

Supplementary Issue: Computational Advances in Cancer Informatics (A)

Empirical Transition Probability Indexing Sparse-Coding Belief Propagation (ETPI-SCoBeP) Genome Sequence Alignment

Aminmohammad Roozgard¹, Nafise Barzigar¹, Shuang Wang², Xiaoqian Jiang²
and Samuel Cheng¹

¹School of Electrical and Computer Engineering, University of Oklahoma, Tulsa, OK, USA. ²Division of Biomedical Informatics, University of California, San Diego, CA, USA.

ABSTRACT: The advance in human genome sequencing technology has significantly reduced the cost of data generation and overwhelms the computing capability of sequence analysis. Efficiency, efficacy, and scalability remain challenging in sequence alignment, which is an important and foundational operation for genome data analysis. In this paper, we propose a two-stage approach to tackle this problem. In the preprocessing step, we match blocks of reference and target sequences based on the similarities between their empirical transition probability distributions using belief propagation. We then conduct a refined match using our recently published sparse-coding belief propagation (SCoBeP) technique. Our experimental results demonstrated robustness in nucleotide sequence alignment, and our results are competitive to those of the SOAP aligner and the BWA algorithm. Moreover, compared to SCoBeP alignment, the proposed technique can handle sequences of much longer lengths.

KEYWORDS: empirical transition probability, indexing, sparse-coding, belief propagation, genome sequence alignment

SUPPLEMENT: Computational Advances in Cancer Informatics (A)

CITATION: Roozgard et al. Empirical Transition Probability Indexing Sparse-Coding Belief Propagation (ETPI-SCoBeP) Genome Sequence Alignment. *Cancer Informatics* 2014;13(S1) 159–165 doi: 10.4137/CIN.S13887.

RECEIVED: April 30, 2014. **RESUBMITTED:** October 9, 2014. **ACCEPTED FOR PUBLICATION:** October 10, 2014.

ACADEMIC EDITOR: JT Efrid, Editor in Chief

TYPE: Methodology

FUNDING: Shuang Wang and Xiaoqian Jiang were funded in part by the NIH grants R00LM011392, R21LM012060, U54HL108460, and the NHGRI (K99HG008175). The authors confirm that the funder had no influence over the study design, content of the article, or selection of this journal.

COMPETING INTERESTS: Authors disclose no potential conflicts of interest.

COPYRIGHT: © the authors, publisher and licensee Libertas Academica Limited. This is an open-access article distributed under the terms of the Creative Commons CC-BY-NC 3.0 License.

CORRESPONDENCE: roozgard@ou.edu

Paper subject to independent expert blind peer review by minimum of two reviewers. All editorial decisions made by independent academic editor. Upon submission manuscript was subject to anti-plagiarism scanning. Prior to publication all authors have given signed confirmation of agreement to article publication and compliance with all applicable ethical and legal requirements, including the accuracy of author and contributor information, disclosure of competing interests and funding sources, compliance with ethical requirements relating to human and animal study participants, and compliance with any copyright requirements of third parties. This journal is a member of the Committee on Publication Ethics (COPE).

I. Introduction

In bioinformatics, sequence alignment is an important way to identify similar regions that might be associated with similar functional and structural relationship between sequences. With the quick growth of genomic data, it is important to develop effective sequence alignment techniques that are scalable. The past decade has witnessed the development of many sequence alignment technologies. Cancers are caused by the collection of genomic sequence changes.¹ Therefore, alignment and analyses of cancer genome sequences provide basics to understand cancer biology, diagnosis, and therapy.

In general, pairwise sequence alignment methods can be classified into local and global approaches. The global

alignment attempts to find the best match between two strings with similar lengths through global optimization. In contrast, the local alignment is usually used to identify regions of similarity between a short query and a longer sequence. Global alignments^{2–5} are less prone to demonstrating false homology as each letter of one sequence is constrained to being aligned to only one letter of the other. Local alignments,^{6–9} on the other hand, can cope with rearrangements between non-syntenic, orthologous sequences by identifying similar regions in sequences; this, however, comes at the expense of a higher false positive rate because of the inability of local aligners to take into account overall conservation maps.¹⁰



A lot of efforts have been made to improve the efficiency and efficacy of sequence alignments. The ClustalW program proposed by Thompson and Larkin^{11,12} uses a multi-stage mechanism to weigh and to align subsequences based on sequence divergences. In addition, sequence annealing technique incrementally builds sequence alignment one at a time by checking whether a single match is consistent with a partial multiple alignment.¹³ Darling et al proposed a hidden Markov model that uses a sum-of-pairs breakpoint score to facilitate the detection of rearrangement breakpoints, when genomes have unequal gene content.¹⁴ Mummer is a highly efficient suffix tree-based matching tool for whole genome alignment as well as incomplete genomes.¹⁵

