

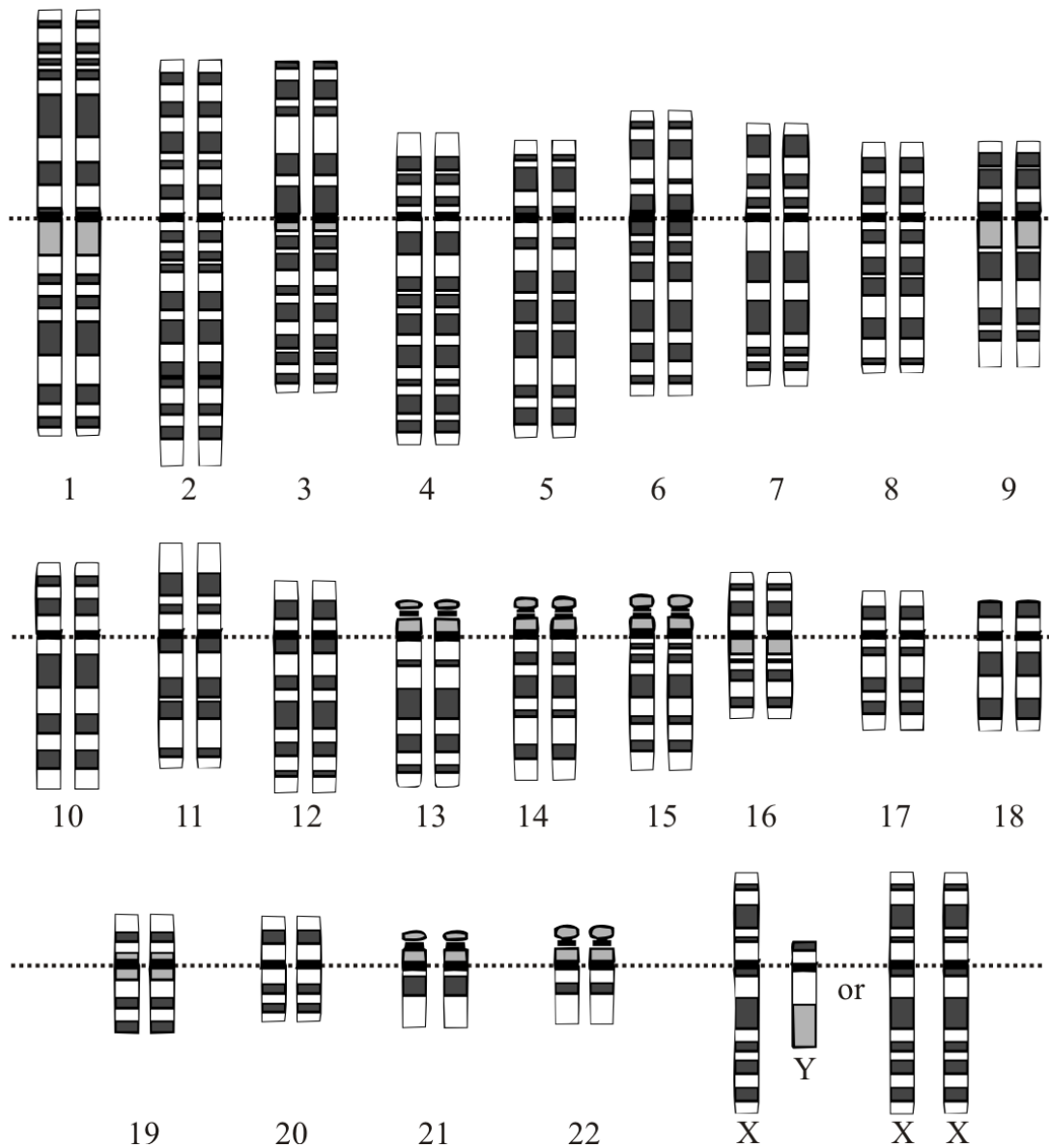
MPEG-G: The Standard for Genomic Information Representation

Jan Voges

Institut für Informationsverarbeitung



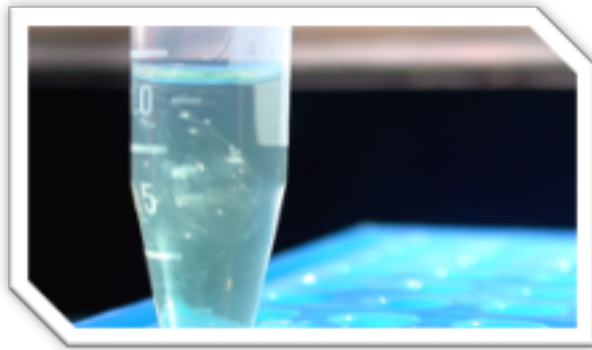
Human karyogram



- Diploid genome
- 46 chromosomes
- 22 pairs of autosomes (1-22)
- 1 allosome pair (XY | XX)
- 4 bases (A, C, G, T)
- ~3 billion base pairs

