

# Telescope: Characterization of the retrotranscriptome by accurate estimation of transposable element expression

Matthew L. Bendall<sup>1,2</sup>, Miguel de Mulder<sup>2</sup>, Luis Pedro Iñiguez<sup>3</sup>, Aarón Lecanda-Sánchez<sup>3</sup>, Marcos Pérez-Losada<sup>1,4,5</sup>, Mario A. Ostrowski<sup>6,7</sup>, R. Brad Jones<sup>2</sup>, Lubbertus C. F. Mulder<sup>8,9</sup>, Gustavo Reyes-Terán<sup>3</sup>, Keith A. Crandall<sup>1,4</sup>, Christopher E. Ormsby<sup>3</sup> and Douglas F. Nixon<sup>2\*</sup>

<sup>1</sup> Computational Biology Institute, Milken Institute School of Public Health, George Washington University, Washington, D.C., USA.

<sup>2</sup> Division of Infectious Diseases, Department of Medicine, Weill Cornell Medical College, New York, N.Y., USA.

<sup>3</sup> Center for Research in Infectious Diseases (CIENI), Instituto Nacional de Enfermedades Respiratorias, Mexico City, Mexico.

<sup>4</sup> Department of Epidemiology and Biostatistics, Milken Institute School of Public Health, George Washington University, Washington, D.C., USA

<sup>5</sup> CIBIO-InBIO, Centro de Investigação em Biodiversidade e Recursos Genéticos, Universidade do Porto, Campus Agrário de Vairão, Vairão 4485-661, Portugal.

<sup>6</sup> Department of Immunology, University of Toronto, Toronto, Ontario, Canada.

<sup>7</sup> Keenan Research Centre for Biomedical Science of St. Michael's Hospital, Toronto, Ontario, Canada.

<sup>8</sup> Department of Microbiology, Icahn School of Medicine at Mount Sinai, New York, New York, USA.

<sup>9</sup> The Global Health and Emerging Pathogens Institute, Icahn School of Medicine at Mount Sinai, New York, New York, USA.

\* Corresponding author

E-mail: [dnixon@med.cornell.edu](mailto:dnixon@med.cornell.edu) (DFN)

# 1    **Abstract**

2    Characterization of Human Endogenous Retrovirus (HERV) expression within the  
3    transcriptomic landscape using RNA-seq is complicated by uncertainty in fragment  
4    assignment because of sequence similarity. We present Telescope, a computational  
5    software tool that provides accurate estimation of transposable element expression  
6    (retrotranscriptome) resolved to specific genomic locations. Telescope directly addresses  
7    uncertainty in fragment assignment by reassigning ambiguously mapped fragments to the  
8    most probable source transcript as determined within a Bayesian statistical model. We  
9    demonstrate the utility of our approach through single locus analysis of HERV expression  
10    in 13 ENCODE cell types. When examined at this resolution, we find that the magnitude  
11    and breadth of the retrotranscriptome can be vastly different among cell types.  
12    Furthermore, our approach is robust to differences in sequencing technology, and  
13    demonstrates that the retrotranscriptome has potential to be used for cell type  
14    identification. Telescope performs highly accurate quantification of the  
15    retrotranscriptomic landscape in RNA-seq experiments, revealing a differential  
16    complexity in the transposable element biology of complex systems not previously  
17    observed. Telescope is available at [github.com/mlbendall/telescope](https://github.com/mlbendall/telescope).

18

# 1 **Author Summary**

2 Almost half of the human genome is composed of Transposable elements (TEs), but their  
3 contribution to the transcriptome, their cell-type specific expression patterns, and their  
4 role in disease remains poorly understood. Recent studies have found many elements to  
5 be actively expressed and involved in key cellular processes. For example, human  
6 endogenous retroviruses (HERVs) are reported to be involved in human embryonic stem  
7 cell differentiation. Discovering which exact HERVs are differentially expressed in  
8 RNA-seq data would be a major advance in understanding such processes. However,  
9 because HERVs have a high level of sequence similarity it is hard to identify which exact  
10 HERV is differentially expressed. To solve this problem, we developed a computer  
11 program which addressed uncertainty in fragment assignment by reassigning  
12 ambiguously mapped fragments to the most probable source transcript as determined  
13 within a Bayesian statistical model. We call this program, “Telescope”. We then used  
14 Telescope to identify HERV expression in 13 well-studied cell types from the ENCODE  
15 consortium and found that different cell types could be characterized by enrichment for  
16 different HERV families, and for locus specific expression. We also showed that  
17 Telescope performed better than other methods currently used to determine TE  
18 expression. The use of this computational tool to examine new and existing RNA-seq  
19 data sets may lead to new understanding of the roles of TEs in health and disease.

