

Speeding Up HMM Decoding and Training by Exploiting Sequence Repetitions*

Shay Mozes**¹, Oren Weimann¹, and Michal Ziv-Ukelson²

¹ MIT Computer Science and Artificial Intelligence Laboratory,
32 Vassar Street, Cambridge, MA 02139, USA.

shaymozes@gmail.com, oweimann@mit.edu

² School of Computer Science, Tel-Aviv University, Tel-Aviv 69978, Israel.
michaluz@post.tau.ac.il

Abstract. We present a method to speed up the dynamic program algorithms used for solving the HMM decoding and training problems for discrete time-independent HMMs. We discuss the application of our method to Viterbi’s decoding and training algorithms [26], as well as to the forward-backward and Baum-Welch [5] algorithms. Our approach is based on identifying repeated substrings in the observed input sequence. We describe three algorithms based alternatively on byte pair encoding (BPE) [24], run length encoding (RLE) and Lempel-Ziv (LZ78) parsing [27]. Compared to Viterbi’s algorithm, we achieve a speedup of $\Omega(r)$ using BPE, a speedup of $\Omega(\frac{r}{\log r})$ using RLE, and a speedup of $\Omega(\frac{\log n}{k})$ using LZ78, where k is the number of hidden states, n is the length of the observed sequence and r is its compression ratio (under each compression scheme). Our experimental results demonstrate that our new algorithms are indeed faster in practice. Furthermore, unlike Viterbi’s algorithm, our algorithms are highly parallelizable.

Key words: HMM, Viterbi, dynamic programming, compression

1 Introduction

Over the last few decades, Hidden Markov Models (HMMs) proved to be an extremely useful framework for modeling processes in diverse areas such as error-correction in communication links [26], speech recognition [7], optical character recognition [2], computational linguistics [20], and bioinformatics [13].

The core HMM-based applications fall in the domain of classification methods and are technically divided into two stages: a training stage and a decoding stage. During the *training* stage, the emission and transition probabilities of an HMM are estimated, based on an input set of observed sequences. This stage is usually executed once as a preprocessing stage and the generated (“trained”) models are stored in a database. Then, a *decoding* stage is run, again and again, in order to classify input sequences. The objective of this stage is to find the most probable sequence of states to have generated each input sequence given each model, as illustrated in Fig. 1.

Obviously, the training problem is more difficult to solve than the decoding problem. However, the techniques used for decoding serve as basic ingredients in solving the training problem. The Viterbi algorithm (VA) [26] is the best known tool for solving the decoding problem. Following its invention in 1967, several other algorithms have been devised for the decoding and training problems, such as the forward-backward and Baum-Welch [5] algorithms. These algorithms are all based on dynamic programs whose running times depend linearly on the length of the observed sequence. The challenge of speeding up VA by utilizing HMM topology was posed in 1997 by Buchsbaum and Giancarlo [7] as a major open problem. In this contribution, we address this open problem by using text compression and present the first provable speedup of these algorithms.

The traditional aim of text compression is the efficient use of resources such as storage and bandwidth. Here, we will compress the observed sequences in order to speed up HMM algorithms. Note

* A preliminary version of this paper appeared in [21].

** Work conducted while visiting MIT

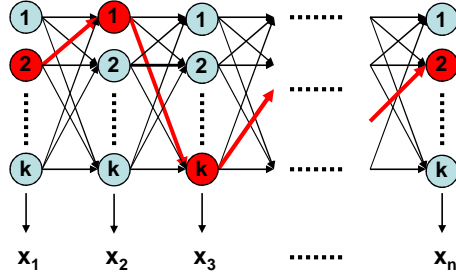


Fig. 1. The HMM on the observed sequence $X = x_1, x_2, \dots, x_n$ and states $1, 2, \dots, k$. The highlighted path is a possible path of states that generate the observed sequence. VA finds the path with highest probability.

that this approach, denoted “acceleration by text-compression”, has been recently applied to some classical problems on strings. Various compression schemes, such as LZ77, LZW-LZ78, Huffman coding, Byte Pair Encoding (BPE) and Run Length Encoding (RLE), were employed to accelerate exact string matching [16, 19, 24], approximate pattern matching [1, 15, 17, 22] and sequence alignment [3, 4, 8, 12, 18]. In light of the practical importance of HMM-based classification methods in state-of-the-art research, and in view of the fact that such techniques are also based on dynamic programming, we set out to answer the following question: can “acceleration by text compression” be applied to HMM decoding and training algorithms?

Our results. In this study we address the above challenge of speeding up HMM dynamic programming algorithms (Viterbi, forward-backward and Baum-Welch) by compression. We compress only in one dimension, the sequence axis, since typically $n \gg k$ and the states are non-repetitive. This compression enables the algorithm to adapt to the data and to utilize its repetitions. We present a basic toolkit of operations that could be further extended beyond this paper and applied to variant HMM-based problems in order to utilize common and repeated substrings. In general, the input sequences can be pre-compressed, as an offline stage, before our algorithms are applied. Such pre-compression, which is usually time-linear in the size of the input sequences, is done in this case not in order to save space but rather as a good investment in preparation for an all-against-all classification scheme in which each input sequence will be decoded many times according to various models and thus it pays off to pre-compress it once and for all.

Let X denote the input sequence and let n denote its length. Let k denote the number of states in the HMM. For any given compression scheme, let n' denote the number of parsed blocks in X and let $r = n/n'$ denote the compression ratio. Our results are as follows.

1. BPE is used to accelerate decoding by a factor of $\Omega(r)$.
2. RLE is used to accelerate decoding by a factor of $\Omega(\frac{r}{\log r})$.
3. Using LZ78, we accelerate decoding by a factor of $\Omega(\frac{\log n}{k})$. Our algorithm guarantees no degradation in efficiency even when $k > \log n$ and is experimentally more than five times faster than VA.
4. The same speedup factors apply to the Viterbi training algorithm.
5. For the Baum-Welch training algorithm, we show how to preprocess a repeated substring of size ℓ once in $O(\ell k^4)$ time so that we may replace the usual $O(\ell k^2)$ processing work for each occurrence of this substring with an alternative $O(k^4)$ computation. This is beneficial for any repeat with λ non-overlapping occurrences, such that $\lambda > \frac{\ell k^2}{\ell - k^2}$.
6. As opposed to VA, our algorithms are highly parallelizable.

