

Fig. S1 Read classification by Bloom-filter vs alignment. 100bp-Illumina reads were simulated by pIRS (v1.1.1) with coverage depths ranging from 10X to 100X from 580 COSMIC (v77) genes and an equal number of non-COSMIC genes randomly selected from RefSeq. Read classification performance (right Y-axis) and run-time (green, left Y-axis) of BBT (v2.1.0) were compared against BWA-MEM (v0.7.12). The classified target of every read was compared against its true gene origin to calculate true (blue) and false (red) positive rates. Benchmarking was done in the same computational environment described in Fig. S5.

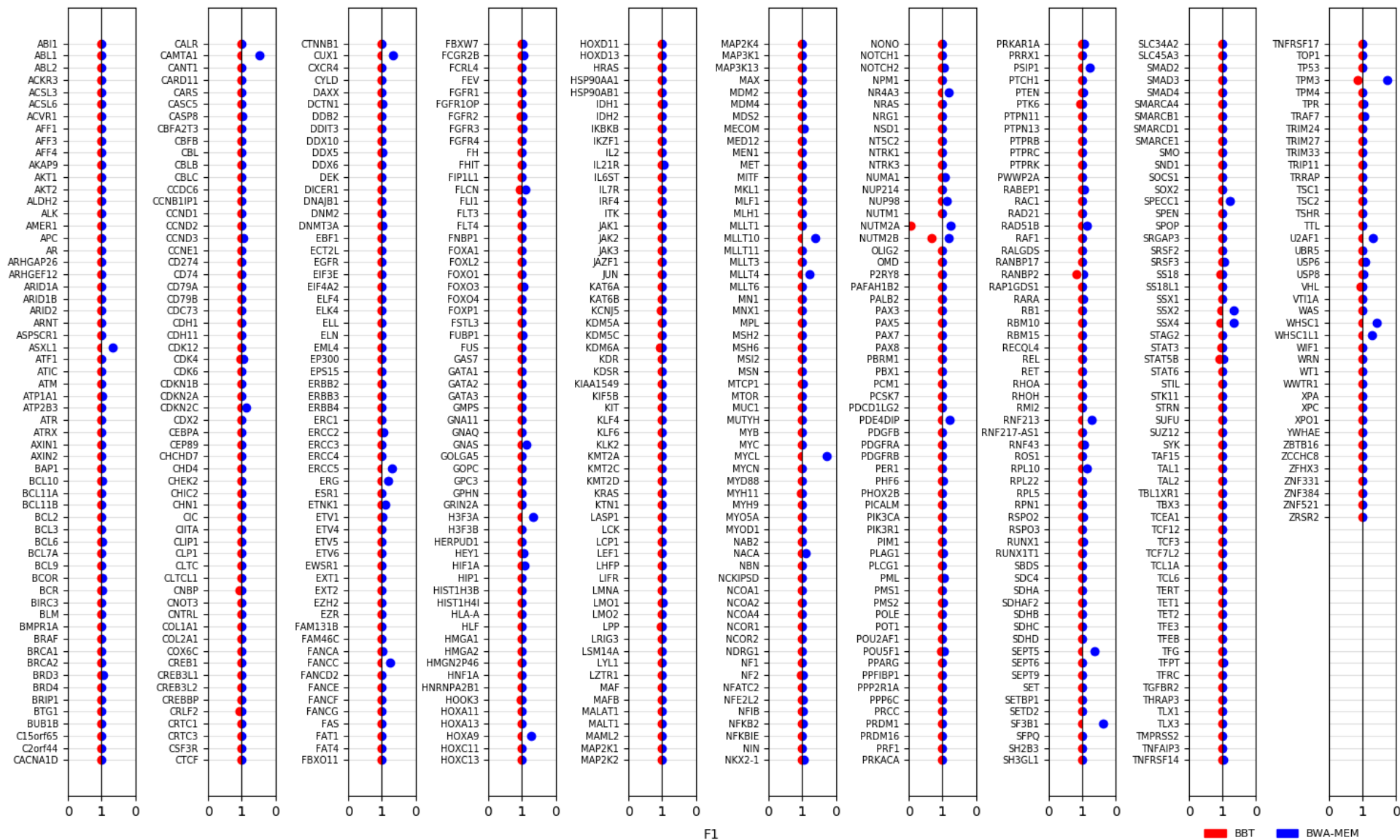


Fig. S2 Per-gene comparison of classification performance by BBT vs BWA-MEM. F1 scores for both methods using the same simulation dataset described in Fig. S1 is calculated for each gene and plotted on the same horizontal line (red=BBT, blue=BWA-mem). The scale on the X-axis for BWA-mem on the right is reversed for easy visual comparison such that higher scores for both methods localize to the middle while lower scores are off-centre. F1 score is calculated as follows: $F1 = \frac{2pr}{p + r}$; where p(precision) = number of correctly assigned reads / total reads assigned; r(recall) = number of correctly assigned reads / total reads simulated (for the gene in question).

