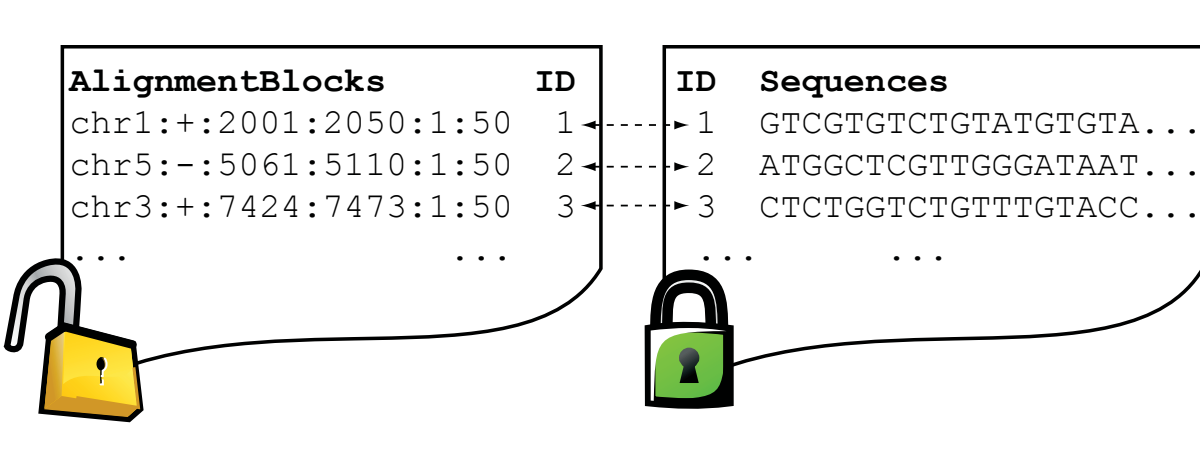


A modular framework for RNA-Seq data analysis using compact, anonymized data summaries

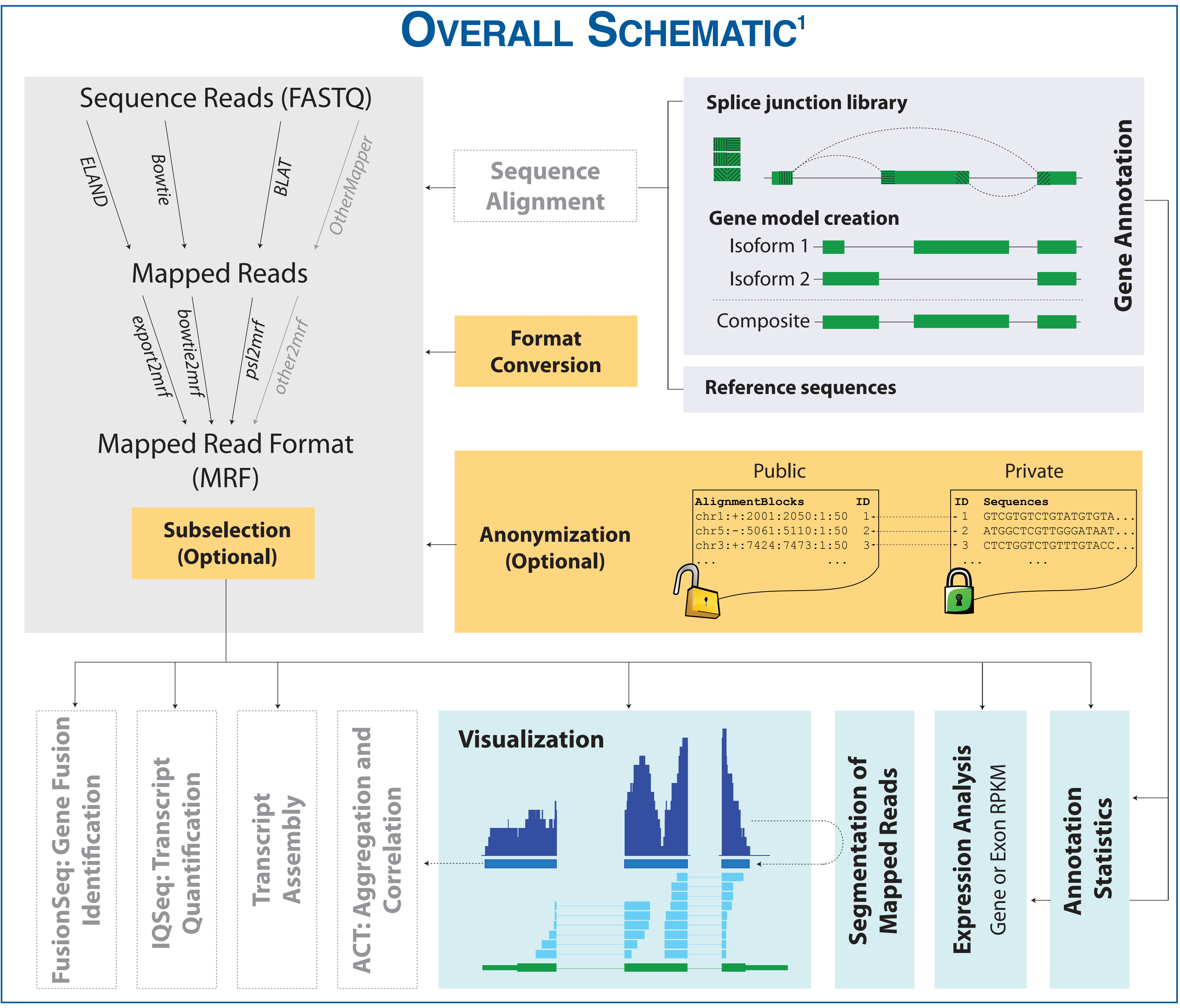
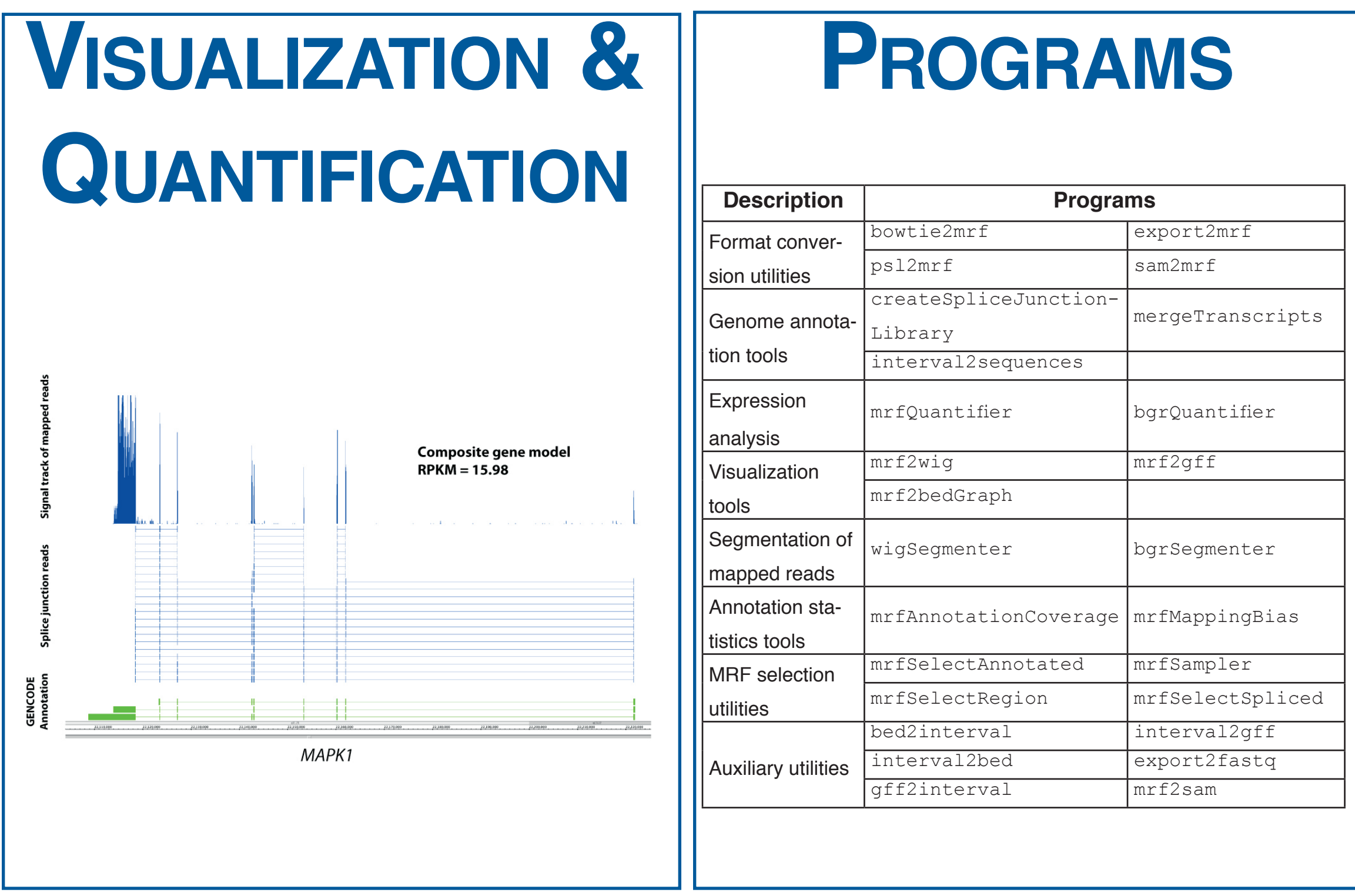
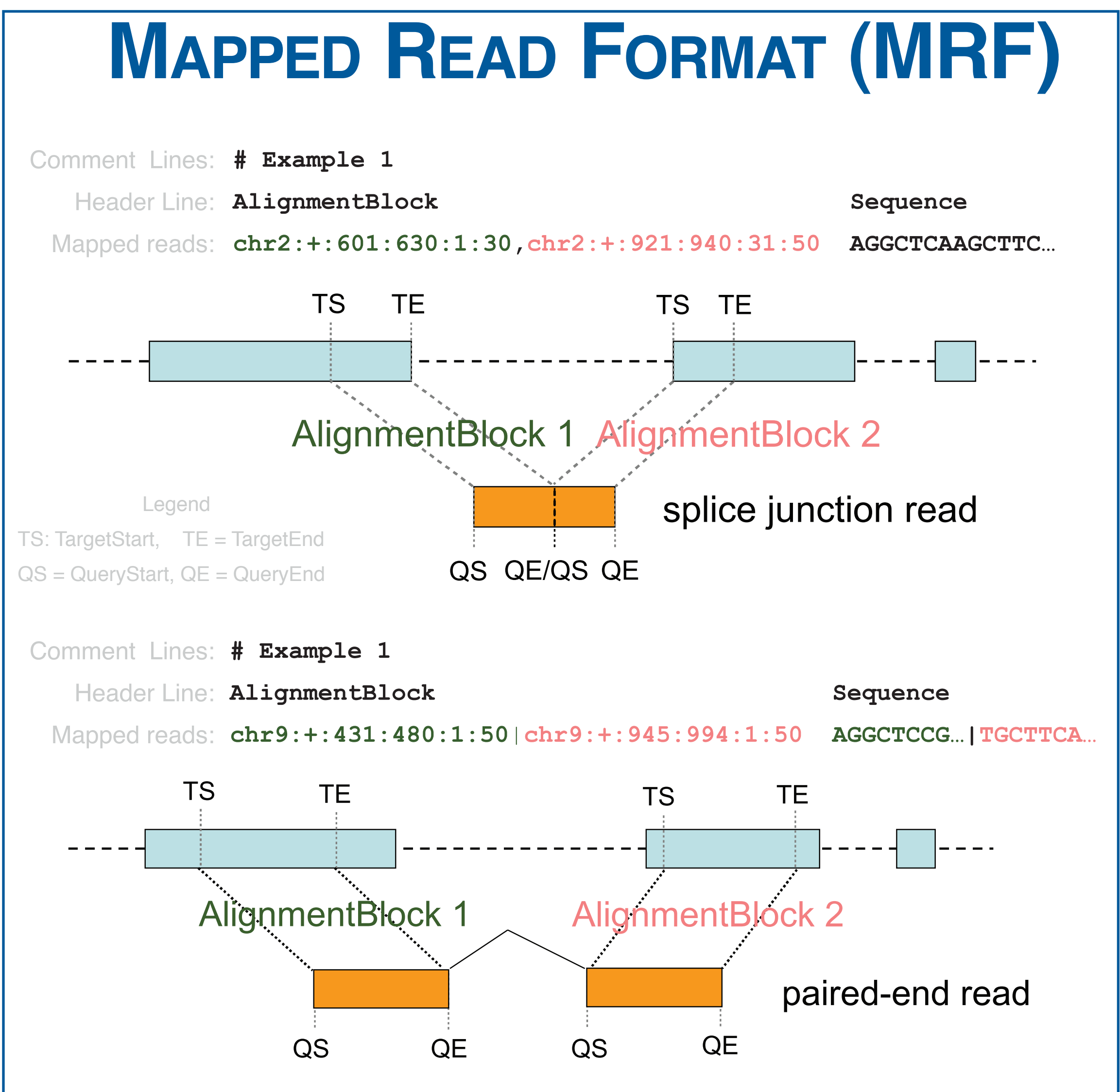


L. HABEGGER^{1,2}, A. SBONER^{1,2}, R. ALEXANDER^{1,2}, T.A. GIANOULIS^{3,4}, J.S. ROZOWSKY², A. AGARWAL^{2,5}, D.Z. CHEN¹, W.M. HUSSAIN^{6,7}, D. PFLUEGER⁶, S. TERRY⁶, F. DEMICHELIS^{6,7}, M.A. RUBIN⁶, M. SNYDER⁸, AND M.B. GERSTEIN^{1,2,5}