DNA sequencing

Human genome: ~ 3 billion base pairs \times 2 bits per base = **~ 750 MB**



+ read-out redundancy
 ~ 500 GB

+ meta information
 ~ 1 TB

+ alignment information
 ~ 1.5 TB

Whole genome sequencing

Chromosome



↓ fragmentation



↓ sequenced

Reads



Alignments

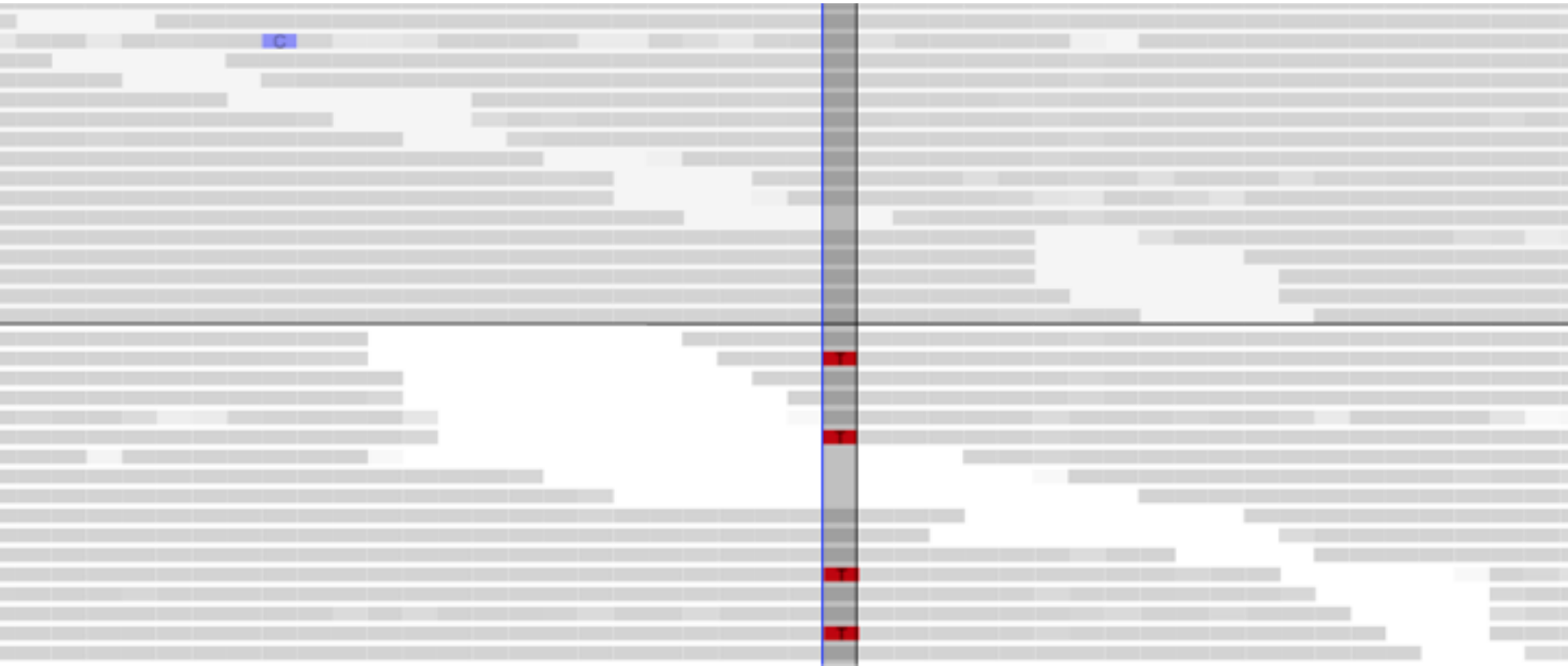
GCTATCAGGCTAGGTTA GTTACAGTGCATGCATA CATACACGTAGCTATACG

↓

Assembly

GCTATCAGGCTAG GTTACAGTGCATGCATA CATACACGTAGCTATACG

Alignment



Evolution of genome sequencing

Sequencing technology

	2009	2017/2018
Cost/genome	\$100k	\$1k
Coverage	~30x	> 200x
Number of reads	~1 billion	> 6 billion
Size of raw sequencing files	~0.25 TB	> 1.5 TB

Storage & transmission infrastructure

	2009	2017/18
Cost/TB	\$100	\$50
Download speed	10 Mbps	100 Mbps

No technology is keeping with the pace of genome sequencing!

What is MPEG-G?

- MPEG-G = International Standard ISO/IEC 23092
- Largest coordinated and international effort addressing the problems and limitations of current technologies
- Paves the road towards a truly efficient and economical handling of genomic information
- Utilizing the latest technologies
- Planned release: end of 2018



- **Interoperable selective access to data in the compressed domain** by means of standard APIs:
 - Genomic region
 - Class of data
 - ...
- On top of compression **higher performance is provided by a specific file format and transport format.**



Support the evolution of a SW/HW ecosystem so that compression technology becomes a commodity for the users.

Benefits provided by MPEG-G

Selective access to
compressed data

Data streaming

Genomic studies
aggregation

Enforcement of
privacy rules

Selective encryption
of sequencing data
and metadata

Annotation and
linkage of genomic
segments in the
compressed domain

Interoperability with
main existing
technologies and
legacy formats

Incremental update
of sequencing data
and metadata

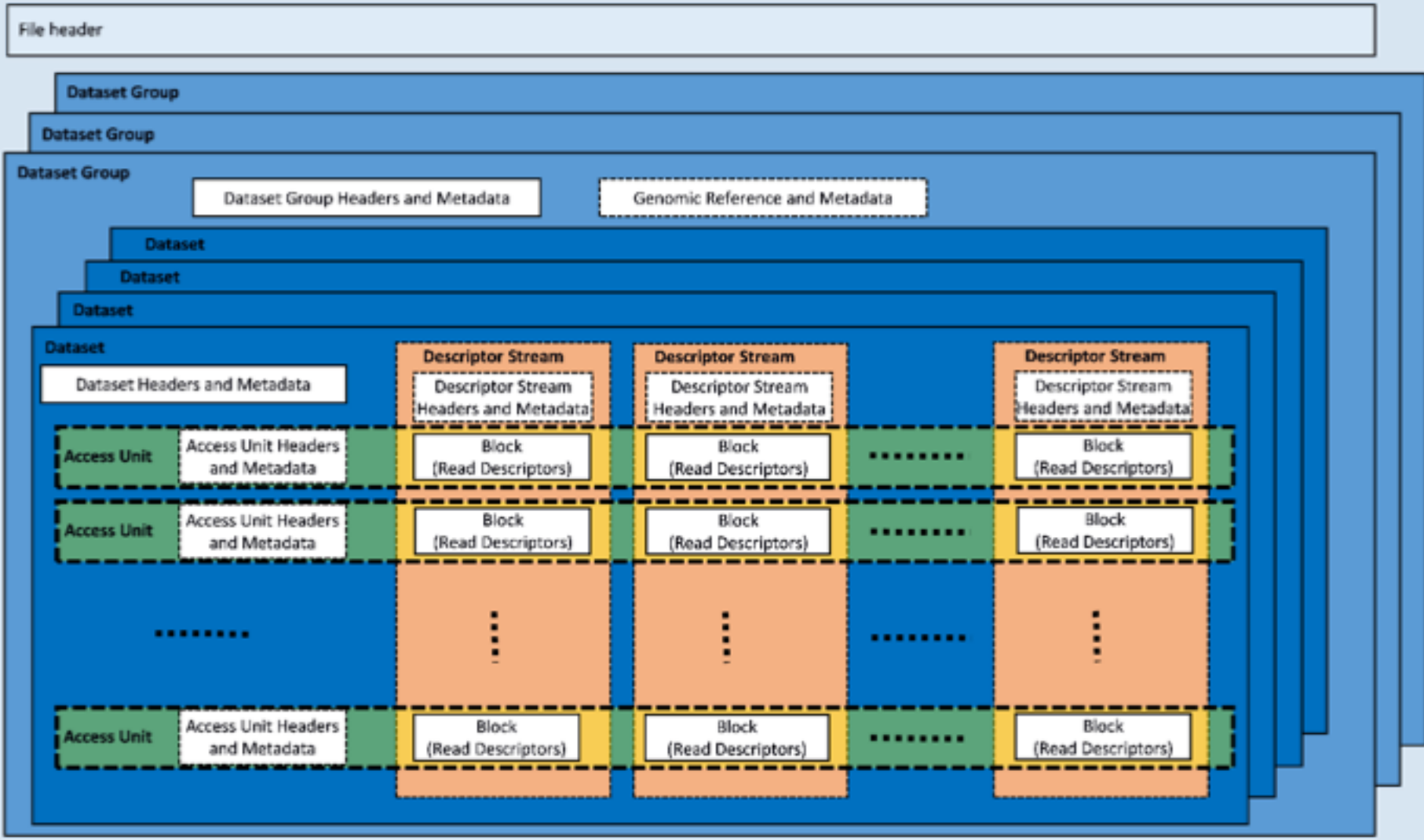
Compressed file
concatenation

- **Part 1: File and Transport Format**
 - The technology to transport and access data
- **Part 2: Genomic Information Representation**
 - The compressed representation
- **Part 3: APIs**
 - The standard interfaces with genomic data applications and legacy formats
- **Part 4: Conformance**
 - The methodology to test compliance with the standard
- **Part 5: Reference Software**
 - The standard support to the implementation of applications

