

# 1 Order-Preserving Pattern Matching 2 Indeterminate Strings

3 **Rui Henriques**

4 INESC-ID and Instituto Superior Técnico, Universidade de Lisboa, Portugal  
5 rmch@tecnico.ulisboa.pt (correspondence)

6 **Alexandre P. Francisco**

7 INESC-ID and Instituto Superior Técnico, Universidade de Lisboa, Portugal

8 **Luís M. S. Russo**

9 INESC-ID and Instituto Superior Técnico, Universidade de Lisboa, Portugal

10 **Hideo Bannai**

11 Department of Computer Science, Kyushu University, Japan

## 12 — Abstract —

13 Given an indeterminate string pattern  $p$  and an indeterminate string text  $t$ , the problem of order-  
14 preserving pattern matching with character uncertainties ( $\mu$ OPPM) is to find all substrings of  $t$   
15 that satisfy one of the possible orderings defined by  $p$ . When the text and pattern are determ-  
16 inate strings, we are in the presence of the well-studied exact order-preserving pattern matching  
17 (OPPM) problem with diverse applications on time series analysis. Despite its relevance, the  
18 exact OPPM problem suffers from two major drawbacks: 1) the inability to deal with indetermin-  
19 ation in the text, thus preventing the analysis of noisy time series; and 2) the inability to deal  
20 with indetermination in the pattern, thus imposing the strict satisfaction of the orders among all  
21 pattern positions.

22 This paper provides the first polynomial algorithm to answer the  $\mu$ OPPM problem when  
23 indetermination is observed on the pattern or text. Given two strings with length  $m$  and  $O(r)$   
24 uncertain characters per string position, we show that the  $\mu$ OPPM problem can be solved in  
25  $O(mr \lg r)$  time when one string is indeterminate and  $r \in \mathbb{N}^+$ . Mappings into satisfiability  
26 problems are provided when indetermination is observed on both the pattern and the text.

27 **2012 ACM Subject Classification** Theory of computation  $\rightarrow$  Pattern matching

28 **Keywords and phrases** Order-preserving pattern matching; Indeterminate string analysis;  
29 Generic pattern matching; Satisfiability

30 **Digital Object Identifier** 10.4230/LIPIcs.CPM.2018.2

## 31 **1** Introduction

32 Given a pattern string  $p$  and a text string  $t$ , the exact order preserving pattern matching  
33 (OPPM) problem is to find all substrings of  $t$  with the same relative orders as  $p$ . The problem  
34 is applicable to strings with characters drawn from numeric or ordinal alphabets. Illustrat-  
35 ing, given  $p=(1,5,3,3)$  and  $t=(5,1,4,2,2,5,2,4)$ , substring  $t[1..4]=(1,4,2,2)$  is reported since it  
36 satisfies the character orders in  $p$ ,  $p[0] \leq p[2] = p[3] \leq p[1]$ . Despite its relevance, the OPPM  
37 problem has limited potential since it prevents the specification of errors, uncertainties or  
38 don't care characters within the text.

39 Indeterminate strings allow uncertainties between two or more characters per position.  
40 Given indeterminate strings  $p$  and  $t$ , the problem of order preserving pattern matching  
41 uncertain text ( $\mu$ OPPM) is to find all substrings of  $t$  with an assignment of values that satisfy  
42 the orders defined by  $p$ . For instance, let  $p=(1,2|5,3,3)$  and  $t=(5,0,1,2|1,2,5,2|3,3|4)$ . The



© Henriques et al., 2018;

licensed under Creative Commons License CC-BY

29th Annual Symposium on Combinatorial Pattern Matching (CPM 2018).

Editors: Gonzalo Navarro, David Sankoff, and Binhai Zhu; Article No. 2; pp. 2:1–2:15

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

43 substrings  $t[1..4]$  and  $t[4..7]$  are reported since there is an assignment of values that preserve  
 44 either  $p[0] < p[1] < p[2] = p[3]$  or  $p[0] < p[2] = p[3] < p[1]$  orderings: respectively  $t[1..4] = (0, 1, 2, 2)$   
 45 and  $t[4..7] = (2, 5, 3, 3)$ .

46 Order-preserving pattern matching captures the structural isomorphism of strings, there-  
 47 fore having a wide-range of relevant applications in the analysis of financial times series,  
 48 musical sheets, physiological signals and biological sequences [32, 39, 36]. Uncertainties of-  
 49 ten occur across these domains. In this context, although the OPPM problem is already a  
 50 relaxation of the traditional pattern matching problem, the need to further handle localized  
 51 errors is essential to deal with noisy strings [33]. For instance, given the stochasticity of  
 52 gene regulation (or markets), the discovery of order-preserving patterns in gene expression  
 53 (or financial) time series needs to account for uncertainties [35, 34]. Numerical indexes of  
 54 amino-acids (representing physiochemical and biochemical properties) are subjected to er-  
 55 rors difficulting the analysis of protein sequences [38]. Another example are ordinal strings  
 56 obtained from the discretization of numerical strings, often having two uncertain characters  
 57 in positions where the original values are near a discretization boundary [33].

58 Let  $m$  and  $n$  be the length of the pattern  $p$  and text  $t$ , respectively. The exact OPPM  
 59 problem has a linear solution on the text length  $O(n + m \lg m)$  based on the Knuth-Morris-  
 60 Pratt algorithm [41, 39, 22]. Alternative algorithms for the OPPM problem have also been  
 61 proposed [21, 12, 20]. Contrasting with the large attention given to the resolution of the  
 62 OPPM problem, to our knowledge there are no polynomial-time algorithms to solve the  
 63  $\mu$ OPPM problem. Naive algorithms for  $\mu$ OPPM assess all possible pattern and text assign-  
 64 ments, bounded by  $O(nr^m)$  when considering up to  $r$  uncertain characters per position.

65 This work proposes the first polynomial algorithms able to answer the  $\mu$ OPPM problem.  
 66 Accordingly, the contributions are organized as follows. First, we show that an indeterminate  
 67 string of length  $m$  order-preserving matches a determinate string with the same length in  
 68  $O(mr \lg r)$  time based on their monotonic properties. Second, and given two indeterminate  
 69 strings with the same size, we provide a linear encoding of the  $\mu$ OPPM into a satisfiability  
 70 formula with properties of interest. Third, given a pattern and text strings with lengths  
 71  $m$  and  $n$ , we show that the  $\mu$ OPPM problem can be solved in linear space and its average  
 72 efficiency boosted under effective filtration procedures.

## 73 2 Background

74 Let  $\Sigma$  be a totally ordered alphabet and an element of  $\Sigma^*$  be a string. The length of a string  
 75  $w$  is denoted by  $|w|$ . The empty string  $\varepsilon$  is a string of length 0. For a string  $w = xyz$ ,  $x$ ,  $y$   
 76 and  $z$  are called a prefix, substring, and suffix of  $w$ , respectively. The  $i$ -th character of a  
 77 string  $w$  is denoted by  $w[i]$  for each  $1 \leq i \leq |w|$ . For a string  $w$  and integers  $1 \leq i \leq j \leq |w|$ ,  $w[i..j]$   
 78 denotes the substring of  $w$  from position  $i$  to position  $j$ . For convenience, let  $w[i..j] = \varepsilon$  when  
 79  $i > j$ .

80 Given strings  $x$  and  $y$  with equal length  $m$ ,  $y$  is said to order-preserving against  $x$   
 81 [41], denoted by  $x \approx y$ , if the orders between the characters of  $x$  and  $y$  are the same, i.e.  
 82  $x[i] \leq x[j] \Leftrightarrow y[i] \leq y[j]$  for any  $1 \leq i, j \leq m$ . A non-empty pattern string  $p$  is said to  
 83 order-preserving match (*op-match* in short) a non-empty text string  $t$  iff there is a position  
 84  $i$  in  $t$  such that  $p \approx t[i - |p| + 1..i]$ . The *order-preserving pattern matching* (OPPM) problem  
 85 is to find all such text positions.

## 2.1 The Problem

Given a totally ordered alphabet  $\Sigma$ , an indeterminate string is a sequence of disjunctive sets of characters  $x[1]x[2]..x[n]$  where  $x[i] \subseteq \Sigma$ . Each position is given by  $x[i]=\sigma_1..\sigma_r$  where  $r \geq 1 \wedge \sigma_i \in \Sigma$ .

