

Palindromic Decompositions with Gaps and Errors

Michał Adamczyk¹, Mai Alzamel², Panagiotis Charalampopoulos², Costas S. Iliopoulos²,
and Jakub Radoszewski^{1,2}

¹ Institute of Informatics, University of Warsaw, Warsaw, Poland,
[michal.adamczyk,jrad]@mimuw.edu.pl

² Department of Informatics, King's College London, London, UK,
[mai.alzamel,panagiotis.charalampopoulos,costas.iliopoulos]@kcl.ac.uk

Abstract

Identifying palindromes in sequences has been an interesting line of research in combinatorics on words and also in computational biology, after the discovery of the relation of palindromes in the DNA sequence with the HIV virus. Efficient algorithms for the factorization of sequences into palindromes and maximal palindromes have been devised in recent years. We extend these studies by allowing gaps in decompositions and errors in palindromes, and also imposing a lower bound to the length of acceptable palindromes.

We first present an algorithm for obtaining a palindromic decomposition of a string of length n with the minimal total gap length in time $\mathcal{O}(n \log n \cdot g)$ and space $\mathcal{O}(n \cdot g)$, where g is the number of allowed gaps in the decomposition. We then consider a decomposition of the string in maximal δ -palindromes (i.e. palindromes with δ errors under the edit or Hamming distance) and g allowed gaps. We present an algorithm to obtain such a decomposition with the minimal total gap length in time $\mathcal{O}(n \cdot (g + \delta))$ and space $\mathcal{O}(n \cdot g)$.

1 Introduction

A palindrome is a symmetric word that reads the same backward and forward. The detection of palindromes is a classical and well-studied problem in computer science, language theory and algorithm design with a lot of variants arising out of different practical scenarios. String and sequence algorithms related to palindromes have long drawn the attention of stringology researchers [4, 12, 19]. Interestingly, in the seminal Knuth-Morris-Pratt paper presenting the well-known string matching algorithm [18], a problem related to palindrome recognition was also considered. In combinatorics on words, for example, studies have investigated the inhabitation of palindromes in Fibonacci words or Sturmian words in general [7, 8]. There is also an interesting conjecture related to periodicity of infinite strings whose every factor can be decomposed into a bounded number of palindromes [10].

In computational biology, palindromes play an important role in regulation of gene activity and other cell processes because these are often observed near promoters, introns and specific untranslated regions. Hairpins (also called complemented palindromes) in the HIV virus are strings of the form $x\bar{x}^R$, where \bar{x}^R is the reverse complement of x , while (even) palindromes are strings of the form xx^R . Algorithms for detecting palindromes can often be adapted to compute hairpins as well. Hence, we can identify palindromes in the DNA sequence and then align the part of the DNA sequence that contains them with the HIV virus.

In the beginnings of algorithmic study of palindromes, Manacher discovered an on-line sequential algorithm that finds all initial palindromes in a string [21]. A string $S[1..n]$ is said to have an initial palindrome of length k if $S[1..k]$ is a palindrome. Later Apostolico et al. observed that the algorithm given by [21] is able to find all maximal palindromic factors in the string in $\mathcal{O}(n)$ time [3]. Gusfield gave another linear-time algorithm to find all maximal palindromes in a string [15]. He also discussed the relation between biological sequences and gapped (separated) palindromes (i.e. strings of the form $xv\bar{x}^R$). Gupta et al. [14] presented

an $\mathcal{O}(n)$ -time algorithm to compute specific classes—based on length constraints—of such palindromes. Algorithms for finding the so-called gapped palindromes were also considered in [11, 19]. (In our study, we consider gaps *between* palindromes, not inside them.)

A problem that gained significant attention recently was decomposing a string into a minimal number of palindromes; any such decomposition is called a palindromic factorization. Fici et al. [9] presented an on-line $\mathcal{O}(n \log n)$ -time algorithm for computing a palindromic factorization of a string of length n . A similar on-line algorithm was presented by I et al. [17] as well as an on-line algorithm with the same time complexity to factorize a string into maximal palindromes. Alatabbi et al. gave an off-line $\mathcal{O}(n)$ -time algorithm for the latter problem [2]. In addition, Rubinchik and Shur [23] devised an $\mathcal{O}(n)$ -sized data structure that helps locate palindromes in a string; they also show how it can be used to compute the palindromic factorization of a string in $\mathcal{O}(n \log n)$ time.

A similar problem, first studied by Galil and Seiferas in [13], asked whether a given string can be decomposed into k palindromes. Galil and Seiferas [13] presented an on-line $\mathcal{O}(n)$ -time algorithm for $k = 1, 2$ and an off-line $\mathcal{O}(n)$ -time algorithm for $k = 3, 4$. In 2014, Kosolobov et al. presented an on-line $\mathcal{O}(kn)$ -time algorithm to decide this for arbitrary k [20].

Our work is a continuation of this line of research, motivated by possible errors and inconsistencies in the biological data. We extend the previous work by introducing a constraint on the length of the palindromes and allowing gaps and errors in the decompositions. By *gaps* we mean regions of the string that are not decomposed into palindromes of sufficient length. We allow *errors* in the palindromes, so that a *palindrome with errors* is a string having a small Hamming or edit distance from an ideal palindrome. We present two approaches for decomposing a string into sufficiently long palindromes; one allowing only gaps in the decomposition and the other allowing both gaps in the decomposition and errors in the palindromes. We first present an algorithm that computes a palindromic decomposition with the minimal total gap length of a string of length n in time $\mathcal{O}(n \log n \cdot g)$ and space $\mathcal{O}(n \cdot g)$, where g is the number of allowed gaps. Secondly, we present an $\mathcal{O}(n \cdot (g + \delta))$ -time and $\mathcal{O}(n \cdot g)$ -space algorithm for the decomposition of a string of length n into maximal palindromes with at most δ errors each, under the Hamming or edit distance, and g allowed gaps. The algorithms can be applied for both standard and complemented palindromes.