Roadmap. The rest of the paper is organized as follows. In section 2 we give a unified presentation of the HMM dynamic programs. We then show in section 3 how these algorithms can be improved by

identifying repeated substrings. Three different implementations of this general idea are presented in section 4. Section 5 discusses the recovery of the optimal state-path. In section 6 we show how to adapt the algorithms to the training problem. A parallel implementation of our algorithms is described in section 7. Experimental results are presented in section 8. We summarize and discuss future work in section 9.

2 Preliminaries

Let Σ denote a finite alphabet and let $X \in \Sigma^n$, $X = x_1, x_2, \dots, x_n$ be a sequence of observed letters. A Markov *model* is a set of k states, along with emission probabilities $e_k(\sigma)$ - the probability to observe $\sigma \in \Sigma$ given that the state is k , and transition probabilities $P_{i,j}$ - the probability to make a transition to state i from state j .

The Viterbi Algorithm. The Viterbi algorithm (VA) finds the most probable sequence of hidden states given the model and the observed sequence. i.e., the sequence of states s_1, s_2, \dots, s_n which maximize

$$\prod_{i=1}^n e_{s_i}(x_i) P_{s_i, s_{i-1}} \quad (1)$$

The dynamic program of VA calculates a vector $v_t[i]$ which is the probability of the most probable sequence of states emitting x_1, \dots, x_t and ending with the state i at time t . v_0 is usually taken to be the vector of uniform probabilities (i.e., $v_0[i] = \frac{1}{k}$). v_{t+1} is calculated from v_t according to

$$v_{t+1}[i] = e_i(x_{t+1}) \cdot \max_j \{P_{i,j} \cdot v_t[j]\} \quad (2)$$

Definition 1 (Viterbi Step). We call the computation of v_{t+1} from v_t a Viterbi step.

Clearly, each Viterbi step requires $O(k^2)$ time. Therefore, the total runtime required to compute the vector v_n is $O(nk^2)$. The probability of the most likely sequence of states is the maximal element in v_n . The actual sequence of states can be then reconstructed in linear time.

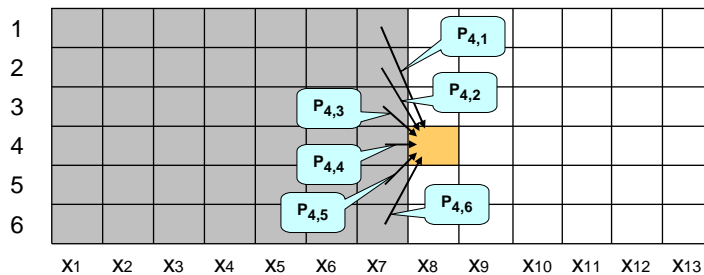


Fig. 2. The VA dynamic program table on sequence $X = x_1, x_2, \dots, x_{13}$ and states 1, 2, 3, 4, 5, 6. The marked cell corresponds to $v_8[4] = e_4(x_8) \cdot \max\{P_{4,1} \cdot v_7[1], P_{4,2} \cdot v_7[2], \dots, P_{4,6} \cdot v_7[6]\}$.

It is useful for our purposes to rewrite VA in a slightly different way. Let M^σ be a $k \times k$ matrix with elements $M_{i,j}^\sigma = e_i(\sigma) \cdot P_{i,j}$. We can now express v_n as:

$$v_n = M^{x_n} \odot M^{x_{n-1}} \odot \dots \odot M^{x_2} \odot M^{x_1} \odot v_0 \quad (3)$$

where $(A \odot B)_{i,j} = \max_k \{A_{i,k} \cdot B_{k,j}\}$ is the so called max-times matrix multiplication. VA computes v_n using (3) from right to left in $O(nk^2)$ time. Notice that if (3) is evaluated from left to right the

computation would take $O(nk^3)$ time (matrix-vector multiplication vs. matrix-matrix multiplication). Throughout, we assume that the max-times matrix-matrix multiplications are done naïvely in $O(k^3)$. Faster methods for max-times matrix multiplication [9] and standard matrix multiplication [25, 11] can be used to reduce the k^3 term. However, for small values of k this is not profitable.

The Forward-Backward Algorithms. The *forward-backward* algorithms are closely related to VA and are based on very similar dynamic programs. In contrast to VA, these algorithms apply standard matrix multiplication instead of max-times multiplication. The forward algorithm calculates $f_t[i]$, the probability to observe the sequence x_1, x_2, \dots, x_t requiring that $s_t = i$ as follows:

$$f_t = M^{x_t} \cdot M^{x_{t-1}} \cdot \dots \cdot M^{x_2} \cdot M^{x_1} \cdot f_0 \quad (4)$$

The backward algorithm calculates $b_t[i]$, the probability to observe the sequence $x_{t+1}, x_{t+2}, \dots, x_n$ given that $s_t = i$ as follows:

$$b_t = b_n \cdot M^{x_n} \cdot M^{x_{n-1}} \cdot \dots \cdot M^{x_{t+2}} \cdot M^{x_{t+1}} \quad (5)$$

Another algorithm which is used in the training stage and employs the forward-backward algorithm as a subroutine, is the Baum-Welch algorithm, to be further discussed in Section 6.

A motivating example. We briefly describe one concrete example from computational biology to which our algorithms naturally apply. CpG islands [6] are regions of DNA with a large concentration of the nucleotide pair CG . These regions are typically a few hundred to a few thousand nucleotides long, located around the promoters of many genes. As such, they are useful landmarks for the identification of genes. The observed sequence (X) is a long DNA sequence composed of four possible nucleotides ($\Sigma = \{A, C, G, T\}$). The length of this sequence is typically a few millions nucleotides ($n \simeq 2^{25}$). A well-studied classification problem is that of parsing a given DNA sequence into CpG islands and non CpG regions. Previous work on CpG island classification used Markov models with either 8 or 2 states ($k = 8$ or $k = 2$) [10, 13].