Researchers also proposed heuristics to accelerate sequence alignment. For example, the bounded sparse dynamic programming (BSDP) is used to support rapid approximation of exhaustive alignment in Slater and Birney.¹⁶ Another heuristic-driven approach, namely FastTree, is a tree-based method that stores profiles of internal nodes in a tree, such that candidate joins can be quickly identified. FastTree is also scalable for handling alignments over 10,000 sequences.^{17,18}

Maximum-likelihood-based approaches like PhyML and RAxML-VI-HPC have been developed as well. PhyML¹⁹ used a hill-climbing algorithm that adjusts tree topology and branch length at each tree modification iteration. RAxML-VI-HPC,²⁰ which stands for randomized accelerated maximum likelihood for high-performance computing, takes advantages of a parallel program to support large-scale genome alignment.

In this paper, we propose a novel alignment method that uses sparse coding²¹ and empirical transition probability to tackle the scalability challenge. Thanks to the sparse representation, our mechanism can handle long sequences with reduced memory footprint. We also leverage belief propagation (BP) to combine local and neighboring information of candidate nucleotides into consideration, and generate matching scores to determine the best match. The rest of this paper is structured as follows. Section II introduces our proposed method. Section III presents our results, including the comparison against SOAP aligner²² and BWA.²³ Finally, we draw our conclusions in Section IV.

II. Proposed Method

In this section, we present our genome indexing and alignment framework in detail, where the proposed method includes three steps: indexing, index matching, and sequence matching. In this paper, we refer to “reference sequence” as the baseline sequence and try to align a “read sequence” against the baseline sequence.

Indexing. The current genome indexing methods generate huge indices before performing the actual alignment to decrease the alignment time.^{24,25} The indexing process can be very time-consuming. In contrast, our proposed indexing technique provides a faster and light-weight alternative for

index generation, which is similar to the big data retrieval systems that were proposed.^{26–28} These indices can reduce the search space and provide an estimation of the read sequence locations in the reference sequence. The proposed genome indexing technique models a nucleotide sequence as a graph by counting the transitions between each pair of nucleotides. To be more specific, as shown in Figure 1, we consider a graph with four states according to the different types of nucleotides and 16 vertices according to all possible transitions between nucleotides. We read the first nucleotide of the sequence and treat it as the initial state. Then, we move from one state to the other state by scanning the next nucleotide repeatedly till the end of the sequence. Afterward, we calculate the number of nucleotide transitions where we count how many times we pass one vertex in the graph and store them in a 4×4 matrix. Finally, we normalize the resulting matrix as follows:

$$I = \begin{matrix} & \begin{matrix} A & C & G & T \end{matrix} \\ \begin{matrix} A \\ C \\ G \\ T \end{matrix} & \begin{pmatrix} k_{aa} & k_{ac} & k_{ag} & k_{at} \\ k_{ca} & k_{cc} & k_{cg} & k_{ct} \\ k_{ga} & k_{gc} & k_{gg} & k_{gt} \\ k_{ta} & k_{tc} & k_{tg} & k_{tt} \end{pmatrix} \end{matrix} \times \frac{1}{\sum_{s,w \in \{a,c,g,t\}} k_{sw}} \quad (1)$$

where k_{sw} is the number that has the S -type nucleotide immediately before the W -type nucleotide.

If the length of a sequence is larger than a given threshold, ie, h , we divide it into subsequences with maximum length of h , where each subsequence will have o nucleotides that overlap with their neighbors. We set $o \geq \frac{h}{2}$ so that each pair of nucleotides