- **Part 1: File and Transport Format**
 - The technology to transport and access data
- **Part 2: Genomic Information Representation**
 - The compressed representation
- **Part 3: APIs**
 - The standard interfaces with genomic data applications and legacy formats
- **Part 4: Conformance**
 - The methodology to test compliance with the standard
- **Part 5: Reference Software**
 - The standard support to the implementation of applications

MPEG-G file format

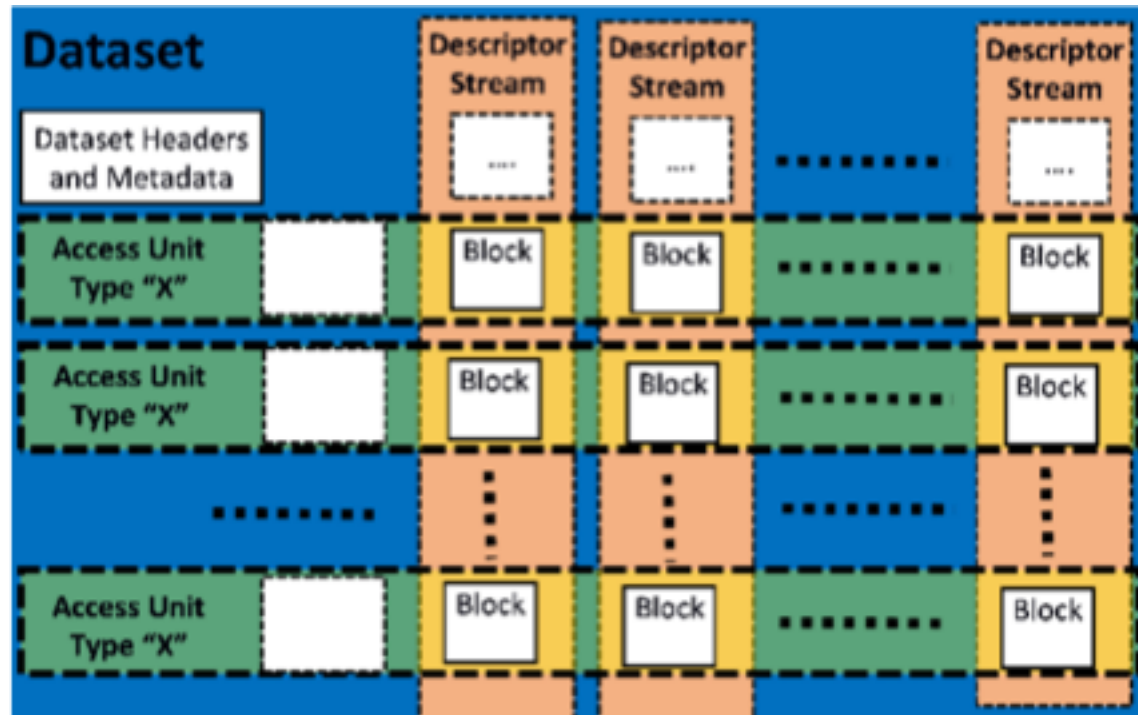
FILE FORMAT



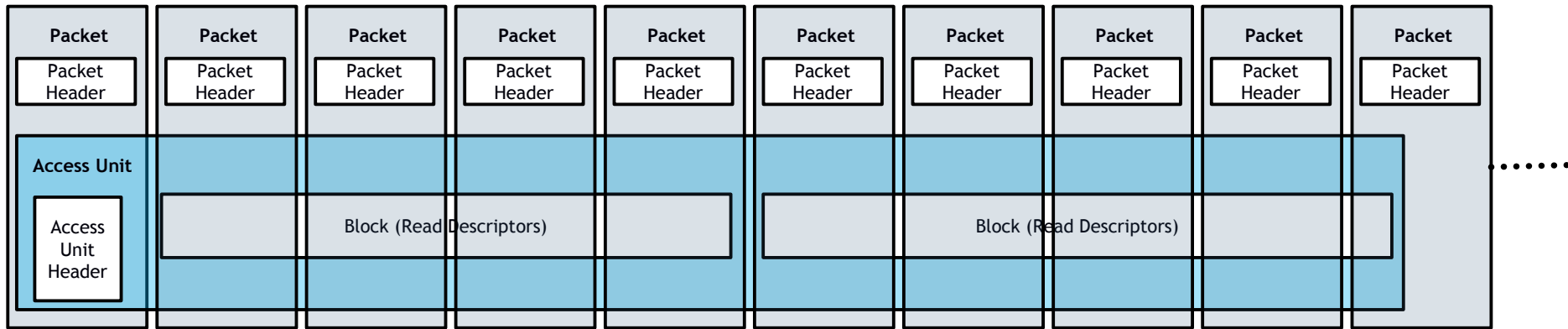
Selective access

The indexing tools embedded in an MPEG-G file enable several types of selective access that can be combined in the same query, e.g.:

- Genomic interval in terms of start to end mapping position on a given reference sequence
- Type of data (i.e., a single data class)



TRANSPORT FORMAT

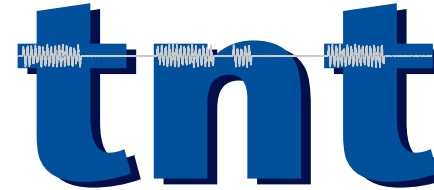


MPEG-G streaming features:

- Packet size adaptation to the channel characteristics/state
- Error detection and support of re-transmission of erroneous/incomplete data for error-free delivery
- Support of out-of-order delivery
- Packet-based filtering of genomic data
- Full convertibility of file and transport formats

- **Part 1: File and Transport Format**
 - The technology to transport and access data
- **Part 2: Genomic Information Representation**
 - The compressed representation
- **Part 3: APIs**
 - The standard interfaces with genomic data applications and legacy formats
- **Part 4: Conformance**
 - The methodology to test compliance with the standard
- **Part 5: Reference Software**
 - The standard support to the implementation of applications

