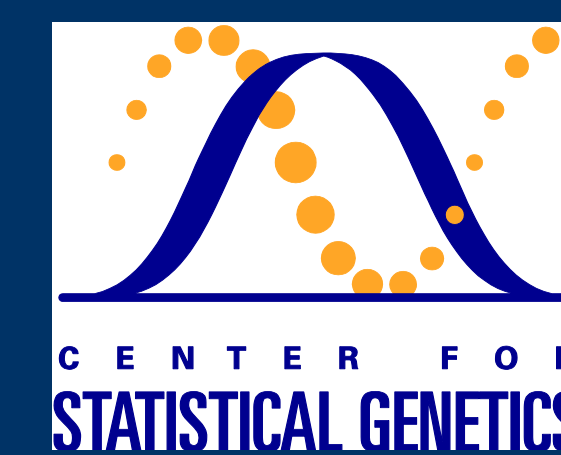


C++ library and tools for next generation sequence data

M. Trost, H.M. Kang, P. Anderson, B. Li, W. Chen, C. Fuchsberger, X. Zhan, A. Tan, G.R. Abecasis
Dept of Biostatistics and Center for Statistical Genetics, Univ Michigan, Ann Arbor, MI, 48109, USA;



INTRODUCTION

In order to handle the increasing volume of next generation sequencing and genotyping data created and developed:

- **C++ Library** – open source, freely available (GPL license), easy to use APIs
 - File/Stream I/O – uncompressed, BGZF, GZIP, stdin, stdout
 - Common file formats – SAM/BAM, FASTQ, GLF
 - Indexed access to BAM files
 - Accessors to get/set values
 - Utility classes, including:
 - Cigar – interpretation and mapping between query and reference
 - Pileup – structured access to data by individual reference position
- **Standalone Tools** – efficiently process and analyze data using this library