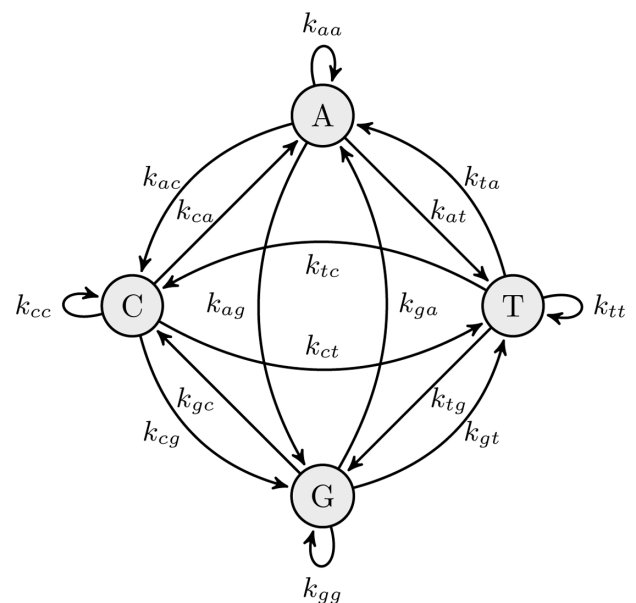


Figure 1. The transition diagram between nucleotides. k_{sw} is the number of appearance of the W -type nucleotide immediately after the S -type nucleotide, where $s, w \in \{a, c, g, t\}$.

can be counted at least twice. For each subsequence i , we count the transition of the nucleotides from the start of the subsequence till its end to reveal the number of different nucleotides that reside beside each other. In Figure 2, an input sequence with $b = 250$ is used to demonstrate the proposed indexing process, where I_i is the calculated index for the input sequence based on the transition graph shown on the left hand side. Finally, we normalize the transition matrix, which will be used to find the approximate location of each subsequence in the next step.

Index matching. The index matching step is designed to find similar indices based on global information of the sequence. We define a symmetric distance function between two index matrices, I and J , as follows: $D_{MSE}(I, J) = \|I - J\|_F$, where $\|\cdot\|_F$ is the Frobenius norm of the matrix.

After generating the indices of the reference sequence and the read sequence, the D_{MSE} distances to all reference sequence indices are calculated, where the top t most similar indices in terms of D_{MSE} are chosen as candidate indices. To find the best matched index, we resort to BP on a factor graph. In this paper, we provide a concise review about the BP algorithm on factor graph for the proposed algorithm. Interested readers can check our earlier publications in Refs. 29–31 for more details about the factor graph design and the BP algorithm.

We apply BP to the factor graph of the test sequence with n candidate nucleotides as the prior knowledge. BP updates the probability of candidate nucleotides based on the probabilities of their neighbors.

Then, the candidate index numbers are fed to a factor graph, and the corresponding D_{MSE} of each of candidates is employed to calculate the initial probability (prior probability) of each candidate. Then, message passing (ie, forward and backward) algorithm is applied to calculate the best match indices. The corresponding subsequences of these indices are used in the next step.

Sequence matching. The sequence matching step is based on sparse coding and BP algorithm. In this step, we use the subsequences that were selected in the previous step to generate an over-complete dictionary. Then, for each nucleotide in the read sequence, we pick n candidate nucleotides using sparse coding. By applying BP to a factor graph, we can obtain the best match for each nucleotide in the read sequence. A detailed description about the sequence matching can be found in our recent publications.^{31,32} A summary of the main procedure for our proposed alignment method is shown in Algorithm 1.

Implementation details.

- $I_i = \text{MakeIndex}(x_i)$ fills the state matrix I_i using the relationship of nucleotides in the subsequence x_i . The subsequence x_i is scanned through all its nucleotides, and the corresponding counts will be stored into the state matrix I_i . For example, k_{gg} in I_i in (2) shows how many times the nucleotide C will be identified, which is next to the nucleotide G in the subsequence x_i . Note that each subsequence x_i has a separate state matrix I_i , where i is the subsequence index.

$$I_i = \begin{matrix} & \begin{matrix} A & C & G & T \end{matrix} \\ \begin{matrix} A \\ C \\ G \\ T \end{matrix} & \begin{pmatrix} k_{aa} & k_{ac} & k_{ag} & k_{at} \\ k_{ca} & k_{cc} & k_{cg} & k_{ct} \\ k_{ga} & k_{gc} & k_{gg} & k_{gt} \\ k_{ta} & k_{tc} & k_{tg} & k_{tt} \end{pmatrix} \end{matrix} \quad (2)$$

- $[c_j, \rho_j] = \text{FindCandidates}(J_j, I, k)$ identifies k candidate state matrices that are highly similar to the test state matrix J_j in I and stores their indices in vector c_j and their probabilities in vector ρ_j . Note that the approach will compute the mean square error (MSE) of the test state matrix J_j with each possible I_i of the reference state matrices and select I_i that has the smallest MSEs.