20

# 1 Introduction

2 Transposable elements (TEs) represent the largest class of biochemically  
3 functional DNA elements in mammalian genomes(Dunham et al. 2012; Kellis et al. 2014)  
4 comprising nearly 50% of the human genome. As many of these transcriptionally active  
5 elements originated as retroelements, we refer to the set of RNA molecules transcribed  
6 from these elements in a population of cells as the retrotranscriptome. The contribution of  
7 the retrotranscriptome to the total transcriptome, cell-type specific expression patterns,  
8 and the role of retroelement transcripts in disease remain poorly understood (Magiorkinis  
9 et al. 2013). Although most TEs are hypothesized to be transcriptionally silent (due to  
10 accumulated mutations), recent studies have found many elements to be actively  
11 expressed and involved in key cellular processes. For example, aberrant expression of  
12 LINE-1 (L1) elements, the most expansive group of TEs, has been implicated in the  
13 pathogenesis of cancer (Wang-Johanning et al. 2003; Tang et al. 2017; Rodić et al. 2015;  
14 Ardeljan et al. 2017), while human endogenous retroviruses (HERVs) are reported to be  
15 involved in human embryonic stem cell differentiation(Grow et al. 2015; Göke et al.  
16 2015) and in the pathogenesis of amyotrophic lateral sclerosis(Li et al. 2015). We, and  
17 others, have shown that HIV-1 infection increases HERV transcription(Garrison et al.  
18 2007; Jones et al. 2012; Ormsby et al. 2012; Contreras-Galindo et al. 2012; Gonzalez-  
19 Hernandez et al. 2014). These lines of evidence therefore indicate that TEs have  
20 important roles in the regulation of human health and disease.

21 The ability to observe and quantify TE expression, especially the specific  
22 genomic locations of active elements, is crucial for understanding the molecular basis  
23 underlying a wide range of conditions and diseases(Flockerzi et al. 2008). Traditional

1 techniques for interrogating the TE transcriptome include quantitative PCR (Muradrasoli  
2 et al. 2006; Rangwala et al. 2009) and RNA expression microarrays (Seifarth et al. 2003;  
3 Pérot et al. 2012; Gnanakkan et al. 2013; Young et al. 2014; Becker et al. 2017).  
4 However, these techniques are unable to discover elements not specifically targeted by  
5 the assay, and may fail to detect rare, previously unknown, or weakly expressed  
6 transcripts. High-throughput RNA sequencing (RNA-seq) promises to overcome many of  
7 these shortcomings, enabling highly sensitive detection of transcripts across a wide  
8 dynamic range. Mathematical and computational approaches for transcriptome  
9 quantification using RNA-seq are well established ((Mortazavi et al. 2008; Marioni et al.  
10 2008), see review (Garber et al. 2011)) and provide researchers with reproducible  
11 analytical pipelines (Trapnell et al. 2010, 2012). Such approaches are highly effective at  
12 quantifying transcripts when sequenced fragments can be uniquely aligned to the  
13 reference genome, since the original genomic template for each transcript can be  
14 unambiguously identified (Trapnell et al. 2013; Conesa et al. 2016). In contrast,  
15 sequencing fragments generated by TEs often have high scoring alignments to many  
16 genomic locations with similar sequences, leading to uncertainty in transcript count  
17 estimates. Approaches that fail to account for these uncertainties may incorrectly estimate  
18 TE abundance and falsely detect significant changes in expression (Trapnell et al. 2013).

19 A growing number of studies are using high-throughput sequencing to  
20 characterize the retrotranscriptome. Three general approaches are used to deal with  
21 challenges of aligning short sequencing reads to repetitive elements. i) “Family-level”  
22 approaches combine read counts across all instances of a TE family, since fragments  
23 mapping to multiple genomic locations can often be uniquely assigned to a single repeat

1 family. ii) “Heuristic” approaches simplify the problem of multi-mapped fragments by  
 2 either discarding ambiguous reads (unique counts) or randomly assigning ambiguous  
 3 reads to one of its best scoring alignments (best counts). Finally, iii) “statistical”  
 4 approaches estimate the most probable assignment of fragments given a statistical model.  
 5 Our approach, Telescope, implements a Bayesian statistical model for reassigning  
 6 ambiguous fragments; previous work that has used statistical approaches include the  
 7 TETranscripts package (Jin et al. 2015) and an ad hoc model implemented by (Santoni et  
 8 al. 2012).

