# Streamlined Genome Sequence Compression using Distributed Source Coding

Shuang Wang, Xiaoqian Jiang

Division of Biomedical Informatics
University of California, San Diego
La Jolla, CA 92093
Email: {shw070, x1jiang}@ucsd.edu

Feng Chen, Lijuan Cui, Samuel Cheng

School of Electrical and Computer Engineering
University of Oklahoma
Tulsa, OK 74135
Email: {achenfengb, lj.cui, samuel.cheng}@ou.edu

## Abstract

We aim at developing a streamlined genome sequence compression algorithm to support alternative miniaturized sequencing devices, which have limited communication, storage, and computation power. Existing techniques that requires heavy-client (encoder side) cannot be applied. To tackle this challenge, we carefully examined distributed source coding (DSC) theory and developed a customized reference-based genome compression protocol to meet the low-complexity need at the client side. Based on the variation between source and reference, our protocol will pick adaptively either syndrome coding or hash coding to compress subsequences of changing code length. Our experimental results showed promising performance of the proposed method when compared with the state of the art algorithm (GRS).

## I. INTRODUCTION

The vision of miniaturized sequencing devices is turning into reality with the emergence of MinION (see Fig.1) by Oxford Nanopore [1]. Such devices are promising in a variety of potential applications, ranging from studying of wildlife and clinical capture of sequenced genes, to food inspection for identifying pathogens. However, such portable devices are commonly subject to the constraints in processing capabilities, power budget, and storage and communication limitations. With these constraints, the traditional view of genome compression architecture as simple decoder and complex encoder needs to be changed. It is urgent to develop novel techniques to satisfy the emerging reality challenges. Data compression methods (for reducing the storage space with significantly lower computational complexity and memory requirements) become crucial for the efficient management of genomic data in portable devices.



Fig. 1. An example of a miniaturised sequencing instrument.

In both situations (with or without reference sequences), traditional genome compression is computationally expensive at the encoder. The complexity is dominated by matching (approximately) repeated patterns of nucleotides—namely Adenine (A), Cytosine (C), Guanine (G) and Thymine (T)—between or within the DNA sequences. These patterns are also accompanied by insertions, deletions and substitutions of single nucleotides.

To date, a number of specialized DNA sequence compression algorithms have been proposed. In the spirit of Ziv and Lempel [2], Grumbach and Tachi [3] proposed the first DNA sequence compressor, Biocompress, to compress the exact repeating patterns with a specially designed Fibonacci coder. The algorithm was then improved in Biocompress-2 [4] by introducing a Markov model for encoding the

non-repeated regions. Chen et al. [5], [6] extended the earlier approach to cover approximated repeats by further exploiting the nature of DNA sequences. Meanwhile, the work in [7] introduced a combined CTW+LZ algorithm for searching approximate repeats and palindrome using hash and dynamic programming. Behzadi and Fessant [8] proposed a dynamic programming approach for the optimal selection of approximate repeats with promising compression efficiency being witnessed. However, such methods are heuristic as the underlying statistics of the sequence patterns are generally ignored. The authors in [9]–[11] proposed to combine the matching and substitution of approximate repeats and a specific normalized maximum likelihood model, obtaining a much higher compression ratio. Subsequently, statistical modeling for predicting the generation of symbols and arithmetic coding for such symbols in DNA sequences were proposed for more efficient compression. Cao *et al.* [12] proposed to estimate the probability distribution of symbols with a panel of "expert" to tackle the approximate repeat problem. Alternatively, finite context models are proposed to capture different aspects of statistical information along the sequence [13], [14], such reference free methods are plagued by their low compression rates (not greater than 6:1) and prohibitive computational consumption for large DNA sets.

Recognizing reference-free architectures do not fully utilizing information, a series of algorithms are proposed to compress sequences by matching approximate repeats with a reference sequence. The RLZ algorithm proposed by Kuruppu *et al.* [15] performed relative Lempel-Ziv compression of DNA sequences with the collection of related sequences. Wang *et al.* [16] proposed the GRS compressor, that is able to compress a sequence using a reference without any additional information. Applying the copy model into the matching of exact repeats in reference sequences, GReEn [17] achieved even larger gains when compared to [15] and [16]. Recently, reference-based algorithms [18], [19] achieved highly efficient compression performance for the fastq data format, by matching and comparing repeated subsequences in the reference sequences. Although the reference-based architectures can achieve hundreds of folds compression, the requirement of reference sequences makes it impractical for miniaturized devices, which have very limited storage space and communication bandwidth.

In this paper, we propose a novel and pioneering architecture for the genome compression application in miniaturized devices with limited processing capabilities, power budget, storage space and communication bandwidth. The contribution of this paper is three-fold.

1) First, to the best of our knowledge, the proposed architecture is the first practical one to meet the demands of miniaturized devices. Motivated by the distributed source coding (DSC) for sensor networks [20], the proposed scheme includes a simplified encoder without having access to reference sequences or communicating with other encoders, and a complex decoder that detects repeated subsequences in the stored reference sequences and decompress the received encoded bits with the specifically designed graphical model. Hence, the proposed compression system can successfully meet the constraints and requirements of the miniaturized sequencing devices.

2) Second, a flexible encoding and decoding mechanism is proposed. Using feedback from the decoder, the encoder transmits either hashes conducting the detection of variable-size exact repeats in decoder or syndromes obtained with low-complexity Slepian-Wolf coding [21] of the non-repeated subsequences. The proposed encoder and decoder perform efficiently by taking consideration of both exact repeats and approximately repeated subsequences (e.g. insertion, deletion and substitution).

3) Third, [21] the syndrome and reference sequences at the decoder, we construct a novel factor graph model to tackle the challenge in detecting insertion, deletion and substitution between the reference and original source. Experimental results show that the proposed architecture can achieve an efficient compression performance with significantly low encoding complexity when compared to the benchmark compressor GRS.