1 Program in Computational Biology and Bioinformatics, Yale University, New Haven, CT, USA
 2 Molecular Biophysics and Biochemistry, Yale University, New Haven, CT, USA
 3 Wyss Institute for Biologically-Inspired Engineering, Harvard University, Boston, MA, USA
 4 Department of Genetics, Harvard Medical School, Boston, MA, USA
 5 Department of Computer Science, Yale University, New Haven, CT, USA
 6 Department of Pathology and Laboratory Medicine, Weill Cornell Medical College, New York, NY, USA
 7 Institute for Computational Biomedicine, Weill Cornell Medical College, New York, NY, USA
 8 Department of Genetics, Stanford University School of Medicine, Stanford, CA, USA

INTRODUCTION

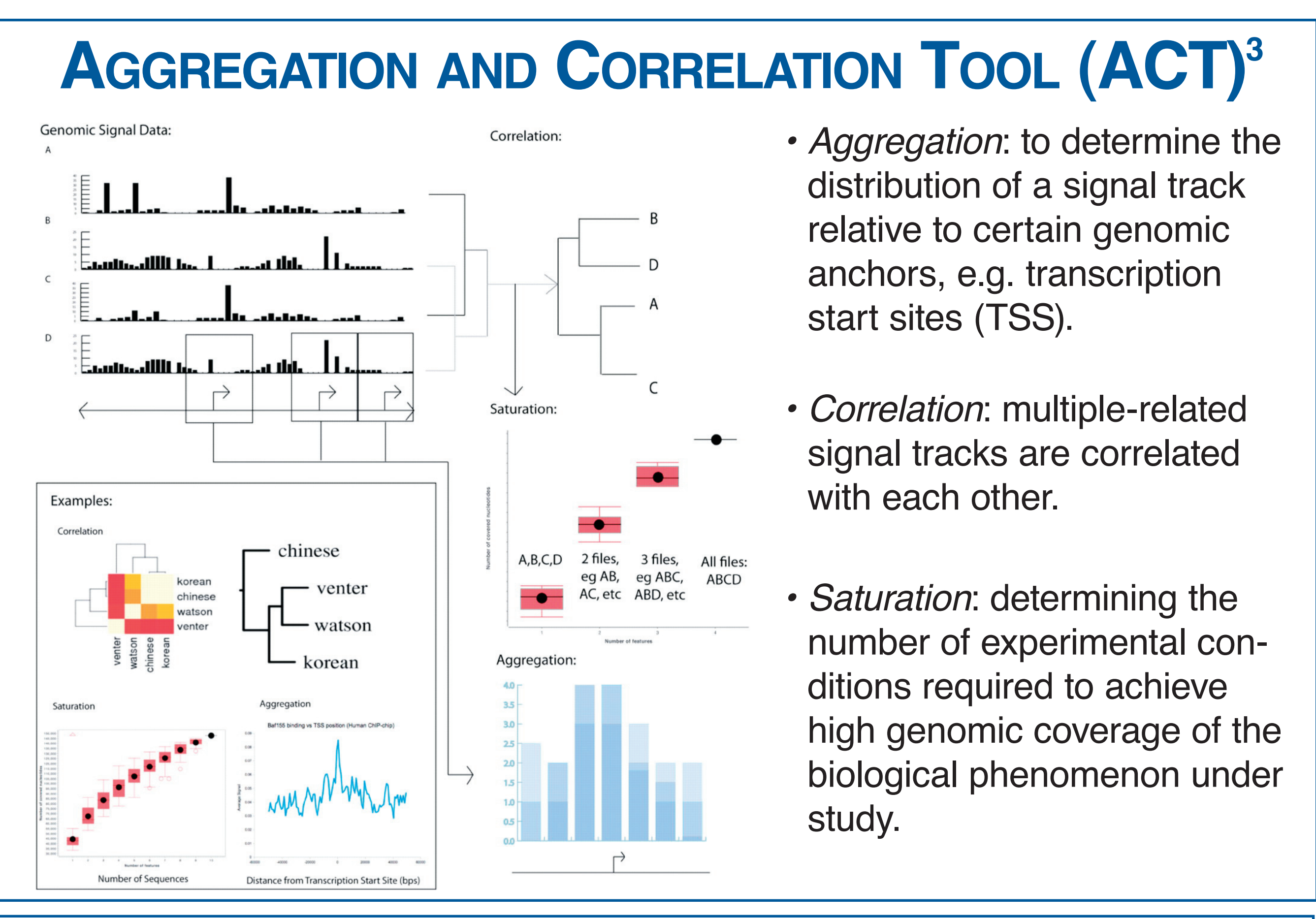
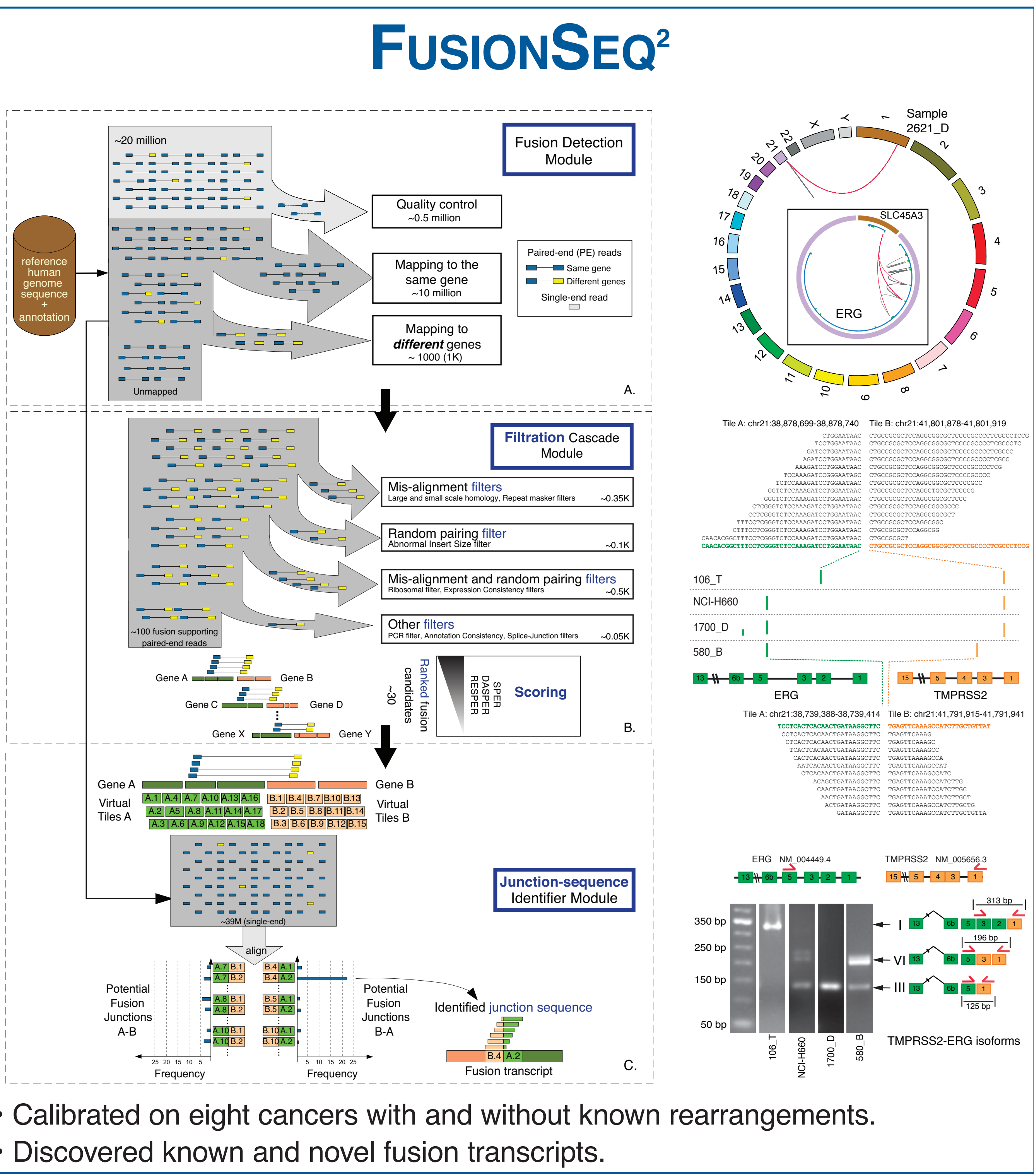
- Next-generation sequencing for functional genomics has given rise to large quantities of sequence information.
- Sequence reads from a specific individual can potentially identify and genetically characterize that person, raising privacy concerns.
