# Matching Patterns with Variables

Florin Manea[1] and Markus L. Schmid[2]

[1]Kiel University, Kiel, Germany, `flm@informatik.uni-kiel.de`
[2]Trier University, Trier, Germany, `MLSchmid@MLSchmid.de`

**Abstract**

A pattern $\alpha$ (i.e., a string of variables and terminals) matches a word $w$, if $w$ can be obtained by uniformly replacing the variables of $\alpha$ by terminal words. The respective matching problem, i.e., deciding whether or not a given pattern matches a given word, is generally NP-complete, but can be solved in polynomial-time for classes of patterns with restricted structure. In this paper we overview a series of recent results related to efficient matching for patterns with variables, as well as a series of extensions of this problem.

## 1 Introduction

A *pattern with variables*, called simply pattern in the context of this work, is a string that consists of *terminal symbols* (e.g., $\mathtt{a}, \mathtt{b}, \mathtt{c}$) and *variables* (e.g., $x_1, x_2, x_3$). The terminal symbols are treated as constants, while the variables are to be uniformly replaced by strings over the set of terminals (i.e., different occurrences of the same variable are replaced by the same string); thus, a pattern is mapped to a terminal word. For example, $x_1\mathtt{ab}x_1x_2\mathtt{c}x_2x_1$ can be mapped to $\mathtt{acabaccaaccaaac}$ and $\mathtt{babbacab}$ by the replacements $(x_1 \to \mathtt{ac}, x_2 \to \mathtt{caa})$ and $(x_1 \to \mathtt{b}, x_2 \to \mathtt{a})$, respectively.

Patterns with variables appear in various areas of theoretical computer science, such as language theory (pattern languages [2]), learning theory (inductive inference [2, 69, 75, 24], PAC-learning [53]), combinatorics on words (word equations [51, 67], unavoidable patterns [61]), pattern matching (generalised function matching [1, 70]), database theory (extended conjunctive regular path queries [5]), and we can also find them in practice in the form of extended regular expressions with backreferences [12, 37, 31], used in programming languages like Perl, Java, Python, etc.

Generally, in all these contexts, patterns with variables are used to model various combinatorial pattern matching questions. For instance, searching for a word $w$ in a text $t$ can be expressed as testing whether the pattern $xwy$ can be mapped to $t$ and testing whether a word $w$ contains a cube is equivalent to testing whether the pattern $xy^3z$ can be mapped to $w$, such that $y$ is not mapped to an empty word. Not only problems of testing whether a given word contains a regularity or a motif of a certain form can be expressed by patterns, but also problems asking whether a word can be factorised in a specifically restricted manner can be modelled in this way. For instance, asking whether $x_1^2x_2^2\ldots x_k^2$ can be mapped to $w$, such that none of the variables $x_i$ are mapped to an empty word, is equivalent to asking whether the word $w$ can be factorised into $k$ non-empty squares.

Unfortunately, deciding whether a given arbitrary pattern can be mapped to a given word, the *matching problem*, is NP-complete [2], whether we ask that the variables are mapped to non-empty words or not. This intractability result severely limits the practical application of patterns. Indeed, in many tasks related to applications of patterns, the matching problem is a necessary step, so the tasks become intractable as well. For instance, this is the case for the task of computing so-called descriptive patterns for finite sets of words (see [2, 35, 36] for more information on descriptive patterns): one cannot solve this problem without solving a series of (general) pattern matching tasks [27]. A more detailed analysis of the complexity of the hardness of the matching problem will be presented in Section 3.

On the other hand, some strong restrictions on the structure of patterns yield subclasses for which the matching problem is tractable (i.e., can be solved in polynomial time). This is clearly

the case of patterns where the number of different variables in the patterns is bounded by a constant, but more sophisticated and general such subclasses can be defined. We will discuss a series of results related to this topic in Sections 4.2, 5.1 and 5.2. In our analysis, the most general class of patterns which allow for a polynomial-time pattern matching problem is defined by establishing a deep connection between strings/patterns and graphs, and considering only patterns which correspond to graphs with bounded structural parameters. As such, the subclass of *patterns with bounded treewidth*. The question of finding classes of patterns which can be matched in polynomial time but do not have bounded treewidth seemed interesting to us. We show a natural construction of such patterns in Section 6.

We continue this survey with a result showing that considering some of the structural parameters, that lead to efficient pattern matching algorithms, as general structural parameters of strings, may lead to remarkable results in other apparently unrelated domains. We show in Section 7 how our results for strings can be used to obtain a state-of-the-art approximation algorithm for computing the *cutwidth of graphs*.

We conclude the survey with a series of extensions. We discuss the problem of *injective pattern matching* as well as the satisfiability problem for word equations with restricted form.

## 2 Basic Definitions

For detailed definitions regarding combinatorics on words we refer to [60].

We denote our *alphabet* by $\Sigma$, the *empty word* by $\varepsilon$, the set of all non-empty words over $\Sigma$ by $\Sigma^+$, the set of all words over $\Sigma$ by $\Sigma^*$, and the *length* of a word $w$ by $|w|$. $(\Sigma^*, \cdot, \varepsilon)$ is the free monoid over $\Sigma$ with *concatenation* as its binary operation, written $\cdot$. For $w \in \Sigma^*$ and every integers $i, j$ with $1 \leq i \leq j \leq |w|$, let $w[i..j] = w[i] \cdots w[j]$, where $w[k]$ represents the *letter on position $k$* and $1 \leq k \leq |w|$. A *period* of $w$ is any positive integer $p$ for which $w[i] = w[i+p]$, for all defined positions. Moreover, in this case, $w$ is said to be *$p$-periodic*. Its *minimal period* is denoted by $per(w)$ and represents the smallest period of $w$. For example, $w = \texttt{abacabacabacabacab}$ has periods 8 and 4; in particular, $per(w) = 4$. A word $w$ is called periodic if $per(w) \leq \frac{|w|}{2}$.

The *concatenation* of $k$ words $w_1, w_2, \ldots, w_k$ is written $\Pi_{i=1,k} w_i$. If $w = w_i$ for all integers $i$ with $1 \leq i \leq k$, this represents the *$k$th power* of $w$, denoted by $w^k$; here, $w$ is a *root* of $w^k$. We can further extend the notion of a power of a word by saying that $w = w[1..per(w)]^{\frac{|w|}{per(w)}}$. We say that $w$ is *primitive* if it cannot be expressed as a power of exponent $\ell$ of any root, where $\ell$ is an integer with $\ell > 1$. Conversely, if $w = v^\ell$ for some integer $\ell > 1$, then $w$ is also called a *repetition*. The infinite repetition $vvv \cdots$ of some word $v$ is denoted $v^\omega$.

For any word $w \in \Sigma^+$ with $w = xyz$, we say that $y$ is a *factor* of $w$. If $x$ is empty, then $y$ is also a *prefix* of $w$, while when $z$ is empty, then $y$ is also a *suffix*. Whenever we have a factor both as a prefix and as a suffix, the factor is said to be a *border* of the word. Furthermore, every word $u = yzx \in \Sigma^+$ is a *conjugate* of $w$. Note that, if $w$ is primitive, so is every conjugate of it. If $w = vu$, then $v^{-1}w = u$.

Let $X = \{x_1, x_2, x_3, \ldots\}$ and call every $x \in X$ a *variable*. For a finite alphabet $\Sigma$ of *terminals* with $\Sigma \cap X = \emptyset$, we define $\mathsf{Pat}_\Sigma = (X \cup \Sigma)^+$ and $\mathsf{Pat} = \bigcup_\Sigma \mathsf{Pat}_\Sigma$. Every $\alpha \in \mathsf{Pat}$ is a *pattern* and every $w \in \Sigma^*$ is a (*terminal*) *word*. Given a word or a pattern $v$, for the smallest sets $B \subseteq \Sigma$ and $Y \subseteq X$ with $v \in (B \cup Y)^*$, we denote $\mathsf{alph}(v) = B$ and $\mathsf{var}(v) = Y$. For any $x \in \Sigma \cup X$ and $\alpha \in \mathsf{Pat}_\Sigma$, $|\alpha|_x$ denotes the number of occurrences of $x$ in $\alpha$; for the sake of convenience, we set $|\alpha|_x = 0$ for every symbol $x$ not in $\Sigma \cup X$. For a pattern $\alpha$, we say that $w = \alpha[i..i + |w|]$ is a maximal terminal factor of $\alpha$ if $\alpha[i-1]$ and $\alpha[i + |w| + 1]$ are either not defined, or are variables.

A *substitution* (*for $\alpha$*) is a mapping $h : \mathsf{var}(\alpha) \to \Sigma^*$. For every $x \in \mathsf{var}(\alpha)$, we say that *$x$ is substituted by $h(x)$* and $h(\alpha)$ denotes the word obtained by substituting every occurrence of a variable $x$ in $\alpha$ by $h(x)$ and leaving the terminals unchanged. We say that the pattern $\alpha$ *matches* $w \in \Sigma^+$ if $h(\alpha) = w$ for some substitution $h : \mathsf{var}(\alpha) \to \Sigma^*$. Substitutions of the form $h : \mathsf{var}(\alpha) \to \Sigma^+$, i.e., the empty word is excluded from the range of the substitution, are also called *non-erasing*; in order to emphasize that the substitution by the empty word is allowed, we also use the term *erasing* substitution.

*Example 1.* Let $\beta = x_1 \texttt{a} x_2 \texttt{b} x_2 x_1 x_2$ be a pattern and let $u = \texttt{bacbabbbbacbb}$ and $v = \texttt{abaabbababab}$ be terminal words. The pattern $\beta$ matches both $u$ and $v$, witnessed by the substitutions $h$ with

$h(x_1) = \mathtt{bacb}$, $h(x_2) = \mathtt{b}$ and $g$ with $g(x_1) = g(x_2) = \mathtt{ab}$, respectively. Moreover, $\beta$ also matches the word $w = \mathtt{acbbcbcb}$ by the *erasing* substitution $h$ with $h(x_1) = \varepsilon$, $h(x_2) = \mathtt{cb}$; it can be easily verified that there is no non-erasing substitution that maps $\beta$ to $w$.

The *matching problem*, denoted by MATCH, is to decide for a given pattern $\alpha$ and word $w$, whether there exists a substitution $h$ with $h(\alpha) = w$. The variant where we are only concerned with non-erasing substitutions is called the *non-erasing case* of the matching problem; we also use the term *erasing-case* in order to emphasize that substitution by the empty word is allowed. Another special variant is the *terminal-free case* of the matching problem, where the input patterns are terminal-free, i.e., they do not contain any occurrences of terminal symbol. We shall briefly discuss some particularities of these different special cases of the matching problem in Section 3. Note that in the sections on efficient algorithms, namely Sections 5.1, 5.2, and 6, we only consider the non-erasing case (with terminal symbols) of the matching problem. The presented results can easily be generalised to the general setting, but we prefer the respective framework for the ease of the presentation.

