

Optimal Packed String Matching

Oren Ben-Kiki* Philip Bille† Dany Breslauer‡ Leszek Gasieniec§
Roberto Grossi¶ Oren Weimann||

November 20, 2011

Abstract

In the packed string matching problem it is assumed that each machine word can accommodate up to α characters, thus an n -character text occupies n/α memory words. We extend the Crochemore-Perrin constant-space $O(n)$ -time string matching algorithm to run in optimal $O(n/\alpha)$ time and even in real-time, achieving a factor α speedup over traditional algorithms that examine each character individually. Benchmarks show that our solution can be efficiently implemented, unlike prior theoretical packed string matching work.

Our algorithm only uses the standard AC^0 instructions of the RAM model plus two specialized AC^0 word-size packed string instructions, i.e. no integer multiplication. The main *string-matching* instruction WSSM is available in Intel's commodity processors, in SSE4.2/AVX Advanced String Operations. The other *maximal-suffix* instruction WSLM is required only during the pattern preprocessing. In the absence of these specialized packed string instructions, we propose alternative theoretically-efficient emulations: using the four Russian's table lookup technique to emulate WSLM and using parallel algorithms techniques to emulate WSSM via standard integer multiplication instructions (integer multiplication is not AC^0 , but is typically highly optimized).

1 Introduction

Hundreds of articles have been published about string matching, exploring the multitude of theoretical and practical facets of this fundamental problem. For an n -character text T

*Intel Research and Development Center, Haifa, Israel.

†Technical University of Denmark, Copenhagen, Denmark.

‡Caesarea Rothchild Institute for Interdisciplinary Applications of Computer Science, University of Haifa, Haifa, Israel. Partially supported by the European Research Council (ERC) Project SFEROT and by the Israeli Science Foundation Grant 686/07, 347/09 and 864/11.

§University of Liverpool, Liverpool, United Kingdom.

¶Dipartimento di Informatica, Università di Pisa, Pisa, Italy. Partially supported by Italian project PRIN AlgoDEEP (2008TFBWL4) of MIUR.

||Computer Science Department, University of Haifa, Haifa, Israel.

and an m -character pattern x , the classical algorithm by Knuth, Morris and Pratt [30] takes $O(n + m)$ time and uses $O(m)$ auxiliary space to find all pattern occurrences in the text, namely, all text positions i , such that $T[i..i + m - 1] = x$. Many other algorithms have been published; some are faster on the average, use only constant auxiliary space, operate in real-time, or have other interesting benefits. In an extensive study, Faro and Lecroq [20] offer an experimental comparative evaluation of some 85 string matching algorithms.

Packed strings. In modern computers, the size of a machine word is typically larger than the size of an alphabet character and the machine level instructions operate on whole words, i.e., 64-bit or longer words vs. 8-bit ASCII, 16-bit UCS, 2-bits biological DNA, 5-bits amino acid alphabets, etc. The *packed string representation* fits multiple characters into one larger word, so that the characters can be compared in bulk rather than individually: if the characters of a string are drawn from an alphabet Σ , then a word of $\omega \geq \log_2 n$ bits fits up to α characters, where the packing factor is $\alpha = \frac{\omega}{\log_2 |\Sigma|} \geq \log_{|\Sigma|} n$. Throughout the paper, we assume that $|\Sigma|$ is a power of two, ω is divisible by $\log_2 |\Sigma|$, and the packing factor α is a whole integer.

Using the packed string representation in the string matching problem is not a new idea and goes back to early string matching papers by Knuth, Morris and Pratt [30, §4] and Boyer and Moore [11, §8-9], to times when hardware character byte addressing was new and often less efficient than word addressing. Since then, several practical solutions that take advantage of the packed string representation have been proposed in the literature [7, 9, 19, 23, 24, 34, 35]. However, none of these algorithms improves over the worst-case $O(n)$ time bounds of the traditional algorithms. On the other hand, any string matching algorithm should take at least $\Omega(n/\alpha)$ time to read a packed text in the worst case, so there remains a gap to fill. Note that on the average, it is not even required to examine all the text characters [11, 37].

Existing work. A significant theoretical step recently taken introduces a few solutions based on either tabulation (a.k.a. “the Four-Russian technique”) or word-level parallelism (a.k.a. “bit-parallelism”). Fredriksson [23, 24] used tabulation and obtained an algorithm that uses $O(n^\varepsilon m)$ space and $O(\frac{n}{\log_{|\Sigma|} n} + n^\varepsilon m + occ)$ time, where occ denotes the number of pattern occurrences and $\varepsilon > 0$ denotes an arbitrary small constant. Bille [10] improved these bounds to $O(n^\varepsilon + m)$ space and $O(\frac{n}{\log_{|\Sigma|} n} + m + occ)$ time. Very recently, Belazzougui [8] showed how to use word-level parallelism to obtain $O(m)$ space and $O(\frac{n}{m} + \frac{n}{\alpha} + m + occ)$ time. Belazzougui’s algorithm uses a number of succinct data structures as well as hashing: for $\alpha \leq m \leq n/\alpha$, his time bound is optimal while space occupancy is not. As reported by the above authors, none of these results is practical. A summary of the known bounds and our new results is given in Table 1, where our result uses two specialized packed string instructions WSSM and WSLM which are described later on.

Time	Space	Reference
$O(\frac{n}{\log_{ \Sigma } n} + n^\varepsilon m + occ)$	$O(n^\varepsilon m)$	Fredriksson [23, 24]
$O(\frac{n}{\log_{ \Sigma } n} + m + occ)$	$O(n^\varepsilon + m)$	Bille [10]
$O(\frac{n}{\alpha} + \frac{n}{m} + m + occ)$	$O(m)$	Belazzougui [8]
$O(\frac{n}{\alpha} + \frac{m}{\alpha} + occ)$	$O(1)$	using WSSM and WSLM

Table 1: Comparison of packed string matching algorithms.