Our first result is a multistage generalization of the algorithm computing a palindromic factorization from [9]. We show that their approach can be extended to handle gaps, a constraint on the length of palindromes, and complemented palindromes as well. In our second result, the first step is the computation of maximal palindromes with at most δ errors each and the second step is a dynamic programming approach that takes as input these maximal palindromes and returns a decomposition of the string with the minimal total gap length. An $\mathcal{O}(n \cdot \delta)$ computation of maximal palindromes under the Hamming distance was presented in [15]. We present an equally efficient approach for computing all maximal palindromes under the edit distance. After our conference publication [1], we have learnt about alternative algorithms computing this type of maximal approximate palindromes [22, 16]. The approach of [22] works in $\mathcal{O}(n \cdot \delta^2)$ time. The algorithm from [16] has same time complexity as ours ($\mathcal{O}(n \cdot \delta)$), however, it is much more complex, as it was designed for a more general problem.

2 Notation and terminology

Let $S = S[1]S[2] \cdots S[n]$ be a *string* of length $|S| = n$ over an alphabet Σ . We consider the case of an integer alphabet; in this case each letter can be replaced by its rank so that the resulting string consists of integers in the range $\{1, \dots, n\}$. For two positions i and j , where $1 \leq i \leq j \leq n$, in S , we denote the *factor* $S[i]S[i+1] \cdots S[j]$ of S by $S[i..j]$. We denote the reverse string of S by S^R , i.e. $S^R = S[n]S[n-1] \cdots S[1]$. The empty string (denoted by ε) is the unique string over Σ of length 0. A string S is said to be a *palindrome* if and only if $S = S^R$. If $S[i..j]$ is a palindrome, the number $\frac{i+j}{2}$ is called the center of $S[i..j]$. Let $S[i..j]$, where $1 \leq i \leq j \leq n$, be a palindromic factor in S . It is said to be a *maximal palindrome* if there is no longer palindrome in S with center $\frac{i+j}{2}$. Note that a maximal palindrome can be a factor of another palindrome.

Definition 1. We say that $S = p_1 p_2 \cdots p_\ell$ is a (maximal) palindromic decomposition of S if all the strings p_i are (maximal) palindromes.

Definition 2. A (maximal) palindromic decomposition of S such that the number of (maximal) palindromes is minimal is called a (maximal) palindromic factorization of S .

Note that any single letter is a palindrome and, hence, every string can always be decomposed into palindromes. However, not every string can be decomposed into maximal palindromes; e.g. consider $S = \text{abaca}$ [2].

Let f be an *involution* on the alphabet Σ , i.e., a function such that $f^2 = \text{id}$. We extend f into a morphism on strings over Σ . We say that a string x is a *generalized palindrome* if $x = f(x^R)$. Two known notions fit this definition:

- If $f = \text{id}$, then a generalized palindrome is a standard palindrome.
- If $\Sigma = \{A, C, G, T\}$ and $f(A) = T$, $f(C) = G$, $f(G) = C$, $f(T) = A$, then a generalized palindrome corresponds to a so-called complemented palindrome [15].

Example 1. The string $A G T A C T T C A T G A$ is a standard palindrome and the string $T A G T C G A C T A$ is a complemented palindrome.

We also consider (generalized) palindromes with errors. Let us recall two well-known metrics on strings. Let u and v be two strings. If $|u| = |v|$, then the *Hamming distance* between u and v is the number of positions where u and v do not match. The *edit (or Levenshtein) distance* between u and v is the minimum number of edit operations (insertions, deletions, substitutions) needed to transform u into v . We say that x is a *generalized δ -palindrome* under the Hamming distance (or the edit distance) if the minimum Hamming distance (edit distance, respectively) from x to any generalized palindrome is *at most* δ .

A generalized palindrome $S[i..j]$ is called *maximal* if there is no longer generalized palindrome with the same center. Similarly, a generalized δ -palindrome $S[i..j]$ under the Hamming/edit distance is called *maximal* if there is no longer generalized δ -palindrome under the same distance measure with the same center.

Example 2. All maximal 0-palindromes/1-palindromes in $GTATCG$ (for $f = \text{id}$) under the Hamming and under the edit distance are as follows:

center	1	1.5	2	2.5	3	3.5	4	4.5	5	5.5	6
0	G	ε	T	ε	TAT	ε	T	ε	C	ε	G
1 under Hamming	G	GT	GTA	TA	GTATC	AT	ATC	TC	TCG	CG	G
1 under edit	G	GT	GTA	GTAT	GTATC	GTATCG	ATC	TC	TCG	CG	G

For instance, the whole string $GTATCG$ is a 1-palindrome under the edit distance, as deleting its fifth letter yields a palindrome $GTATG$.

The computational problems we study can be formally stated as follows.

GENERALIZED PALINDROMIC DECOMPOSITION WITH GAPS

Input: A string S of length n , an involution f , and integers $g, m \geq 1$

Output: A decomposition of S into generalized palindromes with the minimal possible total length of gaps, $\sum_i^q |g_i|$, such that:

- There are at most g gaps, i.e. $q \leq g$
- Each palindrome is of length at least m

GENERALIZED MAXIMAL δ -PALINDROMIC DECOMPOSITION WITH GAPS

Input: A string S of length n , an involution f , and integers $g, m, \delta \geq 1$

Output: A decomposition of S into maximal generalized δ -palindromes with the minimal possible total length of gaps, $\sum_i^q |g_i|$, such that:

- There are at most g gaps, i.e. $q \leq g$
- Each generalized δ -palindrome is of length at least m

We apply several instances of dynamic programming. For simplicity of presentation, we only show how to compute the minimal total length of gaps and omit describing the retrieval of the decomposition itself. To compute the latter, in each of the dynamic programming matrices we would store a pointer to the cell that gave us the minimum value so that we could actually compute the decomposition with the minimal total length of the gaps by backtracing.