The rest of this paper is organized as follows. In Section II, we introduce the proposed DSC based genome compression system, which includes the implementation details of the proposed hash based (exactly) repeated sequence coding with adaptive length and an overview of syndrome based non-repeated
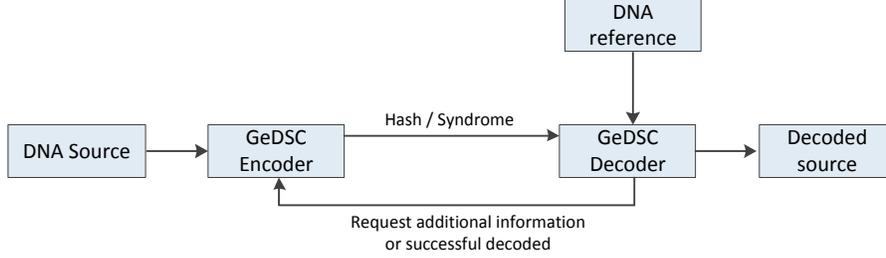
Fig. 2. Workflow of genome compression based on DSC.

sequence coding. Then, in Section III, we focus on the design of the syndromes based non-repeated sequence coding, which can handle the insertion, deletion, and substitution between sources and reference. The experimental results and conclusion remarks can be found in Sections IV and V, respectively.

## II. SYSTEM ARCHITECTURE

The block diagram of the proposed Genome compression framework is depicted in Fig. 2. Suppose that there are two correlated DNA sequences (i.e., source and reference sequences) available at the encoder and decoder, respectively, where the variations between two sequences are modeled by insertion, deletion and substitution. The alphabet of our studied DNA sequence is confined within the set $\{$'$A$', '$C$', '$G$', '$T$', '$N$'$\}$, where '$N$' denotes an unknown base due to a low sequencing quality. Fig. 3 shows the logical flow of the proposed framework, which we will discuss in details.

At the encoder (see the left hand side of Fig. 3), a streaming DNA sequence obtained from the portable sequencer will be first stored in the incoming data buffer for further processing. Second, a sub-sequence $\mathbf{x}_i^L$, which starts with the $i$-th to be compressed base in the source sequence, is extracted from the incoming data buffer, where its length $L$ and the corresponding coding method are decided by the adaptive code length and types selection module. The compressed sequence can be either LDPC Accumulate (LDPCA) syndromes $\mathbf{s}_{\mathbf{x}_i^L} = \mathbf{H}\mathbf{x}_i^L$ or hash bits $\mathbf{h}_{\mathbf{x}_i^L}$ depending on whether variations are presented between the source and the reference sequence, based on the decoder feedback, where $\mathbf{H}$ is the parity check matrix in LDPC codes. Third, the encoded sequence will be temporally stored in the forward data buffer and send to the decoder.

At the decoder (see the right hand side of Fig. 3), the received streaming data in the incoming data buffer will be processed by one of the following modules based on the corresponding data compression mode (i.e., either hash bits or syndromes).

1) For the received hash data $\mathbf{h}_{\mathbf{x}^L}$, it will be compared with the hashes generated from a bunch of sub-sequence candidates $\mathbf{y}_{j+U}^L, \cdots, \mathbf{y}_{j+V}^L$ within the reference sequence for $V - U + 1$ total candidates, where $j$ is the current offset compensated start location, and $U$ and $V$ are predefined lower and upper bounds of the search region for start locations. Then, the comparison result can be further processed as follows.