Our results. We propose an $O(n/\alpha + m/\alpha)$ time string matching algorithm (where the term m/α is kept for comparison with the other results) that is derived from the elegant Crochemore-Perrin [16] algorithm. The latter algorithm takes linear time, uses only constant auxiliary space, and can be implemented in real-time following the recent work by Breslauer, Grossi and Mignosi [12] – benefits that are also enjoyed in our settings. The algorithm has an attractive property that it compares the text characters only moving forward on two wavefronts without ever having to back up, relying on the celebrated Critical Factorization Theorem [13, 31].

We use a *specialized word-size packed string matching instruction* WSSM to anchor the pattern in the text and continue with bulk character comparisons that match the remainder of the pattern. Our reliance on a specialized packed string matching instruction is not far fetched, given the recent availability of such instructions in commodity processors, which has been a catalyst for our work. Our algorithm is easily adaptable to situations where the packed string matching instruction and the bulk character comparison instruction operate on different word sizes. The output occurrences are compactly provided in a bit-mask that can be spelled out as an extensive list of text positions in extra $O(occ)$ time.

Unlike the prior theoretical work, our solution has a cache-friendly sequential memory access without using large external tables or succinct data structures, and therefore, can also be efficiently implemented. The same specialized packed string matching instruction could also be used in other string matching and string processing algorithms, e.g. the Knuth-Morris-Pratt algorithm [28, §10.3.3], but our algorithm has the additional advantages that it also works in real-time and uses only constant auxiliary space. Nonetheless, we expect that algorithms design using non-standard specialized instructions and non-standard models of computation is going to continue and evolve into an exploding area of future research, where the available instructions and the algorithmic design work will cross fertilize.

Model of computation. We adapt the standard word-RAM model with ω -bit words and with only AC^0 instructions (i.e., arithmetic, bitwise and shift operations but no integer multiplication) plus two other specialized AC^0 instructions. The main word-size packed string matching instruction is available in the recent *Advanced String Operations* in *Intel’s Streaming SIMD Extension (SSE4.2)* and *Advanced Vector Extension (AVX) Efficient Accelerated String and Text Processing* instruction set [27, 29]. The other instruction, which is only used in the pattern preprocessing, finds the lexicographically maximum suffix. We adopt the

notation $[d] = \{0, 1, \dots, d - 1\}$. The two specialized instructions are the following:

Word-Size String Matching (wssm): find occurrences of one short pattern x that fits in one word (up to α characters) in a text y that fits in two words (up to $2\alpha - 1$ characters). The output is a binary word Z of $2\alpha - 1$ bits such that its i th bit $Z[i] = 1$ iff $y[i..i + |x| - 1] = x$, for $i \in [2\alpha - 1]$. When $i + |x| - 1 \geq \alpha$, this means that only a prefix of x is matched.

Word-Size Lexicographically Maximum Suffix (wslm): given a packed string x that fits in one word (up to α characters), return position $i \in [\alpha]$ such that $x[i..\alpha - 1]$ is lexicographically maximum among the suffixes in $\{x[j..\alpha - 1] \mid j \in [\alpha]\}$.

Specialized instruction emulation. If these two specialized packed string instructions are not available, then we can emulate them, but our proposed emulations cause a small slowdown as shown in Table 2. While the four Russians’ table lookup technique can be used to emulate either of the two specialized instruction, its space use often makes it impractical and limits the packing factor to $\alpha \leq \log_{|\Sigma|} n$. For the WSSM instruction, we offer a better bit-parallel emulation using the standard word RAM instructions that is built upon techniques that were developed for the *Parallel Random Access Machine* model [15, 36].

Time	Space	Emulation
$O(\frac{n \log \alpha}{\alpha} + occ)$	$O(1)$	bit-parallel WSSM no pre-processing
$O(\frac{n}{\alpha} + \alpha + occ)$	$O(\alpha)$	bit-parallel WSSM pre-processing
$O(\frac{m}{\log_{ \Sigma } n})$	$O(n^\epsilon)$	four Russian WSLM table lookup

Table 2: Word-RAM emulaton of ω -bit WSSM and WSLM instructions.

In Section 2, we start with the reduction of the packed string matching problem using the Crochemore-Perrin algorithm to the two specialized packed string instructions WSSM and WSLM. We then show how to emulate these two specialized packed string instructions in the standard RAM model in Section 3 and report on some experimental results with the WSSM instruction on contemporary processors in Section 4. Conclusions and open problems are given in Section 5.

2 Packed String Matching

In this section we describe how to solve the *packed string matching* problem using the two specialized word-size string matching instructions WSSM and WSLM, and standard word-RAM bulk comparisons of packed strings.

Theorem 2.1 *Packed string matching for a length m pattern and a length n text can be solved in $O(\frac{m}{\alpha} + \frac{n}{\alpha})$ time in the word-RAM extended with constant-time WSSM and WSLM instructions. Listing explicitly the occ text positions of the pattern occurrences takes an additional $O(occ)$ time. The algorithm can be made real-time, and uses just $O(1)$ auxiliary words of memory besides the read-only $\frac{m}{\alpha} + \frac{n}{\alpha}$ words that store the input.*

The algorithm behind Theorem 2.1 follows the classical scheme, in which a text scanning phase is run after the pattern preprocessing. In the following, we first present the necessary background and then describe how to perform the text scanning phase using WSSM, and the pattern preprocessing using WSLM.

2.1 Background

Critical Factorization. Properties of periodic strings are often used in efficient string algorithms. A string u is a *period* of a string x if x is a prefix of u^k for some integer k , or equivalently if x is a prefix of ux . The shortest period of x is called *the period* of x and its length is denoted by $\pi(x)$. A *substring* or a *factor* of a string x is a contiguous block of symbols u , such that $x = x'ux''$ for two strings x' and x'' . A *factorization* of x is a way to break x into a number of factors. We consider factorizations of a string $x = uv$ into two factors: a *prefix* u and a *suffix* v . Such a factorization can be represented by a single integer and is *non-trivial* if neither of the two factors is equal to the empty string.