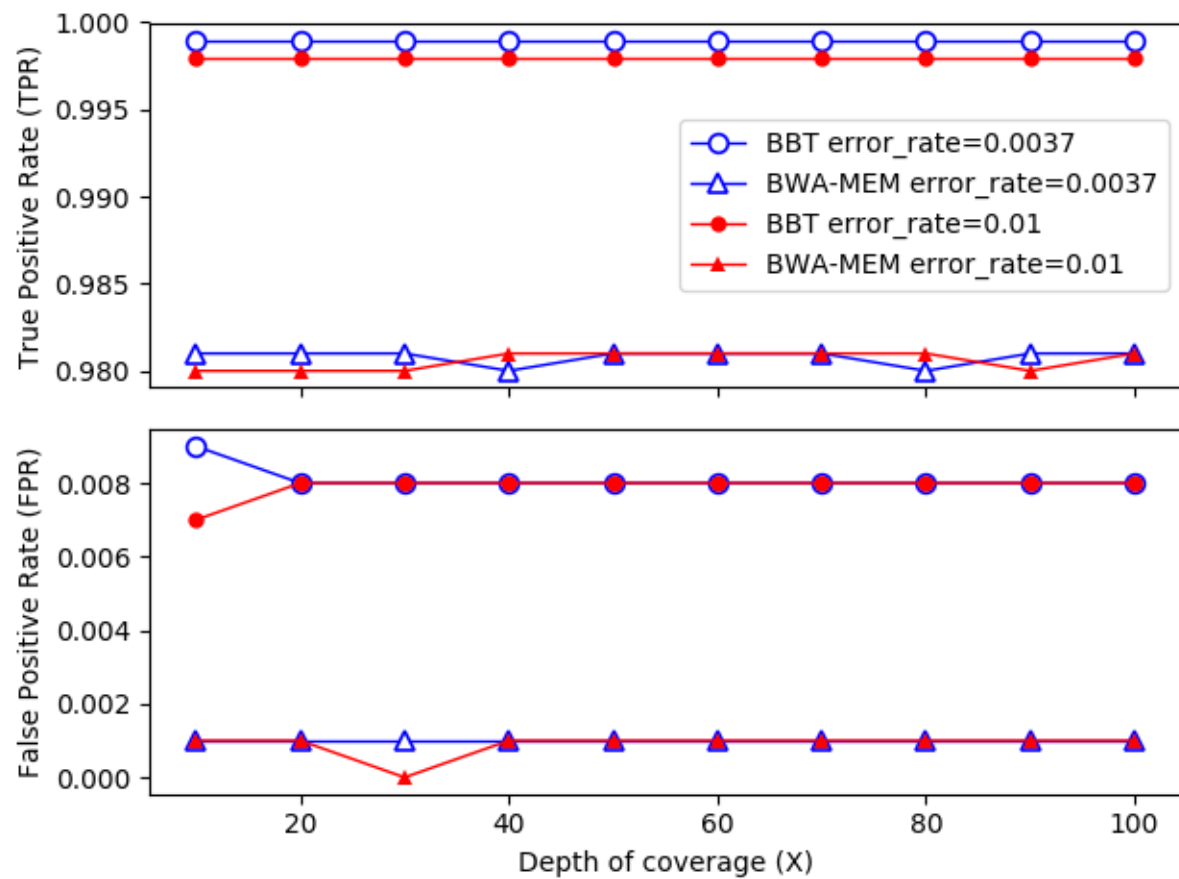


Fig. S3 Effect of sequencing error rate on performance of read classification. Simulation was repeated with increased substitution-error rate of 1% (red and solid data points) and results were compared against the original experiment (Figure S1, 0.37% substitution-error rate, blue and hollow data points). True (top panel) and false positive rates (bottom panel) of both methods, BBT(circle) and BWA-mem(triangle), were plotted against depth of coverage (X-axis).