   a) If a matched hash $\mathbf{h}_{\mathbf{y}_k^L}$ for $k = j + U, \cdots, j + V$ is detected (i.e., $\mathbf{h}_{\mathbf{y}_k^L} = \mathbf{h}_{\mathbf{x}_i^L}$), the next offset compensated start location of the sliding window can be updated as $j = k + L$ (see Fig. 4). Moreover, we claim that $\mathbf{y}_k^L$ will be identical to $\mathbf{x}_i^L$, if $\mathbf{h}_{\mathbf{y}_k^L}$ and $\mathbf{h}_{\mathbf{x}_i^L}$ are matched with each other, which is the fundamental assumption of our proposed system. Intuitively, the aforementioned assumption can be enforced by choosing a strong hash code with a small search region. The experimental results based on sequences [22], [23] with total more than 238 million bases demonstrate that a 16-bit cyclic redundancy check (CRC) hash code with a search region $U = -2$ and $V = 10$ provides a strong assertion of such assumption. In addition, the decoder will inform the success to the encoder and request a longer code length based on a
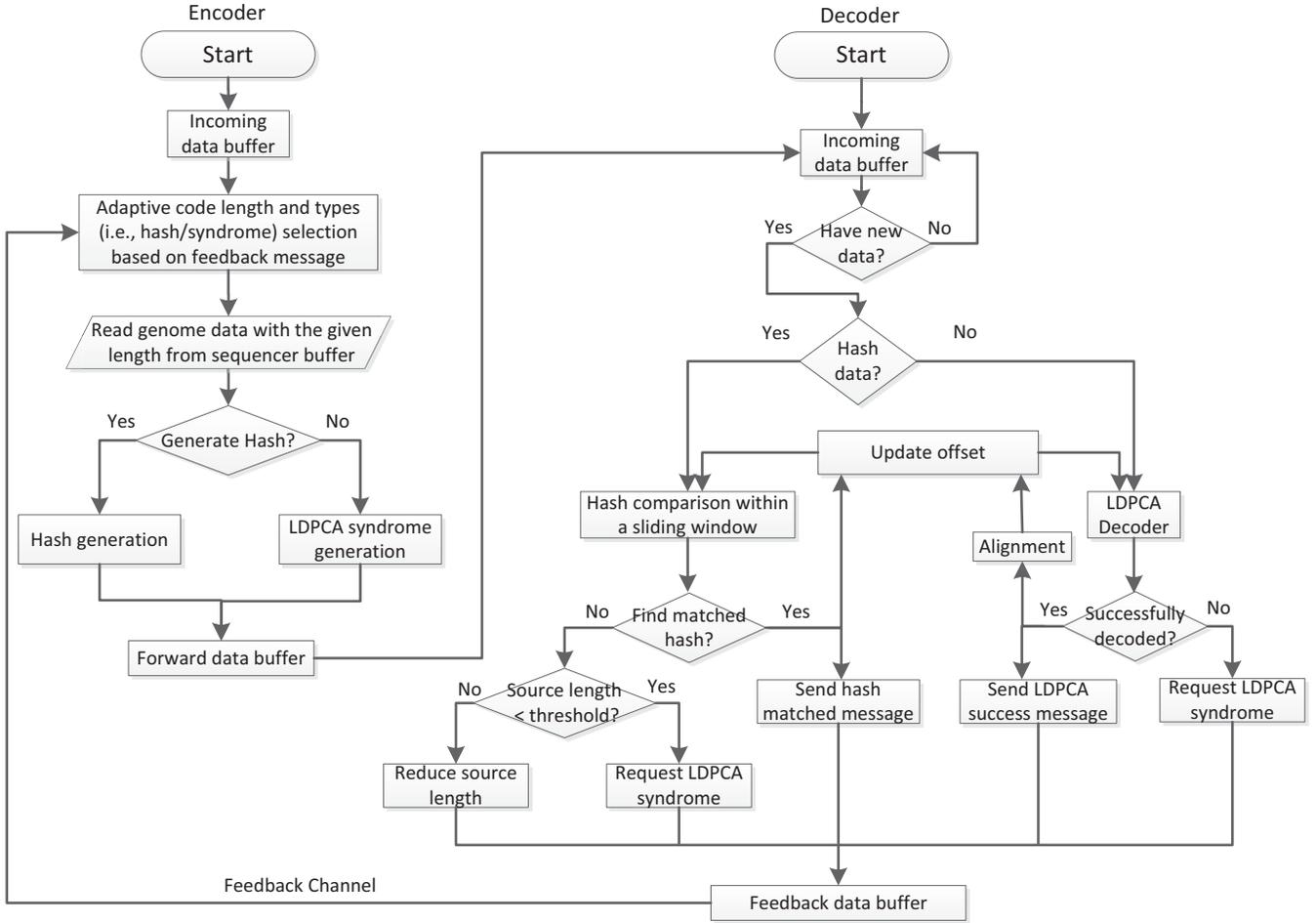
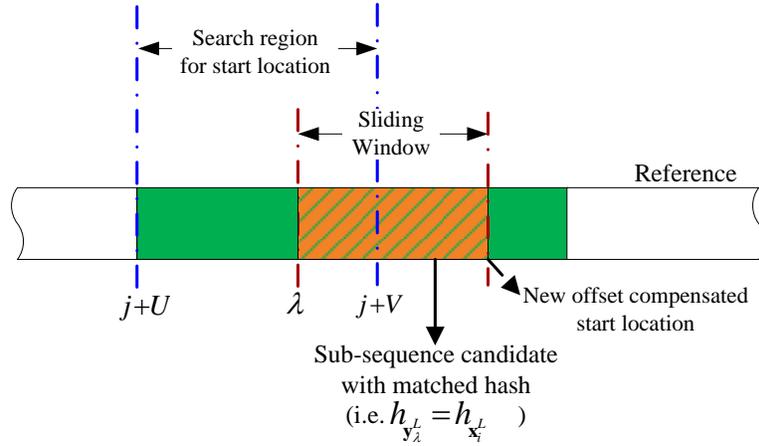Fig. 3. Workflow of genome compression based on DSC.



Fig. 4. The diagram of hash-based coding

predefined protocol as updating $L^{\text{current}} = bL_0$, where $L_0$ is a predefined initial length and the scaling factor $b$ is updated as $b = b + d_b$, $d_b$ is an incremental constant, and $b$ is initialized as 0. For example, at the beginning, $L^{\text{current}} = L_0$, if a matched hash is detected, the adaptive length $L^{\text{current}}$ will be updated as $L^{\text{current}} = d_b L_0$, as the scaling factor $b = 0 + d_b$. Similarly, if $n_h$ number of successively matched hashes are detected, the adaptive length and its corresponding

scale factor will be $L^{\text{current}} = n_h d_b L_0$ and $b = n_h d_b$, respectively.

   b) If no matched hash can be detected, the following two conditions will be checked.

      i) if $L^{\text{current}} = L_0$, the decoder will inform the hash matching failure to the encoder and request syndromes from the encoder for further action.

      ii) Otherwise, the decoder also informs the hash matching failure to the encoder, but requests a shorter code length by setting $L^{\text{current}} = L_0$.

2) For the received syndromes, the decoder will pass the syndrome to the proposed factor graph based LDPCA decoder with the capability of handling deletion, insertion and substitution between the source and the reference (see the next section for more implementation details). The following two conditions will be checked.

   a) If the decoded source $\hat{\mathbf{x}}_i^L$ satisfies both the parity check constraint (i.e., $\mathbf{s}_{\mathbf{x}_i} = H\hat{\mathbf{x}}_i^L$) and the hash constraint (i.e., $\mathbf{h}_{\mathbf{x}_i^L} = \mathbf{h}_{\hat{\mathbf{x}}_i^L}$), the decoder will send an LDPCA success message back to the encoder and update the offset compensated start location $j$ through the Smith-Waterman local alignment between the reference and the decoded source. Moreover, the encoder will send hash codes to the decoder for the next sub-sequence.

   b) Otherwise, the decoder will request additional LDPCA syndromes from the encoder.