For any $P \subseteq \mathsf{Pat}$, the *matching problem for $P$* (or MATCH *for $P$*, for short) is the matching problem, where the input patterns are from $P$. In the sections of this paper we will introduce and discuss several interesting families of patterns.

As we discuss efficient algorithms, it is important to describe the computational model we use in this work. This is the standard unit-cost RAM with logarithmic word size. Also, all logarithms appearing in our time complexity evaluations are in base 2. For the sake of generality, we assume that whenever we are given as input a word $w \in \Sigma^*$ of length $n$, the symbols of $w$ are in fact integers from $\{1, 2, \ldots, n\}$ (i.e., $\Sigma = \mathsf{alph}(w) \subseteq \{1, 2, \ldots, n\}$), and $w$ is seen as a sequence of integers. This is a common assumption in the area of algorithmics on words (see, e.g., the discussion in [52]). Clearly, our algorithmic results hold canonically for constant alphabets, as well.

## 3  The Hardness of the Matching Problem

First, we recall that there are several different variants of the matching problem: the most general case (substitution by the empty word and occurrences of terminals in the patterns are possible), the non-erasing case (with terminal symbols), the terminal-free (erasing) case, and finally the terminal free non-erasing case. As we shall see, these differences do not matter too much if we are only concerned with the matching problem of patterns. However, in other contexts of patterns with variables (e. g., other decision problems, learning theory), these differences are most crucial and we therefore briefly provide some background.

For the class of the so-called *pattern languages*, i.e., the sets of all words that match a pattern, the difference between the erasing and the non-erasing case is important, since these classes of formal languages differ quite substantially with respect to basic decision problems. For example, in the non-erasing case, two patterns describe the same language if and only if the patterns are identical (up to a renaming of variables), while it is open whether the equivalence problem is even decidable in the erasing-case (see, e. g., Section 6 in [68], or [74]). Moreover, the inclusion problem, which is undecidable for both the erasing and the non-erasing case (see [50, 34]), can be decided for terminal-free patterns in the erasing case, while for terminal-free non-erasing patterns the decidability status is open (intuitively speaking, this has to do with the fact that avoidability questions of the form "does pattern $\beta$ necessarily occur in long enough words over a $k$-letter alphabet?" can be expressed as inclusion for two languages given by terminal-free non-erasing patterns). Finally, also whether patterns (or descriptive patterns) can be inferred from positive data strongly depends on whether the erasing or non-erasing case is considered, or whether or not terminal symbols in the patterns are allowed (see [73, 75, 36, 35]).

For the matching problem (note that this corresponds to the membership problem for pattern languages), whether we consider erasing or non-erasing substitution, or whether or not we disallow terminal symbols in the patterns, has little impact on its computational hardness. In fact, that the matching problem for patterns with variables is NP-complete has been independently discovered in different communities and for slightly different problem variants (see, e. g., the introductions of [28, 29] for some remarks on the history of the investigation of the matching problem).

If we consider the most general case, i.e., erasing substitutions and terminals in the patterns, then a hardness-reduction is rather simple. For example, the Boolean formula

$$((v_1, v_2, v_3), (v_2, v_4, v_5), (v_3, v_1, v_3), (v_4, v_1, v_2))$$

in 3-CNF (without negated variables) is 1-in-3 satisfiable (i.e., satisfiable with exactly one literal per clause set to *true*) if and only if the following word $w$ is matched by the pattern $\alpha$:

$$
\begin{array}{ccccccccccccc}
w & = & \texttt{a} & \texttt{b} & \texttt{a} & \texttt{b} & \texttt{a} & \texttt{b} & \texttt{a} \\
\alpha & = & x_1 x_2 x_3 & \texttt{b} & x_2 x_4 x_5 & \texttt{b} & x_3 x_1 x_3 & \texttt{b} & x_4 x_1 x_2
\end{array}
$$

We can further observe that this simple reduction also shows that the matching problem is hard even for binary terminal alphabets and under the restriction that variables are substituted by single symbols (or the empty word) only. This directly raises the questions under which restrictions the matching problem remains hard. For example, a problem instance has a large number of natural parameters (length of the pattern, length of the word, number of variables, number of occurrences per variable, alphabet size, length of words substituted for variables) and in addition to that, it comes in four natural variants resulting from whether we consider the erasing or non-erasing case, and whether or not we allow terminals in the pattern. In the above reduction, the number of variables, the number of occurrences per each variable and the word length are unbounded.

All these numerous restricted problem variants have been thoroughly investigated in [28] and it turns out that the matching problem remains NP-hard under rather strong restrictions. We cite the following result as an example and refer to [28] for further details.

**Theorem 2** ([28])**.** *The erasing case of the matching problem for patterns with variables is* NP-*complete, even if $\Sigma = \{\texttt{a}, \texttt{b}\}$, every variable has at most 2 occurrences and every variable can only be substituted by a single symbol or the empty word.*

This result also holds as stated for terminal-free patterns. In the non-erasing case, however, it holds when the bound on the substitution words is 3 instead of 1, and in the non-erasing and terminal-free case the result holds when additionally the bounds on the occurrences per variable and alphabet size are 3 and 4, respectively.

The only polynomial-time solvable cases of the matching problem obtained by restricting the numerical parameters mentioned above are trivial ones. More precisely, the matching problem can be easily solved for unary alphabets (in this case, we only have to solve an equation in the integers and with integer coefficients, which are given in unary encoding), or if every variable has only one occurrence (the patterns are then *regular*, see Section 4), or if the number of variables or the length of the input word is bounded by a constant (the former is obvious, while the latter, in the erasing case, requires a slightly more careful argument [43]).

In particular, this also points out that Theorem 2 describes some kind of dichotomy, i.e., if we would further restrict the alphabet size, or the maximum number of occurrences per variable to 1, then we would obtain a polynomial-time solvable variant (even if all other parameters are unrestricted); similarly, if we allow variables to be substituted by single symbols only, but not the empty word, then the matching problem becomes efficiently solvable as well (regardless of the alphabet size).

Generally, by brute-force algorithms, the matching problem can be solved in time $|\alpha|^{O(|w|)}$ or $|w|^{O(|\alpha|)}$, making it polynomial-time solvable provided that there is a constant upper bound on $|w|$ or $|\alpha|$ (in fact, a bound on $|\mathsf{var}(\alpha)|$ is sufficient). However, this constant upper bound occurs in the exponent, which means that even for rather low such bounds, say 7, the corresponding polynomial-time algorithms are most likely impractical for larger problem instances. This leads to the question whether exponential-time algorithms are possible whose running-times are such that the exponential part *exclusively* depends on, say $|\mathsf{var}(\alpha)|$, but not on $|w|$, i.e., running-times of the form $f(|\mathsf{var}(\alpha)|) \times g(|\alpha|, |w|)$, where $g$ is a polynomial and $f$ is some *computable* function (exponential, or even double-exponential etc.). Such a running-time is polynomial for upper bounded $|\mathsf{var}(\alpha)|$, but the degree of the polynomial is always the same independent from the actual upper bound. If a problem has an algorithm with such a running-time, then it is called *fixed-parameter tractable* (with respect to the bounded parameter); see the textbooks [23, 30] for more information on parameterised complexity. Whether the matching problem for patterns
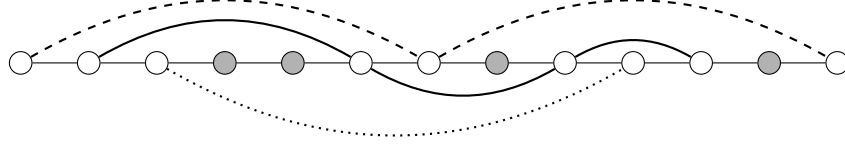
Figure 1: The standard graph representation $G_\alpha^{\mathsf{pat}}$ for $\alpha = x_1 x_2 x_3 \mathtt{bb} x_2 x_1 \mathtt{a} x_2 x_3 x_2 \mathtt{c} x_1$; the dashed, straight and dotted equality edges correspond to occurrences of $x_1$, $x_2$ and $x_3$, respectively; the grey vertices correspond to occurrences of terminal symbols.

with variables allows fixed-parameter tractability for some parameters has been thoroughly investigated in [29]. Although there are some more or less trivial cases of fixed-parameter tractability, the main insight provided by [29] is of a negative nature and can be summarised in the following way.

**Theorem 3** ([29]). *All variants of the matching problem parameterised by $|\alpha|$ are* W[1]-*hard. The erasing case of the matching problem parameterised by $|w|$ is* W[1]-*hard.*[1]

Note that since $|\alpha|$ and $|w|$ are rather general parameters, this result covers other parameters as well, e.g., $|\Sigma|$ or $|\mathsf{var}(\alpha)|$. In the non-erasing case, $|w|$ is an upper bound for $|\mathsf{var}(\alpha)|$; thus, treating $|w|$ as a parameter means that $|\mathsf{var}(\alpha)|$ is also a parameter and therefore the matching problem is fixed-parameter tractable by the obvious brute-force algorithm. We refer to [29] for further such simple fixed-parameter tractable case.

Consequently, even strong restrictions of the obvious numerical parameters of instances of the matching problem, i.e., number of variables, alphabet size, occurrences per variable etc., does not yield interesting efficiently matchable subclasses of patterns with variables. However, as discussed in the next section, looking deeper into the structure of patterns will help.

## 4    Structural Restrictions for Patterns

From an intuitive point of view it is clear that not only the mere length of a pattern or the number of its variables should have an impact on the matching complexity, but also the actual order of the variables. For example, it has been observed rather early in [79] that if the variable occurrences in the patterns are sorted, e.g., as in $x_1 \mathtt{a} x_1 x_2 x_2 \mathtt{ab} x_2 x_2 \mathtt{a} x_3 x_4 \mathtt{c} x_4$, then they can be matched efficiently "from-left-to-right" (more precisely, it is observed in [79] that matching such patterns can be done in logarithmic space).

A systematic investigation of such structural restrictions has been done in the last decade and numerous efficiently matchable subclasses of patterns have been found. In the following, we first present a unifying approach based on graph morphisms and the concept of treewidth. Then, we define and summarise several structural parameters for patterns and respective subclasses of patterns.

### 4.1    Pattern Matching by Graph Morphisms

The following general framework for matching patterns with variables has been developed in [76]. For a pattern $\alpha \in (X \cup \Sigma^*)$, the *standard graph representation* of $\alpha$ is the undirected graph $\mathcal{G}_\alpha^{\mathsf{pat}} = (V_\alpha, E_\alpha)$, where $V_\alpha = \{1, 2, \ldots, |\alpha|\}$ and $E_\alpha = E_\alpha^{\mathsf{equ}} \cup E_\alpha^{\mathsf{nei}}$ with $E_\alpha^{\mathsf{equ}} = \{\{i, i+1\} \mid 1 \le i \le |\alpha|-1\}$ being the set of *neighbour edges* and $E_\alpha^{\mathsf{equ}} = \{\{i, j\} \mid \alpha[i..j] = x\beta x, x \in X, |\beta|_x = 0\}$ being the set of *equality edges* (see Figure 1 for an illustration).