- **genie** = **GEN**omic **I**nformation **E**ncoding
- Joint collaborative effort to produce a standard-compliant open source encoder





genie

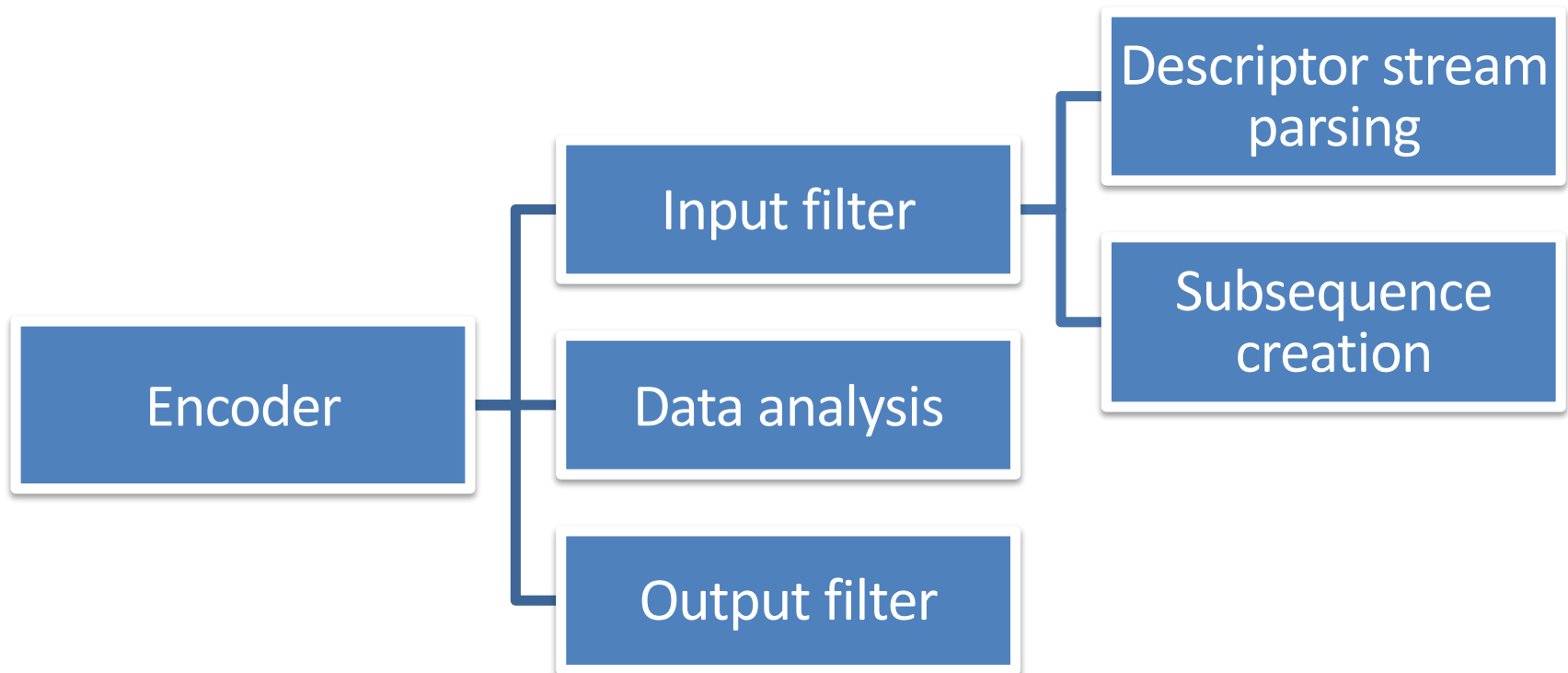
FF/TF
(Part 1)

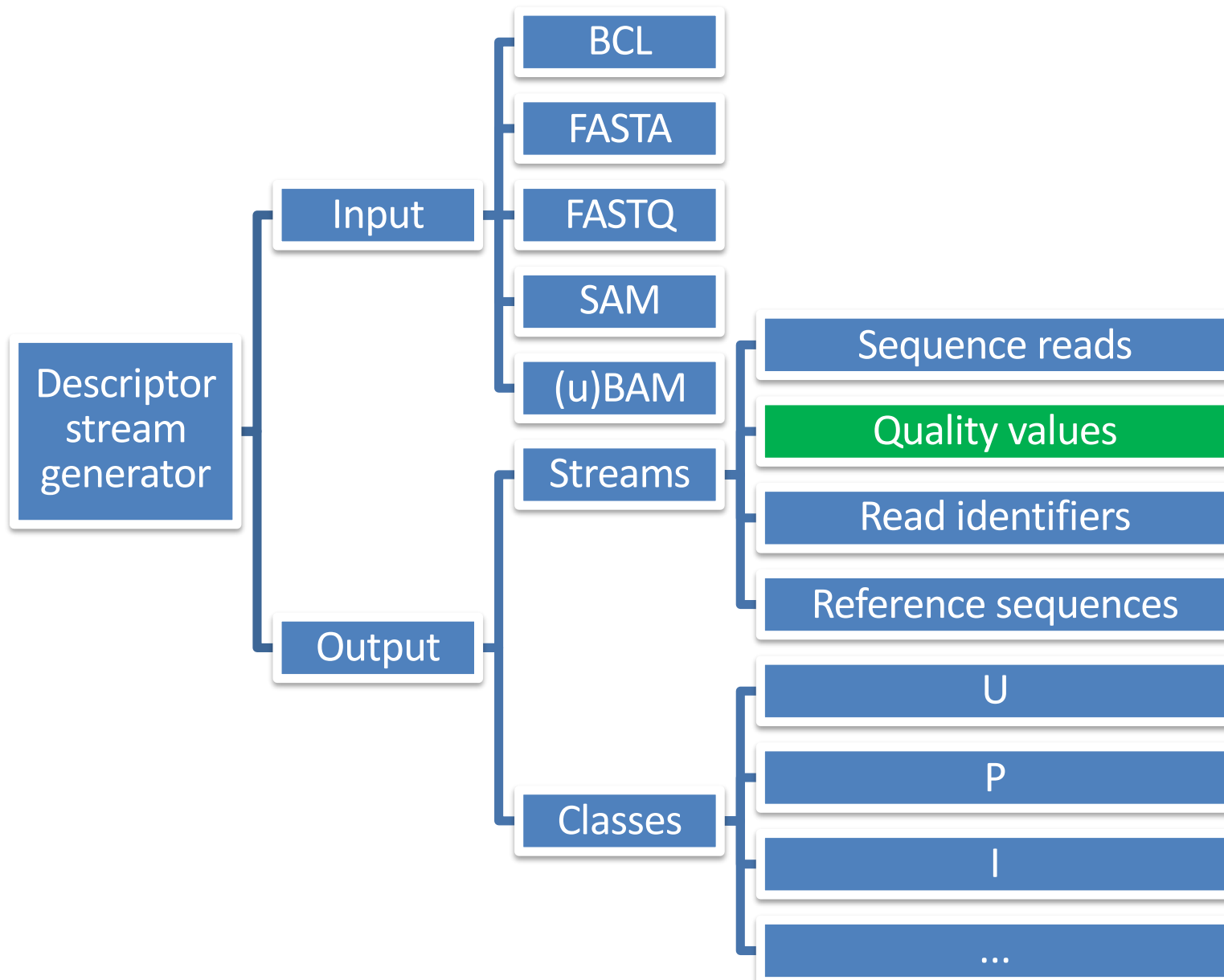
Genomic Information
Representation
(Part 2)