## III. SYNDROME BASED NON-REPEATED SEQUENCE CODING

As previously mentioned in our system architecture, if an exact repeat cannot be identified by hash coding, the decoder will request syndromes from the encoder through a feedback channel. In this section, we introduce the codec design of the proposed syndrome based non-repeated sequence coding.

### A. Syndrome based non-repeated sequence encoding

The first step of the proposed syndrome based non-repeat encoder is to convert DNA data into a binary source, such that they can be compressed under a binary LDPCA encoder. Suppose the following mapping rule for the letters within the alphabet, i.e., 'A' $\rightarrow$ 000, 'C' $\rightarrow$ 001, 'G' $\rightarrow$ 010, 'T' $\rightarrow$ 011, 'N' $\rightarrow$ 100, a DNA subsequence $\mathbf{x}$ can be represented by the corresponding binary vector $\mathbf{x_b}$. For instance, given a DNA subsequence $\mathbf{x} = [\text{'A' 'T' 'G' 'C' 'T' 'N'}]^T$ with length $N = 6$, its corresponding binary vector will be $\mathbf{x_b} = [\mathbf{000\ 011\ 010\ 001\ 011\ 100}]^{\mathbf{T}}$ with length $3N$. Thus, for LDPC based Slepian-Wolf (SW) coding (i.e., lossless DSC), the compressed syndromes will be generated through $\mathbf{s_x} = \mathbf{Hx_b}$, where $\mathbf{H}$ is a sparse parity check matrix with size $M \times 3N$ and $M < 3N$. Thus, the resulting code rate can be expressed as $R = M/N$ bits per base. It is worth mentioning that the computational complexity of the aforementioned encoder is ultra-low, since the only operation is the bit-wise multiplication between the sparse matrix $\mathbf{H}$ and the original source. Moreover, we employ LDPCA codes to implement rate adaptive decoding, where the decoder can incrementally request additional LDPCA syndromes from the encoder through a feedback channel, when facing decoding errors.

### B. Syndrome based non-repeated sequence decoding

To perform syndrome based decoding for non-repeat DNA subsequence $\mathbf{x}$ with the reference sequence as side information $\mathbf{y}$, the key factor is to be able to explore the variations between the source subsequence $\mathbf{x}$ and the reference sequence $\mathbf{y}$, where the variations are modeled by the insertion, deletion, and substitution between the source and reference. Moreover, a substitution can be expressed as an insertion in the source sequence followed by a deletion in the corresponding location in the reference sequence. In this section, we demonstrate that such variations can be effectively estimated through Bayesian inference on graphical models. The graphical model of our proposed syndrome based decoding with variation is depicted in Fig. 5. In Fig. 5, the variable nodes (usually depicted by a circle) denote variables such as source symbol, binary source bits, local offset introduced by variation, and syndromes. Besides, factor nodes (depicted
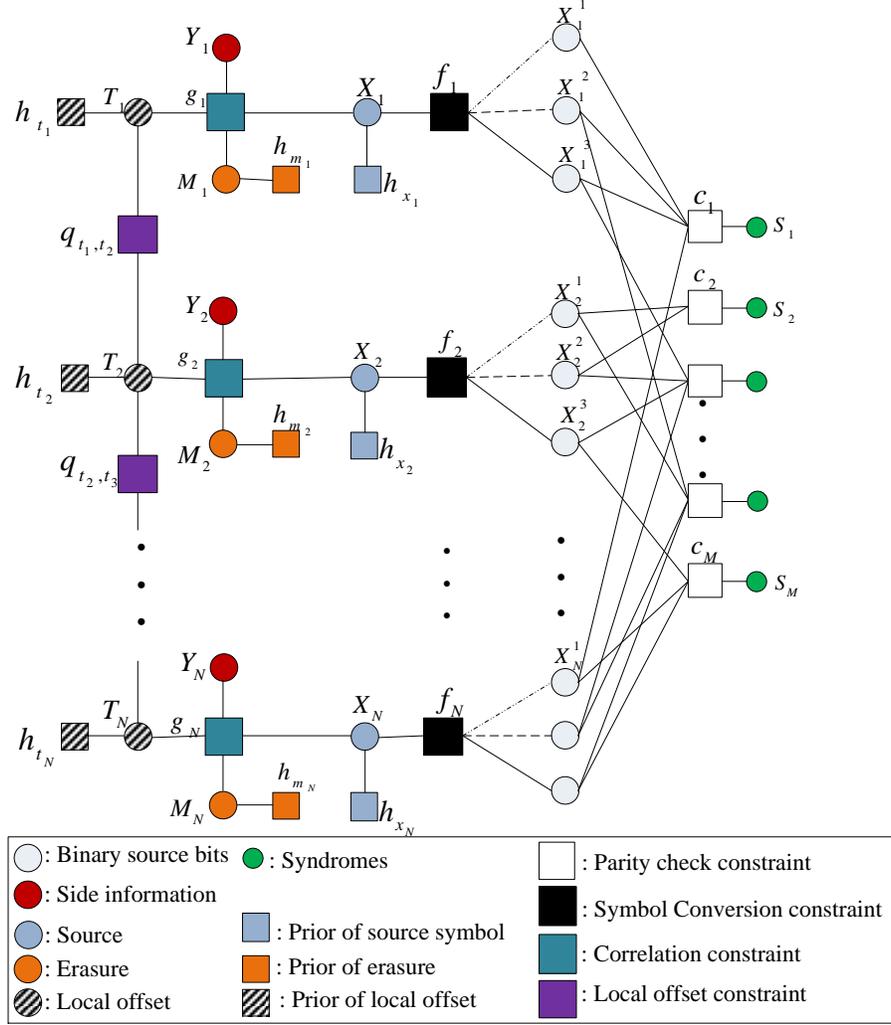
Fig. 5. Factor graph of genome compression based on DSC.

by squares) represent the relationship among the connected variable nodes. In the rest of this section, we will describe how to construct the proposed factor graph for the DNA sequence decoding with variations.