In a similar way, we can also encode words $w \in \Sigma^*$ as graph structures $\mathcal{G}_w^{\mathsf{wo}}$, where every factor $w[i..j]$ of $w$ is represented by a vertex $(i, j)$, equality edges are drawn between $(i, j)$ and $(i', j')$ if $w[i..j] = w[i'..j']$, and neighbour edges if $j + 1 = i'$. It has been shown in [76] that $\alpha$ matches $w$ if and only if there is a graph morphism from $\mathcal{G}_\alpha^{\mathsf{pat}}$ to $\mathcal{G}_w^{\mathsf{wo}}$. Moreover, the concept of the

---

[1]Problems that are hard for the parameterised complexity class W[1] are strongly believed to be not fixed-parameter tractable.
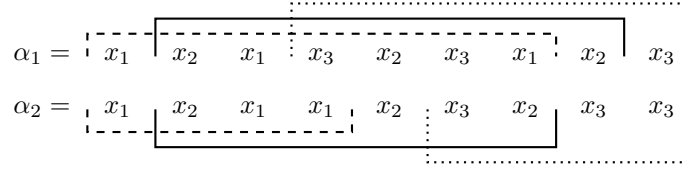
$$\alpha_1 = \quad x_1 \quad x_2 \quad x_1 \quad x_3 \quad x_2 \quad x_3 \quad x_1 \quad x_2 \quad x_3$$

$$\alpha_2 = \quad x_1 \quad x_2 \quad x_1 \quad x_1 \quad x_2 \quad x_3 \quad x_2 \quad x_3 \quad x_3$$

Figure 2: Two pattern $\alpha_1$ and $\alpha_2$ with $\mathsf{scd}(\alpha_1) = 3$ and $\mathsf{scd}(\alpha_2) = 2$. The scopes of variable $x_1$ (dashed line), $x_2$ (straight line) and $x_3$ (dotted line) are highlighted.

treewidth for graphs now also applies to patterns (i. e., the treewidth of a pattern is the treewidth of its standard graph representation), which is of relevance since the graph morphism problem can be solved in polynomial-time provided that the source graphs have bounded treewidth.[2] Consequently, we can conclude the following algorithmic meta-theorem.

**Theorem 4** ([76]). *If a class $P$ of patterns has bounded treewidth, then the matching problem for $P$ can be solved in polynomial-time.*

Due to the generality of the statement of Theorem 4, the polynomial-time matching algorithm that it implies is of little practical value, even for rather simple classes of patterns. On the other hand, its theoretical relevance is demonstrated by the fact that it covers almost all known classes of patterns with a polynomial-time matching problem.[3] After an additional remark regarding [76], we shall briefly define and compare those efficiently matchable classes of patterns in the next subsection.

*Remark* 5. Technically, the matching problem reduces to the morphism problem for (simple) relational structures instead of undirected graphs. However, since we are here only interested in the treewidth of these structures, we can as well only talk about the underlying undirected graphs.

Moreover, the actual meta-theorem of [76] is stronger in the sense that there the treewidth of patterns is not defined with respect to the standard graph representation, but with respect to a slightly more general graph representations (i. e., we allow any way of drawing the equality edges as long as all vertices corresponding to the same variable form a connected component).

## 4.2 Efficiently Matchable Classes of Patterns

The most obvious way to restrict patterns is to limit their number of (repeated) variables or the number of occurrences per variable. In this regard, let $\mathsf{var}_k$ and $\mathsf{repv}_k$ be the class of patterns with at most $k$ variables and with at most $k$ repeated variables, respectively. Due to Theorem 2, we already know that bounding the number of occurrences per variable does not in general yield polynomial-time matchable classes. The only exception are patterns with at most one occurrence per variable, which are called *regular* patterns and are denoted by $\mathsf{reg}$, e. g., $x_1\mathtt{a}x_2\mathtt{bac}x_3\mathtt{a}$ is a regular pattern. Regular patterns have been first considered in [79] and their name is motivated by the fact that the corresponding pattern languages are regular languages.

Next, we define the so-called scope coincidence degree (see [76]). For every $y \in \mathsf{var}(\alpha)$, the *scope of $y$ in $\alpha$* is defined by $\mathrm{sc}_\alpha(y) = \{i, i+1, \ldots, j\}$, where $i$ is the leftmost and $j$ the rightmost occurrence of $y$ in $\alpha$. The scopes of some variables $y_1, y_2, \ldots, y_k \in \mathsf{var}(\alpha)$ *coincide in $\alpha$* if $\bigcap_{1 \leq i \leq k} \mathrm{sc}_\alpha(y_i) \neq \emptyset$. By $\mathsf{scd}(\alpha)$, we denote the *scope coincidence degree* of $\alpha$, which is the maximum number of variables in $\alpha$ such that their scopes coincide, and by $\mathsf{scd}_k$, we denote the class of patterns with scope coincidence degree of at most $k$. See Figure 2 for an example of the scope coincidence degree. An important special class is $\mathsf{scd}_1$, which has been first introduced in [79] as the class of *non-cross* patterns (denoted by $\mathsf{nc}$). Intuitively speaking, the variables in non-cross patterns are sorted, e. g., $x_1\mathtt{a}x_1x_2\mathtt{ba}x_2\mathtt{c}x_3\mathtt{a}x_3x_3$.

Next, we define the locality number, which is a general string-parameter, and which has been first introduced in [15]. A word is *$k$-local* if there exists an order of its symbols such that, if we

---

[2]See [23, 30] for a formal definition of the treewidth.
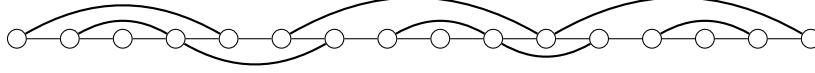[3]See Section 6 for the respective exceptions.

Figure 3: The standard graph representation $G_\alpha^{\mathsf{pat}}$ for $\alpha = x_1x_3x_4x_3x_1x_2x_3x_5\mathsf{b}x_5x_2x_5x_6\mathsf{a}x_6x_2$. By definition, $\alpha$ is mildly entwined. Furthermore, since no vertex is completely "surrounded" by edges, the shown embedding is outer-planar.

*mark* the symbols in the respective order (which is called a *marking sequence*), at each stage there are at most $k$ contiguous blocks of marked symbols in the word. This $k$ is called the *marking number* of that marking sequence. The *locality number* of a word is the smallest $k$ for which that word is $k$-local, or, in other words, the minimum marking number over all marking sequences. For example, the marking sequence $\sigma = (\mathsf{a}, \mathsf{g}, \mathsf{c})$ marks $w = \mathsf{agagcac}$ as follows (marked blocks are illustrated by overlines): $\mathsf{agagcac}$, $\overline{\mathsf{a}}\mathsf{g}\overline{\mathsf{a}}\mathsf{gc}\overline{\mathsf{a}}\mathsf{c}$, $\overline{\mathsf{agag}}\mathsf{c}\overline{\mathsf{a}}\mathsf{c}$, $\overline{\mathsf{agagcac}}$; thus, the marking number of $\sigma$ is 3. In fact, all marking sequences for $w$ have a marking number of 3, except $(\mathsf{g}, \mathsf{a}, \mathsf{c})$, for which it is 2: $\mathsf{a}\overline{\mathsf{g}}\mathsf{a}\overline{\mathsf{g}}\mathsf{cac}$, $\overline{\mathsf{agag}}\mathsf{cac}$, $\overline{\mathsf{agagcac}}$. Thus, the locality number of $w$, denoted by $\mathsf{loc}(w)$, is 2. When we measure the locality number for patterns, we simply ignore all terminal symbols, e.g., $\mathsf{loc}(\mathsf{ab}x_1x_2\mathsf{a}x_1x_2\mathsf{c}x_3x_1\mathsf{a}x_3) = \mathsf{loc}(x_1x_2x_1x_2x_3x_1x_3) = 2$. The class of patterns with locality number at most $k$ is denoted by $\mathsf{loc}_k$.

The next classes have been first considered in [76] and are based on possible nesting structures of variables. For a pattern $\alpha$, we call two variables $x, y \in \mathsf{var}(\alpha)$ *entwined* if $\alpha$ contains $xyxy$ or $yxyx$ as a subsequence. A pattern $\alpha$ is *nested*, if no two variables in $\alpha$ are entwined; the class of nested patterns is denoted by $\mathsf{nest}$. A proper subclass of $\mathsf{nest}$, considered in [15], are the so-called *strongly* nested patterns (denoted by $\mathsf{snest}$), which are inductively defined as follows: any pattern $\alpha \in \mathsf{var}_1$ is strongly nested; if $\alpha_1$ and $\alpha_2$ are strongly nested and variable-disjoint patterns, $x$ is a variable not in $\mathsf{var}(\alpha_1) \cup \mathsf{var}(\alpha_2)$ and $\beta_1, \beta_2 \in (\{x\} \cup \Sigma)^*$, then $\alpha_1\alpha_2$ and $\beta_1\alpha_1\beta_2$ are strongly nested patterns. For example, the pattern $\alpha = x_1x_2\mathsf{a}x_2x_1\mathsf{b}x_3x_4\mathsf{a}x_3$ is strongly nested, whereas $\alpha x_1$ is nested, but not strongly nested anymore.

If, for every $x, y \in \mathsf{var}(\alpha)$, $\alpha = \beta x\gamma_1 y\gamma_2 x\gamma_3 y\delta$ implies $\gamma_2 = \varepsilon$, then $\alpha$ is called *closely entwined*, and a pattern $\alpha$ is *mildly entwined* if it is closely entwined and, for every $x \in \mathsf{var}(\alpha)$, if $\alpha = \beta x\gamma x\delta$ with $|\gamma|_x = 0$, then $\gamma$ is nested. We denote the class of mildly entwined patterns by $\mathsf{ment}$. The main motivation for the somewhat peculiar class of mildly entwined patterns is that mildly entwined patterns are exactly those patterns that have a standard graph representation that is outer-planar (see [76]).[4] It is known that outer-planar graphs have a rather low treewidth of at most 2. Since the concept of outer-planarity generalises to $k$-outer-planarity and $k$-outer-planar graphs have a treewidth of at most $3k - 1$, we can also define the classes $\mathsf{outp}_k$ of $k$-outer-planar patterns (i.e., their standard graph representation is $k$-outer-planar). In this regard note that $\mathsf{outp}_1 = \mathsf{ment}$. See Figure 3 for an example of a mildly-entwined pattern.

It can be easily verified that all of the pattern classes defined above have bounded treewidth; thus, by application of Theorem 4, they can be matched efficiently. For some of them this upper bound on the treewidth is rather low (e.g., $\mathsf{reg}$, $\mathsf{nc}$, $\mathsf{ment}$), while for those classes obtained by bounding a structural parameter, e.g., $\mathsf{repv}_k$, $\mathsf{scd}_k$, $\mathsf{loc}_k$, the bound on the treewidth also grows with this parameter. Figure 4 shows how these pattern classes relate to each other and how they form infinite hierarchies within the class of all patterns (denoted by $\mathsf{Pat}$).

