

# Approximate String Matching with SIMD

FERNANDO J. FIORI<sup>1</sup>, WALTTERI PAKALÉN<sup>2</sup> AND JORMA TARHIO<sup>2</sup>

<sup>1</sup>*Departamento de Ciencias de la Computación, Facultad de Ciencias Exactas, Ingeniería y Agrimensura, Universidad Nacional de Rosario, Pellegrini 250, S2000BTP Rosario, Argentina*

<sup>2</sup>*Department of Computer Science, Aalto University, P.O.B. 15400, FI-00076 Aalto, Finland*

*Corresponding author: Mr Fernando Fiori. E-mail: fioriff@gmail.com*

**We consider the  $k$  mismatches version of approximate string matching for a single pattern and multiple patterns. For these problems, we present new algorithms utilizing the single instruction multiple data (SIMD) instruction set extensions for patterns of up to 32 characters. We apply SIMD computation in three ways: in counting of mismatches, in comparison of substrings and in calculation of fingerprints. We show the competitiveness of the new algorithms by practical experiments.**

*Keywords:: approximate string;  $k$  mismatches; Hamming distance; SIMD.*

*Received 22 February 2020; Revised 30 August 2020; Accepted 24 September 2020*

*Handling editor: Professor Prudence Wong*

This is an extended version of a conference paper presented in The Prague Stringology Conference, 2017.

## 1.. INTRODUCTION

The *string matching problem* is defined as follows: given a pattern  $P = p_0 \cdots p_{m-1}$  and a text  $T = t_0 \cdots t_{n-1}$  over an alphabet  $\Sigma$  of size  $\sigma$ , find all the occurrences of  $P$  in  $T$ . In this paper, we consider the  $k$  mismatches variation of the problem where  $P'$  is an occurrence of  $P$ , if  $|P'| = |P|$  holds and  $P'$  has at most  $k$  mismatches with  $P$ . The mismatch distance of two strings of equal length is also called the *Hamming distance*.

There are numerous good solutions for the  $k$  mismatches problem, see e.g. Navarro's survey [1]. In this article, we introduce new algorithms for the problem. Besides the single pattern problem, we also consider the multiple pattern variation where there are  $r$  patterns  $P^0, \dots, P^{r-1}$  to search where  $P^i = p_0^i \cdots p_{m_i-1}^i$ . Our solutions utilize single instruction multiple data (SIMD) instruction set extensions [2]. SIMD is a type of parallel architecture that allows one instruction to be operated on multiple data items at the same time. We apply SIMD computation in three ways: in counting of mismatches, in comparison of substrings and in calculation of fingerprints a.k.a. hash values. Although our emphasis is on the practical efficiency of the algorithms, we analyze the time complexities of our algorithms. We show the competitiveness of the new algorithms by practical experiments. Our new algorithms for the single pattern problem are faster than reference methods in most cases tested, and our multiple pattern algorithm

outperforms Fredriksson and Navarro's algorithm [3] with a wide margin.

Algorithms for the  $k$  mismatches problem have several applications such as virus and intrusion detection [4], spelling [5], speech recognition [6], optical character recognition [7], handwriting recognition [7], text retrieval under synonym or thesaurus expansion [8], matching of nucleotide sequences in metagenomics in computational biology [9, 10], among others. Moreover, some singlepattern approximate search algorithms resort to multipattern searching of pattern pieces [11]. Depending on the application,  $r$  may vary from a few to thousands.

We will refer to bitwise operations AND and left shift with the same symbols as in the programming language C: '&' and '<<', respectively.

The rest of the paper is organized as follows. Section 2 reviews earlier solutions. Section 3 introduces SIMD computation and the SIMD techniques applied. Sections 4 and 5 describe the new algorithms, Section 6 presents the results of practical experiments and Section 7 concludes the article.

## 2.. EARLIER SOLUTIONS

There are many algorithms for the string matching with  $k$  mismatches problem. Most of them solve the single pattern variation, whereas only a few exist for the multiple pattern counterpart. Naively, an algorithm for the single pattern variation can be extended to solve the multiple pattern variation by executing it separately for each pattern. In the following, we

will review earlier solutions to these variations. None of the earlier solutions introduced here applies SIMD computation.

### 2.1.. Single string matching with $k$ mismatches

A naive algorithm works as follows: starting at the first text position, it compares character by character with the pattern until it detects more than  $k$  mismatches or until it makes  $m$  comparisons. If there are not more than  $k$  mismatches in this text window of  $m$  characters, it reports an occurrence. Then it moves the window one position to the right, and repeats the comparison. This algorithm works in  $O(mn)$  time in the worst case and in  $O(kn)$  time on average if individual characters in  $P$  and  $T$  are chosen independently and uniformly from the alphabet  $\Sigma$ . We will make this same assumption in all subsequent considerations on average time complexity.

Baeza-Yates and Gonnet [12] presented Shift-Add (SA), the first bit-parallel algorithm for the  $k$  mismatches problem. Shift-Add works in linear time for short patterns  $m \leq w/\lceil \log_2(k+1) \rceil + 1$ , where  $w$  is the width of the computer word. Shift-Add is still competitive for short patterns and large  $k$  [13]. Āurian *et al.* [14] presented two variations of SA, TuSA and TwSA, which process the alignment window backwards.

Approximate Boyer–Moore (ABM) by Tarhio and Ukkonen [15] is a generalization of the Boyer–Moore–Horspool algorithm [16] to approximate string matching. In ABM, shifting is based on a  $q$ -gram,  $q = k + 1$ . Liu *et al.* [17] tuned ABM for small alphabets. Their algorithm applies wider  $q$ -grams and is called FFAST (short for a Fast Algorithm for Approximate STring matching). Salmela *et al.* [18] designed an enhanced version of FFAST. This algorithm called EF has a faster preprocessing phase than FFAST and it counts a part of mismatches during preprocessing.

Approximate BNNDM (ABNNDM) by Navarro and Raffinot [19] is based on the BNNDM algorithm [19] for exact string matching. BNNDM simulates the suffix automaton of the reversed pattern with bit-parallelism. ABNNDM as well as ABM, FFAST, EF and TwSA achieve a sublinear running time on average in the case of favorable problem parameters.

The Baeza-Yates–Perleberg algorithm (BYP) [11] is based on a partitioning scheme, where the pattern is divided into substrings of length  $l = \lfloor m/(k+1) \rfloor$ . At least one of these subpatterns is exactly present in an approximate occurrence of the pattern. In the preprocessing phase, it splits the pattern into subpatterns of length  $l$ , and then it performs a multiple exact string matching search of these subpatterns. Whenever one of the subpatterns is found, it checks if there is an approximate pattern match with Ukkonen’s dynamic algorithm [20].

Besides practically oriented algorithms mentioned above, there are other solutions to the  $k$  mismatches problem: the kangaroo method [21, 22], the algorithms based on the fast Fourier transform and marking [23–25] and the  $O(nk^2 \log k/m + n \text{ polylog } m)$  solution presented by Clifford *et al.* [26], which is the best theoretical result.

### 2.2.. Multiple string matching with $k$ mismatches

A naive algorithm for multiple patterns is a natural extension of the naive algorithm for the single pattern case presented above. It performs the single pattern algorithm  $r$  times, which yields  $O(rmn)$  running time in the worst case and  $O(rkn)$  time on average.

The first non-naive algorithm for multiple approximate pattern matching was presented by Muth and Manber [27] for  $k = 1$ . In preprocessing, it uses an *Ad-hoc* function to obtain hash values of all the strings of length  $l - 1$  that result from taking each character out of each prefix of size  $l$  of every pattern, where  $l$  is chosen empirically depending on problem parameters. Then it proceeds in the same way for each text window of  $l$  characters, and naively verifies each hash coincidence. It detects insertions and deletions of characters as well as substitutions as errors, but it can easily be adapted to only consider mismatches by a small change in the way it checks each occurrence. The algorithm has an average preprocessing time complexity of  $O(lr)$  and an average search time complexity of  $O(ln)$ . Its average total time is then  $O(l(r+n))$ , which is rather independent of the number of patterns to search because  $n$  is usually much larger than  $r$ . As for the case when  $k > 1$ , it can be modified to compute all the strings that result from taking  $k$  characters out of each pattern prefix and text window. Given that prefixes and windows have length  $l$ , it needs to generate  $\frac{l!}{k!(l-k)!}$  strings for each text position and each pattern. Its complexity then limits its usage to only very small  $k$  in practice.

Later, Baeza-Yates and Navarro [28] designed an algorithm which is based on a similar partition scheme as in the BYP algorithm. This algorithm addresses the problem of searching occurrences of patterns in a text within an *edit distance* of at most  $k$ . They split every pattern into subpatterns of length  $l = \lfloor m/(k+1) \rfloor$  and perform an exact multiple pattern search using an extension of Sunday’s algorithm [29]. Whenever there is a subpattern occurrence, they check for the entire pattern with an approximate single pattern matching algorithm. This process is done by applying a *hierarchical piece verification* and employing a Non-deterministic Finite Automaton which exploits bit-parallelism.

The fastest algorithm to date is Fredriksson and Navarro’s algorithm [3], which is optimal on average. It places a window over the text, in which  $q$ -grams are read in a backwards order. Whenever an occurrence is impossible, the window is shifted past the read  $q$ -grams. The average complexity of the algorithm is  $O((k + \log_\sigma(rm))n/m)$  for  $\beta < 1/2 - O(1/\sqrt{\sigma})$ , where  $\beta = k/m$  is the difference ratio.

## 3.. SIMD TECHNIQUES

SIMD [2] is a type of parallel architecture that allows one instruction to be operated on multiple data items at the same time. Initially, SIMD was used in multimedia, especially in

processing images or audio files. SIMD instructions have since found applications in other areas such as cryptography. Recently, they have also been applied to string matching [30–34].

Streaming SIMD extensions (SSE) comprise of SIMD instruction sets supported by modern processors that allow computation on vectors of length 16 bytes in the case of SSE2 and 32 bytes in the case of AVX2. With AVX-512, one can process 64 bytes. The instructions operate on such vectors stored in special registers. As one instruction is performed on all the data in these vectors, it is considered SIMD computation.

Next, we describe our techniques to use SIMD in the new algorithms. In the descriptions, the SSE2 instructions are listed for 16 bytes (= 128 bits). There are corresponding AVX instructions for 32 bytes (= 256 bits). AVX-512 is not considered in this article. We assume that a byte represents one character.

### 3.1.. Counting of mismatches

Counting mismatches is a usual operation in approximate string matching. It can be done with the instructions *simd-cmpeq*( $x, y$ ) and *simd-popcount*( $x$ ) explained below. In practice, we also need the instruction *simd-load*( $x$ ), which uses an intrinsic function of the compiler. It is formally defined as follows:

`__m128i _mm_loadu_si128(__m128i * x).`

This function loads 16 bytes from the address  $x$  to a SIMD register given as the left-hand side of an assignment statement.