We first study the parity check constraint imposed by the received syndromes, where $s_1, \cdots, s_M$, the realization of variable node $S_l$, $l = 1, \cdots, M$, denotes the received syndromes in Fig. 5. Similar to the standard LDPC codes, the factor nodes $c_l$, $l = 1, \cdots, M$, take into account the parity check constraints, where the corresponding factor function can be expressed as

$$c_l(\mathbf{x}_{c_l}, s_l) = \begin{cases} 1, & \text{if } s_l \oplus \bigoplus \mathbf{x}_{c_l} = 0, \\ 0, & \text{otherwise.} \end{cases} \tag{1}$$

where $\mathbf{x}_{c_l}$ denotes the set of neighbors of the factor node $c_l$, and $\bigoplus \mathbf{x}_{c_l}$ denotes the binary sum of all elements of the set $\mathbf{x}_{c_l}$.

Moreover, $x_i^1, x_i^2, x_i^3$, the realization of variable node $X_i^r$ with $i = 1, \cdots, N$, $r = 1, 2, 3$, are the binary representation for the $i$-th base $x_i$ in the DNA sequence according to the mapping rule introduced in the Section III-A, where the mapping rule is captured by the factor node $f_i$, $i = 1, \cdots, N$ with corresponding factor function as follows

$$f_i(x_i, x_i^1, x_i^2, x_i^3) = \begin{cases} 1, & \text{if } \mathbf{map}(x_i^1, x_i^2, x_i^3) = x_i, \\ 0, & \text{otherwise.} \end{cases} \tag{2}$$

where $\mathbf{map}(\bullet)$ denotes the mapping from the binary bits "$\bullet$" to a letter in the alphabet, e.g. the output of $\mathbf{map}(0, 1, 1)$ corresponds to the letter "T".

Moreover, since the alphabet is not uniformly distributed in an arbitrary DNA sequence, the prior distribution for the alphabet is captured by the factor node $h_{x_i}$, where learning prior through training DNA sequences will be discussed shortly in the results section.

Now, we introduce an additional erasure variable node $M_i$ to capture the variation between reference $y_i$ and source $x_i$, where the variable $m_i = 1$ indicates the presence of variations, $m_i = 0$ means the existence of matches $y_{i+t_i} = x_i$ and $t_i = -T, \cdots, T$ are all possible local offsets within the search region $[-T, T]$. Moreover, the corresponding prior distribution of variable $m_i$ is captured by the factor node $h_{m_i}$ with factor function defined as

$$h_{m_i}(m_i) = \begin{cases} 1 - p_e & \text{if exist matches } y_{i+t_i} = x_i \\ p_e, & \text{otherwise,} \end{cases} \tag{3}$$

where $p_e$ can be learnt through training DNA sequence.

For the existence of matches $y_{i+t_i} = x_i$, the local offset $t_i$ is captured by the variable node $T_i$ and its corresponding prior is represented by the factor node $h_{t_i}$ with $h_{t_i}(t_i) = p_{t_i}$, where $p_{t_i}$ can be learnt through training DNA sequences. Furthermore, as the local offsets between adjacent DNA bases do not vary significantly in our assumption, it is expected that adjacent variables $t_i$ will not differ much in value. Such characteristic is captured by the additional factor node $q_{t_i, t_{i+1}}$, where the corresponding factor function is defined as

$$q_{t_i, t_{i+1}}(t_i, t_{i+1}, \alpha) = \frac{\alpha}{2} e^{-\alpha|t_i - t_{i+1}|}, \tag{4}$$

where $\alpha$ is the scale parameter of the Laplace distribution.

The factor node $g_i$ and its corresponding factor function $g_i(m_i, y_i, x_i, t_i)$ are introduced to combine the impact imposed by the side information $y_i$, erasure $m_i$ and local offset $t_i$. For $m_i = 0$, the factor function can be expressed as,

$$g_i(m_i = 0, y_i, x_i, t_i) = \begin{cases} 1 & \text{if } y_{i+t_i} = x_i, \\ 0, & \text{otherwise.} \end{cases} \tag{5}$$

For $m_i = 1$, the variable nodes $T_i$ and $Y_i$ will be disconnected from the factor node $g_i$. Therefore, the simplified factor function $g_i(m_i = 1, x_i) = 1$ can be used to take the impact of erasure into account.

Finally, in the context of Bayesian inference, estimating unknown variables corresponds to the evaluations of their marginal distribution, which can be efficiently achieved by performing message passing on a factor graph. With the factor graph defined above in Fig. 5, the original DNA sequence can be recovered by the estimated posterior distribution of each source $x_i$.

## IV. RESULTS

Two genome sequence data - The Arabidopsis Information Resrouce (TAIR) [24] and The Institute for Genomic Research (TIGR) [25] are adopted in our experiments. These two database are collected by professional groups or institutes, and have been widely used by research communities.