9 Here, we introduce Telescope, a tool which provides accurate estimation of TE  
 10 expression resolved to specific genomic locations. Our approach directly addresses  
 11 uncertainty in fragment assignment by reassigning ambiguously mapped fragments to the  
 12 most probable source transcript as determined within a Bayesian statistical model. We  
 13 implement our approach using a descriptive statistical model of the RNA-seq process and  
 14 use an iterative algorithm to optimize model parameters. We use Telescope to investigate  
 15 the expression of HERVs in cell types from the ENCODE consortium.

# 1   **Results**

2

## 3   **Telescope: Single locus resolution of transposable element expression**

4            Resolution of transposable element (including those of human endogenous  
5   retroviruses, HERVs) expression from RNA-seq data sets has been complicated by the  
6   many similarities of these repetitive elements. Telescope is a computational pipeline  
7   program that solves the problem of ambiguously aligned fragments by assigning each  
8   sequenced fragment to its most likely transcript of origin. We assume that the number of  
9   fragments generated by a transcript is proportional to the amount of transcript present in  
10   the sample; thus, the most likely source template for a randomly selected fragment is a  
11   function of its alignment uncertainty and the relative transcript abundances. Telescope  
12   describes this relationship using a Bayesian mixture model where the estimated  
13   parameters include the relative transcript abundances and the latent variables define the  
14   possible source templates for each fragment (Francis et al. 2013).

15            The first step in this approach is to independently align each fragment to the  
16   reference genome; the alignment method should search for multiple valid alignments for  
17   each fragment and report all alignments that meet or exceed a minimum score threshold  
18   (Fig 1A). Next, alignments are tested for overlap with known TE transcripts; transcript  
19   assignments for each fragment are weighted by the score of the corresponding alignment  
20   (Fig 1B and 1C). In our test cases, we typically find that less than 50% of the fragments  
21   aligning to TEs can be uniquely assigned to a single genomic location and many  
22   fragments have more than 20 possible originating transcripts.

**Fig 1. Telescope conceptual overview.** A set of possible genomic locations for each fragment is determined by alignment to the reference genome. In order to find as many high-scoring mappings as possible, we use sensitive local alignment parameters and search for up to 100 alignments for each fragment using bowtie2 (A). Using an annotation containing known TE locations, Telescope intersects the aligned fragments with annotated TE loci (B,C). The set of alignments and corresponding alignment scores for each fragment are used to calculate the expected assignment weights, initially assuming equal expression for all elements (D). The assignment weights estimated in (D) are used to find the maximum likelihood estimate (MLE) for the proportion of each transcript (E). Next, we update the expected assignment weights, now assuming that the MLE represents our best estimate of transcript expression (D,E). The steps in panels (D) and (E) describe an expectation-maximization procedure, and we further refine the assignment weights and MLE by iterating until parameter estimates converge. Telescope produces a report that includes the maximum a posteriori estimate of the transcript proportions and the final number of fragments assigned to each transcript, as well as an updated alignment including the final fragment assignments (F).

Telescope estimates the transcript proportions and expected source templates using an expectation-maximization algorithm. In the expectation step (E-step), the expected value of the source template for each fragment is calculated under current estimates of transcript abundance (Fig 1D). The maximization step (M-step) finds maximum *a posteriori* estimates of the transcript abundance dependent on the expected values from the E-step (Fig 1E). These steps are repeated until parameter estimates



1 converge (Fig 1D and 1E). Telescope reports the proportion of fragments generated by  
2 each transcript and the expected transcript of origin for each fragment (Fig 4F). The final  
3 counts estimated by Telescope correspond to actual observations of sequenced fragments  
4 and are suitable for normalization and differential analysis by a variety of methods. The  
5 software also provides an updated alignment with final fragment assignments that can be  
6 examined using common genome visualization tools. Telescope is available at  
7 [github.com/mlbendall/telescope](https://github.com/mlbendall/telescope).

8

## 9 **Determination of HERV expression in major cell types from the ENCODE** 10 **consortium**

11 To investigate HERV expression in a robust way across a diverse platform of cell  
12 types we relied on publicly available RNA-seq data. The ENCODE data project is an  
13 invaluable source of genomic data from disparate sources and provides the opportunity to  
14 mine the transposable element expression in a setting of maximum genomic information.  
15 We profiled 13 human cell types, including common lines designated by the ENCODE  
16 consortium, as well as primary cell types, and applied our approach to determine HERV  
17 expression across the spectrum of human cell types, including normal or transformed, and  
18 contrasting cell lines with primary cells (Table 1).