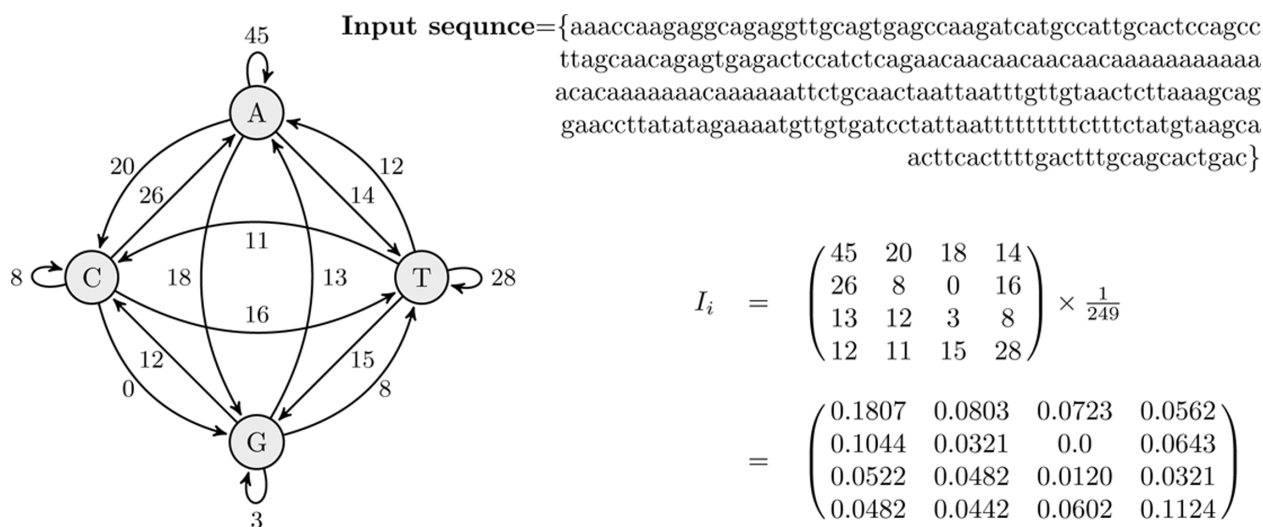


Figure 2. An example of the indexing procedure for a small sample subsequence.



Algorithm 1 Proposed nucleotide sequence alignment algorithm for estimating the location of the input sequence

Inputs: a reference sequence $X \in R^M$, a test sequence $Y \in R^N$, number of the candidate state matrix k , and number of the candidate points n .

Initialize: a 4×4 state matrix I storing the numbers of nucleotide states (2) and nucleotide overlap v .

Fill the reference state matrix I : For each subsequence $x_p \in X$ with v nucleotide overlap in each direction perform:

- $I_i = \text{MakeIndex}(x_p)$

Fill the test state matrix J : For each subsequence, $y_j \in Y$ with v nucleotide overlap in each direction perform:

- $J_j = \text{MakeIndex}(y_j)$
- $[c_p, \rho_p] = \text{FindCandidates}(J_j, I, k)$

Refine the candidate state matrix:

- $\hat{\rho} = \text{BP}(c, \rho)$

Find the corresponding nucleotide in the reference sequence X : For each subsequence, $y_j \in Y$ with v nucleotide overlap in each direction perform:

- $z_j = \text{FindBestSubsequence}(X, y_j, \theta, n)$ (see Refs. 31,32 for more details).

Output: the estimated version of aligned sequence Z

- $\hat{\rho} = \text{BP}(c, \rho)$ models the problem by a factor graph and applies BP³³ to update probability ρ . The updated probability $\hat{\rho}$ can be used to align the reference state matrix index onto the test state matrix index. In our case, we assign a variable node for each test state matrix index and connect each pair of neighboring state matrix indices with a factor node. Also, we introduce one extra factor node to take care of the prior knowledge obtained in the MSE step for each test state matrix index (for more details, see Ref. 32).

- $z_j = \text{FindBestSubsequence}(X, y_j, \theta, n)$ finds the corresponding location for a nucleotide $y_j \in Y$. In this step, the reference nucleotide sequence X and the test nucleotide sequence Y are converted into two integer sequences. Then, an over-complete dictionary is built with all subsequences in the X . We then apply sparse coding followed by using BP to identify the best matches (see Refs.31 and 32 for more details). Note that we used non-overlapped subsequences to build the dictionary. This change decreases the memory usage and the accuracy of the proposed algorithm in comparison with 1D-SCoBeP³¹ but it increases the speed of our alignment algorithm.