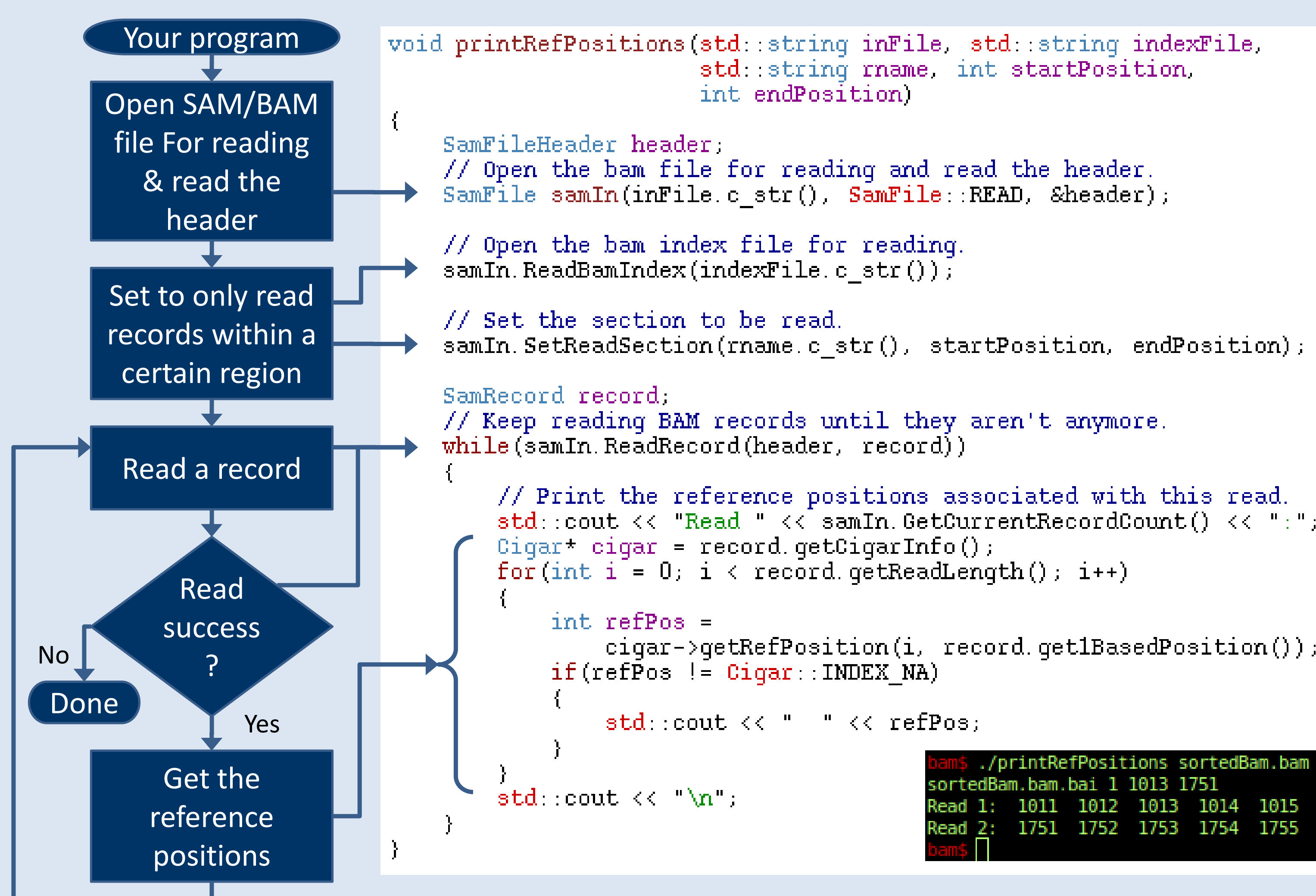
URL FOR LIBRARY & TOOLS

<http://genome.sph.umich.edu/wiki/Software>

- Download of tools and library
- Description and usage information
- Class descriptions, examples, & FAQs