19 **Table 1. ENCODE cell types used in this study**

Cell Type	Description	Karyotype	Lineage	Tissue	Replicates
H1-hESC	Embryonic stem cell	Normal	ICM	ESC	4
GM12878	B-lymphocyte	Normal	mesoderm	blood	4
K562	Leukemia	Cancer	mesoderm	blood	3
HeLa-S3	Cervical carcinoma	Cancer	ectoderm	cervix	3
HepG2	Hepatocellular carcinoma	Cancer	endoderm	liver	3
HUVEC	Umbilical vein endothelial cells	Normal	mesoderm	vessel	3

SK-N-SH	Neuroblastoma	Cancer	ectoderm	brain	1
IMR90	Fetal lung fibroblasts	Normal	endoderm	lung	1
A549	Lung carcinoma	Cancer	endoderm	lung	1
MCF-7	Mammary gland adenocarcinoma	Cancer	ectoderm	breast	2
CD20+	CD20+ B cells	Normal	mesoderm	blood	1
CD14+	CD14+ Monocytes	Normal	mesoderm	blood	1
NHEK	Epidermal keratinocytes	Normal	ectoderm	skin	3

1

2 Over 2.7 billion sequenced fragments aligned to human reference hg38 with  
3 between 23.6% and 46.1% of the fragments in each sample aligning ambiguously to  
4 multiple genomic locations. Telescope intersected the aligned fragments with a set of  
5 14,968 manually curated HERV loci belonging to 60 families (see methods) and  
6 identified over 27 million fragments that appear to originate from HERV proviruses.  
7 Most (80.1%) of these fragments aligned to multiple genomic locations; we used  
8 Telescope to reassign ambiguous fragments to the most likely transcript of origin and  
9 estimate expression at specific HERV loci.

10 We developed genome-wide maps of HERV expression for 8 of the analyzed cell  
11 types that had replicates (Table 1) , and used CIRCOS (Krzywinski 2009) to visualize the  
12 data (Fig 2). The outer track is a bar chart showing the number of HERV loci in 10 Mbp  
13 windows, with the red part of the bar representing the number of loci that are expressed  
14 in one or more cell types. The 8 inner rings show the expression levels (log2 counts per  
15 million (CPM)) of 1365 HERV loci that were expressed at least one of the cell types  
16 examined. Moving from the outer ring to the inner ring are replicates for each of the 8  
17 cell types with duplicates: H1-hESC, GM12878, K562, HeLa-S3, HepG2, HUVEC,  
18 MCF-7, and NHEK.

19

# **Fig 2. Genome-wide maps of locus-specific HERV expression for 8 ENCODE tier 1**

**and 2 cell types.** The outer track is a bar chart showing the number of HERV loci in 10 Mbp windows, with the red part of the bar representing the number of loci that are expressed in one or more cell types. The 8 inner rings show the expression levels (log2 counts per million (CPM)) of 1365 HERV loci that were expressed at least one of the cell types examined. Moving from the outer ring to the inner ring are replicates for each of the 8 cell types with duplicates: H1-hESC, GM12878, K562, HeLa-S3, HepG2, HUVEC, MCF-7, and NHEK.

We found 1365 HERV loci that were expressed in at least one of the cell types (CPM > 0.5). Not all HERVs were expressed in all cell types, some were widely expressed in all cells, whereas others were only expressed in one or more cell type (Fig 2). There is also a spectrum of differential HERV expression, with some HERVs having significantly higher expression than others. On a chromosome by chromosome analysis, there are certain regions of the genome that have minimal HERV expression, while other regions appear dense in HERV expression. There are areas of scarce HERV expression on chromosomes 3, 5, 9, 15, and the Y chromosome (Fig 2). Interestingly, the Y chromosome is host to a greater density of HERV locations, yet they are mostly silent. In contrast, several chromosomes exhibit a greater than expected number of active HERV locations, i.e. chromosome 19 (S1 Fig) and chromosome 6 (S2 Fig).

## **HERV Locus-specific-analysis**

To ascertain, global, family and locus level specific HERV expression, we assessed the number of HERVs expressed in each cell type. All cell types expressed HERVs; the number of expressed loci ranged from 216 (in MCF-7), to 533 (H1-hESC) (Fig 3A). The number and proportion of cell type specific locations (expressed in only one cell) differed among cell types. Nearly half (46.3%) of locations expressed in H1-hESC were not expressed in any other cell type, while 89.3% of locations expressed in MCF-7 were also present in other cell types (Fig 3A). This suggests that regulatory networks are shared among some cell types but not others. We next examined the relative contribution of HERV families to overall HERV transcription and found that different cell types could be characterized by enrichment for different HERV families. For example, HERVH accounted for 91.8% of the transcriptomic output in H1-hESC cells, while HERVE was dominant in K562 cells (24.4%) (Fig 4A). Other families, such as HERVL, were evenly distributed across cell types, both in number of expressed locations and in expression levels (Fig 4B). Resolving the most highly expressed specific locations in each cell type at a locus specific level shows that the distribution of expression varies among cell types. (Fig 3C). For example, HepG2 is characterized by high expression from a single locus, while H1-hESC has many locations that are activated.