Given a factorization $x = uv$, a *local period* of the factorization is defined as a non-empty string p that is consistent with both sides u and v . Namely, (i) p is a suffix of u or u is a suffix of p , and (ii) p is a prefix of v or v is a prefix of p . The shortest local period of a factorization is called *the local period* and its length is denoted by $\mu(u, v)$. A non-trivial factorization $x = uv$ is called a *critical factorization* if the local period of the factorization is of the same length as the period of x , i.e., $\mu(u, v) = \pi(uv)$. See Figure 1.

$\begin{array}{c} \text{a} \mid \text{b a a a b a} \\ \text{b a} \quad \text{b a} \\ \text{(a)} \end{array}$	$\begin{array}{c} \text{a b} \mid \text{a a a b a} \\ \text{a a a b} \quad \text{a a a b} \\ \text{(b)} \end{array}$	$\begin{array}{c} \text{a b a} \mid \text{a a b a} \\ \text{a} \quad \text{a} \\ \text{(c)} \end{array}$
--	--	--

Figure 1: The local periods at the first three non-trivial factorizations of the string `abaaaba`. In some cases the local period overflows on either side; this happens when the local period is longer than either of the two factors. The factorization (b) is a critical factorization with local period `aaab` of the same length as the global period `abaa`.

Crochemore-Perrin algorithm. We assume that the reader is familiar with the Crochemore-Perrin algorithm [16] and its real-time variation Breslauer-Grossi-Mignosi [12] and briefly review these algorithms. Recall that Crochemore and Perrin use Theorem 2.2 to break up the pattern as critical factorization $x = uv$ with non-empty prefix u and suffix v , such that $|u| \leq \pi(x)$.

Theorem 2.2 (*Critical Factorization Theorem, Cesari and Vincent [13, 31]*) *Given any $|\pi(x)| - 1$ consecutive non-trivial factorizations of a string x , at least one is critical.*

Then, they exploit the critical factorization of $x = uv$ by matching the longest prefix z of v against the current text symbols, and using Theorem 2.3 whenever a mismatch is found.

Theorem 2.3 (Crochemore and Perrin [16]) *Let $x = uv$ be a critical factorization of the pattern and let p be any local period at this factorization, such that $|p| \leq \max(|u|, |v|)$. Then $|p|$ is a multiple of $\pi(x)$, the period length of the pattern.*

Precisely, if $z = v$, they show how to declare an occurrence of x . Otherwise, the symbol following z in v is mismatching when compared to the corresponding text symbol, and the pattern x can be *safely* shifted by $|z| + 1$ positions to the right (there are other issues for which we refer the reader to [16]).

To simplify the matter in the rest of the paper, we discuss how to match the pattern suffix v assuming without loss of generality that $|u| \leq |v|$. Indeed, if $|u| > |v|$, the Crochemore-Perrin approach can be simplified as shown in [12]: use two critical factorizations, $x = uv$ and $x' = u'v'$, for a prefix x' of x such that $|x'| > |u|$ and $|u'| \leq |v'|$. In this way, matching both u' and v' suitably displaced by $|x| - |x'|$ positions from matching v , guarantees that x occurs. This fact enables us to focus on matching v and v' , since the cost of matching u' is always dominated by the cost of matching v' , and we do not need to match u . For the sake of discussion, it suffices to consider only one instance, namely, suffix v .

We now give more details on the text processing phase, assuming that the pattern preprocessing phase has correctly found the critical factorization of the pattern x and its period $\pi(x)$, and any additional pattern preprocessing that may be required (Section 2.3).

While other algorithms may be used with the WSSM instruction, the Crochemore-Perrin algorithm is particularly attractive because of its simple text processing. Therefore, it is convenient to assume that the period length and critical factorization are exactly computed in the pattern preprocessing burying the less elegant parts in that phase.

2.2 Text processing

The text processing has complementary parts that handle short patterns and long patterns. A pattern x is *short* if its length is at most α , namely, the packed pattern fits into a single word, and is *long* otherwise. Processing short patterns is immediate with WSSM and, as we shall see, the search for long patterns reduces to that for short patterns.

Short patterns. When the pattern is already short, WSSM is repeatedly used to directly find all occurrences of the pattern in the text.

Lemma 2.4 *There exists an algorithm that finds all occurrences of a short pattern of length $m \leq \alpha$ in a text of length n in $O(\frac{n}{\alpha})$ time using $O(1)$ auxiliary space.*

Proof. Consider the packed text blocks of length $\alpha + m - 1$ that start on word boundaries, where each block overlaps the last $m - 1$ characters of the previous block and the last block might be shorter. Each occurrence of the pattern in the text is contained in exactly one such block. Repeatedly use the WSSM instruction to search for the pattern of length $m \leq \alpha$ in these text blocks whose length is at most $\alpha + m - 1 \leq 2\alpha - 1$. \square

Long patterns. Let x be a long pattern of length $m > \alpha$: occurrences of the pattern in the text must always be spaced at least the period $\pi(x)$ locations apart. We first consider the easier case where the pattern has a long period, namely $m \geq \pi(x) > \alpha$, and so there is at most one occurrence starting within each word.