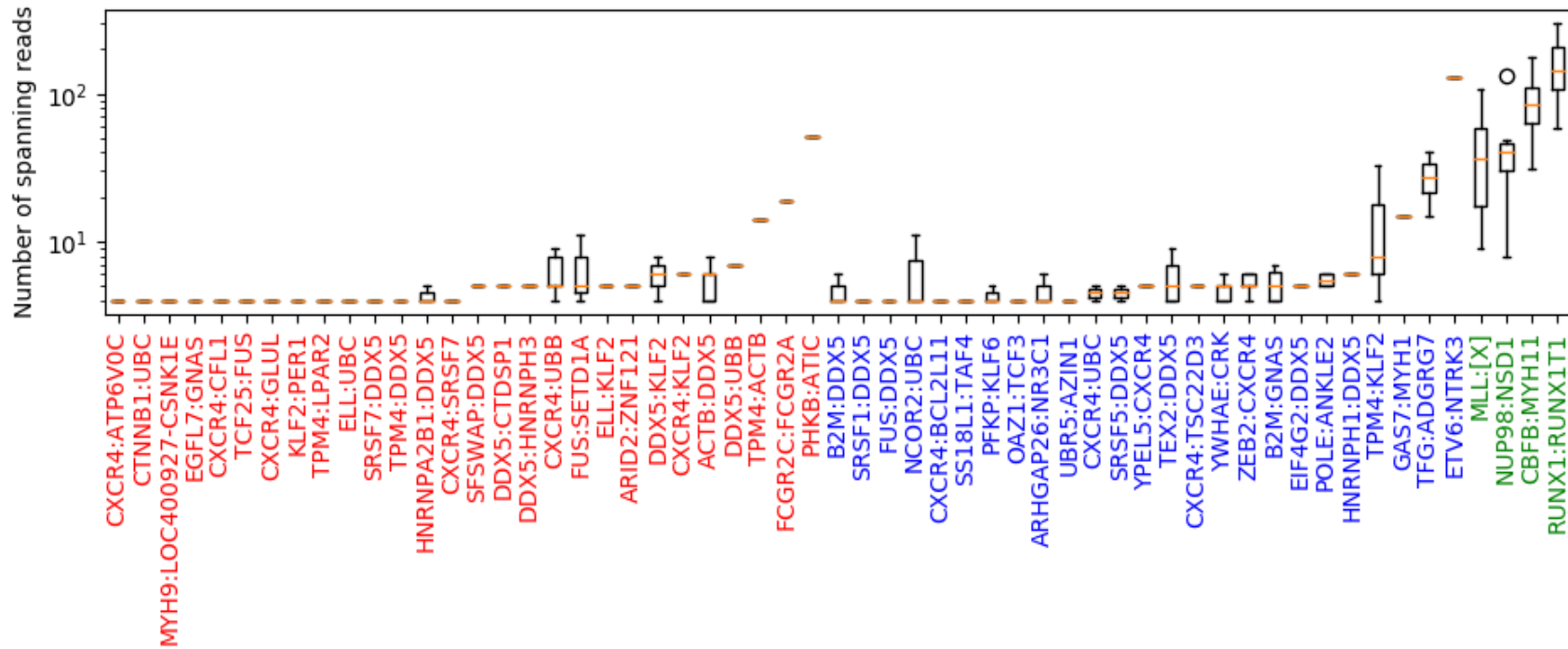


Fig. S4 Support level of gene fusions detected in Leucegene samples. A boxplot is shown to indicate the distribution of the number of junction spanning reads (Y-axis, log-scale) of each fusion event detected in Leucegene samples analyzed in this study (X-axis, grouped by gene names regardless of orientation; *MLL* fusions with different partners in the *MLL*-F cohort were grouped as *MLL*:[X]). Events were colored based on their level of “legitimacy”: green = AML target events (events targeted in Leucegene publications); blue = non-target events with literature support; red = non-target events without literature support. List of non-target events with their associated literature support: B2M:DDX5[1], SRSF5:DDX5[1], FUS:DDX5[1], NCOR2:UBC[2], CXCR4:BCL2L11[1], SS18L1:TAF4[3], PFKP:KLF6[1], OAZ1:TCF3[1], ARHGAP26:NR3C1[1], UBR5:AZIN1[4], CXCR4:UBC[1], SRSF1:DDX5[1], YPEL5:CXCR4[1], TEX2:DDX5[1], CXCR4:TSC22D3[1], YWHAE:CRK[1], ZEB2:CXCR4[1], B2M:GNAS[1], EIF4G2:DDX5[1], POLE:ANKLE2[1], HNRNPH1:DDX5[1], TPM4:KLF2[5], GAS7:MYH1[6], TFG:ADGRG7[7], ETV6:NTRK3[8]