The function *simd-cmpeq*( $x, y$ ) is formally defined as

`_mm_movemask_epi8(_mm_cmpeq_epi8(_mm_loadu_si128(x), _mm_loadu_si128(y)))`

The intrinsic function `_mm_cmpeq_epi8` compares 16 bytes in  $x$  and  $y$  bitwise for equality. It returns a 16-byte value (of type `__m128i`) where each of its bytes equals 0 if the parameters' bytes in that position are different, or 255 (all of its bits set to 1, equal to 0xFF in hexadecimal) otherwise. For example, let  $S_1 = a...abcd$  and  $S_2 = a...aBce$  be strings of length 16 with  $a$  as their first 13 characters, then `_mm_cmpeq_epi8( $S_1, S_2$ ) = 0xFF...FF00FF00`, where the result's first 13 bytes have value 0xFF.

The other intrinsic function used, called `_mm_movemask_epi8`, creates a bitvector from the most significant bit of each of the 16 bytes of the parameter (which is of type `__m128i`), storing the result in the least significant 16 bits of a 32-bit integer value. For example, `_mm_movemask_epi8(0xFF00807F)` is 0xA.

The function *simd-popcount*( $x$ ) counts the number of on bits in a 32-bit integer  $x$  and returns the result as another integer. For example, *simd-popcount*(0xA) returns 2. It is formally defined as follows:

`int _mm_popcnt_u32(x).`

The *simd-cmpeq*( $x, y$ ) function, therefore, makes it possible to compare up to  $b$  characters at the same time, where  $b$  is 16 or 32. The result is a bitvector of the pairwise comparisons. Lastly, a popcount operation on the result tells the number of matching characters. For example,

```
z := simd-cmpeq( $S_1, S_2$ )
z = 0xFFFFA
matches := simd-popcount(z)
matches = 14
```

where  $S_1$  and  $S_2$  are the strings defined in the example of `_mm_cmpeq_epi8`.

### 3.2.. CRC as a fingerprint

There are many filtration methods for approximate string matching. Those methods contain two phases that are usually interleaved. The filtration phase selects match candidates and the checking phase verifies them. The former often entails the calculation of a fingerprint or a hash value from a  $q$ -gram, with which precomputed tables are accessed. Such a calculation can be performed with the *simd-crc*( $x$ ) instruction. A similar instruction was first used by Faro and Külekci [31, 35] in exact string matching.

The function *simd-crc*( $x$ ) first calculates a 32-bit cyclic redundancy checksum (CRC) of a 64-bit value  $x$  with an intrinsic instruction, and then takes the  $\alpha$  least significant bits of the CRC with a mask. It is formally

`_mm_crc32_u64(CRC, x) & mask`

where `_mm_crc32_u64` is an intrinsic function that starting with an unsigned 64-bit integer CRC, accumulates a CRC32 value for the unsigned 64-bit integer  $x$ , and returns it as another 64-bit integer. We use a fixed random constant for CRC and mask is  $2^\alpha - 1$ . Based on our experiments, the best value of  $\alpha$  depends on the problem parameters.

## 4.. IMPROVED SOLUTIONS—SINGLE PATTERN

### 4.1.. Variations of naive

A straightforward approach to string matching with at most  $k$  mismatches is the naive counting of mismatches. Algorithm 1 is the pseudocode of a naive algorithm ANS (short for Approximate Naive with SIMD) which applies SIMD computation. ANS counts the character matches with  $P$  starting from the  $n - m + 1$  first positions of the text for patterns of up to 32 characters. According to our experiments (see Section 6), it is clearly faster than both the classical Shift-Add [12] and TuSA [14].

**Algorithm 1:** ANS  
 $occ \leftarrow 0$   
 $x \leftarrow \text{simd-load}(p_0 \dots p_{m-1})$   
 for  $i \leftarrow 0$  to  $n - m$  do  
    $y \leftarrow \text{simd-load}(t_i \dots t_{i+m-1})$   
    $t \leftarrow \text{simd-cmpeq}(x, y)$   
   if  $\text{simd-popcount}(t) \geq m - k$  then  
      $occ \leftarrow occ + 1$   
 return  $occ$

Note that *simd-cmpeq* takes two 16-byte values (or 32-byte values if AVX2 is used) as arguments. In cases, where  $m < 16$  (or 32 in the case of AVX2), the values of  $x$  and  $y$  can be padded with bytes that differ so there are no false matches reported by it. For example,  $x$  could be padded with zeroes and  $y$  with ones.

There is a way to make ANS even faster. When  $m \leq 16$  holds, we preprocess the condition  $\text{simd-popcount}(t) \geq m - k$  to a boolean array  $D$  for each vector  $t$  of 16 bits. Then the last if statement of ANS is changed to

if  $D[t]$  then  $occ \leftarrow occ + 1$

For longer patterns,  $16 < m \leq 32$ , the last conditional statement becomes

if  $D[t \& \text{mask}]$  then  
   if  $\text{simd-popcount}(t) \geq m - k$  then  
      $occ \leftarrow occ + 1$

where  $\text{mask}$  is  $2^{16} - 1$ . In other words, the first 16 characters of the pattern are tested first. This variation is called ANS2. Pseudocode for its preprocessing phase is shown in Algorithm 2, and pseudocode for its search phase is shown in Algorithms 3 and 4 for pattern sizes of  $m \leq 16$  and  $16 < m \leq 32$ , respectively.

**Algorithm 2:** ANS2 Preprocessing  
 for  $i \leftarrow 0$  to  $2^{16} - 1$  do  
   if  $\text{simd-popcount}(i) \geq m - k$  then  
      $D[i] \leftarrow 1$   
   else  
      $D[i] \leftarrow 0$

**Algorithm 3:** ANS2 Search for  $m \leq 16$   
 $occ \leftarrow 0$   
 $x \leftarrow \text{simd-load}(p_0 \dots p_{m-1})$   
 for  $i \leftarrow 0$  to  $n - m$  do  
    $y \leftarrow \text{simd-load}(t_i \dots t_{i+m-1})$   
    $t \leftarrow \text{simd-cmpeq}(x, y)$   
   if  $D[t]$  then  
      $occ \leftarrow occ + 1$   
 return  $occ$

**Algorithm 4:** ANS2 Search for  $m > 16$   
 $occ \leftarrow 0$   
 $x \leftarrow \text{simd-load}(p_0 \dots p_{m-1})$   
 $\text{mask} \leftarrow 2^{16} - 1$   
 for  $i \leftarrow 0$  to  $n - m$  do  
    $y \leftarrow \text{simd-load}(t_i \dots t_{i+m-1})$   
    $t \leftarrow \text{simd-cmpeq}(x, y)$   
   if  $D[t \& \text{mask}]$  then  
     if  $\text{simd-popcount}(t) \geq m - k$  then  
        $occ \leftarrow occ + 1$   
 return  $occ$

In our test environment (see Section 6), the computation of  $D$  takes about 2 ms, which is tolerable. Note that the preprocessing time would grow exponentially, if  $D$  were extended for wider vectors. Note also that the speed of ANS does not depend on  $k$  for values of  $m$  that fit in a *simd-load* (i.e. 32 characters long). On the other hand, ANS2 performance degrades when  $k$  approaches  $m$  for  $m > 16$  because it is then more probable to apply the *simd-popcount*() function.

Furthermore, its speed also decreases for  $m \leq 16$  when  $k$  is rather large compared to  $m$ . For example, for  $m = 16$  its execution time peaks when  $k$  is around 11 in the DNA alphabet, as shown in Fig. 3. The main reason for this decrement in performance is an increment in the number of CPU conditional branch mispredictions, as analyzed in Section 6.

A way to avoid this dependency on  $k$  for pattern sizes not greater than 16 characters would be to change the last line

if  $D[t]$  then  $occ \leftarrow occ + 1$

into

$occ \leftarrow occ + D[t]$

where  $D$  is the boolean array described before, storing *true* values as ones and *false* values as zeros. This modification would force it to always perform an addition without a conditional statement. It saves time when there are many occurrences, which is more likely to happen when the difference ratio is high and the alphabet size is small. Let us call this variation ANS2b. The pseudocode for its search phase when  $m \leq 16$  is given as Algorithm 5. ANS2b is the same as ANS2 for  $m > 16$ .

**Algorithm 5:** ANS2b Search for  $m \leq 16$   
 $occ \leftarrow 0$   
 $x \leftarrow \text{simd-load}(p_0 \dots p_{m-1})$   
 for  $i \leftarrow 0$  to  $n - m$  do  
    $y \leftarrow \text{simd-load}(t_i \dots t_{i+m-1})$   
    $t \leftarrow \text{simd-cmpeq}(x, y)$   
    $occ \leftarrow occ + D[t]$   
 return  $occ$

Let  $b$  be the number of characters that fit in a single SIMD register used by the SIMD instructions employed in ANS



variations. For example,  $b$  is 16 and 32 for SSE4.2 and AVX2 instructions sets respectively. If ANS algorithm and its variations were extended to work with a general  $m$ , they would have a worst case time complexity in  $O(\lceil m/b \rceil n)$ . Because the naive algorithm has the average time complexity of  $O(kn)$ , a similar analysis on ANS results in the  $O(\lceil k/b \rceil n)$  time in the average case.

Besides the *simd-cmpeq* instruction and other basic SIMD commands, the SIMD architecture comprises of several aggregation operations for string processing. However, they are too slow for the  $k$  mismatches problem on those processors we have tested. Hirvola [36] implemented several algorithms similar to ANS with PCMP and STTNi instructions, but all those algorithms are clearly slower than ANS and TuSA.

#### 4.2.. EF enhanced with SIMD

EF [18] contains a filtration and a checking phase. The checking method can be replaced with ANS2, while the fingerprint computation of the filtration method can be replaced with the CRC fingerprint technique. Algorithm 6 shows the pseudocode of EF.

```

Algorithm 6: EF Search
occ ← 0
s ← m - 1
while s < n do
  /* Filtration phase. */
  f ←  $\sum_{i=0}^{q-1} \text{map}(t_{s-i}) * 4^i$ 
  if M[f] ≤ k then
    /* Checking phase. */
    c ← M[f]
    for i ← 1 to m - q do
      if  $t_{s-q-i+1} \neq p_{m-q-i}$  then
        c ← c + 1
      if c > k then break
    if c ≤ k then occ ← occ + 1
  s ← s + Sq[f]
return occ

```

For each  $q$ -gram  $u_0 \dots u_{q-1}$ , the preprocessing phase of EF computes the Hamming distance with the end of all strict prefixes of the pattern. With this information, a shift table  $S_q$  can be constructed: given a  $q$ -gram  $f$ ,  $S_q[f]$  returns  $m - q - \text{good}_f$ , where  $\text{good}_f$  is the position of the first character of the rightmost  $q$ -gram of the pattern (excluding its last  $q$ -gram) that has at most  $k$  mismatches with  $f$ . If there are several such  $q$ -grams, then the minimum value is stored to  $S_q[f]$ . If there are no pattern  $q$ -grams with at most  $k$  mismatches with  $f$ , then  $S_q[f] = m - q + 1$ , which corresponds to the largest possible shift. This table is used to compute the jump to the next text position to be analyzed.

