3}, \ell^k\})$  using space  $\mathcal{O}(N)$  if  $k = \mathcal{O}(1)$ , where  $N = d\ell$ .*

**Proof** We apply Lemma 4.2 to obtain a hypergraph with  $|V| = \ell$  and  $|E| = d$ . Starting with  $k = 1$  and for growing values of  $k$ , we solve HEAVIEST  $k$ -SECTION until we obtain a solution of weight at least  $z$ , employing either only Lemma 4.3, or Lemmas 4.3, 4.4, 4.5, 4.6 for  $k = 1, 2, 3$  and  $k \geq 4$ , respectively. We obtain  $\mathcal{O}(N + \min\{N^{k/3}, \ell^k\})$  time and  $\mathcal{O}(N)$  space.  $\square$

### 5 Exact Algorithms for a Bounded Number $m$ of Query Strings

Recall that masking a potential match  $(q_1, q_2)$  in record linkage can be performed by solving PMDM twice and superimposing the wildcards (see Sect. 1). In this section, we consider the following generalized version of PMDM to perform the masking simultaneously. The advantage of this approach is that it minimizes the final number of wildcards in  $q_1$  and  $q_2$ .

**MULTIPLE PATTERN MASKING FOR DICTIONARY MATCHING (MPMDM)**

**Input:** A dictionary  $\mathcal{D}$  of  $d$  strings, each of length  $\ell$ , a collection  $\mathcal{M}$  of  $m$  strings, each of length  $\ell$ , and a positive integer  $z$ .

**Output:** A smallest set  $K \subseteq \{1, \dots, \ell\}$  such that, for every  $q$  from  $\mathcal{M}$ ,  $q_K = q \otimes K$  matches at least  $z$  strings from  $\mathcal{D}$ .

Let  $N = d\ell$ . We show the following theorem.

**Theorem 5.1** *MPMDM can be solved in time  $\mathcal{O}(N + \min\{N^{k/3}z^{m-1}, \ell^k\})$  if  $k = \mathcal{O}(1)$  and  $m = \mathcal{O}(1)$ , where  $N = d\ell$ .*

We use a generalization of HEAVIEST  $k$ -SECTION in which the weights are  $m$ -tuples that are added and compared component-wise, and we aim to find a subset  $K$  such that the weight of  $H[K]$  is at least  $(z, \dots, z)$ . An analogue of Lemma 4.3 holds without any alterations, which accounts for the  $\mathcal{O}(N + \ell^k)$ -time algorithm. We adapt the proof of Lemma 4.6 as follows. The branching remains the same, but we have to tweak the final step, that is, what happens when we are in Case 2. For  $m = 1$  we could simply select a number of largest weights, but for  $m > 1$  multiple criteria need to be taken into consideration. All in all, the problem reduces to a variation of the classic Multiple-Choice Knapsack problem [47], which we solve using dynamic programming.

The variation of the classic Multiple-Choice Knapsack problem is as follows.

$\kappa$  HEAVIEST VECTORS ( $\kappa$ -HV)

**Input:** A collection  $\mathcal{T}$  of  $t$  vectors from  $\mathbb{Z}_{\geq 0}^m$ , a vector  $x$  from  $\{0, \dots, z\}^m$ , for a positive integer  $z$ , and an integer  $\kappa \in \{0, \dots, t\}$ .

**Output:** Compute  $\kappa$  elements of  $\mathcal{T}$  (if they exist) such that if  $y$  is their component-wise sum,  $y[i] \geq x[i]$  for all  $i \in \{1, \dots, m\}$ .

The exact reduction from Case 2 is as follows: the set  $\mathcal{T}$  contains weights of subsequent nodes  $v \in V \setminus X$  (defined as the sums of weights of edges  $Y \cup \{v\} \in E$  for  $Y \subseteq X$ ), so  $t \leq |V|$ ,  $x$  is  $(z, \dots, z)$  minus the sum of weights of all edges  $e \in E$  such that  $e \subseteq X$ , and  $\kappa = k - |X|$ .

The solution to  $\kappa$ -HV is a rather straightforward dynamic programming.

**Lemma 5.2** For  $\kappa, m = \mathcal{O}(1)$ ,  $\kappa$ -HV can be solved in time  $\mathcal{O}(t \cdot z^{m-1})$ .

**Proof** We apply dynamic programming. Let  $\mathcal{T} = v_1, \dots, v_t$ . We compute an array  $A$  of size  $\mathcal{O}(t\kappa z^{m-1})$  such that, for  $i \in \{0, \dots, t\}$ ,  $j \in \{0, \dots, \kappa\}$  and  $v \in \{0, \dots, z\}^{m-1}$ ,

$$A[i, j, v] = \max\{a : \exists S \subseteq \{v_1, \dots, v_i\}, |S| = j, \sum_{u \in S} u = (v, a)\},$$

where  $(v, a)$  denotes the operation of appending element  $a$  to vector  $v$ . From each state  $A[i, j, v]$  we have two transitions, depending on whether  $v_{i+1}$  is taken to the subset or not. Each transition is computed in  $\mathcal{O}(m) = \mathcal{O}(1)$  time. This gives time  $\mathcal{O}(t \cdot z^{m-1})$  in total.

The array is equipped with a standard technique to recover the set  $S$  (parents of states). The final answer is computed by checking, for each vector  $v \in \{0, \dots, z\}^{m-1}$  such that  $v[i] \geq x[i]$ , for all  $i = 1, \dots, m - 1$ , if  $A[t, \kappa, v] \geq x[m]$ . □

Overall, we pay an additional  $\mathcal{O}(z^{m-1})$  factor in the complexity of handling of Case 2, which yields the complexity of Theorem 5.1.

## 6 A Data Structure for PMDM Queries

We next show algorithms and data structures for the PMDM problem under the assumption that  $2^\ell$  is reasonably small. We measure space in terms of  $w$ -bit machine words, where  $w = \Omega(\log(d\ell))$ , and focus on showing space versus query-time trade-offs for answering  $q, z$  PMDM queries over  $\mathcal{D}$ . A summary of the complexities of the data structures is shown in Table 1. Specifically, algorithm SMALL- $\ell$  and data structure SIMPLE are used as building blocks in the more involved data structure SPLIT underlying the following theorem.

**Theorem 6.1** There exists a data structure that answers  $q, z$  PMDM queries over  $\mathcal{D}$  in time  $\mathcal{O}(2^{\ell/2}(2^{\ell/2} + \tau)\ell)$  w.h.p. and requires space  $\mathcal{O}(2^\ell d^2/\tau^2 + 2^{\ell/2}d)$ , for any  $\tau \in [1, d]$ .

**Algorithm SMALL- $\ell$ :**  $\mathcal{O}(d\ell)$  Space,  $\mathcal{O}(2^\ell \ell + d\ell)$  Query Time.

Our algorithm is based on the Fast zeta/Möbius transform [28, Theorem 10.12]. No data structure on top of the dictionary  $\mathcal{D}$  is stored. In the query algorithm, we initialize an integer array  $A$  of size  $2^\ell$  with zeros. For an  $\ell$ -bit vector  $m$ , by  $K_m \subseteq \{1, \dots, \ell\}$  let us denote the set of the positions of set bits of  $m$ . Now for every possible  $\ell$ -bit vector  $m$  we want to compute the number of strings in  $\mathcal{D}$  that match  $q_{K_m} = q \otimes K_m$ .

To this end, for every string  $s \in \mathcal{D}$ , we compute the set  $K$  of positions in which  $s$  and  $q$  differ. For  $m$  that satisfies  $K = K_m$ , we increment  $A[m]$ , where  $m$  is the integer representation of the bit vector. This computation takes  $\mathcal{O}(d\ell)$  time and  $\mathcal{O}(1)$  extra space. Then we apply a folklore dynamic-programming-based approach to compute an integer array  $B$ , which is defined as follows:

$$B[m] = \sum_{j \in S(m)} A[j], \text{ where } S(m) = \{j \in [1, 2^\ell] : K_j \subseteq K_m\}.$$

In other words,  $B[m]$  stores the number of strings from  $\mathcal{D}$  that match  $q_{K_m}$ .

We provide a description of the folklore algorithm here for completeness. Consider a vector (mask)  $m$ . Let  $S(m, i)$  consist of the subsets of  $m$  which do not differ from  $m$  but (possibly) in the rightmost  $i$  bits, and

$$B[m, i] = \sum_{j \in S(m, i)} A[j].$$

Clearly,  $S(m)$  is equal to  $S(m, \ell)$ , and hence  $B[m]$  is equal to  $B[m, \ell]$ . The following equation is readily verified (in the first case, since the  $i$ th bit of  $m$  is 0, no element of  $S(m)$  can have the  $i$ th bit set):

$$B[m, i] = \begin{cases} B[m, i - 1] & \text{if the } i\text{th bit of } m \text{ is 0,} \\ B[m, i - 1] \text{ OR } B[m \text{ XOR } 2^i, i - 1] & \text{if the } i\text{th bit of } m \text{ is 1.} \end{cases}$$

By OR and XOR we denote the standard bitwise operations. Overall, there are  $\mathcal{O}(2^\ell \ell)$  choices for  $m$  and  $i$ . We can compute  $B[\cdot, \cdot]$  column by column, in constant time per entry, thus obtaining an  $\mathcal{O}(2^\ell \ell)$ -time algorithm. We can limit the space usage to  $\mathcal{O}(2^\ell)$  by discarding column  $i$  when we are done computing column  $i + 1$ .