Lemma 2.5 *There exists an algorithm that finds all occurrences of a long-period long pattern of length $m \geq \pi(x) \geq \alpha$, in a text of length n in $O(\frac{n}{\alpha})$ time using $O(1)$ auxiliary space.*

Proof. The Crochemore-Perrin algorithm can be naturally implemented using the WSSM instruction and bulk character comparisons. Given the critical factorization $x = uv$, the algorithm repeatedly searches using WSSM for an occurrence of a prefix of v of length $\min(|v|, \alpha)$ starting in each packed word aligned with v , until such an occurrence is discovered. If more than one occurrence is found starting within the same word, then by Lemma 2.3, only the first such occurrence is of interest. The algorithm then uses the occurrence of the prefix of v to anchor the pattern within the text and continues to compare the rest of v with the aligned text and then compares the pattern prefix u , both using bulk comparison of words containing α packed characters. Bulk comparisons are done by comparing words; in case of a mismatch the mismatch position can be identified using bitwise XOR operation, and then finding the most significant set bit.

A mismatch during the attempt to verify the suffix v allows the algorithm to shift the pattern ahead until v is aligned with the text after the mismatch. A mismatch during the attempt to verify u , or after successfully matching u , causes the algorithm to shift the pattern ahead by $\pi(x)$ location. In either case the time adds up to only $O(\frac{n}{\alpha})$. \square

When the period of the pattern is shorter than the word size, that is $\pi(x) \leq \alpha$, there may be several occurrences of the pattern starting within each word. The algorithm is very similar to the long period algorithm above, but with special care to efficiently manipulate the bit-masks representing all the occurrences.

Lemma 2.6 *There exists an algorithm that finds all occurrences of a short-period long pattern of length m , such that $m > \alpha > \pi(x)$, in a text of length n in $O(\frac{n}{\alpha})$ time using $O(1)$ auxiliary space.*

Proof. Let p be the prefix of x of length $\pi(x)$, and write $x = p^r p'$, where p' is a prefix of p . If we can find the maximal runs of consecutive ps inside the text, then it is easy to locate the occurrences of x . To this end, let $k \leq r$ be the maximum positive integer such that $k \cdot \pi(x) \leq \alpha$ while $(k + 1) \cdot \pi(x) > \alpha$. Note that there cannot exist two occurrences of p^k that are completely inside the same word.

We examine one word w of the text at a time while maintaining the current run of consecutive ps spanning the text word w' preceding w . We apply WSSM to p^k and $w'w$, and take the rightmost occurrence of p^k whose matching substring is completely inside $w'w$. We have two cases: either that occurrence exists and is aligned with the current run of ps , and so we extend it, or we close the current run and check whether p' occurs soon after. The

latter case arises when there is no such an occurrence of p^k , or it exists but is not aligned with the current run of ps . Once all the maximal runs of consecutive occurrences of ps are found (some of them are terminated by p') for the current word w , we can decide by simple arithmetics whether $x = p^r p'$ occurs on the fly. \square

Real-time algorithm. As mentioned in Section 2.1, the Crochemore-Perrin algorithm can be implemented in real time using two instances of the basic algorithm with carefully chosen critical factorizations [12]. Since we are following the same scheme here, our algorithm reports the output bit-mask of pattern occurrences ending in each text word in $O(1)$ time after reading the word. Thus, we can obtain a real-time version as claimed in Theorem 2.1.

2.3 Pattern preprocessing

Given the pattern x , the pattern preprocessing of Crochemore-Perrin produces the period length $\pi(x)$ and a critical factorization $x = uv$ (Section 2.1): for the latter, they show that v is the lexicographically maximum suffix in the pattern under either the regular alphabet order or its inverse order, and use the algorithm by Duval [18]. The pattern preprocessing of Breslauer, Grossi and Mignosi [12] uses Crochemore-Perrin preprocessing, and it also requires to find the prefix x' of x such that $|x'| > |u|$ and its critical factorization $x' = u'v'$ where $|u'| \leq |v'|$. Our pattern preprocessing requires to find the period π' for the first α characters in v (resp., those in v'), along with the longest prefix of v (resp., v') having that period. We thus end up with only the following two problems:

1. Given a string x , find its lexicographically maximum suffix v (under the regular alphabet order or its inverse order).
2. Given a string $x = uv$, find its period $\pi(x)$ and the period of a prefix of v .

When $m = O(\frac{n}{\alpha})$, which is probably the case in many situations, we can simply run the above algorithms in $O(m)$ time to solve the above two problems. We focus here on the case when $m = \Omega(\frac{n}{\alpha})$, for which we need to give a bound of $O(\frac{m}{\alpha})$ time.

Lemma 2.7 *Given a string x of length m , its lexicographically maximum suffix v can be found in $O(\frac{m}{\alpha})$ time.*

Proof. Duval’s algorithm [18] is an elegant and simple linear-time algorithm that can be easily adapted to find the lexicographically maximum suffix. It maintains two positions i and j , one for the currently best suffix and the other for the current candidate. Whenever there is a mismatch after matching k characters ($x[i+k] \neq x[j+k]$), one position is “defeated” and the next candidate is taken. Its implementation in word-RAM is quite straightforward, by comparing α characters at a time, except when the interval $[\min(i, j), \max(i, j) + k]$ contains less than α positions, and so everything stays in a single word: in this case, we can potentially perform $O(\alpha)$ operations for the $O(\alpha)$ characters (contrarily to the rest, where we perform $O(1)$ operations). We show how to deal with this situation in $O(1)$ time. We employ WSLM,

and let w be the suffix thus identified in the word. We set i to the position of w in the original string x , and j to the first occurrence of w in x after position i (using WSSM). If j does not exist, we return i as the position of the lexicographically maximum suffix; otherwise, we set $k = |w|$ and continue by preserving the invariant of Duval’s algorithm. \square

Lemma 2.8 *The preprocessing of a pattern of length m takes $O(\frac{m}{\alpha})$ time.*

3 Word-Size Instruction Emulation

Our algorithm uses two specialized word-size packed string matching instructions, WSSM and WSLM, that are assumed to take $O(1)$ time. In the circuit complexity sense both are AC^0 instructions, which are easier than integer multiplication that is not AC^0 , since integer multiplication can be used to compute the parity [25]. Recall that the class AC^0 consist of problems that admit polynomial size circuits of depth $O(1)$, with Boolean **and/or** gates of unbounded fan-in and **not** gates only at the inputs.