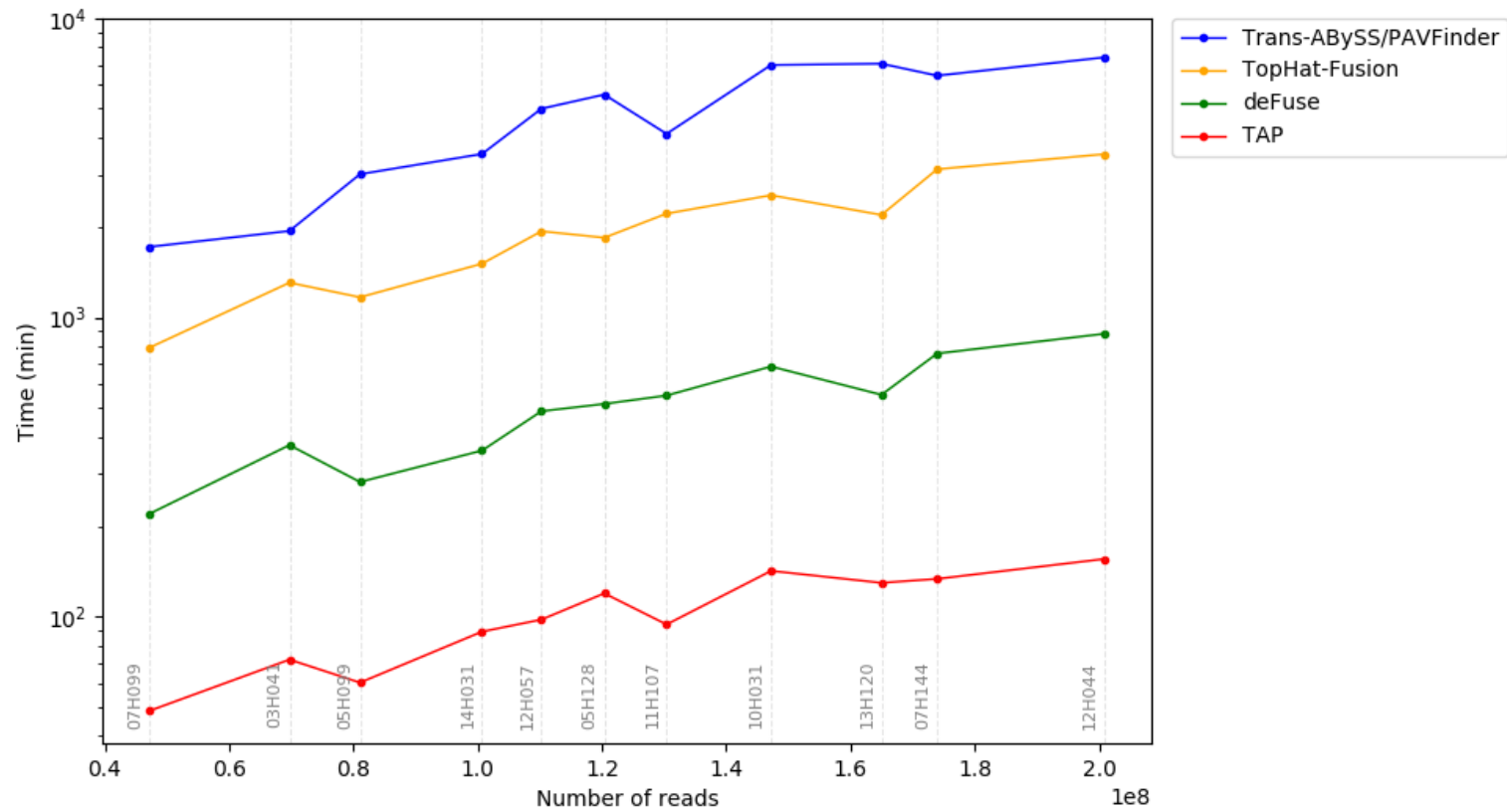


Fig. S5 Benchmarking of TAP and other fusion callers. TAP, Trans-ABYSS/PAVFinder, TopHat-Fusion (v2.1.0), and deFuse (v0.8.0) were run on eleven Leucegene samples with sequencing depths ranging from 50 to 200 million read pairs. Target gene set for TAP comprises 580 COSMIC (v77) genes; the other 3 tools /pipelines were run on entire transcriptome. TopHat-Fusion and deFuse were run with default parameters and references as per their instructions. The second step of TopHat-Fusion (tophat-fusion-post) failed to finish on these samples as it requires in excess of 380GB which is the upper limit of memory available on the testing machines. All benchmarking were performed on Dell C6320 Intel E5 2.2Ghz 48-core 2.2Ghz machines with 384GB RAM running CentOS 6.7.

Table S1 Alignment features used by PAVFinder for classifying various types of transcriptomic structural variants. Types of alignments, alignment strands relative to both the genome and transcript, gene loci involved are used to elucidate the underlying rearrangement each contig with aberrant alignment reconstructs.

Event type	c2g alignment type	Gene loci mapped	c2g strands	Gene strands	Breakpoints flush with exon boundaries
gene fusion	Split	2	NA	same	yes*
read-through	Split/Gapped	2	same	same	yes
internal tandem duplication (ITD)	Split/Gapped/Partial**	1	same	same	NA
partial tandem duplication (PTD)	Split	1	same	same	yes
small insertion	Gapped	1	NA	same	NA
small deletion	Gapped	1	NA	same	NA

c2g strands: reference genome strand contig sequence is aligned to

Gene strands: strand of c2g alignment relative to gene strand

Split alignment: contig is split into 2 segments aligned to non-adjacent locations of the genome (chimeric alignment)

Gapped alignment: contig is aligned to a single genomic but with gap(s) either in the contig or target (genome) sequence not corresponding to introns

NA: not applicable

* Default requirement for fusion calling; parameter can be turned off to allow non-exon-bound fusions to be called

** Contig harboring ITD breakpoints may lead to partial c2g alignments, in which case the unaligned portion will be aligned against the target transcript (determined by c2g) sequence for prediction of an ITD event

Table S2 Block-vs-exon alignment characteristics used by PAVFinder to identify various classes of novel splice variants. Alignment of an assembled transcriptomic contig against the reference genome (c2g) yields alignment “blocks”, with each block corresponding to an exon when a contig reconstructs a reference transcript. When novel splicing not observed in the annotation happens, the colinearity of blocks-vs-exons is disrupted, yielding clues for classifying different types of splice variants.

