# The GeneROOT Project Status and Plans

Fons Rademakers

CERN openlab Chief Research Officer
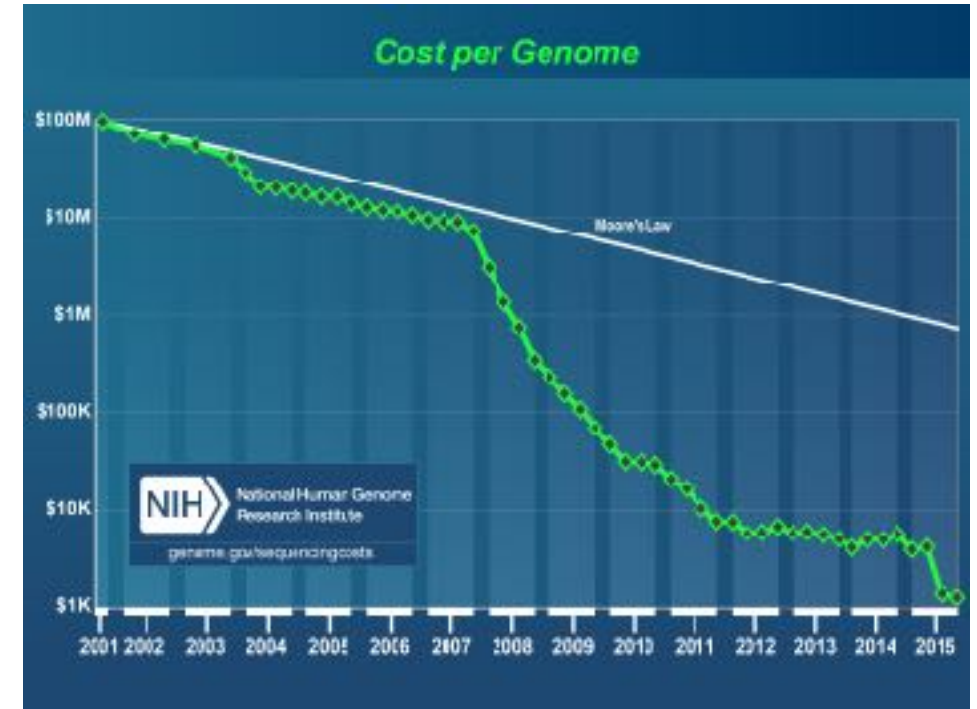
# GeneROOT - Using ROOT for Handling Genomics Data

# King College London - TwinsUK Project

- Collaboration between the KCL and CERN openlab

- Try to optimize genomics data storage and processing using HEP tools

- Working on a local 400TB copy of TwinsUK data
  - 750 Monozygotic twins
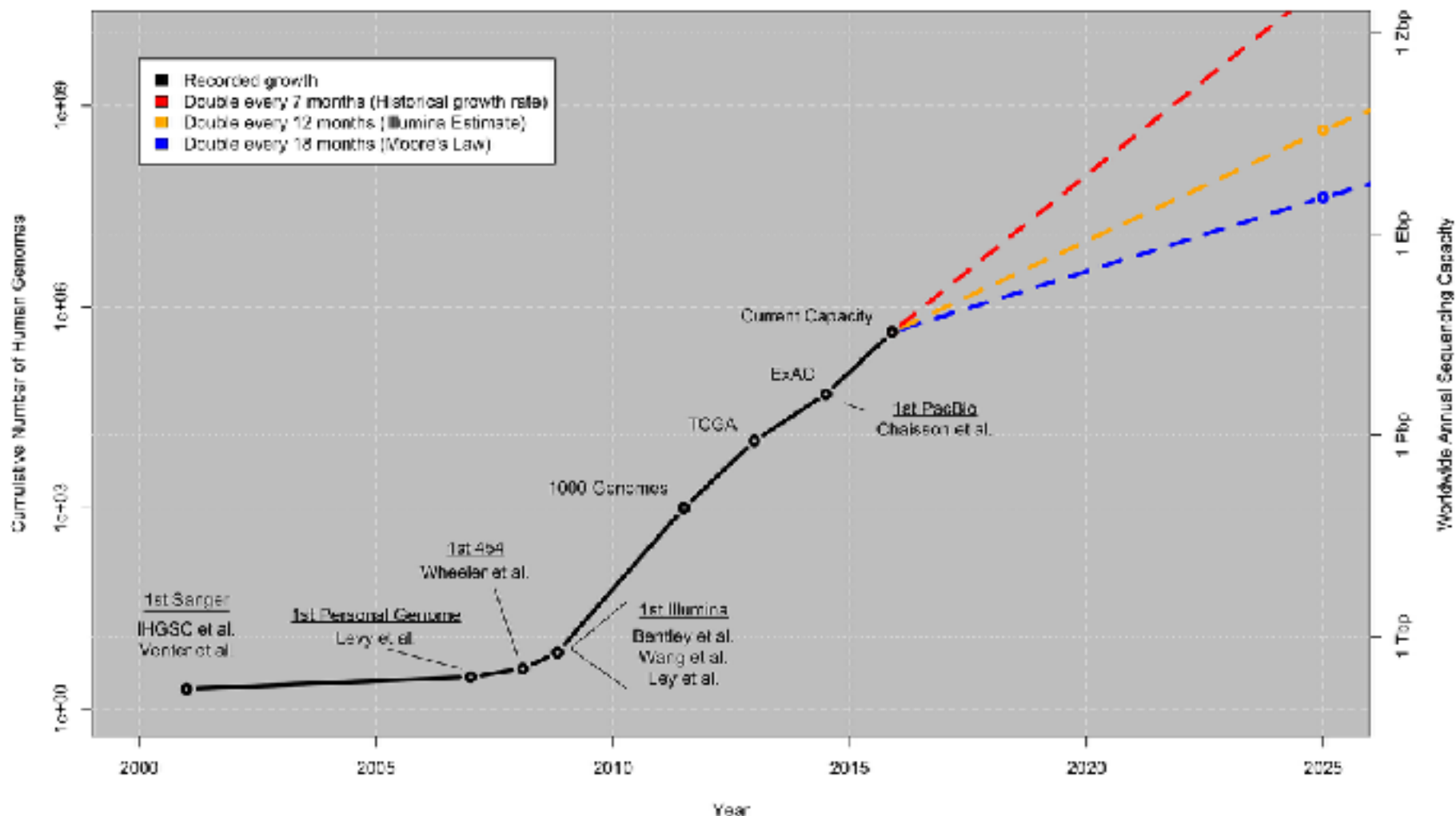  - 900 Dizygotic twins
  - 138 Singletons