**Fig 3. Overall HERV expression patterns.** (A) Number of HERV elements that are expressed for each cell type; expressed loci have CPM > 0.5 in the majority of replicates. The darker section of the bar corresponds to expressed loci that are unique to cell type, while the lighter part is expressed in other cell types. (B) The proportion of mapped RNA-seq fragments that are generated from HERV transcripts in each of eight replicated

1 cell types. Each point is one replicate; boxplot shows the median and first and third  
2 quartiles. (C) Top 10 most highly expressed loci for each cell type. Height of the bar is  
3 average CPM of all replicates with error bars representing the standard error calculated  
4 from replicates CPM values.

5  
6 **Fig 4. Family-level HERV expression profiles using Telescope.** Family-level HERV  
7 expression profiles were computed from locus-specific profiles (generated by Telescope)  
8 by summing expression across all locations within each family. (A) The proportion of  
9 fragments assigned to each HERV family relative to the total amount of HERV  
10 expression. Families that account for at least 5% of total HERV expression in at least one  
11 cell type are shown, with the remaining families in “other”. (B) Number of expressed  
12 HERV loci and fragment counts per million mapped fragments (CPM) for selected  
13 HERV families.

# 14 15 **HERV expression profiles generated by Telescope are cell type specific**

16 Previous work has suggested that estimates of HERV expression are highly  
17 sensitive to sequencing technology used, and differences due to sequencing technology  
18 can obscure biological differences due to cell type (Haase et al. 2015). Since aligning  
19 shorter fragments (i.e. single-end reads) tends to produce more ambiguously mapping  
20 fragments compared to longer fragments, we hypothesized that Telescope (which  
21 resolves ambiguity) would create HERV expression profiles that are robust to differences  
22 in sequencing technology. Hierarchical clustering of all 30 polyA RNA-seq HERV  
23 profiles shows that replicates from the same cell type cluster most closely with other

1 samples from the same cell type, regardless of the sequencing technology used (Fig 5A).  
 2 Clusters for all cell types had significant support using multiscale bootstrap resampling  
 3 (approximately unbiased (AU) > 95%). Principal component analysis (PCA) also  
 4 indicates that cell type, not sequencing technology, is associated with the strongest  
 5 differences among expression profiles. The first principal component, accounting for  
 6 44% of the total variance in the data, separates H1-hESC samples from all other samples  
 7 (Fig 5B). The second and third components further separate the samples into the other 12  
 8 cell types, and capture 13% and 10% of the total variance, respectively. Interestingly, the  
 9 second component separates blood-derived cell types (K562, GM12878, CD20+ and  
 10 CD14+) from the other cell types, suggesting that cells derived from the same tissue may  
 11 share similarities in HERV expression profiles.

12

13 **Fig 5. Cell type characterization based on HERV expression profiles using**  
 14 **unsupervised learning and linear models.** Unsupervised learning and linear modeling  
 15 were used to identify patterns in HERV expression profiles generated by Telescope for  
 16 30 polyA RNA-seq datasets from 13 cell types. (A) Similarities among normalized  
 17 expression profiles were explored using hierarchical cluster analysis. Supporting p-values  
 18 were based on 1000 multiscale bootstrap replicates and calculated using Approximately  
 19 Unbiased (AU, red) and Bootstrap probability (BP, green) approaches. (B) Principal  
 20 component analysis (PCA) of normalized expression profiles. The first component  
 21 accounts for 44% of the variance in the data, and is plotted against component 2 and 3,  
 22 which account for 13% and 10% of the variance, respectively. (C) Heatmap of the  
 23 number of HERV elements found to be significantly differentially expressed (DE) among

each pair of cell types. Significance was determined using cutoffs for the false discovery rate ( $FDR < 0.1$ ) and log2 fold change ( $abs(LFC) > 1.0$ ). Yellow indicates low numbers of differentially expressed elements, while blue indicates high numbers.

We further explored differences among cell types using differential expression (DE) analysis. Pairwise contrasts between cell types were performed to determine the number of significant DE loci ( $FDR < 0.1$ ,  $abs(LFC) > 1.0$ ) (Fig 5C). As found in the unsupervised analysis, HERV expression in H1-hESC was drastically different from other cell types, with between 578 and 1127 significantly DE loci.