Given an indeterminate string  $x$ , a *valid assignment*  $\$x$  is a (determinate) string with a single character at position  $i$ , denoted  $\$x[i]$ , contained in the  $x[i]$  set of characters, i.e.  $\$x[1] \in x[1], \dots, \$x[m] \in x[m]$ . For instance, the indeterminate string  $(1|3, 3|4, 2|3, 1|2)$  has  $2^4$  valid assignments. Given an indeterminate position  $x[i] \subseteq \Sigma$ ,  $\$x_j[i]$  is the  $j^{\text{th}}$  ordered value of  $x[i]$  (e.g.  $\$x_0[i]=1$  for  $x[i]=1|2$ ). Given an indeterminate string  $x$ , let a *partially assigned* string  $\S x$  be an indeterminate string with an arbitrary number of uncertain characters removed, i.e.  $\S x[1] \subseteq x[1], \dots, \S x[m] \subseteq x[m]$ .

Given a determinate string  $x$  of length  $m$ , an indeterminate string  $y$  of equal length is said to be *order-preserving* against  $x$ , identically denoted by  $x \approx y$ , if there is a valid assignment  $\$y$  such that the relative orders of the characters in  $x$  and  $\$y$  are the same, i.e.  $x[i] \leq x[j] \Leftrightarrow \$y[i] \leq \$y[j]$  for any  $1 \leq i, j \leq m$ . Given two indeterminate strings  $x$  and  $y$  with length  $m$ ,  $y$  preserves the orders of  $x$ ,  $x \approx y$ , if exists  $\$y$  in  $y$  that respects the orders of a valid assignment  $\$x$  in  $x$ .

A non-empty indeterminate pattern string  $p$  is said to order-preserving match (*op-match* in short) a non-empty indeterminate text string  $t$  iff there is a position  $i$  in  $t$  such that  $p \approx t[|p|+1..i]$ . The problem of *order-preserving pattern matching with character uncertainties* ( $\mu$ OPPM) problem is to find all such text positions.

To understand the complexity of the  $\mu$ OPPM problem, let us look to its solution from a naive stance yet considering state-of-the-art OPPM principles. The algorithmic proposal by Kubica et al. [41] is still up to this date the one providing a lowest bound,  $O(n+q)$ , where  $q=m$  for alphabets of size  $m^{O(1)}$  ( $q=m \lg m$  otherwise). Given a determinate string  $x$  of length  $m$ , an integer  $i$  ( $1 \leq i < m$ ) is said in the context of this work to be an *order-preserving border* of  $x$  if  $x[1..i] \approx x[m-i+1..m]$ . In this context, given a pattern string  $p$ , the orders between the characters of  $p$  are used to linearly infer the order borders. The order borders can then be used within the Knuth-Morris-Pratt algorithm to find op-matches against a text string  $t$  in linear time [41].

Given a determinate string  $p$  of length  $m$  and an indeterminate string  $t$  of length  $n$ , the previous approach is a direct candidate to the  $\mu$ OPPM problem by decomposing  $t$  in all its possible assignments,  $O(r^n)$ . Since determinate assignments to  $t$  are only relevant in the context of  $m$ -length windows, this approach can be improved to guarantee a maximum of  $O(r^m)$  assignments at each text position. Despite its simplicity, this solution is bounded by  $O(nr^m)$ . This complexity is further increased when indetermination is also considered in the pattern, stressing the need for more efficient alternatives.

## 2.2 Related work

The exact OPPM problem is well-studied in literature. Kubica et al. [41], Kim et al. [39] and Cho et al. [22] presented linear time solutions on the text length by respectively combining order-borders, rank-based prefixes and grammars with the Knuth-Morris-Pratt (KMP) algorithm [40]. Cho et al. [21], Belazzougui et al. [12], and Chhabra et al. [20] presented  $O(nm)$  algorithms that show a sublinear average complexity by either combining bad character heuristics with the Boyer-Moore algorithm [13] or applying filtration strategies. Recently, Chhabra et al. [18] proposed further principles to solve OPPM using word-size packed string matching instructions to enhance efficiency.

133 In the context of numeric strings, multiple relaxations to the exact pattern matching  
 134 problem have been pursued to guarantee that approximate matches are retrieved. In norm  
 135 matching [7, 44, 2, 47], matches between numeric strings occur if a given distance threshold  
 136  $f(x, y) \leq \theta$  is satisfied. In  $(\delta, \gamma)$ -matching [14, 26, 24, 23, 42, 43, 45], strings are matched if the  
 137 maximum difference of the corresponding characters is at most  $\delta$  and the sum of differences  
 138 is at most  $\gamma$ .

139 In the context of nominal strings, variants of the pattern matching task have also been  
 140 extensively studied to allow for don't care symbols in the pattern [37, 25, 9], transposition-  
 141 invariant [42], parameterized matching [11, 6], less than matching [1], swapped matching  
 142 [3, 46], gaps [15, 16, 31], overlap matching [5], and function matching [4, 8].

143 Despite the relevance of the aforementioned contributions to answer the exact order-  
 144 preserving pattern matching and generic pattern matching, they cannot be straightforwardly  
 145 extended to efficiently answer the  $\mu$ OPPM problem.

### 147 **3 Polynomial time $\mu$ OPPM for equal length pattern and text**

148 *Section 3.1* introduces the first efficient algorithm to solve the  $\mu$ OPPM problem when one  
 149 string is indeterminate ( $r \in \mathbb{N}^+$ ). *Section 3.2* discusses the existence of efficient solvers when  
 150 both strings are indeterminate. Based on the reducibility of the graph coloring problem to  
 151 the formulations proposed in *Section 3.2*, we hypothesize that op-matching indeterminate  
 152 strings with an arbitrary number of uncertain characters per position ( $r \in \mathbb{N}^+$ ) is in class  
 153 **NPC**. The proof of this intuition is, nevertheless, considered out of the scope, being regarded  
 154 as future work.

#### 155 **3.1 $O(mr \lg r)$ time $\mu$ OPPM when one string is indeterminate**

156 Given a determinate string  $x$  of length  $m$ , there is a well-defined permutation of positions,  
 157  $\pi$ , that specifies a non-monotonic ascending order of characters in  $x$ . For instance, given  
 158  $x=(1,4,3,1)$ , then  $x[0]=x[3]<x[2]<x[1]$  and  $\pi=(0,3,2,1)$ . Given a determinate string  $y$  with  
 159 the same length,  $y$  op-matches  $x$  if it  $y$  satisfies the same  $m-1$  orders. For instance, given  
 160  $x=(1,4,3,1)$  and  $y=(2,5,4,3)$ ,  $x$  orders are not preserved in  $y$  since  $y[0] \neq y[3] < y[2] < y[1]$ .

161 The monotonic properties can be used to answer  $\mu$ OPPM when one string is indetermin-  
 162 ate. Given an indeterminate string  $y$ , let  $x_\pi$  and  $y_\pi$  be the permuted strings in accordance  
 163 with  $\pi$  orders in  $x$ . To handle equality constraints, positions in  $y_\pi$  with identical characters  
 164 in  $x_\pi$  can be intersected, producing a new string  $y'_\pi$  with  $s$  length ( $s \leq m$ ). Illustrating, given  
 165  $x=(4,1,4,2)$  and  $y=(2|7, 2, 7|8, 1|4|8)$ , then  $\pi=(1,3,0,2)$ ,  $x_\pi=(1,2,4,4)$ ,  $y_\pi=(2, 8|4|1, 7|2, 8|7)$   
 166 and  $y'_\pi=(y_\pi[0], y_\pi[1], y_\pi[2] \cap y_\pi[3])=(2, 8|4|1, 7)$ . To handle monotonic inequalities,  $y'_\pi[i]$   
 167 characters can be concatenated in descending order to compose  $z=y'_\pi[0]y'_\pi[1]..y'_\pi[s]$  and the  
 168 orders between  $x$  and  $y$  verified by testing if the longest increasing subsequence (LIS) [29]  
 169 of  $z$  has  $s$  length. In the given example,  $z=(2, 8, 4, 1, 7)$ , and the LIS of  $z=(2, 8, 4, 1, 7)$  is  
 170  $w=(2,4,7)$ . Since  $|w|=|y'_\pi|=3$ ,  $y$  op-matches  $x$ .

