

Mobile Genomics: Tools and Techniques for Tackling Transposons

Kathryn O'Neill¹, David Brocks², and Molly Gale Hammell¹

¹ Simons Center for Quantitative Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, 11724, USA.

² The Weizmann Institute of Science, Department of Computer Science and Applied Mathematics, Rehovot, IL

Abstract

Next generation sequencing (NGS) approaches have fundamentally changed the types of questions that can be asked about gene function and regulation. With the goal of approaching truly genome-wide quantifications of all the interaction partners and downstream effects of particular genes, these quantitative assays have allowed for an unprecedented level of detail in exploring biological interactions. However, many challenges remain in our ability to accurately describe and quantify the interactions that take place in those hard to reach and extremely repetitive regions of our genome comprised mostly of transposable elements (TEs). Tools dedicated to TE derived sequences have lagged behind, making the inclusion of these sequences in genome-wide analyses difficult. Recent improvements, both computational and experimental, allow for the better inclusion of TE sequences in genomic assays and a renewed appreciation for the importance of TE biology. This review will discuss the recent improvements that have been made in the computational analysis of TE derived sequences as well as the areas where such analysis still proves difficult.

Keywords

Transposable elements; Computational Genomics; Retrotransposons; Single cell Analysis

Introduction

While several types of genomic repeated sequences exist, the largest fraction of the human genome, approximately half, is comprised of transposable elements (TEs)¹, though some groups estimate much larger TE fractions². These TEs, often called transposons or jumping genes, are DNA sequences that have, or once had, the ability to mobilize within the genome, either directly or through an RNA intermediate. TEs are present, to varying degrees, in the genomes of all known types of organisms, both prokaryotic and eukaryotic, with some species showing more genomic transposons than host sequences³. Several excellent reviews have discussed the many and varied types of TEs⁴⁻⁶. Briefly, TEs come in two major types. Class I TEs, also called retrotransposons, first transcribe an RNA copy that is then reverse transcribed to cDNA before inserting elsewhere in the genome. Class II TEs, also called DNA transposons, directly excise themselves from one location before reinsertion. In the human genome, the vast majority of TEs are of the Class I, retrotransposon type. Nearly all human TEs have lost the ability to fully mobilize⁷⁻⁹, with the human specific LINE-1 element (L1HS) being the only fully autonomous TE with the ability to generate new transposition events

to date. However, most TEs have retained some level of functionality, including the ability to direct their own transcription. Thus, transcriptome-wide sequencing assays, like RNA-seq, frequently include transposon-derived transcripts among the set of expressed sequences. Moreover, some transposon transcripts have been coopted to play a role in host function, particularly during early development, such that some expressed transposon transcripts have been shown to be necessary for proper cell differentiation and maintenance of identity¹⁰⁻¹⁴. In addition to their roles in general cellular function, several types of transposons have become intricately entangled within gene regulatory networks¹⁵, contributing both to cis regulatory sequences¹⁶⁻¹⁸ as well as general chromatin environments¹⁹⁻²¹. For this reason, it is paramount we consider the contribution of repetitive elements as we unravel the genomic and epigenomic landscapes that control gene expression.

Properly accounting for repetitive regions in most genomics analysis settings requires special considerations for the challenges presented by the number of nearly-identical transposon sequences dispersed throughout our genomes. Thus, reads derived from these regions are frequently discarded in most sequencing data analysis protocols due to the difficulty in properly assigning TE-derived reads to the correct locus of origin. Few packages explicitly support inclusion of repeats and some intentionally discard reads from these regions, as discussed in a recent review²². Of the packages designed to address TEs, many tools focus on the detection of novel TE insertions or TE-associated genomic rearrangements. Few tools are developed specifically to address regulatory and transcriptional activity of TEs in common assays, such as RNA-seq, chromatin immunoprecipitation sequencing (ChIP-seq), cross-linking immunoprecipitation sequencing (CLIP-seq), and small RNA-seq (sRNA-seq). In this review, we seek to provide an overview of the packages that explicitly support the inclusion of TE sequences in differential expression and binding analyses, and the strides which have been made to improve our ability to resolve ambiguously mapped reads in genomics analysis.

Annotation and de novo detection

A well assembled and annotated genome is the foundation for effective analysis, as all subsequent analyses discussed below require a reference genome as well as a map of gene and TE positions. While many genomes have near complete assemblies, and extensive annotation, the quality of both tends to drop over repeat rich regions for the same reasons discussed above: ambiguity in placing near-identical sequence reads from highly similar copies of related transposons. This ambiguity leads to non-contiguous and erroneous chromosomal assembly, which will feed forward into any genomics analyses using these assemblies²³. Genome assembly has benefitted immensely from long read sequencing technologies, particularly in the context of highly repetitive centromeric regions and in nested repeating elements^{24,25}. While these long read technologies are improving the reference genomes used to map new datasets, one caveat is that transposons are often polymorphic within populations, such that each new sample sequenced is expected to have many non-reference transposon-associated insertions, deletions, and other structural variants that may be rare or private^{26,27}.

Once a high-quality assembly is constructed, the process of annotation may begin. Many curated annotation databases have been developed for identifying repeat elements. For an in depth review of annotation practices and existing repositories please refer to the review by Goerner-Potvin et al.²² Here the distinction between TE-, genome-, and polymorphism-focused annotation repositories is emphasized in addition to a list of software for de novo insertion detection. The most widely used database of TE consensus sequences is RepBase²⁸, which provides the sequences with which genome-specific annotation files are constructed. These annotation files are available through the University of California Santa Cruz Genome Browser(UCSC) and RepeatMasker²⁹. While new RepBase consensus sequences require a subscription, several open databases for repeat annotation are available in addition to UCSC including: RepetDB³⁰, ERVdb³¹, Dfam³², TREP³³, SPTEdb³⁴, ConTEdb³⁵, and mips-REdat³⁶. The ideal database for analysis will vary depending on model organism and transposable elements of interest, as some databases are species and TE type specific.

Mapping

After the construction of a well annotated reference genome, one is faced with the task of mapping experimental data to the appropriate reference. Even with a perfectly annotated and constructed genome, ambiguously mapped sequencing reads still present a challenging problem. One of the first approaches to address this problem, designed for RNA-seq analysis, was to probabilistically assign multi-mapped reads to regions that also show a higher density of uniquely mapped reads, i.e. reads with a single best genomic alignment under the mapping software's heuristics³⁷. However, this was a highly gene centric model that was primarily focused on host gene expression, and was not explicitly intended for estimating expression from TE loci. Moreover, this approach is biased toward regions that have some uniquely mappable content. Unfortunately, the most recently integrated TE insertions are also the least likely to be uniquely mappable, and are thus the most likely to be lost or underestimated by these methods. To highlight this, Figure 1 displays the estimated mappability of several different types of TEs in the human genome, with a specific emphasis on younger types of TEs shown to be active in the human genome³⁸. Mappability in this plot was defined as the inverse of the number of times a simulated 76bp paired end read mapped to the genome, allowing three mismatches. Mappability was scored per nucleotide with the score assigned to the first nucleotide of the read. This track was procured from an in-depth analysis performed by Sexton and Han which considers the many parameters that contribute to the mappability of a particular sequence including the mapping software chosen and the length of the sequenced read³⁹. These analyses still return to the same basic theme displayed in Figure 1: mappability rates vary for different types of transposons, and the most recently inserted transposons are the most likely to be discarded by standard analyses that rely on uniquely mapped reads. In other words, the transposons that present the most problems in genomics analyses are precisely those that are more likely to be functional in terms of: carrying fully functional promoters, encoding for functional proteins, and, rarely, mobilizing within the genome. In addition, many older elements with degraded versions of these components have been recycled to play roles in cis regulatory architecture⁴⁰.

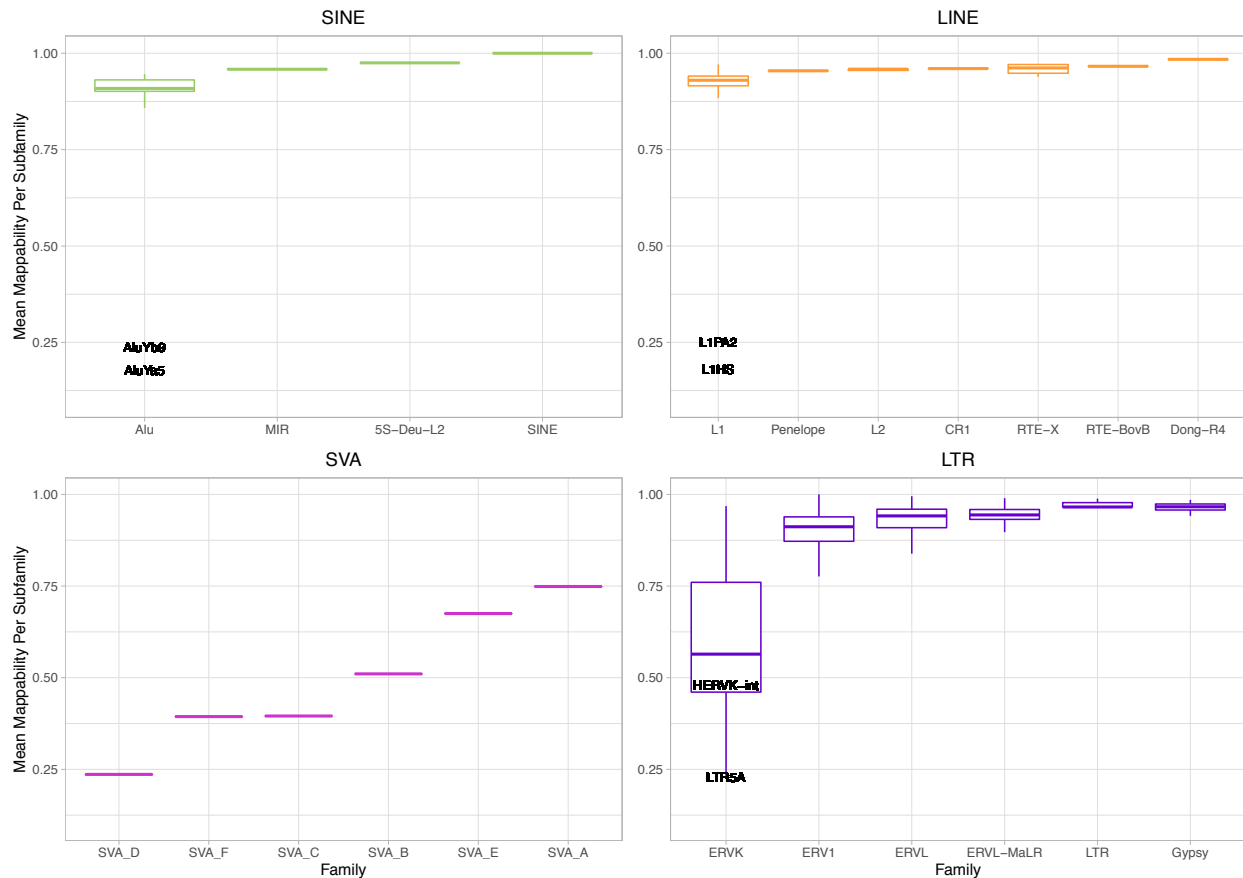


Figure 1 | Estimated mean mappability for different types of TEs in the human genome. Mappability tracks from the analysis by Sexton and Han for hg38 were used to construct mean mappability estimates (average probability that a pair of 76bp reads would map uniquely to a genomic instance of that TE). These were then aggregated by subfamily (L1HS is a human specific subfamily of the LINE class). Some TEs have accumulated enough mutations across each locus that nearly all copies are uniquely mappable. Very recently inserted, and/or still active TEs, show the lowest mappability rates with many copies still very close to the consensus sequence (e.g. Alu and SVA types). In contrast, many older SINE and LINE TEs have high mappability rates and can easily be assessed using only uniquely aligning reads with standard analysis procedures. Mappability was calculated by counting number of times a 76bp paired end read (242-mer with an internal gap of 100 nt) would map within the genome at a particular nucleotide where that nucleotide was the beginning of a 242-mer.