Finally, we asked whether heuristic approaches for TE quantification would be sufficient to identify cell type specific signal in the data or whether these approaches would be sensitive to other variables. We performed the same unsupervised analyses with HERV expression profiles obtained using unique and best counting approaches. Hierarchical clustering using unique count expression profiles produced a very similar topology to that found using Telescope (S3 Fig). Replicate samples were properly clustered by cell type, though GM12878 and HUVEC had slightly less bootstrap support. In contrast, clustering with the best count profiles did not recover all cell type clusters; two HeLa-S3 samples clustered with H1-hESC, while the third was more similar to A549 cells (S3 Fig). There was also less support for several clusters, including one cell type cluster (NHEK) that did not meet the 95% threshold.

## Performance of Telescope compared to current methods

In order to examine the sensitivity and biases of computational approaches for quantifying TE expression, we designed simulation experiments with known expression values. Earlier studies have suggested that the HERV-K(HML-2) family (hereafter referred to as HML-2) is expressed in human tissue and may be relevant to human health (Hohn et al. 2013; Grow et al. 2015; Li et al. 2015; Weiss 2016). Furthermore, its relatively few family members (~90 distinct genomic loci (Subramanian et al. 2011)) and high nucleotide identity make HML-2 a good model for studying TE expression. Here, we report on the performance of each method to detect locus-specific expression of HML-2 by simulating RNA-seq fragments. We simulated 25 independent datasets, each simulation consisted of 10 randomly chosen HML-2 loci and a fragment count, which could be interpreted as an expression value. We used the following TE quantification approaches for estimating locus specific HML-2 expression: unique counts, best counts, RepEnrich (Criscione et al. 2014), TETranscripts (Jin et al. 2015) and Telescope.

**Fig 6. Comparison of performance results for unique counts, best counts, RepEnrich, TETranscripts and Telescope using simulated data.** 25 RNA-seq samples were simulated, each sample consisted of 10 randomly chosen HML-2 loci with each having different expression value. The possible expression values are shown along the x-axis, 0 represent all HML-2 not expressed, red dashed line represents the expected expression value. A box plot representing the count distribution from each expression value is plotted. The resulted count from the different counting method from each simulated expression value per sample is plotted over the boxplot. Counting methods tested: (A) unique count, (B) best count, (C) RepEnrich, (D) TETranscript, (E) Telescope.



1 (F) The precision and recall for each sample simulated as well as the mean of both are  
2 shown for all methods.

3

4 The greatest strength of the unique counts approach was the low false detection  
5 rate since across all 25 simulations (41.5K simulated fragments), only 6 fragments were  
6 incorrectly assigned. However, unique counts consistently underestimated expression  
7 levels with ~60% of all estimates (151 out of 250) missing at least 50% of the true  
8 expression (Fig 6A). One striking example of this underestimation was for  
9 HML2\_5q33.3; this locus did not generate any fragment that could be counted by unique  
10 counts despite being expressed in 5 simulations. We presume that the underestimation is  
11 a direct consequence of discarding ambiguously mapped reads, as the unique counts  
12 discarded 62.6% of the simulated fragments.

13 In contrast to unique counts, the best counts approach offers greater sensitivity  
14 (Fig 5F). Instead of discarding ambiguously mapped fragments, all the fragments are  
15 used, resulting in more accurate expression estimates (Fig 5B). The majority of fragments  
16 were assigned to the true source transcript, while incorrect hits account for 14.2% of the  
17 total fragments. These off-target assignments resulted in false detection of unexpressed  
18 loci in each simulation, representing a major drawback of this approach (Fig 5B). Despite  
19 the high sensitivity of best counts, we conclude that the high number of incorrect  
20 detections outweighs the possible advantages of using this approach.

21 In order to make a direct comparison between Telescope and RepEnrich, which  
22 quantifies TE families at the family level, we modified RepEnrich annotations to give  
23 each locus a unique “family” name. We found that expression was underestimated by this

1 approach for all expressed HML-2 elements (Fig 5C). We attribute this bias to the large  
2 number of fragments that were discarded by this approach, 38.8%. Nevertheless, the  
3 number of discarded reads is fewer than the unique counts approach where 62.6% of all  
4 fragments were omitted. Another negative aspect of RepEnrich is that 12.1 % of the  
5 counted fragments were assigned to non-expressed TE.

6 We tested TETranscripts on our simulated data with the TE counting method set  
7 under “multi” mode. We provided the algorithm with our HERV annotation thus results  
8 could be compared with Telescope. TETranscripts showed a better performance than  
9 RepEnrich as shown in Jin *et al.* (2015) involving multiple mapped fragments assignment  
10 but 21.7% of all fragments were assigned incorrectly to a non-expressed HML-2. Based  
11 on the precision and recall of all methods tested TETranscripts performed as the third  
12 best TE single locus counting method on the simulated data (S4 Fig).