Thus, overall, the (query) time required by algorithm SMALL- $\ell$  is  $\mathcal{O}(\ell 2^\ell + d\ell)$ , the data structure space is  $\mathcal{O}(d\ell)$ , and the extra space is  $\mathcal{O}(2^\ell)$ .

We now present SIMPLE, an auxiliary data structure, which we will apply later on to construct DS SPLIT, a data structure with the space/query-time trade-off of Theorem 6.1.

**DS SIMPLE:**  $\mathcal{O}(2^\ell d)$  Space,  $\mathcal{O}(2^\ell \ell)$  Query Time. We initialize an empty set  $\mathcal{Q}$ . For each possible subset of  $\{1, \dots, \ell\}$  we do the following. We mask the corresponding positions in all strings from  $\mathcal{D}$  and then sort the masked strings lexicographically. By iterating over the lexicographically sorted list of the masked strings, we count how many copies of each distinct (masked) string we have in our list. We insert each such (masked) string to  $\mathcal{Q}$  along with its count. After processing all  $2^\ell$  subsets, we construct



**Table 1** Basic complexities of the data structures from Sect. 6

Data structure	Space	Query time
Algorithm SMALL- $\ell$	$\mathcal{O}(d\ell)$	$\mathcal{O}(2^\ell \ell + d\ell)$
DS SIMPLE	$\mathcal{O}(2^\ell d)$	$\mathcal{O}(2^\ell \ell)$
DS SPLIT, any $\tau$	$\mathcal{O}(2^\ell d^2/\tau^2 + 2^{\ell/2}d)$	$\mathcal{O}(2^{\ell/2} \cdot (2^{\ell/2} + \tau)\ell)$
DS SPLIT for $\tau = 2^{\ell/4}\sqrt{d}$	$\mathcal{O}(2^{\ell/2}d)$	$\mathcal{O}(2^\ell \ell + 2^{3\ell/4}\sqrt{d}\ell)$

a compacted trie for the strings in  $\mathcal{Q}$ ; each leaf corresponds to a unique element of  $\mathcal{Q}$ , and stores this element’s count. The total space occupied by this compacted trie is thus  $\mathcal{O}(2^\ell d)$ . Upon an on-line query  $q$  (of length  $\ell$ ) and  $z$ , we apply all possible  $2^\ell$  masks to  $q$  and read the count for each of them from the compacted trie in  $\mathcal{O}(\ell)$  time per mask. Next, we show how to decrease the exponential dependency on  $\ell$  in the space complexity when  $2^\ell = o(d)$ , incurring extra time in the query.

**DS SPLIT:**  $\mathcal{O}(2^\ell d^2/\tau^2 + 2^{\ell/2}d)$  Space,  $\mathcal{O}(2^{\ell/2} \cdot (2^{\ell/2} + \tau)\ell)$  Query Time, for any  $\tau$ . This trade-off is relevant when  $\tau = \omega(\sqrt{d})$ ; otherwise the DS SIMPLE is better.

We split each string  $p \in \mathcal{D}$  roughly in the middle, to prefix  $p_L$  and suffix  $p_R$ ; specifically,  $p = p_L p_R$  and  $|p_L| = \lceil \ell/2 \rceil$ . We create dictionaries  $\mathcal{D}_L = \{p_L : p \in \mathcal{D}\}$  and  $\mathcal{D}_R = \{p_R : p \in \mathcal{D}\}$ . Let us now explain how to process  $\mathcal{D}_L$ ; we process  $\mathcal{D}_R$  analogously. Let  $\lambda = \lceil \ell/2 \rceil$ . We construct DS SIMPLE over  $\mathcal{D}_L$ . This requires space  $\mathcal{O}(2^{\ell/2}d)$ . Let  $\tau$  be an input parameter, intuitively used as the minimum frequency threshold. For each of the possible  $2^\lambda$  masks, we can have at most  $\lfloor d/\tau \rfloor$  (masked) strings with frequency at least  $\tau$ . Over all masks, we thus have at most  $2^\lambda \lfloor d/\tau \rfloor$  such strings, which we call  $\tau$ -frequent. For every pair of  $\tau$ -frequent strings, one from  $\mathcal{D}_L$  and one from  $\mathcal{D}_R$ , we store the number of occurrences of their concatenation in  $\mathcal{Q}$  using a compacted trie as in DS SIMPLE. This requires space  $\mathcal{O}(2^\ell d^2/\tau^2)$ .

Consider  $\mathcal{D}_L$ . For each mask  $i$  and each string  $p_L \in \mathcal{D}_L$ , we can afford to store the list of all strings in  $\mathcal{D}_L$  that match  $p_L \otimes i$ . Note that we have computed this information when sorting for constructing DS SIMPLE over  $\mathcal{D}_L$ . This information requires space  $\mathcal{O}(2^{\ell/2}d)$ . Thus, DS SPLIT requires  $\mathcal{O}(2^\ell d^2/\tau^2 + 2^{\ell/2}d)$  space overall.

Let us now show how to answer an on-line  $q, z$  query. Let  $q = q_L q_R$  with  $|q_L| = \lceil \ell/2 \rceil$ . We iterate over all possible  $2^\ell$  masks.

For a mask  $i$ , let  $q' = q \otimes i$ . We split  $q'$  into two halves,  $q'_L$  and  $q'_R$  with  $q' = q'_L q'_R$  and  $|q'_L| = \lceil \ell/2 \rceil$ . First, we check whether each of  $q'_L$  and  $q'_R$  is  $\tau$ -infrequent using the DS SIMPLE we have constructed for  $\mathcal{D}_L$  and  $\mathcal{D}_R$ , respectively, in time  $\mathcal{O}(\ell)$ . We have the following two cases (inspect also Fig. 3).

- If both halves are  $\tau$ -frequent, then we can read the frequency of their concatenation using the stored compacted trie in time  $\mathcal{O}(\ell)$ .
- Else, at least one of the two halves is  $\tau$ -infrequent. Assume without loss of generality that  $q'_L$  is  $\tau$ -infrequent. Let  $\mathcal{F}$  be the dictionary consisting of at most  $\tau$  strings from  $\mathcal{D}_R$  that correspond to the right halves of strings in  $\mathcal{D}_L$  that match  $q'_L$ . Naïvely counting how many elements of  $\mathcal{F}$  match  $q'_R$  could require  $\Omega(\tau \ell)$  time, and thus  $\Omega(2^\ell \tau \ell)$  overall. Instead, we apply algorithm SMALL- $\ell$  on  $q_R$  and

$\mathcal{D}$	$\mathcal{D}_L$	$\mathcal{D}_R$	Counts of SIMPLE		Compacted trie of SPLIT	
			$\mathcal{D}_L$	$\mathcal{D}_R$		
abba	ab	ba	a* 4	b* 2	***	a**
acba	ac	ba	*b 2	*a 3		
acca	ac	ca	*c 2	...		
abac	ab	ac	** 4	** 4	*	a

**Fig. 3** Let  $\tau = 3$ . If both  $q'_L$  and  $q'_R$  are 3-frequent (we check this using the counts of DS SIMPLE), then we read the count for  $q'_L q'_R$  from the compacted trie of DS SPLIT. If  $q'_L$  is 3-infrequent, then we apply SMALL- $\ell$  on  $q_R$  and on the dictionary consisting of at most  $\tau = 3$  strings from  $\mathcal{D}_R$  corresponding to the right halves of strings in  $\mathcal{D}_L$  that match  $q'_L$

**Table 2** Basic complexities of the data structures from Sect. 6 for  $d = 2^{2\ell}$

Data structure	Space	Query time
Algorithm SMALL- $\ell$	$\mathcal{O}(2^{2\ell}\ell)$	$\mathcal{O}(2^{2\ell}\ell)$
DS SIMPLE	$\mathcal{O}(2^{3\ell})$	$\mathcal{O}(2^\ell\ell)$
DS SPLIT for $\tau = 2^{5\ell/4}$	$\mathcal{O}(2^{5\ell/2})$	$\mathcal{O}(2^{7\ell/4}\ell)$

$\mathcal{F}$ . The crucial point is that if we ever come across  $q'_L$  again (for a different mask on  $q$ ), we will not need to do anything. We can maintain whether  $q'_L$  has been processed by temporarily marking the leaf corresponding to it in DS SIMPLE for  $\mathcal{D}_L$ . Thus, overall, we perform the SMALL- $\ell$  algorithm  $\mathcal{O}(2^{\ell/2})$  times, each time in  $\mathcal{O}((2^{\ell/2} + \tau)\ell)$  time. This completes the proof of Theorem 6.1.