APIs
(Part 3)

Descriptor
stream
generator

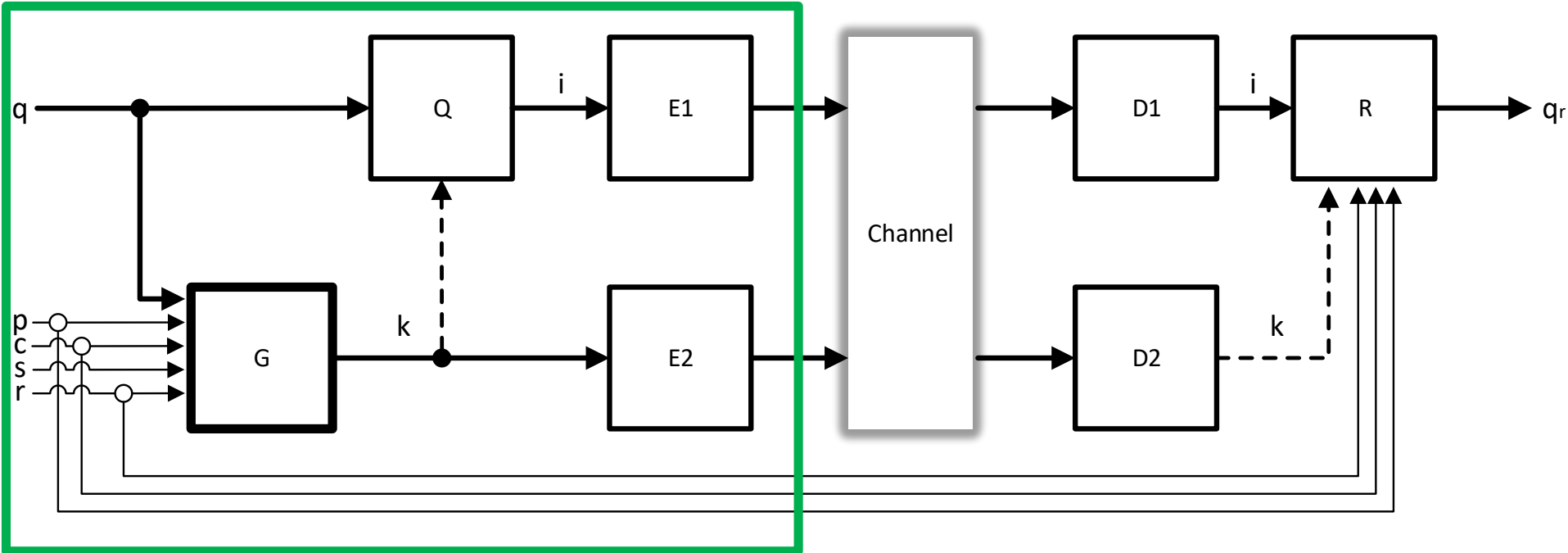
Encoder





Coding of quality values

Year	Tool
2011	<ul style="list-style-type: none">• SlimGene (Kozanitis et al.)
2012	<ul style="list-style-type: none">• SCALCE (Hach et al.)
2013	<ul style="list-style-type: none">• QualComp (Ochoa et al.)• BEETL (Janin et al.)• Fastqz (Bonfield et al.)
2014	<ul style="list-style-type: none">• Illumina's binning• P-/R-Block (Cánovas et a.l.)
2015	<ul style="list-style-type: none">• Quartz (Yu et al.)• QVZ (Malysa et al.)
2016	<ul style="list-style-type: none">• Crumble (Bonfield)
2017	<ul style="list-style-type: none">• CALQ (Voges et al.)
2018	<ul style="list-style-type: none">• QSComp (Voges et al.)

