

**INTERNATIONAL ORGANISATION FOR STANDARDISATION  
ORGANISATION INTERNATIONALE DE NORMALISATION  
ISO/IEC JTC 1/SC 29/WG 11  
CODING OF MOVING PICTURES AND AUDIO**

**ISO/IEC JTC1/SC29/WG11 MPEG2017/M40271**  
**April 2017, Hobart, AU**

**Source:** Requirements

**Title:** Use cases for an efficient genomic information representation

**Authors:** Nicolas Guex, Christian Iseli, Thierry Schuepbach, Ioannis Xenarios (SIB/Vital-IT), Daniele Renzi, Giorgio Zoia (GenomSys), Claudio Alberti, Marco Mattavelli (EPFL), Jaime Delgado, Silvia Llorente, Daniel Naro (Distributed Multimedia Applications Group – Universitat Politècnica de Catalunya)

## **Table of Contents**

1	Purpose .....	2
2	Terminology .....	2
3	Use cases .....	3
3.1	Selective access to aligned data .....	3
3.2	Files concatenation .....	5
3.3	Genomic studies aggregation .....	6
3.4	Data streaming .....	7
3.5	Checking action conformity with privacy rules .....	8
3.6	Multiple alignments .....	10
4	Conclusions .....	10

## 1 Purpose

This document reports a collection of different use cases involving handling of genomic data; the document intends to clearly define modes of usage in terms of functionality of the sub-system and interactions with other (sub-) systems or actors during usage.

These use cases originate from the current procedures and needs of the daily genomic data processing activity that takes place at several bioinformatics data processing centers.

## 2 Terminology

Term	Definition
Alignment	A sequence read mapped on a reference DNA sequence
BAM	Compressed binary version of SAM
CRAM	GIR that includes SAM + Compression configuration
FASTA	GIR that includes header and sequence reads (nucleotides sequence)
FASTQ	GIR that includes FASTA + Quality Scores
GIR	Genomic Information Representation
Indel	An additional or missing nucleotide in a DNA sequence with respect to a reference DNA sequence.
MAF	Mutation Annotation Format. File format used to mark the genes and other biological features in a DNA sequence
Mate pairs	Two reads from the same (long) DNA strand extracted by sequencing machines. The orientation is the opposite of paired ends.
Paired ends	Couple of reads produced by the same (short) DNA fragment by sequencing both ends. The orientation is the opposite of mate pairs.
Quality score	A quality score is assigned to each nucleotide base call in automated sequencing processes. It expresses the base-call accuracy.
Read header	Each sequence read stored in FASTA and FASTQ format starts with a textual field called “header” containing a sequence identifier and an optional description
SAM	GIR that is human readable and includes FASTQ + Alignment and analysis information
Sequence read	The readout, by a specific technology more or less prone to errors, of a continuous part of a segment of DNA extracted from an organic sample

Some of the concepts described in this document are currently supported by tools used by MPEG members to process large datasets of genomic data on production systems used for scientific research. Other concepts have not been implemented yet but are required in order to enable usage scenarios currently not supported by existing formats such as SAM and CRAM.

The current Use Case were identified and described by:

- the Swiss Institute of Bioinformatics (SIB) and Vital-IT
- EPFL
- GenomSys
- DMAG-UPC

### 3 Use cases

#### 3.1 Selective access to aligned data

The types of selective access to aligned data described in this Use Case are important in genomic analysis. The new genomic information representation structure should support partitioning and efficient querying of data according to the criteria listed in the actions below.