**Efficient Construction** For completeness, we next show how to construct DS SPLIT in  $\mathcal{O}(d\ell \log(d\ell) + 2^\ell d\ell + 2^\ell \ell d^2/\tau^2)$  time. We preprocess  $\mathcal{D}$  by sorting its letters in  $\mathcal{O}(d\ell \log(d\ell))$  time. The DS SIMPLE for  $\mathcal{D}_L$  and  $\mathcal{D}_R$  can then be constructed in  $\mathcal{O}(2^{\ell/2}d\ell)$  time. We then create the compacted trie for pairs of  $\tau$ -frequent strings. For each of the  $2^\ell$  possible masks, say  $i$ , and each string  $p \in \mathcal{D}$ , we split  $p' = p \otimes i$  in the middle to obtain  $p'_L$  and  $p'_R$ . If both  $p'_L$  and  $p'_R$  are  $\tau$ -frequent then  $p'$  will be in the set of strings for which we will construct the compacted trie for pairs of  $\tau$ -frequent strings. The counts for each of those strings can be read in  $\mathcal{O}(\ell)$  time from a DS SIMPLE over  $\mathcal{D}$ , which we can construct in time  $\mathcal{O}(2^\ell d\ell)$ —this data structure is then discarded. The compacted trie construction requires time  $\mathcal{O}(2^\ell \ell d^2/\tau^2)$ .

**Comparison of the Data Structures** DS SIMPLE has lower query time than algorithm SMALL- $\ell$ . However, its space complexity can be much higher. DS SPLIT can be viewed as an intermediate option. For  $\tau$  as in Table 1, it has lower query time than algorithm SMALL- $\ell$  for  $d = \omega(2^{3\ell/2})$ , while keeping moderate space complexity. DS SPLIT always has higher query time than DS SIMPLE, but its space complexity is lower by a factor of  $2^{\ell/2}$ . For example, for  $d = 2^{2\ell}$  we get the complexities shown in Table 2.

Let us now discuss why our data structure results cannot be directly obtained using the same data structures as for the problem Dictionary Matching with  $k$ -wildcards

(see Sect. 1 for the problem definition). Conceivably, one could construct such a data structure, and then iterate over all subsets of  $\{1, \dots, \ell\}$ , querying for the masked string. Existing data structures for dictionary matching with wildcards (cf. [12, Table 1], [52], and [34]), that allow querying a pattern with at most  $\ell$  wildcards, have:

- (a) Either  $\Omega(\min\{\sigma^\ell, d\})$  query time, thus yielding  $\Omega(2^\ell \cdot \min\{\sigma^\ell, d\})$  query time for our problem, and space  $\Omega(d\ell)$ , a trade-off dominated by the SMALL- $\ell$  algorithm (cf. our Table 1);
- (b) Or  $\Omega(\ell)$  query time, thus yielding  $\Omega(2^\ell \ell)$  query time for our problem, and  $\Omega(d\ell \log^\ell \log(d\ell))$  space, a trade-off dominated by the DS SIMPLE (cf. our Table 1).

### 7 Approximation Algorithm for PMDM

Clearly, PMDM is at least as hard as PMDM-SIZE because it also outputs the positions of the wildcards (set  $K$ ). Thus, PMDM is also NP-hard. In what follows, we show existence of a polynomial-time approximation algorithm for PMDM whose approximation factor is given with respect to  $d$ . Specifically, we show the following approximation result for PMDM.

**Theorem 7.1** *For any constant  $\epsilon > 0$ , there is an  $\mathcal{O}(d^{1/4+\epsilon})$ -approximation algorithm for PMDM, whose running time is polynomial in  $N$ , where  $N = d\ell$ .*

Our result is based on a reduction to the Minimum Union (MU) problem [20], which we define next.

**MINIMUM UNION (MU)**

**Input:** A collection  $\mathcal{S}$  of  $d$  sets over a universe  $U$  and a positive integer  $z \leq d$ .

**Output:** A collection  $\mathcal{T} \subseteq \mathcal{S}$  with  $|\mathcal{T}| = z$  such that the size of  $\cup_{S \in \mathcal{T}} S$  is minimized.

To illustrate the MU problem, consider an instance of it where  $U = \{1, 2, 3, 4, 5\}$ ,  $\mathcal{S} = \{\{1\}, \{1, 2, 3\}, \{1, 3, 5\}, \{3\}, \{3, 4, 5\}, \{4\}, \{4, 5\}, \{5\}\}$ , with  $d = |\mathcal{S}| = 8$ , and  $z = 4$ . Then  $\mathcal{T} = \{\{3\}, \{3, 4, 5\}, \{4\}, \{4, 5\}\}$  is a solution because  $|\mathcal{T}| = z = 4$  and  $|\cup_{S \in \mathcal{T}} S| = 3$  is minimum. The MU problem is NP-hard and the following approximation result is known.

**Theorem 7.2** ([20]) *For any constant  $\epsilon > 0$ , there is an  $\mathcal{O}(d^{1/4+\epsilon})$ -approximation algorithm for MU, whose running time is polynomial in the size of  $\mathcal{S}$ .*

We next describe the reduction that leads to our result.

**Theorem 7.3** *PMDM can be reduced to MU in time polynomial in  $N$ .*

**Proof** We reduce the PMDM problem to MU in polynomial time as follows. Given any instance  $\mathcal{I}_{\text{PMDM}}$  of PMDM, we construct an instance  $\mathcal{I}_{\text{MU}}$  of MU in time  $\mathcal{O}(d\ell)$  by performing the following steps:

1. The universe  $U$  is set to  $\{1, \dots, \ell\}$ .

2. We start with an empty collection  $\mathcal{S}$ . Then, for each string  $s_i$  in  $\mathcal{D}$ , we add member  $S_i$  to  $\mathcal{S}$ , where  $S_i$  is the set of positions where string  $q$  and string  $s_i$  have a mismatch. This can be done trivially in time  $\mathcal{O}(d\ell)$  for all strings in  $\mathcal{D}$ .
3. Set the  $z$  of the MU problem to the  $z$  of the PMDM problem.

Thus, the total time  $\mathcal{O}(d\ell)$  needed for Steps 1 to 3 above is clearly polynomial in the size of  $\mathcal{I}_{\text{PMDM}}$ .

**Claim** For any solution  $\mathcal{T}$  to  $\mathcal{I}_{\text{MU}}$  and any solution  $K$  to  $\mathcal{I}_{\text{PMDM}}$ , such that  $\mathcal{I}_{\text{MU}}$  is obtained from  $\mathcal{I}_{\text{PMDM}}$  using the above three steps, we have  $|K| = |\cup_{S \in \mathcal{T}} S|$ .

**Proof of Claim** Let  $\mathcal{F} \subseteq \mathcal{D}$  consist of  $z$  strings that match  $q_K$ . Further, let the set  $\mathcal{F}^*$  consist of the elements of  $\mathcal{S}$  corresponding to strings in  $\mathcal{F}$ . We have  $|\cup_{S \in \mathcal{T}} S| \leq |\cup_{S \in \mathcal{F}^*} S| \leq |K|$ .

Now, let  $C = \cup_{S \in \mathcal{T}} S$ . Then,  $q_C = q \otimes C$  matches at least  $z$  strings from  $\mathcal{D}$  and hence  $|K| \leq |C| = |\cup_{S \in \mathcal{T}} S|$ . □

To conclude the proof, it remains to show that given a solution  $\mathcal{T}$  to  $\mathcal{I}_{\text{MU}}$  we can obtain a solution  $K$  to  $\mathcal{I}_{\text{PMDM}}$  in time polynomial in the size of  $\mathcal{I}_{\text{MU}}$ . This readily follows from the proof of the above claim: it suffices to set  $K = \cup_{S \in \mathcal{T}} S$ . □

We can now prove Theorem 7.1.

**Proof of Theorem 7.1** The reduction in Theorem 7.3 implies that there is a polynomial-time approximation algorithm for PMDM. In particular, Theorem 7.2 provides an approximation guarantee for MU that depends on the number of sets of the input  $\mathcal{S}$ . In Step 2 of the reduction of Theorem 7.3, we construct *one* set for the MU instance per *one* string of the dictionary  $\mathcal{D}$  of the PMDM instance. Also, from the constructed solution  $\mathcal{T}$  to the MU instance, we obtain a solution  $K$  to the PMDM instance by simply substituting the positions of  $q$  corresponding to the elements of the sets of  $\mathcal{T}$  with wildcards. This construction implies the approximation result of Theorem 7.1 that depends on the size of  $\mathcal{D}$ . □

Applying Theorem 7.1 to solve PMDM is not practical, as in real-world applications, such as those in Sect. 1,  $d$  is typically in the order of thousands or millions [24, 30, 35, 71].