171 **► Theorem 1.** *Given a determinate string  $x$  and an indeterminate string  $y$ , let  $x_\pi$  and  $y_\pi$  be*  
 172 *the sorted strings in accordance with  $\pi$  order of characters in  $x$ . Let the positions with equal*  
 173 *characters in  $x_\pi$  be intersected in  $y_\pi$  to produce a new indeterminate string  $y'_\pi$ . Consider  $z_i$*   
 174 *to be a string with  $y'_\pi[i]$  characters in descending order and  $z=z_1z_2..z_m$ , then:*

$$175 \quad |w| = |y'_\pi| \Leftrightarrow y \approx x \quad \text{where } w = \text{longest increasing subsequence in } z \quad (1)$$

176 **Proof.** ( $\Rightarrow$ ) If the length of the longest increasing subsequence (LIS),  $|w|$ , equals the number  
 177 of monotonic relations in  $x$ ,  $|y'_\pi|$ , then  $y \approx x$ . By sorting characters in descending order  
 178 per position, we guarantee that at most one character per position in  $y'_\pi$  appears in the  
 179 LIS (respecting monotonic orders in  $x$  given  $y'_\pi$  properties). By intersecting characters  
 180 in positions of  $y$  with identical characters in  $x$ , we guarantee the eligibility of characters  
 181 satisfying equality orders in  $x$ , otherwise empty positions in  $y'_\pi$  are observed and the LIS  
 182 length is less than  $|y'_\pi|$ . ( $\Leftarrow$ ) If  $|w| < |y'_\pi|$ , there is no assignment in  $y$  that op-matches  $x$   
 183 due to one of two reasons: 1) there are empty positions in  $y'_\pi$  due to the inability to satisfy  
 184 equalities in  $x$ , or 2) it is not possible to find a monotonically increasing assignment to  $y'_\pi$   
 185 and, given the properties of  $y'_\pi$ ,  $y_\pi$  cannot preserve the orders of  $x_\pi$ .  $\blacktriangleleft$

186 Solving the LIS task on a string of size  $n$  is  $O(n \lg n)$  [29] where  $n=|z|=O(rm)$ . In addi-  
 187 tion, set intersection operations are performed  $O(m)$  times on sets with  $O(r)$  size, which  
 188 can be accomplished in  $O(rm \lg r)$  time. As a result, the  $\mu$ OPPM problem with one inde-  
 189 terminate string can be solved in  $O(rm \lg(rm))$ .

190 Given the fact that the candidate string for the LIS task has properties of interest, we  
 can improve the complexity of this calculus (*Theorem 2*) in accordance with Algorithm 1.

---

**Algorithm 1:**  $O(mr \lg r)$   $\mu$ OPPM algorithm with one indeterminate string

---

```

Input: determinate  $x$ , indeterminate  $y$  ( $|x|=|y|=m$ )
 $\pi \leftarrow \text{sortedIndexes}(x);$  //  $O(m)$  if  $|\Sigma| = m_r^{O(1)}$  ( $O(m \lg m)$  otherwise)
 $x_\pi \leftarrow \text{permute}(x, \pi);$  //  $O(m + mr)$ 
 $y_\pi \leftarrow \text{permute}(y, \pi);$ 
 $j \leftarrow 0; y'_\pi[0] \leftarrow \{y_\pi[0]\};$ 
foreach  $i \in 1..m-1$  do //  $O(mr \lg r)$ 
  | if  $x_\pi[i] = x_\pi[i-1]$  then  $y'_\pi[j] \leftarrow y'_\pi[j] \cup \{y_\pi[i]\};$  //  $O(r \lg r)$ 
  | else  $j \leftarrow j+1; y'_\pi[j] \leftarrow \{y_\pi[i]\};$ 
 $s \leftarrow |y'_\pi|, \text{nextMin} \leftarrow -\infty;$ 
foreach  $i \in 0..s-1$  do //  $O(mr)$ 
  |  $\text{nextMin} \leftarrow \min\{a \mid a \in y'_\pi[i], a > \text{nextMin}\};$  //  $O(r)$ 
  | if  $\nexists \text{nextMin}$  then return false};
return true};

```

---

191

192  $\blacktriangleright$  **Theorem 2.**  $\mu$ OPPM two strings of length  $m$ , one being indeterminate, is in  $O(mr \lg r)$   
 193 time, where  $r \in \mathbb{N}^+$ .

194 **Proof.** In accordance with Algorithm 1,  $\mu$ OPPM is bounded by the verification of equalities,  
 195  $O(mr \lg r)$  [27]. Testing inequalities after set intersections can be linearly performed on the  
 196 size of  $y$ ,  $O(mr)$  time, improving the  $O(mr \lg(mr))$  bound given by the LIS calculus.  $\blacktriangleleft$

197 The analysis of Algorithm 1 further reveals that the  $\mu$ OPPM problem with one inde-  
 198 terminate string requires linear space in the text length,  $O(mr)$ .

### 199 3.2 $\mu$ OPPM with indeterminate pattern and text

200 As indetermination in real-world strings is typically observed between pairs of characters  
 201 [33], a key question is whether  $\mu$ OPPM on two indeterminate strings is in class **P** when  
 202  $r=2$ . To explore this possibility, new concepts need to be introduced. In OPPM research,  
 203 character orders in a string of length  $m$  can be decomposed in 3 sequences with  $m$  unit sets:

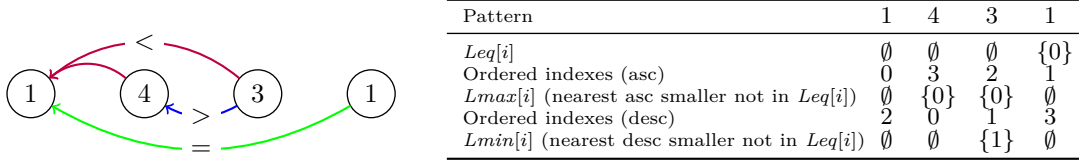
204  $\blacktriangleright$  **Definition 3.** For  $i=0, \dots, m-1$ :

- 205  $\blacksquare$   $Leq_x[i] = \{\max\{k : k < i, x[i] = x[k]\}\}$  ( $\emptyset$  if there is no eligible  $k$ )
- 206  $\blacksquare$   $Lmax_x[i] = \{\max\{\arg\max_k\{x[k] : k < i, x[i] > x[k]\}\}\}$  ( $\emptyset$  if there is no eligible  $k$ )

207 ■  $Lmin_x[i] = \{\max\{\text{argmin}_k\{x[k] : k < i, x[i] < x[k]\}\}\}$  ( $\emptyset$  if there is no eligible  $k$ )

208  $Leq$ ,  $Lmax$  and  $Lmin$  capture  $=$ ,  $>$  and  $<$  relationships between each character  $x[i]$  in  
 209  $x$  and the closest preceding character  $x[k]$ . These orders can be inferred in linear time  
 210 for alphabets of size  $m^{O(1)}$  and in  $O(m \lg m)$  time for other alphabets by answering the  
 211 “all nearest smaller values” task on the sorted indexes [41]. Fig.1 depicts  $Leq$ ,  $Lmax$  and  
 212  $Lmin$  for  $x=(1,4,3,1)$ . Given determinate strings  $x$  and  $y$ ,  $A=Leq_x[t+1]$ ,  $B=Lmax_x[t+1]$  and  
 213  $C=Lmin_x[t+1]$ , if  $x[1..t] \approx y[1..t]$  then:

214  $x[1..t+1] \approx y[1..t+1] \Leftrightarrow \forall a \in A (y[t+1] = y[a]) \wedge \forall b \in B (y[t+1] > y[b]) \wedge \forall c \in C (y[t+1] < y[c])$ . (2)



■ **Figure 1** Orders identified for  $p=(1,4,3,1)$  in accordance with Kubica et al. [41].

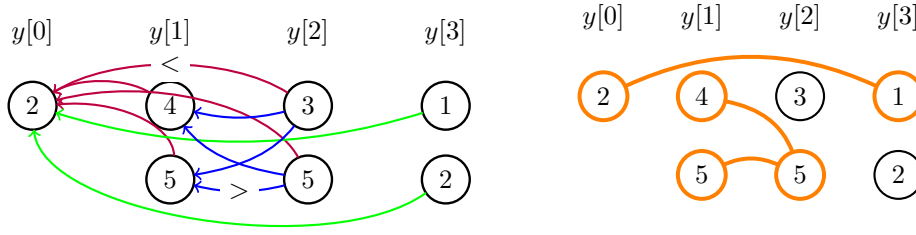
215 When allowing uncertainties between pairs of characters, previous research on the OP  
 216 PM problem cannot be straightforwardly extended due to the need to trace  $O(2^m)$  assignments  
 217 on indeterminate strings.