# Rapidly Increasing Amount of Genomics Data

- Next generation Sequencing (NGS)
  - Dramatic increase in the amount of data
  - Improved data confidence
- NGS is enabler for more sophisticated research questions in Genomics



**Issue: Leaps in sequencing technology have outperformed advances in computing**

Growth of DNA sequencing data both in terms sequenced human genomes and total sequencing capacity

Z. D. Stephens, S. Y. Lee, F. Faghri, R. H. Campbell, C. Zhai, M. J. Efron, R. Iyer, M. C. Schatz, S. Sinha, and G. E. Robinson, "Big data: Astronomical or genomical?", Plos biology, vol. 13, no. 7, e1002195, 2015.

# SAM - Sequence Alignment/Map

- Text-based format for storing biological sequences aligned to a reference one



- SAM - Plaintext
- BAM - Binary Compressed

Challenge: data size is extremely large (about 500GB for a single human genome)

Genomic data size will overtake LHC data in the coming years

# SAM Example

```
@SQ      SN:chrM LN:16571
@SQ      SN:chrX LN:155270560
SOLEXA-1GA-2_2  0   chr1  10145  25  36M  *  0  0  AACCCCTAACCCTAACCCTAACCCTA  hhhhHcWhhHTghcKA_ONhAAEEBZ
SOLEXA-1GA-2_2  0   chr1  10148  25  36M  *  0  0  CCCTAACCCTAACCCTAACCCTAACC  hbfhhhXUYhT_ULZdLRTKNIMIKG  NM:i:0
SOLEXA-1GA-2_2  16  chr1  10149  25  36M  *  0  0  CCAAACACTAACCCTAACCCTAACCC  ><<>B@>?>?D>>?B?D>DBC?E@BDH  NM:i:1  X1:i:1
SOLEXA-1GA-2_2  0   chr1  10150  25  36M  *  0  0  CTAACCCTAAACCTAACCCTAACCCT  hhW_X]MXNOHQQWMILHGIFMJGJ
```

The file consists of a header section (lines starting with a @)
and an alignment section with 11 mandatory fields + several optional fields

A human genome SAM file consist of about 1 billion lines in the alignement section and is about 500GB large

# BAM - Binary Alignment/Map

- BGZF is block compression implemented on top of the standard gzip file format

- File is gunzip compatible

- Each 64KB block of data is compressed and added to the file

- Using an .bai index file, random access is supported in the BAM file

# ROOT Framework

Don't reinvent the wheel



Most of the challenges with massive data processing are common to HEP. **ROOT** has decades of experience in software design and optimisation

# Conversion - SAM to ROOT

We convert from TAB-delimited SAM format to a ROOT TTree



ROOT handles compression, memory buffers, datatypes, endian-ness, etc.

The columns are defined using a C++ class.

# Performance - SAM to ROOT

There is a tradeoff between compression and read/write speed for 100GB file.