**Sanity Check** We remark that Theorem 3.1 (reduction from  $k$ -CLIQUE to  $k$ -PMDM) and Theorem 7.1 (approximation algorithm for PMDM) do not contradict the inapproximability results for the maximum clique problem (see Sect. 3), since our reduction from  $k$ -CLIQUE to  $k$ -PMDM cannot be adapted to a reduction from maximum clique to PMDM-SIZE.

**Two-Way Reduction** Chlamtáč et al. [20] also show that their polynomial-time  $\mathcal{O}(d^{1/4+\epsilon})$ -approximation algorithm for MU is tight under a plausible conjecture for the so-called Hypergraph Dense vs Random problem. In what follows, we also show that approximating the MU problem can be reduced to approximating PMDM in polynomial time and hence the same tightness result applies to PMDM.

**Theorem 7.4** MU can be reduced to PMDM in time polynomial in the size of  $\mathcal{S}$ .

**Proof** Let  $||\mathcal{S}||$  denote the total number of elements in the  $d$  members of  $\mathcal{S}$ . We reduce the MU problem to the PMDM problem in polynomial time as follows. Given any instance  $\mathcal{I}_{\text{MU}}$  of MU, we construct an instance  $\mathcal{I}_{\text{PMDM}}$  of PMDM by performing the following steps:

1. Sort the union of all elements of members of  $\mathcal{S}$ , assign to each element  $j$  a unique rank  $\text{rank}(j) \in \{1, \dots, |U|\}$ , and set  $\ell = |U|$ . This can be done in  $\mathcal{O}(|\mathcal{S}| \log |\mathcal{S}|)$  time.
2. Set the query string  $q$  equal to the string  $a^\ell$  of length  $\ell$ . For each set  $S_i$  in  $\mathcal{S}$ , construct a string  $s_i = a^\ell$ , set  $s_i[\text{rank}(j)] := b$  if and only if  $j \in S_i$ , and add  $s_i$  to dictionary  $\mathcal{D}$ . This can be done in  $\mathcal{O}(d\ell)$  time.
3. Set the  $z$  of the PMDM problem equal to the  $z$  of the MU problem. This can be done in  $\mathcal{O}(1)$  time.

Thus, the total time  $\mathcal{O}(d\ell \log(d\ell))$  needed for Steps 1 to 3 above is clearly polynomial in the size of  $\mathcal{I}_{\text{MU}}$  as  $\ell \leq ||\mathcal{S}||$ .

A proof of the following claim is analogous to that of Claim 7.

**Claim** For any solution  $\mathcal{T}$  to  $\mathcal{I}_{\text{MU}}$  and any solution  $K$  to  $\mathcal{I}_{\text{PMDM}}$ , such that  $\mathcal{I}_{\text{PMDM}}$  is obtained from  $\mathcal{I}_{\text{MU}}$  using the above three steps, we have  $|K| = |\cup_{S \in \mathcal{T}} S|$ .

To conclude the proof, it remains to show that, given a solution  $K$  to  $\mathcal{I}_{\text{PMDM}}$ , we can obtain a solution  $\mathcal{T}$  to  $\mathcal{I}_{\text{MU}}$  in time polynomial in the size of  $\mathcal{I}_{\text{PMDM}}$ . It suffices to pick  $z$  sets in  $\mathcal{S}$  that are subsets of  $K$ . Their existence is guaranteed by construction, because such sets correspond to the at least  $z$  strings in  $\mathcal{D}$  that have  $b$  in a subset of the positions in  $K$ . This selection can be done naively in  $\mathcal{O}(|\mathcal{S}|)$  time. Finally, the above claim guarantees that they indeed form a solution to  $\mathcal{I}_{\text{MU}}$ . □

## 8 A Greedy Heuristic for PMDM

We design a heuristic called GREEDY  $\tau$ -PMDM that solves PMDM and which, for a given constant  $\tau \geq 1$ , iteratively applies Theorem 4.1 (see Sect. 4), for  $k = 1, \dots, \tau$ . Intuitively, the larger the  $\tau$ , the more effective—but the slower—GREEDY  $\tau$ -PMDM is. Specifically, in iteration  $i = 1$ , we apply Theorem 4.1 for  $k = 1, \dots, \tau$  and check whether there are at least  $z$  strings from  $\mathcal{D}$  that can be matched when at most  $k$  wildcards are substituted in the query string  $q$ . If there are, we return the minimum such  $k$  and terminate. Clearly, by Theorem 4.7, the returned solution  $K_1$  is an optimal solution to PMDM. Otherwise, we proceed into the next iteration,  $i = 2$ . In this iteration, we construct string  $q_{K_1} = q \otimes K_1$  and apply Theorem 4.1, for  $k = 1, \dots, \tau$ , to  $q_{K_1}$ . This returns a solution  $K_2$  telling us whether there are at least  $z$  strings from  $\mathcal{D}$  that can be matched with  $q_{K_2} = q_{K_1} \otimes K_2$ . If there are, we return  $K_1 \cup K_2$ , which is now a (sub-optimal) solution to PMDM, and terminate. Otherwise, we proceed into iteration  $i = 3$ , which is analogous to iteration  $i = 2$ . Note that GREEDY  $\tau$ -PMDM always terminates with some (sub-optimal) solution  $K_1 \cup K_2 \cup \dots \cup K_j$ , for some  $j \leq \lceil \ell/\tau \rceil$ . Namely, in the worst case, it returns set  $\{1, \dots, \ell\}$  after  $\lceil \ell/\tau \rceil$  iterations and  $q_{\{1, \dots, \ell\}}$  matches all strings in  $\mathcal{D}$ . The reason why GREEDY  $\tau$ -PMDM does not guarantee finding an optimal solution

to PMDM is that at iteration  $i$  we fix the positions of wildcards based on solution  $K_1 \cup \dots \cup K_{i-1}$ , whereas some of those positions might not belong to the global optimal solution.

Since  $\tau = \mathcal{O}(1)$ , the time complexity of GREEDY  $\tau$ -PMDM is  $\mathcal{O}((N + N^{\tau/3})\ell)$ : each iteration takes time  $\mathcal{O}(N + N^{\tau/3})$  by Theorem 4.1, and then there are no more than  $\lceil \ell/\tau \rceil = \mathcal{O}(\ell)$  iterations. The space complexity of GREEDY  $\tau$ -PMDM is  $\mathcal{O}(N)$ . The hypergraph  $H = (V, E)$  used in the implementation of Theorem 4.1 has edges of size up to  $k$ . If every string in  $\mathcal{D}$  has more than  $k$  mismatches with  $q$ , then all edges in  $H$  have size larger than  $k$ . In this case, we preprocess the hypergraph  $H$ , as detailed below. The objective is to remove selected nodes and edges from  $H$ , so that it has at least one edge of size up to  $k$  and then apply GREEDY  $\tau$ -PMDM.

**Hypergraph Preprocessing** Let us now complete the description of our heuristic, by describing the hypergraph preprocessing. We want to ensure that hypergraph  $H = (V, E)$  has at least one edge of size up to  $k$  so that GREEDY  $\tau$ -PMDM can be applied. To this end, if there is no edge of size up to  $k$  at some iteration, then we add some nodes into the partial solution with the following heuristic.

1. We assign a score  $s(u)$  to each node  $u$  in  $V$  using the function:

$$s(u) = |E_u| \cdot \frac{\sum_{e \in E_u} w(e)}{\sum_{e \in E_u} |e|},$$

where  $E_u = \{e \in E : u \in e\}$  and  $w(e)$  is the edge weight.

2. Then, we add the node with maximum score from  $H$  (breaking ties arbitrarily) into the partial solution and update the edges accordingly.

These two steps are repeated until there is at least one edge of size up to  $k$ ; this takes  $\mathcal{O}(d\ell^2)$  time. After that, we add the removed nodes into the current solution  $K_k$  and use the resulting hypergraph to apply GREEDY  $\tau$ -PMDM.

The intuition behind the above process is to add nodes which appear in many short edges (so that we mask few positions) with large weight (so that the masked positions greatly increase the number of matched strings). We have also tried a different scoring function  $s'(u) = \sum_{e \in E_u} \frac{w(e)}{|e|}$  instead of  $s(u)$ , but the results were worse, and thus not reported.

## 9 Experimental Evaluation

**Methods** We compared the performance of our heuristic GREEDY  $\tau$ -PMDM (henceforth, GR  $\tau$ ), for the values  $\tau \in [3, 5]$ , for which its time complexity is subquadratic in  $N$ , to the following two algorithms:

- BASELINE (henceforth, BA). In iteration  $i$ , BA adds a node of hypergraph  $H$  into  $K$  and updates  $H$  according to the preprocessing described in Sect. 8. If at least  $z$  strings from  $\mathcal{D}$  match the query string  $q$  after the positions in  $q$  corresponding to  $K$  are replaced with wildcards, BA returns  $K$  and terminates; otherwise, it proceeds

into iteration  $i + 1$ . Note that BA generally constructs a *suboptimal* solution  $K$  to PMDM and takes  $\mathcal{O}(d\ell^2)$  time.

