# Tools and best practices for retrotransposon analysis using high-throughput sequencing data

_____

Aurélie Teissandier[1,2,3,4], Nicolas Servant[1,23,*], Emmanuel Barillot[1,2,3] and Deborah Bourc'his[1,4] *

## Supplementary Methods

**Mapping parameters**

The following parameters were used for each mapper and mode.

**Unique mode :**

- Bowtie v1.0.0 : -m 1 -e 400 -v 3 --chunkmbs 100 -p 4 -I 0 -X 1000 --nomaqround -y --best –strata
- Novoalign v3.2.11 : -o SAM -r None -i 1000,200 -F STDFQ
- STAR v2.5.2b : --runThreadN 4 --outSAMtype BAM SortedByCoordinate --runMode alignReads --outFilterMultimapNmax 1 --outFilterMismatchNmax 3 --alignEndsType EndToEnd --alignIntronMax 1 --alignMatesGapMax 350
- Bowtie2 v2.1.0 : -N 1  -X 1000 -p 4 + post filtering keeping alignments when alignment score (AS tag) is higher than second valid alignment (XS tag).
- BWA aln v0.7.15 : -t 4 -n 3 + post filtering keeping alignments when XT tag is equal to U (meaning Unique alignment)
- BWA mem v0.7.15 : -t 4 -c 50000 -T 20 + post filtering keeping alignments when XT tag is equal to U (meaning Unique alignment)


**Random mode :**

- Bowtie v1.0.0 : -M 1 -e 400 -v 3 --chunkmbs 100 -p 4 -I 0 -X 1000 --nomaqround -y --best --strata
- Novoalign v3.2.11 : -o SAM -r Random -i 1000,200 -F STDFQ
- STAR v2.5.2b : --runThreadN 4 --outSAMtype BAM SortedByCoordinate --runMode alignReads --outFilterMultimapNmax 1000 --outSAMmultNmax 1 --outFilterMismatchNmax 3 --outMultimapperOrder Random --winAnchorMultimapNmax 1000 --alignEndsType EndToEnd --alignIntronMax 1 --alignMatesGapMax 350
- Bowtie2 v2.1.0 : -N 1  -X 1000 -p 4
- BWA aln v0.7.15 : -t 4 -n 3

36 &bull; BWA mem v0.7.15 : -t 4 -c 50000 -T 20

37

38 **<u>Multi-hit mode :</u>**

39 &bull; Bowtie v1.0.0 : -a -m 5000 -e 400 -v 3 --chunkmbs 100 -p 4 -I 0 -X 1000 --
40  nomaqround -y --best --strata
41 &bull; Novoalign v3.2.11 : -o SAM -r All 5000 -i 1000,200 -F STDFQ
42 &bull; STAR v2.5.2b : --runThreadN 4 --outSAMtype BAM SortedByCoordinate --runMode
43 alignReads --outFilterMultimapNmax 1000 --outFilterMismatchNmax 3 --
44 outMultimapperOrder Random --winAnchorMultimapNmax 1000 --alignEndsType EndToEnd
45 --alignIntronMax 1 --alignMatesGapMax 350
46 &bull; Bowtie2 v2.1.0 : -N 1 -k 5000 -X 1000 -p 4

47

48 Due to complex TE content in the mouse genome, parameters were adapted using STAR
49 when genome-wide libraries were mapped.

50 **<u>Mouse genome-wide mapping :</u>**

51 &bull; STAR v2.5.2b unique mode : --runThreadN 4 --outSAMtype BAM Unsorted --
52 runMode alignReads --outFilterMultimapNmax 5000 --outFilterMismatchNmax 3 --
53 outMultimapperOrder Random --winAnchorMultimapNmax 5000 --alignEndsType EndToEnd
54 --alignIntronMax 1 --alignMatesGapMax 350 --seedSearchStartLmax 30 --
55 alignTranscriptsPerReadNmax 30000 --alignWindowsPerReadNmax 30000  --
56 alignTranscriptsPerWindowNmax 300 --seedPerReadNmax 3000 --seedPerWindowNmax
57 300 --seedNoneLociPerWindow 1000

58 &bull; STAR v2.5.2b random mode : --runThreadN 4 --outSAMtype BAM Unsorted --
59 runMode alignReads --outFilterMultimapNmax 5000 --outSAMmultNmax 1 --
60 outFilterMismatchNmax 3 --outMultimapperOrder Random --winAnchorMultimapNmax 5000 -
61 -alignEndsType EndToEnd --alignIntronMax 1 --alignMatesGapMax 350 --
62 seedSearchStartLmax 30 --alignTranscriptsPerReadNmax 30000 --
63 alignWindowsPerReadNmax 30000  --alignTranscriptsPerWindowNmax 300 --
64 seedPerReadNmax 3000 --seedPerWindowNmax 300 --seedNoneLociPerWindow 1000