- RSEQtools address these issues with *Mapped Read Format (MRF)*, a compact data summary format for short, long and paired-end read alignments.
- RSEQtools also provide a set of programs for RNA-seq data analysis.



RUN STATISTICS (ILLUMINA GAII DATA)

Purpose	Program	Time to process	File sizes (uncompressed)	Notes
Alignment + Conversion	ELAND2	~1 day	Export: 4.2Gb	Processing of one flow cell (8 lanes)
	export2mrf	~ 2 min	MRF: 400Mb	Number of mapped reads: ~12M
Quantification	mrfQuantifier	~ 45 sec	Gene expression values: 3.5Mb	Gencode composite gene models (~22,000)
Visualization	mrf2wig	~ 2 min	One WIG file per chromosome: 1Mb - 150Mb	Signal track of mapped reads normalized per million mapped reads
	mrf2gff	~ 45 sec	One GFF file per chromosome: 100Kb - 16Mb	To visualize splice junction reads

1. L. Habegger et al. "RSEQtools: a modular framework to analyze RNA-Seq data using compact, anonymized data summaries." *Bioinformatics* 27, no. 2 (2011): 281-283.
 2. A. Sboner et al. "FusionSeq: a modular framework for finding gene fusions by analyzing Paired-End RNA-Sequencing data." *Genome Biology* 11, no. 10 (2010): R104.
 3. J. Jee et al. "ACT: aggregation and correlation toolbox for analyses of genome tracks." *Bioinformatics* 27, no. 8 (2011): 1152-1154.



• RSEQtools: <http://rseqtools.gersteinlab.org>
 • FusionSeq: <http://mseq.gersteinlab.org/fusionseq>
 • ACT: <http://act.gersteinlab.org>
 • IQSeq: <http://code.google.com/p/iqseq/>