13 Finally, we tested Telescope’s ability to reassign ambiguous alignments and  
14 estimate locus-specific fragment counts using the same simulation data. On average  
15 57.8% of the simulated fragments aligned to multiple loci, which need to be reassigned.  
16 Telescope reassigned ambiguously mapped fragments to the expected transcript of origin  
17 according to our model (see methods) and reported the final number of fragments aligned  
18 to each TE. The estimated levels of expression, calculated with Telescope, from each  
19 HML-2 resembled closely the simulated levels (Fig 5E). Only 3 HML-2 loci that were  
20 simulated to be expressed did not present any fragments and 99.9% of all simulated  
21 fragments were counted with the Telescope approach.

22 Of all methods considered here, Telescope had the highest rate of precision and  
23 recall from all other counting methods tested (Fig 5F). In contrast to the best counts

1 approach, the second best (S4 Fig), Telescope assigned only 15 fragments were assigned  
 2 to genomic locations that were not expressed, while 5871 fragments were assigned  
 3 incorrectly by best counts. Deviations of Telescope estimates from true expression levels,  
 4 as measured by F1-score, was the highest of all approaches (S4 Fig). These simulation  
 5 results demonstrate that Telescope resolves ambiguously aligned fragments and produces  
 6 unbiased estimates of TE expression that are robust to sequencing error.

# 1 DISCUSSION

2 High-throughput RNA sequencing has enabled the simultaneous characterization  
3 and quantification of an entire transcriptome with remarkable resolution and sensitivity.  
4 Current studies have primarily focused on protein-coding transcripts, with greater  
5 attention being given to non-coding and micro RNAs in recent years. The transposable  
6 elements represent another major biochemically active group of transcripts and are  
7 increasingly recognized as important regulators in complex biological systems and  
8 disease yet have been largely ignored in the literature. We present a novel software  
9 program, Telescope, that can be used to mine new or existing RNA-seq datasets to  
10 accurately quantify the expression of TEs. The key advantage of our approach is the  
11 capability to localize TE expression to an exact chromosomal location.

12 As TEs are repetitive elements located throughout the genome, existing programs  
13 have limitations in performing accurate alignments from RNA fragments because of  
14 sequence similarity. The management of alignment uncertainty has been approached in  
15 several ways. The unique count approach discards fragments that align ambiguously, but  
16 this approach underestimates or fails to detect gene expression. The best count method  
17 assigns each fragment to the source template with the best scoring alignment, but this  
18 underestimates TEs that are truly expressed and spuriously detects those that are in fact  
19 absent. While the family method of alignment mitigates uncertainty in alignments by  
20 classification according to repeat family, this method does not locate the genomic site of  
21 TE transcription. Our approach, Telescope, reassigns fragments to the most likely  
22 originating transcript using Bayesian mixture model that relating the relative transcript  
23 abundances to the possible source templates for each fragment. This approach thus

1 resolves ambiguously aligned fragments and results in accurate quantification of TE loci  
2 for differential analysis.

3 Telescope will have widespread utility in other settings. Studies on TE expression  
4 have become prominent in studies of embryonic stem cell development (Grow et al.  
5 2015)(Göke et al. 2015), neural cell plasticity (Muotri et al. 2010; Gage and Muotri  
6 2012), oncogenesis (Wang-Johanning et al. 2003; Rakoff-Nahoum et al. 2006; Takahashi  
7 et al. 2008; Tang et al. 2017; Rodić et al. 2015; Ardeljan et al. 2017), psychiatric and  
8 neurological disorders(Perron et al. 2012; Christensen 2016; Mortelmans et al. 2016) and  
9 autoimmune diseases (Nexø et al. 2015; Hanke et al. 2016). As the breadth of knowledge  
10 on TEs expands, expression profiling of TEs using Telescope will allow scientists to  
11 discover unique and collective TE transcripts involved in the biology of complex  
12 systems.

13

# 1    **Methods**

## 2    **Fragment reassignment mixture model**

3            Telescope implements a generative model of RNA-seq relating the probability of  
4    observing a sequenced fragment to the proportions of fragments originating from each  
5    transcript. Formally, let  $F = [f_1, f_2, \dots, f_N]$  be the set of  $N$  observed sequencing  
6    fragments. We assume these fragments originate from  $K$  annotated transcripts in the  
7    transcriptome  $T = [t_1, t_2, \dots, t_K]$ . In practice, annotations fail to identify all possible  
8    transcripts that generate fragments, thus we include an additional category,  $t_0$ , for  
9    fragments that cannot be assigned to annotated transcripts. Let  $G = [G_1, G_2, \dots, G_N]$   
10   represent the true generating transcripts for  $F$ , where  $G_i \in T$  and  $G_i = t_j$  if  $f_i$  originates  
11   from  $t_j$ . Since the process of generating  $F$  from  $T$  cannot be directly observed, the true  
12   generating transcripts  $G$  are considered to be “missing” data. The objective of our model  
13   is to estimate the proportions of  $T$  by learning the generating transcripts of  $F$ .