- BRUTEFORCE (henceforth, BF). In iteration  $i$ , BF applies Lemma 4.3 in the process of obtaining an *optimal* solution  $K$  of size  $i = k$  to PMDM. In particular, it checks whether at least  $z$  strings from  $\mathcal{D}$  match the query string  $q$ , after the  $i$  positions in  $q$  corresponding to  $K$  are replaced with wildcards. If the check succeeds, BF returns  $K$  and terminates; otherwise, it proceeds into iteration  $i + 1$ . BF takes  $\mathcal{O}(k(2\ell)^k + dk)$  time (see Lemma 4.3).

Since—as mentioned in Sect. 1—there are no existing algorithms for addressing PMDM, in the evaluation we used our own baseline BA. We have implemented all of the above algorithms in C++. Our implementations are freely available at <https://bitbucket.org/pattern-masking/pmdm/>.

**Datasets** We used the North Carolina Voter Registration database [57] (NCVR); a standard benchmark dataset for record linkage [23, 30, 45, 73]. NCVR is a collection of 7,736,911 records with attributes such as *Forename*, *Surname*, *City*, *County*, and *Gender*. We generated 4 subcollections of NCVR: (I) FS is comprised of all 952,864 records having *Forename* and *Surname* of total length  $\ell = 15$ ; (II) FCi is comprised of all 342,472 records having *Forename* and *City* of total length  $\ell = 15$ ; (III) FCiCo is comprised of all 342,472 records having *Forename*, *City*, and *County* of total length  $\ell = 30$ ; and (IV) FSCiCo is comprised of all 8,238 records having *Forename*, *Surname*, *City* and *County* of total length  $\ell = 45$ .

We also generated a synthetic dataset, referred to as SYN, using the IBM Synthetic Data Generator [41], a standard tool for generating sequential datasets [29, 40]. SYN contains a collection of  $6 \cdot 10^6$  records, each of length  $\ell = 50$ , over an alphabet of size  $|\Sigma| = 10$ . We also generated subcollections of SYN comprised of:  $x \cdot 10^6$  arbitrarily selected records; the length- $y$  prefix of each selected record. We denote each resulting dataset by  $\text{SYN}_{x,y}$ .

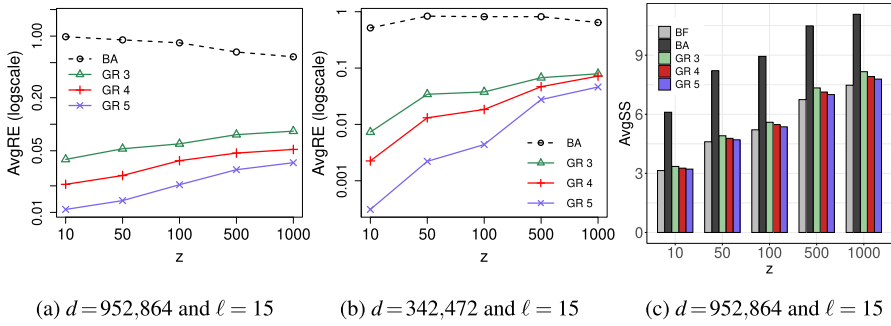
**Comparison Measures** We evaluated the *effectiveness* of the algorithms using:

**AvgRE** An Average Relative Error measure, computed as  $\text{avg}_{i \in [1, 1000]} \frac{k_i - k_i^*}{k_i^*}$ , where  $k_i^*$  is the size of the optimal solution produced by BF, and  $k_i$  is the size of the solution produced by one of the other tested algorithms. Both  $k_i^*$  and  $k_i$  are obtained by using, as query  $q_i$ , a record of the input dictionary selected uniformly at random.

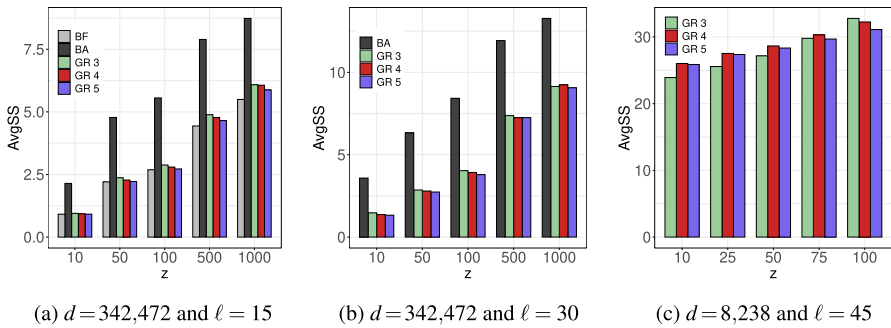
**AvgSS** An Average Solution Size measure computed as  $\text{avg}_{i \in [1, 1000]} k_i^*$  for BF and  $\text{avg}_{i \in [1, 1000]} k_i$  for any other algorithm.

We evaluated *efficiency* by reporting  $\text{avg}_{i \in [1, 1000]} t_i$ , where  $t_i$  is the elapsed time of a tested algorithm to obtain a solution for query  $q_i$  over the input dictionary.

**Execution Environment** In our experiments we used a PC with Intel Xeon E5-2640@2.66GHz and 160GB RAM running GNU/Linux, and a gcc v.7.3.1 compiler at optimization level -O3.



**Fig. 4** AvgRE (in logscale) versus  $z$  computed for **a** FS and **b** FCI; **c** AvgSS versus  $z$  for FS



**Fig. 5** AvgSS versus  $z$  for (a) FCI. (b) FCI Co (BF did not produce results for any  $z$  within 48 h), and (c) FSCiCo (BF and BA did not produce results for any  $z$  within 48 h. The results of GR for  $z > 100$  are omitted because AvgSS  $> 40$  which is close to  $\ell = 45$ )

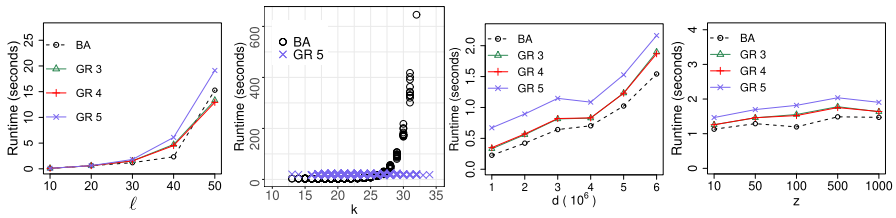
**Effectiveness** Figs. 4 and 5a show that GR produced nearly-optimal solutions, significantly outperforming BA. In Fig. 4a, the solutions of GR 3 were no more than 9% worse than the optimal, while those of BA were up to 95% worse. In Fig. 5a, the average solution size of BF was 5.4 versus 5.9 and 9, for the solution size of GR 3 and BA, respectively.

In Figs. 5b and c, we examined the effectiveness of GR for larger  $\ell$  values. Figure 5b shows that the solution size of GR 3 was at least 31% and up to 60% smaller than that of BA on average, while Fig. 5c shows that the solution of GR 3 was comparable to that of GR 4 and 5. We omit the results for BF from Figs. 5b and c and those for BA from Fig. 5c, as these algorithms did not produce results for all queries within 48 h, for any  $z$ . This is because, unlike GR, BF does not scale well with  $\ell$  and BA does not scale well with the solution size, as we will explain later.

Note that increasing  $\tau$  generally increases the effectiveness of GR as it computes more positions of wildcards per iteration. However, even with  $\tau = 3$ , it remains competitive to BF.