$M$  is another precomputed table, which gives the Hamming distance of a  $q$ -gram against the last  $q$ -gram of the pattern. Whenever  $M[t_{s-q+1} \dots t_s] > k$  holds, the algorithm shifts forward without processing the alignment window fur-

ther. Both the tables are accessed with the fingerprint  $f \leftarrow \sum_{i=0}^{q-1} \text{map}(t_{s-i}) * 4^i$ , where the function *map* maps each DNA character to an integer in  $\{0, 1, 2, 3\}$ .

Algorithm 7 is the pseudocode of EFS (short for EF with SIMD) for  $m \leq 16$ . Like EF, EFS is intended for small alphabets like the DNA alphabet. The array  $D$  is computed in the same way as for ANS2. For longer patterns,  $16 < m \leq 32$ , the required change is the same as in the case of ANS2.

```

Algorithm 7: EFS search for  $m \leq 16$ 
occ ← 0
x ← simd-load(p0 ... pm-1)
s ← m - 1
while s < n do
  /* Filtration phase. */
  f ← simd-crc(ts-q+1 ... ts)
  if M[f] ≤ k then
    /* Checking phase. */
    y ← simd-load(ts-m+1 ... ts)
    t ← simd-cmpeq(x, y)
    if D[t] then
      occ ← occ + 1
  s ← s + Sq[f]
return occ

```

#### 4.3.. BYP enhanced with SIMD

The original BYP [11] looks for *exact* occurrences of  $k + 1$  subpatterns of the pattern in the text:  $k + 1 - (m \bmod (k + 1))$  of length  $l$  and  $m \bmod (k + 1)$  of length  $l + 1$ , where  $l = \lfloor m/(k + 1) \rfloor$ . To achieve this, we employed a tuned version of the MEPSM algorithm [35] for exact multiple string matching. MEPSM reports subpattern occurrences, which are later verified by ANS2. Let us call the total algorithm BYPS.

MEPSM computes CRC fingerprints for  $q$ -grams of each subpattern, where  $q \leq l$  is a parameter of MEPSM. The information about which  $q$ -gram the fingerprint belongs to is stored in a table. Afterwards, during the search, the algorithm looks for matching fingerprints of  $q$ -grams in the text. Whenever a subpattern candidate is found, it is naively verified and forwarded to further processing in case of a match. After each  $q$ -gram analysis, the algorithm shifts forward by  $\text{shift} = l - q + 1$  characters.

BYPS can then be summarized with the following steps: (i) in the preprocessing phase, we split every pattern into  $k + 1$  subpatterns as explained at the beginning of this section. Then we compute the CRC fingerprint of the first  $\text{shift}$   $q$ -grams of each subpattern. The fingerprint is used to access a table that stores information about which subpattern it was computed from. (ii) In the search, we compute the fingerprint of a  $q$ -gram in the text, with which we fetch the corresponding information from the table. We perform a shift of  $\text{shift}$  characters in the text after analyzing a  $q$ -gram. Here  $\text{shift}$  is the maximum number of characters we can skip. (iii) For every subpattern associated with the fingerprint, we naively check if it *exactly* appears at

this point. If it does, a possible approximate occurrence of the pattern is detected. (iv) Every time a match candidate of the pattern is found, we use an *approximate* single pattern matching algorithm to verify it.

Note that the last  $q$ -gram of a subpattern of length  $l + 1$  is not preprocessed. Let  $A$  be a subpattern of length  $l + 1$  and let  $A$  appear exactly in the text.  $A$  is recognized even if the algorithm reads the last  $q$ -gram of  $A$ , because it then reads also the first  $q$ -gram of  $A$ , as the offset between those  $q$ -grams is precisely  $l - q + 1$ , the step of the algorithm.

To avoid reverifying an occurrence, we keep track of the position up to which the text has already been analyzed. In this way, we do not need to check for occurrences starting at positions to the left of the last one analyzed. This would be trivial if MEPSM reported exact subpattern occurrences corresponding to ordered approximate pattern matches, but it does not. This has been solved by executing the approximate single pattern matching algorithm in a larger window. If there is an occurrence at position  $x$  in the text of a subpattern that begins at position  $sp$  of the pattern, we check for an approximate occurrence of the pattern from position  $x - (m - l)$  to  $x + m - sp$ . Thus, once an occurrence of the pattern has been found, a newer occurrence will never precede it positionally.

We tuned MEPSM by setting the  $q$ -grams as large as possible, that is  $q = \min(l, 8)$ . The maximum value for  $q$  is 8 because we need to compute the fingerprint of  $q$ -grams with *simd-crc*( $x$ ). To do so, it uses a SIMD instruction, which means it can take up to a 64-bit value  $x$  as input, i.e. 8 bytes. This causes fewer fingerprint collisions, but on the other hand, larger  $q$  reduces shifts between alignments ( $shift = l - q + 1$ ). However, this trade-off showed to be satisfactory in practice, especially in the case of small subpatterns.

This tuning is then employed in combinations of  $m$  and  $k$  that make the length of subpatterns searched by MEPSM smaller than 8 characters (i.e.  $l < 8$ ). On the other hand, the minimum value admitted for  $q$  (and thus for  $l$ ) is 4 because the number of collisions while hashing would quickly rise otherwise, critically slowing down the algorithm.

We also present a variation of BYPS where we use the ANS2b algorithm for searching approximate occurrences of the pattern in the neighborhood of exact subpattern matches. In this algorithm called BYPSb, we also replaced the naive exact subpattern check done by MEPSM after finding a  $q$ -gram matching fingerprint in the text, with a SIMD-based comparison function *simd-memcmp*, which employs function *simd-cmpeq*. Function *simd-memcmp* is shown in Algorithm 8 and works for strings of equal size  $|s|$  up to 32 characters<sup>1</sup> using AVX2 instructions. Its time complexity is in  $O(\lceil |s|/b \rceil)$ , where  $b$  is the number of characters that fit in a SIMD register of the extension set available. Given that we are using AVX2 where  $b = 32$ , then it can be said it works in constant time.

The preprocessing and search phases of BYPSb for cases when  $l < 8$  are shown in Algorithms 9 and 10, respectively. For cases when  $l \geq 8$  the original MEPSM calculation of  $q$  is used. Symbols  $\langle \rangle$  are used to denote the beginning and end of a list of elements, and  $\#$  denotes list concatenation.

We implemented another version of BYPSb which completely skips the step of exactly verifying the occurrence of a subpattern and assumes a match every time there is a fingerprint match in the text. Given that  $q = \min(l, 8)$ , this step would not be necessary in cases when  $q = l$  if the hash function worked perfectly. Note that if a wrong subpattern match assumption is made, the efficiency of the algorithm is not affected because ANS2b analyzes the window afterwards. Let us call it BYPSc.

**Algorithm 8:** *simd-memcmp*( $s_1, s_2$ )  
 $eqmask \leftarrow (1 << |s_1|) - 1$   
 $x \leftarrow \text{simd-load}(s_1)$   
 $y \leftarrow \text{simd-load}(s_2)$   
 $mask \leftarrow \text{simd-cmpeq}(x, y)$   
 return  $(eqmask \sim \& \sim mask) = eqmask$

**Algorithm 9:** BYPSb Preprocessing  
 ANS2bPreprocessing( $m, k$ )  
 $l \leftarrow \lfloor m/(k+1) \rfloor$   
 $rem \leftarrow m \bmod l$   
 $q \leftarrow \min(l, 8)$   
 $shift \leftarrow l - q + 1$   
 for  $f \leftarrow 0$  to  $2^{16} - 1$  do  
 $H[f] \leftarrow \langle \rangle$   
  
 /\* For every subpattern of length  $l$ , calculate the  
 fingerprint of its first  $shift$   $q$ -grams. \*/  
 for  $subpat \leftarrow 0$  to  $m - l - rem * (l + 1)$  step  $l$  do  
 for  $j \leftarrow 0$  to  $shift - 1$  do  
 $f \leftarrow \text{simd-crc}(p_{subpat+j} \sim \dots \sim p_{subpat+j+q-1})$   
 $H[f] \leftarrow \langle (subpat, j) \rangle \# H[f]$   
  
 /\* Do the same for subpatterns of length  $l + 1$ . \*/  
 for  $subpat \leftarrow m - rem * (l + 1)$  to  $m - (l + 1)$  step  $l + 1$  do  
 for  $j \leftarrow 0$  to  $shift - 1$  do  
 $f \leftarrow \text{simd-crc}(p_{subpat+j} \sim \dots \sim p_{subpat+j+q-1})$   
 $H[f] \leftarrow \langle (subpat, j) \rangle \# H[f]$

In the worst case, BYPSb time complexity is equal to executing ANS2b over all the text plus the work needed to check the exact occurrences of all corresponding subpatterns at each of the  $O(n/shift)$  text alignments for each fingerprint match. Assuming the worst case for the hash function in which all text fingerprints match every preprocessed fingerprint of the  $shift$   $q$ -grams of each of the  $k + 1$  subpatterns, its complexity is in

$$O(\lceil m/b \rceil n + (k + 1) \cdot shift \lceil l/b \rceil n / shift) =$$

$$O(\lceil m/b \rceil n + k \lceil l/b \rceil n), \text{ for } k > 0 \quad (1)$$

<sup>1</sup> Its extension to general sizes is straightforward.