III. Experimental Results

We designed our experiments based on the work in Darling et al.¹⁴ to evaluate the proposed method for aligning the nucleotide sequences and to compare it with SOAP aligner,²² BWA,²³ and 1D-SCoBeP.³¹ We considered the problem of aligning a sequence of human nucleotides from the National Center for Biotechnology Information³⁴ and Cancer Genomics Hub.³⁵

To evaluate the performance of our approach, we conducted two sets of tests on the nucleotide sequences. In the first set, we selected 50 short subsequences of human genomes and then used SOAP aligner, BWA, 1D-SCoBeP, and the proposed method to find the location of selected subsequence nucleotide in the human chromosome. All the four algorithms successfully passed this test. We created 20 shuffled subsequences of the reference sequence as follows: for each read sequence R , we cut it into five pieces p_1, p_2, p_3, p_4 , and p_5 . Then, we switched p_2 with p_4 . Therefore, we converted a read sequence $R = [p_1, p_2, p_3, p_4, p_5]$ into a new read sequence $\hat{R} = [p_1, p_4, p_3, p_2, p_5]$.

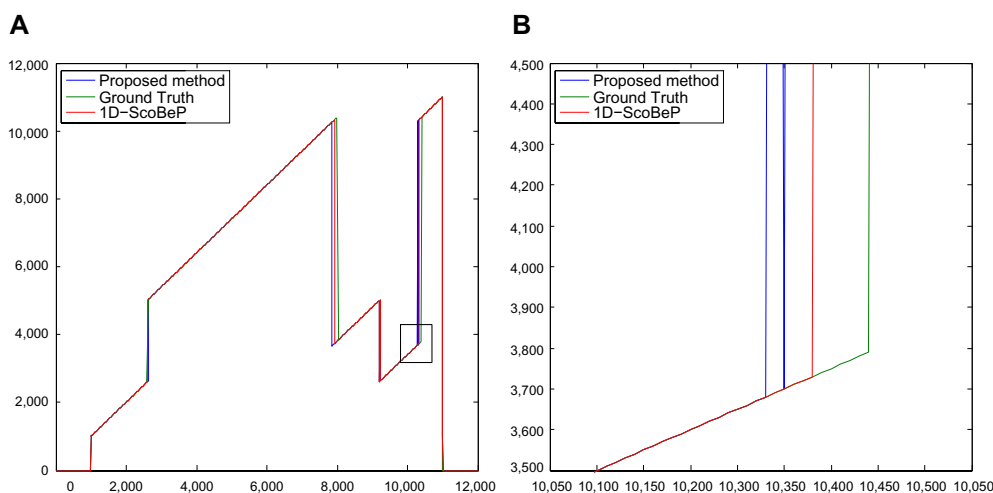


Figure 3. The results of proposed method for non-collinear nucleotide sequence alignment are shown. (A) Comparison among alignment results of the ground truth, 1D-SCoBeP,³¹ and the proposed method. (B) Magnified black square in (A) shows the gap between the proposed method, ground truth, and 1D-SCoBeP³¹ on the jump point. The x-axis and y-axis are the index numbers of the original genome sequences and the shuffled genome sequences, respectively.

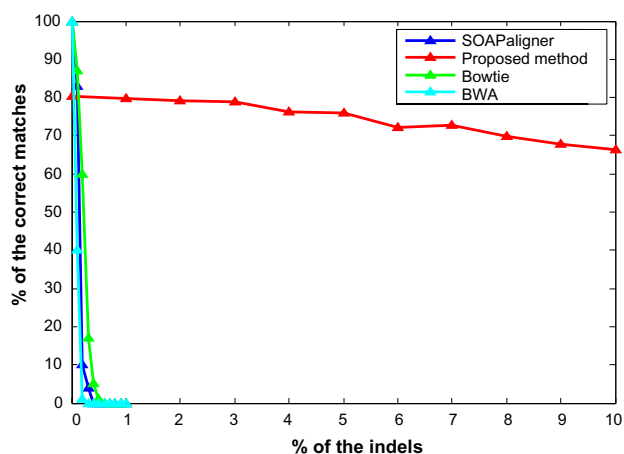


Figure 4. Accuracy of BWA,²³ SOAP aligner,²² and the proposed method in the presence of different Indel rates, where the testing genome sequences were obtained from Ref. 34.