In a sense, Figure 4 is a "tractability map" for the matching problem of patterns with variables. For the classes that have low treewidth, we can expect matching algorithms that are rather efficient. On the other hand, these classes are quite restricted (compared to the full class of patterns) and are most likely only applicable for very special pattern matching tasks. An obvious approach to matching general patterns would be to first perform a preprocessing that identifies a "low class" of the tractability map that contains the input pattern and then uses the most efficient algorithm for matching it. In this regard, it is even an asset that most of the different efficiently matchable classes and hierarchies of classes are incomparable: it is possible that an input pattern has a very large locality number of 100, but can nevertheless be matched efficiently, because its standard graph representation is 2-outerplanar; on the other hand, a pattern could have a large scope coincidence degree and a large number of variables,

---

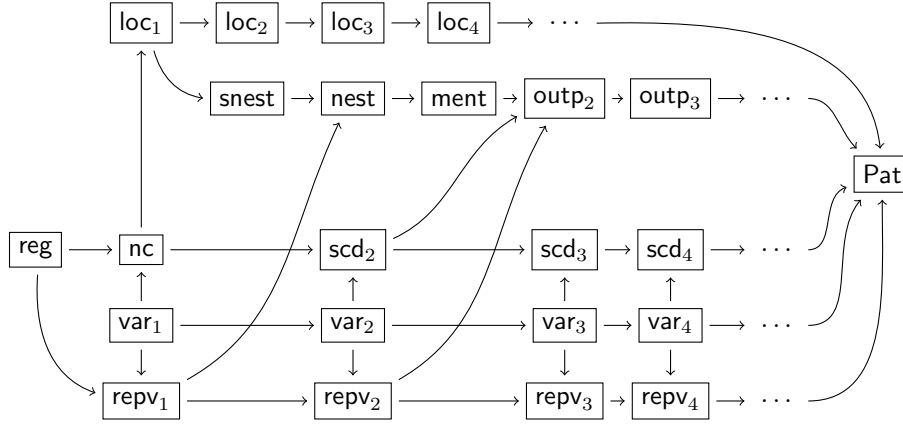[4] A graph is outer-planar if it has a planar embedding with all vertices lying on the outer face.

Figure 4: An overview of efficiently matchable classes of patterns. By $A \to B$, we denote $A \subset B$; pairs without arrow are incomparable. Note that $\mathsf{nc} = \mathsf{scd}_1$ and $\mathsf{ment} = \mathsf{outp}_1$.

but at the same time a very low locality number. It might even be a worthwhile research task to experimentally analyse a large corpus of (random) patterns with respect to the classes of the tractability map in which they are contained.

*Remark* 6. Bounding the structural parameters defined above yields polynomial-time matchable classes of patterns; thus, the question arises whether the matching problem is also fixed-parameter tractable with respect to those parameters. However, Theorem 3 already states that this is most likely not the case for parameter $|\mathsf{var}(\alpha)|$, and since $|\mathsf{var}(\alpha)|$ is an upper bound for the number of repeated variables, the scope coincidence degree, the outer-planarity and the locality number of $\alpha$, it is also highly unlikely that we can achieve fixed-parameter tractability with respect to those parameters.

### 4.3 Computing Structural Parameters for Patterns

Since the structural restrictions of patterns surveyed above are all meant to be exploited algorithmically, the task of checking them (or computing the respective parameters) is an important issue. In this regard, note that in general computing the treewidth of a graph is an NP-hard problem and it is also not known whether it can be computed efficiently for standard graph representations of patterns. This also emphasises the importance of *easily computable* parameters that are bounding the treewidth of a pattern and also points out why the value of Theorem 4 is of a theoretical nature that provides guidance in finding such restrictions with higher practical relevance. Restrictions like the regularity, the non-cross condition, number of (repeated) variables and the different nesting properties can be easily checked for. Moreover, also the scopes of a pattern and therefore its scope coincidence degree can be efficiently computed, and the smallest $k$ for which a graph is $k$-outerplanar can also be computed in polynomial time (for more details see [77]). On the other hand, computing the locality number seems more difficult and it was left open in [15] whether or not is is hard to compute. This gap was closed in [13] where it was shown that computing the locality number is NP-hard, but fixed-parameter tractable (if the locality number or $|\Sigma|$ is considered a parameter); in addition, approximation of the locality number has also been investigated in [13] (note that these result will be discussed in more detail in Section 7).

## 5   Faster Pattern Matching

In this section we will overview some efficient matching algorithms developed for various classes of patterns, some defined already in the previous sections, and some defined via some other natural structural restrictions. Most of the result of this paper were shown in [26, 15, 16].

## 5.1 Patterns with Low Scope Coincidence Degree

We start with several definitions. The *one-variable blocks* in a pattern are maximal contiguous blocks of occurrences of the same variable. A pattern $\alpha$ with $m$ one-variable blocks can be written as $\alpha = w_0 \Pi_{i=1,m}(z_i^{k_i} w_i)$ with $z_i \in \mathsf{var}(\alpha)$ for $i \in \{1, 2, \ldots, m\}$ and $z_i \neq z_{i+1}$, whenever $w_i = \varepsilon$ for $i \in \{1, 2, \ldots, m-1\}$. The number of one-variable blocks is a natural complexity measure that we will consider.

*Example* 7. The pattern $\alpha = x_1 x_2 x_2 \mathsf{a} x_2 x_2 x_2 x_3 \mathsf{a} x_3 x_2 x_2 x_3 x_3$ has the following 7 one-variable blocks: $x_1, x_2 x_2, x_2 x_2 x_2, x_3, x_3, x_2 x_2, x_3 x_3$.

As discussed in the previous sections, prominent subclasses of patterns for which MATCH can be solved in polynomial time are the classes of patterns with a bounded number of (repeated) variables ($\mathsf{var}_k$ and $\mathsf{repv}_k$), of regular patterns ($\mathsf{reg}$), of non-cross patterns ($\mathsf{nc}$), and of patterns with a bounded scope coincidence degree ($\mathsf{scd}_k$). However, the known respective algorithms are rather poor considering their running times. For example, for $\mathsf{var}_k$, the matching problem can be solved in $O(\frac{mn^{k-1}}{(k-1)!})$, where $m$ and $n$ are the lengths of the pattern and the word (see [45]). For patterns with a scope coincidence degree of at most $k$, an $O(mn^{2(k+3)}(k+2)^2)$ time algorithm can be derived using the general matching technique described by Theorem 4, where $m$ and $n$ are the lengths of the pattern and the word, respectively, and the proof that the matching problem for non-cross patterns is in P (see [79]) leads to an $O(n^4)$-time algorithm. Hence, for all these classes, we consider the following refinement of the problem of showing that the matching problem for a class of patterns is in P.

*Problem* 8. Let $K$ be a class of patterns for which the matching problem can be solved in polynomial time. Find an efficient algorithm that solves the matching problem for $K$.

The main class of patters considered in the following is that of patterns with bounded scope coincidence degree, and its subclasses.

If the scope coincidence degree is bounded by 1, i.e., non-cross patterns, we can decide whether a pattern $\alpha$ having $m$ one-variable blocks matches a word $w$ of length $n$ in $O(mn \log n)$ time. This result can be achieved via a general dynamic programming approach, which tries to match prefixes of the pattern $\alpha$ to the prefixes of the word $w$. This general approach is rather standard but the big gain is that it can be implemented efficiently by a a detailed combinatorial analysis of the possible matches between the one-variable blocks occurring in $\alpha$ and factors of $w$. For instance, if the shortest factor of $\alpha$ containing all occurrences of a variable $x$ starts with a one-variable block containing at least two occurrences the variable $x$, we can efficiently find the matches of this factor by exploiting a major result from [14], which states that the primitively rooted squares contained in a word of length $n$ can be listed optimally in $O(n \log n)$. As each match for a factor starting with two occurrences of a variable starts with a primitively rooted square, the respective matches can be found efficiently. The result regarding primitively rooted squares can be extended to show that, given a word $w$ of length $n$ and a word $v$ with length shorter than $n$, the word $w$ contains $O(n \log n)$ factors of the form $uvu$ with $uv$ primitive, and all these factors can be found optimally in $O(n \log n)$ time. This allows us to find efficiently the matches for one-variable that the shortest factor of $\alpha$ which contains all occurrences of $x$ and starts with $xvx$, for all choices of a variable $x$ such that $v$ is a non-empty terminal string.

**Theorem 9** ([26]). *The matching problem for $\mathsf{nc}$ is solvable in $O(mn \log n)$ time, where $w$ is the input word of length $n$ and $m$ is the number of one-variable blocks occurring in the pattern.*

Two particular subclasses of non-cross patterns are of interest: the regular patterns $\mathsf{reg}$ and the one-variable patterns $\mathsf{var}_1$ (see also Figure 4). It is not hard to show that regular patterns can be matched in linear time $O(|\alpha| + |w|)$, by iteratively using the Knuth-Morris-Pratt algorithm to identify greedily the terminal factors occurring in the pattern, in their orders of occurrences. All factors of a word $w$ that match a given regular pattern $\alpha$ can be detected in linear time too.

More interesting is the case of one-variable patterns. The simplest example of one-variable patterns are the repetitions, i.e., patterns of the form $x^k$. Checking whether a word is a match for a pattern $x^k$ can be done in linear time. Moreover, a compact representation of all periodic factors of a word $w$ can be also obtained in linear time by identifying the (at most $|w|$) so-called runs inside $w$ [4]. With this, a compact representation of occurrences of $x^k$ in $w$ can also be

obtained in linear time. More complex one-variable patterns are the pseudo-repetitions (see [42, 39, 41] and the references therein). These are patterns from $\{x, x^R\}^*$, where $x^R$ is a variable that is always substituted by the reverse image of the string substituting $x$. Checking whether a string matches a given pseudo-repetition can be done in linear time [42]. The following general result can be shown for one-variable patterns, see [57]. Given a pattern $\alpha = v_1 x v_2 x \cdots v_{r-1} x v_r$ such that $x$ is a variable and $v_1, v_2, \ldots, v_r$ are terminal strings, a compact representation of all $P$ instances of $\alpha$ in the input string $w$ of length $n$ can be computed in $\mathrm{O}(rn)$ time, so that one can report those occurrences in $\mathrm{O}(P)$ time. The same result holds also for the case when some of the occurrences of $x$ in such a pattern are replaced by $x^R$. It is worth noting that using this algorithm to find the factors of a given word that match the shortest factor of $\alpha$ containing all occurrences of a variable $x$ inside a non-cross pattern in our approach for matching nc does not lead to a faster matching algorithm in that case.