3 Palindromic decomposition with gaps

In this section we develop an efficient solution to the GENERALIZED PALINDROMIC DECOMPOSITION WITH GAPS problem. It is based on several transformations of the algorithm for computing a palindromic factorization by Fici et al. [9]. For a string S of length n this algorithm works in $\mathcal{O}(n \log n)$ time. The algorithm consists of two steps:

1. Let P_j be the sorted list of starting positions of all palindromes ending at position j in S . This list may have size $\mathcal{O}(j)$. However, it follows from combinatorial properties of palindromes that the sequence of consecutive differences in P_j is non-increasing and contains at most $\mathcal{O}(\log j)$ distinct values. Let $P_{j,\Delta}$ be the maximal sublist of P_j containing elements whose predecessor in P_j is smaller by exactly Δ . Then there are $\mathcal{O}(\log j)$ such sublists in P_j . Hence, P_j can be represented by a set G_j of size $\mathcal{O}(\log j)$ which consists of triples of the form (i, Δ, k) that represent $P_{j,\Delta} = \{i, i + \Delta, \dots, i + (k-1)\Delta\}$. The triples are sorted according to decreasing values of Δ and all starting positions in each triple are greater than in the previous one. Fici et al. show that G_j can be computed from G_{j-1} in $\mathcal{O}(\log j)$ time.
2. Let $PL[j]$ denote the number of palindromes in a palindromic factorization of $S[1..j]$. Fici et al. show that it can be computed via a dynamic programming approach, using all palindromes from G_j in $\mathcal{O}(\log j)$ time. Their algorithm works as follows. Let $PL_\Delta[j]$ be the minimum number of palindromes we can decompose $S[1..j]$ in, provided that we use a palindrome from $(i, \Delta, k) \in G_j$. Then $PL_\Delta[j]$ can be computed in constant time using $PL_\Delta[j - \Delta]$ based on the fact that if $(i, \Delta, k) \in G_j$ and $k \geq 2$, then $(i, \Delta, k-1) \in G_{j-\Delta}$. Exploiting this fact, $PL_\Delta[j]$ can be computed by only considering $PL_\Delta[j - \Delta]$ and the shortest palindrome in (i, Δ, k) . Finally, we compute $PL[j]$ from all such $PL_\Delta[j]$ values.

In Appendix A we show for completeness that the same approach works for generalized palindromes for any involution f .

To solve the GENERALIZED PALINDROMIC DECOMPOSITION WITH GAPS problem, we first need to modify each of the triples in G_j to reflect the length constraint (m). More precisely, due to the length constraint, in each G_j some triples will disappear completely, and at most one triple will get *trimmed* (i.e. the parameter k will be decreased).

Our algorithm then computes an array $MG[1..n][0..g]$ such that $MG[j][q]$ is the minimum possible total length of gaps in a palindromic decomposition of $S[1..j]$, provided that there are at most q gaps. Simultaneously, our algorithm computes an auxiliary array $MG'[1..n][0..g]$ such that $MG'[j][q]$ is the minimum possible total length of gaps up to position j provided that this position belongs to a gap: at most the q -th one.

For $j > 0$ and $q \geq 0$ we have the following formula:

$$MG[j][q] = \min(MG'[j][q], \min_{\Delta} \{MG_\Delta[j][q]\})$$

where $MG_{\Delta}[j][q]$ is the partial minimum computed only using generalized palindromes from $(i, \Delta, k) \in G_j$. The formula means: either we have a gap at position j , or we use a generalized palindrome ending at position j . We also set $MG[0][q] = 0$ for any $q \geq 0$.

We compute $MG_{\Delta}[j][q]$ for $(i, \Delta, k) \in G_j$ using the same approach as Fici et al. [9] used for PL_{Δ} , ignoring the triples that disappear due to the length constraint. If there is a triple that got trimmed, then the corresponding triple at position $j - \Delta$ (from which we reuse the values in the dynamic programming) must have got trimmed as well. More precisely, if the triple (i, Δ, k) is trimmed to (i, Δ, k') at position j , then at position $j - \Delta$ there is a triple $(i, \Delta, k - 1)$ which is trimmed to $(i, \Delta, k' - 1)$; that is, by the same number of generalized palindromes. Consequently, to compute $MG_{\Delta}[j][q]$ from $MG_{\Delta}[j - \Delta][q]$, we need to include one additional generalized palindrome (the shortest one in the triple) just as in Fici et al.'s approach.

Example 3. Consider the string AACCAACCAACCAACCA, $f = \text{id}$, and let $m = 7$.

A	A	C	C	A	A	C	C	A	A	C	C	A	A	C	C	A	A
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18

Then $G_{18} = \{(1, \infty, 1), (5, 4, 4), (18, 1, 1)\}$, where:

- $(1, \infty, 1)$ represents the whole string,
- $(5, 4, 4)$ represents $\{\text{AACCAACCAACCA}, \text{AACCAACCA}, \text{AACCA}, \text{AA}\}$ which will get trimmed by 2 palindromes due to the length constraint, becoming $(5, 4, 2)$,
- $(18, 1, 1)$ represents $\{\text{A}\}$ and disappears.

Now looking at position $j - \Delta = 18 - 4 = 14$ for the trimmed group, we had $(5, 4, 3) \in G_{14}$ representing $\{\text{AACCAACCA}, \text{AACCA}, \text{AA}\}$, and this also gets trimmed by 2 palindromes, becoming $(5, 4, 1)$.

Finally, for $j > 0$ and $q > 0$ we compute MG' using the following formula:

$$MG'[j][q] = \min(MG'[j - 1][q], MG[j - 1][q - 1]) + 1.$$

The first case corresponds to continuing the gap from position j , whereas the second to using a generalized palindrome finishing at position $j - 1$ or a gap finishing at position $j - 1$ (the latter will be suboptimal). Here the border cases are $MG'[j][0] = \infty$ for $j \geq 0$ and $MG'[0][q] = \infty$ for $q > 0$.

Thus we arrive at the complete solution to the problem.