Event	Block-vs-exon alignment
skipped exon	adjacent contig blocks mapped to non-adjacent exons
novel exon	adjacent exons mapped to non-adjacent contig blocks
novel intron	adjacent contig blocks mapped to single exon with outer boundaries flush
retained intron	adjacent exons mapped to single block with outer boundaries flush
novel splice donor	novel block boundary corresponding to the splice donor of the matching exon
novel splice acceptor	novel block boundary corresponding to the splice acceptor of the matching exon

All the splice variants require the novel block(s) to be flanked by the canonical splicing motif (GT-AG).

Table S3 Software and command lines used in TAP and benchmark experiments

software	version	command
PIRS	1.1.1	pirs simulate -i <reference.fa> -s </pIRS_1111/Profiles/Base-Calling_Profiles/humNew.PE100.matrix.gz> -b </pIRS_1111/Profiles/InDel_Profiles/phixv2.InDel.matrix> -d </pIRS_1111/Profiles/GC-depth_Profiles/humNew.gcdep_100.dat> -x <coverage_depth> -m 150 -Q 33 -c 1 -o <output_prefix> [-e <error_rate>]
biobloommaker	2.1.0	biobloommaker -F -f 0.005 -p cancer_census -a 4 -m -S "00000100110010101001001100000010110111010110001011101001101111101001101111101010000111011000000001010101110010100000010000000000000000000010001100011000110010010010010001100011000110011100100000000001000000101001110101010000000011011100001010111111101100101110100011010111011010000011001001010100110010000" cancer_census/*.fa
TAP	0.4.2	tap.py <output_prefix> <output_directory> --bf </path/biobloommaker/output.bf> --fq </path/xxx_1.fastq.gz> </path/xxx_2.fastq.gz> --k 32 62 --readlen <read_length> --nprocs <number_processes> --params </path/to/tap.cfg>
biobloomcategorizer	2.1.0	biobloomcategorizer -fq -i -p <output_prefix> -a 2 -t <num_threads> -e -f </path/biobloommaker/output.bf> </path/xxx_1.fastq> </path/xxx_2.fastq>
Trans-ABYSS	1.5.4	transabyss --kmer <kmer_size> --pe </path/xxx_1.fastq.gz> </path/xxx_2.fastq.gz> -outdir <output_directory> --name <prefix_name> --cleanup 3 transabyss-merge -mink 32 -maxk 62 --prefixes k32 k62 -length 100 %s --out <output_file> --force
GMAP	2014-12-18	gmap -D </path/gmap_index/> -d hg19 <input.fasta> -t <num_threads> -f samse -n 0 -x 10
PAVFinder	0.4.2	find_sv_transcriptome.py --gbam </path/c2g.bam> --tbam </path/c2t.bam> --transcripts_fasta </path/transcripts.fasta> --genome_index </path/gmap_index/> --r2c </path/r2c.bam> --nproc <number_process> </path/contigs.fasta> </path/refGene.sorted.gtf.gz> </path/hg19.fa> <output_dir> map_splice.py </path/c2g.bam> </path/contigs.fasta> </path/refGene.sorted.gtf.gz> </path/hg19.fa> <output_dir> --r2c </path/r2c.bam> --nproc <number_process> --suppl_annot