3 Exploiting Repeated Substrings in the Decoding Stage

Consider a substring $W = w_1, w_2, \dots, w_\ell$ of X , and define

$$M(W) = M^{w_\ell} \odot M^{w_{\ell-1}} \odot \dots \odot M^{w_2} \odot M^{w_1} \quad (6)$$

Intuitively, $M_{i,j}(W)$ is the probability of the most likely path starting with state j , making a transition into some other state, emitting w_1 , then making a transition into yet another state and emitting w_2 and so on until making a final transition into state i and emitting w_ℓ .

In the core of our method stands the following observation, which is immediate from the associative nature of matrix multiplication.

Observation 1. *We may replace any occurrence of $M^{w_\ell} \odot M^{w_{\ell-1}} \odot \dots \odot M^{w_1}$ in eq. (3) with $M(W)$.*

The application of observation 1 to the computation of equation (3) saves $\ell - 1$ Viterbi steps *each* time W appears in X , but incurs the additional cost of computing $M(W)$ once.

An intuitive exercise. Let λ denote the number of times a given word W appears, in non-overlapping occurrences, in the input string X . Suppose we naïvely compute $M(W)$ using $(|W| - 1)$ max-times matrix multiplications, and then apply observation 1 to all occurrences of W before running VA. We gain some speedup in doing so if

$$\begin{aligned} (|W| - 1)k^3 + \lambda k^2 &< \lambda |W| k^2 \\ \lambda &> k \end{aligned} \quad (7)$$

Hence, if there are at least k non-overlapping occurrences of W in the input sequence, then it is worthwhile to naïvely precompute $M(W)$, regardless of its size $|W|$.

Definition 2 (Good Substring). We call a substring W good if we decide to compute $M(W)$.

We can now give a general four-step framework of our method:

- (I) *Dictionary Selection:* choose the set $D = \{W_i\}$ of good substrings.
- (II) *Encoding:* precompute the matrices $M(W_i)$ for every $W_i \in D$.
- (III) *Parsing:* partition the input sequence X into consecutive good substrings $X = W_{i_1}W_{i_2} \cdots W_{i_{n''}}$ and let X' denote the compressed representation of this parsing of X , such that $X' = i_1i_2 \cdots i_{n''}$.
- (IV) *Propagation:* run VA on X' , using the matrices $M(W_i)$.

The above framework introduces the challenge of how to select the set of good substrings (step I) and how to efficiently compute their matrices (step II). In the next section we show how the BPE, RLE and LZ78 compression schemes can be applied to address this challenge. In practice, the choice of the appropriate compression scheme should be made according to the nature of the observed sequences. For example, genomic sequences tend to compress well with BPE [24] and ????????? [?] are well compressed by RLE. LZ78 guarantees asymptotic compression for any sequence and is useful in cases such as the CpG islands [6] identification problem in DNA sequences [10, 13], where k is smaller than $\log n$.

Another challenge is how to parse the sequence X (step III) in order to maximize acceleration. We show that, surprisingly, this optimal parsing may differ from the initial parsing induced by the selected compression scheme. To our knowledge, this feature was not applied by previous “acceleration by compression” algorithms.

Throughout this paper we focus on computing path probabilities rather than the paths themselves. The actual paths can be reconstructed in linear time as described in section 5.

4 Three different implementations of the general framework

4.1 Acceleration via Byte-Pair Encoding

In this section byte pair encoding is utilized to accelerate the Viterbi decoding computations by a factor of $\Omega(r)$, where n' is the number of phrases in the BPE-compressed sequence X and $r = n/n'$ is the BPE compression ratio of the sequence. The corresponding pre-processing term for encoding is $O(\Sigma'k^3)$, where Σ' denotes the number of character codes in the extended alphabet.

Byte pair encoding [14, 23, 24] is a simple form of data compression in which the most common pair of consecutive bytes of data is replaced with a byte that does not occur within that data. This operation is repeated until either all new characters are used up or no pair of consecutive two characters appears frequently in the substituted text. For example, the input string $ABABCABCD$ could be BPE encoded to $XYXD$, by applying the following two substitution operations: First $AB \rightarrow X$, yielding $XXCXCD$, and then $XC \rightarrow Y$. A substitution table, which stores for each character code the replacement it represents, is required to rebuild the original data. The compression ratio of BPE has been shown to be about 30% for biological sequences [24].

The compression time of BPE is $O(\Sigma'n)$. Alternatively, one could follow the approach of Shibata et al. [23, 24], and construct the substitution-table offline, during system set-up, based on a sufficient set of representative sequences. Then, using this pre-constructed substitution table, the sequence parsing can be done in time linear in the total length of the original and the substituted text. Let $\sigma \in \Sigma'$ and let W_σ denote the word represented by σ in the BPE substitution table. The four-step framework described in section 3 is applied as follows.

- (I) *Dictionary Selection:* all words appearing in the BPE substitution table are good substrings, i.e. $D = \{W_\sigma\}$ for all $\sigma \in \Sigma'$.

(II) *Encoding*: if σ is a substring obtained via the substitution operation $AB \rightarrow \sigma$ then

$$M(W_\sigma) = M(W_A) \odot M(W_B).$$

So each matrix can be computed by multiplying two previously computed matrices.

(III) *Parsing*: given an input sequence X , apply BPE parsing as described in [23, 24].

(IV) *Propagation*: run VA on X' , using the matrices $M(W_i)$ as described in section 3.