When considering general patterns with bounded scope coincidence degree, one can show, using a similar dynamic programming approach as in the case of non-cross patterns, that the matching problem for $\mathsf{scd}_k$ is solvable in $\mathrm{O}(\frac{mn^{2k}}{((k-1)!)^2})$ time, where $n$ is the length of the input word and $m$ is, again, the number of one-variable blocks occurring in the pattern. One should note that in this case it seems hard to use the combinatorial insights used for non-cross patterns (thus, the $\log n$ factor is replaced by an $n$ factor in the evaluation of the time complexity), but, still, this algorithm is significantly faster than the previously known solution.

**Theorem 10** ([26]). *The matching problem for $\mathsf{scd}_k$ is solvable in $\mathrm{O}\left(\frac{mn^{2k}}{((k-1)!)^2}\right)$ time, where $w$ is the input word of length $n$ and $m$ is the number of one-variable blocks occurring in the pattern.*

Next we consider the classes $\mathsf{repv}_k$. For the basic case of $k = 1$, the matching problem can be solved in $\mathrm{O}(n^2)$ time, where $n$ is the length of the input word. The idea of this algorithm is to guess the length $\ell$ of the repeated variable $x$, and then to partition the suffix array of the input word into clusters, such that all suffixes in a cluster start with the same factor of length $\ell$. Essentially, in a match between the pattern and the word, where $x$ is mapped to a factor of length $\ell$, the positions where the factors matching $x$ occur in the input word belong to the same cluster. Using this idea, the desired complexity is then reached, again via dynamic programming.

**Theorem 11** ([26]). *The matching problem for $\mathsf{repv}_k$ is solvable in quadratic time.*

Further, one can use this result to show that the matching problem for the general class of patterns $\mathsf{repv}_k$ is solvable in $\mathrm{O}(\frac{n^{2k}}{((k-1)!)^2})$ time. This algorithm is better than the one that could have been obtained by using the fact that patterns with at most $k$ repeated variables have the scope coincidence degree bounded by $k + 1$, and then directly applying our previous algorithm solving the matching problem for $\mathsf{scd}_{k+1}$.

**Theorem 12** ([26]). *The matching problem for $\mathsf{repv}_k$ is solvable in $\mathrm{O}\left(\frac{n^{2k}}{((k-1)!)^2}\right)$ time, where $n$ is the length of the input word.*

Note that the classes of non-cross patterns and of patterns with a bounded scope coincidence degree or with a bounded number of repeated variables are of special interest, since for them we can compute so-called descriptive patterns (see [2, 79]) in polynomial time. A pattern $\alpha$ is *descriptive* (with respect to, say, non-cross patterns) for a finite set $S$ of words if it can generate all words in $S$ and there exists no other non-cross pattern that describes the elements of $S$ in a better way. Computing a descriptive pattern, which is NP-complete in general, means to infer a pattern common to a finite set of words, with applications for inductive inference of pattern languages (see [69]). For example, our algorithm for computing non-cross patterns can be used in order to obtain an algorithm that computes a descriptive non-cross pattern in time $\mathrm{O}(\sum_{w \in S}(m^2|w|\log|w|))$, where $m$ is the length of a shortest word of $S$ (see [27] for details).

The algorithms, except the ones for the basic cases of regular and non-cross patterns and patterns with only one repeated variable, still have an exponential dependency on the number of repeated variables or the scope coincidence degree. Therefore, only for very low constant bounds on these parameters can these algorithms be considered efficient. Naturally, finding a polynomial time algorithm for which the degree of the polynomial does not depend on the

number of repeated variables or on the scope coincidence degree would be desirable. However, by Remark 6 such algorithms are very unlikely.

Finally we recall a result regarding *gapped repeats and palindromes*. A gapped repeat (palindrome) is an instance of a terminal-free pattern $xyx$ (respectively, $xyx^R$). For $\alpha \geq 1$, an $\alpha$-gapped repeat in a word $w$ is a factor $uvu$ of $w$ such that $|uv| \leq \alpha|u|$; the two factors $u$ in such a repeat are called arms, while the factor $v$ is called gap. Such a repeat is called maximal if its arms cannot be extended simultaneously with the same symbol to the right or, respectively, to the left. In a sense, $\alpha$-gapped repeats are instances of the pattern $xyx$ where length constraints are imposed on the strings that substitute $x$ and $y$. In [40] it was shown that the number of maximal $\alpha$-gapped repeats that may occur in a word is upper bounded by $18\alpha n$. Using this, an algorithm finding all the maximal $\alpha$-gapped repeats of a word in $O(\alpha n)$ was defined; this result is optimal, in the worst case, as there are words that have $\Theta(\alpha n)$ maximal $\alpha$-gapped repeats. Comparable results were developed for the case of $\alpha$-gapped palindromes, i.e., factors $uvu^R$ with $|uv| \leq \alpha|u|$. On the one hand, these results were relevant as they provided optimal algorithms for the identification of $\alpha$-gapped repeats and palindromes, and closed an open problem from [55, 56] (see also [40] and the references therein for more on gapped repeats and palindromes). On the other hand, they point towards the study of MATCH for patterns with (linear) length constraints on the images of the variables.

## 5.2    Patterns with Low Locality Number

Intuitively, the notion of $k$-locality (already introduced in Section 4.2) involves marking the variables in the pattern in some arbitrary order until all the variables are marked. The pattern is $k$-local if this marking can be done while never creating more than $k$ marked blocks. Variables which only occur adjacent to those which are already marked can be marked "for free" – without creating any new blocks, and thus a valid marking sequence allows a sort-of parsing of the pattern whilst maintaining a degree of closeness (locality) to the parts already parsed. The notion of $k$-locality was introduced and further analysed in [15]. With respect to pattern matching, the main result proven in that paper is the following:

**Theorem 13** ([15])**.** MATCH *for* $\mathsf{loc}_k$ *can be decided in* $O(mkn^{\max(3k-1,2k+1)})$ *time, where* $m$ *is the length of the input pattern and* $n$ *is the length of the input word.*

To solve the matching problem for $\mathsf{loc}_k$ we use the following idea. Using a simple dynamic programming approach we can show that, given a pattern $\beta \in (X \cup \Sigma)^*$ of length $m$, we can decide in $O(m^{2k}k)$ time whether $\beta \in \mathsf{loc}_k$, and if the answer is positive, we can produce in the same time a marking sequence witnessing that $\beta$ is $k$-local. As such we can keep track of the marked factors in the pattern, while executing the marking according to the computed marking sequence. We also need now to keep track to which factors of the input word the marked factors correspond. Then we try to assign every new variable so that it fits nicely around the already matched factors. This is done efficiently using a data structure from [57], mentioned also above: given a word $w$ and a one-variable pattern $\gamma$ (so, $|\mathsf{var}(\gamma)| = 1$), one can produce a compact representation of all the $g$ factors of $w$ matching $\gamma$ in $O(|\gamma||w|)$ time; moreover, we can obtain all the $g$ factors of $w$ matching $\gamma$ in $O(|g|)$ time. This allows us to test efficiently which factors of $w$ match any of the one-variable blocks of $\beta$, and, ultimately, to assign a value to each variable. In comparison to the algorithm from [76] for patterns of bounded treewidth, which firstly constructs relational structures from $\alpha$ and $w$, and solves the homomorphism problem on these relational structures (see Section 4.1), the above algorithm exploits directly the locality structure present in the patterns. The advantage of this more focussed approach is that it allows for a considerable improvement in the required time, reducing the exponent of $n$ from $4k + 4$ to $3k - 1$.

# 6    Efficient Pattern Matching Beyond Bounded Treewidth

In [16] the authors tried to identify classes of patterns that do not have bounded treewidth but can still be matched in polynomial time. The idea behind defining such classes was relatively simple: consider generalised repetitions of patterns.

One simple observation is that, if we can match patterns from a class $\mathcal{C}$ in polynomial time, then we can also match repetitions of these patterns in polynomial time: if we wish to check whether $\alpha^k$ matches a word $w$, where $\alpha$ is chosen from the class $\mathcal{C}$ for which we can solve MATCH efficiently, then we can firstly check whether $w = v^k$ for some word $v$, and then check whether $\alpha$ matches $v$, so we can also match $\alpha^k$ efficiently. Moreover, it can be observed that most parameters that lead to efficiently matchable classes, e. g., the scope coincidence degree or locality, are defined independently from the terminal symbols, i. e., via the word obtained after removing all terminals, which shall be called *skeleton* in the following (e. g., the skeleton of $x_1\mathtt{a}x_2\mathtt{b}\mathtt{a}x_1x_2\mathtt{b}$ is $x_1x_2x_1x_2$). As a result, it is possible that a pattern, that is *not* a repetition of any $\alpha \in \mathcal{C}$, has nevertheless a skeleton that is a repetition of a skeleton from $\mathcal{C}$. For example, $\mathtt{a}x_1(x_2)^3x_3\mathtt{b}x_3x_1(x_2)^2\mathtt{b}x_2\mathtt{a}(x_3)^2$ is not a repetition of a non-cross pattern, but its skeleton $(x_1(x_2)^3(x_3)^2)^2$ is. In [16] it is shown that, for some important classes $\mathcal{C}$ of patterns, including $\mathsf{loc}_k$ and $\mathsf{scd}_k$, for constant $k$, the polynomial time solvability of MATCH does not only extend from $\mathcal{C}$ to exact repetitions, but also to such skeleton-repetitions, called $\mathcal{C}$-*repetitions*.

**Theorem 14** ([16])**.** *For $\mathcal{C} \in \{\mathsf{nc}, \mathsf{reg}, \mathsf{loc}_k, \mathsf{scd}_k\}$, solving the matching problem for the class of $\mathcal{C}$-repetitions can be done in polynomial time.*

It is interesting to note that the general treewidth-based framework of polynomial time matching of patterns does not seem to cover a very simple and natural aspect: repetitions of the same pattern. More precisely, if $\mathcal{C}$ is one of the known efficiently matchable classes of patterns, then a repetition $\alpha^k$ for some $\alpha \in \mathcal{C}$ is usually not in $\mathcal{C}$ anymore. In fact, it can be shown that even for patterns $\alpha$ with bounded and very low treewidth, the treewidth of repetitions $\alpha^k$ can be unbounded.

**Theorem 15** ([16])**.** *Let $\mathcal{C}$ be a class of patterns that contains $\mathsf{reg}$. Then the class of $\mathcal{C}$-repetitions contains patterns with arbitrarily large treewidth.*

In particular, the previous theorem holds for the class $\mathsf{reg}$ of regular patterns, arguably the simplest class allowing an unbounded number of variables (note that patterns with a constant number of variables can trivially be matched in polynomial-time). In the same paper it is shown that if the notion of repetition is relaxed further, by considering a setting where the order in which the variables appear is no longer constrained at all (i.e., considering *abelian repetitions* instead of repetitions), then the matching problem is NP-complete. This holds even in the minimal case when the number of repetitions is restricted to two, and that the pattern which is repeated is regular.