**Algorithm 10:** BYPSb Search

```

lastpos ← m - 1
occ ← 0
for i ← 0 to n - q step shift do
    /* Calculate the fingerprint of the q-gram
       starting at the current text position. */
    f ← simd-crc(ti ~ ... ~ ti+q-1)

    /* Analyze each preprocessed q-gram with fingerprint f.
       Each of these q-grams starts at position j inside
       the subpattern that starts at position subpat
       in the pattern. */
    forall (subpat, j) in H[f] do

        /* Compare the subpattern with the corresponding
           text window for equality. */
        if simd-memcmp(
            ti-j ~ ... ~ ti-j+l-1,
            psubpat ~ ... ~ psubpat+l-1
        ) then

            /* Look for an approximate pattern match in
               the corresponding text window. */
            begin ← max{i - j - (m - l), ~ lastpos - m + 1}
            end ← min{i - j - subpat + m - 1, ~ n - 1}
            if end - begin + 1 ≥ m then
                lastpos ← end
                occ ← occ + ANS2bSearch(tbegin ~ ... ~ tend, p, k)
return occ

```

Now let us analyze the average case of BYPSb. As in the worst case, we will split this calculation into two pieces:  $E_{exact}$ , the expected work needed to obtain all exact subpattern matches in the processed text alignments, and  $E_{approx}$ , the total average amount of time demanded by ANS2b searches. In this way, the average time complexity of BYPSb can be written as:

$$O(E_{exact} + E_{approx}) \quad (2)$$

In first place,  $E_{exact}$  corresponds to traversing the text comparing  $l$  characters with at most  $k + 1$  subpatterns using *simd-memcmp* at most  $shift$  times at each of the  $O(n/shift)$  alignments. Then

$$\begin{aligned} E_{exact} &\in O((k + 1)shift \lceil l/b \rceil n / shift) \\ &= O(k \lceil l/b \rceil n), \text{ for } k > 0 \end{aligned} \quad (3)$$

Note that it is a pessimistic bound as we are assuming no hash filtering, otherwise *simd-memcmp* is only called whenever there is a hash occurrence.

On the other hand, we could optimistically assume a perfect hash function. Moreover, if we also simplify the problem conditions by assuming that all  $(k + 1)shift$   $q$ -grams analyzed inside the pattern are independent (note that they might overlap and perhaps appear very often in certain real-life patterns) and are uniformly distributed over the set of possible  $q$ -grams built with characters from  $\Sigma$ , we get that the expected number

of  $q$ -gram matches after analyzing a  $q$ -gram fingerprint at a given text position is  $E_{qmatches} = (k + 1)shift / \sigma^q$ . We analyze  $O(n/shift)$  text alignments and the expected work needed by *simd-memcmp* to compare  $l$  text characters with a subpattern is  $\lceil l/b \rceil$ , so we get

$$\begin{aligned} E_{exact} &\in O\left(E_{qmatches} \lceil l/b \rceil \frac{n}{shift}\right) \\ &= O\left(\frac{(k + 1)shift}{\sigma^q} \lceil l/b \rceil \frac{n}{shift}\right) \\ &= O\left(\frac{k \lceil l/b \rceil n}{\sigma^q}\right), \text{ for } k > 0 \end{aligned} \quad (4)$$

For the calculation of  $E_{approx}$ , it is important to note that it cannot be greater than the average complexity of ANS2b  $O(\lceil k/b \rceil n)$  because it is applied at most over all the text. Now let us call the expected number of exact subpattern occurrences in the text  $E_{subpatterns}$ , and the expected work needed to check for an approximate occurrence of the pattern each time a subpattern is found  $E_{pattern}$ . Then we can rewrite  $E_{approx}$  as follows:

$$E_{approx} \in O(\min(\lceil k/b \rceil n, E_{subpatterns} \times E_{pattern})) \quad (5)$$

The probability of exactly matching a subpattern in a given alignment is  $1/\sigma^l$ , which coincides with the expected number of matches for this subpattern in one alignment. As we have at most  $k + 1$  different subpatterns, then by linearity of expectation we get that the expected number of subpattern matches in a given alignment is at most  $(k + 1)/\sigma^l$ . Also, BYPSb analyzes  $O(n/shift)$  alignments of subpatterns in the text. So we can conclude that

$$E_{subpatterns} \in O\left(\frac{n}{shift} \frac{k + 1}{\sigma^l}\right) \quad (6)$$

On the other hand,  $E_{pattern}$  corresponds to the expected work needed by ANS2b to look for an approximate occurrence of the pattern in a window of size at most  $2m$  characters. Then, we can affirm that

$$E_{pattern} \in O(\lceil k/b \rceil m) \quad (7)$$

Replacing Equations 6 and 7 in 5, we obtain

$$\begin{aligned} E_{approx} &\in O\left(\min\left(\lceil k/b \rceil n, \frac{n}{shift} \frac{k + 1}{\sigma^l} \lceil k/b \rceil m\right)\right) \\ &= O\left(\min\left(\lceil k/b \rceil n, \frac{\lceil k/b \rceil kmn}{\sigma^l(l - q + 1)}\right)\right), \text{ for } k > 0 \end{aligned} \quad (8)$$

Finally, if we assume no hash filtering, we can replace Equations 3 and 8 in 2 and get an average time complexity of

BYPSb in

$$O\left(k\lceil l/b \rceil n + \min\left(\lceil k/b \rceil n, \frac{\lceil k/b \rceil kmn}{\sigma^l \text{shift}}\right)\right) \quad (9)$$

$$= O(k\lceil l/b \rceil n), \text{ for } k > 0$$

Otherwise, under the assumption of perfect hashing and from Equation 4, we get a total average complexity of BYPSb in

$$O\left(\frac{k\lceil l/b \rceil n}{\sigma^q} + \min\left(\lceil k/b \rceil n, \frac{\lceil k/b \rceil kmn}{\sigma^l \text{shift}}\right)\right)$$

$$= O\left(\frac{k\lceil l/b \rceil n}{\sigma^q} + \min\left(\lceil k/b \rceil n, \frac{\lceil k/b \rceil kmn}{\sigma^{l(l-q+1)}}\right)\right), \quad (10)$$

for  $k > 0$

From this analysis, it can be predicted a faster execution of BYPSb for lower difference ratios (because  $l = \lfloor m/(k+1) \rfloor$  and  $q = \min(l, 8)$ ) and for larger alphabets.

As regards BYPSc, we omit the step of using *simd-cmpeq* after a fingerprint coincidence, so from Equation 3 we get

$$E_{\text{exact}} \in O(n/\text{shift}) = O(n/(l-q+1)) \quad (11)$$

Now for the calculation of  $E_{\text{approx}}$  of BYPSc, we consider the expected number of exact subpattern occurrences  $E_{\text{subpatterns}}$  to be the same as the expected number of fingerprint coincidences. There are at most  $(k+1)\text{shift}$  different fingerprints, so if we assume a perfect hashing function and follow an analogous reasoning as that used in Equation 6, we get

$$E_{\text{subpatterns}} \in O\left(\frac{n}{\text{shift}} \frac{(k+1)\text{shift}}{\sigma^q}\right)$$

$$= O\left(\frac{nk}{\sigma^q}\right), \text{ for } k > 0 \quad (12)$$

Under the assumption of no hash filtering, every hash value would be the same, so we would expect to assume that all the preprocessed  $q$ -grams of all subpatterns exactly appear at each text position analyzed. Then, we get

$$E_{\text{subpatterns}} \in O\left(\frac{n}{\text{shift}} (k+1)\text{shift}\right)$$

$$= O(nk), \text{ for } k > 0 \quad (13)$$

Given that  $E_{\text{pattern}}$  does not change in BYPSc, if we assume perfect hashing, we can replace Equations 7 and 12 in Equation

5 and obtain

$$E_{\text{approx}} \in O\left(\min\left(\lceil k/b \rceil n, \frac{n}{\text{shift}} \frac{(k+1)\text{shift}}{\sigma^q} \lceil k/b \rceil m\right)\right) \quad (14)$$

$$= O\left(\min\left(\lceil k/b \rceil n, \frac{\lceil k/b \rceil kmn}{\sigma^q}\right)\right), \text{ for } k > 0$$

Otherwise, assuming no hash filtering, from Equation 13 we get

$$E_{\text{approx}} \in O(\min(\lceil k/b \rceil n, n(k+1)\lceil k/b \rceil m))$$

$$= O(\lceil k/b \rceil n) \quad (15)$$

Finally, assuming no hash filtering, we replace Equations 11 and 15 in Equation 2 and get a time complexity of BYPSc in

$$O\left(\frac{n}{\text{shift}} + \lceil k/b \rceil n\right) = O(\lceil k/b \rceil n) \quad (16)$$

On the other hand, if we assume perfect hashing, from Equation 14 we obtain the following average time complexity

$$O\left(\frac{n}{\text{shift}} + \min\left(\lceil k/b \rceil n, \frac{\lceil k/b \rceil kmn}{\sigma^q}\right)\right)$$

$$= O\left(\frac{n}{l-q+1} + \min\left(\lceil k/b \rceil n, \frac{\lceil k/b \rceil kmn}{\sigma^q}\right)\right), \quad (17)$$

for  $k > 0$

## 5.. IMPROVED SOLUTION—MULTIPLE PATTERNS

We extend the BYPS algorithm to work with multiple patterns. The new algorithm MBYPS works as follows (changes respect to BYPS have been highlighted in bold red):

- (i) In the preprocessing phase, we split every pattern into subpatterns of length  $l$  or  $l+1$  in the same way as in BYPS, where  $l = \lfloor m/(k+1) \rfloor$ . Then, we compute the CRC fingerprint of  $\text{shift}$   $q$ -grams of each subpattern, where  $q \leq l$  is a tuned parameter of the MEPSM algorithm, and  $\text{shift} = l - q + 1$ . The fingerprint is used to access a table that stores information about which subpattern of which pattern it was computed from.
- (ii) In the search, we compute the fingerprint of a  $q$ -gram in the text, with which we fetch the corresponding information from the table. We perform a shift of  $\text{shift}$  characters in the text after analyzing each  $q$ -gram, which is the maximum number of characters we can skip.
- (iii) For every subpattern associated with the fingerprint, we naively check if it *exactly* appears at this point. If it does, a possible approximate occurrence of the corresponding pattern is detected.
- (iv) Every time a match candidate of a pattern is found, we use an *approximate* single pattern matching algorithm to verify it.



For the phase of exact multiple string matching, we use our tuned version of MEPSM, as described in Section 4.3, except that now we extend our tuning of the  $q$  calculation to all subpattern lengths. For the phase of approximate single string matching, we use ANS2 for  $m \leq 32$ . For longer patterns, another algorithm should be used.

Following the same reasoning as in BYPSb, we also consider a variation of MBYPS where ANS2b is used instead of ANS2 for searching approximate pattern occurrences and *simd-memcmp* is used for exactly comparing strings, which we call MBYPSb. Algorithms 11 and 12 show pseudocode for MBYPSb preprocessing and search phases respectively. Changes respect to BYPSb have been highlighted in red. Note that we are comparing MBYPSb with BYPSb version for  $l < 8$  because MBYPSb uses this same tuning for calculating  $q$  for all values of  $l$ .

Furthermore, we implemented another version of MBYPSb analogous to BYPSb that skips the step of exactly comparing a subpattern with a text window of  $l$  characters after a fingerprint match. We call it MBYPSc.