Time and Space Complexity Analysis Step I was already taken care of during the system-setup stage and therefore does not count in the analysis. Step II is implemented as an offline, preprocessing stage that is independent of the observed sequence X but dependent on the training model. It can therefore be conducted once in advance, for all sequences to come. The time complexity of this stage is $O(|\Sigma'|k^3)$ since each matrix is computed by one max-times matrix multiplication in $O(k^3)$ time. The space complexity is $O(|\Sigma'|k^2)$. Since we assume Step III is conducted in pre-compression, the compressed decoding algorithm consists of Step IV. In the propagation step (IV), given an input sequence of size n , we run VA using at most n' matrices, where n' denotes the number of blocks in the parsed text. Since each VA step takes k^2 time, the time complexity of this step is $O(n'k^2)$. Thus, the time complexity of the BPE-compressed decoding is $O(n'k^2)$ and the speedup, compared to the $O(nk^2)$ time of VA, is $\Omega(r)$.

4.2 Acceleration via Run-length Encoding

In this section we obtain an $\Omega(\frac{r}{\log r})$ speedup for decoding an observed sequence with run-length compression ratio r . A string S is *run-length encoded* if it is described as an ordered sequence of pairs (σ, i) , often denoted “ σ^i ”. Each pair corresponds to a *run* in S , consisting of i consecutive occurrences of the character σ . For example, the string $aaabbbcccc$ is encoded as $a^3b^2c^6$. Run-length encoding serves as a popular image compression technique, since many classes of images (e.g., binary images in facsimile transmission or for use in optical character recognition) typically contain large patches of identically-valued pixels. The four-step framework described in section 3 is applied as follows.

(I) *Dictionary Selection*: for every $\sigma \in \Sigma$ and every $i = 1, 2, \dots, \log n$ we choose σ^{2^i} as a *good substring*.

(II) *Encoding*: since $M(\sigma^{2^i}) = M(\sigma^{2^{i-1}}) \odot M(\sigma^{2^{i-1}})$, we can compute the matrices using repeated squaring.

(III) *Parsing*: Let $W_1W_2 \dots W_{n'}$ be the RLE of X , where each W_i is a run of some $\sigma \in \Sigma$. X' is obtained by further parsing each W_i into at most $\log |W_i|$ good substrings of the form σ^{2^j} .

(IV) *Propagation*: run VA on X' , as described in Section 3.

Time and Space Complexity Analysis. The offline preprocessing stage consists of steps I and II. The time complexity of step II is $O(|\Sigma|k^3 \log n)$ by applying max-times repeated squaring in $O(k^3)$ time per multiplication. The space complexity is $O(|\Sigma|k^2 \log n)$. This work is done offline once, during the training stage, in advance for all sequences to come. Furthermore, for typical applications, the $O(|\Sigma|k^3 \log n)$ term is much smaller than the $O(nk^2)$ term of VA.

Steps III and IV both apply one operation per occurrence of a good substring in X' : step III computes, in constant time, the index of the next parsing-comma, and step IV applies a single Viterbi step in k^2 time. Since $|X'| = \sum_{i=1}^{n'} \log |W_i|$, the complexity is

$$\sum_{i=1}^{n'} k^2 \log |W_i| = k^2 \log(|W_1| \cdot |W_2| \cdot \dots \cdot |W_{n'}|) \leq k^2 \log((n/n')^{n'}) = O(n'k^2 \log \frac{n}{n'}).$$

Thus, the speedup compared to the $O(nk^2)$ time of VA is $\Omega(\frac{n}{\log \frac{n}{n'}}) = \Omega(\frac{r}{\log r})$.

4.3 Acceleration via LZ78 Parsing

In this section we obtain an $\Omega(\frac{\log n}{k})$ speedup for decoding, and a constant speedup in the case where $k > \log n$. We show how to use the LZ78 [27] (henceforth LZ) parsing to find good substrings and how to use the incremental nature of the LZ parse to compute $M(W)$ for a good substring W in $O(k^3)$ time.

LZ parses the string X into substrings (LZ-words) in a single pass over X . Each LZ-word is composed of the longest LZ-word previously seen plus a single letter. More formally, LZ begins with an empty dictionary and parses according to the following rule: when parsing location i , look for the longest LZ-word W starting at position i which already appears in the dictionary. Read one more letter σ and insert $W\sigma$ into the dictionary. Continue parsing from position $i + |W| + 1$. For example, the string “AACGACG” is parsed into four words: A, AC, G, ACG. Asymptotically, LZ parses a string of length n into $O(hn/\log n)$ words [27], where $0 \leq h \leq 1$ is the entropy of the string. The LZ parse is performed in linear time by maintaining the dictionary in a trie. Each node in the trie corresponds to an LZ-word. The four-step framework described in section 3 is applied as follows.

- (I) *Dictionary Selection*: the good substrings are all the LZ-words in the LZ-parse of X .
- (II) *Encoding*: construct the matrices incrementally, according to their order in the LZ-trie, $M(W\sigma) = M(W) \odot M^\sigma$.
- (III) *Parsing*: X' is the LZ-parsing of X .
- (IV) *Propagation*: run VA on X' , as described in section 3.

Time and Space Complexity Analysis. Steps I and III were already conducted offline during the pre-processing compression of the input sequences (in any case LZ parsing is linear). In step II, computing $M(W\sigma) = M(W) \odot M^\sigma$, takes $O(k^3)$ time since $M(W)$ was already computed for the good substring W . Since there are $O(n/\log n)$ LZ-words, calculating the matrices $M(W)$ for all W s takes $O(k^3n/\log n)$. Running VA on X' (step IV) takes just $O(k^2n/\log n)$ time. Therefore, the overall runtime is dominated by $O(k^3n/\log n)$. The space complexity is $O(k^2n/\log n)$.

The above algorithm is useful in many applications, such as CpG island classification, where $k < \log n$. However, in those applications where $k > \log n$ such an algorithm may actually slow down VA.

We next show an adaptive variant that is guaranteed to speed up VA, regardless of the values of n and k . This graceful degradation retains the asymptotic $\Omega(\frac{\log n}{k})$ acceleration when $k < \log n$.

