

Systematic evaluation of spliced alignment programs for RNA-seq data

Pär G Engström^{1,13}, Tamara Steijger¹, Botond Sipos¹, Gregory R Grant^{2,3}, André Kahles^{4,5}, The RGASP Consortium⁶, Gunnar Räscher^{4,5}, Nick Goldman¹, Tim J Hubbard⁷, Jennifer Harrow⁷, Roderic Guigo^{8,9} & Paul Bertone^{1,10-12}

High-throughput RNA sequencing is an increasingly accessible method for studying gene structure and activity on a genome-wide scale. A critical step in RNA-seq data analysis is the alignment of partial transcript reads to a reference genome sequence. To assess the performance of current mapping software, we invited developers of RNA-seq aligners to process four large human and mouse RNA-seq data sets. In total, we compared 26 mapping protocols based on 11 programs and pipelines and found major performance differences between methods on numerous benchmarks, including alignment yield, basewise accuracy, mismatch and gap placement, exon junction discovery and suitability of alignments for transcript reconstruction. We observed concordant results on real and simulated RNA-seq data, confirming the relevance of the metrics employed. Future developments in RNA-seq alignment methods would benefit from improved placement of multimapped reads, balanced utilization of existing gene annotation and a reduced false discovery rate for splice junctions.

examines the density of independent reads at those loci. Many algorithms also consider base-call quality scores and use sophisticated indexing schemes to decrease runtime.

Here we assess the performance of 26 RNA-seq alignment protocols on real and simulated human and mouse transcriptomes. We adopted a competitive evaluation model applied in other areas of bioinformatics¹¹⁻¹⁴. Developers were invited to run their software and submit results for evaluation as part of the RNA-seq Genome Annotation Assessment Project (RGASP). Programs included six spliced aligners GSNAP⁷, MapSplice⁴, PALMapper⁸, ReadsMap, STAR⁹ and TopHat^{5,6}) and four alignment pipelines (GEM³, PASS¹⁵, GSTRUCT and BAGET). GSTRUCT is based on GSNAP, whereas BAGET uses a contiguous DNA aligner to map reads to the genome as well as to exon junction sequences derived from reference gene annotation. For comparison, the contiguous aligner SMALT was also tested. SMALT can map reads in a split manner, but it lacks several features of dedicated spliced aligners, such as precise determination of exon-intron boundaries. We demonstrate that choice of align-