As regards the time complexity of MBYPSb, there are two differences with the analysis done for BYPSb. One is that MBYPSb looks for exact occurrences of at most  $(k+1)r$  subpatterns instead of just  $k+1$  as in BYPSb. The second one is that a subpattern occurrence can now belong to many patterns, so it needs to perform at most  $r$  ANS2b searches in each window of at most  $2m$  characters. In this way, by multiplying both terms of the sum of Equation 1 by  $r$ , we can conclude that MBYPSb time complexity in the worst case is in

$$O(r\lceil m/b \rceil n + rk\lceil l/b \rceil n), \text{ for } k > 0 \quad (18)$$

As regards the average time complexity of MBYPSb, if we assume no hash filtering,  $E_{exact}$  from Equation 3 is increased by a factor of  $r$

$$E_{exact} \in O(rk\lceil l/b \rceil n), \text{ for } k > 0 \quad (19)$$

On the other hand, if we assume a perfect hash function and that all  $q$ -grams analyzed are independent and uniformly distributed over the set of possible  $q$ -grams, we get that  $E_{qmatches} = r(k+1)shift/\sigma^q$  since there are now  $r(k+1)shift$  preprocessed  $q$ -grams. Then, following the same reasoning from Equation 4, we get

$$\begin{aligned} E_{exact} &\in O\left(E_{qmatches}\lceil l/b \rceil \frac{n}{shift}\right) \\ &= O\left(\frac{r(k+1)shift}{\sigma^q}\lceil l/b \rceil \frac{n}{shift}\right) \quad (20) \\ &= O\left(\frac{rk\lceil l/b \rceil n}{\sigma^q}\right), \text{ for } k > 0 \end{aligned}$$

The increment in the number of subpatterns also affects the average time complexity of ANS2b applied over all the text, because now it is employed at most once for each of the  $r$  patterns. Consequently, from Equation 5 we get

$$\begin{aligned} E_{approx} &\in \\ O\left(\min\left(r\lceil k/b \rceil n, E_{subpatterns} \times E_{pattern}\right)\right) \end{aligned} \quad (21)$$

This factor of  $r$  is propagated to  $E_{subpatterns}$  too, because now we have at most  $r(k+1)$  different subpatterns. So, from Equation 6 it can be deduced that

$$E_{subpatterns} \in O\left(\frac{n}{shift} \frac{(k+1)r}{\sigma^l}\right) \quad (22)$$

Since now the algorithm needs to perform at most  $r$  ANS2b searches for every subpattern match, Equation 7 becomes

$$E_{pattern} \in O(\lceil k/b \rceil mr) \quad (23)$$

Then, replacing Equations 22 and 23 in 21, and then replacing it and Equation 19 in 2 we obtain that MBYPSb search phase average time complexity if there is no hash filtering is in

$$\begin{aligned} &O\left(rk\lceil l/b \rceil n + \right. \\ &\left. \min\left(r\lceil k/b \rceil n, \frac{n}{shift} \frac{(k+1)r}{\sigma^l} \lceil k/b \rceil mr\right)\right) \quad (24) \\ &= O(rk\lceil l/b \rceil n), \text{ for } k > 0 \end{aligned}$$

Or, assuming perfect hashing, from Equation 20, we obtain a total average time complexity of MBYPSb in

$$\begin{aligned} &O\left(\frac{rk\lceil l/b \rceil n}{\sigma^q} + \min\left(r\lceil k/b \rceil n, \frac{r^2 k\lceil k/b \rceil mn}{\sigma^l shift}\right)\right) \\ &= O\left(\frac{rk\lceil l/b \rceil n}{\sigma^q} + \min\left(r\lceil k/b \rceil n, \frac{r^2 k\lceil k/b \rceil mn}{\sigma^l(l-q+1)}\right)\right), \quad (25) \\ &\text{for } k > 0. \end{aligned}$$

Note that we pessimistically assumed in Equation 23 that each subpattern found appears in all patterns. Nonetheless, this analysis implies that the lower the difference ratio (keep in mind that  $l = \lfloor m/(k+1) \rfloor$ , so it becomes larger), the bigger the alphabet and the smaller the set of patterns, the faster MBYPSb will perform.

As regards MBYPSc, following the same reasoning as in Section 4.3 and the analysis done for MBYPSb, if we assume



**Algorithm 12:** MBYPSb Search

```

occ ← 0
lastpospat ← m - 1 ~  $\forall \sim pat$ 
for i ← 0 to n - q step shift do
    /* Calculate the fingerprint of the q-gram
       starting at the current text position. */
    f ← simd-crc( $t_i \sim \dots \sim t_{i+q-1}$ )

    /* Analyze each preprocessed q-gram with fingerprint f.
       Each of these q-grams starts at position j inside
       the subpattern that starts at position subpat
       in pattern pat. */
    forall (pat, subpat, j) in H[f] do

        /* Compare the subpattern with the corresponding
           text window for equality. */
        if simd-memcmp(
             $t_{i-j} \sim \dots \sim t_{i-j+l-1}$ ,
             $p_{subpat}^{pat} \sim \dots \sim p_{subpat+l-1}^{pat}$ 
        ) then

            /* Look for an approximate pattern match in
               the corresponding text window. */
            begin ← max{i - j - (m - l), ~ lastpospat - m + 1}
            end ← min{i - j - subpat + m - 1, ~ n - 1}
            if end - begin + 1 ≥ m then
                lastpospat ← end
                occ ← occ + ANS2bSearch( $t_{begin} \sim \dots \sim t_{end}$ ,  $p^{pat}$ , k)
return occ

```

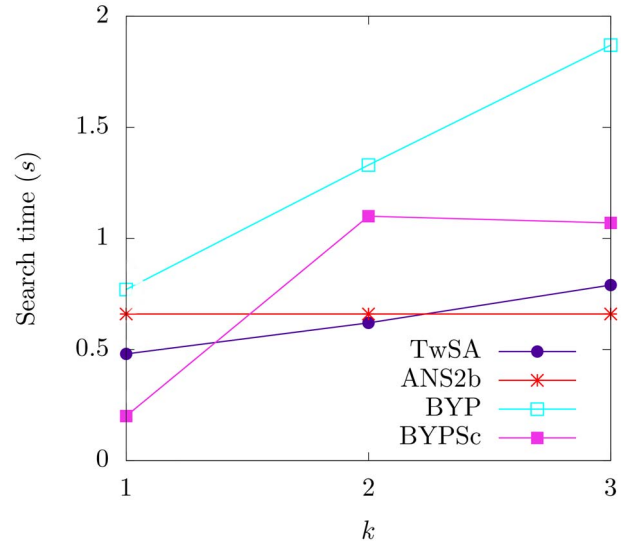
time because we address the *online* version of the problem, so it becomes insignificant compared with search time for an arbitrarily long text.

We also ran tests on a protein sequence (from the Human sequence genome, 3.1 MiB) taken from the SMART corpus. We obtained similar results to those reported in this chapter for the English alphabet for single string matching and the DNA alphabet in the case of multiple string matching. We omit these results to avoid redundancy.

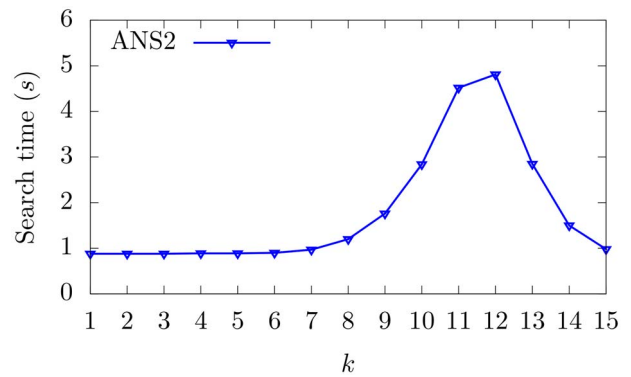
### 6.1.. Single pattern matching with up to $k$ mismatches

Algorithms that were compared for this problem are as follows: SA: Shift-Add [12]; TuSA: Tuned Shift-Add [14]; TwSA: Two-way Shift-Add [14]; EF: Enhanced FFAST [18]; EFS: Enhanced FFAST with SIMD; ANS, ANS2 and ANS2b: Approximate Naive enhanced with SIMD; BYP: the original Baeza-Yates-Perleberg algorithm [11]; BYPS, BYPSb and BYPSc: Baeza-Yates-Perleberg enhanced with SIMD.

According to tests by Hirvola [36], TwSA is the best for the English data. According to tests by Salmela *et al.* [18], EF is the best for DNA data. Our results are shown in Table 1 with the best times highlighted. In Figures 1 and 2, there are plots of search times of the fastest variations of tested algorithms as a function of  $k$  for  $m = 16$  in both alphabets. ANS2b and BYPSb/BYPSc were almost always the fastest for all parameter combinations on both DNA and English. They were



**FIGURE 2.** Search times of the fastest algorithms tested for single pattern matching as a function of  $k$  for  $m = 16$  in the English alphabet.



**FIGURE 3.** Search times of ANS2 as a function of  $k$  for  $m = 16$  in the DNA alphabet.

only surpassed by TwSA in the English alphabet when  $m = 16$  and  $k = 2$ . BYPSb and BYPSc were the best for cases with a low difference ratio. BYPSb performed better in DNA whereas BYPSc did so on English. On the other hand, ANS2b worked better than BYPSb/BYPSc when the difference ratio was rather high.

It can be observed a clear improvement in the performance of SIMD-based algorithms compared to their non-SIMD original versions, specially in the case of BYP and BYPSb/BYPSc. The most notorious change is seen in the DNA alphabet, where BYPSc got an average speedup of x8.5 over BYP across all our tests, and a maximum speedup of x24.5 in the case of  $k = 1$  and  $m = 32$ . For the English alphabet, BYPSc got an average improvement of x3, with a maximum of x4.4 in the case of  $k = 1$  and  $m = 32$ .

ANS, ANS2 and ANS2b work for all possible values of  $k$ . ANS did so at an almost constant speed independent of the

**TABLE 1** Search times (in seconds) of algorithms for single approximate pattern matching with up to  $k$  mismatches ran 100 times with 100 different patterns.