Theorem 1. *The GENERALIZED PALINDROMIC DECOMPOSITION WITH GAPS problem can be solved in $\mathcal{O}(n \log n \cdot g)$ time and $\mathcal{O}(n \cdot g)$ space.*

4 Computing maximal palindromes with errors

Recall that all maximal (standard) palindromes in a string can be computed in $\mathcal{O}(n)$ time by Manacher's [6, 21] and Gusfield's [15] algorithms. These algorithms perform different computations for odd- and for even-length palindromes. Recall that we defined the centers of odd-length palindromes as integers and the centers of even-length palindromes as odd multiples of $\frac{1}{2}$.

Gusfield's algorithm [15] applies Longest Common Extension (LCE) Queries in the string $T = S\$S^R$, where $\$ \notin \Sigma$ is a sentinel character. An $LCE(i, j)$ query returns the length of the longest common prefix of the suffixes $T[i..|T|]$ and $T[j..|T|]$. For example, to compute the length of the maximal even-length palindrome centered between positions i and $i + 1$, the algorithm computes $LCE(i + 1, 2n + 2 - i)$ in T . Recall that LCE queries in a string (over an integer alphabet) can be answered in $\mathcal{O}(1)$ time after linear-time preprocessing [5].

Gusfield's approach can be easily adapted to generalized palindromes: it suffices to apply LCE-queries on $T = S\$f(S^R)$. To further simplify the description of this approach, we introduce the Longest Gapped

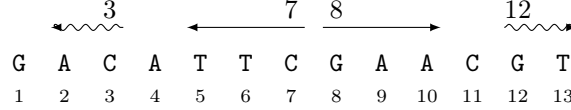


Figure 1: To find the longest complemented 1-palindrome under the Hamming distance centered at position 7.5 in $S = \text{GACATTCGAACGT}$, it suffices to ask two LGPal-queries: $\text{LGP}al(7, 8) = 3$ finds the first mismatch, and $\text{LGP}al(3, 12)$ extends the 1-palindrome after the mismatch. Note that each of these LGPal-queries is equivalent to an appropriate LCP-query in $S\$f(S^R)$.

Palindrome (LGP)al Queries, such that $\text{LGP}al(i, j)$ is the maximum k such that $f(S[i - k + 1 \dots i]^R) = S[j \dots j + k - 1]$; see Fig. 1. As we have already noticed, LGPal-queries are equivalent to LCE-queries in $T = S\$f(S^R)$.

It is known (see [15]) that all maximal generalized δ -palindromes under the Hamming distance can be computed in $O(n \cdot \delta)$ time via at most δ applications of the LGPal-query for each possible center position. Below we show how to compute maximal generalized δ -palindromes under the edit distance within the same time complexity.

Recall that if u is a generalized δ -palindrome under the edit distance, then there exists a generalized palindrome v such that the minimal number of edit operations (insertion, deletion, substitution) required to transform u to v is at most δ . The following simple observation shows that we can restrict ourselves to deletions and substitutions only, which we call in what follows the *restricted edit operations*. Intuitively, instead of inserting at position i a character to match the character at position $|u| - i + 1$, we can delete the character at position $|u| - i + 1$.

Observation 1. *Let u be a generalized δ -palindrome and v a generalized palindrome such that the edit distance between u and v is minimal. Then there exists a generalized palindrome v' such that the number of restricted edit operations needed to transform u to v' is equal to the edit distance between u and v .*

We can extend a maximal generalized δ -palindrome $S[i \dots j]$ to a maximal generalized $(\delta + 1)$ -palindrome in three ways; either ignore the letter $S[i - 1]$ and then perform an LGPal-query, or ignore the letter $S[j + 1]$ and then perform an LGPal-query, or ignore both and then perform the LGPal-query. More formally:

Definition 3. Assume that $S[i \dots j]$ is a generalized δ -palindrome. Then we say that each of the factors $S[i' \dots j']$ for:

- $i' = i - 1 - d, j' = j + d$, where $d = \text{LGP}al(i - 2, j + 1)$
- $i' = i - d, j' = j + 1 + d$, where $d = \text{LGP}al(i - 1, j + 2)$
- $i' = i - 1 - d, j' = j + 1 + d$, where $d = \text{LGP}al(i - 2, j + 2)$

is an *extension* of $S[i \dots j]$. If the index i' is smaller than 1 or the index j' is greater than $|S|$, the corresponding extension is not possible. We also say that $S[i \dots j]$ can be extended to any of the three strings $S[i' \dots j']$.

Clearly, the extensions of a generalized δ -palindrome are always generalized $(\delta + 1)$ -palindromes.

To facilitate the case of δ -palindromes being prefixes or suffixes of the text, we also introduce the following *border-reductions* for $S[i \dots j]$ being a generalized δ -palindrome:

- If $i = 1$, a border reduction leads to $S[1 \dots j - 1]$.
- If $j = n$, a border reduction leads to $S[i + 1 \dots n]$.

If any of the reductions is possible, we also say that $S[i \dots j]$ can be border-reduced to the corresponding strings. As previously, border-reductions of a generalized δ -palindrome are always generalized $(\delta + 1)$ -palindromes.



Figure 2: If the outermost restricted edit operation on $S[i' \dots j']$ is a substitution (from letter X to letter Y), then $S[i' \dots j']$ is an extension of the third type of the maximal generalized $(\delta - 1)$ -palindrome $S[i \dots j]$.

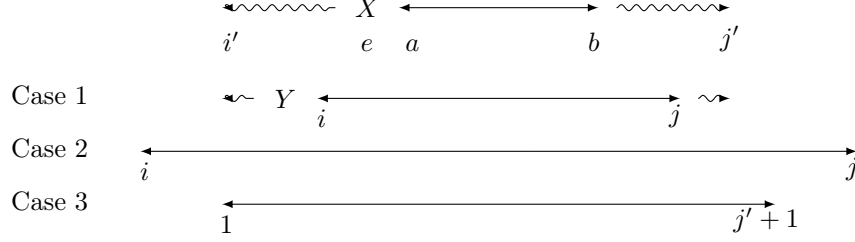


Figure 3: Three cases resulting when the outermost edit operation on $S[i' \dots j']$ is a deletion of a character X .