Most genome alignment software is aware of the difficulties posed by ambiguously mapped reads, and thus provide extensive parameter sets designed to allow the user to choose the number of alignments considered for each sequenced read. This includes standard genome mapping software applicable to genome resequencing studies as well as ChIP-seq based studies of protein-DNA binding, such as BWA⁴¹, bowtie⁴², and Novoalign (<http://novocraft.com/>). For RNA-seq aligners there are two approaches, those

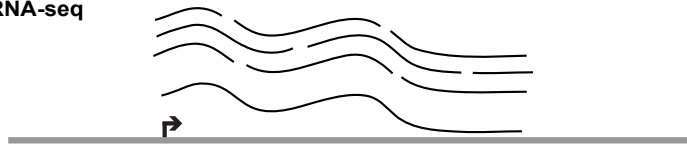
that align to reference transcriptomes and those that align to genomes. Transcriptome methods like kallisto⁴³ and salmon⁴⁴ perform pseudoalignments with transcript derived k-mers and can attempt to build the reference transcriptome from the RNA-seq data itself. Salmon can be specified to report unmapped reads, kallisto does not include this option. While pseudoalignment is very fast, computationally less intensive, and helpful in organisms without a reference genome, it can be complicated in the context of repetitive elements, where all of the caveats that make genome assembly difficult (discussed above) would also apply to de novo transcriptome assembly. With regard to genome based RNA-seq aligners, there are a number of packages available including: STAR⁴⁵, HISAT2⁴⁶, GSNAP⁴⁷, Novoalign, RUM⁴⁸, Minimap2⁴⁹ and others⁵⁰. In the context of sRNA-seq data, short read genome-based aligners (BWA⁴¹, bowtie⁴² and SCRAM⁵¹) that do not consider splice junctions tend to work as well or better than RNA-seq tailored algorithms, with SCRAM being specifically designed for small RNA analysis pipelines. Another approach to improve mappability would be to incorporate long-read sequencing methods, as longer reads contain more information and can serve as a way to reduce ambiguity in the context of RNA-seq. Many of the previous aligners like STAR, HISAT2 and GSNAP have been applied to long-read sequencing data after error correction⁵² and have been shown to work well. In addition, algorithms like BLASR⁵³, GraphMap⁵⁴, rHAT⁵⁵, LAMSA⁵⁶, Kart⁵⁷, NGLMR⁵⁸, and lordFAST⁵⁹ have been developed specifically to address the increased length and error rates associated with long read technologies.

Some tools designed to improve mapping rates for repetitive regions work after an initial analysis with one of the tools listed above. These standalone tools can use alignment files as input and then attempt to statistically redistribute the ambiguous reads based on distributions of neighboring alignments. One such algorithm is MMR⁶⁰ which iteratively redistributes ambiguously mapped reads across their respective loci to maximize smoothness of multimapped read distribution in the context of unique reads, or reduce the variance in coverage. Another is a Gibbs sampling method⁶¹ which uses stochastic redistribution of multimapped reads, normalized to the background distribution, in order to iteratively search for the most likely locus of origin. This type of iterative statistical technique for optimal assignment of reads to the correct loci has been picked up and elaborated on by several different groups, and represents a theme throughout the review. While it does not employ the statistical redistribution of reads, CoCo⁶² is a package which corrects and salvages multimapped reads by taking into consideration nested genomic architecture, a common feature associated with TEs.

Analysis

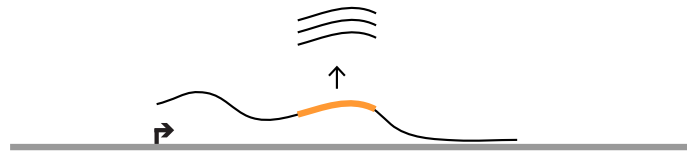
The next step in a general NGS sequencing analysis pipeline is to annotate and quantify those reads which mapped to the genome. The mapping profiles will vary widely based on molecular context of the sequencing library. Each type of NGS data comes with its own challenges in the context of highly repetitive elements. The remaining sections will go through analysis strategies for each of the most common NGS data types in detail. The tools in these sections are listed for reference on Figure 2, where they are grouped by the experimental assays used to generate the data. Table 1 gives references and links to the software for all tools described.

RNA-seq



RSEM
RepEnrich
TETranscripts
TETools
SalmonTE
ERVmap
LIONS
SQUIRE
TeXP
Telescope

sRNA-seq



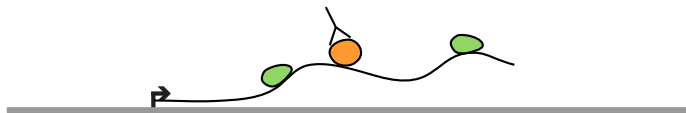
miRDeep2
ShortStack
PiPipes
Chimira
unitas
Oasis 2
TEsmall

ChIP-seq



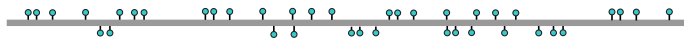
CSEM
MOSAICS
LONUT
DROMPA
Perm-seq
MapRRCon
Crunch

RIP-seq/CLIP-seq



CLIPSeqTools
PROBER
CLAM

DNA methylation-seq



TEPID
EpiTEome

Figure 2 | Published tools available for including repetitive regions in several common genomics analysis protocols. These have been divided into those that are geared toward RNA expression analysis (RNA-seq), small RNA expression analysis (sRNA-seq), genome and chromatin binding factors (ChIP-seq), RNA-binding factors (RIP/CLIP-seq), and DNA methylation analysis (DNA methylation-seq) A table describing these tools (Table 1) also provides links and references for the software and associated publications.

RNA-seq

RNA-seq for expression analysis is one of the most well studied areas in genomics, and this is also reflected in the diversity of tools available for analysis of transcripts from repetitive regions. RNA-seq data derived from short-read sequencing platforms is comprised of small fragments, derived from short single- or paired-end reads tiled across the region of a transcript of origin. Of the tools which have been developed to facilitate transcriptional analysis of repetitive elements, here we will focus on those which take into consideration ambiguously mapped reads. How to address ambiguously mapped reads is an old problem in genome science particularly when using older sequencing technologies from which reads were much shorter (~36 nt) than what we currently consider a short read (~150 nt). These early RNA-seq packages were largely gene centric, as investigation of repetitive elements with these earlier technologies was (and remains) a challenge. However, the basic principles for probabilistic redistribution of ambiguously mapped reads emerged at this time. The first strategies employed a single-step multimapped read redistribution based on the number of uniquely mapped reads at each locus.³⁷ This was followed quickly by an expectation maximization (EM) algorithm to iteratively estimate the most likely expression levels of gene transcripts based on relative counts of unique and multimapped reads⁶³. In addition to probabilistic redistribution of reads, packages like Cufflinks⁶⁴ and HTseq⁶⁵ have multimapper modes where ambiguously mapped reads are weighted by the relative number of genomic alignments (as $1/n$, where n is the number of potential alignments in the genome). The package Scavenger⁶⁶ considers multimapped reads and uses an intermediate consensus assignment with remapping to rescue unmapped reads. Differences in strategies used to address multimapped reads and their associated limitations are outlined in detail by Treangen and Salzberg²³.

As interest broadened to begin investigating transposon expression through RNA-seq explicitly, several packages were developed to handle transposons separately from the rest of the transcriptome. Among the first TE-centric packages was RepEnrich⁶⁷ which functions by creating repetitive element pseudochromosomes, which are a series of contigs that represent all of the genomic instances of each transposon subfamily annotated in RepeatMasker, concatenated onto a single region. These subfamily pseudochromosomes were then used to identify reads that mapped only to one subfamily of transposons, such as the human specific LINE element L1Hs, even if the exact generating locus was still ambiguous. This was able to separate the level of uncertainty to finer detail, such that reads could be described as: unique in the genome, unique to a particular subfamily, or ambiguously mapping to multiple types of transposons. Similar to

RepEnrich, TETools⁶⁸ is another transcript quantification method which uses a detailed annotation file or 'rosette' to facilitate quantification from TE derived reads, and which again aggregates reads at the subfamily level. TeXP⁶⁹ is a package which focuses on LINE-1 elements specifically and models spurious genome transcription to more accurately quantify LINE-1 expression. TETranscripts⁷⁰ was the first TE-centric algorithm to implement statistical read redistribution to handle multimapped reads. TETranscripts uses an expectation maximization algorithm to find the most likely distribution of ambiguously mapped TE-derived RNA-seq reads, and also includes expression estimates for both host genes and transposable elements in the output. After TETranscripts, other packages have been developed to expand the methods used for statistical read redistribution including MMR⁶⁰ and SalmonTE,⁷¹ with SalmonTE being unique in its use of a pseudoalignment strategy from the authors of the original Salmon⁴⁴ package in order to bypass the mapping step typically used in RNA-seq analysis. Yanagi⁷² expands on this pseudoalignment strategy by mapping to a segmented version of the transcriptome to reduce ambiguity of mapping.

In the packages described above, quantification was performed at the subfamily level, as determining the specific expressed genomic loci within a subfamily is quite difficult for transposable elements that are close to the consensus sequence. However, several newer packages have been released to address the need for locus specific quantification of TE derived transcripts. TE-centric packages include SINEsFIND⁷³, and ERVmap⁷⁴ which are specialized for their respective TE family of interest. Two pipelines used genome guided de novo transcriptome assembly with Trinity⁷⁵ to quantify TE expression at a locus specific level: TECandidates⁷⁶ and a pipeline described by Guffanti et al.⁷⁷ More recently, SQUIRE⁷⁸ (Software for Quantifying Interspersed Repeat Expansion), and Telescope⁷⁹ adapted the EM-based read redistribution strategies described above to infer originating loci of ambiguously mapped reads, using uniquely mapped reads surrounding the locus to guide the EM read redistribution.

One of the motivating reasons to study transposable elements is for their influence over regulatory networks in our genome. To address this specifically, a final type of RNA-seq analysis package has been released at the interface of gene-centric and TE-centric models. LIONS⁸⁰ is a novel package which detects novel fusion events that connect TE promoter sequences to downstream coding gene sequences. These chimeric TE/gene transcripts represent one of the many ways that TE promoter elements might affect regulation of adjacent genes.