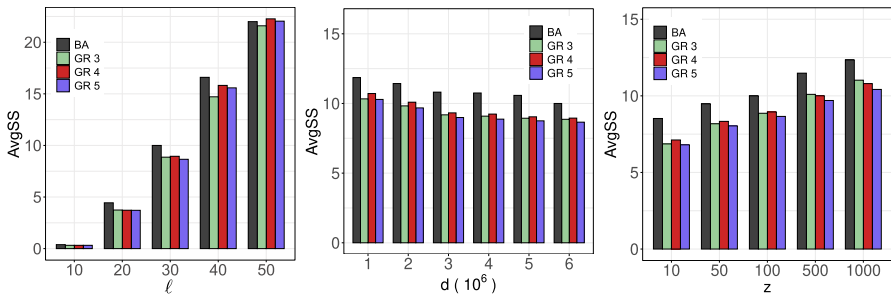
**Efficiency** Having shown that GR produced nearly-optimal solutions, we now show that it is comparable in terms of efficiency or faster than BA for synthetic data. (BF was at least two orders of magnitude slower than the other methods on average and thus we





(a)  $d = 6 \cdot 10^6$  and  $z = 100$  (b)  $\ell = 50$  and  $z = 100$  (c)  $\ell = 30$  and  $z = 100$  (d)  $d = 6 \cdot 10^6$  and  $\ell = 30$

**Fig. 6** Efficiency versus (a)  $\ell$  for  $SYN_{6,\ell}$ , (b)  $k$  for  $SYN$ , (c)  $d$  for  $SYN_{x,30}$ ,  $x \in \{1 \cdot 10^6, 2 \cdot 10^6, \dots, 6 \cdot 10^6\}$ , and (d)  $z$  for  $SYN_{6,30}$ . The results of BF are omitted, because it was slower than other methods by at least two orders of magnitude on average



(a)  $d = 6 \cdot 10^6$ ,  $z = 100$  (b)  $\ell = 30$ ,  $z = 100$  (c)  $d = 6 \cdot 10^6$ ,  $\ell = 30$

**Fig. 7** AvgSS versus (a)  $\ell$  for  $SYN_{6,\ell}$ , (b)  $d$  for  $SYN_{x,30}$ ,  $x \in \{1 \cdot 10^6, 2 \cdot 10^6, \dots, 6 \cdot 10^6\}$ , and (c)  $z$  for  $SYN_{6,30}$

omit its results.) The results for NCVR were qualitatively similar (omitted). Figure 6a shows that GR spent, on average, the same time for a query as BA did. However, it took significantly (up to 36 times) less time than BA for queries with large solution size  $k$ . This can be seen from Fig. 6b, which shows the time each query with solution size  $k$  took; the results for GR 3 and 4 were similar and thus omitted. The reason is that BA updates the hypergraph every time a node is added into the solution, which is computationally expensive when  $k$  is large. Figures 6c and d show that all algorithms scaled sublinearly with  $d$  and with  $z$ , respectively. The increase with  $d$  is explained by the time complexity of the methods. The slight increase with  $z$  is because  $k$  gets larger, on average, as  $z$  increases (see Fig. 7c next, in which we also show the average solution size for the experiments in Figs. 6a, c). GR 3 and 4 performed similarly to each other, being faster than GR 5 in all experiments as expected: increasing  $\tau$  from 3 or 4 to 5 trades-off effectiveness for efficiency.

**Average Solution Size** Figs. 7a–c show the average solution size in the experiments of Fig. 6a, c, and d, respectively. Observe that the results are analogous to those obtained using the NCVR datasets: GR outperforms BA significantly. Also, observe in Fig. 7c that the solution size of each tested algorithm gets larger, on average, as  $z$  increases.

**Summary** We have presented an extensive experimental evaluation demonstrating the effectiveness and efficiency of the proposed heuristic on real-world datasets used in record linkage, as well as on synthetic datasets. In the experiments that we have performed, the proposed heuristic: (I) found nearly-optimal solutions for varying values of  $z$  and  $\ell$ , even when applied with a small  $\tau$ ; and (II) scaled as predicted by the complexity analysis, requiring fewer than 3 s for  $d = 6 \cdot 10^6$  in all tested cases. Our experimental results suggest that our methods can inspire solutions in large-scale real-world systems, where no sophisticated algorithms for PMDM are being used.

## 10 Open Questions

The following questions of theoretical nature remain unanswered:

1. Can we improve on the exact  $\mathcal{O}(d\ell + (d\ell)^{k/3})$ -time algorithm presented in Sect. 4 for PMDM and  $k = \mathcal{O}(1)$  using fast matrix multiplication [4, 53] or show that algebraic techniques cannot help, e.g., via the  $(c, k)$ -HYPERCLIQUE problem?
2. Can we improve on the  $\mathcal{O}(2^{\ell/2}(2^{\ell/2} + \tau)\ell)$ -time and  $\mathcal{O}(2^\ell d^2/\tau^2 + 2^{\ell/2}d)$ -space trade-off presented in Sect. 6 for the data structure answering  $q, z$  PMDM queries?

**Acknowledgements** Panagiotis Charalampopoulos was partially supported by the Israel Science Foundation grants 592/17 and 810/21. Huiping Chen was supported by a CSC Scholarship. Part of the work was performed when this author was at King's College London. Grigorios Loukides was supported in part by the Leverhulme Trust RPG-2019-399 project. Nadia Pisanti is supported by NextGeneration EU programme PNRR ECS00000017 Tuscany Health Ecosystem, by MUR PRIN 2022YRB97K PINC, and by the ALPACA EU project (GA 956229). Jakub Radoszewski was supported by the Polish National Science Center, grant number 2018/31/D/ST6/03991. The authors acknowledge the use of the following HPC infrastructure: King's Computational Research, Engineering and Technology Environment (CREATE).

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

## References

1. Abboud, A., Backurs, A., Williams, V.V.: If the current clique algorithms are optimal, so is Valiant's parser. *SIAM J. Comput.* **47**(6), 2527–2555 (2018). <https://doi.org/10.1137/16M1061771>
2. Afshani, P., Nielsen, J.S.: Data structure lower bounds for document indexing problems. In: 43rd International Colloquium on Automata, Languages and Programming (ICALP 2016), LIPIcs, vol. 55, pp. 93:1–93:15 (2016). <https://doi.org/10.4230/LIPIcs.ICALP.2016.93>
3. Aho, A.V., Corasick, M.J.: Efficient string matching: an aid to bibliographic search. *Commun. ACM* **18**(6), 333–340 (1975). <https://doi.org/10.1145/360825.360855>
4. Alman, J., Williams, V.V.: A refined laser method and faster matrix multiplication. In: Marx, D. (ed.) Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms, SODA 2021, Virtual Conference, January 10–13, 2021, pp. 522–539. SIAM, Philadelphia (2021). <https://doi.org/10.1137/1.9781611976465.32>

5. Applebaum, B.: Pseudorandom generators with long stretch and low locality from random local one-way functions. *SIAM J. Comput.* **42**(5), 2008–2037 (2013). <https://doi.org/10.1137/120884857>
6. Arimura, H., Uno, T.: An efficient polynomial space and polynomial delay algorithm for enumeration of maximal motifs in a sequence. *J. Comb. Optim.* **13**(3), 243–262 (2007). <https://doi.org/10.1007/s10878-006-9029-1>
7. Bailey, D.R., Battle, A.J., Gomes, B.A., Nayak, P.P.: Estimating Confidence for Query Revision Models. U.S. Patent US7617205B2 (granted to Google) (2009)
8. Bailey, M., Cole, C., Henderson, M., Massey, C.: How well do automated linking methods perform? Lessons from U.S. historical data. NBER Working Papers 24019, National Bureau of Economic Research, Inc (2017). <https://doi.org/10.3386/w24019>
9. Battaglia, G., Cangelosi, D., Grossi, R., Pisanti, N.: Masking patterns in sequences: a new class of motif discovery with don't cares. *Theor. Comput. Sci.* **410**(43), 4327–4340 (2009). <https://doi.org/10.1016/j.tcs.2009.07.014>
10. Belazzougui, D.: Faster and space-optimal edit distance “1” dictionary. In: 20th Annual Symposium on Combinatorial Pattern Matching (CPM 2009), Lecture Notes in Computer Science, vol. 5577, pp. 154–167. Springer, New York (2009). [https://doi.org/10.1007/978-3-642-02441-2\\_14](https://doi.org/10.1007/978-3-642-02441-2_14)
11. Belazzougui, D., Venturini, R.: Compressed string dictionary search with edit distance one. *Algorithmica* **74**(3), 1099–1122 (2016). <https://doi.org/10.1007/s00453-015-9990-0>
12. Bille, P., Gørtz, I.L., Vildhøj, H.W., Vind, S.: String indexing for patterns with wildcards. *Theory Comput. Syst.* **55**(1), 41–60 (2014). <https://doi.org/10.1007/s00224-013-9498-4>
13. Borodin, A., Ostrovsky, R., Rabani, Y.: Lower bounds for high dimensional nearest neighbor search and related problems. In: 31st ACM Symposium on Theory of Computing (STOC 1999), pp. 312–321 (1999). <https://doi.org/10.1145/301250.301330>
14. Brodal, G.S., Venkatesh, S.: Improved bounds for dictionary look-up with one error. *Inf. Process. Lett.* **75**(1–2), 57–59 (2000). [https://doi.org/10.1016/S0020-0190\(00\)00079-X](https://doi.org/10.1016/S0020-0190(00)00079-X)
15. Calabro, C., Impagliazzo, R., Paturi, R.: The complexity of satisfiability of small depth circuits. In: Parameterized and Exact Computation, 4th International Workshop (IWPEC 2009), Lecture Notes in Computer Science, vol. 5917, pp. 75–85. Springer, New York (2009). [https://doi.org/10.1007/978-3-642-11269-0\\_6](https://doi.org/10.1007/978-3-642-11269-0_6)
16. Chan, H., Lam, T.W., Sung, W., Tam, S., Wong, S.: Compressed indexes for approximate string matching. *Algorithmica* **58**(2), 263–281 (2010). <https://doi.org/10.1007/s00453-008-9263-2>
17. Charalampopoulos, P., Chen, H., Christen, P., Loukides, G., Pisanti, N., Pissis, S.P., Radoszewski, J.: Pattern Masking for Dictionary Matching. In: Ahn, H.K., Sadakane, K. (eds.) 32nd International Symposium on Algorithms and Computation (ISAAC 2021), Leibniz International Proceedings in Informatics (LIPIcs), vol. 212, pp. 65:1–65:19. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl, Germany (2021). <https://doi.org/10.4230/LIPIcs.ISAAC.2021.65>
18. Charikar, M., Indyk, P., Panigrahy, R.: New algorithms for subset query, partial match, orthogonal range searching, and related problems. In: 29th International Colloquium on Automata, Languages and Programming (ICALP 2002), pp. 451–462 (2002). [https://doi.org/10.1007/3-540-45465-9\\_39](https://doi.org/10.1007/3-540-45465-9_39)
19. Chen, J., Huang, X., Kanj, I.A., Xia, G.: Strong computational lower bounds via parameterized complexity. *J. Comput. Syst. Sci.* **72**(8), 1346–1367 (2006). <https://doi.org/10.1016/j.jcss.2006.04.007>
20. Chlamtáč, E., Dinitz, M., Makarychev, Y.: Minimizing the union: tight approximations for small set bipartite vertex expansion. In: 28th ACM-SIAM Symposium on Discrete Algorithms (SODA 2017), pp. 881–899 (2017). <https://doi.org/10.1137/1.9781611974782.56>
21. Chlamtáč, E., Dinitz, M., Konrad, C., Kortsarz, G., Rabanca, G.: The densest k-subhypergraph problem. *SIAM J. Discrete Math.* **32**(2), 1458–1477 (2018). <https://doi.org/10.1137/16M1096402>
22. Christen, P.: Data Matching - Concepts and Techniques for Record Linkage, Entity Resolution, and Duplicate Detection. Data-Centric Systems and Applications. Springer, Heidelberg (2012). <https://doi.org/10.1007/978-3-642-31164-2>
23. Christen, P., Gayler, R.W., Tran, K.N., Fisher, J., Vatsalan, D.: Automatic discovery of abnormal values in large textual databases. *J. Data Inf. Qual.* (2016). <https://doi.org/10.1145/2889311>
24. Christen, P., Ranbaduge, T., Schnell, R.: Linking Sensitive Data. Springer, Heidelberg (2020). <https://doi.org/10.1007/978-3-030-59706-1>
25. Cohen-Addad, V., Feuilloley, L., Starikovskaya, T.: Lower bounds for text indexing with mismatches and differences. In: 30th ACM-SIAM Symposium on Discrete Algorithms (SODA 2019), pp. 1146–1164 (2019). <https://doi.org/10.1137/1.9781611975482.70>