<b>UCId</b>	GI-UC01
<b>Use Case Name</b>	Selective Access to Aligned Data
<b>Description</b>	Genomic Sequencing Information that has been aligned is accessed in many different ways during a typical analysis session.
<b>Actors</b>	Analyst
<b>Assumptions</b>	The sequencing data to be analysed have been aligned prior to the current session and stored in a suitable format
<b>Actions</b>	<ol style="list-style-type: none"><li>1. The analyst opens the file</li><li>2. The analyst filters unmapped reads, cuts them from the current file and saves them in a new file.</li><li>3. The analyst filters duplicate reads and removes them from the file.</li><li>4. The analyst filters reads with no mismatches and remove them from the current session view</li><li>5. The analyst filters reads with mismatches only (no indels) and decodes reads with a number of mismatches of 3 or less</li><li>6. The analyst filters reads with mismatches, indels and soft clips and decode reads with a distance of X or less</li><li>7. The analyst browses the selected read subsets and selects regions of interest in the genome being analysed</li><li>8. The analyst labels the selected sub-regions and saves the results in a new file.</li></ol>
<b>Inclusion or Extensions</b>	The analyst reloads the previously saved file and can easily decode and access only the sub-regions selected in the previous session.
<b>Issues</b>	

The typical actions of an analysis session described above have precise rationales for the bioinformaticians. They are made more explicit in the following table:

	Description	Rationale
0	Efficiently retrieve both reads to a pair and reconstruct a FASTQ file equivalent to the original input	Useful to remap with a different tool, or visualize aligned pairs on a genomic browser
1	Access reads with no mismatches	Useful to determine the coverage of any particular nucleotide or genomic region (CNV, RNAseq) but not as a source for the discovery of new variants
2	Access reads with mismatches only	They contain the most useful information for SNP calling or to identify premature stop codons
2.1	Decode only reads with a number of mismatches below/above a given threshold	Used in variant calling to filter reads with too many mismatches which may indicate poor sequencing quality and / or incorrect mapping
3	Access reads with indels	They contain the most useful information to detect frameshifts, or to identify chromosome rearrangements
4	Access unmapped reads only	Useful to detect potential contamination or infections
5	Group genomic sub-regions uniquely identified with user-defined criteria	
5.1	Count the reads (or pairs in the case of paired-end experiments) that fall within each defined region	Used in gene expression analysis
5.2	Extract all the reads with mismatches or indels that fall within a specific region	Used to perform a local re-alignment or assembly of the reads to improve SNP calling or indel detection
6	Remove/filter duplicates	Crucial to prevent over-counting PCR effects during expression analyses
7	Perform the query on several genomic datasets (i.e. the results of a sequencing run)	Used in population genetics, for example to count the penetrance of a given trait / SNP
8	Remove or properly weight reads that map equally well at more than one location	Used to normalize gene expression analysis
9	Access pair-end reads whose ends map on distinct chromosomes	Used to detect cancer chromosome rearrangements and also used in HiC experiments

### 3.2 Files concatenation

It should be possible to concatenate compressed files without the need to decompress and recompress them. This is particularly useful to transmit, merge and aggregate small subparts of a sequencing experiment, such as gene panels or regions that cover a few SNPs of interest. It would also be very useful in the common case where a library is re-sequenced in order to get additional sequencing depth. It however causes potential issues with duplicate reads (see issues)

<b>UCId</b>	GI-UC02
<b>Use Case Name</b>	Files concatenation
<b>Description</b>	Concatenate compressed files without the need to decode and re-encode them
<b>Actors</b>	Analyst
<b>Assumptions</b>	
<b>Actions</b>	<ol style="list-style-type: none"><li>1. The analyst obtains a genomic sample of raw reads from a human individual</li><li>2. The analyst performs alignment and compression against reference sequences representing the 24 human chromosomes (22 + X and Y)</li><li>3. The alignment is performed in parallel per each chromosome and the output is a compilation of 24 compressed files. The analyst opens the 24 files and concatenates them into a single compressed bitstream very quickly with no need of decompression and re-compression of each individual file.</li></ol>
<b>Inclusion or Extensions</b>	<p>Other examples where files concatenation is useful is when different “experiences” are performed on the same genomic sequence data and the results shall be concatenated. For example alignment of the same raw data using different configurations of the same tool or different alignment approaches.</p> <p>Another application is the common case where a library is re-sequenced in order to get additional sequencing depth.</p>
<b>Issues</b>	

### 3.3 Genomic studies aggregation

It should be possible to encapsulate in the same compressed file several related studies that can be separately accessible. Additionally transversal queries on all the compressed files should be possible (e.g. “select chr1 of all compressed samples”). This is particularly useful when a study is performed on large populations of individuals of the same species or when the same individual is sequenced/analyzed several times during his life time.