Small RNA-seq

Cells regulate transposable element expression using multiple strategies. The most potent silencers of TEs in germline cells are small RNAs (sRNAs) of the PIWI-interacting RNA (piRNA) class⁸¹. In somatic tissues, two additional classes of small RNAs contribute to TE silencing: short interfering RNAs (siRNAs) derived from expressed transposon transcripts⁸¹ and the more recently described 3' tRNA derived fragments (3' tRFs)⁸². Therefore, it is integral to the study of transposon biology to consider sRNAs and

accurately quantify their production. To this end, several packages have been released to investigate sRNA species, which prove particularly challenging when derived from repetitive loci in the genome as they are short in length, typically between 18-36 nucleotides. Packages like MiRdeep2⁸³, ShortStack⁸⁴, PiPipes⁸⁵, Chimira⁸⁶, sRNAtoolbox⁸⁷, Oasis 2⁸⁸, and Manatee⁸⁹ have been developed to detect specific types of sRNA loci in the genome and quantify their differential expression. While microRNAs (miRNAs) are not known to play a large role in transposon regulation, a large fraction of miRNAs and other known TE regulatory sRNAs are present in multiple copies in the genome, making TE-focused strategies for multimapped read resolution useful, even for non TE-derived sRNAs. Statistical techniques, including machine learning, have already been extensively employed in the arena of piRNA prediction, a critical step for the ultimate quantification of piRNA reads accumulation in packages like piRNAPredictor⁹⁰, Piano⁹¹, and a k-mer based method described by Zhang et al.⁹² ShortStack after publication was updated to include Butter⁹³ which now performs statistical redistribution of multimapped reads.

These methods described above have largely considered sRNA classes separately, however several packages including Uritas⁹⁴ and TEsSmall⁹⁵ have strived to consider sRNA classes comprehensively to facilitate proper normalization of heterogeneous sRNA libraries, and to facilitate differential expression analysis across classes while taking into consideration ambiguously mapped reads.

While several iterative statistical methods have been employed in the study of sRNAs for annotation and target prediction⁹⁶, there is still much room for improvement in the handling of ambiguously mapped reads for small RNA expression analysis. Many of these issues have been nicely reviewed by Bousios et al.⁹⁷ particularly in the context of plants whose genomes are highly enriched in TEs and where sRNAs form a large component of the TE silencing machinery. Briefly, the chief challenge for applying probabilistic read redistribution algorithms for sRNA loci is that many types of sRNAs accumulate as very short transcripts cut from larger precursors. Often the precursors are rapidly processed and/or would not be caught by sRNA library preparation protocols. For miRNAs, for example, typically only the guide and passenger strands are detected in sRNA-seq libraries, leaving only two short ~22 nucleotide sRNAs and few surrounding reads from the precursor transcript to help guide decisions about the true originating locus. Thus, some loci may be more amenable to statistical inference algorithms, while others need additional assays in order to determine the precise source of sRNA biogenesis.

IP-seq (ChIP, CLIP, and RIP)

In this section, we have grouped together multiple disparate genomics data types that all involve immunoprecipitation-based steps in order to find protein binding sites in nucleic acids. These data can be derived from chromatin bound factors (ChIP-seq) or RNA-binding proteins (CLIP-seq/RIP-seq), but are grouped here as IP-seq because of the similar challenges these data types present for computational analysis pipelines. Typically, the published pipelines for IP-seq data analysis begin by discarding

multimapped reads in order to achieve higher specificity and resolution for the protein binding sites. This can be troublesome when studying proteins which bind to regions rich in repetitive elements. For example, H3K9me3 histone markers are known to be enriched in constitutive heterochromatin⁹⁸, a region of the genome highly enriched in repeat elements. Therefore, when calling H3K9me3 peaks using only uniquely mapped reads, the actual enrichment above background levels may be significantly higher than what is reported, skewing the estimates of background levels and discarding many truly bound regions. While this is a known issue for heterochromatin binding proteins, recent surveys of DNA- and RNA- factors have shown that transposon-derived regulatory elements form a significant fraction of both transcription factor binding sites^{99,18} as well as RNA-binding protein recognition elements^{100,101}.

For ChIP-seq based datasets, it is important to acknowledge the differences and difficulties associated with attempting to detect binding elements for chromatin binding factors and marked histones that typically bind broadly over large areas (broad peaks) as compared to transcription factors, which typically display sharp, narrow peaks. H3K9me3 typically shows a broad peak profile, as these histone marks are found on nucleosomes spread across wide stretches of chromatin. This distribution warrants a different detection strategy than that used for a typical transcription factor, such as MYC, which might occupy narrow binding regions, on the order of ~50-150 nucleotides in a typical assay. This is particularly relevant when these different peaks occur in repetitive genomic regions. The larger the bound region, the more likely it is that some of that genomic sequence will be uniquely mappable, which can guide the inference about read accumulation in adjacent sequences.

To address multimapped reads specifically, packages like the peak caller CSEM¹⁰² have used expectation maximization to redistribute ambiguously mapped ChIP-seq reads based on the distribution of surrounding uniquely mapped reads. Due to the reliance on uniquely mappable reads, these methods function best on broader peaks because they query a larger region, which may be more likely to contain uniquely mappable content. LONUT¹⁰³ calls a set of unique peaks and a set of non-unique peaks, then aggregates both call sets together to remove any redundancy. MOSAiCS¹⁰⁴, while not specifically developed to handle repetitive regions, recommends using the CSEM algorithm as a pre-processing step in order to include multimapped reads. DROMPA¹⁰⁵ and Crunch¹⁰⁶ take into account multimapped reads using a simple 1/n fractional distribution strategy. Crunch subsequently places a large emphasis on motif prediction and annotation. The analysis pipeline MapRRcon¹⁰⁷ uses unique and multimapped reads, but resolves the issue of multimapped read ambiguity by calling peaks on the consensus sequence of transposon subfamilies.

There is still significant room for progress in the arena of ChIP-seq analysis in repetitive regions. It is still difficult to call narrow peaks in repetitive regions, due to the lack of sufficient reads surrounding the locus of interest to guide the inference algorithms. Perm-seq¹⁰⁸ addresses this issue by using the orthogonal dataset of DNAase hypersensitivity profiling for better resolution in repetitive regions of the genome. As sufficient reference datasets become available in multiple cell types and conditions, this may make this

strategy feasible as a general method. In contrast, while broad peak callers tend to include more information within the locus of interest to help guide inference across repetitive regions, the data from these methods tend to have a lower signal-to-noise ratio, such that improvement of broad peak callers generally is still an active area of computational development.

The problems described above in the context of ChIP-seq analysis are compounded in the context of CLIP- and RIP-seq datasets, where one must also normalize for differences in the expression level of the bound transcript substrates. If the bound transcripts contain repetitive regions, or are entirely composed of repetitive elements, one must first find a way to accurately distribute ambiguous reads among the input transcriptome dataset before calling enriched binding sites in particular transcripts. CLIPper¹⁰⁹ was one of the first CLIP-seq pipelines, but was restricted to uniquely mapped reads only. CLIPSeqTools¹¹⁰ is a CLIP- analysis pipeline which randomly assigns ambiguously mapped reads to one of their candidate mapping loci. CLAM¹¹¹ uses expectation maximization algorithms, as described above, to redistribute ambiguously mapped reads between expressed transcripts, but the algorithm works only on the alignment file and does not include information about enriched peaks in its statistical weights. PROBer¹¹² has been developed as a general purpose algorithm for detecting sites of RNA binding or modification (termed ‘toeprint’ profiling) and includes an algorithm for handling multimapped reads using a Gibbs sampler approach to iteratively infer a single “best” alignment for each read. While PROBer does include steps to handle multimapped reads, it was not developed specifically for TEs, and thus has not been tested on highly repetitive regions, such as TEs that are very close to the consensus.

DNA methylation-seq

We have detailed several methods to assess differential expression, and protein binding in the context of repetitive elements. However, a critical component to the understanding of transposon biology is the analysis of DNA methylation as it is the main mechanism by which transposons are transcriptionally silenced long term¹¹³. To assess DNA methylation, particularly the 5-methylcytosine (5-mC) modification, several techniques have been developed and compared¹¹⁴. In brief, the most common method to assess DNA methylation is bisulfite sequencing: whole genome DNA sequencing following bisulfite conversion of all non-methylated cytosine residues to uracil. Bisulfite sequencing based methods can be non-directional¹¹⁵, or directional¹¹⁶ allowing one to reduce the ambiguity of strand of origin. One of the first analysis pipelines developed for high-throughput bisulfite sequencing was in *Arabidopsis*¹¹⁶ and analysis was performed in conjunction with sRNA-seq datasets. In this pipeline, ambiguously mapped reads were discarded by mapping to a repeat masked version of the genome, a technique once commonly used in animal systems to reduce mapping ambiguity in the context of bisulfite induced C>T conversions¹¹⁷. Bisulfite sequencing analysis differs significantly from other analysis pipelines in that often two reference genomes are used, one which contains converted cytosines in addition to the original reference genome. In this context, what are considered ambiguous reads are those reads which map to both the converted and unconverted reference genomes. This compounds the difficulty of assigning multimapped reads, such that many published bisulfite sequencing software packages choose not to

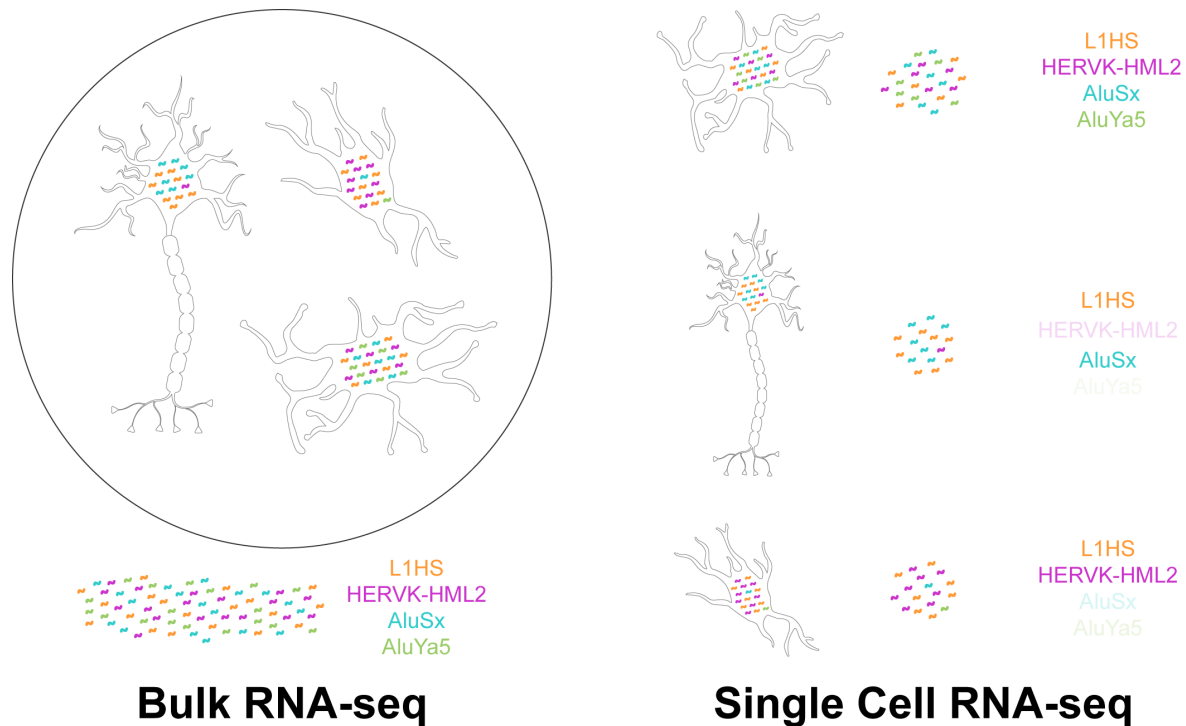
include multimapped reads to avoid this confounded ambiguity (see Table 1). The most commonly used pipelines for bisulfite sequencing reads including BSMAP¹¹⁸, Bismark¹¹⁹, MOABS¹²⁰, and BS-Seeker3¹²¹, none of which include probabilistic handling of multimapper reads. For a more comprehensive list of non-TE-specific methylation pipelines, please see the review by Adusumalli et al.¹²² and the supplemental material of a recently published pipeline, bicycle¹²³. Here, confounding between ambiguity in bisulfite conversion rates, non-reference polymorphisms, and read non-uniqueness can complicate the statistical tests used to determine if a site in the genome is differentially methylated. Thus, this represents an area of computational genomics that could benefit greatly from further development.