218 ► **Lemma 4.** *Given a determinate string  $x$ , an indeterminate string  $y$ , and the singleton*  
 219 *sets  $A=Leq_x[t+1]$ ,  $B=Lmax_x[t+1]$  and  $C=Lmin_x[t+1]$  containing a position in  $1..t$ . If*  
 220  *$x[1..t] \approx y[1..t]$  is verified on a specific assignment of  $y$  characters, denoted  $\$y$ , then:*

221  $x[1..t+1] \approx y[1..t+1] \Leftrightarrow \exists_{\$y[t+1] \in \$y[t+1]} \forall_{a \in A} \exists_{\$y[a] \in \$y[a]} \forall_{b \in B} \exists_{\$y[b] \in \$y[b]} \forall_{c \in C} \exists_{\$y[c] \in \$y[c]}$   
 222  $\$y[t+1] = \$y[a] \wedge \$y[t+1] > \$y[b] \wedge \$y[t+1] < \$y[c]$

223 **Proof.** ( $\Rightarrow$ ) In accordance with  $Leq$ ,  $Lmax$  and  $Lmin$  definition, for any  $a \in A$ ,  $b \in B$  and  $c \in C$   
 224 we have  $x[t+1]=x[a]$ ,  $x[t+1]>x[b]$  and  $x[t+1]<x[c]$ . If there is an assignment to  $y[1..t+1]$  in  
 225  $\$y$  that preserves the orders of  $x[1..t+1]$ , then for each  $a \in A$ ,  $b \in B$  and  $c \in C$   $\$y[t+1]=\$y[a]$ ,  
 226  $\$y[t+1]>\$y[b]$  and  $\$y[t+1]<\$y[c]$  (where  $\$y[t+1] \in \$y[t+1]$ ,  $\$y[a] \in \$y[a]$ ,  $\$y[b] \in \$y[b]$ ,  
 227  $\$y[c] \in \$y[c]$ ). ( $\Leftarrow$ ) We need to show that  $x[1..t+1] \approx y[1..t+1]$ . Since  $x[1..t] \approx y[1..t]$ ,  
 228 for  $i < t$ ,  $\exists_{\$y[i] \in \$y[i], \$y[t+1] \in \$y[t+1]}: x[t+1]>x[i] \Leftrightarrow \$y[t+1]>\$y[i]$ . Assuming  $x[t+1]>x[i]$   
 229 for some  $i \in \{1..t\}$ : by the definition of  $Lmax$ ,  $\forall_{b \in B} x[b]>x[i]$ ; by the order-isomorphism of  
 230  $x[1..t]$  and  $\$y[1..t]$  in  $\$y[1..t]$ , there is  $\$y[i] \in \$y[i]$  and  $\$y[b] \in \$y[b]$  that  $\forall_{b \in B} \$y[b]>\$y[i]$ ;  
 231 and by the assumption of the lemma,  $\forall_{b \in B} \$y[t+1]>\$y[b]$ ; hence  $\$y[t+1]>\$y[i]$ . Similarly,  
 232  $x[t+1]<x[i]$  (and  $x[t+1]=x[i]$ ) implies  $\$y[t+1]<\$y[i]$  (and  $\$y[t+1]=\$y[i]$ ), yielding the  
 233 stated equivalence. ◀

234 Given two strings of equal length, the  $\mu$ OPPM problem can be schematically represented  
 235 according to the identified order restrictions. Fig.2 represents restrictions on the indeterm-  
 236 inate string  $y=(2,4|5,3|5,1|2)$  in accordance with the observed orders in  $x=(1,4,3,1)$ . The  
 237 left side edges are placed in accordance with *Lemma 4* and capture assessments on the or-  
 238 ders between pairs of characters. The right side edges capture incompatibilities detected  
 239 after the assessments, i.e. pairs of characters that cannot be selected simultaneously (for  
 240 instance,  $y[0]=2$  and  $y[3]=1$ , or  $y[1]=4$  and  $y[2]=5$ ). For the given example, there are two  
 241 valid assignments,  $\$y_1=(2,4,3,2)$  and  $\$y_2=(2,5,3,2)$ , that satisfy  $x[0]=x[3]<x[2]<x[1]$ , thus  
 242  $y$  op-matches  $x$ .



■ **Figure 2** Schematic representation of the pairwise ordering restrictions for text  $y=(2,4|5,3|5,1|2)$  and pattern  $x=(1,4,3,1)$ . In the left side, all order verifications are represented, while in the right side only the order conflicts are signaled (e.g.  $y[1]=4$  cannot be selected together with  $y[2]=5$ ).

243 To verify whether there is an assignment that satisfies the identified ordering restrictions,  
244 we propose the reduction of  $\mu$ OPPM problem to a Boolean satisfiability problem.

245 Given a set of Boolean variables, a formula in conjunctive normal form is a conjunction of  
246 clauses, where each clause is a disjunction of literals, and a literal corresponds to a variable  
247 or its negation. Let a 2CNF formula be a formula in the conjunctive normal form with at  
248 most two literals per clause. Given a CNF formula, the *satisfiability* (SAT) problem is to  
249 verify if there is an assigning of values to the Boolean variables such that the CNF formula  
250 is satisfied.

251 ► **Theorem 5.** *The  $\mu$ OPPM problem over two strings of equal length, one being indeterm-*  
252 *inate, can be reduced to a satisfiability problem with the following CNF formula:*

$$\begin{aligned}
 253 \quad \phi = & \bigwedge_{i=0}^{m-1} \left( \bigvee_{\$y[i] \in y[i]} z_{i, \$y[i]} \right) \wedge \bigwedge_{i=0}^{m-1} \left( \bigwedge_{\$y[i] \in y[i]} \bigwedge_{j \in Leq[i], \$y[j] \in y[j]} (\neg z_{i, \$y[i]} \vee \neg z_{j, \$y[j]} \vee \$y[i] = \$y[j]) \right) \\
 & \wedge \bigwedge_{\$y[i] \in y[i]} \bigwedge_{\substack{j \in Lmax[i] \\ \$y[j] \in y[j]}} (\neg z_{i, \$y[i]} \vee \neg z_{j, \$y[j]} \vee \$y[i] > \$y[j]) \wedge \bigwedge_{\$y[i] \in y[i]} \bigwedge_{\substack{j \in Lmin[i] \\ \$y[j] \in y[j]}} (\neg z_{i, \$y[i]} \vee \neg z_{j, \$y[j]} \vee \$y[i] < \$y[j])
 \end{aligned}
 \tag{3}$$

255 **Proof.** Let us show that if  $x$  op-matches  $y$  then  $\phi$  is satisfiable, and if  $x$  does not op-match  
256  $y$  then  $\phi$  is not satisfiable. ( $\Rightarrow$ ) When  $x \approx y$ , there is an assignment of values to  $y$ ,  $\$y$ ,  
257 that satisfy the orderings of  $x$ .  $\phi$  is satisfiable if there is at least one variable assigned  
258 to true per clause  $\bigvee_{\$y[i] \in y[i]} z_{i, \$y[i]}$  given conflicts  $\neg z_{i, \$y[i]} \vee \neg z_{j, \$y[j]}$ . As conflicts do not  
259 prevent the existence of a valid assignment (by assumption), then  $\exists_{\$y} \wedge_{i \in \{0..m-1\}} z_{i, \$y[i]}$  and  
260  $\phi$  is satisfiable. ( $\Leftarrow$ ) When  $x$  does not op-match  $y$ , there is no assignment of values  $\$y \in y$   
261 that can satisfy the orders of  $x$ . Per formulation, the conflicts  $\neg z_{i, \$y[i]} \vee \neg z_{j, \$y[j]}$  prevent the  
262 satisfiability of one or more clauses  $\bigvee_{\$y[i] \in y[i]} z_{i, \$y[i]}$ , leading to a non-satisfiable formula. ◀

263 If the established  $\phi$  formula is satisfiable, there is a Boolean assignment to the variables  
264 that specify an assignment of characters in  $y$ ,  $\$y$ , preserving the orders of  $x$  (as defined  
265 by *Leq*, *Lmax* and *Lmin*). Otherwise, it is not possible to select an assignment  $\$y$  op-  
266 matching  $x$ .  $\phi$  has at most  $r \times m$  variables,  $\{z_{i, \sigma} \mid i \in \{0..m-1\}, \sigma \in \Sigma\}$ . The Boolean value  
267 assigned to a variable  $z_{i, \sigma}$  simply defines that the associated character  $\sigma$  from  $y[i]$  can be  
268 either considered (when true) or not (when false) to compose a valid assignment  $\$y$  that  
269 op-matches the given determinate string  $x$ . The reduced (3) formula is composed of two  
270 major types of clauses:  $\bigvee_{\$y[i] \in y[i]} z_{i, \$y[i]}$ , and  $(\neg z_{i, \$y[i]} \vee \neg z_{j, \$y[j]} \vee \mathbf{bool})$  where **bool** is either  
271 given by  $\$y[i] = \$y[j]$ ,  $\$y[i] < \$y[j]$  or  $\$y[i] > \$y[j]$ . Clauses of the first type specify the need to  
272 select at least one character per position in  $y$  to guarantee the presence of valid assignments.