## 7 From Locality to Graph Parameters

Following the ideas of Section 3 we explore further the connection between string and graph parameters. The main idea behind such a connection is to reach it by "flattening" a graph into a sequential form, or by "inflating" a string into a graph, so that algorithmic techniques available for each one of these become applicable for the other one as well. In this section, following [13], we are concerned with certain structural parameters (and the problems of computing them) for graphs and strings: the *cutwidth* $\mathsf{cw}(G)$ of a graph $G$ (i. e., the maximum number of "stacked" edges if the vertices of a graph are drawn on a straight line), the *pathwidth* $\mathsf{pw}(G)$ of a graph $G$ (i. e., the minimum width of a tree decomposition the tree structure of which is a path), and the *locality number* $\mathsf{loc}(\alpha)$ of a string $\alpha$ (explained in more detail in Section 4.2). By CUTWIDTH, PATHWIDTH and LOC, we denote the corresponding natural decision problems (i. e., decide whether a given graph has a pathwidth/cutwidth, or a given string has a locality number of at most $k$, for given $k$) and with the prefix MIN, we refer to the minimisation variants. The two former graph-parameters are very classical. Pathwidth is a simple (yet still hard to compute) subvariant of treewidth, which measures how much a graph resembles a path. The problems PATHWIDTH and MINPATHWIDTH are intensively studied (in terms of exact, parameterised and approximation algorithms) and have numerous applications (see the surveys and textbook [9, 54, 7]). CUTWIDTH is the best known example of a whole class of so-called *graph layout problems* (see the survey [20, 71] for detailed information), which are studied since the 1970s and were originally motivated by questions of circuit layouts.

In comparison, the locality number seems a rather simple parameter directly defined on strings, but, however, it bounds the treewidth of the string (in the sense defined in Section 4.1), and the corresponding marking sequences can be seen as instructions for a dynamic programming algorithm for matching the pattern. In this way, it resembles a bit to the way the pathwidth and treewidth of graphs are used in algorithmic settings. Moreover, compared to other "tractability-parameters" of strings, it seems to cover best the treewidth of a string, but it also cannot be efficiently computed compared to the other simpler parameters.

Going more into detail, for LOC, exact exponential-time algorithms are not hard to be devised [15] but whether it can be solved in polynomial-time, or whether it is at least fixed-parameter tractable was left open in the paper where this measure was introduced. On the other hand, PATHWIDTH and CUTWIDTH are known NP-complete problems, fixed-parameter tractable with respect to parameter $\mathsf{pw}(G)$ or $\mathsf{cw}(G)$, respectively (even with "linear" fpt-algorithms with running-time $g(k)\,\mathrm{O}(n)$ [8, 10, 80]). With respect to approximation, their minimisation variants have received a lot of attention, mainly because they yield (like many other graph parameters) general algorithmic approaches for numerous graph problems, i.e., a good linear arrangement or path-decomposition can often be used to design a dynamic programming (or even divide and conquer) algorithm for other problems. The best known approximation algorithms for the problems MINPATHWIDTH and MINCUTWIDTH (with approximations ratios of $\mathrm{O}(\sqrt{\log(\mathsf{opt})}\log(n))$ and $\mathrm{O}(\log^2(n))$, respectively) follow from approximations of vertex separators (see [25]) and edge separators (see [58]), respectively.

There are two natural approaches to represent a word $\alpha$ over alphabet $\Sigma$ as a graph $G_\alpha = (V_\alpha, E_\alpha)$: (1) $V_\alpha = \{1, 2, \ldots, |\alpha|\}$ and the edges are somehow used to represent the actual symbols (note that this is the case for the standard graph representation of patterns defined in Section 4.1), or (2) $V_\alpha = \Sigma$ and the edges are somehow used to represent the positions of $\alpha$. A reduction of type (2) can be defined such that $|E_\alpha| = \mathrm{O}(|\alpha|)$ and $\mathsf{cw}(G_\alpha) = 2\,\mathsf{loc}(\alpha)$, and a reduction of type (1) can be defined such that $|E_\alpha| = \mathrm{O}(|\alpha|^2)$ and $\mathsf{loc}(\alpha) \leq \mathsf{pw}(G_\alpha) \leq 2\,\mathsf{loc}(\alpha)$. Since these reductions are parameterised reductions and also allow to transfer approximation results, one may conclude that LOC is fixed-parameter tractable if parameterised by $|\Sigma|$ (note that for parameter $|\Sigma|$ a simple, but less efficient fpt-algorithm is trivially obtained by simply enumerating all marking sequences) or by the locality number, and also that there is a polynomial-time $\mathrm{O}(\sqrt{\log(\mathsf{opt})}\log(n))$-approximation algorithm for MINLOC.

In addition, one can represent an arbitrary multi-graph $G = (V, E)$ by a word $\alpha_G$ over alphabet $V$ with $|\alpha_G| = |E|$ and $\mathsf{cw}(G) = \mathsf{loc}(\alpha)$. This describes a Turing-reduction from CUTWIDTH to LOC which also allows to transfer approximation results between the minimisation variants. As a result, LOC is NP-complete. Finally, by plugging together the reductions from MINCUTWIDTH to MINLOC and from MINLOC to MINPATHWIDTH, one obtains a reduction which transfers approximation results from MINPATHWIDTH to MINCUTWIDTH, which yields an $\mathrm{O}(\sqrt{\log(\mathsf{opt})}\log(n))$-approximation algorithm for MINCUTWIDTH. This result from [13] improved, for the first time since 1999, the best approximation for CUTWIDTH from [58]. Interestingly, this improvement appeared as a side-product of relating string-parameters with graph-parameters.

**Theorem 16** ([13])**.** *There is an* $\mathrm{O}(\sqrt{\log(\mathbf{opt})}\log(h))$*-approximation algorithm (running in polynomial time) for* MINCUTWIDTH *on multigraphs with $h$ edges. In particular, this yields an* $\mathrm{O}(\sqrt{\log(\mathbf{opt})}\log(n))$*-approximation algorithm for* MINCUTWIDTH *for graphs.*

Moreover, this approach allows also for establishing a direct connection between cutwidth and pathwidth, which preserves the good algorithmic properties, and has not yet been reported in the literature so far. This is rather surprising, since CUTWIDTH and PATHWIDTH have been jointly investigated in the context of exact and approximation algorithms, especially in terms of balanced vertex and edge separators. We think that a reason for overlooking this connection might be that it is less obvious on the graph level and becomes more apparent if linked via the string parameter of locality, emphasising, as such, the value of such mixed approaches.

# 8 Extensions

## 8.1 Injectivity

In our setting, the substitutions that map variables to words are *not* required to be injective, i. e., different variables can be mapped to the same word. However, the requirement of injectivity is natural in some contexts. For example, in the pattern matching community, the first mentioning of pattern matching with variables concerns the case where variables have to be substituted by single symbols and in an injective way. More precisely, this *parameterised pattern matching* was introduced in [3] to formalise the problem of detecting code clones (i. e., we want to find code segments that are created by copying some code blocks and renaming program variables (this renaming will be injective, since otherwise the semantic of the code might change)). More generally speaking, the injectivity condition is appropriate whenever we know a priori that different variables should always refer to different words (e. g., when matching the pattern

$$x_1 \text{ name: } y \text{ ; address: } z \text{ ; } x_2 \text{ name: } y \text{ ; address: } z \, x_3$$

in order to check whether there is a repetition of some name-address data tuple, then it is likely that we can assume injectivity).

Depending on the actual variant, the injectivity condition can make the matching problem harder or easier. In [26], it is shown that it is NP-hard to decide for a given word $w$ and an integer $k$ whether $w$ can be factorised into at least $k$ *pairwise different* factors. This immediately implies that the injectivity condition makes the matching problem NP-hard even for the "trivial" pattern class $\{x_1 x_2 \ldots x_n \mid n \geq 1\}$ (note that this is even a subset of the class reg of regular patterns). On the other hand, if we have an upper bound on $|\Sigma|$ and $\max\{|h(x)| \mid x \in X\}$ (recall that this case is still NP-hard even for bounds 2 and 1, respectively; see Theorem 2) then also the total number of possible substitution words is bounded; thus, the injectivity condition bounds the total number of variables and therefore the matching problem becomes tractable (see [28]). A similar observation can be made with respect to fixed-parameter tractability if we parameterise by $|\Sigma|$ and $\max\{|h(x)| \mid x \in X\}$ (see [29]).

## 8.2 Word Equations

A *word equation* is an equality $\alpha = \beta$, where $\alpha$ and $\beta$ are patterns with variables, e. g., $\alpha = x_1 \mathsf{ab} x_2$ and $\beta = \mathsf{a} x_1 x_2 \mathsf{b}$ define the equation $x_1 \mathsf{ab} x_2 = \mathsf{a} x_1 x_2 \mathsf{b}$. A *solution* to an equation $\alpha = \beta$ is a substitution $h : (\mathsf{var}(\alpha) \cup \mathsf{var}(\beta)) \to \Sigma^*$ (in the sense defined in Section 2) that satisfies $h(\alpha) = h(\beta)$. For the example equation from above, the solutions are the substitutions $h$ with $h(x_1) = \mathsf{a}^k$, for $k \geq 0$, and $h(x_2) = \mathsf{b}^\ell$, for $\ell \geq 0$.

The study of word equations (or the existential theory of equations over free monoids) is an important topic found at the intersection of algebra and computer science, with significant connections to, e.g., combinatorial group or monoid theory [64, 63, 21], unification [78, 46, 47]), and, more recently, data base theory [33, 32].

The central computational problem for word equations is the satisfiability problem, i. e., the problem of deciding whether a given word equation $\alpha = \beta$ has a solution or not. In this regard, the matching problem for patterns with variables describes just the special case of the satisfiability problem for word equations where one side of the equation is a terminal word, e. g., $x_1 \mathsf{ab} x_1 x_2 \mathsf{c} x_2 x_1 = \mathsf{babbacab}$ is an instance of the matching problem already mentioned in the introduction, phrased as a word equation. Consequently, the satisfiability problem is intractable, even for very strongly restricted cases (see Theorems 2 and 3). Also note that it has been shown in [22] that the solvability problem remains NP-hard if every variable has at most two occurrences in $\alpha\beta$ (called *quadratic* equations), but the proof of [22] actually talks about the matching problem for patterns with at most two occurrences per variable.

While the matching problem for patterns with variables is trivially decidable, it is not at all obvious how to solve the satisfiability problem for word equations. In fact, the question whether it is decidable was initially approached with the expectation that it will be answered in the negative. It was, however, shown to be decidable by Makanin [65] (see Chapter 12 of [62] for a survey). Later it was shown that the satisfiability problem is in PSPACE by Plandowski [72]; a new proof of this result was obtained in [49], based on a new simple technique called recompression. There

14

are also cases when the satisfiability problem is tractable. For instance, word equations with only one variable can be solved in linear time in the size of the equation, see [48]; equations with two variables can be solved in time $O(|\alpha\beta|^5)$, see [19].