For example when aggregating several datasets from the same individual/species, the available reference sequences (e.g. human genome) should be stored only once with a significant gain in storage efficiency in case of studies on large populations (up to several thousand individuals) or when sequencing several times a single individual.

<b>UCId</b>	GI-UC03
<b>Use Case Name</b>	Genomic Studies Aggregation
<b>Description</b>	Encapsulate in the a compressed file several related studies that are maintained separately accessible
<b>Actors</b>	Analyst 1, Analyst 2
<b>Assumptions</b>	The analyst has performed various types of analysis (such as alignment, variant calling, gene expression analysis) on different datasets of genomic sequence data. The results are in the form of compressed genomic bitstreams containing sequence data and related annotations/metadata.
<b>Actions</b>	<ol style="list-style-type: none"> <li>1. Analyst 1 opens the several files he has available from previous analysis sessions. <ol style="list-style-type: none"> <li>a. The reference sequences used for alignment should be stored only a minimum number of times (one in case of single species).</li> </ol> </li> <li>2. Analyst 1 selects the previous studies that he/she wants to aggregate in a single file</li> <li>3. Analyst 1 saves the aggregation and transmits the file to Analyst 2</li> <li>4. Analyst 2 opens the file he received from Analyst 1</li> <li>5. Analyst 2 browse one or more specific dataset of the aggregation and all the related results</li> <li>6. Analyst 2 identifies an interesting result and browse all the datasets of the aggregation for same or similar types of results.</li> <li>7. Analyst 2 specifies a matching function on one region of the available</li> </ol>

	<p>datasets for results such as:</p> <ol style="list-style-type: none"> <li>i. variant calling on chr2 of all dataset</li> <li>ii. retrieve all dataset where a specific mutation is present above a given threshold</li> <li>iii. count the number of genes in all datasets and retrieve only those above a given threshold</li> </ol>
<b>Inclusion or Extensions</b>	<p>Each aggregated dataset should be looked as one dimension of a “hypercube” with the possibility to randomize access according to any (set of) dimensions. This can be seen as well in terms of "experiences" (i.e. with several concatenated files) where the "hypercube" can be browsed according to dimensions such as genomic region, experience, individual, etc.</p>
<b>Issues</b>	

### 3.4 Data streaming

<b>UCId</b>	GI-UC04
<b>Use Case Name</b>	Data Streaming
<b>Description</b>	Transmission of compressed data from a datacenter for analysis with real-time access
<b>Actors</b>	Sequencing operator, Analyst
<b>Assumptions</b>	
<b>Actions</b>	<ol style="list-style-type: none"> <li>1. At the sequencing facility the operator launches a sequencing process</li> <li>2. Raw sequence data are produced and stored on a local storage device</li> <li>3. Raw sequence data are compressed and the operator starts a streaming session to an analysis center that subscribed to receive the data, as soon as data are available on the local storage device</li> <li>4. The analyst at the analysis center starts receiving data from the open link with the sequencing facility</li> <li>5. The analyst starts analyzing the received data as soon as they are received, even before sequencing and streaming are completed. For</li> </ol>

	<p>example:</p> <ol style="list-style-type: none"> <li>a. alignment and variant calling against an available reference</li> <li>b. if a mutation is detected above a given threshold an alarm is triggered (this could provide the required result before the sequencing process is completed)</li> </ol> <ol style="list-style-type: none"> <li>6. The analyst forward some portions of the data being received to a subcontractor for further analysis and/or specific refinements <ol style="list-style-type: none"> <li>a. These data may be selected based on e.g. some specific classification such as non-mapping data or some specific given threshold as above, etc.</li> </ol> </li> <li>7. The analyst terminates his/her analysis after all data have been received from the sequencing facility</li> <li>8. The analyst waits for further results from the subcontractor and aggregates them with his/her results into a single file.</li> </ol>
<b>Inclusion or Extensions</b>	
<b>Issues</b>	