273 The remaining clauses specify ordering constraints between characters. If an inequality,  
 274 such as  $\$y[i] > \$y[j]$ , is assessed as true, the associated clause is removed. Otherwise,  
 275  $(\neg z_{i,\sigma_1} \vee \neg z_{j,\sigma_2})$  is derived, meaning that these  $\sigma_1$  and  $\sigma_2$  characters should not be selected  
 276 simultaneously since they do not satisfy the orders defined by a given pattern. For instance,  
 277 the pairs of characters in orange from Fig.2 should not be simultaneously selected due to  
 278 order conflicts. To this end,  $(\neg z_{0,2} \vee \neg z_{3,1})$  and  $(\neg z_{1,4} \vee \neg z_{2,5})$  clauses need to be included to  
 279 verify if  $y \approx x$ . Considering  $y=(2,4|5,4|5,1|2)$  and  $x=(1,4,3,1)$ , schematically represented  
 280 in Fig.2, the associated CNF formula is:

$$281 \quad \phi = z_{0,2} \wedge (z_{1,4} \vee z_{1,5}) \wedge (z_{2,4} \vee z_{2,5}) \wedge (z_{3,1} \vee z_{3,2}) \wedge (\neg z_{0,2} \vee \neg z_{3,1}) \wedge (\neg z_{1,4} \vee \neg z_{2,5})$$

282 ► **Theorem 6.** *Given two strings of length  $m$ , one being indeterminate with  $r=2$ , the  $\mu$ OPPM*  
 283 *problem can be reduced to a 2SAT problem with a CNF formula with  $O(m)$  size.*

284 **Proof.** Given *Theorem 5* and the fact that the reduced CNF formula has at most two literals  
 285 per clause –  $\phi$  is a composition of  $\vee_{\$y[i] \in y[i]} z_{i,\$y[i]}$  clauses with  $|y[i]| \in \{1, 2\}$  and  $(\neg z_{i,\$y[i]} \vee$   
 286  $\neg z_{j,\$y[j]} \vee \text{bool})$  clauses –  $\mu$ OPPM with  $r=2$  and one indeterminate string is reducible to  
 287 2SAT. The reduced formula has at most  $10m$  clauses with 2 literals each, being linear in  $m$ :

- 288 ■ [clauses that impose the selection of at least one character per position in  $y$ ] Since  $y$   
 289 has  $m$  positions, and each position is either determinate (unitary clause) or defines an  
 290 uncertainty between a pair of characters, there are  $m$  clauses and at most  $2m$  literals;
- 291 ■ [clauses that define the ordering restrictions between two variables] A position in the  
 292 indeterminate string  $y[i]$  needs to satisfy at most two order relations. Considering that  
 293  $i$ ,  $Leq[i]$ ,  $Lmax[i]$  and  $Lmin[i]$  specify uncertainties between pairs of characters, there  
 294 are up to 12 restrictions per position: 4 ordering restrictions between characters in  $y[i]$   
 295 and  $y[Leq[i]]$ ,  $y[Lmax[i]]$  and  $y[Lmin[i]]$ . Whenever the order between two characters is  
 296 not satisfied, a clause is added per position, leading to at most  $12m$  clauses. ◀

298 ► **Theorem 7.** *The  $\mu$ OPPM between determinate and indeterminate strings of equal length*  
 299 *can be solved in linear time when  $r=2$ .*

300 **Proof.** Given the fact that a 2SAT problem can be solved in linear time  $[10]^*$ , this proof  
 301 directly derives from *Theorem 6* as it guarantees the soundness of reducing  $\mu$ OPPM ( $r=2$ )  
 302 to a 2SAT problem with a CNF formula with  $O(m)$  size. ◀

303 As the size of the mapped CNF formula  $\phi$  is  $O(m)$  and the a valid algorithm to verify  
 304 its satisfiability would require the construction of a graph with  $O(m)$  nodes and edges, the  
 305 required memory for the target  $\mu$ OPPM problem is  $\Theta(m)$ .

306 When moving from one to two indeterminate strings, previous contributions are insuf-  
 307 ficient to answer the  $\mu$ OPPM problem. In this context, the  $Leq$ ,  $Lmax$  and  $Lmin$  vectors  
 308 need to be redefined to be inferred from an indeterminate string:

309 ► **Definition 8.**

- 310 ■  $Leq_x[i|j]=\{k : k < i, \exists_p \$x_j[i]=\$x_p[k]\}$  ( $\emptyset$  if there is no eligible  $k$ ), for  $i=0, \dots, m-1$

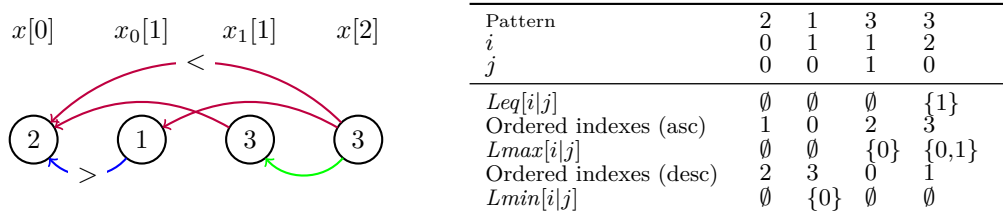
---

\*2SAT problems have linear time and space solutions on the size of the input formula. Consider for instance the original proposal [10], the formula  $\phi$  is modeled by a directed graph  $G=(V, E)$ , with two nodes per variable  $z_i$  in  $\phi$  ( $z_i$  and  $\neg z_i$ ) and two directed edges for each clause  $z_i \vee z_j$  (the equivalent implicative forms  $\neg z_i \Rightarrow z_j$  and  $\neg z_j \Rightarrow z_i$ ). Given  $G$ , the strongly connected components (SCCs) of  $G$  can be discovered in  $O(|V| + |E|)$ . During the traversal if a variable and its complement belong to the same SCC, then the procedure stops as  $\phi$  is determined to be unsatisfiable. Given the fact that both  $V=O(m)$  and  $E=O(m)$  by *Lemma 6*, this procedure is  $O(m)$  time and space.



- 311 ■  $Lmax_x[i|j]=\{k : k < i, \exists_p \$x_j[i] > \$x_p[k]\}$  ( $\emptyset$  if there is no eligible  $k$ ), for  $i=0, \dots, m-1$
- 312 ■  $Lmin_x[i|j]=\{k : k < i, \exists_p \$x_j[i] < \$x_p[k]\}$  ( $\emptyset$  if there is no eligible  $k$ ), for  $i=0, \dots, m-1$

313 Fig.3 schematically represents the order relationships of  $x=(2, 1|3, 3)$  and the associated  
 314  $Leq$ ,  $Lmax$  and  $Lmin$  vectors. In this scenario,  $x[2]$  needs to be verified not only against  
 315  $x_0[1]$  but also against  $x_1[1]$  in case  $x_0[1]$  is disregarded.



■ **Figure 3** Order relationships of  $x=(2, 1|3, 3)$  and the corresponding  $Lmax$  and  $Lmin$  vectors.

316 ► **Corollary 9.** Given  $Leq$ ,  $Lmax$  and  $Lmin$  (Def.8), there are  $O((rm)^2)$  order relationships  
 317 when  $r \in \mathbb{N}^+$  since each character in a given position establishes at most  $O(m)$  relationships  
 318 with characters in preceding positions.

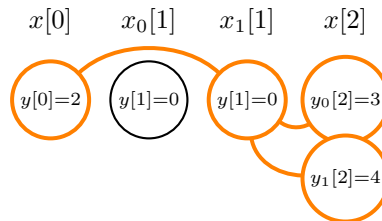
319 ► **Lemma 10.** Given indeterminate strings  $x$  and  $y$ , let  $A_j=Leq_x[t+1|j]$ ,  $B_j=Lmax_x[t+1|j]$   
 320 and  $C_j=Lmin_x[t+1|j]$  (Def.8) be the orders associated with  $\$x_j[t+1]$ . If  $x[1..t] \approx y[1..t]$  is  
 321 verified on a partial assignment of  $y$  characters, denoted by  $\$y$ , then:

$$322 \quad x[1..t+1] \approx y[1..t+1] \Leftrightarrow \exists_{j \in \{0,1\}} \exists_{\$y[t+1] \in \$y[t+1]} \forall_{a \in A_j, b \in B_j, c \in C_j}$$

$$323 \quad \exists_{\$y[a] \in \$y[a], \$y[b] \in \$y[b], \$y[c] \in \$y[c]} (\$y[t+1] = \$y[a] \wedge \$y[t+1] > \$y[b] \wedge \$y[t+1] < \$y[c])$$

324 **Proof.** ( $\Rightarrow$ ) Similar to the proof of Lemma 4, yet  $A$ ,  $B$  and  $C$  conditional to  $x[t+1]$  (Def.3)  
 325 are now given by  $A_j$ ,  $B_j$  and  $C_j$  conditional to  $x_j[t+1]$  (Def.8). If there is an assignment  
 326 to  $y[1..t+1]$  in  $\$y$  that preserves one of the possible orders in  $x[1..t+1]$ , then for any  
 327  $a \in A_j$ ,  $b \in B_j$  and  $c \in C_j$ :  $\$y[t+1]=\$y[a]$ ,  $\$y[t+1]>\$y[b]$  and  $\$y[t+1]<\$y[c]$  (where  
 328  $\$y[t+1] \in \$y[t+1]$ ,  $\$y[a] \in \$y[a]$ ,  $\$y[b] \in \$y[b]$ ,  $\$y[c] \in \$y[c]$ ).