	$m = 8$			$m = 16$			$m = 24$			$m = 32$			
$k$	1	2	3	1	2	3	1	2	3	1	2	3	$\Sigma$
SA	1.68	1.69	1.69	1.68	1.69	1.68	1.68	$_{-a}$	$_{-a}$	1.68	$_{-a}$	$_{-a}$	DNA
TuSA	1.31	1.31	1.31	1.31	1.31	1.31	1.31	$_{-a}$	$_{-a}$	1.31	$_{-a}$	$_{-a}$	
TwSA	1.62	2.13	2.40	0.81	1.07	1.29	0.55	$_{-a}$	$_{-a}$	0.41	$_{-a}$	$_{-a}$	
ANS	1.25	1.25	1.25	1.25	1.25	1.25	1.35	1.35	1.34	1.34	1.34	1.34	
ANS2	0.86	0.91	1.16	0.88	0.88	0.88	$_{-b}$	$_{-b}$	$_{-b}$	$_{-b}$	$_{-b}$	$_{-b}$	
ANS2b	<b>0.75</b>	<b>0.75</b>	<b>0.75</b>	0.78	<b>0.78</b>	<b>0.78</b>	1.05	1.05	<b>1.07</b>	1.05	1.05	1.07	
EF	1.46	2.28	3.90	0.70	1.04	1.75	0.49	0.77	1.37	0.39	0.64	1.20	
EFS	1.41	2.14	3.92	0.67	0.97	1.59	0.48	0.71	1.24	0.38	0.59	1.10	
BYP	4.18	10.13	14.82	3.56	5.32	8.23	3.62	5.16	6.36	3.67	4.99	5.83	
BYPS	1.60	$_{-a}$	$_{-a}$	0.35	1.56	1.93	0.25	0.42	1.60	0.19	0.35	0.48	
BYPSb	1.43	$_{-a}$	$_{-a}$	<b>0.30</b>	1.36	1.84	<b>0.20</b>	<b>0.35</b>	1.42	0.18	<b>0.28</b>	<b>0.40</b>	
BYPSc	1.16	$_{-a}$	$_{-a}$	0.37	1.10	1.42	0.42	0.63	1.18	<b>0.15</b>	0.65	0.84	
SA	1.47	1.47	1.47	1.47	1.47	1.47	1.47	$_{-a}$	$_{-a}$	1.47	$_{-a}$	$_{-a}$	English
TuSA	1.14	1.14	1.14	1.14	1.14	1.14	1.14	$_{-a}$	$_{-a}$	1.14	$_{-a}$	$_{-a}$	
TwSA	0.83	1.17	1.53	0.48	<b>0.62</b>	0.79	0.33	$_{-a}$	$_{-a}$	0.26	$_{-a}$	$_{-a}$	
ANS	1.09	1.09	1.09	1.09	1.09	1.09	1.17	1.17	1.17	1.17	1.17	1.17	
ANS2	0.75	0.75	0.76	0.75	0.75	0.75	$_{-b}$	$_{-b}$	$_{-b}$	$_{-b}$	$_{-b}$	$_{-b}$	
ANS2b	<b>0.65</b>	<b>0.65</b>	<b>0.65</b>	0.66	0.66	<b>0.66</b>	0.91	0.91	<b>0.91</b>	0.91	0.91	0.91	
BYP	1.19	2.29	3.24	0.77	1.33	1.87	0.54	0.92	1.31	0.49	0.80	1.08	
BYPS	1.37	$_{-a}$	$_{-a}$	0.27	1.43	1.42	0.16	0.28	1.44	0.17	0.20	0.28	
BYPSb	1.19	$_{-a}$	$_{-a}$	0.24	1.25	1.23	<b>0.14</b>	0.25	1.27	0.15	<b>0.18</b>	0.25	
BYPSc	1.04	$_{-a}$	$_{-a}$	<b>0.20</b>	1.10	1.07	0.15	<b>0.23</b>	1.12	<b>0.13</b>	<b>0.18</b>	<b>0.24</b>	

<sup>a</sup> Algorithm not designed to work in this case. <sup>b</sup> Same as ANS2b.

value of  $k$ . The speed of ANS2 decreased as  $k$  grew, which was most noticeable when  $m = 8$  in DNA because of the high number of occurrences. As regards ANS2b, its performance was independent from  $k$  for  $m \leq 16$  as explained in Section 4.1. It showed a clear improvement compared to ANS2 surpassing it in all tests, even in cases of few occurrences (low difference ratios). Timings for ANS2 when  $m > 16$  have been omitted as it works in the same way as ANS2b.

We also ran tests of ANS2 for all possible values of  $k$  for the DNA alphabet and  $m = 16$ , getting timings plotted in Figure 3. By using the perf tool<sup>3</sup>, we gathered information about those runs. In particular, the percentage of conditional branches that were wrongly predicted and the number of instructions it took each execution to finish. The results are reported in Figures 4 and 5, respectively, as a function of  $k$ . It is easy to see that the peak in execution time of ANS2 is based on an increment in the number of branch mispredictions. The number of instructions executed by the program also rose as  $k$  grew because there were more approximate pattern occurrences so it needed to

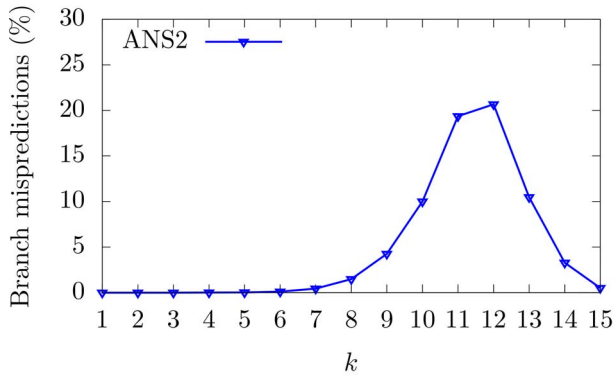
perform more additions to the occurrence counter. However, it is not as meaningful as the variation in the amount of branch mispredictions. The same testing was performed on ANS2b, which showed a constant behavior for all  $k$ . Its execution time, number of instructions executed and percentage of branch mispredictions were independent of the value of  $k$ . This characteristic is also evident in its search times in Table 1.

SA, TuSA and TwSA are limited to small values of  $k$  for long patterns. For example, they only work for  $k = 1$  in the case of  $m \in \{24, 32\}$ . Furthermore, the speed of TwSA degraded as  $k$  grows. EF, EFS, BYP, BYPS, BYPSb and BYPSc exhibited a similar behavior, with  $k$  affecting their speed. Despite this, the growth of  $k$  can be tolerated given that  $m$  is large enough, i.e. when we have small difference ratios. Some timings of BYPS and its variations have been omitted because they do not work for  $l = \lfloor m/(k+1) \rfloor < 4$ .

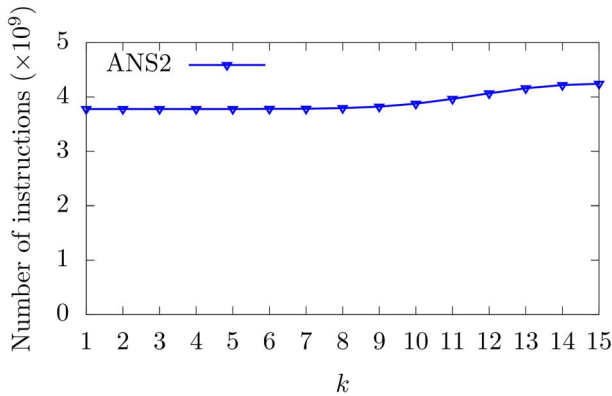
BYPS was also tested for longer patterns. According to our experiments and following the same line stated by Baeza-Yates and Perleberg [11], BYP and BYPS and its variations obtained their best results for low difference ratios. As regards BYPSb, we found out that it was always faster than BYPS. On the other hand, BYPSc sometimes worked faster than BYPSb

<sup>3</sup> Performance analysis tool for Linux. [https://perf.wiki.kernel.org/index.php/Main\\_Page](https://perf.wiki.kernel.org/index.php/Main_Page)





**FIGURE 4.** Percentage of branch prediction misses of the total of conditional branches executed in ANS2 as a function of  $k$  for  $m = 16$  in the DNA alphabet.



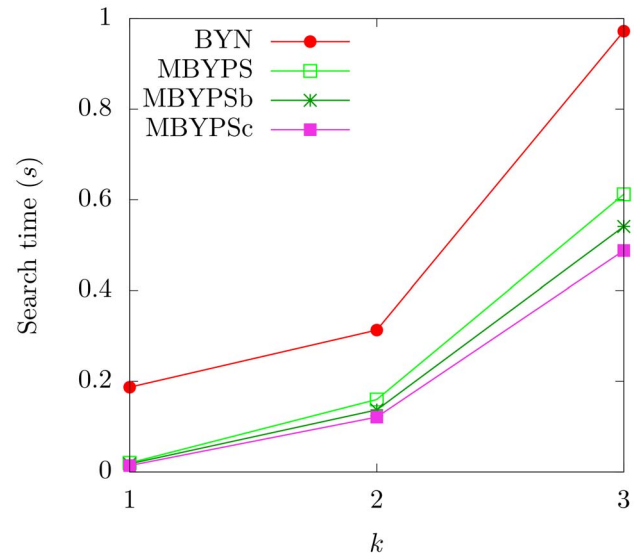
**FIGURE 5.** Number of instructions executed in a run of ANS2 as a function of  $k$  for  $m = 16$  in the DNA alphabet.

(particularly in the English alphabet) while otherwise it was even slower than BYPS. Note the performance improvement of BYPS and its variations respect to original BYP, specially in cases of low difference ratios.

We also tested our algorithms for  $k = 0$  (the times are not shown). BYPSc was the fastest one in this setting. We ran BYPSc against EPSM [31] which is one of the most efficient exact algorithms. BYPSc was only 10–15% slower than EPSM in our experiments.

## 6.2.. Multiple pattern matching with up to $k$ mismatches

We used sets of 10, 100 and 1000 patterns and small  $k$  for testing algorithms for the multiple pattern variation of the problem. Algorithms that were tested are MM: the Muth–Manber algorithm working under Hamming distance [27]; FN: Fredriksson and Navarro’s algorithm [3]; BYN: the original Baeza-Yates–Navarro algorithm which detects occurrences under *edit distance* [28]; HBYN: the Baeza-Yates–Navarro algorithm adapted to Hamming distance;



**FIGURE 6.** Search times of the fastest algorithms tested for multiple pattern matching as a function of  $k$  for  $m = 16$  and  $r = 100$  in the DNA alphabet.

MBYPS, MBYPSb and MBYPSc: Multiple-pattern Baeza-Yates–Perleberg enhanced with SIMD.

Since BYN works under *edit distance*, it has been necessary to modify it to only detect mismatches. We modified its automaton, obtaining a similar one to that used in Shift-Add algorithm. It is worth mentioning that HBYN was never significantly slower than the original BYN. This means HBYN can be fairly compared to other algorithms since it addresses the same problem as them (unlike BYN), and it was not degraded in performance.

As regards the Fredriksson–Navarro algorithm, thorough testing was performed in order to choose the best parameters for each case. For DNA we obtained the same tuning mentioned in [3] as the best configuration.