		</path/acembly.sorted.gtf.gz>
BWA MEM	0.7.12	bwa mem -t <number_threads> </path/reference/bwa/index> </path/xxx_1.fastq.gz> </path/xxx_2.fastq.gz> samtools view -bhS - -o <output.bam>
TopHat-Fusion	2.1.0	tophat -o <output_directory> -p 12 --fusion-search --keep-fasta-order --bowtie1 --no-coverage-search -r 0 --mate-std-dev 80 --max-intron-length 100000 --fusion-min-dist 100000 --fusion-anchor-length 13 --fusion-ignore-chromosomes chrM </path/bowtie/hg19/index/> </path/xxx_1.fastq> </path/xxx_2.fastq>
deFuse	0.6.1	defuse_run.pl -c </path/config.txt> -d </path/to/defuse_ref/> -o <output_directory> -1 </path/xxx_1.fastq> -2 </path/xxx_2.fastq> -n <name> -p 8

Reference:

1. Marincevic-Zuniga Y, Dahlberg J, Nilsson S, Raine A, Nystedt S, Lindqvist CM, Berglund EC, Abrahamsson J, Cavelier L, Forestier E *et al*: **Transcriptome sequencing in pediatric acute lymphoblastic leukemia identifies fusion genes associated with distinct DNA methylation profiles.** *Journal of hematology & oncology* 2017, **10**(1):148.
2. Li Y, Heavican TB, Vellichirammal NN, Iqbal J, Guda C: **ChimeRScope: a novel alignment-free algorithm for fusion transcript prediction using paired-end RNA-Seq data.** *Nucleic acids research* 2017, **45**(13):e120.
3. **Atlas of Genetics and Cytogenetics in Oncology and Haematology** [<http://AtlasGeneticsOncology.org>]
4. Ruffle F AJ, Boureux A *et al.*: **New chimeric RNAs in acute myeloid leukemia [version 1; referees: 1 approved].** *F1000Research* 2017, **6**(ISCB Comm J):1302.
5. Roberts KG, Morin RD, Zhang J, Hirst M, Zhao Y, Su X, Chen SC, Payne-Turner D, Churchman ML, Harvey RC *et al*: **Genetic alterations activating kinase and cytokine receptor signaling in high-risk acute lymphoblastic leukemia.** *Cancer cell* 2012, **22**(2):153-166.
6. Yoshihara K, Wang Q, Torres-Garcia W, Zheng S, Vegesna R, Kim H, Verhaak RG: **The landscape and therapeutic relevance of cancer-associated transcript fusions.** *Oncogene* 2015, **34**(37):4845-4854.
7. Chase A, Ernst T, Fiebig A, Collins A, Grand F, Erben P, Reiter A, Schreiber S, Cross NC: **TFG, a target of chromosome translocations in lymphoma and soft tissue tumors, fuses to GPR128 in healthy individuals.** *Haematologica* 2010, **95**(1):20-26.
8. Lannon CL, Sorensen PH: **ETV6-NTRK3: a chimeric protein tyrosine kinase with transformation activity in multiple cell lineages.** *Seminars in cancer biology* 2005, **15**(3):215-223.