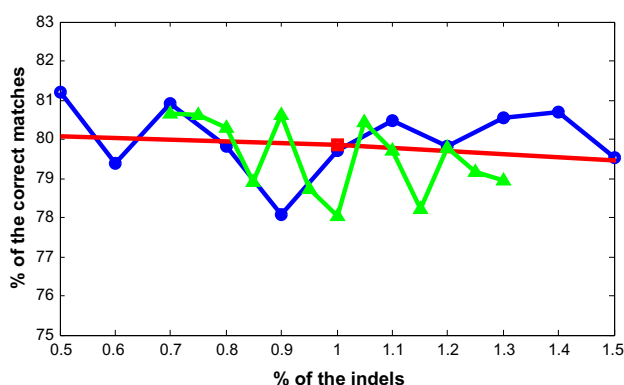


Figure 5. The percentage of successful alignments in the presence of 0.5–1.5% indels. The green line is the percentage of successful alignments, where the rate of the indels is changing between 0.7 and 1.3% at equal step of 0.05%. The blue line is the percentage of successful alignments, where the rate of the indels is changing between 0.5 and 1.5% at equal step of 0.1%, and the red line is the same as in Figure 4. Each point represents 10^5 random site selection with same indel rate. Note that the genome sequences used in this study were obtained from the National Center for Biotechnology Information and the Cancer Genomics Hub.^{34, 35}

Figure 3 shows that the result of the 1D-SCoBeP and the proposed method show a better performance with a gap of 100–120 nucleotides away from the ground truth. Since we were using non-overlapped subsequences for the dictionary generation, the gap between the proposed method and the ground truth was larger than those reported in 1D-SCoBeP.³¹ In our experiments, the following parameters were used: the number of candidate points $n = 3$, the sparsity factor $k = 3$, and the dictionary column size $a = 200$.

To evaluate the robustness of the proposed method, we generate indices for long human genome sequences (ie, 5×10^8 nucleotides), where $h = 10000$ and $\sigma = 5000$. Moreover, we synthesized insertion, deletion, and mutation (ie, indel) in these sequences. For indel rate, we picked 10^5

Table 1. Percentage of successful alignments.

% OF THE INDELS	ACCURACY OF RED LINE	ACCURACY OF BLUE LINE	ACCURACY OF GREEN LINE
0.00	80.33	–	–
0.50	–	81.19	–
0.60	–	79.38	–
0.70	–	80.90	80.64
0.75	–	–	80.63
0.80	–	79.82	80.29
0.85	–	–	78.90
0.90	–	78.09	80.63
0.95	–	–	78.74
1.00	79.85	79.71	78.04
1.05	–	–	80.43
1.10	–	80.49	79.70
1.15	–	–	78.22
1.20	–	79.81	79.78
1.25	–	–	79.16
1.30	–	80.54	78.94
1.40	–	80.70	–
1.50	–	79.52	–
2.00	79.09	–	–
3.00	78.90	–	–
4.00	76.33	–	–
5.00	75.90	–	–
6.00	72.09	–	–
7.00	72.86	–	–
8.00	69.79	–	–
9.00	67.87	–	–
10.00	66.42	–	–

number of subsequences with size of 10^4 nucleotides. Then, we randomly modified a certain number of nucleotides (based on the indel rate) and aligned them with the references. We counted the number of times the alignment location and real subsequence location (ie, ground truth) are matched, where the accuracy is defined as the count of the successfully aligned sequences over total number of the subsequences. Figure 4 shows the accuracy of alignment of the proposed method, BWA, and SOAP aligner in the presence of the different indel rates. The proposed method showed similar accuracies even when we increased the indel rate to 3%. Moreover, the proposed algorithm still showed more than 75% accuracy even after we modified 5% of the nucleotides in our selected subsequences. In contrast, the accuracy of the BWA and SOAP aligner decreased sharply as the indel rates increase.

We investigate the impact of small indel rate in the range from 0.5 to 1.5% in Figure 5. In this figure, we showed the



accuracy of 1% indels in red for the dataset used in Figure 4 as reference. To verify our result, we repeat the experiments with different indel steps and different read locations, and present the results in green and blue, respectively. Note that each point in this figure was obtained from the evaluation over 10^5 read sequences. There are slight variations among the curves because of statistical deviation. The summary of the indel rate accuracy is shown in Table 1.