SAMPLE C++ CODE USING LIBRARY

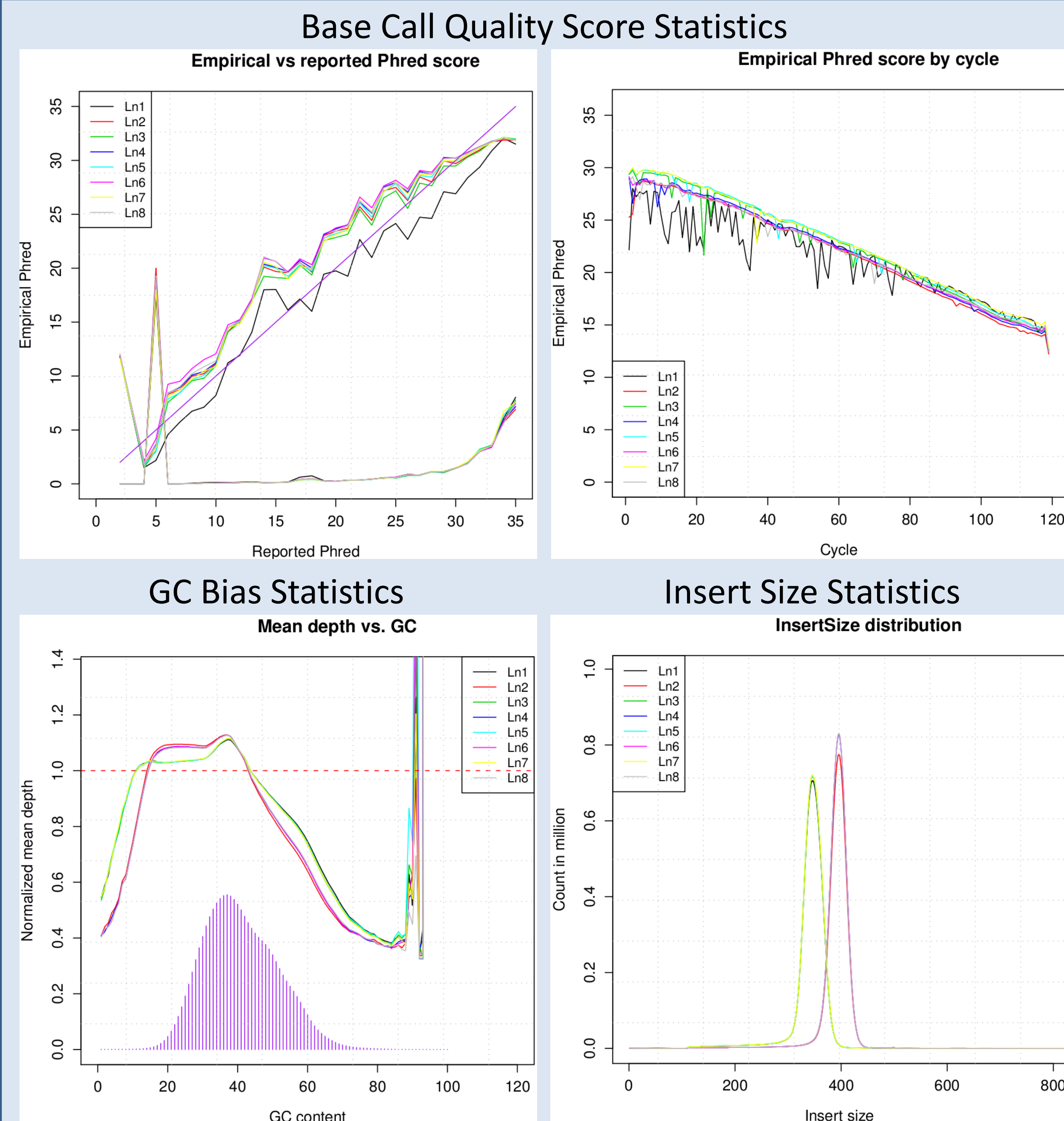
- Print reference positions for BAM file reads that overlap the specified positions.



