

GEne COmpression Z (GECOZ) Formats.

The GECOZ file formats is a set of file formats for efficient storing of genetic information.

It includes the Reference Genome file format (*.gcz) which contains compressed genome reference as a chain of blocks that consist of headers and nucleotide sequences.

This way it represents a compact FASTA format representation based on the Burrows–Wheeler Transformed (BWT) sequence. The BWT sequence is stored as a Huffman Shaped Wavelet Tree (HSWT) which automatically provides Huffman compression. Given that compression of genetic information doesn't improve much with higher order entropy encoding ($H_0 = 1.974$, $H_4 = 1.910$ bit/nct), further HSWT nodes compression via RLE, RRR found inefficient (Grossi, Vitter, & Xu, 2011). Using pre-calculated ranks for HSWT nodes extract and locate operations may be performed in a constant time. The implemented Succinct Suffix Array index is a variation of FM-index and is broadly used for indexing genomic information in many bioinformatics tools such as BWA (Li & Durbin, 2009) and Bowtie2 (Langmead & Salzberg, 2012).

In addition to the compressed Reference Genome *.gcz file, there is a corresponding Sparse Suffix Array index file (*.gcs). The size of the index file depends of the sampling rate. The reasonable parameters used by other bioinformatics tools are 16, 32, 64 and 128 (Ferragina, González, Navarro, & Venturini, 2009).

GECOZ SEQUENCE REFERENCE FORMAT

The *.gcz file consists of a sequence blocks. Any block corresponds to continuous sequence[s] (i.e. chromosome[s] or contigs) and consists of a header and the data section (Table 1).

"GecozBWT"	"magic" string	8 bytes
version	1	1 byte
size	block size	8 bytes
len	the sequence string S length	8 bytes
headers	FASTA header1	len (string)
	terminating \0	1 byte
	FASTA header2	len (string)
	terminating \0	1 byte
	headers terminating \0	1 byte
Deflate Length Table	RFC 1951 3.2.7 encoded Huffman table	
HSWT	Huffman Shaped Wavelet Tree, a sequence of bit vectors.	

Table 1. BWT Data block format

The block may contain more than one sequence. Every FASTA header is terminated by the '\0' byte. The last FASTA header is terminated by the double '\0'.

The data block contains the Deflate Length Table (as described in RFC 1951 3.2.7) and the HSWT nodes placed one after another using pre-ordered left-to-right traversal algorithm. The Huffman Codes are also generated according RFC 1951.

For fast O(1) extract/locate operations HSWT nodes include pre-calculated ranks:

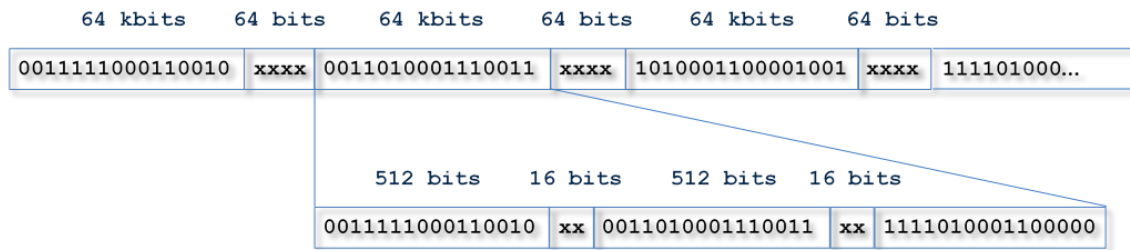


Figure 1. Ranked bit vector format

GECOZ SUCCINCT SUFFIX ARRAY INDEX FORMAT

The Succinct Suffix Array data is stored in the *.gcz file and consists of a ranked bit vector of stored Suffix Array indexes followed by the indexes themselves.

"GecozSSA"	"magic" string	8 bytes
version	1	1 byte
len	the index size	8 bytes
hash	the headers hash	8 bytes
SSA bit vector	bit vector that marks stored SSA indexes	
SSA indexes	wavelet table (bit vectors) that stores permuted SSA indexes	

Table 2. SSA Index block format

Bit vectors use the same format as on the Figure 1. SSA indexes are stored as a wavelet table where each row holds a bit vector of sorted index bit values (Table 3). Find and get index operations take $\log(n) \times \text{rank}()$ time, where n is a number of indexes.

06 ₀	10 ₀	11 ₀	04 ₀	08 ₀	14 ₀	00 ₀	12 ₀	02 ₀	05 ₀	16 ₁	03 ₀	07 ₀	01 ₀	09 ₀	15 ₀	13 ₀
00110	01010	01011	00100	01000	01110	00000	01100	00010	00101	10000	00011	00111	00001	01001	01111	01101
006 ₀	010 ₁	011 ₁	004 ₀	008 ₁	014 ₁	000 ₀	012 ₁	002 ₀	005 ₀	003 ₀	007 ₀	001 ₀	009 ₁	015 ₁	013 ₁	016 ₀
006 ₁	004 ₁	000 ₀	002 ₀	005 ₁	003 ₀	007 ₁	001 ₀	110 ₀	111 ₀	108 ₀	114 ₁	112 ₁	109 ₀	115 ₁	113 ₁	016 ₀
000 ₀	002 ₁	003 ₁	001 ₀	106 ₁	104 ₀	105 ₀	107 ₁	010 ₁	011 ₁	008 ₀	009 ₀	114 ₁	112 ₀	115 ₁	113 ₀	016 ₀
000 ₀	001 ₁	102 ₀	103 ₁	004 ₀	005 ₁	106 ₀	107 ₁	008 ₀	009 ₁	110 ₀	111 ₁	012 ₀	013 ₁	114 ₀	115 ₁	016 ₀

Table 3. Suffix Array permutation.

Vectors' bits (blue) are used for the next vector sorting (red bits).

References

- Ferragina, P., González, R., Navarro, G., & Venturini, R. (2009, February). Compressed Text Indexes: From Theory to Practice. *J. Exp. Algorithmics*, 13, 12:1.12--12:1.31.
doi:10.1145/1412228.1455268
- Grossi, R., Vitter, J. S., & Xu, B. (2011). Wavelet Trees: From Theory to Practice. *Proceedings of the 2011 First International Conference on Data Compression, Communications and Processing* (pp. 210-221). Washington: IEEE Computer Society.
doi:10.1109/CCP.2011.16
- Langmead, B., & Salzberg, S. L. (2012, April). Fast gapped-read alignment with Bowtie 2. *Nat Meth*, 9, 357-359. Retrieved from <http://dx.doi.org/10.1038/nmeth.1923>
- Li, H., & Durbin, R. (2009, July). Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics (Oxford, England)*, 25(14), 1754-1760.
doi:10.1093/bioinformatics/btp324