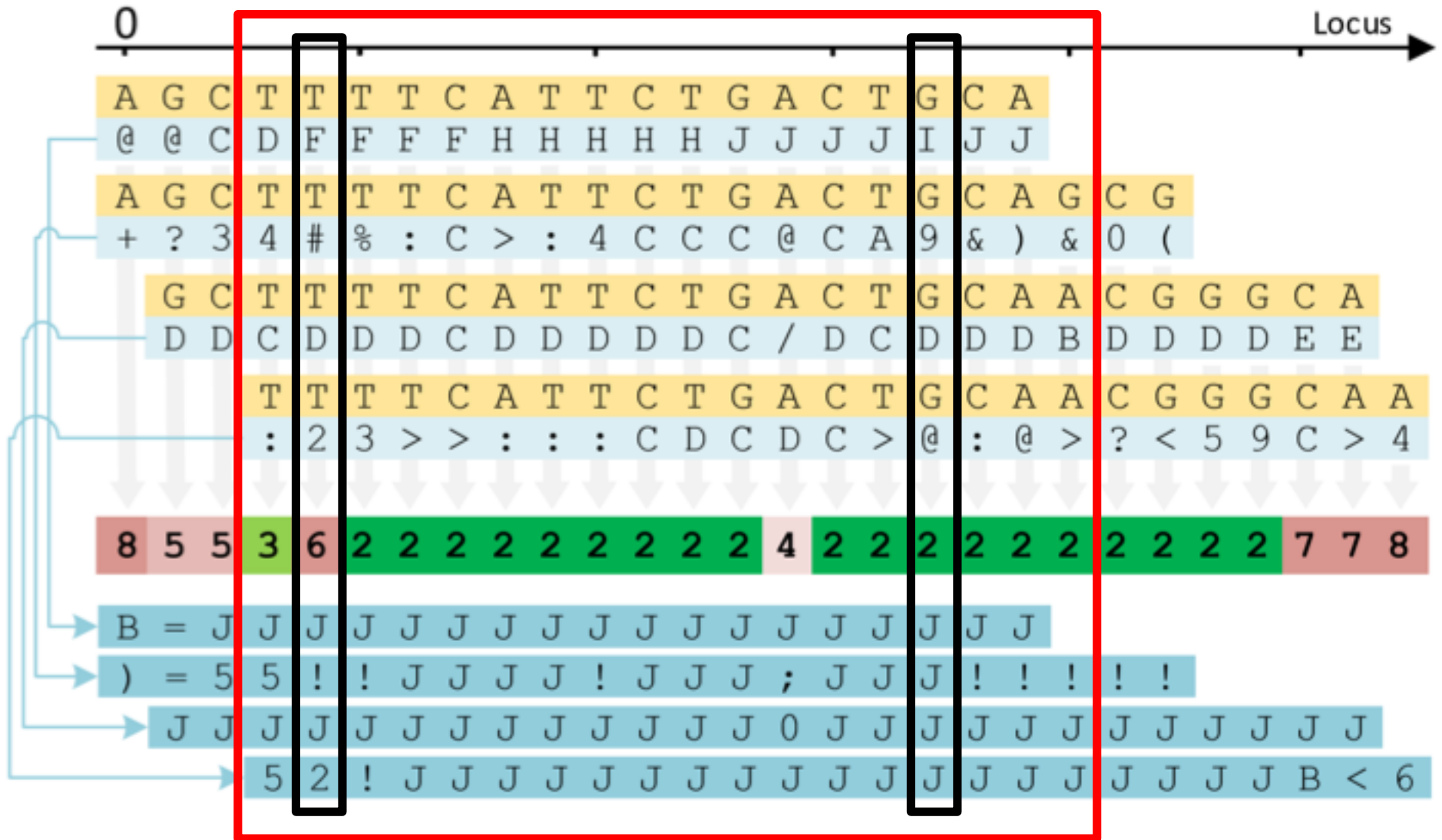


Legend

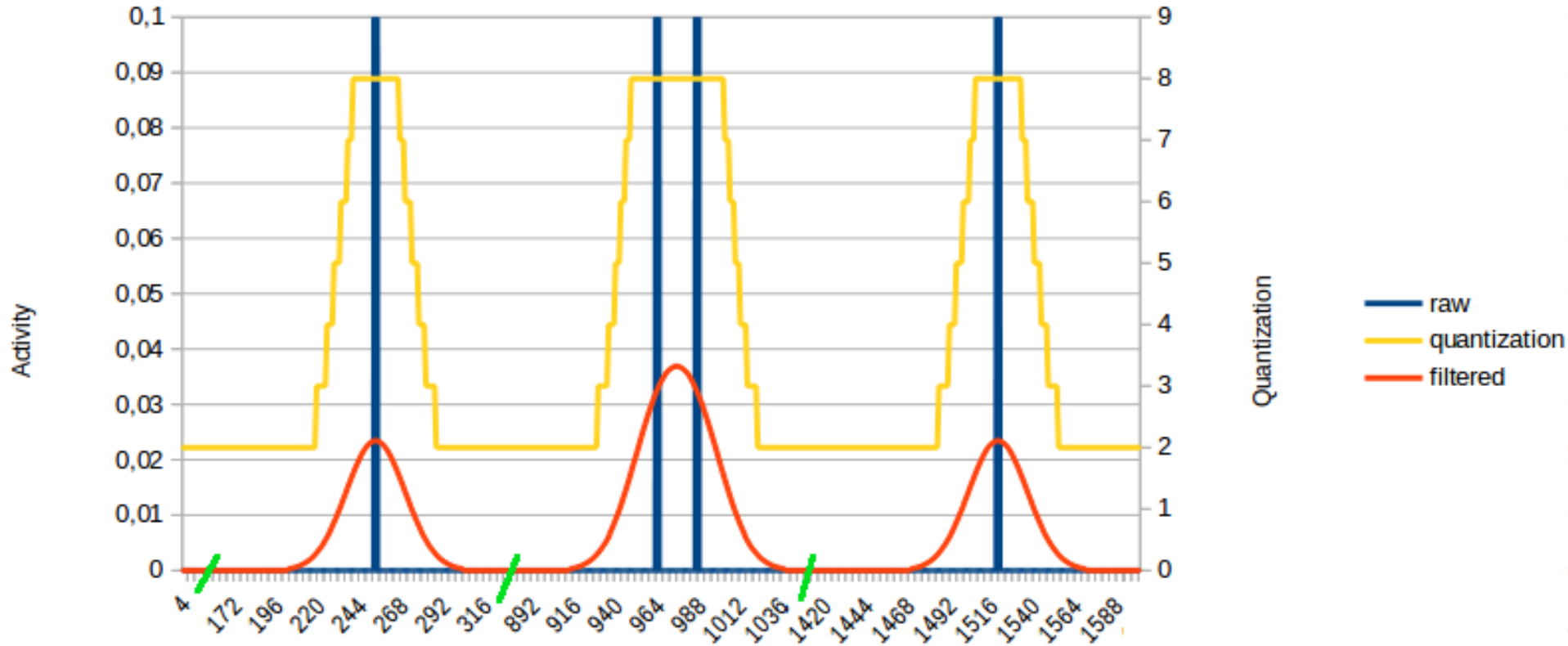
q	quality values
q _r	reconstructed quality values
p, c, s, r	side information (mapping positions, CIGAR strings, donor sequences, reference sequence(s))
k	quantizer indexes
i	quantization indexes