26. Cole, R., Gottlieb, L., Lewenstein, M.: Dictionary matching and indexing with errors and don't cares. In: 36th ACM Symposium on Theory of Computing (STOC 2004), pp. 91–100 (2004). <https://doi.org/10.1145/1007352.1007374>
27. Cuzzocrea, A., Shahriar, H.: Data masking techniques for nosql database security: a systematic review. In: 2017 IEEE International Conference on Big Data (BigData 2017), pp. 4467–4473 (2017). <https://doi.org/10.1109/BigData.2017.8258486>
28. Cygan, M., Fomin, F.V., Kowalik, L., Lokshtanov, D., Marx, D., Pilipczuk, M., Pilipczuk, M., Saurabh, S.: Parameterized Algorithms. Springer, New York (2015). <https://doi.org/10.1007/978-3-319-21275-3>
29. Ding, B., Lo, D., Han, J., Khoo, S.C.: Efficient mining of closed repetitive gapped subsequences from a sequence database. In: 25th IEEE International Conference on Data Engineering (ICDE), pp. 1024–1035 (2009). <https://doi.org/10.1109/ICDE.2009.104>
30. Durham, E.A., Kantarcioglu, M., Xue, Y., Tóth, C., Malin, B.: Composite bloom filters for secure record linkage. *IEEE Trans. Knowl. Data Eng.* **26**(12), 2956–2968 (2014). <https://doi.org/10.1109/TKDE.2013.91>
31. Federico, M., Pisanti, N.: Suffix tree characterization of maximal motifs in biological sequences. *Theor. Comput. Sci.* **410**(43), 4391–4401 (2009). <https://doi.org/10.1016/j.tcs.2009.07.020>
32. Fredman, M.L., Komlós, J., Szemerédi, E.: Storing a sparse table with  $O(1)$  worst case access time. *J. ACM* **31**(3), 538–544 (1984). <https://doi.org/10.1145/828.1884>
33. Fung, B.C.M., Wang, K., Chen, R., Yu, P.S.: Privacy-preserving data publishing: a survey of recent developments. *ACM Comput. Surv.* **42**(4), 1–53 (2010). <https://doi.org/10.1145/1749603.1749605>
34. Gawrychowski, P., Lewenstein, M., Nicholson, P.K.: Weighted ancestors in suffix trees. In: Algorithms - 22th Annual European Symposium (ESA 2014), Lecture Notes in Computer Science, vol. 8737, pp. 455–466. Springer (2014). [https://doi.org/10.1007/978-3-662-44777-2\\_38](https://doi.org/10.1007/978-3-662-44777-2_38)
35. Gollapudi, S., Jeong, S., Ntoulas, A., Papatrinos, S.: Efficient query rewrite for structured web queries. In: 20th ACM International Conference on Information and Knowledge Management (CIKM 2011), pp. 2417–2420 (2011). <https://doi.org/10.1145/2063576.2063981>
36. Grossi, R., Pietracaprina, A., Pisanti, N., Pucci, G., Upfal, E., Vandin, F.: MADMX: a strategy for maximal dense motif extraction. *J. Comput. Biol.* **18**(4), 535–545 (2011). <https://doi.org/10.1089/CMB.2010.0177>
37. Grossi, R., Menconi, G., Pisanti, N., Trani, R., Vind, S.: Motif trie: An efficient text index for pattern discovery with don't cares. *Theor. Comput. Sci.* **710**, 74–87 (2018). <https://doi.org/10.1016/j.tcs.2017.04.012>
38. Hastad, J.: Clique is hard to approximate within  $n^{1-\epsilon}$ . *Acta Math.* **182**, 105–142 (1999). <https://doi.org/10.1007/BF02392825>
39. Herzog, T.N., Scheuren, F.J., Winkler, W.E.: Data Quality and Record Linkage Techniques. Springer, New York (2007)
40. I, T., Enokuma, Y., Bannai, H., Takeda, M.: General algorithms for mining closed flexible patterns under various equivalence relations. In: Machine Learning and Knowledge Discovery in Databases, pp. 435–450 (2012). [https://doi.org/10.1007/978-3-642-33486-3\\_28](https://doi.org/10.1007/978-3-642-33486-3_28)
41. IBM Synthetic Data Generator for Itemsets and Sequences. <https://github.com/zakimjz/IBMGenerator> (2020)
42. Ilyas, I.F., Beskales, G., Soliman, M.A.: A survey of top-k query processing techniques in relational database systems. *ACM Comput. Surveys* (2008). <https://doi.org/10.1145/1391729.1391730>
43. Impagliazzo, R., Paturi, R.: On the complexity of k-SAT. *J. Comput. Syst. Sci.* **62**(2), 367–375 (2001). <https://doi.org/10.1006/jcss.2000.1727>
44. Jayram, T.S., Khot, S., Kumar, R., Rabani, Y.: Cell-probe lower bounds for the partial match problem. *J. Comput. Syst. Sci.* **69**(3), 435–447 (2004). <https://doi.org/10.1016/j.jcss.2004.04.006>
45. Karapiperis, D., Gkoulalas-Divanis, A., Verykios, V.S.: Summarizing and linking electronic health records. *Distrib. Parallel Datab.* (2019). <https://doi.org/10.1007/s10619-019-07263-0>
46. Karp, R.M.: Reducibility among combinatorial problems. In: 50 Years of Integer Programming 1958-2008 - From the Early Years to the State-of-the-Art, pp. 219–241. Springer (2010). [https://doi.org/10.1007/978-3-540-68279-0\\_8](https://doi.org/10.1007/978-3-540-68279-0_8)
47. Kellerer, H., Pferschy, U., Pisinger, D.: The Multiple-Choice Knapsack Problem, pp. 317–347. Springer, Berlin, Heidelberg (2004). [https://doi.org/10.1007/978-3-540-24777-7\\_11](https://doi.org/10.1007/978-3-540-24777-7_11)
48. Konda, P., et al.: Technical perspective: toward building entity matching management systems. *SIGMOD Record* **47**(1), 33–40 (2018). <https://doi.org/10.1145/3277006.3277015>