Since DNA methylation is a critical mechanism by which transposons are silenced, several groups have used new methods to improve methylation analysis for TEs. TEPID¹²⁴ and epiTEome¹²⁵ were designed to improve analysis of TE methylation levels by including the analysis of split reads that cross junctions between TEs and uniquely mappable genome regions. An approach employed to assess the low mappability of young transposable elements, like L1-Ta, in the human genome was repurposed to align bisulfite reads to a consensus sequence as described in Shukla et al.¹²⁶. One interesting method to improve methylation analysis is to first rigorously determine the average bisulfite conversion rates genome-wide, then use this as a parameter to tease apart mapping ambiguities from differences in conversion rates, as done by Noshay et al.¹²⁷. Despite these improvements, DNA methylation analysis is still a difficult bioinformatic challenge that would benefit from further study.

Single Cell RNA-seq

All of the software described above has been geared towards genomics datasets generated from bulk tissue samples. However, bulk profiling of heterogeneous cell populations only provides averages that obscure underlying variability of TE expression across cell types, as illustrated in Figure 3. This problem is further amplified when aggregating transcriptional signal across numerous loci within high copy-number TE families. It remains largely unknown how TE de-repression varies between individual cells, what factors drive such differences, and how this variability might affect cellular phenotypes. Single cell RNA-sequencing (scRNA-seq) promises to answer some of those questions and has already redefined our knowledge about tissue composition and gene regulatory networks¹²⁸. While its broad application has so far been largely restricted to the study of gene activity patterns, a few pioneering studies have utilized first-generation protocols to identify TE expression dynamics across single pre-implantation embryonic cells^{129,130}. Those early efforts were largely limited by small cell numbers, high sequencing burden per cell, and lack of molecular barcode counts to estimate true transcriptional output, thus preventing broad-scale adaptation. Since then, the increasing demand in single cell transcriptome data has seen an unprecedented expansion of available scRNA-seq protocols with considerably improved throughput, robustness, and error-rates¹³¹. One such publication was by Guo et al.¹³² where the number of cells were scaled up allowing for investigation into TE dynamics in spermatogenesis.

Figure 3 | Comparison of bulk RNA-seq versus single cell RNA-seq. Heterogeneity in expression profiles across cell types is masked by bulk sequencing methods. Transposable element (TE) expression may vary across cell types, between cells of the same type, and within the same cells across time. Single cell methods are necessary to reveal this heterogeneity, but software for single cell data analysis is not currently optimized for handling TEs.



Despite such experimental advancements, inherent design principles of scRNA-seq protocols that cooperate with the well-known challenges of TE transcriptome analysis have so far prevented their common application for the study of TE expression at single cell resolution (**Fig 4**). For example, many popular methods quantify RNA molecules at the 3' end of polyadenylated mRNAs^{133–137} and therefore depend on accurate reference models to bridge the gap between polyadenylation sites and the corresponding transcript isoform and/or promoter. This is problematic for TE-derived transcripts, which are generally poorly annotated in many species. While protocols with full-length transcript coverage might alleviate some of those problems, the naïve assignment of reads to the nearest TE interval can still lead to erroneous assignment, misattribution of intronic reads from unprocessed pre-mRNAs, and hence misinterpretation of TE de-repression. Full-length protocols additionally suffer from higher sequencing burden, often lack of unique molecular identifiers to account for PCR duplicates, and potentially higher background TE read coverage due to intronic signal originating from pre-mRNAs^{138,139}.

A potential solution to minimize misattribution problems are 5' end based scRNA-seq protocols that incorporate a template switch oligo (TSO) towards the start of transcription initiation^{140,141}. Although incomplete processing and premature TSO incorporation during library preparation might vary between transcripts and cells, such protocols have already been successfully used to map alternative transcription start sites between individual cells¹⁴². Importantly, a recent study also demonstrated its utility to quantify unexpected variability in TE promoter activity between thousands of single cancer cells following epigenetic therapy¹⁴³. However, the problem of premature TSO incorporation, combined with the pervasive nature of TEs, and technical noise inherent to all current scRNA-seq protocols requires dedicated strategies to mitigate the danger of spurious estimates of TE cell-to-cell variation. To the best of our knowledge, no peer-reviewed computational pipeline currently combines such features with the reliable quantification of TEs at single cell resolution, but unpublished efforts already aim to facilitate TE single cell analysis for a wide array of available scRNA-seq protocols (<https://tanaylab.github.io/Repseq/>). With the continuous methodological advancements and the increasing interest in TE biology, we anticipate a rapid progress towards the routine quantification of TEs in individual cells that will be accompanied by the discovery of unprecedented heterogeneity in TE transcription patterns.

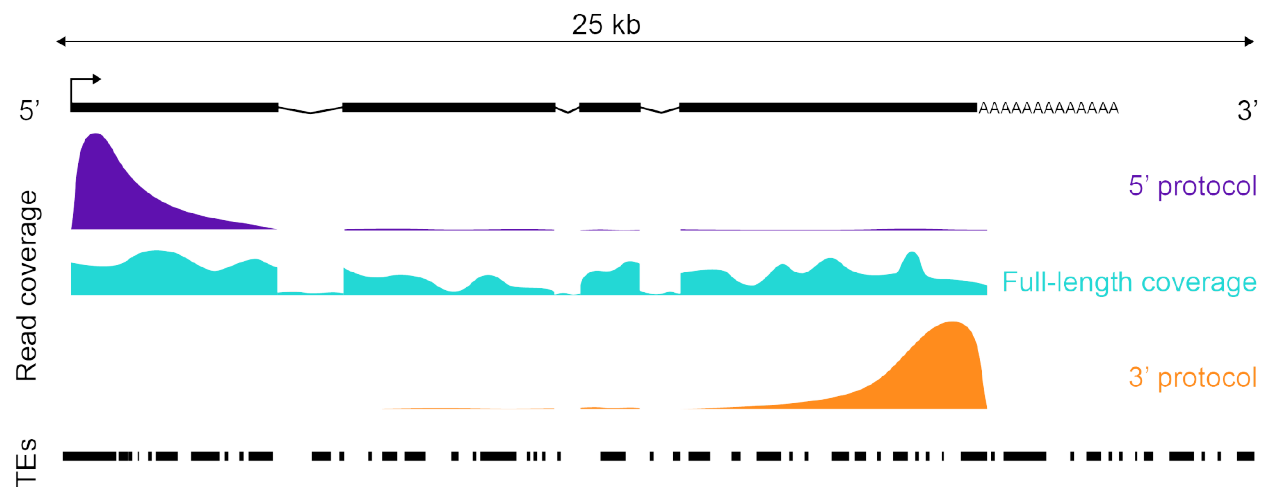


Figure 4 | Impact of different RNA-seq library strategies on read coverage along a TE-derived transcript with exon/intron structure. (top). TE intervals are shown at the bottom. Briefly, scRNA-seq protocols that offer full-length transcript coverage provide the best means to identify full length transcribed TEs in a locus-specific manner, but this method suffers from noise due to intronic TEs in host genes that might be mistaken for expressed TE transcripts as well as the inability to barcode individual mRNA molecules. 5' and 3' based protocols allow for barcodes that enable mRNA molecule counting, with 5' protocols also offering the ability to detect TE transcripts originating from proper TE promoters.

Conclusions

What is now doable?

The last years have seen a general improvement in sequencing read length, making it possible to study the majority of transposable elements in a genome wide fashion. For particularly young and less diverged families, we have discussed at length the strides made in genome biology to address the difficulties of treating ambiguously mapped sequencing fragments for differential expression and binding analyses. In the context of highly repetitive regions of the genome these difficulties are compounded, particularly for the most active TEs, which remain close to their consensus sequence and thus are the most difficult to map. The greatest progress has been made with RNA-seq data analysis, as we have progressed from using simple fractional assignments of multimapped reads within genes to approaching true locus specific resolution in the most repetitive regions of the genome – such as the L1HS subfamily, active Alu families, and composite SVA elements. Progress has been made in the realm of sRNA analysis as these improved algorithms for RNA-seq analysis have now been incorporated into sRNA-seq data analysis pipelines. In immunoprecipitation based assays, for ChIP- and CLIP-seq datasets, efforts have been made to use probabilistic read redistribution for peaks within repetitive regions, but challenges remain.

What is still hard?

sRNA-seq data contains a large proportion of multimapped reads, and while significant effort has been put forth to leverage advanced iterative statistical methods for novel sRNA discovery and target prediction, these methods have not been as widely applied to sRNA-seq transcript quantification. This may be attributed to the tight distribution of sRNA reads across their mapping loci, making it difficult to garner locus-specific information from adjacent reads. Moreover, these much shorter reads (18-30 nt) are intrinsically less unique in the genome than longer sequences.

In ChIP-seq data, the expected profile of read distributions can vary widely from the typically tall, narrow peaks associated with most transcription factor profiles or RNA-binding proteins to the broader, shorter, and noisier peaks associated with some marked histones, such as H3K9me3. Algorithms have been developed to address both types of ChIP-seq profiles. Yet the lines between these categories can be blurred, and there is a large tradeoff between the window size in peak calling and the ability to use uniquely mapped reads to probabilistically reassign all other reads to a particular locus. One area of active research for broader regions would be to incorporate multimapped reads into segmentation models which allow for the detection of changes in peak landscape, as opposed to simply calling the absence or presence of individual peaks.

Single cell RNA-seq (scRNA-seq) represents one of the newest genomic assays to be used for TE expression profiling, and as such, remains an area of greatest need for improvements in software packages specifically designed to handle the complexities inherent in TE genomics. Efforts are already underway, but as yet no published software packages for scRNA-seq are available. That said, many standard scRNA-seq packages could be adapted for this use, as in the example protocol described above. However, as discussed in detail, differences in the experimental protocols used to generate scRNA-

seq libraries will have a large impact upon the interpretability of the data, and this is particularly problematic for TE expression analysis.

Two types of analysis which largely do not include multimapped reads are assays for transposase accessible chromatin using sequencing (ATAC-seq)¹⁴⁴ and Hi-C¹⁴⁵, an extension of chromosome conformation capture (3C). The read distributions for ATAC-seq data greatly resemble those of ChIP-seq and this analysis encounters similar computational difficulties when studying repetitive regions of the genome. Fortunately, as this analysis is similar to ChIP-seq there has already been significant effort which could be incorporated into ATAC-seq analysis. Adapting Hi-C pipelines to take into account multimapped reads is still a difficult task as this type of analysis already requires the resolution of chimeric reads representing genomic proximity. mHiC¹⁴⁶ has been developed to address this issue, but the relative sensitivity to highly repetitive transposon regions is unclear. Significant work has been done using these methods to address the role of transposons in genome architecture and the transition from the embryonic cell state to early embryonic like cells^{147–149}. These analyses can only improve as better methods for handling repetitive reads are included.