4.4 An improved algorithm

Recall that given $M(W)$ for a good substring W , it takes k^3 time to calculate $M(W\sigma)$. This calculation saves k^2 operations each time $W\sigma$ occurs in X in comparison to the situation where only $M(W)$ is computed. Therefore, in step I we should include in D , as good substrings, only words that appear as a prefix of at least k LZ-words. Finding these words can be done in a single traversal of the trie. The following observation is immediate from the prefix monotonicity of occurrence tries.

Observation 2. *Words that appear as a prefix of at least k LZ-words are represented by trie nodes whose subtrees contain at least k nodes.*

In the previous case it was straightforward to transform X into X' , since each phrase p in the parsed sequence corresponded to a good substring. Now, however, X does not divide into just good substrings and it is unclear what is the optimal way to construct X' (in step III). Our approach for constructing X' is to first parse X into all LZ-words and then apply the following greedy parsing to each LZ-word W : using the trie, find the longest good substring $w' \in D$ that is a prefix of W , place a parsing comma immediately after w' and repeat the process for the remainder of W .

Time and Space Complexity Analysis. The improved algorithm utilizes substrings that guarantee acceleration (with respect to VA) so it is therefore faster than VA even when $k = \Omega(\log n)$. In addition, in spite of the fact that this algorithm re-parses the original LZ partition, the algorithm still guarantees an $\Omega(\frac{\log n}{k})$ speedup over VA as shown by the following lemma.

Lemma 1. *The running time of the above algorithm is bounded by $O(k^3n/\log n)$.*

Proof. The running time of step II is at most $O(k^3n/\log n)$. This is because the size of the entire LZ-trie is $O(n/\log n)$ and we construct the matrices, in $O(k^3)$ time each, for just a subset of the trie nodes. The running time of step IV depends on the number of new phrases (commas) that result from the re-parsing of each LZ-word W . We next prove that this number is at most k for each word.

Consider the first iteration of the greedy procedure on some LZ-word W . Let w' be the longest prefix of W that is represented by a trie node with at least k descendants. Assume, contrary to fact, that $|W| - |w'| > k$. This means that w'' , the child of w' , satisfies $|W| - |w''| \geq k$, in contradiction to the definition of w' . We have established that $|W| - |w'| \leq k$ and therefore the number of re-parsed words is bounded by $k + 1$. The propagation step IV thus takes $O(k^3)$ time for each one of the $O(n/\log n)$ LZ-words. So the total time complexity remains $O(k^3n/\log n)$. \square

Based on Lemma 1, and assuming that steps I and III are pre-computed offline, the running time of the above algorithm is $O(nk^2/e)$ where $e = \Omega(\max(1, \frac{\log n}{k}))$. The space complexity is $O(k^2n/\log n)$.

5 Optimal state-path recovery

In this section we show how our decoding algorithms can trace back the optimal path, within the same space complexity and in $O(n)$ time. To do the traceback, VA keeps, along with the vector v_t (see eq. (2)), a vector of the maximizing arguments of eq. (2), namely:

$$u_{t+1}[i] = \operatorname{argmax}_j \{P_{i,j} \cdot v_t[j]\} \quad (8)$$

It then traces the states of the most likely path in reverse order. The last state s_n is simply the largest element in v_n , $\operatorname{argmax}_j \{v_n[j]\}$. The rest of the states are obtained from the vectors u by $s_{t-1} = u_t[s_t]$. We use exactly the same mechanism in the propagation step (IV) of our algorithm. The problem is that in our case, this only retrieves the states on the boundaries of good substrings but not the states within each good substring. We solve this problem in a similar manner.

Note that in all of our decoding algorithms every good substring W is such that $W = W_A W_B$ where both W_A and W_B are either good substrings or single letters. In LZ78-accelerated decoding, W_B is a single letter, when using RLE $W_A = W_B = \sigma^{|W|/2}$, and with BPE $W_A, W_B \in \Sigma'$. For this reason, we keep, along with the matrix $M(W)$, a matrix $R(W)$ whose elements are:

$$R(W)_{i,j} = \operatorname{argmax}_k \{M(W_A)_{i,k} \odot M(W_B)_{k,j}\} \quad (9)$$

Now, for each occurrence of a good substring $W = w_1, w_2, \dots, w_\ell$ we can reconstruct the most likely sequence of states s_1, s_2, \dots, s_ℓ as follows. From the partial traceback, using the vectors u , we know the two states s_0 and s_ℓ , such that s_0 is the most likely state immediately before w_1 was generated and s_ℓ is the most likely state when w_ℓ was generated. We find the intermediate states by recursive application of the computation $s_{|W_A|} = R(W)_{s_0, s_\ell}$.

Time and Space Complexity Analysis In all compression schemes, the overall time required for tracing back the most likely path is $O(n)$. Storing the matrices R does not increase the basic space complexity, since we already stored the similar-sized matrices $M(W)$.

6 The Training Problem

In the training problem we are given as input the number of states in the HMM and an observed training sequence X . The aim is to find a set of model parameters θ (i.e., the emission and transition probabilities) that maximize the likelihood to observe the given sequence $P(X|\theta)$. The most commonly used training algorithms for HMMs are based on the concept of Expectation Maximization. This is an iterative process in which each iteration is composed of two steps. The first step solves the decoding problem given the current model parameters. The second step uses the results of the decoding process to update the model parameters. These iterative processes are guaranteed to converge to a local maximum. It is important to note that since the dictionary selection step (I) and the parsing step (III) of our algorithm are independent of the model parameters, we only need run them once, and repeat just the encoding step (II) and the propagation step (IV) when the decoding process is performed in each iteration.