While either instruction can be emulated using the four Russians’ technique, table lookup limits the packing factor and has limited practical value for two reasons: it sacrifices the constant auxiliary space and has no more cache friendly access. We focus here on the easier and more useful main instruction WSSM and propose efficient bit parallel emulations in the word-RAM, relying on integer multiplication for fast Boolean convolutions.

Lemma 3.1 *After a preprocessing of $O(\omega)$ time, the $\omega/\log \log W$ -bit WSSM and WSLM instructions can be emulated in $O(1)$ time on a ω -bit word RAM.*

3.1 Bit-parallel emulation of wssm

String matching problems under *general matching relations* were classified in [32, 33] into easy and hard problems, where easy problems are equivalent to string matching and are solvable in $O(n + m)$ time, and hard problems are at least as hard as one or more Boolean convolutions, that are solved using *FFT* and integer convolutions in $O(n \log m)$ time [2, 22]. To efficiently emulate the WSSM instruction we introduce *two layers* of increased complexity: first, we observe that the problem can also be solved using Boolean convolutions, and then, we use the powerful, yet standard, integer multiplication operation, that resembles integer convolutions, to emulate Boolean convolutions. In the circuit complexity sense Boolean convolution is AC^0 , and therefore, is easier than integer multiplication.

String matching and boolean convolution via integer multiplication. Consider the text $t = t_0 \cdots t_{n-1}$ and the pattern $p = p_0 \cdots p_{m-1}$ where $p_i, t_i \in \Sigma$. Our goal is to compute the occurrence vector of p in t . This is a vector c so that $c_k = 1$ iff $t_{k+i} = p_i$ for all $i \in [m]$.

Since each character is encoded in $\log_2 |\Sigma|$ bits, we view t and p as *binary* vectors of length $n \log_2 |\Sigma|$ and $m \log_2 |\Sigma|$ respectively. In the occurrence vector c we will then only

consider positions $k = 0, \log_2 |\Sigma|, \log_2 |\Sigma|, \dots$ (all the other positions correspond to partial characters and will be discarded). In general, we have

$$c_k = \bigwedge_{i=0, \dots, m-1} (t_{k+i} = p_i) = \overline{\left(\bigvee_{i=0, \dots, m-1} (t_{k+i} \wedge \overline{p_i}) \right)} \vee \left(\bigvee_{i=0, \dots, m-1} (\overline{t_{k+i}} \wedge p_i) \right).$$

Define the *convolution operator* $a \star b$ for binary vectors $a = a_0 \cdots a_{n-1}$ and $b = b_0 \cdots b_{m-1}$ to be

$$(a \star b)_k = \bigvee_{i=0, \dots, \min\{m-1, n-k-1\}} (a_{i+k} \wedge b_i).$$

The occurrence vector c can then be computed by taking the n least significant bits from $(t \star \overline{p}) \vee (\overline{t} \star p)$. This is illustrated in Figure 2. We now explain how to compute $t \star \overline{p}$, computing $\overline{t} \star p$ is done similarly.

Notice that $(t \star \overline{p})_k = 1$ iff aligning the pattern p to the text t starting at position k has *at least one* mismatch location where t has a 1 and p has a 0. We will instead require that $(t \star \overline{p})_k = x$ iff there are *exactly* x such mismatches. This way, we can compute $t \star \overline{p}$ using standard integer multiplication $t \times \overline{p}$. This is because with the left shift operator \ll we have:

$$t \star \overline{p} = \bigvee_{i=0, \dots, \min\{m-1, n-k-1\}} [(t \ll i) \times \overline{p}_i] = t \times \overline{p}.$$

The only problem with this method is that the number of mismatches x might be as large as m . To account for this, we pad each digit of t and \overline{p} with $L = \lceil \log m \rceil$ zeros, and think of each group of L bits as a *field*. Since we are adding up at most m numbers the fields would not overflow. Thus, performing the integer multiplication on the padded strings gives fields with value x when the number of mismatches is x .

Adding the two convolutions $t \star \overline{p}$ and $\overline{t} \star p$ together, we get the overall number of mismatches, and we need to identify the fields with value 0 (these correspond to occurrences, i.e., no mismatches). In other words, if we use padded vectors $t', \overline{t}', p',$ and \overline{p}' , we can compute $r = (t' \times \overline{p}') + (\overline{t}' \times p')$ and set $c_k = 0$ if and only if the the corresponding field in r is non-zero.

We use the constant time word-RAM bit techniques in Fich [21] to pad and compact. Note that in each field with value f we have that $0 - f$ is either 0, or borrows from the next field 1s on the left side. Take a mask with 1 in each field at the least significant bit, and subtract our integer r from this mask. We get that only zero fields have 0 in their most significant bit. Next, Boolean AND with the mask to keep the most significant bit in each field. Finally, shift right to the least significant bit in the field.

The only caveat in the above “string matching via integer multiplication” is its need for padding, thus extending the involved vectors by a factor of $L = \Theta(\log m) = O(\log w)$. We now have to use L machine words which incurs a slowdown of $\Omega(L)$. We next show how to reduce the required padding from L to $\log \log \alpha$.

Padding the pattern 101 and the text 01101010 (padding bits are in gray) we get that:

$$p = 010001, t = 0001010001000100$$

$$\bar{p} = 000100, \bar{t} = 0100000100010001$$

Doing standard integer multiplication on these vectors we get that:

$$p \times \bar{t} = 1000101001000100001$$

$$\bar{p} \times t = 0000101000100010000$$

Adding these we get the mismatch vector:

$$(p \times \bar{t}) + (\bar{p} \times t) = 1\ 00\ 10\ 10\ 00\ 11\ 00\ 11\ 00\ 01$$

Replacing each field (two bits) by the number it holds gives:

$$(p \times \bar{t}) + (\bar{p} \times t) = 1022030301$$

Taking the $n = 8$ least significant bits gives the mismatch vector 22030301.