Lemma 1. *Given a maximal generalized δ -palindrome $S[i' \dots j']$ with $\delta > 0$, there exists a maximal generalized $(\delta - 1)$ -palindrome $S[i \dots j]$ which can be extended or border-reduced to $S[i' \dots j']$.*

Proof. Consider a shortest sequence of restricted edit operations that transforms $u = S[i' \dots j']$ into a generalized palindrome v . Let us consider the position where we perform a restricted edit operation that is closest to i' or j' . Assume w.l.o.g. that this position—denote it by e —is not further to i' than to j' .

Assume first that this edit operation is a substitution. Then $S[i \dots j]$, for $i = e + 1$ and $j = j' - (e + 1 - i')$, is a generalized $(\delta - 1)$ -palindrome (the witness generalized palindrome is the corresponding factor of v); see Fig. 2. Moreover, it is a maximal generalized $(\delta - 1)$ -palindrome, as otherwise $S[e] = S[i - 1]$ would be equal to $f(S[j + 1])$, which means that the substitution at the position e would not be necessary. This completes the proof in this case.

Now assume that the edit operation at the position e was a deletion. Let $a = e + 1$ and $b = j' - (e - i')$. Again, we see that clearly $S[a \dots b]$ is a generalized $(\delta - 1)$ -palindrome. If it is maximal, then we are done. Otherwise, consider the maximal generalized $(\delta - 1)$ -palindrome $S[i \dots j]$ centered at the same position as $S[a \dots b]$ ($a - i = j - b > 0$). Now we have three cases; see Fig. 3.

1. If $j \leq j'$, then we can obtain $S[i' \dots j']$ by an extension (of the first type) of $S[i \dots j]$; i.e. ignoring the letter $S[i - 1]$.
2. If $j > j'$, then we have that $S[i' \dots j' + 1]$ is a generalized $(\delta - 1)$ -palindrome. If, additionally, $i' > 1$, then $S[i' - 1 \dots j' + 1]$ is a generalized δ -palindrome, which contradicts the maximality of $S[i' \dots j']$.
3. Finally, if $j > j'$ and $i' = 1$, then $i = 1$, $j = j' + 1$. Hence, $S[i' \dots j']$ obtained from $S[i \dots j]$ by a border-reduction.

This completes the proof of the lemma. \square

The combinatorial characterization of Lemma 1 yields an algorithm for generating all maximal generalized d -palindromes, for all centers and subsequent $d = 0, \dots, \delta$. Maximal generalized 0-palindromes are computed using Gusfield's approach (LGPal-queries). For a given $d < \delta$, we consider all the maximal generalized d -palindromes and try to extend each of them in all three possible ways (and border-reduce, if possible). This way we obtain a number of generalized $(d + 1)$ -palindromes amongst which, by Lemma 1, are all maximal generalized $(d + 1)$ -palindromes. To exclude the non-maximal ones, we group the generalized $(d + 1)$ -palindromes by their centers (in $\mathcal{O}(n)$ time via bucket sort) and retain only the longest one for each center. We arrive at the following intermediate result.

Lemma 2. *Under the edit distance, all maximal generalized δ -palindromes in a string of length n can be computed in $\mathcal{O}(n \cdot \delta)$ time and $\mathcal{O}(n)$ space.*

5 Maximal palindromic decomposition with gaps and errors

Let \mathcal{F} be a set of factors of the text $S[1..n]$. In this section we develop a general framework that allows to decompose S into factors from \mathcal{F} , allowing at most g gaps. We call such a factorization a (g, \mathcal{F}) -factorization of S . Our goal is to find a (g, \mathcal{F}) -factorization of S that minimizes the total length of gaps. We aim at the time complexity $\mathcal{O}((n + |\mathcal{F}|) \cdot g)$ and space complexity $\mathcal{O}(n \cdot g + |\mathcal{F}|)$.

In our solution we use dynamic programming to compute two arrays, similar to the ones used in Section 3:

$MG[1..n][0..g]$: $MG[j][q]$ is the minimum total length of gaps in a (q, \mathcal{F}) -factorization of $S[1..j]$.

$MG'[1..n][0..g]$: $MG'[j][q]$ is the minimum total length of gaps in a (q, \mathcal{F}) -factorization of $S[1..j]$ for which the position j belongs to a gap.

We use the following formulas, for $j > 0$ and $q > 0$:

$$MG[j][q] = \min(MG'[j][q], \min_{S[a..j] \in \mathcal{F}} MG[a-1][q])$$

$$MG'[j][q] = \min(MG[j-1][q-1], MG'[j-1][q]) + 1$$

The border cases are exactly the same as in Section 3.

Clearly, the space complexity of this solution is $\mathcal{O}(n \cdot g + |\mathcal{F}|)$. Let us analyze its time complexity. Fix $q \in \{0, \dots, g\}$. The number of transitions using the factors from \mathcal{F} in the dynamic programming is $|\mathcal{F}|$ in total, as each factor is used only for the position j where it ends. Hence, the formulas for $MG[j][q]$ take $\mathcal{O}(n \cdot g + |\mathcal{F}| \cdot g)$ time to evaluate. Computing the $MG'[j][q]$ values takes $\mathcal{O}(n \cdot g)$ time. Thus we arrive at the desired time complexity of $\mathcal{O}((n + |\mathcal{F}|) \cdot g)$.

We apply this approach to maximal generalized δ -palindromes in each of the considered metrics (see the classic result from [15] for the Hamming distance and Lemma 2 for the edit distance) to obtain the following result.