The computational complexity of proposed is mainly determined by the following five steps: (1) indexing, (2) index matching, (3) extracting subsequence nucleotides as features and constructing the dictionary, (4) finding candidate nucleotides via sparse coding, and (5) applying BP. Assume that the sizes of the read and reference sequences are N and M nucleotides, respectively. The time required to create indexes is $O(M + N)$, in order to scan whole read and reference sequences. The number of reference sequence indexes is $I_M = \frac{M}{b} + \frac{2oM}{b^2} = O\left(\frac{M}{b}\right)$ and similarly, $I_N = O\left(\frac{N}{b}\right)$. Therefore, the time for the index matching is $O\left(\frac{M}{b}\right) \times O\left(\frac{N}{b}\right) = O\left(\frac{MN}{b^2}\right)$. After index matching step, the size of search space reduces from M to $\bar{M} = I_s \times b$, where I_s is the number of selected indexes I_s , and b is the size of each index. The time required for feature extraction will be $O(a(\bar{M} + N))$ where a is the size of the vector of extracted features for each nucleotide. The dictionary construction step involves the normalization of each column, which requires $O(a\bar{M})$ amount of time. Thus, the total time complexity of the first step is $O(a(\bar{M} + N))$. In the next step, the time complexity of subspace pursuit (SP) is $O(\log(f)a\bar{M})$,³⁶ where f is the number of iterations for searching the sparse vector. As we have to repeat the process to find candidate points for all N feature vectors, the time complexity of finding candidate points by SP is $O(\log(f)a\bar{M}N)$. Then, the time complexity of BP in our factor graph is $O(vn^2\bar{M})$, where v is the number of iterations before converging, and n is the number of candidates in each variable node. Finally, the time complexity of proposed method will be $O(MN + \log(f)a\bar{M}N + vn^2\bar{M})$.

IV. Conclusion

In this paper, we proposed a sparse coding and BP-based method for indexing and alignment genome sequences. The proposed method builds a transition matrix based on the neighboring nucleotides of an input sequence and then reduces the search space by selecting the top K most similar subsequences based on their distances. The proposed algorithm selects candidate nucleotides by using sparse coding with an overcomplete dictionary, which was constructed from the nucleotides of reference sequence in the indexing step. BP algorithm is then applied to select the best matches. Through experimental results, we showed that the proposed algorithms are comparable to SOAP aligner,²² BWA,²³ and 1D-SCoBeP³¹ in terms of the alignment accuracy. In addition, the proposed

method is robust to insertions, deletions, and mutations in the genome sequences when comparing with SOAP aligner and BWA. Finally, the proposed method is able to process much longer sequences than our previous 1D-SCoBeP approach.

Author Contributions

Conceived and designed the experiments: AR, NB, SW, XJ, SC. Analyzed the data: AR, NB, SW, XJ, SC. Wrote the first draft of the manuscript: AR, NB, SW, XJ, SC. Contributed to the writing of the manuscript: AR, NB, SW, XJ, SC. Agree with manuscript results and conclusions: AR, NB, SW, XJ, SC. Jointly developed the structure and arguments for the paper: AR, NB, SW, XJ, SC. Made critical revisions and approved final version: AR, NB, SW, XJ, SC. All authors reviewed and approved of the final manuscript.