329 ( $\Leftarrow$ ) We need to show that  $x[1..t+1] \approx y[1..t+1]$ . Since  $x[1..t] \approx y[1..t]$ , it is sufficient  
 330 to prove that for  $i \leq t$ : exists  $\$x[i] \in \$x[i]$ ,  $\$x[t+1] \in \$x[t+1]$ ,  $\$y[i] \in \$y[i]$ ,  $\$y[t+1] \in$   
 331  $\$y[t+1]$  such that  $\$x[t+1]=\$x[i] \Leftrightarrow \$y[t+1]=\$y[i]$ ,  $\$x[t+1]>\$x[i] \Leftrightarrow \$y[t+1]>\$y[i]$  and  
 332  $\$x[t+1]<\$x[i] \Leftrightarrow \$y[t+1]<\$y[i]$ . This results from Def.8, the order-isomorphism property  
 333 and Lemma 4. ◀



■ **Figure 4** Conflicts when op-matching  $y=(2, 0, 3|4)$  against  $x=(2, 1|3, 3)$ .

334 Fig.4 represents encountered restrictions when op-matching  $x=(2, 1|3, 3)$  against  $y=(2, 0, 3|4)$ .

335 The right side edges capture the detected incompatibilities, i.e. pairs of characters that can-  
 336 not be selected simultaneously. For the given example, there are 2 valid assignments –  
 337  $\$y_1=(2,0,3)$  and  $\$y_2=(2,0,4)$  – satisfying  $\$x_0[1]<\$x_0[0]<\$x_0[2]$ , thus  $x \approx y$ .

338 To verify whether there is an assignment that satisfies the identified ordering restrictions,  
 339 Theorem 11 extends the previously introduced SAT mapping given by (3).

340 ► **Theorem 11.** *Given  $Leq$ ,  $Lmax$  and  $Lmin$  (Def.8),  $\mu OPPM$  problem over two indeterm-  
 341 inate strings of equal length can be reduced to a satisfiability problem with the following CNF  
 342 formula:*

$$\begin{aligned}
 & \bigwedge_{i=0}^{m-1} \bigvee_{\substack{\$y[j] \in y[j] \\ \$x[j] \in x[j]}} z_{i, \$x[i], \$y[i]} \wedge \bigwedge_{i=0}^{m-1} \bigwedge_{\substack{\$y[j] \in y[j] \\ \$x[j] \in x[j]}} \left( \bigwedge_{j \in Leq[i]} \bigwedge_{\substack{\$y[j] \in y[j] \\ \$x[j] \in x[j]}} (\neg z_{i, \$x[i], \$y[i]} \vee \neg z_{j, \$x[j], \$y[j]} \vee \$y[i] = \$y[j]) \right) \\
 & \bigwedge_{j \in Lmax[i]} \bigwedge_{\substack{\$y[j] \in y[j] \\ \$x[j] \in x[j]}} \left( \neg z_{i, \$x[i], \$y[i]} \vee \neg z_{j, \$x[j], \$y[j]} \vee \$y[i] > \$y[j] \right) \wedge \bigwedge_{j \in Lmin[i]} \bigwedge_{\substack{\$y[j] \in y[j] \\ \$x[j] \in x[j]}} \left( \neg z_{i, \$x[i], \$y[i]} \vee \neg z_{j, \$x[j], \$y[j]} \vee \$y[i] < \$y[j] \right)
 \end{aligned}
 \tag{4}$$

345 **Proof.** If  $x \approx y$  then  $\phi$  is satisfiable, and if  $x$  does not op-match  $y$  then  $\phi$  is not satisfiable.

346 ( $\Rightarrow$ ) When  $x$  op-matches  $y$ , there is an assignment of values in  $x$  and  $y$  such that  $\$x \approx \$y$ .  
 347  $\phi$  is satisfiable if there is at least one valid assignment  $z_{i, \$x[i], \$y[i]}$  per  $i^{th}$  position. As  
 348 conflicts  $\neg z_{i, \$x[i], \$y[i]} \vee \neg z_{j, \$x[j], \$y[j]}$  do not prevent the existence of a valid assignment (by  
 349 assumption), one or more variables  $z_{i, \$x[i], \$y[i]}$  can be selected per position.  $\phi$  can then be  
 350 satisfied by fixing a single variable  $z_{i, \$x[i], \$y[i]}$  per  $i^{th}$  position as true and the remaining  
 351 variables as false. ( $\Leftarrow$ ) When  $x$  does not op-match  $y$ , there is no assignment of values  $\$x \in x$   
 352 and  $\$y \in y$  such that  $\$x \approx \$y$ . Per formulation, in the absence of an order-preserving match,  
 353 conflicts will prevent the assignment of at least one variable  $z_{i, \$x[i], \$y[i]}$  per  $i^{th}$  position, thus  
 354 making  $\phi$  formula unsat. ◀

355 If (4) formula is satisfiable, there is a Boolean assignment to the variables such that  
 356 there is an assignment of characters in  $y$ ,  $\$y$ , and in  $x$ ,  $\$x$ , such that both strings op-  
 357 match. Otherwise, it is not possible to select assignments such that  $x \approx y$ . Given  $r=2$ , the  
 358 established  $\phi$  formula has at most  $4m$  variables,  $\{z_{i, \sigma_1, \sigma_2} \mid i \in \{0..m-1\}, \sigma_1, \sigma_2 \in \Sigma\}$ . The  
 359 Boolean values assigned to these variables define whether characters  $\sigma_1 \in x[i]$  and  $\sigma_2 \in y[i]$   
 360 belong to an op-match. The reduced formula is composed of two major types of clauses:

- 361 ■ (4.1) at least one combination of characters,  $\$x[i]$  and  $\$y[i]$ , should be selected per  $i^{th}$   
 362 position;
- 363 ■ (4.2) clauses specify ordering constraints between pairs of characters  $\sigma_1 \in y[i]$  and  
 364  $y[Leq[i]]$ ,  $y[Lmax[i]]$  and  $y[Lmin[i]]$ . If the inequalities  $\$y[i]=\$y[j]$ ,  $\$y[i]>\$y[j]$  and  
 365  $\$y[i]<\$y[j]$  are assessed as false, these leads to clauses of the form  $(\neg z_{i, \sigma_1} \vee \neg z_{j, \sigma_2})$ ,  
 366 meaning that these characters should not be selected simultaneously in the given posi-  
 367 tions (see Fig.4).

368 To instantiate the proposed mapping, consider  $x=(2, 1|3, 3)$  and  $y=(2, 0, 3|4)$ , schemat-  
 369 ically represented in Fig.3. The associated CNF formula is:

$$\begin{aligned}
 \phi &= z_{0,2,2} \wedge (z_{1,1,0} \vee z_{1,3,0}) \wedge (z_{2,3,3} \vee z_{2,3,4}) // (4.1) \\
 & \wedge (\neg z_{0,2,0} \vee \neg z_{1,3,0}) \wedge (\neg z_{1,3,0} \vee \neg z_{2,3,3}) \wedge (\neg z_{1,3,0} \vee \neg z_{2,3,4}) // 4.2
 \end{aligned}$$

372 ► **Theorem 12.** *The  $\mu$ OPPM problem for two indeterminate strings of equal length is reducible into a satisfiability problem over a CNF formula with  $O((mr)^2)$  size.*  
 373

374 **Proof.** The reduced formula (4) is in the two conjunctive normal form (CNF) with at most  
 375  $4m$  clauses given by (4.1) and a maximum of  $O(mr)$  orders per position (*Corollary 9*),  
 376 totalling at most  $O((mr)^2)$  order conflicts between characters (4.2). ◀

377 Given the unique properties of the above satisfiability formulation, effective backtracking  
 378 in accordance with (4.1), as well as dedicated conflict pruning principles derived from (4.2),  
 379 can be considered to develop efficient SAT solvers able to solve the  $\mu$ OPPM problem.