Given the fact that there are many structural restrictions of patterns that yield tractability (with respect to the matching problem, see Section 4), the question naturally arises how the complexity of the satisfiability problem for word equations (which are essentially equations of patterns) behaves if these restrictions are applied to word equations. More precisely, while each class of patterns with NP-hard matching problem yields a class of word equations with NP-hard satisfiability problem, the hardness of the satisfiability problem for equations with sides in some efficiently matchable class of patterns is no longer immediate. An investigation of that question was initiated in [66], where the following results were obtained. Firstly, the satisfiability problem for non-cross word equations (i. e., word equations for which both sides are non-cross) remains NP-hard. In particular, solving non-cross equations $\alpha = \beta$ where each variable occurs at most three times, at most twice in $\alpha$ and exactly once in $\beta$, is NP-hard (note that this constitutes the first NP-hardness result for word equations that is not a direct conclusion from a hardness result for the matching problem). Secondly, the satisfiability of one-repeated variable equations (i. e., at most one variable occurs more than once in $\alpha\beta$, but arbitrarily many other variables occur only once) having at least one non-repeated variable on each side, was shown to be trivially in P.

In [18], it is shown that it is (still) NP-hard to solve regular ordered word equations. More precisely, these are word equations where each side is a regular pattern and the order of the variables in both sides is the same (it is, however, possible that some variables only occur on one side of the equation), e. g., $x_1\mathtt{a}x_2\mathtt{ba}x_3x_4 = \mathtt{b}x_1x_3\mathtt{aa}x_4$ is a regular ordered word equation. They are particular cases of both quadratic equations and non-cross equations, so the reductions showing the hardness of solving these more general equations do not carry over. In particular, note that the class of regular patterns is arguably the most simple class of patterns in terms of their matching complexity (see Section 5.1).

The respective hardness reduction relied on some deep word-combinatorics ideas. As a first step, a reachability problem for a certain type of (regulated) string rewriting systems was introduced, and showed it is NP-complete. This was achieved via a reduction from the problem 3-Partition [38], which is strongly NP-complete. Then it was shown that this reachability problem can be reduced to the satisfiability of regular-ordered word equations; in this reduction the applications of the rewriting rules of the system were encoded into the periods of the words assigned to the variables in a solution to the equation. The main technicality was to make sure to only use one occurrence of each variable per side, and moreover to even have the variables in the same order in both sides. This result exhibits the arguably structurally-simplest class of word equations for which the satisfiability problem is NP-hard.

The main open problem in the area of word equations remains, even for simple subclasses such as regular equations or quadratic equations, to show that the satisfiability problem of word equations of the respective types is in NP (note that this was already explicitly posed as an open question for the class of quadratic word equations in [22]).

# 9 Conclusions

In this work we tried to survey several results related to the problem of matching patterns with variables, that seem important to us. While this work is clearly not exhaustive, it is aimed to offer a basic understanding of the problems and state of the art in this area.

From an algorithmic point of view, the results we covered provide a wide variety of classes of patterns with variables, for which MATCH can be efficiently solved. Moreover, as explained in Section 4.3, it is usually easy to check whether a pattern belongs to one of these classes. So, putting it all together, one could use the following approach when trying to match a pattern, rather than just using an exponential time algorithm (based, e.g., on general SMT-solvers, or on the theory of string solving [6, 81]). First, check if the pattern belongs to one of the classes for which efficient matching algorithms are known and, then, use this algorithm; only use a general algorithm when no customised one can be applied. Identifying more natural pattern classes for which MATCH can be solved efficiently appears, as such, as a rather useful task. Following the practically motivated challenges that arise from the area of string solving, one could also

try to find efficient matching algorithms for various classes of patterns, enhanced with various constraints: regular constraints, length constraints, etc.

As an important part of this survey deals with polynomial time algorithms, it is natural to also ask whether they are optimal or not. This kind of questions are the focus of the area of fine-grained complexity (see, e.g., the survey [11] and the citations therein). It would be interesting to see, using tools from this area, whether one can show lower bounds for the MATCH problems for different classes of patterns.

In the light of the results from [13], it seems that exploring the connections between string parameters and parameters for other classes of objects could lead to some interesting results in both worlds. So, it also seems like an interesting challenge to explore what the structural parameters of strings that we explored here (and maybe some other new ones) mean when various other types of data are represented as strings, and what consequences can be derived from such a representation.

Finally, the area of word equations abounds with open problems. As mentioned, it is not even clear whether the satisfiability of regular or quadratic equations is in NP. So even if we restrict to equations with structurally simple left and right hand sides, the complexity of solving equations is not known. Such problems become even more involved when we consider equations with various types of constraints (e.g., length or regular). For instance, the decidability of general word equations with length constraints is a long standing open problem, but it is already an interesting open question for simpler cases (once again: regular or quadratic equations); see, e.g., [44, 17, 59], and the references therein. It seems interesting to us whether some of the ideas used in matching patterns can be transferred to solving (simplified) word equations, with or without constraints.

# References

[1] Amihood Amir and Igor Nor. Generalized function matching. *Journal of Discrete Algorithms*, 5:514–523, 2007.

[2] Dana Angluin. Finding patterns common to a set of strings. *Journal of Computer and System Sciences*, 21:46–62, 1980.

[3] Brenda S. Baker. Parameterized pattern matching: Algorithms and applications. *Journal of Computer and System Sciences*, 52:28–42, 1996.

[4] Hideo Bannai, Tomohiro I, Shunsuke Inenaga, Yuto Nakashima, Masayuki Takeda, and Kazuya Tsuruta. The "runs" theorem. *SIAM J. Comput.*, 46(5):1501–1514, 2017.

[5] Pablo Barceló, Leonid Libkin, Anthony W. Lin, and Peter T. Wood. Expressive languages for path queries over graph-structured data. *ACM Transactions on Database Systems*, 37, 2012.

[6] Clark Barrett, Christopher L Conway, Morgan Deters, Liana Hadarean, Dejan Jovanović, Tim King, Andrew Reynolds, and Cesare Tinelli. CVC4. In *International Conference on Computer Aided Verification*, pages 171–177. Springer, 2011.

[7] Hans L. Bodlaender. A tourist guide through treewidth. *Acta Cybern.*, 11(1-2):1–21, 1993.

[8] Hans L. Bodlaender. A linear-time algorithm for finding tree-decompositions of small treewidth. *SIAM J. Comput.*, 25(5):1305–1317, 1996. `doi:10.1137/s0097539793251219`.

[9] Hans L. Bodlaender. A partial $k$-arboretum of graphs with bounded treewidth. *Theor. Comput. Sci.*, 209(1-2):1–45, 1998. `doi:10.1016/S0304-3975(97)00228-4`.

[10] Hans L. Bodlaender. Fixed-parameter tractability of treewidth and pathwidth. In Hans L. Bodlaender, Rod Downey, Fedor V. Fomin, and Dániel Marx, editors, *The Multivariate Algorithmic Revolution and Beyond*, volume 7370 of *LNCS*, pages 196–227, 2012.

[11] Karl Bringmann. Fine-Grained Complexity Theory (Tutorial). In Rolf Niedermeier and Christophe Paul, editors, *36th International Symposium on Theoretical Aspects of Computer Science (STACS 2019)*, volume 126 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 4:1–4:7, Dagstuhl, Germany, 2019. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik. URL: `http://drops.dagstuhl.de/opus/volltexte/2019/10243`, `doi:10.4230/LIPIcs.STACS.2019.4`.

[12] Cezar Câmpeanu, Kai Salomaa, and Sheng Yu. A formal study of practical regular expressions. *International Journal of Foundations of Computer Science*, 14:1007–1018, 2003.

[13] Katrin Casel, Joel D. Day, Pamela Fleischmann, Tomasz Kociumaka, Florin Manea, and Markus L. Schmid. Graph and string parameters: Connections between pathwidth, cutwidth and the locality number. *CoRR, to appear in Proc. ICALP 2019,*, abs/1902.10983, 2019. URL: `http://arxiv.org/abs/1902.10983`, `arXiv:1902.10983`.

[14] Maxime Crochemore. An optimal algorithm for computing the repetitions in a word. *Information Processing Letters*, 12(5):244–250, 1981.

[15] Joel D. Day, Pamela Fleischmann, Florin Manea, and Dirk Nowotka. Local patterns. In *37th IARCS Annual Conference on Foundations of Software Technology and Theoretical Computer Science, FSTTCS 2017*, pages 24:1–24:14, 2017.

[16] Joel D. Day, Pamela Fleischmann, Florin Manea, Dirk Nowotka, and Markus L. Schmid. On matching generalised repetitive patterns. In *Proc. Developments in Language Theory - 22nd International Conference, DLT 2018*, volume 11088 of *Lecture Notes in Computer Science*, pages 269–281. Springer, 2018.

[17] Joel D. Day, Vijay Ganesh, Paul He, Florin Manea, and Dirk Nowotka. The satisfiability of word equations: Decidable and undecidable theories. In *Proc. 12th International Conference Reachability Problems, RP 2018*, volume 11123 of *Lecture Notes in Computer Science*, pages 15–29. Springer, 2018.

[18] Joel D. Day, Florin Manea, and Dirk Nowotka. The hardness of solving simple word equations. In *Proc. MFCS 2017*, volume 83 of *LIPIcs*, pages 18:1–18:14, 2017.

[19] R. Dąbrowski and W. Plandowski. Solving two-variable word equations. In *Proc. 31th International Colloquium on Automata, Languages and Programming, ICALP 2004*, volume 3142 of *Lecture Notes in Computer Science*, pages 408–419, 2004.

[20] Josep Díaz, Jordi Petit, and Maria Serna. A survey of graph layout problems. *ACM Comput. Surv.*, 34(3):313–356, September 2002. `doi:10.1145/568522.568523`.

[21] V. Diekert, A. Jez, and M. Kufleitner. Solutions of word equations over partially commutative structures. In *Proc. 43rd International Colloquium on Automata, Languages, and Programming, ICALP 2016*, volume 55 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 127:1–127:14, 2016.

[22] V. Diekert and J. M. Robson. On quadratic word equations. In *Proc. 16th Annual Symposium on Theoretical Aspects of Computer Science, STACS 1999*, volume 1563 of *Lecture Notes in Computer Science*, pages 217–226, 1999.

[23] Rodney G. Downey and Michael R. Fellows. *Fundamentals of Parameterized Complexity*. Texts in Computer Science. Springer, 2013.

[24] Thomas Erlebach, Peter Rossmanith, Hans Stadtherr, Angelika Steger, and Thomas Zeugmann. Learning one-variable pattern languages very efficiently on average, in parallel, and by asking queries. *Theoretical Computer Science*, 261:119–156, 2001.

[25] Uriel Feige, MohammadTaghi HajiAghayi, and James R. Lee. Improved approximation algorithms for minimum weight vertex separators. *SIAM J. Comput.*, 38(2):629–657, 2008. `doi:10.1137/05064299x`.