	signal
	control signal
	side signal

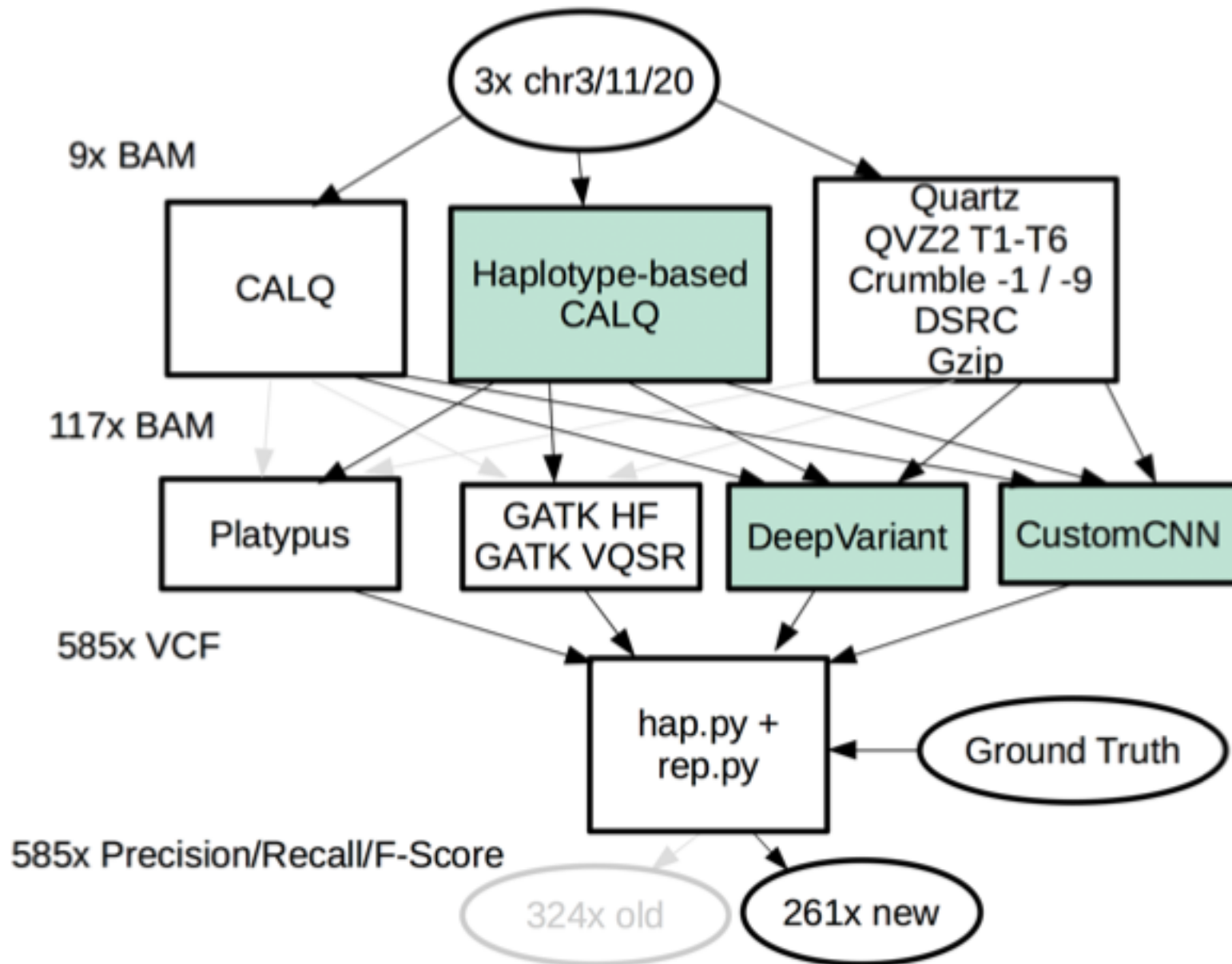
CALQ



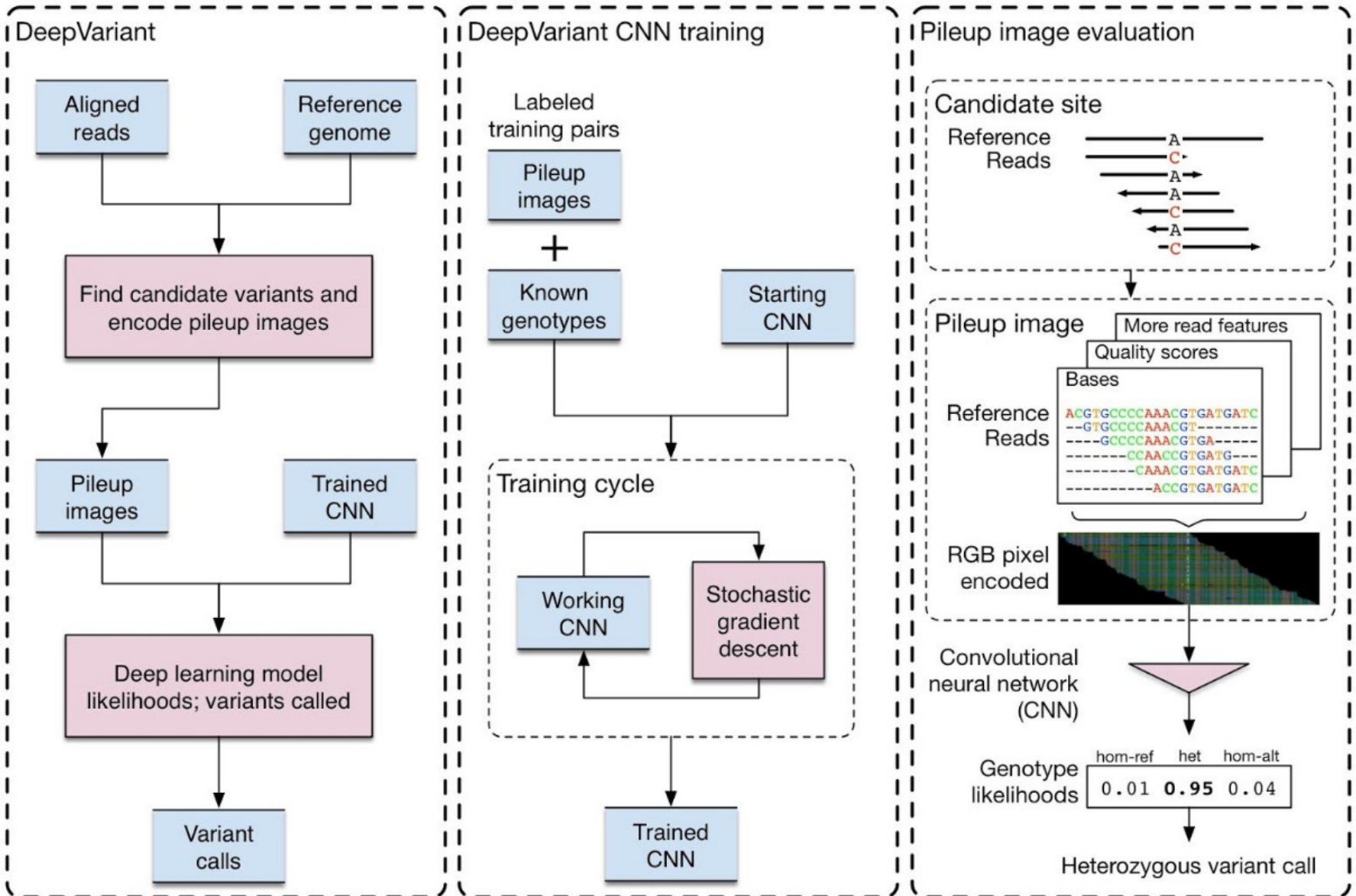
CALQ v2: sequence activity and quantization