Figure 2: An example of searching for the pattern $p = 101$ of length $m = 3$ in the text $t = 01101010$ of length $n = 8$. The mismatch vector is 22030301. i.e., aligning p to t at position 0 gives two mismatches, at position 1 also gives two mismatches, at position 2 there are no mismatches (this is an occurrence) etc.

Sparse convolutions via deterministic samples. A *deterministic sample (DS)* for a pattern with period length π is a collection of at most $\lceil \log \pi \rceil$ pattern positions, such that any two occurrence candidate text locations that match the pattern at the *DS* must be at least π locations apart [36]. To see that a *DS* exists, take π consecutive occurrence candidates. Any two candidates must have at least one mismatch position; add one such position to the *DS* and keep only the remaining minority candidates, removing at least half of the remaining candidates. After at most $\lceil \log \pi \rceil$ iterations, there remains only one candidate and its *DS*. Moreover, if the input characters are expanded into $\log_2 |\Sigma|$ bits, then the *DS* specifies only $\lceil \log \pi \rceil$ bits, rather than characters. Candidates can be eliminated via Boolean convolutions with the two bit vectors representing the 0s and 1s in the *DS*, that is, sparse Boolean vectors with at most $\lceil \log \pi \rceil$ set bits. The period π , the *DS*, and the other required masks and indices are precomputed in $O(\omega)$ time.

Consider now how we performed string matching via integer multiplication in the previous paragraph. Then, the padding in the bitwise convolution construction can be now reduced to only $L' = \lceil \log \log \pi + 1 \rceil$ bits instead of L bits, leading to convolutions of shorter $O(\omega \log \log \pi) = O(\omega \log \log \omega)$ bit words and slowdown of only $O(\log \log \omega)$ time. Using ω -bit words and $O(\omega)$ -time preprocessing, we can treat $O(\omega / \log \log \omega)$ bits in $O(1)$ time using multiplication, thus proving Lemma 3.1.

Preprocessing the pattern The algorithm above requires the period length of the pattern π (see Section 2.3) and certain deterministic samples to solve the string matching problem. We can take the pattern prefix of length $\omega/\log\omega$ and match it against itself and find its period using one multiplication as above. If the period is small, we can verify it and find if and where it terminates. If not, see if there is any value and if possible to work out an algorithm that takes constant operations. This might give a hint for better string matching algorithm.

We can separate out the preprocessing. But first we should solve the string matching and then maybe we can get preprocessing in constant time even.

Observe that even that the pattern preprocessing is done in constant time and can be computed repeatedly, the preprocessing has only to be done once since the WSSM instruction is only used with one pattern string, the whole pattern if the pattern is short ($m \leq \alpha$) or the prefix of the critical factorization tail if the pattern is long.

3.2 Four Russians' table lookup technique

The specialized WSSM and WSLM instruction can both be emulated using the *four Russian's table lookup technique* [6]. The time and space required to create a size n lookup table limit the packing factor α , typically to $\alpha \leq \log_{|\Sigma|} n$. The lookup table size can be reduced further to n^ϵ limiting the packing factor further to $\alpha \leq \epsilon \log_{|\Sigma|} n$ and increasing the time by an ϵ multiplicative factor. In practice, a table lookup implementation would also benefit from small table sizes that fit in the fast cache memory, limiting the size of such lookup tables even further. Thus, using large lookup tables would typically make the algorithm impractical, sacrificing both the small space and the sequential cache friendly access benefits of our algorithm. However, when the alphabet size and the table sizes are small, table lookup may be beneficial.

We outline here only how the four Russian's table lookup technique emulates the WSLM instruction, which is particularly simple. The WSSM bit-parallel emulation above is better than the equivalent table lookup emulation of WSSM. There is some circumstantial evidence that the WSLM emulation is harder than the WSSM emulation, since in the parallel random access machine model the best maximum suffix algorithms take significantly more time than string matching [5, 17, 26].

We create a lookup table $\text{maximum-suffix}[x]$ that gives the length of the lexicographically maximum suffix, for all strings x represented as integers in base $\log_2 |\Sigma|$, with their most significant digit as the first character. Using this definition, short strings are padded with "0" characters on the left and therefore have the same maximum suffix as long strings, with the exception of the 0 entry that has the maximum suffix that is the whole string and its length is the length of the whole string which must be specified separately. The table is created by appending a new first non-zero character on the left of a shorter string and comparing the whole new string lexicographically as one integer to the maximum suffix of the string before adding the first new character using appropriate shift operations. The time to create the maximum-suffix lookup table is clearly linear in its size.

4 wssm on contemporary commodity processors

We conducted benchmarks of the packed string matching instructions that are available in the "Efficient Accelerated String and Text Processing: Advanced String Operations" part of the *Streaming SIMD Extension (SSE4.2) and the Advanced Vector Extension (AVX)* on Intel Sandy Bridge processors [27, 29] and consulted Intel's Optimization Reference Manual [28], both indicate remarkable performance. The instruction *Packed Compare Explicit Length Strings Return Mask (PCMPESTRM, Equal Ordered Aggregation)* produces a bit mask that is suitable for short patterns and the similar but slightly faster instruction *Packed Compare Explicit Length Strings Return Index (PCMPESTRI)* produces only the index of the first occurrence, which is suitable for our longer pattern algorithm.¹ These instructions support WSSM with 8-bit or 16-bit characters and with up to 128-bit long pattern and text strings. There are currently no WSLM equivalent instructions available.