Theorem 2. *The GENERALIZED MAXIMAL δ -PALINDROMIC DECOMPOSITION WITH GAPS problem under the Hamming distance or the edit distance can be solved in $\mathcal{O}(n \cdot (g + \delta))$ time and $\mathcal{O}(n \cdot g)$ space.*

Example 4. Consider the following string¹ of length 92:

GGACTCGGCTTGCTGAGGTGCACACAGCAAGAGGCGAGAGCGGCGACTGGTGAGTACGCCAAATTTTACTAGCGGAGGCTAGAAGGAGAGA

We have used our implementation of the algorithm from Theorem 2 to compute the decomposition of the string into maximal complemented 3-palindromes of length at least 14 under the **edit distance** with at most 4 gaps ($g = 4$, $\delta = 3$, $m = 14$) with the minimal total gap length:

[GGACTCG] GCTTGCTGAGGTGCACACAGCAAGA [GGCGAGAGC] GGCGACTGGTGAGTACGCC [AAATTTTG]
ACTAGCGGAGGCTAGA [AGGAGAGA]

The gaps are given in square brackets. Edit operations are underlined, with deletes additionally given in italics. The gaps have total length 32.

In comparison, the optimal decomposition of this string under the **Hamming distance** with the same parameters ($g = 4$, $\delta = 3$, $m = 14$) uses four gaps of total length 46.

¹See <http://www.cesshiv1.org/disview.php?accession=AB220944>

6 Conclusions

We have presented two algorithms for finding palindromic decompositions: one allowing gaps and the other allowing both gaps in the decomposition and errors in palindromes. The first algorithm shows that (somewhat surprisingly) Fici et al.’s algorithm [9] for finding an exact palindromic factorization can be extended to handle gaps, a constraint on the palindromes length, and complements in palindromes as well. In the second algorithm we decompose a string into maximal palindromes with errors; the most involved part here was computing all such maximal palindromes under the edit distance.

In the problems that were defined in the beginning, the objective was to minimize the total length of gaps, allowing a certain number of gaps. However, the approaches that were presented in this paper can be used to solve different variants of the problems, like minimizing only the total number of gaps or maximizing the total length of palindromes, regardless of the number of gaps.

An open question is to efficiently compute decompositions into palindromes that may contain errors and are not necessarily maximal. This problem seems to be hard, as δ -palindromes do not have such a strong combinatorial structure as palindromes without errors.

Acknowledgments

We would like to warmly thank G. Fici, T. Gagie, J. Kärkkäinen, and D. Kempa for providing us Figures 4, 5, and 6 (that we adapted and present in Appendix A) that they first presented in [9].

References

- [1] Michał Adamczyk, Mai Alzamel, Panagiotis Charalampopoulos, Costas S. Iliopoulos, and Jakub Radoszewski. Palindromic decompositions with gaps and errors. In *Computer Science - Theory and Applications - 12th International Computer Science Symposium in Russia, CSR 2017, Proceedings (accepted for publication)*, 2017.
- [2] Ali Alatabbi, Costas S Iliopoulos, and Mohammad Sohel Rahman. Maximal palindromic factorization. In *Stringology*, pages 70–77, 2013.
- [3] Alberto Apostolico, Dany Breslauer, and Zvi Galil. Parallel detection of all palindromes in a string. *Theor. Comput. Sci.*, 141(1):163–173, 1995.
- [4] Dany Breslauer and Zvi Galil. Finding all periods and initial palindromes of a string in parallel. *Algorithmica*, 14(4):355–366, 1995.
- [5] Maxime Crochemore, Christophe Hancart, and Thierry Lecroq. *Algorithms on Strings*. Cambridge University Press, 2007.
- [6] Maxime Crochemore and Wojciech Rytter. *Jewels of Stringology*. World Scientific, 2003.
- [7] Xavier Droubay. Palindromes in the Fibonacci word. *Inf. Process. Lett.*, 55(4):217–221, 1995.
- [8] Xavier Droubay and Giuseppe Pirillo. Palindromes and Sturmian words. *Theor. Comput. Sci.*, 223(1-2):73–85, 1999.
- [9] Gabriele Fici, Travis Gagie, Juha Kärkkäinen, and Dominik Kempa. A subquadratic algorithm for minimum palindromic factorization. *J. Discr. Alg.*, 28(C):41–48, September 2014.
- [10] A.E. Frid, S. Puzynina, and L.Q. Zamboni. On palindromic factorization of words. *Adv. Appl. Math.*, 50(5):737–748, May 2013.

- [11] Yuta Fujishige, Michitaro Nakamura, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. Finding gapped palindromes online. In Veli Mäkinen, Simon J. Puglisi, and Leena Salmela, editors, *Combinatorial Algorithms - 27th International Workshop, IWOCA 2016, Proceedings*, volume 9843 of *LNCS*, pages 191–202. Springer, 2016.
- [12] Zvi Galil. Real-time algorithms for string-matching and palindrome recognition. In *Proceedings of the Eighth Annual ACM Symposium on Theory of Computing*, pages 161–173. ACM, 1976.
- [13] Zvi Galil and Joel Seiferas. A linear-time on-line recognition algorithm for “palstar”. *J. ACM*, 25(1):102–111, January 1978.
- [14] Shivika Gupta, Rajesh Prasad, and Sunita Yadav. Searching gapped palindromes in dna sequences using dynamic suffix array. *Indian Journal of Science and Technology*, 8(23):1, 2015.
- [15] Dan Gusfield. *Algorithms on Strings, Trees, and Sequences: Computer Science and Computational Biology*. Cambridge University Press, New York, NY, USA, 1997.
- [16] Ping-Hui Hsu, Kuan-Yu Chen, and Kun-Mao Chao. Finding all approximate gapped palindromes. *Int. J. Found. Comput. Sci.*, 21(6):925–939, 2010.
- [17] Tomohiro I, Shiho Sugimoto, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. Computing palindromic factorizations and palindromic covers on-line. In Alexander S. Kulikov, Sergei O. Kuznetsov, and Pavel A. Pevzner, editors, *Combinatorial Pattern Matching - 25th Annual Symposium, CPM 2014. Proceedings*, volume 8486 of *Lecture Notes in Computer Science*, pages 150–161. Springer, 2014.
- [18] Donald E Knuth, James H Morris, Jr, and Vaughan R Pratt. Fast pattern matching in strings. *SIAM journal on computing*, 6(2):323–350, 1977.
- [19] Roman Kolpakov and Gregory Kucherov. Searching for gapped palindromes. *Theor. Comput. Sci.*, 410(51):5365–5373, 2009.
- [20] Dmitry Kosolobov, Mikhail Rubinchik, and Arseny M. Shur. Pal^k is linear recognizable online. In *Proc. 41th Int. Conf. on Theory and Practice of Computer Science (SOFSEM 2015)*, pages 289–301, 2015.
- [21] Glenn Manacher. A new linear-time “on-line” algorithm for finding the smallest initial palindrome of a string. *Journal of the ACM (JACM)*, 22(3):346–351, 1975.
- [22] Alexandre H. L. Porto and Valmir Carneiro Barbosa. Finding approximate palindromes in strings. *Pattern Recognition*, 35(11):2581–2591, 2002.
- [23] Mikhail Rubinchik and Arseny M. Shur. EERTREE: an efficient data structure for processing palindromes in strings. In *Combinatorial Algorithms - 26th International Workshop, IWOCA 2015*, LNCS, pages 321–333, 2015.