49. Kum, H.C., Ragan, E.D., Ilangovan, G., Ramezani, M., Li, Q., Schmit, C.: Enhancing privacy through an interactive on-demand incremental information disclosure interface: applying privacy-by-design to record linkage. In: Fifteenth USENIX Conference on Usable Privacy and Security, pp. 175–189 (2019). <https://doi.org/10.5555/3361476.3361489>
50. Kum, H.C., Krishnamurthy, A., Machanavajjhala, A., Reiter, M.K., Ahalt, S.: Privacy preserving interactive record linkage (PPIRL). *J. Am. Med. Inform. Assoc.* **21**(2), 212–220 (2014). <https://doi.org/10.1136/amiajnl-2013-002165>
51. Kumar, P.S., Arasada, P., Jammalamadaka, R.C.: Systems and Methods for Generating Search Query Rewrites. U.S. Patent US10108712B2 (granted to ebay) (2018)
52. Lewenstein, M., Nekrich, Y., Vitter, J.S.: Space-efficient string indexing for wildcard pattern matching. In: 31st Symposium on Theoretical Aspects of Computer Science (STACS 2014), pp. 506–517 (2014). <https://doi.org/10.4230/LIPIcs.STACS.2014.506>
53. Le Gall, F.: Powers of tensors and fast matrix multiplication. In: Nabeshima, K., Nagasaka, K., Winkler, F., Szántó, Á. (eds.) International Symposium on Symbolic and Algebraic Computation, ISSAC '14, Kobe, Japan, July 23–25, 2014, pp. 296–303. ACM (2014). <https://doi.org/10.1145/2608628.2608664>
54. Lewenstein, M., Munro, J.I., Raman, V., Thankachan, S.V.: Less space: Indexing for queries with wildcards. *Theor. Comput. Sci.* **557**, 120–127 (2014). <https://doi.org/10.1016/j.tcs.2014.09.003>
55. Lincoln, A., Williams, V.V., Williams, R.R.: Tight hardness for shortest cycles and paths in sparse graphs. In: 29th ACM-SIAM Symposium on Discrete Algorithms (SODA 2018), pp. 1236–1252 (2018). <https://doi.org/10.1137/1.9781611975031.80>
56. Miltersen, P.B., Nisan, N., Safra, S., Wigderson, A.: On data structures and asymmetric communication complexity. *J. Comput. Syst. Sci.* **57**(1), 37–49 (1998). <https://doi.org/10.1006/jcss.1998.1577>
57. North Carolina Voter Registration database (dataset ncvoter\_Statewide.zip). <https://dl.ncsbe.gov/?prefix=data/> (2020)
58. Papadakis, G., Skoutas, D., Thanos, E., Palpanas, T.: Blocking and filtering techniques for entity resolution: a survey. *ACM Comput. Surv.* (2020). <https://doi.org/10.1145/3377455>
59. Pătraşcu, M.: Unifying the landscape of cell-probe lower bounds. *SIAM J. Comput.* **40**(3), 827–847 (2011). <https://doi.org/10.1137/09075336X>
60. Pătraşcu, M., Thorup, M.: Higher lower bounds for near-neighbor and further rich problems. *SIAM J. Comput.* **39**(2), 730–741 (2009). <https://doi.org/10.1137/070684859>
61. Pisanti, N., Crochemore, M., Grossi, R., Sagot, M.: A basis of tiling motifs for generating repeated patterns and its complexity for higher quorum. In: 28th International Symposium on Mathematical Foundations of Computer Science 2003 (MFCS), Lecture Notes in Computer Science, vol. 2747, pp. 622–631. Springer, New York (2003). [https://doi.org/10.1007/978-3-540-45138-9\\_56](https://doi.org/10.1007/978-3-540-45138-9_56)
62. Pisanti, N., Crochemore, M., Grossi, R., Sagot, M.: Bases of motifs for generating repeated patterns with wild cards. *IEEE/ACM Trans. Comput. Biol. Bioinf.* **2**(1), 40–50 (2005). <https://doi.org/10.1109/TCBB.2005.5>
63. Ragan, E.D., Kum, H.C., Ilangovan, G., Wang, H.: Balancing privacy and information disclosure in interactive record linkage with visual masking. In: ACM Conference on Human Factors in Computing Systems (CHI 2018) (2018). <https://doi.org/10.1145/3173574.3173900>
64. Rivest, R.L.: Partial-match retrieval algorithms. *SIAM J. Comput.* **5**(1), 19–50 (1976). <https://doi.org/10.1137/0205003>
65. Samarati, P., Sweeney, L.: Generalizing data to provide anonymity when disclosing information (abstract). In: Proceedings of the Seventeenth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS 1998), p. 188. Association for Computing Machinery (1998). <https://doi.org/10.1145/275487.275508>
66. Samarati, P., Sweeney, L.: Protecting Privacy When Disclosing Information:  $k$ -Anonymity and Its Enforcement Through Generalization and Suppression. Tech. rep, Computer Science Laboratory, SRI International (1998)
67. Samarati, P.: Protecting respondents' identities in microdata release. *IEEE Trans. Knowl. Data Eng.* **13**(6), 1010–1027 (2001). <https://doi.org/10.1109/69.971193>
68. Santos, R.J., Bernardino, J., Vieira, M.: A data masking technique for data warehouses. In: 15th International Database Engineering and Applications Symposium (IDEAS 2011), pp. 61–69 (2011). <https://doi.org/10.1145/2076623.2076632>
69. Secure critical data with Oracle Data Safe (white paper). <https://www.oracle.com/a/tech/docs/dbsec/data-safe/wp-security-data-safe.pdf> (2020)

70. Sweeney, L.: Computational disclosure control: a primer on data privacy protection. Ph.D. thesis, Massachusetts Institute of Technology, Cambridge, MA, USA (2001). <http://hdl.handle.net/1721.1/8589>
71. Tan, Z., Xu, C., Jiang, M., Yang, H., Wu, X.: Query rewrite for null and low search results in ecommerce. In: SIGIR Workshop On eCommerce, CEUR Workshop Proceedings, vol. 2311 (2017)
72. Tao, Y.: Entity matching with active monotone classification. In: 37th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems (PODS 2018), pp. 49–62 (2018). <https://doi.org/10.1145/3196959.3196984>
73. Vatsalan, D., Christen, P.: Scalable privacy-preserving record linkage for multiple databases. In: 23rd ACM International Conference on Information and Knowledge Management (CIKM 2014), pp. 1795–1798 (2014). <https://doi.org/10.1145/2661829.2661875>
74. Vatsalan, D., Sehili, Z., Christen, P., Rahm, E.: Privacy-preserving record linkage for Big Data: Current approaches and research challenges. In: Handbook of Big Data Technologies, pp. 851–895. Springer, New York (2017). <https://doi.org/10.1007/978-3-319-49340-4>
75. Weiner, P.: Linear pattern matching algorithms. In: 14th Annual Symposium on Switching and Automata Theory (SWAT 1973), pp. 1–11. IEEE Computer Society (1973). <https://doi.org/10.1109/SWAT.1973.13>
76. Williams, V.V.: On some fine-grained questions in algorithms and complexity. In: 2018 International Congress of Mathematicians (ICM), pp. 3447–3487 (2019). [https://doi.org/10.1142/9789813272880\\_0188](https://doi.org/10.1142/9789813272880_0188)
77. Yao, A.C., Yao, F.F.: Dictionary look-up with one error. *J. Algorithms* **25**(1), 194–202 (1997). <https://doi.org/10.1006/jagm.1997.0875>
78. Zuckerman, D.: Linear degree extractors and the inapproximability of max clique and chromatic number. *Theory Comput.* **3**(1), 103–128 (2007). <https://doi.org/10.4086/toc.2007.v003a006>

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

## Authors and Affiliations

Panagiotis Charalampopoulos<sup>1</sup> · Huiping Chen<sup>2</sup> · Peter Christen<sup>3</sup> ·  
 Grigorios Loukides<sup>4</sup> · Nadia Pisanti<sup>5</sup> · Solon P. Pissis<sup>6,7</sup>  ·  
 Jakub Radoszewski<sup>8</sup>

Panagiotis Charalampopoulos  
 p.charalampopoulos@bbk.ac.uk

Peter Christen  
 peter.christen@anu.edu.au

Grigorios Loukides  
 grigorios.loukides@kcl.ac.uk

Nadia Pisanti  
 nadia.pisanti@unipi.it

Jakub Radoszewski  
 jrad@mimuw.edu.pl

<sup>1</sup> School of Computing and Mathematical Sciences, Birkbeck, University of London, London, UK

<sup>2</sup> School of Computer Science, University of Birmingham, Birmingham, UK

<sup>3</sup> Australian National University, Canberra, Australia

<sup>4</sup> Department of Informatics, King's College London, London, UK

<sup>5</sup> Università di Pisa, Pisa, Italy

- 6 CWI, Amsterdam, The Netherlands
- 7 Vrije Universiteit, Amsterdam, The Netherlands
- 8 Institute of Informatics, University of Warsaw, Warsaw, Poland