What new technology needs to be developed?

Long read sequencing technologies promise to solve many issues inherent in the assays described above. Once the issues with throughput and error rates can be solved, long read sequencing would enable the isolation of entire transcripts and, if correctly barcoded, would also allow for accurately calibrated expression estimates. These technologies could also be combined with antibody based pulldowns and endonuclease based footprinting assays, to accurately call cis-regulatory regions derived from TEs. Finally, long-read genome resequencing assays that sequence through highly repetitive genome regions may allow for better genomic annotations that will benefit all of the applications described above. To this end, not only must new experimental protocols be developed which emphasize longer reads, but new computational pipelines must also be developed to ensure that these long read analysis pipelines properly handle and account for the complications inherent in addressing TE genomics.

Acknowledgments

We would like to acknowledge Ying Jin for her helpful conversations about TE genomics. We apologize for any software we failed to cite in this review pertinent to the study of repetitive elements in genome biology. Schematic images were adapted with permission from Servier Medical Art (<https://smart.servier.com>). MGH receives support from the Rita Allen Foundation and the Chan-Zuckerberg Initiative. KO acknowledges funding from the National Science Foundation Graduate Research Fellowship and NIH training grant (2T32GM065094-16). DB is supported by an EMBO long-term fellowship (ALTF 900-2017).

Competing Interests

The authors have no competing interests to declare.

Author Contributions

All authors contributed to the writing and editing of this manuscript.

References

- 1 Pace II JK, Feschotte C. The evolutionary history of human DNA transposons : Evidence for intense activity in the primate lineage. *Genome Res* 2007;422–32. <https://doi.org/10.1101/gr.5826307.422>.
- 2 de Koning APJ, Gu W, Castoe TA, Batzer MA, Pollock DD. Repetitive Elements May Comprise Over Two-Thirds of the Human Genome. *PLoS Genet* 2011;7:. <https://doi.org/10.1371/journal.pgen.1002384>.
- 3 Schnable PS, Ware D, Fulton RS, Stein JC, Wei F, Pasternak S, *et al.* The B73 Maize Genome: Complexity, Diversity, and Dynamics. *Science (80-)* 2009;326:1112–5. <https://doi.org/10.1038/nature03895>.
- 4 Feschotte C, Pritham EJ. DNA Transposons and the Evolution of Eukaryotic Genomes. *Annu Rev Genet* 2007;41:331–68. <https://doi.org/10.1146/annurev.genet.40.110405.090448>.
- 5 Slotkin RK, Martienssen R. Transposable elements and the epigenetic regulation of the genome. *Nat Rev Genet* 2007;8:272–85. <https://doi.org/10.1038/nrg2072>.
- 6 Levin HL, Moran J V. Dynamic interactions between transposable elements and their hosts. *Nat Rev Genet* 2011;12:615–27. <https://doi.org/10.1038/nrg3030>.
- 7 Boissinot S, Chevret P, Furano A V. L1 (LINE-1) Retrotransposon Evolution and Amplification in Recent Human History. *Mol Biol Evol* 2000;17:915–28.
- 8 Sheen F, Sherry ST, Risch GM, Robichaux M, Nasidze I, Stoneking M, *et al.* Reading between the LINEs : Human Genomic Variation Induced by LINE-1 Retrotransposition. *Genome Res* 2000;10:1496–508. <https://doi.org/10.1101/gr.149400.6>.
- 9 Brouha B, Schustak J, Badge RM, Lutz-Prigge S, Farley AH, Moran J V, *et al.* Hot L1s account for the bulk of retrotransposition in the human population. *Proc Natl Acad Sci* 2003;100:5280–5. <https://doi.org/10.1073/pnas.0831042100>.
- 10 Lu X, Sachs F, Ramsay L, Jacques P-É, Göke J, Bourque G, *et al.* The retrovirus HERVH is a long noncoding RNA required for human embryonic stem cell identity. *Nat Struct Mol Biol* 2014;21:423–5. <https://doi.org/10.1038/nsmb.2799>.
- 11 Grow EJ, Flynn RA, Chavez SL, Bayless NL, Wossidlo M, Wesche DJ, *et al.* Intrinsic retroviral reactivation in human preimplantation embryos and pluripotent cells. *Nature* 2015;522:221–5. <https://doi.org/10.1038/nature14308>.
- 12 Jachowicz JW, Bing X, Pontabry J, Bošković A, Rando OJ, Torres-Padilla M-E. LINE-1 activation after fertilization regulates global chromatin accessibility in the early mouse embryo. *Nat Genet* 2017;49:1502–10. <https://doi.org/10.1038/ng.3945>.
- 13 Rowe HM, Kapopoulou A, Corsinotti A, Fasching L, Macfarlan TS, Tarabay Y, *et al.* TRIM28 repression of retrotransposon-based enhancers is necessary to preserve transcriptional dynamics in embryonic stem cells. *Genome Res* 2013;23:452–61. <https://doi.org/10.1101/gr.147678.112>.
- 14 Percharde M, Lin C-J, Yin Y, Guan J, Peixoto GA, Bulut-Karslioglu A, *et al.* A

- LINE1-Nucleolin Partnership Regulates Early Development and ESC Identity. *Cell* 2018;**174**:391–405. <https://doi.org/10.1016/j.cell.2018.05.043>.
- 15 Chuong EB, Elde NC, Feschotte C. Regulatory activities of transposable elements: from conflicts to benefits. *Nat Rev Genet* 2017;**18**:71–86. <https://doi.org/10.1038/nrg.2016.139>.
- 16 Imbeault M, Helleboid P-Y, Trono D. KRAB zinc-finger proteins contribute to the evolution of gene regulatory networks. *Nature* 2017;**543**:550–4. <https://doi.org/10.1038/nature21683>.
- 17 Chuong EB, Elde NC, Feschotte C. Regulatory evolution of innate immunity through co-option of endogenous retroviruses. *Science (80-)* 2016;**351**:1083–8.
- 18 Sundaram V, Cheng Y, Ma Z, Li D, Xing X, Edge P, *et al*. Widespread contribution of transposable elements to the innovation of gene regulatory networks. *Genome Res* 2014;**24**:1–15. <https://doi.org/10.1101/gr.168872.113>.
- 19 Pontis J, Planet E, Offner S, Turelli P, Duc J, Coudray A, *et al*. Hominoid-Specific Transposable Elements and KZFPs Facilitate Human Embryonic Genome Activation and Control Transcription in Naive Human ESCs. *Cell Stem Cell* 2019;**24**:724–735.e5. <https://doi.org/10.1016/j.stem.2019.03.012>.
- 20 Venuto D, Bourque G. Identifying co-opted transposable elements using comparative epigenomics. *Dev Growth Differ* 2018;**60**:53–62. <https://doi.org/10.1111/dgd.12423>.
- 21 Raviram R, Rocha PP, Luo VM, Swanzey E, Miraldi ER, Chuong EB, *et al*. Analysis of 3D genomic interactions identifies candidate host genes that transposable elements potentially regulate. *Genome Biol* 2018;**19**:1–19. <https://doi.org/10.1186/s13059-018-1598-7>.
- 22 Goerner-Potvin P, Bourque G. Computational tools to unmask transposable elements. *Nat Rev Genet* 2018;**19**:688–704. <https://doi.org/10.1038/s41576-018-0050-x>.
- 23 Treangen TJ, Salzberg SL. Repetitive DNA and next-generation sequencing : computational challenges and solutions. *Nat Rev Genet* 2012;**13**:. <https://doi.org/10.1038/nrg3117>.
- 24 Jain M, Olsen HE, Turner DJ, Stoddart D, Bulazel K V, Paten B, *et al*. Linear assembly of a human centromere on the Y chromosome. *Nat Biotechnol* 2018;**36**:321–3. <https://doi.org/10.1038/nbt.4109>.
- 25 Gordon D, Huddleston J, Chaisson MJP, Hill CM, Kronenberg ZN, Munson KM, *et al*. Long-read sequence assembly of the gorilla genome. *Science (80-)* 2016;**352**:. <https://doi.org/10.1126/science.aae0344>.
- 26 Wong TN, Miller CA, Jotte MRM, Bagegni N, Baty JD, Schmidt AP, *et al*. Cellular stressors contribute to the expansion of hematopoietic clones of varying leukemic potential. *Nat Commun* 2018;**9**:1–10. <https://doi.org/10.1038/s41467-018-02858-0>.
- 27 Chaisson MJP, Sanders AD, Zhao X, Malhotra A, Porubsky D, Rausch T, *et al*. Multi-platform discovery of haplotype-resolved structural variation in human genomes. *Nat Commun* 2019;**10**:1–16. <https://doi.org/10.1038/s41467-018-08148-z>.
- 28 Bao W, Kojima KK, Kohany O. Repbase Update , a database of repetitive elements in eukaryotic genomes. *Mob DNA* 2015;**6**:4–9.

- <https://doi.org/10.1186/s13100-015-0041-9>.
- 29 Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, *et al*. The Human Genome Browser at UCSC. *Genome Res* 2002;**12**:996–1006. <https://doi.org/10.1101/gr.229102>.
- 30 Amselem J, Cornut G, Choisne N, Alaux M, Alfama-Depauw F, Jamilloux V, *et al*. RepetDB : a unified resource for transposable element references. *Mob DNA* 2019;**10**:4–11. <https://doi.org/https://doi.org/10.1186/s13100-019-0150-y>.
- 31 Gifford RJ, Blomberg J, Coffin JM, Fan H, Heidmann T, Mayer J, *et al*. Nomenclature for endogenous retrovirus (ERV) loci. *Retrovirology* 2018;**15**:1–11. <https://doi.org/10.1186/s12977-018-0442-1>.
- 32 Hubley R, Finn RD, Clements J, Eddy SR, Jones TA, Bao W, *et al*. The Dfam database of repetitive DNA families. *Nucleic Acids Res* 2016;**44**:81–9. <https://doi.org/10.1093/nar/gkv1272>.
- 33 Wicker T, Matthews DE, Keller B. TREP : a database for Triticeae repetitive elements. *Trends Plant Sci* 2002;**7**:561–2.
- 34 Yi F, Jia Z, Xiao Y, Ma W, Wang J. SPTEdb : a database for transposable elements in salicaceous plants. *Database* 2018:1–8. <https://doi.org/10.1093/database/bay024>.
- 35 Yi F, Ling J, Xiao Y, Zhang H, Ouyang F, Wang J. ConTEdb : a comprehensive database of transposable elements in conifers. *Database* 2018:1–7. <https://doi.org/10.1093/database/bay131>.
- 36 Nussbaumer T, Martis MM, Roessner SK, Pfeifer M, Bader KC, Sharma S, *et al*. MIPS PlantsDB : a database framework for comparative plant genome research. *Nucleic Acids Res* 2013;**41**:1144–51. <https://doi.org/10.1093/nar/gks1153>.
- 37 Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B. Mapping and quantifying mammalian transcriptomes by RNA-Seq. *Nat Methods* 2008;**5**:621–8. <https://doi.org/10.1038/NMETH.1226>.
- 38 Mills RE, Bennett EA, Iskow RC, Devine SE. Which transposable elements are active in the human genome? *Trends Genet* 2007;**23**:183–91. <https://doi.org/10.1016/j.tig.2007.02.006>.
- 39 Sexton CE, Han M V. Paired-end mappability of transposable elements in the human genome. *Mob DNA* 2019;**10**:1–11. <https://doi.org/https://doi.org/10.1186/s13100-019-0172-5>.
- 40 Feschotte C. Transposable elements and the evolution of regulatory networks. *Nat Rev Genet* 2008;**9**:397–405.
- 41 Li H, Durbin R. Fast and accurate long-read alignment with Burrows – Wheeler transform. *Bioinformatics* 2010;**26**:589–95. <https://doi.org/10.1093/bioinformatics/btp698>.
- 42 Langmead B. Aligning short sequencing reads with Bowtie. *Curr Protoc Bioinforma* 2010:1–24. <https://doi.org/10.1002/0471250953.bi1107s32>.
- 43 Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq quantification. *Nat Biotechnol* 2016;**34**:525–7. <https://doi.org/10.1038/nbt.3519>.
- 44 Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. Salmon provides fast and bias-aware quantification of transcript expression. *Nat Methods* 2017;**14**:417–9. <https://doi.org/10.1038/nmeth.4197>.
- 45 Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, *et al*. STAR:

- ultrafast universal RNA-seq aligner. *Bioinformatics* 2013;**29**:15–21. <https://doi.org/10.1093/bioinformatics/bts635>.
- 46 Kim D, Paggi JM, Park C, Bennett C, Salzberg SL. Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nat Biotechnol* 2019;**37**:907–15. <https://doi.org/10.1038/s41587-019-0201-4>.
- 47 Wu TD, Reeder J, Lawrence M, Becker G, Brauer MJ. GMAP and GSNAP for Genomic Sequence Alignment: Enhancements to Speed, Accuracy, and Functionality. In: Mathé E, Davis S, editors. *Stat. Genomics Methods Protoc*. New York, NY: Springer New York; 2016. p. 283–334.
- 48 Grant GR, Farkas MH, Pizarro AD, Lahens NF, Schug J, Brunk BP, *et al*. Comparative analysis of RNA-Seq alignment algorithms and the RNA-Seq unified mapper (RUM). *Bioinformatics* 2011;**27**:2518–28. <https://doi.org/10.1093/bioinformatics/btr427>.
- 49 Li H. Sequence analysis Minimap2 : pairwise alignment for nucleotide sequences. *Bioinformatics* 2018;**34**:3094–100. <https://doi.org/10.1093/bioinformatics/bty191>.
- 50 Baruzzo G, Hayer KE, Kim EJ, Camillo B Di, Fitzgerald GA, Grant GR. Simulation-based comprehensive benchmarking of RNA-seq aligners. *Nat Methods* 2017;**14**:135–9. <https://doi.org/10.1038/nmeth.4106>.
- 51 Fletcher SJ, Boden M, Mitter N, Carroll BJ. SCRAM : a pipeline for fast index-free small RNA read alignment and visualization. *Bioinformatics* 2018;**34**:2670–2. <https://doi.org/10.1093/bioinformatics/bty161>.
- 52 Krizanovic K, Echchiki A, Roux J, Sikic M. Evaluation of tools for long read RNA-seq splice-aware alignment. *Bioinformatics* 2018;**34**:748–54. <https://doi.org/10.1093/bioinformatics/btx668>.
- 53 Chaisson MJ, Tesler G. Mapping single molecule sequencing reads using basic local alignment with successive refinement (BLASR): application and theory. *BMC Bioinformatics* 2012;**13**:.
- 54 Sovic I, Sikic M, Wilm A, Fenlon SN, Chen S, Nagarajan N. Fast and sensitive mapping of nanopore sequencing reads with GraphMap. *Nat Commun* 2016;**7**:. <https://doi.org/10.1038/ncomms11307>.
- 55 Liu B, Guan D, Teng M, Wang Y. rHAT : fast alignment of noisy long reads with regional hashing. *Bioinformatics* 2015;**32**:1625–31. <https://doi.org/10.1093/bioinformatics/btv662>.
- 56 Liu B, Gao Y, Wang Y. LAMSA: fast split read alignment with long approximate matches. *Bioinformatics* 2017;**33**:192–201. <https://doi.org/10.1093/bioinformatics/btw594>.
- 57 Lin H, Hsu W. Kart: a divide-and-conquer algorithm for NGS read alignment. *Bioinformatics* 2017;**33**:2281–7. <https://doi.org/10.1093/bioinformatics/btx189>.
- 58 Sedlazeck FJ, Rescheneder P, Smolka M, Fang H, Nattestad M, Haeseler A Von, *et al*. Accurate detection of complex structural variations using single-molecule sequencing. *Nat Methods* 2018;**15**:461–8. <https://doi.org/10.1038/s41592-018-0001-7>.
- 59 Haghshenas E, Sahinalp SC, Hach F. lordFAST : sensitive and Fast Alignment Search Tool for LOnoisy Read sequencing Data. *Bioinformatics* 2019;**35**:20–7. <https://doi.org/10.1093/bioinformatics/bty544>.
- 60 Kahles A, Behr J, Räscht G. MMR : a tool for read multi-mapper resolution.

- Bioinformatics* 2016;**32**:770–2. <https://doi.org/10.1093/bioinformatics/btv624>.
- 61 Wang J, Huda A, Lunyak V V, Jordan IK. A Gibbs sampling strategy applied to the mapping of ambiguous short-sequence tags. *Bioinformatics* 2010;**26**:2501–8. <https://doi.org/10.1093/bioinformatics/btq460>.
- 62 Deschamps-Francoeur G, Boivin V, Sherif AE, Scott MS. CoCo: RNA-seq read assignment correction for nested genes and multimapped reads. *Bioinformatics* 2019:1–9. <https://doi.org/10.1093/bioinformatics/btz433>.
- 63 Li B, Ruotti V, Stewart RM, Thomson JA, Dewey CN. RNA-Seq gene expression estimation with read mapping uncertainty. *Bioinformatics* 2010;**26**:493–500. <https://doi.org/10.1093/bioinformatics/btp692>.
- 64 Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, Baren MJ Van, *et al.* Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol* 2010;**28**:516–20. <https://doi.org/10.1038/nbt.1621>.
- 65 Anders S, Pyl PT, Huber W. HTSeq — a Python framework to work with high-throughput sequencing data. *Bioinformatics* 2015;**31**:166–9. <https://doi.org/10.1093/bioinformatics/btu638>.
- 66 Yang A, Tang JYS, Troup M, Ho JWK. Scavenger: A pipeline for recovery of unaligned reads utilising similarity with aligned reads. *F1000 Res* 2019:1–20.
- 67 Criscione SW, Zhang Y, Thompson W, Sedivy JM, Neretti N. Transcriptional landscape of repetitive elements in normal and cancer human cells. *BMC Genomics* 2014;**15**:1–17. <https://doi.org/10.1186/1471-2164-15-583>.
- 68 Lerat E, Fablet M, Modolo L, Lopez-Maestre H, Vieira C. TETOOLS facilitates big data expression analysis of transposable elements and reveals an antagonism between their activity and that of piRNA genes. *Nucleic Acids Res* 2017;**45**:1–12. <https://doi.org/10.1093/nar/gkw953>.
- 69 Navarro FC, Hoops J, Bellfy L, Cerveira E, Zhu Q, Zhang C, *et al.* TeXP: Deconvolving the effects of pervasive and autonomous transcription of transposable elements. *PLOS Comput Biol* 2019;**15**:1–19. <https://doi.org/10.1371/journal.pcbi.1007293> August.
- 70 Jin Y, Tam OH, Paniagua E, Hammell M. Tetranscripts: A package for including transposable elements in differential expression analysis of RNA-seq datasets. *Bioinformatics* 2015;**31**:3593–9. <https://doi.org/10.1093/bioinformatics/btv422>.
- 71 Jeong H-H, Yalamanchili HK, Guo C, Shulman JM, Liu Z. An Ultra-Fast and Scalable Quantification Pipeline for Transposable Elements from next Generation Sequencing Data.
- 72 Gunady MK, Mount SM, Bravo HC. Fast and interpretable alternative splicing and differential gene-level expression analysis using transcriptome segmentation with Yanagi. *bioRxiv* 2018:1–23. <https://doi.org/10.1101/364281>.
- 73 Carnevali D, Conti A, Pellegrini M, Dieci G. Whole-genome expression analysis of mammalian-wide interspersed repeat elements in human cell lines. *Dna Res* 2017;**24**:59–69. <https://doi.org/10.1093/dnares/dsw048>.
- 74 Tokuyama M, Kong Y, Song E, Jayewickreme T, Kang I, Iwasaki A. ERVmap analysis reveals genome-wide transcription of human endogenous retroviruses. *Proc Natl Acad Sci* 2018;**115**:. <https://doi.org/10.1073/pnas.1814589115>.
- 75 Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, *et al.*

- De novo transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis. *Nat Protoc* 2013;**8**:1494–512. <https://doi.org/10.1038/nprot.2013.084>.
- 76 Valdebenito-Maturana B, Riadi G. TEcandidates: prediction of genomic origin of expressed transposable elements using RNA-seq data. *Bioinformatics* 2018;**34**:3915–6. <https://doi.org/10.1093/bioinformatics/bty423>.
- 77 Guffanti G, Bartlett A, Klengel T, Klengel C, Hunter R, Glinsky G, *et al*. Novel Bioinformatics Approach Identifies Transcriptional Profiles of Lineage-Specific Transposable Elements at Distinct Loci in the Human Dorsolateral Prefrontal Cortex. *Mol Biol Evol* 2018;**35**:2435–53. <https://doi.org/10.1093/molbev/msy143>.
- 78 Yang WR, Ardeljan D, Pacyna CN, Payer LM, Burns KH. SQUIRE reveals locus-specific regulation of interspersed repeat expression. *Nucleic Acids Res* 2019;**47**:1–16. <https://doi.org/10.1093/nar/gky1301>.
- 79 Bendall ML, Mulder M De, Iñiguez LP, Lecanda-Sánchez A, Pérez-Losada M, Ostrowski MA, *et al*. Telescope : Characterization of the retrotranscriptome by accurate estimation of transposable element expression. *PLOS Comput Biol* 2019:1–25.
- 80 Babaian A, Thompson IR, Lever J, Gagnier L, Karimi MM, Mager DL. LIONS: analysis suite for detecting and quantifying transposable element initiated transcription from RNA-seq. *Bioinformatics* 2019:1–3. <https://doi.org/10.1093/bioinformatics/btz130>.
- 81 Malone CD, Hannon GJ. Small RNAs as Guardians of the Genome. *Cell* 2009;**136**:656–68. <https://doi.org/10.1016/j.cell.2009.01.045>.
- 82 Schorn AJ, Gutbrod MJ, LeBlanc C, Martienssen R. LTR-Retrotransposon Control by tRNA-Derived Small RNAs. *Cell* 2017;**170**:61–71. <https://doi.org/10.1016/j.cell.2017.06.013>.
- 83 Friedländer MR, Mackowiak SD, Li N, Chen W, Rajewsky N. miRDeep2 accurately identifies known and hundreds of novel microRNA genes in seven animal clades. *Nucleic Acids Res* 2012;**40**:37–52. <https://doi.org/10.1093/nar/gkr688>.
- 84 Axtell MJ. ShortStack: Comprehensive annotation and quantification of small RNA genes. *RNA* 2013;**19**:740–51. <https://doi.org/10.1261/rna.035279.112>.
- 85 Han BW, Wang W, Zamore PD, Weng Z. piPipes: A set of pipelines for piRNA and transposon analysis via small RNA-seq, RNA-seq, degradome-and CAGE-seq, ChIP-seq and genomic DNA sequencing. *Bioinformatics* 2015;**31**:593–5. <https://doi.org/10.1093/bioinformatics/btu647>.
- 86 Vitsios DM, Enright AJ. Chimira: Analysis of small RNA sequencing data and microRNA modifications. *Bioinformatics* 2015;**31**:3365–7. <https://doi.org/10.1093/bioinformatics/btv380>.
- 87 Rueda A, Barturen G, Lebrón R, Gómez-Martín C, Alganza Á, Oliver JL, *et al*. sRNAtoolbox: an integrated collection of small RNA research tools. *Nucleic Acids Res* 2015;**43**:W467–73. <https://doi.org/10.1093/nar/gkv555>.
- 88 Rahman RU, Gautam A, Bethune J, Sattar A, Fiosins M, Magruder DS, *et al*. Oasis 2: Improved online analysis of small RNA-seq data. *BMC Bioinformatics* 2018;**19**:1–10. <https://doi.org/10.1186/s12859-018-2047-z>.
- 89 Handzlik JE, Tastsoglou S, Vlachos IS, Hatzigeorgiou AG. Manatee: detection