A Generalized palindromic factorization

In this section we show that the approach of Fici et al. [9] works for generalized palindromes for any involution f . The following auxiliary lemma extends the combinatorial properties of standard palindromes used in [9] (see Lemmas 1-3 therein) to generalized palindromes. Recall that a string y is called a *border* of a string x if it is both a prefix and a suffix of x . A number p is called a *period* of x if $x[i] = x[i+p]$ for all $i = 1, \dots, |x| - p$. It is well known that x has a period p iff it has a border of length $|x| - p$; see [5, 6].

Lemma 3. (a) *Let y be a suffix of a generalized palindrome x . Then y is a border of x iff y is a generalized palindrome.*

(b) Let x be a string with a border y such that $|x| \leq 2|y|$. Then x is a generalized palindrome iff y is a generalized palindrome.

(c) Let y be a proper suffix of a generalized palindrome x . Then $|x| - |y|$ is a period of x iff y is a generalized palindrome. In particular, $|x| - |y|$ is the smallest period of x iff y is the longest generalized palindromic proper suffix of x .

Proof. (a) Let y' be the prefix of x of length $|y|$. As x is a generalized palindrome, $y' = f(y^R)$. (\Rightarrow) If y is a border of x , then $y = y' = f(y^R)$, so y is a generalized palindrome. (\Leftarrow) If y is a generalized palindrome, then $y' = f(y^R) = y$, so y is a border of x .

(b) (\Rightarrow) From (a), if x is a generalized palindrome and y is its border, then y is a generalized palindrome. (\Leftarrow) If y is a generalized palindrome, $f(x^R)$ has a border $f(y^R) = y$. This border covers the whole string $f(x^R)$ and is the same as the border of x , so $x = f(x^R)$ and x indeed is a generalized palindrome.

(c) This is a consequence of part (a) and the relation between borders and periods of a string. \square

=The crucial combinatorial property of standard palindromes used in Step 1 of the algorithm in Section 3 is that the sequence of consecutive differences in P_j is non-increasing and contains at most $\mathcal{O}(\log j)$ distinct values. We show that the same observation holds for generalized palindromes; this follows from the next lemma, parts (1) and (2). The proof of Lemma 4 follows exactly the lines of the proof of the corresponding Lemma 4 in [9]; Figures 4 and 5 (thanks to the courtesy of the authors of [9]) are included for illustration.

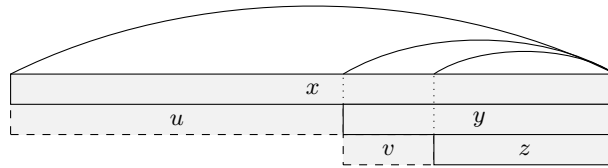


Figure 4: Setting in Lemma 4.

Lemma 4. Let x be a generalized palindrome, y the longest generalized palindromic proper suffix of x , and z the longest generalized palindromic proper suffix of y . Let u and v be strings such that $x = uy$ and $y = vz$. Then:

- (1) $|u| \geq |v|$;
- (2) if $|u| > |v|$ then $|u| > |z|$;
- (3) if $|u| = |v|$ then $u = v$.

Proof. See Figure 4 for an illustration.

(1) By Lemma 3(c), $|u| = |x| - |y|$ is the smallest period of x , and $|v| = |y| - |z|$ is the smallest period of y . Since y is a factor of x , either $|u| > |y| > |v|$ or $|u|$ is a period of y too, and thus it cannot be smaller than $|v|$.

(2) By Lemma 3(a), y is a border of x and thus v is a prefix of x . Let w be a string such that $x = vw$. Then z is a border of w and $|w| = |zu|$ (see Figure 5). Since we assume $|u| > |v|$, we must have $|w| > |y|$. Suppose to the contrary that $|u| \leq |z|$. Then $|w| = |zu| \leq 2|z|$, and by Lemma 3(b), w is a generalized palindrome. But this contradicts y being the longest generalized palindromic proper suffix of x .

(3) In the proof of (2) we saw that v is a prefix of x , and so is u by definition. Thus $u = v$ if $|u| = |v|$. \square

We have thus shown that, also in case of generalized palindromes, the set P_j can be compactly represented by a set G_j , as described in Section 3. To complete Step 1 of the algorithm, we need to show that G_j can be computed from G_{j-1} in $\mathcal{O}(\log j)$ time. For this, just as in [9], we show that each triple $(i, \Delta, k) \in G_{j-1}$ will be either eliminated or replaced by $(i - 1, \Delta, k)$ in G_j . The proof exploits part (3) of Lemma 4.