SAM/BAM TOOLS

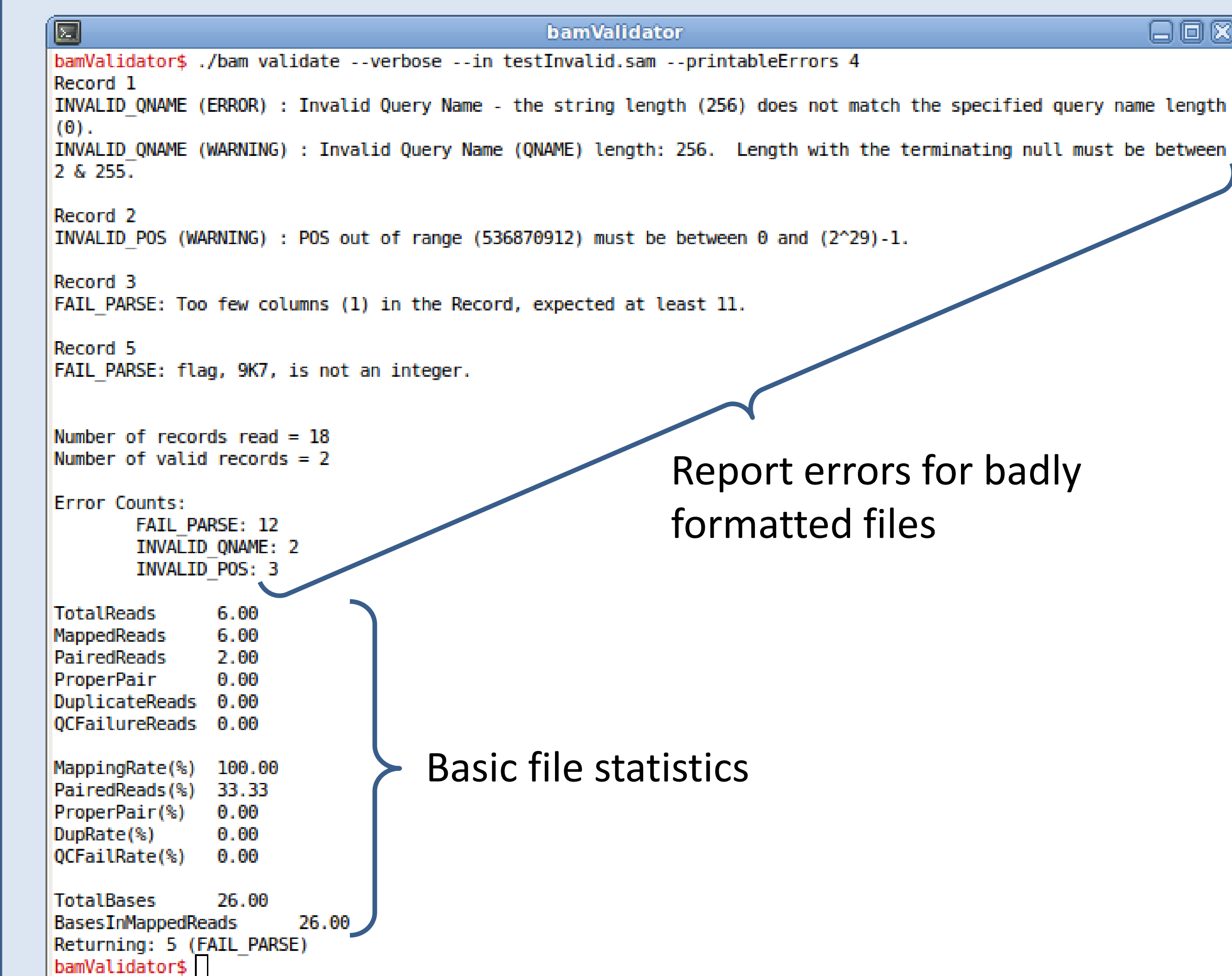
- **VerifyBamID** – Check sample identities for contamination/sample swap
 - Genotype concordance based detection
 - Estimate based on population allele frequencies without genotype data
- **Validate** – Check file format & print statistics
- **Convert** – Convert between SAM & BAM
- **SplitChromosome** – Split into 1 file per Chromosome
- **WriteRegion** – Write only reads in the specified region
- **Filter** – Soft clip ends with too high mismatch % and mark unmapped if quality of mismatches is too high
- **PolishBam** – Add/Update header lines & add RG tag to each record
- **QPLOT** - Calculate & plot summary statistics
- **RGMergeBam** – Merge sorted BAM files adding Read Groups
- **SplitBam** – Split into 1 file per Read Group
- **TrimBam** – Trim end of reads, changing read ends to 'N' & quality to '!'
- **Pileup** – Pileup every base or just bases in specified region and write VCF
- **Deduper** – Mark or remove duplicates
- **Recalibrator** – Resource-efficient tool, which recalibrates base qualities based on an adaptive logistic regression model

QPLOT SAMPLE OUTPUT



- Additional plots include: Depth Distribution, Depth Coverage, and Empirical Q20 Count

BAM VALIDATE SAMPLE OUTPUT



ADDITIONAL TOOLS

- **FastQValidator** – Check format of FASTQ file
 - Reports errors for badly formatted files
 - Reports Base Composition Statistics (%reads at each read index)
- **VcfGenomeStat** – Print flanking sequences and how often they appear for input VCF file

Additional support including VCF and association analysis will be provided in the near future.

Support from National Institutes of Health (U01HG00521401 & HG005552)
References: Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R; 1000 Genome Project Data Processing Subgroup. The Sequence Alignment/Map (SAM) format and SAMtools. *Bioinformatics*, 2009;25, 2078-2079. [PMID: 19505943]