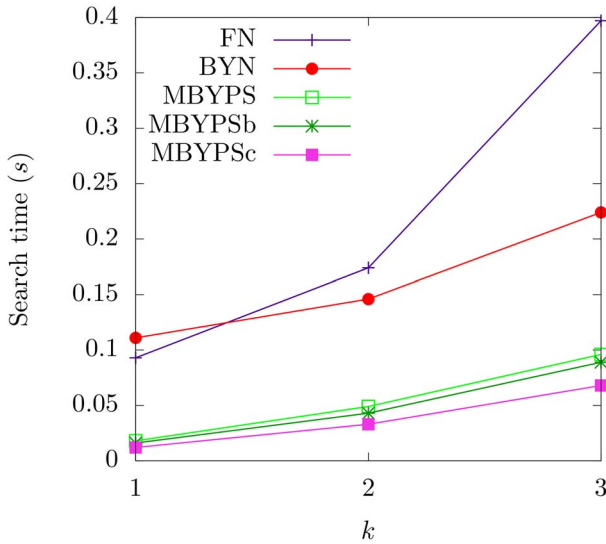
The results are shown in Table 2 with the best times highlighted. In Figures 6 and 7, there are plots of search times of the fastest tested algorithms as a function of  $k$  for  $m = 16$  and  $r = 100$  in both alphabets. Some timings of the mention algorithms have been omitted: MM was designed to work only for  $k = 1$ ; HBYN uses an adaptation of SA, which is limited to small values of  $k$  for long patterns and only works for  $k = 1$  in the case of  $m \in \{24, 32\}$ .

MBYPSb and MBYPSc almost always outperformed all other algorithms, only beaten by MM when  $m = 8$ ,  $k = 1$  and  $r = 1000$  in the English alphabet. In general, there is a larger difference in execution times for lower difference ratios and a larger alphabet (i.e. English). It is worth observing the same performance improvement of SIMD-based algorithms as in the single pattern search test results. For example, for the DNA alphabet and sets of 100 patterns, MBYPSc got an average and a maximum speedup of x88 and x151 compared to MM, x9 and

**TABLE 2** Search times (in seconds) of algorithms for multiple approximate pattern matching with up to  $k$  mismatches.

$k$	$m=8$	$m = 16$			$m = 24$			$m = 32$			$r$	$\Sigma$
	1	1	2	3	1	2	3	1	2	3		
MM	0.105	0.108	$_{-}^a$	$_{-}^a$	0.109	$_{-}^a$	$_{-}^a$	0.115	$_{-}^a$	$_{-}^a$	10	DNA
FN	0.246	0.033	0.249	0.873	0.012	0.015	0.028	0.008	0.010	0.016		
BYN	0.138	0.115	0.132	0.196	0.114	0.131	0.135	0.117	0.126	0.140		
HBYN	0.136	0.115	0.133	0.189	0.113	$_{-}^a$	$_{-}^a$	0.110	$_{-}^a$	$_{-}^a$		
MBYPS	0.043	0.017	0.029	0.087	0.004	0.017	0.022	<b>0.002</b>	0.006	0.018		
MBYPSb	0.038	0.016	0.027	0.076	0.004	0.016	0.021	<b>0.002</b>	0.006	0.016		
MBYPSc	<b>0.028</b>	<b>0.012</b>	<b>0.020</b>	<b>0.060</b>	<b>0.003</b>	<b>0.012</b>	<b>0.016</b>	<b>0.002</b>	<b>0.005</b>	<b>0.013</b>		
MM	0.708	0.738	$_{-}^a$	$_{-}^a$	0.730	$_{-}^a$	$_{-}^a$	0.758	$_{-}^a$	$_{-}^a$	100	
FN	2.215	0.306	1.977	9.952	0.023	0.062	0.164	0.015	0.027	0.065		
BYN	0.510	0.187	0.313	0.972	0.187	0.205	0.278	0.186	0.198	0.220		
HBYN	0.494	0.191	0.308	0.898	0.192	$_{-}^a$	$_{-}^a$	0.188	$_{-}^a$	$_{-}^a$		
MBYPS	0.300	0.020	0.160	0.612	<b>0.005</b>	0.022	0.084	<b>0.003</b>	0.009	0.025		
MBYPSb	0.261	0.018	0.137	0.541	<b>0.005</b>	0.021	0.077	<b>0.003</b>	<b>0.007</b>	0.022		
MBYPSc	<b>0.221</b>	<b>0.014</b>	<b>0.121</b>	<b>0.488</b>	<b>0.005</b>	<b>0.017</b>	<b>0.066</b>	0.005	0.009	<b>0.019</b>		
MM	4.763	4.931	$_{-}^a$	$_{-}^a$	4.985	$_{-}^a$	$_{-}^a$	5.151	$_{-}^a$	$_{-}^a$	1000	
FN	22.283	3.094	21.133	97.438	0.183	0.678	1.999	0.111	0.301	0.668		
BYN	3.855	0.334	1.816	8.742	0.319	0.395	1.244	0.325	0.357	0.435		
HBYN	3.649	0.335	1.720	7.931	0.308	$_{-}^a$	$_{-}^a$	0.324	$_{-}^a$	$_{-}^a$		
MBYPS	2.811	0.054	1.526	6.994	0.018	0.077	0.766	0.014	0.034	0.102		
MBYPSb	2.463	0.049	1.438	6.496	<b>0.013</b>	0.070	0.744	<b>0.012</b>	<b>0.024</b>	0.093		
MBYPSc	<b>2.216</b>	<b>0.041</b>	<b>1.313</b>	<b>5.808</b>	0.033	<b>0.065</b>	<b>0.681</b>	0.035	0.059	<b>0.090</b>		
MM	0.053	0.054	$_{-}^a$	$_{-}^a$	0.053	$_{-}^a$	$_{-}^a$	0.052	$_{-}^a$	$_{-}^a$	10	English
FN	0.048	0.016	0.027	0.052	0.010	0.015	0.023	0.006	0.012	0.017		
BYN	0.061	0.041	0.050	0.055	0.036	0.045	0.059	0.032	0.041	0.047		
HBYN	0.061	0.042	0.050	0.056	0.037	$_{-}^a$	$_{-}^a$	0.031	$_{-}^a$	$_{-}^a$		
MBYPS	0.022	0.015	0.016	0.016	0.003	0.015	0.015	<b>0.002</b>	0.005	0.015		
MBYPSb	0.020	0.013	0.015	0.015	0.003	0.013	0.014	<b>0.002</b>	<b>0.004</b>	0.013		
MBYPSc	<b>0.015</b>	<b>0.010</b>	<b>0.012</b>	<b>0.011</b>	<b>0.002</b>	<b>0.010</b>	<b>0.011</b>	<b>0.002</b>	<b>0.004</b>	<b>0.010</b>		
MM	0.068	0.075	$_{-}^a$	$_{-}^a$	0.072	$_{-}^a$	$_{-}^a$	0.066	$_{-}^a$	$_{-}^a$	100	
FN	0.273	0.093	0.174	0.397	0.063	0.095	0.159	0.048	0.069	0.105		
BYN	0.146	0.111	0.146	0.224	0.104	0.115	0.134	0.093	0.109	0.117		
HBYN	0.146	0.110	0.146	0.218	0.094	$_{-}^a$	$_{-}^a$	0.096	$_{-}^a$	$_{-}^a$		
MBYPS	0.063	0.018	0.049	0.096	0.005	0.018	0.032	0.003	<b>0.006</b>	0.019		
MBYPSb	0.054	0.016	0.043	0.089	<b>0.004</b>	0.016	0.029	<b>0.002</b>	<b>0.006</b>	0.017		
MBYPSc	<b>0.040</b>	<b>0.012</b>	<b>0.033</b>	<b>0.068</b>	<b>0.004</b>	<b>0.013</b>	<b>0.024</b>	0.003	<b>0.006</b>	<b>0.014</b>		
MM	<b>0.242</b>	0.272	$_{-}^a$	$_{-}^a$	0.286	$_{-}^a$	$_{-}^a$	0.274	$_{-}^a$	$_{-}^a$	1000	
FN	3.235	0.679	2.299	7.809	0.368	0.910	1.887	0.270	0.613	1.047		
BYN	0.616	0.208	0.563	1.512	0.186	0.238	0.413	0.178	0.212	0.258		
HBYN	0.602	0.203	0.541	1.436	0.182	$_{-}^a$	$_{-}^a$	0.174	$_{-}^a$	$_{-}^a$		
MBYPS	0.372	0.037	0.332	0.920	0.014	0.048	0.192	0.010	0.025	0.060		
MBYPSb	0.327	0.032	0.299	0.823	<b>0.011</b>	<b>0.041</b>	0.162	<b>0.008</b>	<b>0.019</b>	<b>0.051</b>		
MBYPSc	0.251	<b>0.029</b>	<b>0.225</b>	<b>0.688</b>	0.024	0.045	<b>0.152</b>	0.023	0.039	0.061		

<sup>a</sup> Algorithm not designed to work in this case.



**FIGURE 7.** Search times of the fastest algorithms tested for multiple pattern matching as a function of  $k$  for  $m = 16$  and  $r = 100$  in the English alphabet.

x22 compared to FN, and x16 and x38 compared to HBYN. As regards the English alphabet with  $r = 100$ , MBYPSc average and maximum speed-ups are x12 and x22 compared to MM, x9 and x16 compared to FN, and x12 and x32 compared to HBYN.

MBYPSc was more dominant than MBYPSb, as opposed to the case of BYPSc compared to BYPSb, although it sometimes fell behind it in cases of large  $r$  and a large subpattern length ( $l \geq 8$ ), which corresponds to low difference ratios because  $l$  is  $\lfloor m/(k+1) \rfloor$ .

As regards other algorithms, it is worth mentioning that: (i) Muth–Manber’s algorithm remained competitive only for large sets of small patterns in the English alphabet, specifically for  $r = 1000$  and  $m = 8$  in our tests, showing its tolerance for large sets of patterns; (ii) Fredriksson–Navarro’s algorithm sometimes behaved similarly to MBYPS in the DNA alphabet for small sets of long patterns, although it never outperformed the latter; (iii) Baeza–Yates–Navarro’s algorithm adapted to Hamming distance was always surpassed by MBYPS or its variations, but it was the fastest non-SIMD-based algorithm in many cases of high difference ratios or large sets of patterns.

A curious case worth analyzing is present when an algorithm was run for  $r_1$  and  $r_2$  where  $r_2 > r_1$ , and keeping the same  $m$ ,  $k$  and alphabet, the algorithm yielded search times  $e_1$  and  $e_2$ , respectively, such that  $e_2 > \frac{r_2}{r_1} e_1$ . For example, in Table 2, we can see that MBYPSc took 0.488 s to search pattern occurrences for  $m = 16$ ,  $k = 3$  and  $r = 100$  in the DNA alphabet, whereas it needed 5.808 s to perform the same task for an  $r = 1000$ , around twelve times its previous time.