- and quantification of small non-coding RNAs from next-generation sequencing data. *bioRxiv* 2019. <https://doi.org/10.1101/662007>.
- 90 Li D, Luo L, Zhang W, Liu F, Luo F. A genetic algorithm-based weighted ensemble method for predicting transposon-derived piRNAs. *BMC Bioinformatics* 2016;**17**:1–11. <https://doi.org/10.1186/s12859-016-1206-3>.
- 91 Wang K, Liang C, Liu J, Xiao H, Huang S, Xu J, *et al.* Prediction of piRNAs using transposon interaction and a support vector machine. *BMC Bioinformatics* 2014;**15**:1–8. <https://doi.org/10.1186/s12859-014-0419-6>.
- 92 Zhang Y, Wang X, Kang L. A k-mer scheme to predict piRNAs and characterize locust piRNAs. *Bioinformatics* 2011;**27**:771–6. <https://doi.org/10.1093/bioinformatics/btr016>.
- 93 Axtell MJ. Butter: High-precision genomic alignment of small RNA-seq data. *bioRxiv* 2014:1–16. <https://doi.org/10.1101/007427>.
- 94 Gebert D, Hewel C, Rosenkranz D. unitas: The universal tool for annotation of small RNAs. *BMC Genomics* 2017;**18**:1–14. <https://doi.org/10.1186/s12864-017-4031-9>.
- 95 O'Neill K, Liao WW, Patel A, Hammell M. TEs_{small} identifies small RNAs associated with targeted inhibitor resistance in melanoma. *Front Genet* 2018. <https://doi.org/doi:10.3389/fgene.2018.00461>.
- 96 Hadi LHA, Lin QXXL, Minh TT, Loh M, Ng HK, Salim A, *et al.* miREM: an expectation-maximization approach for prioritizing miRNAs associated with gene-set. *BMC Bioinformatics* 2018;**19**:1–8. <https://doi.org/10.1186/s12859-018-2292-1>.
- 97 Bousios A, Gaut BS, Darzentas N. Considerations and complications of mapping small RNA high-throughput data to transposable elements. *Mob DNA* 2017;**8**:1–13. <https://doi.org/10.1186/s13100-017-0086-z>.
- 98 Zhang T, Cooper S, Brockdorff N. The interplay of histone modifications – writers that read. *EMBO Rep* 2015;**16**:1467–81.
- 99 Cao Y, Chen G, Wu G, Zhang X, McDermott J, Chen X, *et al.* Widespread roles of enhancer-like transposable elements in cell identity and long-range genomic interactions. *Genome Res* 2019;**29**:1–13. <https://doi.org/10.1101/gr.235747.118>.
- 100 Attig J, Agostini F, Gooding C, Chakrabarti AM, Singh A, Haberman N, *et al.* Heteromeric RNP Assembly at LINEs Controls Lineage-Specific RNA Processing. *Cell* 2018;**174**:1067–81. <https://doi.org/10.1016/j.cell.2018.07.001>.
- 101 Kelley DR, Hendrickson DG, Tenen D, Rinn JL. Transposable elements modulate human RNA abundance and splicing via specific RNA-protein interactions. *Genome Biol* 2014;**15**:1–16.
- 102 Chung D, Kuan PF, Li B, Sanalkumar R, Liang K, Bresnick EH, *et al.* Discovering Transcription Factor Binding Sites in Highly Repetitive Regions of Genomes with Multi-Read Analysis of ChIP-Seq Data. *PLoS Comput Biol* 2011;**7**:. <https://doi.org/10.1371/journal.pcbi.1002111>.
- 103 Wang R, Hsu H, Blattler A, Wang Y, Lan X, Wang Y, *et al.* LOcating Non-Unique matched Tags (LONUT) to Improve the Detection of the Enriched Regions for ChIP-seq Data. *PLoS One* 2013;**8**:1–10. <https://doi.org/10.1371/journal.pone.0067788>.
- 104 Sun G, Chung D, Liang K, Keles S. Statistical Analysis of ChIP-seq Data with MOSAiCS. In: Shomron N, editor. *Deep Seq. Data Anal.* New York, NY: Springer

- Science+Business Media; 2013. p. 193–212.
- 105 Nakato R, Itoh T, Shirahige K. DROMPA: easy-to-handle peak calling and visualization software for the computational analysis and validation of ChIP-seq data. *Genes to Cells* 2013;**18**:589–601. <https://doi.org/10.1111/gtc.12058>.
- 106 Berger S, Pachkov M, Arnold P, Omid S, Kelley N, Salatino S, *et al.* Crunch: Integrated processing and modeling of ChIP-seq data in terms of regulatory motifs. *Genome Res* 2019:1–24. <https://doi.org/10.1101/gr.239319.118>.
- 107 Sun X, Wang X, Tang Z, Grivainis M, Kahler D, Yun C, *et al.* Transcription factor profiling reveals molecular choreography and key regulators of human retrotransposon expression. *Proc Natl Acad Sci* 2018;**115**:. <https://doi.org/10.1073/pnas.1722565115>.
- 108 Zeng X, Li B, Welch R, Rojo C, Zheng Y, Dewey CN, *et al.* Perm-seq : Mapping Protein-DNA Interactions in Segmental Duplication and Highly Repetitive Regions of Genomes with Prior-Enhanced Read Mapping. *PLoS Comput Biol* 2015;**11**:1–23. <https://doi.org/10.1371/journal.pcbi.1004491>.
- 109 Lovci MT, Ghanem D, Marr H, Arnold J, Gee S, Parra M, *et al.* Rbfox proteins regulate alternative mRNA splicing through evolutionarily conserved RNA bridges. *Nat Struct Mol Biol* 2013;**20**:1434–42. <https://doi.org/10.1038/nsmb.2699>.
- 110 Maragkakis M, Alexiou P, Nakaya T, Mourelatos Z. CLIPSeqTools — a novel bioinformatics CLIP-seq analysis suite. *RNA* 2016;**22**:1–9. <https://doi.org/10.1261/rna.052167.115>.
- 111 Zhang Z, Xing Y. CLIP-seq analysis of multi-mapped reads discovers novel functional RNA regulatory sites in the human transcriptome. *Nucleic Acids Res* 2017;**45**:9260–71. <https://doi.org/10.1093/nar/gkx646>.
- 112 Li B, Tambe A, Aviran S, Pachter L. PROBer Provides a General Toolkit for Analyzing Sequencing-Based Toeprinting Assays. *Cell Syst* 2017;**4**:568–74. <https://doi.org/10.1016/j.cels.2017.04.007>.
- 113 Deniz Ö, Frost JM, Branco MR. Regulation of transposable elements by DNA modifications. *Nat Rev Genet* 2019;**20**:. <https://doi.org/10.1038/s41576-019-0106-6>.
- 114 Harris RA, Wang T, Coarfa C, Nagarajan RP, Hong C, Downey SL, *et al.* Comparison of sequencing-based methods to profile DNA methylation and identification of monoallelic epigenetic modifications. *Nat Biotechnol* 2010;**28**:. <https://doi.org/10.1038/nbt.1682>.
- 115 Cokus SJ, Feng S, Zhang X, Chen Z, Merriman B, Haudenschild CD, *et al.* Shotgun bisulphite sequencing of the Arabidopsis genome reveals DNA methylation patterning. *Nature* 2008;**452**:215–9. <https://doi.org/10.1038/nature06745>.
- 116 Lister R, Malley RCO, Tonti-Filippini J, Gregory BD, Berry CC, Millar AH, *et al.* Highly Integrated Single-Base Resolution Maps of the Epigenome in Arabidopsis. *Cell* 2008;**133**:523–36. <https://doi.org/10.1016/j.cell.2008.03.029>.
- 117 Lister R, Pelizzola M, Downen RH, Hawkins RD, Hon G, Tonti-Filippini J, *et al.* Human DNA methylomes at base resolution show widespread epigenomic differences. *Nature* 2009;**462**:315–22. <https://doi.org/10.1038/nature08514>.
- 118 Xi Y, Li W. BSMAP : whole genome bisulfite sequence MAPping program. *BMC Bioinformatics* 2009;**10**:1–9. <https://doi.org/10.1186/1471-2105-10-232>.