With ZLIB compression -> 4 times faster
With LZMA compression conversion -> 15% smaller

# View - Random Access

- We also need to be able to view the information as fast as possible
- ROOT columnar structure allow us to just look at the chromosome and position columns to optimize performance

# Indexing For Fast Access

- Improve read speed by using an index

  - For random access BAM format needs a BAI index file

  - For ROOT implemented a RAMIndex, sparse index in combination with fast columnar ROOT file scanning

  - 16 bytes per entry, compressed

```cpp
class RAMIndex {
private:
   typedef std::pair<int,int> Key_t;          // refid (of rname) and pos
   typedef std::map<Key_t,Long64_t> Index_t;  // map of Key_t and TTree entry number

   Index_t fIndex;
   …
};
```

# View - Performance



Median read speed for ramtools across different parameter settings

# View - Performance

## Viewing 3 different regions from a 9-18 GB file

# Future Work

- Increase file sample - The study was done on only one base SAM file (subsampled to 10% and 1%). Conversion, indexing and views should be run on different SAM files so sample bias is reduced

- Extended format comparison - The current work compared only to the well established BAM format. However, comparisons to f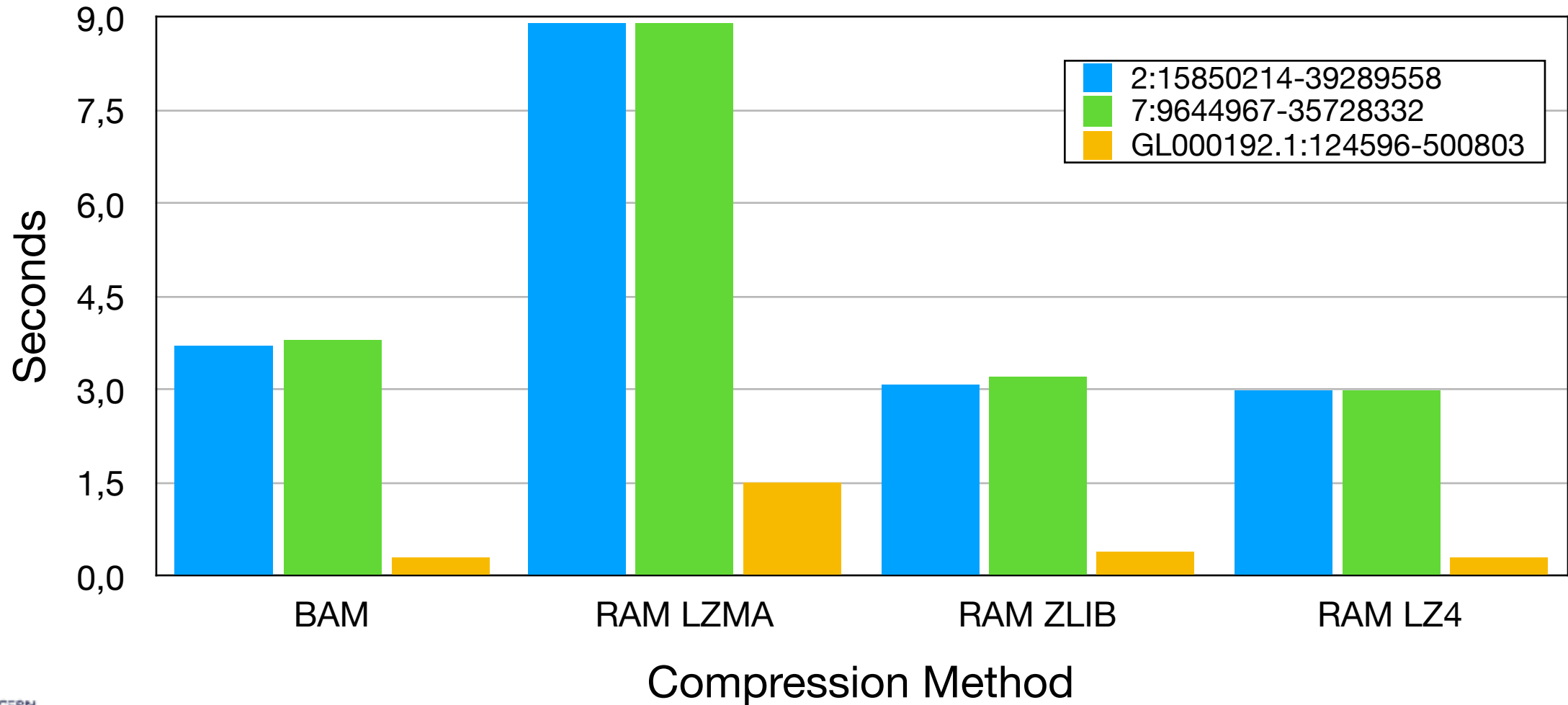ormats such as CRAM should be made to see the relative performances in compression rate and speed and read speed, respectively

- Support for additional operations - Add to ramtools support for sort, merge, split or stats.

- FASTQ/FASTA conversion - The raw sequencer data format. Simpler but closely resembling SAM. Most of the TTree advantages apply also to this type of data. In fact, recent research in formats like LFQC has many analogies to how ROOT branches are used to compress data fields separately, optimizing compression and read speed.

# CERN openlab GeneROOT Technology Transfer Benefits

- Additional use case for CERN's ROOT technology

- Return flow of know-how benefiting the ROOT User community

- CERN openlab provides a doctoral student who gains valuable experience

- Entire Omics community would benefit from improved analysis tools to handle rapidly growing amounts of data

- Project is joint effort between CERN openlab, CERN medical applications group, King's College London

- Big thank you to my 2017 openlab summer student Jose Javier Gonzalez Ortiz