[26] Henning Fernau, Florin Manea, Robert Mercas, and Markus L. Schmid. Pattern matching with variables: Fast algorithms and new hardness results. In *32nd International Symposium on Theoretical Aspects of Computer Science, STACS 2015*, pages 302–315, 2015.

[27] Henning Fernau, Florin Manea, Robert Mercas, and Markus L. Schmid. Revisiting shinohara's algorithm for computing descriptive patterns. *Theoretical Computer Science*, 733:44–54, 2018.

[28] Henning Fernau and Markus L. Schmid. Pattern matching with variables: A multivariate complexity analysis. *Inf. Comput.*, 242:287–305, 2015.

[29] Henning Fernau, Markus L. Schmid, and Yngve Villanger. On the parameterised complexity of string morphism problems. *Theory Comput. Syst.*, 59(1):24–51, 2016.

[30] Jörg Flum and Martin Grohe. *Parameterized Complexity Theory*. Texts in Theoretical Computer Science. An EATCS Series. Springer, 2006.

[31] Dominik D. Freydenberger. Extended regular expressions: Succinctness and decidability. *Theory of Computing Systems*, 53:159–193, 2013.

[32] Dominik D. Freydenberger. A logic for document spanners. *Theory of Computing Systems*, Sep 2018. URL: https://doi.org/10.1007/s00224-018-9874-1.

[33] Dominik D. Freydenberger and Mario Holldack. Document spanners: From expressive power to decision problems. *Theory of Computing Systems*, 62(4):854–898, 2018.

[34] Dominik D. Freydenberger and Daniel Reidenbach. Bad news on decision problems for patterns. *Inf. Comput.*, 208(1):83–96, 2010.

[35] Dominik D. Freydenberger and Daniel Reidenbach. Existence and nonexistence of descriptive patterns. *Theor. Comput. Sci.*, 411(34-36):3274–3286, 2010.

[36] Dominik D. Freydenberger and Daniel Reidenbach. Inferring descriptive generalisations of formal languages. *J. Comput. Syst. Sci.*, 79(5):622–639, 2013.

[37] Jeffrey E. F. Friedl. *Mastering Regular Expressions*. O'Reilly, Sebastopol, CA, third edition, 2006.

[38] Michael R. Garey and David S. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman & Co., New York, NY, USA, 1979.

[39] P. Gawrychowski, F. Manea, and D. Nowotka. Testing generalised freeness of words. In *STACS 2014*, volume 25 of *LIPIcs*, pages 337–349. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2014.

[40] Pawel Gawrychowski, Tomohiro I, Shunsuke Inenaga, Dominik Köppl, and Florin Manea. Tighter bounds and optimal algorithms for all maximal $\alpha$-gapped repeats and palindromes - finding all maximal $\alpha$-gapped repeats and palindromes in optimal worst case time on integer alphabets. *Theory Comput. Syst.*, 62(1):162–191, 2018.

[41] Pawel Gawrychowski, Florin Manea, Robert Mercas, and Dirk Nowotka. Hide and seek with repetitions. *J. Comput. Syst. Sci.*, 101:42–67, 2019. URL: https://doi.org/10.1016/j.jcss.2018.10.004, doi:10.1016/j.jcss.2018.10.004.

[42] Pawel Gawrychowski, Florin Manea, Robert Mercas, Dirk Nowotka, and Catalin Tiseanu. Finding pseudo-repetitions. In *30th International Symposium on Theoretical Aspects of Computer Science, STACS 2013, February 27 - March 2, 2013, Kiel, Germany*, volume 20 of *LIPIcs*, pages 257–268, 2013.

[43] M. Geilke and S. Zilles. Learning relational patterns. In *Proc. 22nd International Conference on Algorithmic Learning Theory, ALT 2011*, volume 6925 of *Lecture Notes in Computer Science*, pages 84–98, 2011.

[44] Simon Halfon, Philippe Schnoebelen, and Georg Zetzsche. Decidability, complexity, and expressiveness of first-order logic over the subword ordering. In *Proc. 32nd Annual ACM/IEEE Symposium on Logic in Computer Science, LICS 2017*, pages 1–12. IEEE Computer Society, 2017.

[45] Oscar H. Ibarra, Ting-Chuen Pong, and Stephen M. Sohn. A note on parsing pattern languages. *Pattern Recognition Letters*, 16:179–182, 1995.

[46] J. Jaffar. Minimal and complete word unification. *Journal of the ACM*, 37(1):47–85, 1990.

[47] A. Jez. Context unification is in PSPACE. In *Proc. 41st International Colloquium on Automata, Languages, and Programming, ICALP 2014*, volume 8573 of *Lecture Notes in Computer Science*, pages 244–255. Springer, 2014.

[48] A. Jeż. One-variable word equations in linear time. *Algorithmica*, 74:1–48, 2016.

[49] A. Jeż. Recompression: A simple and powerful technique for word equations. *Journal of the ACM*, 63, 2016.

[50] Tao Jiang, Arto Salomaa, Kai Salomaa, and Sheng Yu. Decision problems for patterns. *J. Comput. Syst. Sci.*, 50(1):53–63, 1995.

[51] Juhani Karhumäki, Wojciech Plandowski, and Filippo Mignosi. The expressibility of languages and relations by word equations. *Journal of the ACM*, 47:483–505, 2000.

[52] Juha Kärkkäinen, Peter Sanders, and Stefan Burkhardt. Linear work suffix array construction. *Journal of the ACM*, 53:918–936, 2006.

[53] Michael J. Kearns and Leonard Pitt. A polynomial-time algorithm for learning $k$-variable pattern languages from examples. In *Proceedings of the Second Annual Workshop on Computational Learning Theory, COLT 1989, Santa Cruz, CA, USA, July 31 - August 2, 1989.*, pages 57–71, 1989.

[54] Ton Kloks. *Treewidth, Computations and Approximations*, volume 842 of *Lecture Notes in Computer Science*. Springer, 1994. `doi:10.1007/BFb0045375`.

[55] Roman Kolpakov and Gregory Kucherov. Searching for gapped palindromes. *Theor. Comput. Sci.*, 410(51):5365–5373, 2009.

[56] Roman Kolpakov, Mikhail Podolskiy, Mikhail Posypkin, and Nickolay Khrapov. Searching of gapped repeats and subrepetitions in a word. In *Proc. 25th Annual Symposium Combinatorial Pattern Matching, CPM 2014*, volume 8486 of *Lecture Notes in Computer Science*, pages 212–221. Springer, 2014.

[57] Dmitry Kosolobov, Florin Manea, and Dirk Nowotka. Detecting one-variable patterns. In *String Processing and Information Retrieval - 24th International Symposium, SPIRE 2017, Palermo, Italy, September 26-29, 2017, Proceedings*, pages 254–270, 2017.

[58] Tom Leighton and Satish Rao. Multicommodity max-flow min-cut theorems and their use in designing approximation algorithms. *J. ACM*, 46(6):787–832, November 1999. `doi:10.1145/331524.331526`.

[59] Anthony W. Lin and Rupak Majumdar. Quadratic word equations with length constraints, counter systems, and presburger arithmetic with divisibility. In *Proc. 16th International Symposium Automated Technology for Verification and Analysis, ATVA 2018*, volume 11138 of *Lecture Notes in Computer Science*, pages 352–369. Springer, 2018.

[60] M. Lothaire. *Combinatorics on Words*. Cambridge University Press, 1997.

[61] M. Lothaire. *Algebraic Combinatorics on Words*, chapter 3. Cambridge University Press, Cambridge, New York, 2002.

[62] M. Lothaire. *Algebraic Combinatorics on Words*. Cambridge University Press, Cambridge, New York, 2002.

[63] R. C. Lyndon. Equations in free groups. *Transactions of the American Mathematical Society*, 96:445–457, 1960.

[64] R. C. Lyndon and P. E. Schupp. *Combinatorial Group Theory*. Springer, 1977.

[65] G. S. Makanin. The problem of solvability of equations in a free semigroup. *Matematicheskii Sbornik*, 103:147–236, 1977.

[66] F. Manea, D. Nowotka, and M. L. Schmid. On the solvability problem for restricted classes of word equations. In *Proc. 20th International Conference on Developments in Language Theory, DLT 2016*, volume 9840 of *Lecture Notes in Computer Science*, pages 306–318. Springer, 2016.

[67] Alexandru Mateescu and Arto Salomaa. Finite degrees of ambiguity in pattern languages. *RAIRO Informatique Théoretique et Applications*, 28:233–253, 1994.

[68] Alexandru Mateescu and Arto Salomaa. Aspects of classical language theory. In *Handbook of Formal Languages (1)*, pages 175–251. 1997.

[69] Yen K. Ng and Takeshi Shinohara. Developments from enquiries into the learnability of the pattern languages from positive data. *Theoretical Computer Science*, 397:150–165, 2008.

[70] Sebastian Ordyniak and Alexandru Popa. A parameterized study of maximum generalized pattern matching problems. *Algorithmica*, 75:1–26, 2016.

[71] Jordi Petit. Addenda to the survey of layout problems. *Bulletin of the EATCS*, 105:177–201, 2011. URL: http://eatcs.org/beatcs/index.php/beatcs/article/view/98.

[72] W. Plandowski. An efficient algorithm for solving word equations. In *Proceedings of the 38th Annual ACM Symposium on Theory of Computing, STOC 2006*, pages 467–476, 2006.

[73] Daniel Reidenbach. A non-learnable class of e-pattern languages. *Theor. Comput. Sci.*, 350(1):91–102, 2006.

[74] Daniel Reidenbach. An examination of ohlebusch and ukkonen's conjecture on the equivalence problem for e-pattern languages. *Journal of Automata, Languages and Combinatorics*, 12(3):407–426, 2007.

[75] Daniel Reidenbach. Discontinuities in pattern inference. *Theoretical Computer Science*, 397:166–193, 2008.

[76] Daniel Reidenbach and Markus L. Schmid. Patterns with bounded treewidth. *Information and Computation*, 239:87–99, 2014.

[77] Markus L. Schmid. A note on the complexity of matching patterns with variables. *Inf. Process. Lett.*, 113(19-21):729–733, 2013.

[78] K. U. Schulz. Word unification and transformation of generalized equations. *Journal of Automated Reasoning*, 11:149–184, 1995.

[79] Takeshi Shinohara. Polynomial time inference of pattern languages and its application. In *Proceedings of 7th IBM Symposium on Mathematical Foundations of Computer Science, MFCS*, pages 191–209, 1982.

[80] Dimitrios M. Thilikos, Maria J. Serna, and Hans L. Bodlaender. Cutwidth I: A linear time fixed parameter algorithm. *J. Algorithms*, 56(1):1–24, 2005. doi:10.1016/j.jalgor.2004.12.001.

[81] Yunhui Zheng, Vijay Ganesh, Sanu Subramanian, Omer Tripp, Murphy Berzish, Julian Dolby, and Xiangyu Zhang. Z3str2: an efficient solver for strings, regular expressions, and length constraints. *Formal Methods in System Design*, 50(2-3):249–288, 2017.