#### 380 4 Polynomial time $\mu$ OPPM

382 ► **Lemma 13.** *Given a pattern string of length  $m$  and a text string of length  $n$ , one being  
 383 indeterminate, the  $\mu$ OPPM problem can be solved in  $O(nmr \lg r)$  time.*

384 **Proof.** From *Lemma 7*, verifying if two strings of length  $m$  op-match is in  $O(mr \lg r)$  time  
 385 (indetermination in one string) since at most  $n-m+1$  verifications need to be performed. ◀

386 *Lemma 13* confirms that the  $\mu$ OPPM problem with one indeterminate strings is in class  
 387 **P**. This lemma further triggers the research question “*Is  $O(nmr)$  a tight bound to solve the*  
 388  *$\mu$ OPPM?*”, here left as an open research question.

389 Irrespectively of the answer, the analysis of the average complexity is of complementary  
 390 relevance. State-of-the-art research on the exact OPPM problem shows that the average  
 391 performance of algorithms in  $O(nm)$  time can outperform linear algorithms [20, 17, 19].

392 Motivated by the evidence gathered by these works, we suggest the use of filtration  
 393 procedures to improve the average complexity of the proposed  $\mu$ OPPM algorithm while  
 394 still preserving its complexity bounds. A filtration procedure encodes the input pattern  
 395 and text, and relies on this encoding to efficiently find positions in the text with a high  
 396 likelihood to op-match a given pattern. Despite the diversity of string encodings, simplistic  
 397 binary encodings are considered to be the state-of-the-art in OPPM research [20, 17]. In  
 398 accordance with Chhabra et al. [20], a pattern  $p$  can be mapped into a binary string  $p'$   
 399 expressing increases (1), equalities (0) and decreases (0) between subsequent positions. By  
 400 searching for exact pattern matches of  $p'$  in an analogously transformed text string  $t'$ , we  
 401 guarantee that the verification of whether  $p[0..m-1]$  and  $t[i..i+m-1]$  orders are preserved  
 402 is only performed when exact binary matches occur. Illustrating, given  $p=(3,1,2,4)$  and  
 403  $t=(2,4,3,5,7,1,4,8)$ , then  $p'=(1,0,1,1)$  and  $t'=(1,1,0,1,1,0,1,1,0)$ , revealing two matches  $t'[1..4]$   
 404 and  $t'[4..7]$ : one spurious match  $t[1..5]$  and one true match  $t[4..8]$ .

405 When handling indeterminate strings the concept of increase, equality and decrease needs  
 406 to be redefined. Given an indeterminate string  $x$ , consider  $x'[i]=1$  if  $\max(x[i]) < \min(x[i+1])$ ,  
 407  $x'[i]=0$  if  $\min(x[i]) \geq \max(x[i+1])$ , and  $x'[i]=*$  otherwise. Under this encoding, the pattern  
 408 matching problem is identical under the additional guard that a character in  $p'$  always  
 409 matches a don't care position,  $t'[i]=*$ , and vice-versa. Illustrating, given  $p=(6,2|3,5)$  and  
 410  $t=(3|4,5,6|8,6|7,3,5,4|6,7|8,4)$ , then  $p'=(0,1)$  and  $t'=(11*01*10)$ , leading to one true match  
 411  $t[3..5]$  – e.g.  $\$t[3..5]=(6,3,5)$  – and one spurious match  $t[5..7]$ . Exact pattern matching  
 412 algorithms, such as Knuth-Morris-Pratt and Boyer-Moore, can be adapted to consider don't  
 413 care positions while preserving complexity bounds [40, 13].

414 The properties of the proposed encoding guarantee that the exact matches of  $p'$  in  $t'$   
 415 cannot skip any op-match of  $p$  in  $t$ . Thus, when combining the premises of *Lemma 13* with  
 416 the previous observation, we guarantee that the computed  $\mu$ OPPM solution is sound.

417 The application of this simple filtration procedure prevents the recurring  $O(mr \lg r)$   
 418 verifications  $n-m+1$  times. Instead, the complexity of the proposed method to solve the  
 419  $\mu$ OPPM problem becomes  $O(dmr \lg r + n)$  (when one string is indeterminate) where  $d$  is  
 420 the number of exact matches ( $d \ll n$ ). According to previous work on exact OPPM with  
 421 filtration procedures [20], SBNDM2 and SBNDM4 algorithms [28] (Boyer-Moore variants)  
 422 were suggested to match binary encodings. In the presence of small patterns, Fast Shift-Or  
 423 (FSO) [30] can be alternatively applied [20].

424 A given string text can be read and encoded incrementally from the standard input as  
 425 needed to perform  $\mu$ OPPM, thus requiring  $O(mr)$  space. When filtration procedures are con-  
 426 sidered, the aforementioned algorithms for exact pattern matching require  $O(m)$  space [20],  
 427 thus  $\mu$ OPPM space requirements are bound by substring verifications (*Section 3*):  $O(mr)$   
 428 space when one string is indeterminate and  $O((mr)^2)$  when indetermination is considered  
 429 on both strings.

## 430 5 Concluding remark

431 This work addressed the relevant yet scarcely studied problem of finding order-preserving  
 432 pattern matches on indeterminate strings ( $\mu$ OPPM). We showed that the problem has a  
 433 linear time and space solution when one string is indeterminate. In addition, the  $\mu$ OPPM  
 434 problem (when both strings are indeterminate) was mapped into a satisfiability formula of  
 435 polynomial size and two simple types of clauses in order to study efficient solvers for the  
 436  $\mu$ OPPM problem. Finally, we showed that solvers of the  $\mu$ OPPM problem can be boosted  
 437 in the presence of filtration procedures.

438 **Acknowledgments.** This work was developed in the context of a secondment granted by the BIRDS  
 439 MASC RISE project funded in part by EU H2020 research and innovation programme under the  
 440 Marie Skłodowska-Curie grant agreement no.690941. This work was further supported by national  
 441 funds through Fundação para a Ciência e Tecnologia (FCT) with reference UID/CEC/50021/2013.

---

## 442 References

- 443 1 A. Amir and M. Farach. Efficient 2-dimensional approximate matching of half-rectangular  
 444 figures. *Information and Computation*, 118(1):1 – 11, 1995.
- 445 2 Amihood Amir, Yonatan Aumann, Piotr Indyk, Avivit Levy, and Ely Porat. Efficient  
 446 computations of  $l_1$  and  $l_\infty$  rearrangement distances. *Theoretical Computer Science*,  
 447 410(43):4382 – 4390, 2009.
- 448 3 Amihood Amir, Yonatan Aumann, Gad M. Landau, Moshe Lewenstein, and Noa Lewen-  
 449 stein. Pattern matching with swaps. *Journal of Algorithms*, 37(2):247 – 266, 2000.
- 450 4 Amihood Amir, Yonatan Aumann, Moshe Lewenstein, and Ely Porat. Function matching.  
 451 *SIAM Journal on Computing*, 35(5):1007–1022, 2006.
- 452 5 Amihood Amir, Richard Cole, Ramesh Hariharan, Moshe Lewenstein, and Ely Porat. Over-  
 453 lap matching. *Information and Computation*, 181(1):57 – 74, 2003.
- 454 6 Amihood Amir, Martin Farach, and S. Muthukrishnan. Alphabet dependence in paramet-  
 455 erized matching. *Information Processing Letters*, 49(3):111 – 115, 1994.
- 456 7 Amihood Amir, Ohad Lipsky, Ely Porat, and Julia Umanski. Approximate matching in  
 457 the  $l_1$  metric. In *CPM*, volume 5, pages 91–103. Springer, 2005.
- 458 8 Amihood Amir and Igor Nor. Generalized function matching. *Journal of Discrete Al-*  
 459 *gorithms*, 5(3):514 – 523, 2007. Selected papers from Ad Hoc Now 2005.

