



Reference Based Genome Compression

B. G. Chern, I. Ochoa, A. Manolakos, A. No, K. Venkat and T. Weissman
Information Systems Laboratory, Stanford University

Introduction:

DNA sequencing technology has advanced to a point where storage is becoming the central bottleneck in the acquisition and mining of more data. Large amounts of data are vital for genomics research, and generic compression tools, while viable, cannot offer the same savings as approaches tuned to inherent biological properties.

We propose an algorithm to compress a target genome given a known reference genome. The genomes are related via insertions, deletions and substitutions. The proposed algorithm first generates a mapping from the reference to the target genome, and then compresses this mapping with an entropy coder.

Problem description:

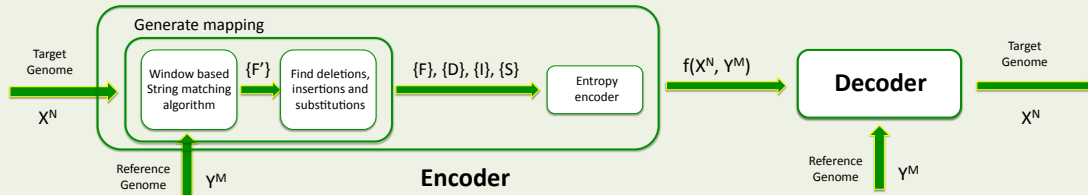
We consider the **problem** of compressing the target genome X^N given the reference genome Y^M as side information, which is available at both the encoder and decoder.

The **encoder** is described by the mapping $f(X^N, Y^M)$, which indicates the compressed version of X^N given Y^M .

The **decoder** is described by the mapping $g(f(X^N, Y^M), Y^M)$, which satisfies $X^N = g(f(X^N, Y^M), Y^M)$.

The **objective** is to construct an encoding/decoding scheme that minimizes the length of the representation $f(., .)$, while ensuring a perfect reconstruction of X^N at the decoder.

Proposed scheme:



1) Window based string matching algorithm (based on LZ77)

Idea: starting from a position in the target sequence, find the longest matching string within a fixed window in the reference sequence.

More formally: Define $Y_{W_k-L_k-L_k}^{W_k+R_k}$ as the window at time k . Assume x_{n_k} has already been encoded, for some n_k . Then, find the largest length l_k s.t. $x_{n_k+i} = y_{i-L_k}$ and $x_{n_k+i+1} = y_{i+1-L_k}$, for some $W_k - L_k \leq i \leq W_k + R_k$ and store the position $p_k = i$, length l_k and the novel symbol $z_k = x_{n_k+i+1}$. Set $n_{k+1} = n_k + l_k + 1$ and repeat the procedure until the target genome is exhausted.

Notice that the set of instructions $F = \{(p, l, z)\}$ suffices to reconstruct the target genome given the reference.

The window is dynamically updating its position and width along the reference so as to efficiently capture the synchronization in traversing along the two sequences.

Example:

REF: A₁ACCT₅GGCAT₁₀TGCAA₁₅CTTAA₂₀CTACC₂₅CGTTC₃₀AATCC₃₅CACTA₄₀...
TARGET: A₁ACCG₅GGCAT₁₀TGTTT₁₅AACGT₂₀ACCCG₂₅TTAAT₃₀CCCCA₃₅CTT₄₀...
 $F = \{(1, 4, G), (6, 7, T), (17, 5, G), (22, 8, A), (32, 5, C), (37, 3, T), \dots\}$

3) Entropy encoder for sets F, S, I and D

Techniques: Delta encoding, Huffman, Golomb codes...

Compression Performance:

Data: hg18 release from the UCSC Genome Browser, the korean genomes KOREF20090131 and KOREF20090224, the genome of a Han Chinese referred to as YH and the Watson JW genome.

Performance: We compare the proposed scheme with Gzip and the algorithm GReEn proposed in [2]. All the sizes are expressed in MB.

Ref	Target	Size	Gzip	GReEn	Our
hg18	JW	2,991	834.8	18.23	6.99
hg18	YH	2,987	832.0	9.85	8.50
KO224	YH	2,987	832.0	10.06	9.54
KO331	JW	2,991	834.8	23.23	10.22

Acknowledgement: The authors would like to thank Itai Sharon, Chris Miller and Golan Yona for helpful discussions.

2) Find substitutions, insertions and deletions and modify F accordingly

Substitutions: $S = \{(p^{(s)}, z^{(s)})\} \longrightarrow S = \{(5, G)\}$
Insertions: $I = \{(p^{(i)}, z^{(i)})\} \longrightarrow I = \{(19, G), (34, C)\}$
Deletions: $D = \{(p^{(d)}, l^{(d)}, z^{(d)})\} \longrightarrow D = \{(27, 1)\}$
Instructions: $F = \{(p, l, z)\} \longrightarrow F = \{(1, 12, T), (17, 23, T)\}$

References:

- [1] B. G. Chern, I. Ochoa, A. Manolakos, A. No, K. Venkat and T. Weissman, "Reference Based Genome Compression", <http://arxiv.org/pdf/1204.1912v1.pdf>
- [2] A. J. Pinho, D. Pratas and S. P. Garcia, "GReEn: a tool for efficient compression of genome resequencing data", Nucleic Acids Research, December 2011.