Faro and Lecroq kindly made their extensive *String Matching Algorithms Research Tool (SMART)* available to us [20]. Benchmarks that we put together in SMART show that for up to 8-character long patterns, the raw packed string matching instructions performed among the top algorithms available in SMART. The Crochemore-Perrin algorithm with packed string matching instructions performed very well on many combinations of input types and longer patterns. These preliminary experimental results must be interpreted cautiously, since on one hand we have implemented the benchmarks very quickly, while on the other hand the existing SMART algorithms could benefit as well from packed string matching instructions and from additional handcrafted machine specific optimization; in fact, a handful of the existing SMART algorithms already use other Streaming SIMD Extension instructions.

Results of our SMART benchmarks with roughly 95 SMART algorithms are summarized in Table 3 and Table 4. Algorithm *SSECP* uses the raw Intel instruction for patterns of length 2, 4, and 8 and the Crochemore-Perrin algorithm for longer patterns. In general, the longer is the pattern and the larger is the variation of the alphabet characters, the algorithms that skip parts of the text have a greater advantage. See Faro and Lecroq's [20] paper and the SMART framework implementation for further details on the listed smart algorithms.

5 Conclusions

We demonstrated how to employ word-size string matching instructions to design optimal packed string matching algorithms in the word-RAM, which are fast both in theory and in practice. There is an array of interesting questions that arise from our investigation.

1. Is it possible to improve our WSSM emulation further towards constant time, including any pattern pre-processing, with only ω -bit words? With only AC^0 operations (i.e. no integer multiplication)?

¹On current generation high end Intel Sandy Bridge processors, 2-cycle throughput and 7- or 8-cycle latency [28, §C.3.1]. We did not evaluate new AMD processors that also realize the packed string matching instructions [3, 4]. Implicit length (null terminated) packed string instruction variants are also available.

	2	4	8	16	32	64	128	256
	SSECP 4.44	SSECP 4.57	UFNDMQ4 4.99	BNDMQ4 4.23	BNDMQ4 3.83	LBNDM 3.91	BNDMQ4 3.71	HASH3 3.24
	SKIP 4.80	RF 5.07	SSECP 5.00	SBNDMQ4 4.31	BNDMQ6 3.86	BNDMQ4 3.94	HASH5 3.83	HASH8 3.81
	SO 4.84	BM 5.33	FSBNDM 5.05	UFNDMQ4 4.31	SBNDMQ4 3.95	SBNDMQ4 3.96	HASH8 3.93	HASH5 3.83
	FNDM 4.94	BNDMQ2 5.46	SBNDMQ2 5.08	UFNDMQ6 4.47	SBNDMQ6 3.97	BNDMQ6 3.97	HASH3 3.94	BNDMQ4 3.91
	FSBNDM 5.03	BF 5.58	BNDMQ2 5.13	SBNDMQ6 4.57	UFNDMQ4 4.00	HASH5 3.98	BNDMQ6 3.97	SBNDMQ4 3.95
				23 SSECP 5.00	27 SSECP 5.29	39 SSECP 4.88	42 SSECP 4.73	38 SSECP 4.58

Table 3: Fastest SMART algorithms on SMART English texts.

	2	4	8	16	32	64	128	256
	SSECP 4.28	SSECP 4.49	BNDMQ2 4.42	SBNDMQ2 4.08	UFNDMQ2 3.75	SBNDMQ4 3.67	BNDMQ4 3.70	HASH5 3.68
	FFS 4.88	SVM1 4.84	SBNDMQ2 4.48	UFNDMQ2 4.08	BNDMQ4 3.79	BNDMQ4 3.72	SBNDMQ4 3.71	BNDMQ4 3.69
	GRASPM 4.93	SBNDMQ4 4.85	SBNDM 4.59	SBNDMQ4 4.10	SBNDMQ4 3.80	UFNDMQ4 3.80	HASH5 3.75	HASH8 3.71
	BR 5.14	BOM2 4.95	SBNDM2 4.59	SBNDM2 4.13	UFNDMQ4 3.80	BNDMQ2 3.89	UFNDMQ4 3.77	LBNDM 3.71
	BWW 5.5.14	EBOM 5.25	UFNDMQ2 4.69	BNDMQ2 4.14	BNDMQ2 3.89	SBNDM2 3.96	HASH8 3.80	SBNDMQ4 3.73
			13 SSECP 5.00	22 SSECP 5.08	35 SSECP 4.77	39 SSECP 4.77	45 SSECP 4.76	48 SSECP 4.77

Table 4: Fastest SMART algorithms on SMART random text with 16-character alphabet.

2. Derive Boyer-Moore style algorithms, that skip parts of the text [11, 14, 34, 37] and may be therefore faster on average, using packed string matching instructions.
3. Extend our specialized packed string instruction results to the dictionary matching problem with multiple patterns [1, 8, 14, 34] and to other string matching problems.
4. Find critical factorizations in linear-time using equality pairwise symbol comparisons (i.e. no alphabet order). Such algorithms could also have applications in our packed string model, possibly eliminating our reliance on the WSLM instruction.
5. Further compare the performance of our new algorithm using hardware packed string matching instructions to existing implementations (e.g. Faro and Lecroq’s [20] SMART and the SSE4.2 platform specific *strstr* in *glibc*). Our WSSM emulation may also be useful in practice in the case of very small alphabets, e.g. 2-bit DNA alphabets [35].

References

- [1] A. Aho and M. Corasick. Efficient string matching: An aid to bibliographic search. *Comm. of the ACM*, 18(6):333–340, 1975.
- [2] A. Aho, J. Hopcroft, and J. Ullman. *The Design and Analysis of Computer Algorithms*. Addison-Wesley, Reading, MA, 1974.