6.1 Viterbi training

The first step of Viterbi training [13] uses VA to find the most likely sequence of states given the current set of parameters (i.e., decoding). Let A_{ij} denote the number of times the state i follows the state j in the most likely sequence of states. Similarly, let $E_i(\sigma)$ denote the number of times the letter σ is emitted by the state i in the most likely sequence. The updated parameters are given by:

$$P_{ij} = \frac{A_{ij}}{\sum_{i'} A_{i'j}} \text{ and } e_i(\sigma) = \frac{E_i(\sigma)}{\sum_{\sigma'} E_i(\sigma')} \quad (10)$$

Note that the Viterbi training algorithm does not converge to the set of parameters that maximizes the likelihood to observe the given sequence $P(X|\theta)$, but rather the set of parameters that locally maximizes the contribution to the likelihood from the most probable sequence of states [13]. It is easy to see that the time complexity of each Viterbi training iteration is $O(k^2n + n) = O(k^2n)$ so it is dominated by the running time of VA. Therefore, we can immediately apply our compressed decoding algorithms from section 4 to obtain a better running time per iteration.

6.2 Baum-Welch training

The Baum-Welch training algorithm converges to a set of parameters that maximizes the likelihood to observe the given sequence $P(X|\theta)$, and is the most commonly used method for model training. Recall the forward-backward matrices: $f_t[i]$ is the probability to observe the sequence x_1, x_2, \dots, x_t requiring that the t 'th state is i and that $b_t[i]$ is the probability to observe the sequence $x_{t+1}, x_{t+2}, \dots, x_n$ given that the t 'th state is i . The first step of Baum-Welch calculates $f_t[i]$ and $b_t[i]$ for every $1 \leq t \leq n$ and every $1 \leq i \leq k$. This is achieved by applying the forward and backward algorithms to the input data in $O(nk^2)$ time (see eqs. (4) and (5)). The second step recalculates A and E according to

$$\begin{aligned} A_{i,j} &= \sum_t P(s_t = j, s_{t+1} = i | X, \theta) \\ E_i(\sigma) &= \sum_{t|x_t=\sigma} P(s_t = i | X, \theta) \end{aligned} \quad (11)$$

where $P(s_t = j, s_{t+1} = i | X, \theta)$ is the probability that a transition from state j to state i occurred in position t in the sequence X , and $P(s_t = i | X, \theta)$ is the probability for the t 'th state to be i in the sequence X . These probabilities are calculated as follows using the matrices $f_t[i]$ and $b_t[i]$ that were computed in the first step.

$P(s_t = j, s_{t+1} = i | X, \theta)$ is given by the product of the probabilities to be in state j after emitting x_1, x_2, \dots, x_t , to make a transition from state j to state i , to emit x_{t+1} at state i and to emit the rest of X given that the state is i :

$$P(s_t = j, s_{t+1} = i | X, \theta) = \frac{f_t[j] \cdot P_{i,j} \cdot e_i(x_{t+1}) \cdot b_{t+1}[i]}{P(X | \theta)} \quad (12)$$

where the division by $P(X | \theta) = \sum_i f_n[i]$ is a normalization by the probability to observe X given the current model parameters. Similarly

$$P(s_t = i | X, \theta) = \frac{f_t[i] \cdot b_t[i]}{P(X | \theta)}. \quad (13)$$

Finally, after the matrices A and E are recalculated, Baum-Welch updates the model parameters according to (10).

We next describe how to accelerate the Baum-Welch algorithm. It is important to notice that, in the first step of Baum-Welch, our algorithms to accelerate VA (sections 4.2 and 4.3) can be used to accelerate the forward-backward algorithms by simply replacing the max-times matrix multiplication with regular matrix multiplication. However, the accelerated forward-backward algorithms will only calculate f_t and b_t on the boundaries of good substrings. In what follows, we explain how to solve this problem and speed up the second step of Baum-Welch as well. We focus on updating the matrix A , updating the matrix E can be done in a similar fashion.

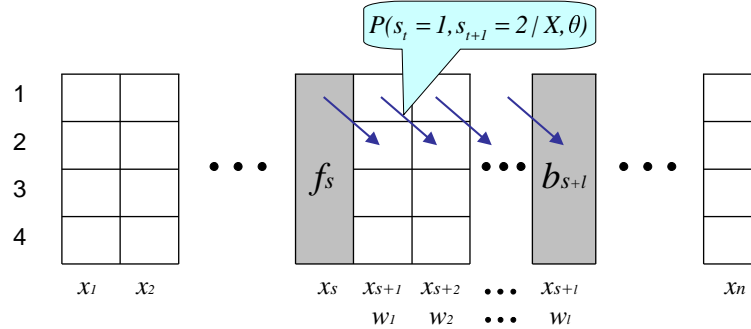


Fig. 3. The contribution of an occurrence of the good string W to A_{12} , computed as the sum of the arrow probabilities.

We observe that when accumulating the contribution of some appearance of a good substring W to A , Baum-Welch performs $O(k^2|W|)$ operations, but updates at most k^2 entries (the size of A). Therefore, we may gain a speedup by precalculating the contribution of each good substring to A and E . More Formally, let $W = w_1 w_2 \dots w_\ell$ be a substring of the observed sequence X starting s characters from the beginning (i.e., $W = x_{s+1} x_{s+2} \dots x_{s+\ell}$ as illustrated in Fig. 3). According to (11) and (12), the contribution of this occurrence of W to A_{ij} is:

$$\begin{aligned} & \sum_{t=s}^{s+\ell-1} \frac{f_t[j] \cdot P_{i,j} \cdot e_i(x_{t+1}) \cdot b_{t+1}[i]}{P(X | \theta)} \\ &= \frac{1}{P(X | \theta)} \sum_{t=0}^{\ell-1} f_{t+s}[j] \cdot P_{ij} \cdot e_j(w_{t+1}) \cdot b_{t+s+1}[i] \end{aligned}$$