### 3.5 Checking action conformity with privacy rules

<b>UCId</b>	GI-UC05
<b>Use Case Name</b>	Checking action conformity with privacy rules
<b>Description</b>	The system checks if the action the user intends to perform is compliant with the rules expressed by the owner of the data.
<b>Actors</b>	Data consumer (analyst, practitioner, service provider, forensic scientist).
<b>Assumptions</b>	<p>The sequenced data are linked to privacy rules, defining which conditions the owner imposes on its access and use.</p> <p>The system might operate locally at the user's premises or remotely (for example on the data custodian's repository)</p>
<b>Actions</b>	1. The data consumer submits the relevant information about the action s/he wants to be performed:

	<ul style="list-style-type: none"> <li>- the intended use (not for further use, commercial use, relatives search, forensics, scientific purposes)</li> <li>- her/his identity as contact information and role (analyst, researcher, healthcare, forensic, not specified)</li> <li>- the action to perform (genetic analysis for diagnostic, analysis for a research project, view information, paternity test, forensic use)</li> <li>- to whom the result will be forwarded (contact information and role)</li> <li>- the set of genomic regions where the action will be performed</li> </ul> <ol style="list-style-type: none"> <li>2. The system checks if the provided user information can be verified</li> <li>3. The system checks if the provided privacy rules can be verified</li> <li>4. The system checks if the provided information meet the privacy rules, including the different uses regulated (not for further use, not for commercial use, not for relatives search use, not for forensics use, forensics use, open access, only for scientific purposes, ...), and performs necessary actions, such as informing the owner or requesting a further confirmation from the owner</li> <li>5. Depending on the result of privacy rules evaluation, the system performs the requested action</li> <li>6. The system notifies the owner about the result, if requested.</li> </ol>
<b>Inclusion or Extensions</b>	In order to enforce the application of the privacy rules, data have to be encrypted and the process has to request the necessary information for decryption.
<b>Issues</b>	

### 3.6 Multiple alignments

It is necessary to support and properly weight reads that map equally well at more than one location. This is used e.g. to normalize gene expression analysis.

The current existing format supporting Multiple Alignments is the Multiple Alignment Format (MAF, <https://cgwb.nci.nih.gov/FAQ/FAQformat.html>).

<b>UCId</b>	GI-UC06
<b>Use Case Name</b>	Multiple alignments
<b>Description</b>	Many genes are duplicated in the genome, and there are many gene families where differences between genes from the same family at the DNA level cannot be determined from short reads, if at all. In such cases, it is common practice to spread the weight of reads where there are multiple mapping positions among all possible positions/genes. In some cases, those reads can also be excluded from the downstream analysis
<b>Actors</b>	Analyst, alignment tool
<b>Assumptions</b>	A scoring scheme is adopted to measure the <i>quality</i> of each alignment. Existing scoring schemes use: <ul style="list-style-type: none"><li>• points to each aligned base</li><li>• probability (normalized to 1 among all alignments)</li></ul>
<b>Actions</b>	<ol style="list-style-type: none"><li>1. The analyst opens files containing non-aligned sequencing data</li><li>2. The analyst selects the scoring scheme in case of multiple alignments</li><li>3. The analyst runs the first alignment pass</li><li>4. When multiple mapping positions are found, each alignment is assigned a score according to the selected scheme</li><li>5. The analyst saves the alignment results in a suitable format allowing the storage of sensible information for all multiple alignment (incl. MAF)</li></ol>
<b>Inclusion or Extensions</b>	In case of multiple pass alignments, an additional integer number is associated to each alignment: usually pass 1 indicates the “strongest” alignment and additional passes indicate the “weaker” ones
<b>Issues</b>	It would be important to store the configuration of the tool or toolset used to generate the multiple alignments. This should probably be encoded as metadata associated with the alignment information.

## 4 Conclusions

Use Cases presented in this document should help the specification of a new file format and genomic data representation and partitioning presented in this document. They have been proposed validated by organizations working on genomic data such as the Swiss Institute of Bioinformatics.

Currently the used compressors are general purpose entropy coders, but the development of compressors exploiting the nature of genomic information would certainly provide better performance.