- 460 **9** Alberto Apostolico. Algorithms and theory of computation handbook. chapter General  
461 Pattern Matching, pages 15–15. Chapman & Hall/CRC, 2010.
- 462 **10** Bengt Aspvall, Michael F Plass, and Robert Endre Tarjan. A linear-time algorithm for  
463 testing the truth of certain quantified boolean formulas. *Information Processing Letters*,  
464 8(3):121–123, 1979.
- 465 **11** Brenda S Baker. A theory of parameterized pattern matching: algorithms and applications.  
466 In *ACM symposium on Theory of computing*, pages 71–80. ACM, 1993.
- 467 **12** Djamal Belazzougui, A. Pierrot, M. Raffinot, and Stéphane Vialette. Single and multiple  
468 consecutive permutation motif search. In *Int. Symposium on Algorithms and Computation*,  
469 pages 66–77. Springer, 2013.
- 470 **13** Robert S Boyer and J Strother Moore. A fast string searching algorithm. *Communications*  
471 *of the ACM*, 20(10):762–772, 1977.
- 472 **14** Emiliios Cambouropoulos, M. Crochemore, C. Iliopoulos, L. Mouchard, and Yoan Pinzon.  
473 Algorithms for computing approximate repetitions in musical sequences. *Int. Journal of*  
474 *Computer Mathematics*, 79(11):1135–1148, 2002.
- 475 **15** Domenico Cantone, Salvatore Cristofaro, and Simone Faro. An efficient algorithm for  $\delta$ -  
476 approximate matching with  $\alpha$ -bounded gaps in musical sequences. In *IW on Experimental*  
477 *and Efficient Algorithms*, pages 428–439. Springer, 2005.
- 478 **16** Domenico Cantone, Salvatore Cristofaro, and Simone Faro. On tuning the  $(\delta, \alpha)$ -sequential-  
479 sampling algorithm for  $\delta$ -approximate matching with alpha-bounded gaps in musical se-  
480 quences. In *ISMIR*, pages 454–459, 2005.
- 481 **17** Domenico Cantone, Simone Faro, and M Oguzhan Külekci. An efficient skip-search ap-  
482 proach to the order-preserving pattern matching problem. In *Stringology*, pages 22–35,  
483 2015.
- 484 **18** Tamanna Chhabra, Simone Faro, M. Oğuzhan Külekci, and Jorma Tarhio. Engineering  
485 order-preserving pattern matching with simd parallelism. *Softw. Pract. Exper.*, 47(5):731–  
486 739, May 2017.
- 487 **19** Tamanna Chhabra, M Oguzhan Külekci, and Jorma Tarhio. Alternative algorithms for  
488 order-preserving matching. In *Stringology*, pages 36–46, 2015.
- 489 **20** Tamanna Chhabra and Jorma Tarhio. A filtration method for order-preserving matching.  
490 *Information Processing Letters*, 116(2):71 – 74, 2016.
- 491 **21** Sukhyeun Cho, Joong Chae Na, Kunsoo Park, and Jeong Seop Sim. Fast order-preserving  
492 pattern matching. In *Combinatorial Optimization and Applications*, pages 295–305.  
493 Springer, 2013.
- 494 **22** Sukhyeun Cho, Joong Chae Na, Kunsoo Park, and Jeong Seop Sim. A fast algorithm for  
495 order-preserving pattern matching. *Information Processing Letters*, 115(2):397–402, 2015.
- 496 **23** Peter Clifford, Raphaël Clifford, and Costas Iliopoulos. Faster algorithms for  $\delta, \gamma$ -matching  
497 and related problems. In *Annual Symposium on Combinatorial Pattern Matching*, pages  
498 68–78. Springer, 2005.
- 499 **24** Raphaël Clifford and C Iliopoulos. Approximate string matching for music analysis. *Soft*  
500 *Computing-A Fusion of Foundations, Methodologies and Applications*, 8(9):597–603, 2004.
- 501 **25** Richard Cole, C. Iliopoulos, T. Lecroq, W. Plandowski, and Wojciech Rytter. On special  
502 families of morphisms related to  $\delta$ -matching and don’t care symbols. *Information Processing*  
503 *Letters*, 85(5):227–233, 2003.
- 504 **26** Maxime Crochemore, Costas S Iliopoulos, Thierry Lecroq, Wojciech Plandowski, and Wo-  
505 jciech Rytter. Three heuristics for delta-matching: delta-bm algorithms. In *CPM*, pages  
506 178–189. Springer, 2002.

- 507 **27** Erik D Demaine, Alejandro López-Ortiz, and J Ian Munro. Adaptive set intersections,  
508 unions, and differences. In *In Proceedings of the 11th Annual ACM-SIAM Symposium on*  
509 *Discrete Algorithms (SODA)*. Citeseer, 2000.
- 510 **28** Branislav Ďurian, Jan Holub, Hannu Peltola, and Jorma Tarhio. Improving practical exact  
511 string matching. *Information Processing Letters*, 110(4):148–152, 2010.
- 512 **29** Michael L. Fredman. On computing the length of longest increasing subsequences. *Discrete*  
513 *Mathematics*, 11(1):29 – 35, 1975.
- 514 **30** Kimmo Fredriksson and Szymon Grabowski. Practical and optimal string matching. In  
515 *SPIRE*, volume 3772, pages 376–387. Springer, 2005.
- 516 **31** Kimmo Fredriksson and Szymon Grabowski. Efficient algorithms for pattern matching  
517 with general gaps, character classes, and transposition invariance. *Information Retrieval*,  
518 11(4):335–357, 2008.
- 519 **32** Xianping Ge. Pattern matching in financial time series data. *final project report for ICS*,  
520 278, 1998.
- 521 **33** Rui Henriques. *Learning from High-Dimensional Data using Local Descriptive Models*. PhD  
522 thesis, Instituto Superior Tecnico, Universidade de Lisboa, Lisboa, 2016.
- 523 **34** Rui Henriques, Cláudia Antunes, and Sara C. Madeira. Methods for the efficient discovery  
524 of large item-indexable sequential patterns. In *New Frontiers in Mining Complex Patterns*,  
525 volume 8399 of *LNCS*, pages 100–116. Springer International Publishing, 2014.
- 526 **35** Rui Henriques and Sara C Madeira. Bicspam: flexible biclustering using sequential pat-  
527 terns. *BMC bioinformatics*, 15(1):130, 2014.
- 528 **36** Rui Henriques and Ana Paiva. Seven principles to mine flexible behavior from physiological  
529 signals for effective emotion recognition and description in affective interactions. In *PhyCS*,  
530 pages 75–82, 2014.
- 531 **37** Jan Holub, W.F. Smyth, and Shu Wang. Fast pattern-matching on indeterminate strings.  
532 *Journal of Discrete Algorithms*, 6(1):37 – 50, 2008. Selected papers from AWOCA 2005.
- 533 **38** Shuichi Kawashima and Minoru Kanehisa. Aaindex: amino acid index database. *Nucleic*  
534 *acids research*, 28(1):374–374, 2000.
- 535 **39** Jinil Kim, Peter Eades, Rudolf Fleischer, Seok-Hee Hong, Costas S Iliopoulos, Kunsoo  
536 Park, Simon J Puglisi, and Takeshi Tokuyama. Order-preserving matching. *Theoretical*  
537 *Computer Science*, 525:68–79, 2014.
- 538 **40** Donald E Knuth, James H Morris, Jr, and Vaughan R Pratt. Fast pattern matching in  
539 strings. *SIAM journal on computing*, 6(2):323–350, 1977.
- 540 **41** Marcin Kubica, Tomasz Kulczyński, Jakub Radoszewski, Wojciech Rytter, and Tomasz  
541 Waleń. A linear time algorithm for consecutive permutation pattern matching. *Information*  
542 *Processing Letters*, 113(12):430–433, 2013.
- 543 **42** Inbok Lee, Raphaël Clifford, and Sung-Ryul Kim. Algorithms on extended  $(\delta, \gamma)$ -matching.  
544 *Computational Science and Its Applications-ICCSA 2006*, pages 1137–1142, 2006.
- 545 **43** Inbok Lee, Juan Mendivelso, and Yoan J Pinzón.  $\delta\gamma$ -parameterized matching. In *Internat-*  
546 *ional Symposium on String Processing and Information Retrieval*, pages 236–248. Springer,  
547 2008.
- 548 **44** Ohad Lipsky and Ely Porat. Approximate matching in the  $\infty$  metric. *Information Pro-*  
549 *cessing Letters*, 105(4):138 – 140, 2008.
- 550 **45** Juan Mendivelso, Inbok Lee, and Yoan J Pinzón. Approximate function matching under  
551  $\delta$ -and  $\gamma$ -distances. In *SPIRE*, pages 348–359. Springer, 2012.

- 552 **46** S Muthukrishnan. New results and open problems related to non-standard stringology. In  
553 *Combinatorial Pattern Matching*, pages 298–317. Springer, 1995.
- 554 **47** Ely Porat and Klim Efremenko. Approximating general metric distances between a pattern  
555 and a text. In *ACM-SIAM symposium on Discrete algorithms*, pages 419–427. SIAM, 2008.