Evaluation of lossy quality value compression



Variant calling using neural networks

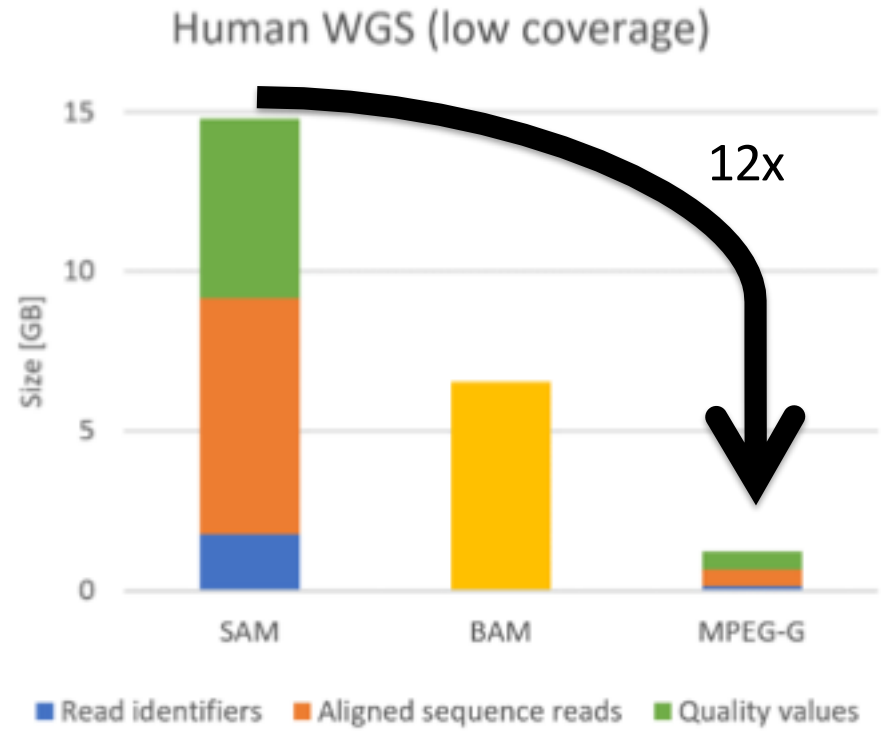
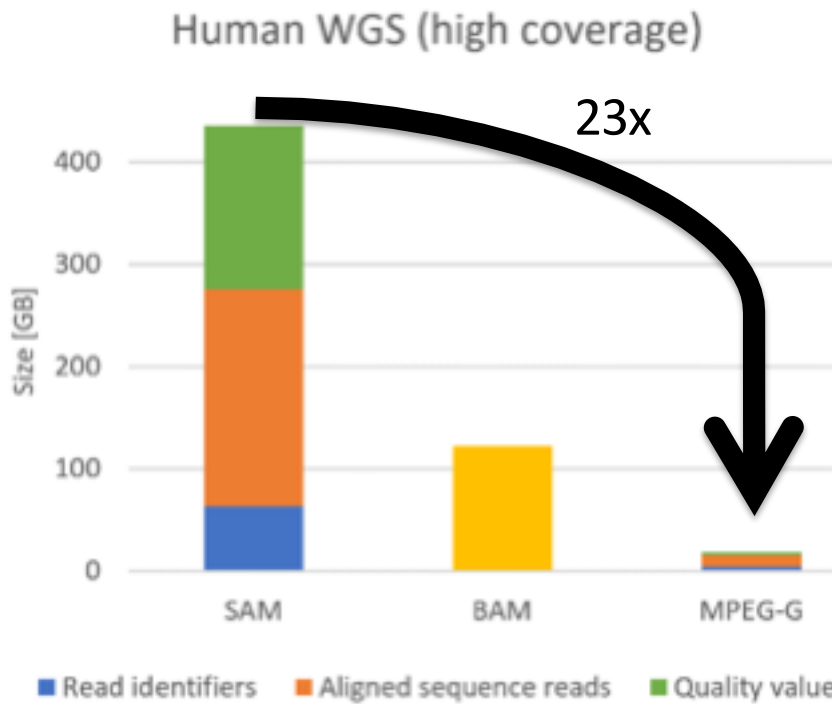


RD diagram for quality values

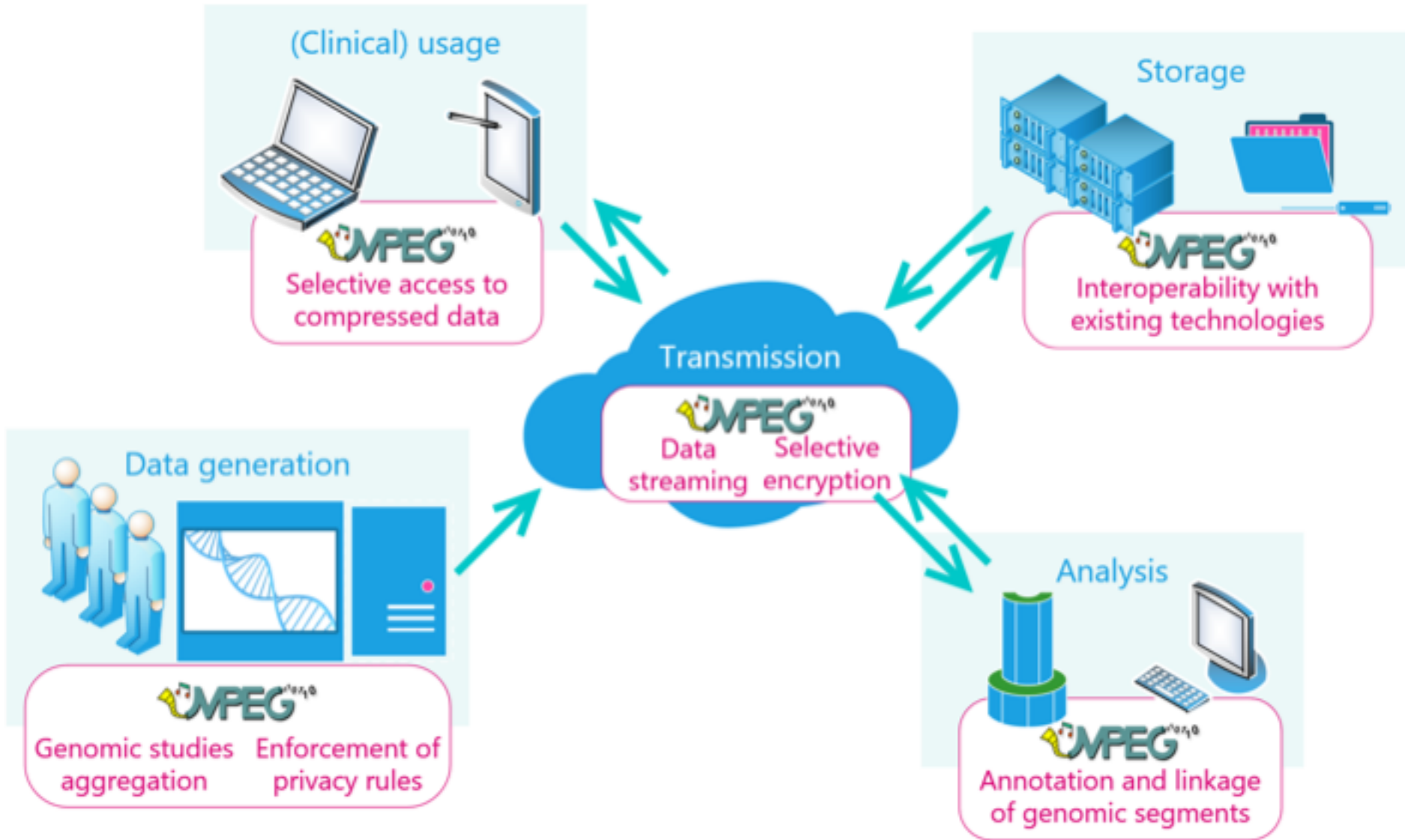
Average F-score difference w.r.t. original data versus bits per quality value



MPEG-G performance



A genomic ecosystem fueled by MPEG-G



And if we try to guess the future of genomic data:

- The technology to compress genomic information will **change over time**
- Genomic information compression performance will **improve over time**
- The MPEG-G Systems technologies will evolve and improve, but **the main functionality will stay** and support the evolution of analysis application

A new logo?