Two different techniques to address this issue have been used in FN in [3] which are called *pattern grouping* and *pattern clustering*. The former simply consists of dividing the pattern

set into smaller sets and search for them separately. The latter is an extension of the first one where the way subsets are divided are beneficial for the searching mechanism, spending more time in the preprocessing phase. Unfortunately, these improvements would not work for MBYPS or its variations, as the reason for its speed decrement relies on its filtering technique. When the number of subpatterns to search starts to grow, the same subpatterns begin to appear in many patterns. So finding them in the text would actually not filter much, since we still need to verify an approximate occurrence of all patterns that contain them. In conclusion, applying pattern grouping in any of these algorithms would only make them go through the text many times, just increasing their search time.

## 7.. CONCLUDING REMARKS

We have demonstrated that simple SIMD solutions are competitive in searching for approximate single pattern matches within the Hamming distance for patterns of at most 32 characters. Our experiments show that SIMD-based algorithms are the fastest options for small number of mismatches ( $k \leq 3$ ) and sets of at most 1000 patterns in the DNA and English alphabets using real-life texts. It is important to note that although SIMD-based algorithms presented in this paper have worse theoretical time complexities than classic algorithms, they perform better in practice.

In Sections 4.1 and 4.2, we showed that the algorithms for naive counting of mismatches can be used as a checking method for single pattern filtration algorithms. Meanwhile, the fingerprint calculation of a filtration method can be replaced with the CRC fingerprint technique of Section 3.2.

We have also presented an effective way of using the SIMD techniques for approximate multiple string matching in Section 5. The resulting algorithm is substantially faster than the previous most competitive algorithm across multiple alphabets.

In the future, we will modify our algorithms to work with AVX-512. It may be possible to achieve better speed-ups because compare and mask instructions have been merged into one operation. In addition, we will apply SIMD instruction to algorithms for approximate string matching under the edit distance allowing insertions and deletions.

## 7.. DATA AVAILABILITY

No new data were generated or analyzed in support of this research. All source code used in the experiments is available at <https://github.com/ffiori/hamming>.

## REFERENCES

- [1] Navarro, G. (2001) A guided tour to approximate string matching. *ACM Comput. Surv.*, 33, 31–88.

- [2] Intel. Intel (R) 64 and IA-32 Architectures Software Developer's Manual. <https://software.intel.com/content/www/us/en/develop/articles/intel-sdm.html>. 8 August 2020.
- [3] Fredriksson, K. and Navarro, G. (2005) Average-optimal single and multiple approximate string matching. *ACM J. Exp. Algorithmics*, 9, 1. 4.
- [4] Kumar, S. and Spafford, E. (1994) A pattern-matching model for intrusion detection. *Proc. 17th National Computer Security Conference*, 1994, pp. 11–21, NIST, Baltimore, MD.
- [5] Kukich, K. (1992) Techniques for automatically correcting words in text. *ACM Comput. Surv.*, 24, 377–439.
- [6] Dixon, N. and Martin, T. (1979) *Automatic Speech and Speaker Recognition*. IEEE Press, Piscataway, NJ.
- [7] Elliman, D.G. and Lancaster, I.T. (1990) A review of segmentation and contextual analysis techniques for text recognition. *Pattern Recognit.*, 23, 337–346.
- [8] Baeza-Yates, R.A. and Ribeiro-Neto, B.A. (1999) *Modern Information Retrieval*. ACM Press, New York, NY.
- [9] Faro, S. and Pappalardo, E. (2010) Ant-CSP: An ant colony optimization algorithm for the closest string problem. In van Leeuwen, J., Muscholl, A., Peleg, D., Pokorný, J., Rumpe, B. (eds) *SOFSEM 2010: Theory and Practice of Computer Science, 36th Conference on Current Trends in Theory and Practice of Computer Science, Spindleruv Mlýn, Czech Republic, January 23–29, 2010. Proceedings*, Lecture Notes in Computer Science (Vol. 5901), pp. 370–381. Springer, Berlin.
- [10] Gog, S., Karhu, K., Kärkkäinen, J., Mäkinen, V. and Välimäki, N. (2012) Multi-pattern matching with bidirectional indexes. In Gudmundsson, J., Mestre, J., Viglas, T. (eds) *Computing and Combinatorics - 18th Annual International Conference, COCOON 2012, Sydney, Australia, August 20–22, 2012. Proceedings, Lecture Notes in Computer Science* (Vol. 7434), pp. 384–395. Springer, Berlin.
- [11] Baeza-Yates, R.A. and Perleberg, C.H. (1996) Fast and practical approximate string matching. *Inf. Process. Lett.*, 59, 21–27.
- [12] Baeza-Yates, R.A. and Gonnet, G.H. (1992) A new approach to text searching. *Commun. ACM*, 35, 74–82.
- [13] Grabowski, S. and Fredriksson, K. (2008) Bit-parallel string matching under Hamming distance in  $O(n[m/w])$  worst case time. *Inf. Process. Lett.*, 105, 182–187.
- [14] Durian, B., Chhabra, T., Ghuman, S.S., Hirvola, T., Peltola, H. and Tarhio, J. (2014) Improved two-way bit-parallel search. In Holub, J., Zdárek, J. (eds) *Proceedings of the Prague Stringology Conference 2014, Prague, Czech Republic, September 1–3, 2014*, pp. 71–83. Department of Theoretical Computer Science, Faculty of Information Technology, Czech Technical University in Prague, Prague.
- [15] Tarhio, J. and Ukkonen, E. (1993) Approximate Boyer-Moore string matching. *SIAM J. Comput.*, 22, 243–260.
- [16] Horspool, R.N. (1980) Practical fast searching in strings. *Softw. Pract. Exp.*, 10, 501–506.
- [17] Liu, Z., Chen, X., Borneman, J. and Jiang, T. (2005) A fast algorithm for approximate string matching on gene sequences. In Apostolico, A., Crochemore, M., Park, K. (eds) *Proceedings of Combinatorial Pattern Matching, 16th Annual Symposium, CPM 2005, Jeju Island, Korea, June 19–22, 2005, Lecture Notes in Computer Science* (Vol. 3537), pp. 79–90. Springer, Berlin.
- [18] Salmela, L., Tarhio, J. and Kalsi, P. (2010) Approximate Boyer-Moore string matching for small alphabets. *Algorithmica*, 58, 591–609.
- [19] Navarro, G. and Raffinot, M. (2000) Fast and flexible string matching by combining bit-parallelism and suffix automata. *ACM J. Exp. Algorithmics*, 5, 4.
- [20] Ukkonen, E. (1985) Finding approximate patterns in strings. *J. Algorithms*, 6, 132–137.
- [21] Galil, Z. and Giancarlo, R. (1986) Improved string matching with  $k$  mismatches. *SIGACT News*, 17, 52–54.
- [22] Landau, G.M. and Vishkin, U. (1986) Efficient string matching with  $k$  mismatches. *Theor. Comput. Sci.*, 43, 239–249.
- [23] Abrahamson, K.R. (1987) Generalized string matching. *SIAM J. Comput.*, 16, 1039–1051.
- [24] Amir, A., Lewenstein, M. and Porat, E. (2004) Faster algorithms for string matching with  $k$  mismatches. *J. Algorithms*, 50, 257–275.
- [25] Fredriksson, K. and Grabowski, S. (2013) Exploiting word-level parallelism for fast convolutions and their applications in approximate string matching. *Eur. J. Comb.*, 34, 38–51.
- [26] Clifford, R., Fontaine, A., Porat, E., Sach, B. and Starikovskaya, T. (2016) The  $k$ -mismatch problem revisited. In Krauthgamer, R. (ed) *Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2016, Arlington, VA, USA, January 10–12, 2016*, SIAM, pp. 2039–2052.
- [27] Muth, R. and Manber, U. (1996) Approximate multiple strings search. In Hirschberg, D.S., Myers, E.W. (eds) *Proceedings of Combinatorial Pattern Matching, 7th Annual Symposium, CPM 96, Laguna Beach, California, USA, June 10–12, 1996, Lecture Notes in Computer Science* (Vol. 1075), pp. 75–86. Springer, Berlin.
- [28] Baeza-Yates, R. and Navarro, G. (2002) New and faster filters for multiple approximate string matching. *Random Struct Algorithms*, 20, 23–49.
- [29] Sunday, D. (1990) A very fast substring search algorithm. *Commun. ACM*, 33, 132–142.
- [30] Chhabra, T., Faro, S., Külekci, M.O. and Tarhio, J. (2017) Engineering order-preserving pattern matching with SIMD parallelism. *Softw. Pract. Exp.*, 47, 731–739.
- [31] Faro, S. and Külekci, M.O. (2013) Fast packed string matching for short patterns. In Sanders, P., Zeh, N. (eds) *Proceedings of the 15th Meeting on Algorithm Engineering and Experiments, ALENEX 2013, New Orleans, Louisiana, USA, January 7, 2013*, SIAM, pp. 113–121.
- [32] Külekci, M.O. (2009) Filter based fast matching of long patterns by using SIMD instructions. In Holub, J., Zdárek, J. (eds) *Proceedings of the Prague Stringology Conference 2009, Prague, Czech Republic, August 31 – September 2, 2009*, pp. 118–128. Prague Stringology Club, Department of Computer Science and Engineering, Faculty of Electrical Engineering, Czech Technical University in Prague, Prague.
- [33] Ladra, S., Pedreira, O., Duato, J. and Brisaboa, N.R. (2012) Exploiting SIMD instructions in current processors to improve classical string algorithms. In Morzy, T., Härder, T., Wrembel, R. (eds) *Advances in Databases and Information Systems – 16th East European Conference, ADBIS 2012, Poznań, Poland, September 18–21, 2012. Proceedings, Lecture Notes in Computer Science* (Vol. 7503), pp. 254–267. Springer, Berlin.



- [34] Tarhio, J., Holub, J. and Giaquinta, E. (2017) Technology beats algorithms (in exact string matching). *Softw. Pract. Exp.*, 47, 1877–1885.
- [35] Faro, S. and Külekci, M.O. (2013) Towards a very fast multiple string matching algorithm for short patterns. In Holub, J., Zdárek, J. (eds) *Proceedings of the Prague Stringology Conference 2013, Prague, Czech Republic, September 2–4, 2013*, pp. 78–91. Department of Theoretical Computer Science, Faculty of Information Technology, Czech Technical University in Prague, Prague.
- [36] Hirvola, T. (2016) Bit-parallel approximate string matching under Hamming distance. M.Sc. thesis, Aalto University, Espoo, Finland.
- [37] Hume, A. and Sunday, D. (1991) Fast string searching. *Softw. Pract. Exp.*, 21, 1221–1248.