REFERENCES

1. Meyerson M, Gabriel S, Getz G. Advances in understanding cancer genomes through second-generation sequencing. *Nat Rev Genet.* 2010;11(10):685–96.
2. Needleman SB, Wunsch CD. A general method applicable to the search for similarities in the amino acid sequence of two proteins. *J Mol Biol.* 1970;48(3):443–53.
3. Morgenstern B. Dialign 2: improvement of the segment-to-segment approach to multiple sequence alignment. *Bioinformatics.* 1999;15(3):211–8.
4. Bray N, Dubchak I, Pachter L. Avid: a global alignment program. *Genome Res.* 2003;13(1):97–102.
5. Brudno M, Do CB, Cooper GM, et al. Lagan and multi-lagan: efficient tools for large-scale multiple alignment of genomic DNA. *Genome Res.* 2003;13(4):721–31.
6. Smith TF, Waterman MS. Comparison of biosequences. *Adv Appl Math.* 1981;2(4):482–9.
7. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol Biol.* 1990;215(3):403–10.
8. Brudno M, Morgenstern B. Fast and sensitive alignment of large genomic sequences. In: Proceedings Bioinformatics Conference, 2002. IEEE Computer Society. IEEE; 2002; Stanford; CA; 138–47
9. Schwartz S, Kent WJ, Smit A, et al. Human-ouse alignments with blastz. *Genome Res.* 2003;13(1):103–7.
10. Brudno M, Malde S, Poliakov A, et al. Glocal alignment: finding rearrangements during alignment. *Bioinformatics.* 2003;19(1):i54–62.
11. Thompson JD, Higgins DG, Gibson TJ. Clustal w: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res.* 1994;22(22):4673–80.
12. Larkin M, Blackshields G, Brown N, et al. Clustal w and clustal x version 2.0. *Bioinformatics.* 2007;23(21):2947–8.
13. Schwartz AS, Pachter L. Multiple alignment by sequence annealing. *Bioinformatics.* 2007;23(2):e24–9.
14. Darling AE, Mau B, Perna NT. Progressive mauve: multiple genome alignment with gene gain, loss and rearrangement. *PLoS One.* 2010;5(6):e11147.
15. Kurtz S, Phillippy A, Delcher AL, et al. Versatile and open software for comparing large genomes. *Genome Biol.* 2004;5(2):R12.
16. Slater GS, Birney E. Automated generation of heuristics for biological sequence comparison. *BMC Bioinformatics.* 2005;6(1):31.
17. Price MN, Dehal PS, Arkin AP. Fasttree: computing large minimum evolution trees with profiles instead of a distance matrix. *Mol Biol Evol.* 2009;26(7):1641–50.
18. Price MN, Dehal PS, Arkin AP. Fasttree 2—approximately maximum-likelihood trees for large alignments. *PLoS One.* 2010;5(3):e9490.
19. Guindon S, Gascuel O. A simple, fast, and accurate algorithm to estimate large phylogenies by maximum likelihood. *Syst Biol.* 2003;52(5):696–704.
20. Stamatakis A. Raxml-vi-hpc: maximum likelihood-based phylogenetic analyses with thousands of taxa and mixed models. *Bioinformatics.* 2006;22(21):2688–90.
21. Olshausen BA, Field DJ. Sparse coding with an overcomplete basis set: a strategy employed by v1? *Vision Res.* 1997;37(23):3311–25.
22. Li R, Yu C, Li Y, et al. Soap2: an improved ultrafast tool for short read alignment. *Bioinformatics.* 2009;25(15):1966–7.
23. Li H, Durbin R. Fast and accurate short read alignment with burrows-wheeler transform. *Bioinformatics.* 2009;25(14):1754–60.
24. Liu C-M, Wong T, Wu E, et al. Soap3: ultra-fast gpu-based parallel alignment tool for short reads. *Bioinformatics.* 2012;28(6):878–9.



25. Vouzis PD, Sahinidis NV. Gpu-blast: using graphics processors to accelerate protein sequence alignment. *Bioinformatics*. 2011;27(2):182–8.
26. Carson C, Thomas M, Belongie S, Hellerstein JM, Malik J. Blobworld: a system for region-based image indexing and retrieval. *Visual Information and Information Systems*. CA, USA: Springer; 1999:509–17.
27. Doermann D. The indexing and retrieval of document images: a survey. *Comput Vis Image Underst*. 1998;70(3):287–98.
28. Liu C, Wechsler H. Robust coding schemes for indexing and retrieval from large face databases. *IEEE Trans Image Process*. 2000;9(1):132–7.
29. Roozgard A, Barzigar N, Cheng S, Verma P. Dense image registration using sparse coding and belief propagation. 5th International Conference on Signal Processing and Communication Systems (ICSPCS), Honolulu, HI, USA. Dec 2011. 1–5.
30. Roozgard A, Barzigar N, Cheng S, Verma P. Medical image registration using sparse coding and belief propagation. In: Engineering in Medicine and Biology Society (EMBC), 2012. Annual International Conference of the IEEE; August 2012; San Diego, CA; 1141–4.
31. Roozgard A, Barzigar N, Wang S, Jiang X, Ohno-Machado L, Cheng S. Nucleotide sequence alignment using sparse coding and belief propagation. In: Engineering in Medicine and Biology Society (EMBC), 2013. 35th Annual International Conference of the IEEE. IEEE; 2013; Osaka; 588–91.
32. Barzigar N, Roozgard A, Cheng S, Verma P. Scobep: Dense image registration using sparse coding and belief propagation. *J Visual Comm Image Represent*. 2012;24:137–47.
33. Kschischang FR, Frey BJ, Loeliger H-A. Factor graphs and the sum-product algorithm. *IEEE Trans Inf Theory*. 2001;47(2):498–519.
34. National Cancer Institute. The national center for biotechnology information. Available at: <http://www.ncbi.nlm.nih.gov/genome?db=genome>, January 2013. [Online].
35. University of California, Sanata Cruz. Cancer genomics hub. Available at: <https://cghub.ucsc.edu/datasets/benchmark-download.html>. January 2013. [Online].
36. Dai W, Milenkovic O. Subspace pursuit for compressive sensing: closing the gap between performance and complexity DTIC Document. Technical Report 2008.