65 &bull; STAR v2.5.2b multi-hit mode : --runThreadN 4 --outSAMtype BAM Unsorted --
66 runMode alignReads --outFilterMultimapNmax 1 --outFilterMismatchNmax 3 --
67 outMultimapperOrder Random --winAnchorMultimapNmax 5000 --alignEndsType EndToEnd
68 --alignIntronMax 1 --alignMatesGapMax 350 --seedSearchStartLmax 30 --
69 alignTranscriptsPerReadNmax 30000 --alignWindowsPerReadNmax 30000  --
70 alignTranscriptsPerWindowNmax 300 --seedPerReadNmax 3000 --seedPerWindowNmax
71 300 --seedNoneLociPerWindow 1000

72

73 **Quantification comparison**

74 The following tools and parameters were used :

75 *repEnrich:*

76 bowtie --chunkmbs 10000 -p 4 -t -m 1 -S --max multimap.fastq -1 R1.fastq -2 R2.fastq >

77 unique.sam ; samtools view –bS unique.sam > unique.bam ; samtools sort -o

78 unique_sort.bam unique.bam ; samtools index unique_sort.bam ; RepEnrich.py rmskFile

79 outDir nameSample genomeFolder multimap_1.fastq --fastqfile2 multimap_2.fastq

80 unique_sort.bam --cpus 4 --pairedend TRUE --is_bed TRUE --allcountmethod TRUE

81 *TEtools:*

82 TEcount.py -rosette rmskFile -column 2 -TE_fasta rmskFile.fa -count TE.count -RNA

83 R1.fastq -RNApair R2.fastq -bowtie2

84 *TEtranscripts:*

85 STAR --readFilesIn R1.fastq R2.fastq --runThreadN 4 --outSAMtype BAM Unsorted --

86 runMode alignReads --outFilterMultimapNmax 100 --winAnchorMultimapNmax 100 --

87 outFilterMismatchNmax 3 --alignEndsType EndToEnd --alignIntronMax 1 --

88 alignMatesGapMax 350 ; samtools  sort -n -o multiple_sort.bam multiple.bam ; TEtranscripts

89 -t multiple_sort.bam  -c multiple_sort.bam –TE rmsk.gtf --stranded no --format BAM –mode

90 uniq|multi --GTF refGene.gtf --project myProject ;

91 *SQuIRE:*

92 squire Map -1 R1.fastq -2 R2.fastq -o outSquire -f genomeSquire -r 100 -p 4 ; squire Count -

93 m outSquire -f genomeSquire -r 100 -p 4 -o outSquire -c genomeCleanSquire

94 *FeatureCounts Unique Alignments*:

95 featureCounts -F SAF -T 1 -s 0 -p -a rmsk.SAF -o outfeatureCounts.txt Input.bam

96 *FeatureCounts Random Alignments*:

97 featureCounts -M -F SAF -T 1 -s 0 -p -a rmsk.SAF -o outfeatureCounts.txt Input.bam

98 *FeatureCounts Multiple Alignments*:

99 featureCounts -M --fraction -F SAF -T 1 -s 0 -p -a rmsk.SAF -o outfeatureCounts.txt

100 Input.bam

101

# Supplementary Figure Legends

**Supplementary Figure 1: Comparison of mapper efficiency with human simulated data.** **(A)** True Positive (TP) rate versus mapping percentage with chromosome 1 of the human genome. The dots are the average values of three independent simulated libraries. SE and PE refer to single end and paired end, respectively. **(B)** Use memory, run time and size of the BAM file with chromosome 1 of the human genome. The error bars correspond to standard deviation from three independent simulated libraries.


**Supplementary Figure 2: Comparison of the methods for the quantification of human retrotransposon families**. **(A)** Comparison of the estimated abundance versus the true abundance for different quantification methods using human simulated TE-derived library. An R-squared value ($R^2$) was calculated to evaluate the correlation of estimated values between simulated values **(B)** Comparison of the estimated abundance versus the true abundance for TEtools and when randomly reported reads are used for the TE quantification with FeatureCounts (*FeatureCounts Random alignments*). A PE genome-wide library (10X coverage) was simulated using the human genome with STAR for the mapping.


**Supplementary Figure 3: Impact of read depth in TE families quantification**. **(A)** Estimated abundance for different quantification methods and true abundance (Simulated counts) using 5X, 10X, 25X, 50X and 100X coverage on specific mouse TE families. Only these TE families were used for the quantification. **(B)** Same as in A), with specific human TE families.


**Supplementary Figure 4: Mappability of the different human retrotransposon families**. **(A)** True Positive (TP) rate versus mapping percentage per TE family using STAR and paired-end library and human simulated TE-derived reads. Black triangle represents the True Positive rate and percentage of mapping for the entire simulated library **(B)** Mapping percentage versus age of L1Md families. Dot colors represent the True Positive (TP) rate. Ages are obtained from previously published divergence analysis study (25) **(C)** Gain of True Positive in percentage versus gain of mapping in percentage when PE library are used in comparison to SE library.