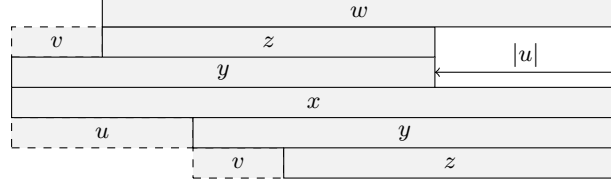


Figure 5: Proof of Lemma 4(2): if $|u| > |v|$ and $|u| \leq |z|$ then w is a generalized palindromic proper suffix of x longer than y .

Lemma 5. *Let p_i and p_{i+1} be two consecutive elements of $P_{j-1, \Delta}$. Then $p_i - 1 \in P_j$ iff $p_{i+1} - 1 \in P_j$.*

Proof. By definition, $p_{i+1} - p_i = \Delta$, and the predecessor of p_i in P_j is $p_{i-1} = p_i - \Delta$. The strings $x = S[p_{i-1} \dots j - 1]$, $y = S[p_i \dots j - 1]$, and $z = S[p_{i+1} \dots j - 1]$ form the situation of Lemma 4(3). Hence, $S[p_i - 1] = S[p_{i+1} - 1] = c$. Thus, $p_i - 1 \in P_j$ iff $S[j] = f(c)$ iff $p_{i+1} - 1 \in P_j$. \square

After this transformation, one might need to update pairs of adjacent triples in G_j because the gaps between them might have changed. This simple process is explained in detail in [9] and takes only $\mathcal{O}(\log j)$ additional time.

As for Step 2 of the algorithm, it suffices to show that the following combinatorial observation holds for generalized palindromes. Again we follow the lines of the proof from [9]; Figure 6 (thanks to the courtesy of the authors of [9]) is included for illustration.

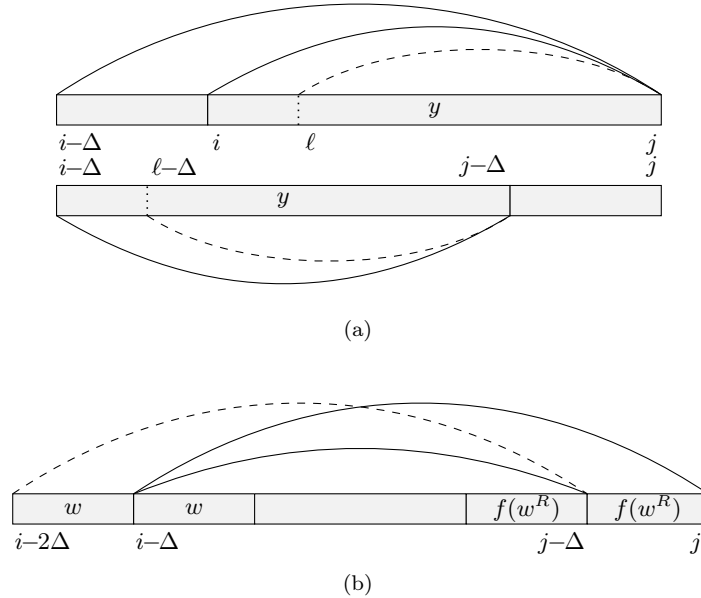


Figure 6: Proof of Lemma 6. (a) $\ell \in P_j$ iff $\ell - \Delta \in P_{j-\Delta}$ for all $\ell \in [i \dots j]$. (b) If $i - 2\Delta \in P_{j-\Delta}$ then $S[i - 2\Delta \dots j]$ is a generalized palindrome.

Lemma 6. *If $(i, \Delta, k) \in G_j$ and $k \geq 2$, then $(i, \Delta, k - 1) \in G_{j-\Delta}$.*

Proof. By definition, $(i, \Delta, k) \in G_j$ is equivalent to saying that $P_{j, \Delta} = \{i, i + \Delta, \dots, i + (k - 1)\Delta\}$, and we need to show that $P_{j-\Delta, \Delta} = \{i, i + \Delta, \dots, i + (k - 2)\Delta\}$. We will show first that $P_{j-\Delta, \Delta} \cap [i - \Delta + 1 \dots j - \Delta] = \{i, i + \Delta, \dots, i + (k - 2)\Delta\}$ and then that $P_{j-\Delta, \Delta} \cap [1 \dots i - \Delta] = \emptyset$.

Since $y = S[i \dots j]$ and $x = S[i - \Delta \dots j]$ are generalized palindromes and y is the longest proper border of x (by Lemma 3(a)), $S[i - \Delta \dots j - \Delta] = y = S[i \dots j]$. Thus for all $\ell \in [i \dots j]$, $\ell \in P_j$ iff $\ell - \Delta \in P_{j-\Delta}$ (see Figure 6a). In particular, the consecutive differences in both cases are the same and for all $\ell \in [i + 1 \dots j]$, $\ell \in P_{j,\Delta}$ iff $\ell - \Delta \in P_{j-\Delta,\Delta}$. Thus $P_{j-\Delta,\Delta} \cap [i - \Delta + 1 \dots j - \Delta] = \{i, i + \Delta, \dots, i + (k - 2)\Delta\}$.

We still need to show that $P_{j-\Delta,\Delta} \cap [1 \dots i - \Delta] = \emptyset$, which is true if and only if $i - 2\Delta \notin P_{j-\Delta}$. Suppose to the contrary that $S[i - 2\Delta \dots j - \Delta]$ is a generalized palindrome and let $w = S[i - 2\Delta \dots i - \Delta - 1]$. Then $S[j - 2\Delta + 1 \dots j - \Delta] = f(w^R)$. Since $z = S[i - \Delta \dots j - \Delta]$ and $S[i - \Delta \dots j]$ are generalized palindromes too, we have that $S[i - \Delta \dots i - 1] = w$ and $S[j - \Delta + 1 \dots j] = f(w^R)$. Finally, since z is a generalized palindrome, $S[i - 2\Delta \dots j] = wzf(w^R)$ is a generalized palindrome (see Figure 6b). This implies that $i - 2\Delta \in P_j$ and thus $i - \Delta \in P_{j,\Delta}$, which is a contradiction. \square