The TAIR maintains a database of genetic and molecular biology data for the model higher plant Arabidopsis thaliana [24]. In this experiment, we test TAIR8 [22] dataset and TAIR9 [23] dataset, where each dataset contains five chromosomes with over 238 million bases in total. Moreover, the genome of TAIR9 is used for testing our compression performance with TAIR8 as reference only available at the

7

decoder. For this experiment, all the hyper-parameters are initialized as follows, the initial code length $L_0 = 528$, the incremental constant $b_d = 3$, the scale parameter of Laplace distribution $\alpha = 1$, the maximum local offset search region $T = 4$ and the erasure probability $p_e = 0.01$. The proposed codec is implemented in MATLAB and evaluated on an Intel 3.0GHz CPU.
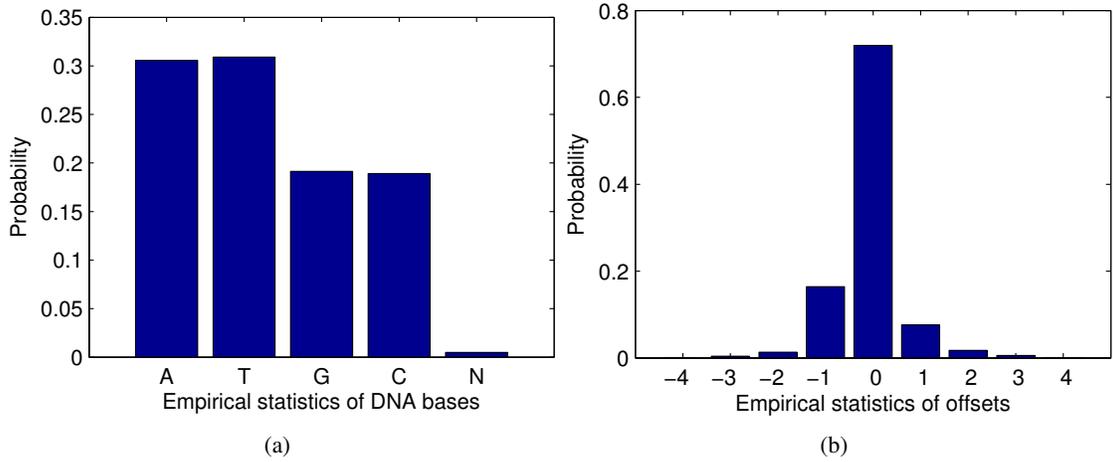


Fig. 6. The empirical statistics of (a) the DNA bases {'A', 'T', 'G', 'C', 'N'} and these of (b) the local offsets with the range from $-4$ to $4$.

First, the empirical marginal statistics of the DNA bases {'A', 'T', 'G', 'C', 'N'} and these of the local offsets $t_i$ within the range from $-4$ to $4$ are shown in Figs. 6(a) and 6(b), respectively, which will be used as the priors in the syndrome based non-repeated sequence decoding. In Figs. 6(a), we verify the assumption that the alphabets of DNA sequences are usually non-uniformly distributed. Moreover, Fig. 6(b) depicts that the maximum local offset with $T = 4$ is sufficiently large for capturing shifts between the reference and the source.
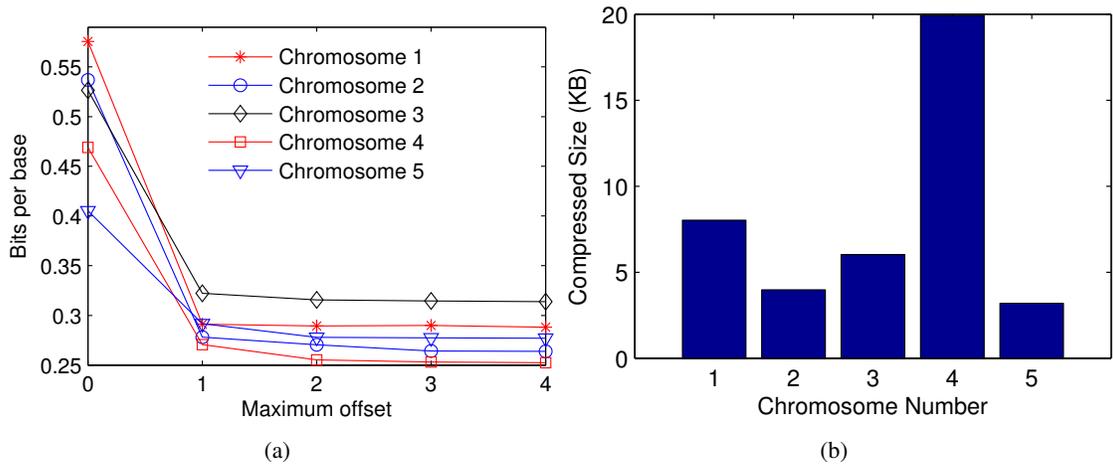


Fig. 7. Compression performance of the proposed codec on TAIR dataset, (a) the average code rates vs. the different maximum local offsets in syndrome coding; (b) the overall compression performance (i.e., hash bits + syndromes) for all 5 chromosomes.

Fig. 7(a) illustrates the relationship between the average code rates and the different maximum local offsets in syndrome coding based on all 5 chromosomes. In Fig. 7(a), we can see that the code rates decrease as the maximum local offsets increase, due to the fact that a larger maximum local offset offers a wider search region for exploring the reference. However, a larger maximum local offset may also result in a higher decoding complexity. Fig. 7(b) shows the overall compression performance (i.e., hash bits +