- [3] AMD. *AMD64® Architecture Programmers Manual Volume 4: 128-Bit and 256-Bit Media Instructions*. AMD Corporation, 2011.
- [4] AMD. *Software Optimization Guide for AMD Family 15h Processors*. AMD Corporation, 2011.
- [5] A. Apostolico and M. Crochemore. Optimal Canonization of All Substrings of a String. *Inf. Comput.*, 95(1):76–95, 1991.
- [6] V. Arlazarov, E. Dinic, M. Kronrod, and I. Faradzev. On economic construction of the transitive closure of a directed graph. *Soviet Math. Dokl.*, 11:1209–1210, 1970.
- [7] R. A. Baeza-Yates. Improved string searching. *Softw. Pract. Exper.*, 19(3):257–271, 1989.
- [8] D. Belazzougui. Worst Case Efficient Single and Multiple String Matching in the RAM Model. In *Proceedings of the 21st International Workshop On Combinatorial Algorithms (IWOCA)*, pages 90–102, 2010.
- [9] M. Ben-Nissan and S. T. Klein. Accelerating Boyer Moore searches on binary texts. In *Proceedings of the 12th International Conference on implementation and Application of Automata (CIAA)*, pages 130–143, 2007.
- [10] P. Bille. Fast searching in packed strings. *J. Discrete Algorithms*, 9(1):49–56, 2011.
- [11] R. Boyer and J. Moore. A fast string searching algorithm. *Comm. of the ACM*, 20:762–772, 1977.
- [12] D. Breslauer, R. Grossi, and F. Mignosi. Simple Real-Time Constant-Space String Matching. In R. Giancarlo and G. Manzini, editors, *CPM*, volume 6661 of *Lecture Notes in Computer Science*, pages 173–183. Springer, 2011.
- [13] Y. Césari and M. Vincent. Une caractérisation des mots périodiques. *C.R. Acad. Sci. Paris*, 286(A):1175–1177, 1978.
- [14] B. Commentz-Walter. A string matching algorithm fast on the average. In *Proc. 6th International Colloquium on Automata, Languages, and Programming*, Lecture Notes in Computer Science, pages 118–132. Springer-Verlag, Berlin, Germany, 1979.
- [15] M. Crochemore, Z. Galil, L. Gasieniec, K. Park, and W. Rytter. Constant-Time Randomized Parallel String Matching. *SIAM J. Comput.*, 26(4):950–960, 1997.
- [16] M. Crochemore and D. Perrin. Two-way string-matching. *J. ACM*, 38(3):651–675, 1991.
- [17] J. W. Daykin, C. S. Iliopoulos, and W. F. Smyth. Parallel RAM Algorithms for Factoring Words. *Theor. Comput. Sci.*, 127(1):53–67, 1994.

- [18] J. Duval. Factorizing Words over an Ordered Alphabet. *J. Algorithms*, 4:363–381, 1983.
- [19] S. Faro and T. Lecroq. Efficient pattern matching on binary strings. In *Proceedings of the 35th Conference on Current Trends in Theory and Practice of Computer Science (SOFSEM)*, 2009.
- [20] S. Faro and T. Lecroq. The exact string matching problem: a comprehensive experimental evaluation report. Technical Report 0810.2390, arXiv, Cornell University Library, 2011. <http://arxiv.org/abs/1012.2547>.
- [21] F. E. Fich. Constant time operations for words of length w . Technical report, University of Toronto, 1999. <http://www.cs.toronto.edu/~faith/algs.ps>.
- [22] M. Fischer and M. Paterson. String matching and other products. In *Complexity of Computation*, pages 113–125. American Mathematical Society, 1974.
- [23] K. Fredriksson. Faster string matching with super-alphabets. In *Proceedings of the 9th International Symposium on String Processing and Information Retrieval (SPIRE)*, pages 44–57, 2002.
- [24] K. Fredriksson. Shift-or string matching with super-alphabets. *IPL*, 87(4):201–204, 2003.
- [25] M. L. Furst, J. B. Saxe, and M. Sipser. Parity, circuits, and the polynomial-time hierarchy. *Mathematical Systems Theory*, 17(1):13–27, 1984.
- [26] C. S. Iliopoulos and W. F. Smyth. Optimal Algorithms for Computing the canonical form of a circular string. *Theor. Comput. Sci.*, 92(1):87–105, 1992.
- [27] Intel. *Intel® SSE4 Programming Reference*. Intel Corporation, 2007.
- [28] Intel. *Intel® 64 and IA-32 Architectures Optimization Reference Manual*. Intel Corporation, 2011.
- [29] Intel. *Intel® Advanced Vector Extensions Programming Reference*. Intel Corporation, 2011.
- [30] D. Knuth, J. Morris, and V. Pratt. Fast pattern matching in strings. *SIAM J. Comput.*, 6:322–350, 1977.
- [31] M. Lothaire. *Combinatorics on Words*. Addison-Wesley, Reading, MA, U.S.A., 1983.
- [32] S. Muthukrishnan and K. V. Palem. Non-standard stringology: algorithms and complexity. In *Proceedings of the 26th Annual ACM Symposium on Theory of Computing (STOC)*, pages 770–779, 1994.
- [33] S. Muthukrishnan and H. Ramesh. String Matching Under a General Matching Relation. *Inf. Comput.*, 122(1):140–148, 1995.

- [34] G. Navarro and M. Raffinot. A bit-parallel approach to suffix automata: Fast extended string matching. In M. Farach-Colton, editor, *CPM*, volume 1448 of *Lecture Notes in Computer Science*, pages 14–33. Springer, 1998.
- [35] J. Tarhio and H. Peltola. String matching in the DNA alphabet. *Software Practice Experience*, 27:851–861, 1997.
- [36] U. Vishkin. Deterministic sampling - a new technique for fast pattern matching. *SIAM J. Comput.*, 20(1):22–40, 1990.
- [37] A. C.-C. Yao. The complexity of pattern matching for a random string. *SIAM J. Comput.*, 8(3):368–387, 1979.