However, by equation (4) we have

$$\begin{aligned} f_{t+s} &= M^{x_{t+s}} \cdot M^{x_{t+s-1}} \dots M^{x_1} \cdot f_0 \\ &= M^{x_{t+s}} \cdot M^{x_{t+s-1}} \dots M^{x_{s+1}} \cdot f_s \\ &= M^{w_t} \cdot M^{w_{t-1}} \dots M^{w_1} \cdot f_s \end{aligned}$$

and similarly by equation (5) $b_{t+s+1} = b_{s+\ell} \cdot M^{w_\ell} \cdot M^{w_{\ell-1}} \dots M^{w_{t+2}}$. The above sum thus equals

$$\begin{aligned} & \frac{1}{P(X|\theta)} \sum_{t=0}^{\ell-1} (M^{w_t} M^{w_{t-1}} \dots M^{w_1} \cdot f_s)_j \cdot P_{ij} \cdot e_j(w_{t+1}) \cdot (b_{s+\ell} \cdot M^{w_\ell} M^{w_{\ell-1}} \dots M^{w_{t+2}})_i \\ &= \frac{1}{P(X|\theta)} \sum_{t=0}^{\ell-1} \sum_{\alpha=1}^k (M^{w_t} M^{w_{t-1}} \dots M^{w_1})_{j,\alpha} \cdot f_s[\alpha] \cdot P_{ij} \cdot e_j(w_{t+1}) \cdot \sum_{\beta=1}^k b_{s+\ell}[\beta] \cdot (M^{w_\ell} M^{w_{\ell-1}} \dots M^{w_{t+2}})_{\beta,i} \\ &= \sum_{\alpha=1}^k \sum_{\beta=1}^k f_s[\alpha] \cdot b_{s+\ell}[\beta] \cdot \underbrace{\frac{1}{P(X|\theta)} \sum_{t=0}^{\ell-1} (M^{w_t} M^{w_{t-1}} \dots M^{w_1})_{j,\alpha} \cdot P_{ij} \cdot e_j(w_{t+1}) \cdot (M^{w_\ell} M^{w_{\ell-1}} \dots M^{w_{t+2}})_{\beta,i}}_{R_{ij}^{\alpha\beta}} \\ &\equiv \sum_{\alpha=1}^k \sum_{\beta=1}^k f_s[\alpha] \cdot b_{s+\ell}[\beta] \cdot R_{ij}^{\alpha\beta} \end{aligned} \tag{14}$$

Notice that the four dimensional array $R_{ij}^{\alpha\beta}$ can be computed in an encoding step (II) in $O(\ell k^4)$ time and is not dependant on the string context prior to X_s or following $X_{s+\ell}$. Furthermore, the vectors f_s and $b_{s+\ell}$ where already computed in the first step of Baum-Welch since they refer to boundaries of a good substring. Therefore, R can be used according to (14) to update A_{ij} for a single occurrence of W and for some specific i and j in $O(k^2)$ time. So R can be used to update A_{ij} for a single occurrence of W and for every i, j in $O(k^4)$ time. To get a speedup we need λ , the number of times the good substring W appears in X to satisfy:

$$\begin{aligned} \ell k^4 + \lambda k^4 &< \lambda \ell k^2 \\ \lambda &> \frac{\ell k^2}{\ell - k^2} \end{aligned} \tag{15}$$

This is reasonable if k is small. If $\ell = 2k^2$, for example, then we need λ to be greater than $2k^2$. In the CpG islands problem, if $k = 2$ then any substrings of length eight is good if it appears more than eight times in the text.

7 Parallelization

In this section we present a parallel version of our algorithm. We discuss our result for VA, but the same applies for the forward and backward algorithms. The basic recursion of VA in eq. (2) imposes a constraint on a parallel implementation. It is possible to compute the maximum in (2) in parallel for all states i , but the dependency of v_{t+1} on v_t makes it difficult to achieve a sublinear parallel algorithm. Once VA is cast into the form of eq. (3), there is a straightforward parallel implementation which runs in $O(\log n)$ time using $\frac{n}{2} k^4$ processors. The max-times multiplication of two matrices can be done in parallel in $O(1)$ time using k^4 processors so we can calculate the product of n matrices in $O(\log n)$ time. For the LZ compression scheme, if the parsing step (III) of our algorithm is performed in advance, then the same parallel algorithm applied to the sequence of good substrings X' runs even faster in $O(\log(n/\log n))$, though asymptotically this is still $O(\log n)$.

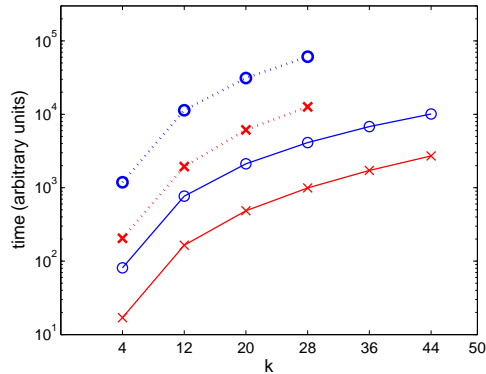


Fig. 4. Comparison of the cumulative running time of steps II and IV of our algorithm (marked x) with the running time of VA (marked o), for different values of k . Time is shown in arbitrary units on a logarithmic scale. Runs on the 1.5Mbp chromosome 4 of *S. cerevisiae* are in solid lines. Runs on the 22Mbp human Y-chromosome are in dotted lines. The roughly uniform difference between corresponding pairs of curves reflects a speedup factor of more than five.

8 Experimental Results

We implemented both our improved LZ-compressed algorithm from subsection 4.4 and classical VA in C++ and compared their execution times on a sequence of approximately 22,000,000 nucleotides from the human Y chromosome and on a sequence of approximately 1,500,000 nucleotides from chromosome 4 of *S. Cerevisiae* obtained from the UCSC genome database. The benchmarks were performed on a single processor of a SunFire V880 server with 8 UltraSPARC-IV processors and 16GB main memory. The implementation is just for calculating the probability of the most likely sequence of states, and does not traceback the optimal sequence itself. As we have seen, this is the time consuming part of the algorithm. We measured the running times for different values of k . In practice we found that the simple implementation of our algorithm (choosing all LZ-words as good substrings) is a little slower than the regular VA implementation. However, an implementation of the refined algorithm (choosing just LZ-words that appear as a prefix of more than k LZ-words) performs roughly five times faster than VA. The fastest variant of our algorithm uses as good substrings all LZ-words that appear as a prefix of more than a threshold of the LZ-words. The optimal threshold is dynamically computed in the parsing step (III) of our algorithm.