syndromes) for all 5 chromosomes in terms of compressed file size. Moreover, Fig. 8 shows a side-by-side comparison of the compression rate and compression time. We can see that both the proposed method and GRS algorithm achieve significant file size reductions (i.e., up to $8252\text{x}$ file size reduction).
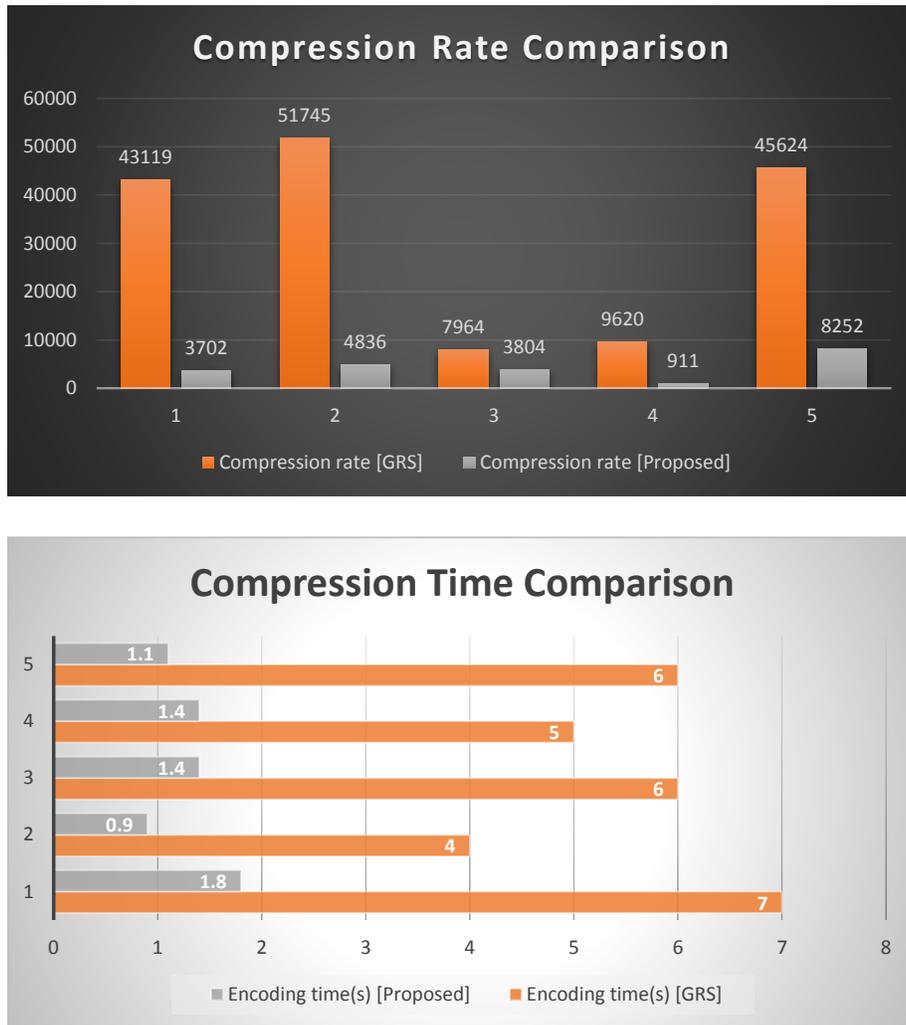


Fig. 8. Performance comparison between GRS and our proposed codec on TAIR dataset.

For the TIGR data, we tested the chromosome $4$ ($35.8\text{MB}$) of the TIGR5 dataset using the chromosome 4 of the TIGR6 as the reference by varying the LDPC code length (i.e., 528, 1056, 1584, 2112 and 2640) as shown in Table I. Moreover, we also implemented GRS method on this dataset, and the result was also listed in Table II as reference. We can see that the compression performance decreases as the LDPC code length increases. The proposed algorithm achieved better compression performance comparing with GRS in this case. It is because that the reference chromosome 4 in TIGR6 includes a significant amount of insertions when comparing with the same chromosome in TIGR 5, where the insertion information in the reference chromosome has no contribution to the size of DSC compressed data.

Our proposed encoder shows a significantly lower encoding complexity. It is worth mentioning that the proposed codec is implemented by MATLAB, where a potential performance boost is highly expected by using more efficient programming languages e.g., C/C++. To the best of our knowledge, this is the first study of DSC based genome compression. There is no doubt that it opens many possibilities for the

9

TABLE I
PERFORMANCES OF OUR PROPOSED METHOD ON CHROMOSOME 4 OF TIGR5 (35.8 MB)

| LDPC Code Length | Compression Size (KB) | Encoding Time (Seconds) | Decoding Time (Seconds) |
|---|---|---|---|
| 528 | 3.68 | 0.04 | 298 |
| 1056 | 4.58 | 0.009 | 596 |
| 1584 | 5.67 | 0.01 | 787 |
| 2112 | 6.97 | 0.01 | 1102 |
| 2640 | 6.72 | 0.08 | 1374 |

TABLE II
PERFORMANCE OF GRS ON CHROMOSOME 4 OF TIGR5 (35.8 MB)

| Compression Size (KB) | Encoding Time (Seconds) | Decoding Time (Seconds) |
|---|---|---|
| 26.34 | 12 | 6 |

portable miniaturized applications in which energy consumption and bandwidth usage are of paramount importance.

## V. CONCLUSION

In this paper, we present a DSC based genome compression architecture. To the best of our knowledge, the proposed framework is the first study of its kind: specially targeted at the low complexity genome encoding for miniaturized devices, which have limited processing capabilities, power budget, storage space and communication bandwidth. Compared to traditional reference based DNA compression algorithm (e.g., GRS), the proposed framework offers ultra-low encoding complexity (non-repeated subsequences are encoded using low complexity DSC encoding), while (exactly) repeated subsequences are compressed through adaptive length hash coding based on the decoder feedback. The customized factor graph based decoder tackles the challenges of detecting insertion, deletion and substitution between the reference and the original source, and it recovers the non-repeated subsequences based on received syndromes. Last but not least, our proposed genome compression framework incorporates LDPCA codes for rate adaptive decoding. Experimental results show that the proposed architecture could achieve an efficient compression performance with significantly lower encoding complexity when compared to the benchmark compressor GRS.