- 119 Krueger F, Andrews SR. Bismark : a flexible aligner and methylation caller for Bisulfite-Seq applications. *Bioinformatics* 2011;**27**:1571–2. <https://doi.org/10.1093/bioinformatics/btr167>.
- 120 Sun D, Xi Y, Rodriguez B, Park HJ, Tong P, Meong M, *et al.* MOABS : model based analysis of bisulfite sequencing data. *Genome Biol* 2014;**15**:1–12. <https://doi.org/10.1186/gb-2014-15-2-r38>.
- 121 Huang KYY, Huang Y-J, Chen P-Y. BS-Seeker3 : ultrafast pipeline for bisulfite sequencing. *BMC Bioinformatics* 2018;**19**:2–5. <https://doi.org/10.1186/s12859-018-2120-7>.
- 122 Adusumalli S, Feroz Mohd Omar M, Soong R, Benoukraf T. Methodological aspects of whole-genome bisulfite sequencing analysis. *Brief Bioinform* 2014;**16**:369–79. <https://doi.org/10.1093/bib/bbu016>.
- 123 Graña O, López-Fernández H, Fdez-Riverola F, Pisano DG, Glez-Peña D. Bicycle : a bioinformatics pipeline to analyze bisulfite sequencing data. *Bioinformatics* 2018;**34**:1414–5. <https://doi.org/10.1093/bioinformatics/btx778>.
- 124 Stuart T, Eichten SR, Cahn J, Karpievitch Y V, Borevitz JO, Lister R. Population scale mapping of transposable element diversity reveals links to gene regulation and epigenomic variation. *Elife* 2016;**5**:1–27. <https://doi.org/10.7554/eLife.20777>.
- 125 Daron J, Slotkin RK. EpiTEome : Simultaneous detection of transposable element insertion sites and their DNA methylation levels. *Genome Biol* 2017;**18**:1–10. <https://doi.org/10.1186/s13059-017-1232-0>.
- 126 Shukla R, Upton KR, Muñoz-Lopez M, Gerhardt DJ, Fisher ME, Nguyen T, *et al.* Endogenous Retrotransposition Activates Oncogenic Pathways in Hepatocellular Carcinoma. *Cell* 2013;**153**:101–11. <https://doi.org/10.1016/j.cell.2013.02.032>.
- 127 Noshay JM, Anderson SN, Zhou P, Ji L, Ricci W, Lu Z, *et al.* Monitoring the interplay between transposable element families and DNA methylation in maize. *PLoS Genet* 2019:1–25. <https://doi.org/10.1371/journal.pgen.1008291>.
- 128 Tanay A, Regev A. Scaling single-cell genomics from phenomenology to mechanism. *Nature* 2017;**541**:331–8. <https://doi.org/10.1038/nature21350>.
- 129 Göke J, Lu X, Chan Y, Ng H, Ly L-H, Sachs F, *et al.* Dynamic Transcription of Distinct Classes of Endogenous Retroviral Elements Marks Specific Populations of Early Human Embryonic Cells. *Cell Stem Cell* 2015;**16**:135–41. <https://doi.org/10.1016/j.stem.2015.01.005>.
- 130 Boroviak T, Stirparo GG, Dietmann S, Hernando-Herraez I, Mohammed H, Reik W, *et al.* Single cell transcriptome analysis of human, marmoset and mouse embryos reveals common and divergent features of preimplantation development. *Development* 2018;**145**:1–18. <https://doi.org/10.1242/dev.167833>.
- 131 Ziegenhain C, Vieth B, Parekh S, Reinius B, Guillaumet-Adkins A, Smets M, *et al.* Comparative Analysis of Single-Cell RNA Sequencing Methods. *Mol Cell* 2017;**65**:631–43. <https://doi.org/10.1016/j.molcel.2017.01.023>.
- 132 Guo J, Grow EJ, Mlcochova H, Maher GJ, Lindskog C, Nie X, *et al.* The adult human testis transcriptional cell atlas. *Cell Res* 2018. <https://doi.org/10.1038/s41422-018-0099-2>.
- 133 Jaitin DA, Kenigsberg E, Keren-Shaul H, Elefant N, Franziska P, Zaretsky I, *et al.* Massively Parallel Single-Cell RNA-Seq for Marker-Free Decomposition of Tissues into Cell Types. *Science (80-)* 2014;**343**:776–9.

- <https://doi.org/10.1126/science.1247651>.
- 134 Macosko EZ, Basu A, Satija R, Nemesh J, Shekhar K, Goldman M, *et al*. Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell* 2015;**161**:1202–14. <https://doi.org/10.1016/j.cell.2015.05.002>.
- 135 Hashimshony T, Senderovich N, Avital G, Klochendler A, Leeuw Y De, Anavy L, *et al*. CEL-Seq2 : sensitive highly-multiplexed single-cell RNA-seq. *Genome Biol* 2016;**17**:1–7. <https://doi.org/10.1186/s13059-016-0938-8>.
- 136 Nakamura T, Yabuta Y, Okamoto I, Aramaki S, Yokobayashi S, Kurimoto K, *et al*. SC3-seq : a method for highly parallel and quantitative measurement of single-cell gene expression. *Nucleic Acids Res* 2015;**43**:1–17. <https://doi.org/10.1093/nar/gkv134>.
- 137 Cao J, Packer JS, Ramani V, Cusanovich DA, Huynh C, Daza R, *et al*. Comprehensive single-cell transcriptional profiling of a multicellular organism. *Science (80-)* 2017;**357**:661–7. <https://doi.org/10.1126/science.aam8940>.
- 138 Manno G La, Soldatov R, Zeisel A, Braun E, Hochgerner H, Petukhov V, *et al*. RNA velocity of single cells. *Nature* 2018;**560**:494–8. <https://doi.org/10.1038/s41586-018-0414-6>.
- 139 Ding J, Adiconis X, Simmons SK, Kowalczyk MS, Hession CC, Marjanovic ND, *et al*. Systematic comparative analysis of single cell RNA-sequencing methods. *bioRxiv* 2019.
- 140 Islam S, Kjällquist U, Moliner A, Zajac P, Fan J, Lönnerberg P, *et al*. Highly multiplexed and strand-specific single-cell RNA 5' end sequencing. *Nat Protoc* 2012;**7**:813–28. <https://doi.org/10.1038/nprot.2012.022>.
- 141 Cole C, Byrne A, Beaudin AE, Forsberg EC, Vollmers C. Tn5Prime , a Tn5 based 5' capture method for single cell RNA-seq. *Nucleic Acids Res* 2018;**46**:1–12. <https://doi.org/10.1093/nar/gky182>.
- 142 Karlsson K, Lönnerberg P, Linnarsson S. Alternative TSSs are co-regulated in single cells in the mouse brain. *Mol Syst Biol* 2017;**13**:1–10. <https://doi.org/10.15252/msb.20167374>.
- 143 Brocks D, Chomsky E, Mukamel Z, Lifshitz A, Tanay A. Single cell analysis reveals dynamics of transposable element transcription following epigenetic de-repression. *bioRxiv* 2018.
- 144 Buenrostro J, Wu B, Chang H, Greenleaf W. ATAC-seq: A Method for Assaying Chromatin Accessibility Genome-Wide. *Curr Protoc Mol Biol* 2015;**109**:1–10. <https://doi.org/10.1002/0471142727.mb2129s109.ATAC-seq>.
- 145 Lieberman-Aiden E, Berkum NL Van, Williams L, Imakaev M, Ragoczy T, Telling A, *et al*. Comprehensive Mapping of Long-Range Interactions Reveals Folding Principles of the Human Genome. *Science (80-)* 2009;**326**:289–94.
- 146 Zheng Y, Ay F, Keles S. Generative modeling of multi-mapping reads with mHi-C advances analysis of Hi-C studies. *Elife* 2019:1–29.
- 147 Cao Y, Chen G, Wu G, Zhang X, McDermott J, Chen X, *et al*. Widespread roles of enhancer-like transposable elements in cell identity and long-range genomic interactions. *Genome Res* 2018:40–52. <https://doi.org/10.1101/gr.235747.118>.
- 148 Kruse K, Díaz N, Enriquez-Gasca R, Gaume X, Torres-Padilla M-E, Vaquerizas JM. Transposable elements drive reorganisation of 3D chromatin during early embryogenesis. *bioRxiv* 2019:1–28.

149 Rodriguez-Terrones D, Gaume X, Ishiuchi T, Weiss A, Kopp A, Kruse K, *et al.* A molecular roadmap for the emergence of early-embryonic-like cells in culture. *Nat Genet* 2018;**50**:. <https://doi.org/10.1038/s41588-017-0016-5>.

Sequencing Data format	Software Package	Includes Multimappers	Includes Iterative Statistics	TE Locus-Specific Calls	TE-centric	Publication
RNA	RSEM	y	y	n	n	Li et al. BMC Bioinformatics 2011
RNA	RepEnrich	y	n	n	y	Criscione et al. BMC Genomics 2014
RNA	TETranscripts	y	y	n	y	Jin et al. Bioinformatics 2015
RNA	TETools	y	n	n	y	Lerat et al. Nucleic Acids Research 2017
RNA	Yanagi	y	n	n	n	Gunady et al. bioRxiv 2018
RNA	SINEsFIND	n	n	y	SINEs	Carnevali et al. DNA Research 2017
RNA	SalmonTE	y	y	n	y	Jeong et al. PacBio Symposium on Biocomputing 2018
RNA	ERVmap	y	n	y	ERVs	Tokuyama et al. PNAS 2018
RNA	SQuIRE	y	y	y	y	Yang et al. Nucleic Acids Research 2019
RNA	LIONS	y	n	y	TE Fusions	Babian et al. Bioinformatics 2019
RNA	TeXP	y	n	n	LINE-1	Navarro et al. PLoS Comp Bio 2019
RNA	Telescope	y	y	y	y	Bendall et al. PLoS Comp Bio 2019
RNA	Scavenger	y	n	n	n	Yang et al. F1000Research 2019
sRNA	miRDeep2	y	n	n	n	Friedländer et al. Nucleic Acids Research 2012
sRNA	ShortStack	y	n	y	y	Axtell RNA 2013
sRNA	Butter	y	y	y	y	Axtell bioRxiv 2014
sRNA	PiPipes	y	n	n	y	Han et al. Bioinformatics 2015
sRNA	Chimira	y	n	n	n	Vitosos and Enright Bioinformatics 2015
sRNA	unitas	y	n	n	y	Gebert et al. BMC Genomics 2017
sRNA	Oasis 2	y	n	n	n	Rahman et al. BMC Bioinformatics 2018
sRNA	TEsmall	y	n	n	y	O'Neill et al. Frontiers in Genetics 2018
sRNA	Manatee	y	n	n	n	Handzlik et al. bioRxiv 2019
ChIP	CSEM	y	y	y	n	Chung et al. PLoS Comp Bio 2011
ChIP	MOSAICS	y	y	y	n	Kuan et al. J Am Stat Assoc 2011
ChIP	LONUT	y	n	y	n	Wang et al. PLoS One 2013
ChIP	DROMPA	y	n	y	n	Nakato et al. Genes to Cells 2013
ChIP	Perm-seq	y	y	y	n	Zeng et al. PLoS Comp Bio 2015
ChIP	MapRRCon	y	n	n	LINE1	Sun et al. PNAS 2018
ChIP	Crunch	y	n	y	n	Berger et al. Genome Research 2019
CLIP	CLIPSeqTools	y	n	n	n	Maragkakis et al. RNA 2016
CLIP	CLAM	y	y	n	n	Zhang and Xing Nucleic Acids Research 2017
CLIP	PROBer	y	y	y	n	Li et al. Cell Systems 2017
General Purpose Re-aligner	Gibbs Sampler	y	y	n	n	Wang et al. Bioinformatics 2010
General Purpose Re-aligner	MMR	y	y	n	n	Kahles et al. Bioinformatics 2016
General Purpose Re-aligner	CoCo	y	n	n	n	Deschamps-Francoeur et al. Bioinformatics 2019
Bisulfite	BSMAP	y	n	y	n	Xi et al. BMC Bioinformatics 2009
Bisulfite	Bismark	y	n	y	n	Krueger et al. Bioinformatics 2011
Bisulfite	MOABS	n	n	y	n	Sun et al. Genome Biology 2014
Bisulfite	BS-Seeker3	n	n	y	n	Huang et al. BMC Bioinformatics 2018
Bisulfite	bicycle	n	n	y	n	Graña et al. Bioinformatics 2018
Bisulfite	TEPID	n	n	y	y	Stuart et al. eLife 2016
Bisulfite	epiTEome	n	n	y	y	Daron et al. Genome Biology 2017
HiC	mHiC	y	n	y	y	Zheng et al. 2019 eLife