Unlike the procedure described in section 4.3, this variant parses X into good substrings by applying the greedy procedure from section 4.3 to the entire sequence X , rather than to each LZ-word individually. As we explained in the previous sections we are only interested in the running time of the encoding and propagation steps (II and IV) since the combined parsing/dictionary-selections steps (I and III) may be performed in advance and are not repeated by the training and decoding algorithms. A comparison of the running time of steps II and IV of this variant to the running time of the corresponding calculation by VA is shown in Fig. 4.

As k becomes larger, the optimal threshold and the number of good substrings decreases. Our algorithm performs faster than VA even for surprisingly large values of k . For example, for $k = 60$ our algorithm is roughly three times faster than VA. It is very likely that using better heuristics for identifying good substrings one can get even faster implementations.

9 Conclusions

FILL HERE...

References

1. G. Benson A. Amir and M. Farach. Let sleeping files lie: Pattern matching in Z-compressed files. *Journal of Comp. and Sys. Sciences*, 52(2):299–307, 1996.
2. O. Agazzi and S. Kuo. HMM based optical character recognition in the presence of deterministic transformations. *Pattern recognition*, 26:1813–1826, 1993.
3. A. Apostolico, G.M. Landau, and S. Skiena. Matching for run length encoded strings. *Journal of Complexity*, 15:1:4–16, 1999.
4. O. Arbell, G. M. Landau, and J. Mitchell. Edit distance of run-length encoded strings. *Information Processing Letters*, 2001.
5. L.E. Baum. An inequality and associated maximization technique in statistical estimation for probabilistic functions of a Markov process. *Inequalities*, 3:1–8, 1972.
6. A.P. Bird. Cpg-rich islands as gene markers in the vertebrate nucleus. *Trends in Genetics*, 3:342–347, 1987.
7. A. L. Buchsbaum and R. Giancarlo. Algorithmic aspects in speech recognition: An introduction. *ACM Journal of Experimental Algorithms*, 2:1, 1997.
8. H. Bunke and J. Csirik. An improved algorithm for computing the edit distance of run length coded strings. *Information Processing Letters*, 54:93–96, 1995.
9. T.M. Chan. All-pairs shortest paths with real weights in $O(n^3/\log n)$ time. In *Proc. 9th Workshop on Algorithms and Data Structures*, pages 318–324, 2005.
10. G.A. Churchill. Hidden Markov chains and the analysis of genome structure. *Computers Chem.*, 16:107–115, 1992.
11. D. Coppersmith and S. Winograd. Matrix multiplication via arithmetical progressions. *Journal of Symbolic Computation*, 9:251–280, 1990.
12. M. Crochemore, G. Landau, and M. Ziv-Ukelson. A sub-quadratic sequence alignment algorithm for unrestricted cost matrices. In *Proc. 13th Annual ACM/SIAM Symposium on Discrete Algorithms*, pages 679–688, 2002.
13. R. Durbin, S. Eddy, A. Krieh, and G. Mitcheson. *Biological Sequence Analysis*. Cambridge University Press, 1998.
14. P. Gage. A new algorithm for data compression. *The C Users Journal*, 12(2), 1994.
15. J. Karkkainen, G. Navarro, and E. Ukkonen. Approximate string matching over Ziv-Lempel compressed text. *Proc. 11th Annual Symposium On Combinatorial Pattern Matching (CPM)*, LNCS 1848:195–209, 2000.
16. J. Karkkainen and E. Ukkonen. Lempel-Ziv parsing and sublinear-size index structures for string matching. *Proc. Third South American Workshop on String Processing (WSP)*, pages 141–155, 1996.
17. J. Karkkainen and E. Ukkonen. Lempel-Ziv parsing and sublinear-size index structures for string matching. In *In Proceeding of the third South American Workshop on String Processing (WSP)*, pages 141–155, 1996.
18. V. Makinen, G. Navarro, and E. Ukkonen. Approximate matching of run-length compressed strings. *Proc. 12th Annual Symposium On Combinatorial Pattern Matching (CPM)*, LNCS 1645:1–13, 1999.
19. U. Manber. A text compression scheme that allows fast searching directly in the compressed file. *Proc. 5th Annual Symposium On Combinatorial Pattern Matching (CPM)*, LNCS 2089:31–49, 2001.
20. C. Manning and H. Schütze. *Statistical Natural Language Processing*. The MIT Press, 1999.
21. S. Mozes, O. Weimann, and M. Ziv-Ukelson. Speeding up hmm decoding and training by exploiting sequence repetitions. *Proc. 18th Annual Symposium On Combinatorial Pattern Matching (CPM)*, 2007. To appear.
22. G. Navarro, T. Kida, M. Takeda, A. Shinohara, and S. Arikawa. Faster approximate string matching over compressed text. *Proc. Data Compression Conference (DCC)*, pages 459–468, 2001.
23. Y. Shibata, T. Kida, S. Fukamachi, M. Takeda, A. Shinohara, T. Shinohara, and S. Arikawa. Byte Pair encoding: A text compression scheme that accelerates pattern matching. *Technical Report DOI-TR-161, Department of Informatics, Kyushu University*, 1999.
24. Y. Shibata, T. Kida, S. Fukamachi, M. Takeda, A. Shinohara, T. Shinohara, and S. Arikawa. Speeding up pattern matching by text compression. *Lecture Notes in Computer Science*, 1767:306–315, 2000.
25. V. Strassen. Gaussian elimination is not optimal. *Numerische Mathematik*, 13:354–356, 1969.
26. A. Viterbi. Error bounds for convolutional codes and an asymptotically optimal decoding algorithm. *IEEE Transactions on Information Theory*, IT-13:260–269, 1967.
27. J. Ziv and A. Lempel. On the complexity of finite sequences. *IEEE Transactions on Information Theory*, 22(1):75–81, 1976.