## ACKNOWLEDGMENT

## REFERENCES

[1] (2012, Feb.) MinION: A miniaturised sensing instrument. Oxford Nanopore Tech. [Online]. Available: http://www.nanoporetech.com/technology/minion-a-miniaturised-sensing-instrument

[2] J. Ziv and A. Lempel, "Compression of individual sequences via variable-rate coding," *IEEE Trans. Inform. Theory*, vol. 24, no. 5, pp. 530–536, Sep. 1978.

[3] S. Grumbach and F. Tahi, "Compression of DNA sequences," in *Proc. Data Compression Conf.*, Snowbird, Utah, USA, Mar. 1993, pp. 340–350.

[4] ——, "A new challenge for compression algorithms: Genetic sequences," *J. Inf. Process. Manage.*, vol. 30, no. 6, pp. 876–887, Nov. 1994.

[5] X. Chen, S. Kwong, and M. Li, "A compression algorithm for DNA sequences," *IEEE Eng. Med. Biol. Mag.*, vol. 20, no. 4, pp. 61–66, Jul. 2001.

[6] X. Chen, M. Li, B. Ma, and J. Tromp, "DNACompress: fast and effective DNA sequence compression," *Bioinformatics*, vol. 18, no. 12, pp. 1696–1698, Dec. 2002.

[7] T. Matsumoto, K. Sadakane, and H. Imai, "Biological sequence compression algorithms," *Genome Informatics*, vol. 11, pp. 43–52, Dec. 2000.

[8] B. Behzadi and F. L. Fessant, "DNA compression challenge revisited: A dynamic programming approach," *Combinatorial Pattern Matching*, vol. 3537, pp. 85–96, Jun. 2005.

[9] I. Tabus, G. Korodi, and J. Rissanen, "DNA sequence compression using the normalized maximum likelihood model for discrete regression," in *Proc. Data Compression Conf.*, Snowbird, Utah, USA, Mar. 2003, pp. 253–262.

[10] G. Korodi, I. Tabus, J. Rissane, and J. Astola, "DNA sequence compression - Based on the normalized maximum likelihood model," *IEEE Signal Processing Mag.*, vol. 24, no. 1, pp. 47–53, Jan. 2007.

[11] G. Korodi and I. Tabus, "Normalized maximum likelihood model of order-1 for the compression of DNA sequences," in *Proc. Data Compression Conf.*, Snowbird, Utah, USA, Mar. 2007, pp. 33–42.

[12] M. D. Cao, T. I. Dix, L. Allison, and C. Mears, "A simple statistical algorithm for biological sequence compression," in *Proc. Data Compression Conf.*, Snowbird, Utah, USA, Mar. 2007, pp. 43–52.

[13] A. J. Pinho, A. Neves, C. Bastos, and P. Ferreira, "DNA coding using finite-context models and arithmetic coding," in *Proc. IEEE Int. Conf. Acoustics, Speech and Signal Process*, Taipei, Apr. 2009, pp. 1693–1696.

[14] D. Pratas and A. Pinho, "Compressing the human genome using exclusively Markov models," in *5th Int. Conf. on Practical Applications of Comput. Biol. Bioinformatics (PACBB 2011)*, Mar. 2011, pp. 213–220.

[15] S. Kuruppu, S. Puglisi, and J. Zobel, "Relative lempel-ziv compression of genomes for large-scale storage and retrieval," *String Process. and Inf. Retrieval*, vol. 6393, pp. 201–206, Oct. 2010.

[16] C. Wang and D. Zhang, "A novel compression tool for efficient storage of genome resequencing data," *Nucleic Acids Res.*, vol. 39, no. 7, p. e45, Apr. 2011.

[17] A. Pinho, D. Pratas, and S. Garcia, "GReEn: A tool for efficient compression of genome resequencing data," *Nucleic Acids Res.*, vol. 40, no. 4, p. e27, Feb. 2012.

[18] C. Kozanitis, C. Saunders, S. Kruglyak, V. Bafna, and G. Varghese, "Compressing genomic sequence fragments using SlimGene," *J. Comput. Biol.*, vol. 18, no. 3, pp. 401–413, Mar. 2011.

[19] M. H. Y. Fritz, R. Leinonen, G. Cochrane, and E. Birney, "Efficient storage of high throughput DNA sequencing data using reference-based compression," *Genome Res.*, vol. 21, pp. 734–740, May 2011.

[20] Z. Xiong, A. D. Liveris, and S. Cheng, "Distributed source coding for sensor networks," *IEEE Signal Processing Mag.*, vol. 21, no. 5, pp. 80–94, Sep. 2004.

[21] S. S. Pradhan and K. Ramchandran, "Distributed source coding using syndromes (DISCUS): Design and construction," *IEEE Trans. Inform. Theory*, vol. 49, no. 3, pp. 626–643, Mar. 2003.

[22] E. Huala, A. Dickerman, M. Garcia-Hernandez, D. Weems, L. Reiser, F. LaFond, D. Hanley, D. Kiphart, M. Zhuang, W. Huang *et al.*, "The arabidopsis information resource (tair): a comprehensive database and web-based information retrieval, analysis, and visualization system for a model plant," *Nucleic acids research*, vol. 29, no. 1, pp. 102–105, 2001.

[23] S. Rhee, W. Beavis, T. Berardini, G. Chen, D. Dixon, A. Doyle, M. Garcia-Hernandez, E. Huala, G. Lander, M. Montoya *et al.*, "The arabidopsis information resource (tair): a model organism database providing a centralized, curated gateway to arabidopsis biology, research materials and community," *Nucleic acids research*, vol. 31, no. 1, pp. 224–228, 2003.

[24] Tair. [Online]. Available: http://www.arabidopsis.org/

[25] Tigr rice genome annotation. [Online]. Available: http://www.arabidopsis.org/