14    As described above, the alignment stage identifies one or more possible alignments for  
15    each fragment, along with corresponding alignment scores. Let  $q_i = [q_{i0}, q_{i1}, \dots, q_{iK}]$  be  
16    the set of mapping qualities for fragment  $f_i$ , where  $q_{ij} = \Pr(f_i | G_i = t_j)$  represents the  
17    conditional probability of observing  $f_i$  assuming it was generated from  $t_j$ ; we calculate  
18    this by scaling the raw alignment score by the maximum alignment score observed for the  
19    data. We write the likelihood of observing uniquely aligned fragment  $f_u$  as a function of  
20    the conditional probabilities  $q_u$  and the relative expression of each transcript for all  
21    possible generating transcripts  $G_u$

$$\Pr(f_u | \pi, q_u) = \sum_{j=0}^K \pi_j q_{uj}$$

1 where  $\boldsymbol{\pi} = [\pi_0, \pi_1, \dots, \pi_K]$  represents the fraction of observed fragments originating from  
2 each transcript. Note that  $q_{uj} = 0$  for all transcripts that are not aligned by  $f_u$ . For non-  
3 unique fragments, we introduce an additional parameter in the above likelihood to  
4 reweight each ambiguous alignment among the set of possible alignments. The  
5 probability of observing ambiguous fragment  $f_a$  is given by

$$6 \quad \Pr(f_a | \boldsymbol{\pi}, \boldsymbol{\theta}, \mathbf{q}_a) = \sum_{j=0}^K \pi_j \theta_j q_{aj}$$

7 where  $\boldsymbol{\theta} = [\theta_0, \theta_1, \dots, \theta_K]$  is a reassignment parameter representing the fraction of non-  
8 unique reads generated by each transcript.

9 Using these probabilities of observing ambiguous and unique fragments, we  
10 formulate a mixture model describing the likelihood of the data given parameters  $\boldsymbol{\pi}$  and  
11  $\boldsymbol{\theta}$ . The  $K$  mixture weights in the model are given by  $\boldsymbol{\pi}$ , the proportion of all fragments  
12 originating from each transcript. To account for uncertainty in the initial fragment  
13 assignments, let  $\mathbf{x}_i = [x_{i0}, x_{i1}, \dots, x_{iK}]$  be a set of partial assignment (or membership)  
14 weights for fragment  $f_i$ , where  $\sum_{j=0}^K x_{ij} = 1$  and  $x_{ij} = 0$  if  $f_i$  does not align to  $t_j$ . We  
15 assume that  $\mathbf{x}_i$  is distributed according to a multinomial distribution with success  
16 probability  $\boldsymbol{\pi}$ . Intuitively,  $x_{ij}$  represents our confidence that  $f_i$  was generated by  
17 transcript  $t_j$ . In order to simplify our notation, we introduce an indicator variable  $\mathbf{y} =$   
18  $[y_1, y_2, \dots, y_N]$  where  $y_i = 1$  if  $f_i$  is ambiguously aligned and  $y_i = 0$  otherwise. The  
19 complete data likelihood is

$$20 \quad L(\boldsymbol{\pi}, \boldsymbol{\theta} | \mathbf{x}, \mathbf{q}, \mathbf{y}) \propto \prod_{i=1}^N \prod_{j=0}^K [\pi_j \theta_j^{y_i} q_{ij}]^{x_{ij}}$$

21

# 1 **Parameter estimation and fragment reassignment by EM**

2 Telescope iteratively optimizes the likelihood function using an expectation-  
3 maximization algorithm (Dempster et al. 1977). First, the parameters  $\pi$  and  $\theta$  are  
4 initialized by assigning equal weight to all transcripts. In the expectation step, we  
5 compute the expected values of  $x_i$  under current estimates of the model parameters. The  
6 expectation is given by the posterior probability of  $x_i$ :

$$7 \quad E[x_{ij}] = \frac{\pi_j \theta_j^{y_i} q_{ij}}{\sum_{k=0}^K \pi_k \theta_k^{y_i} q_{ik}}$$

8 In the M-step we calculate the maximum a posteriori (MAP) estimates for  $\pi$  and  $\theta$

$$9 \quad \hat{\pi}_j = \frac{\sum_{i=1}^N E[x_{ij}] + a_j}{N + \sum_{k=0}^K a_k} \quad \text{and} \quad \hat{\theta}_j = \frac{\sum_{i=1}^N E[x_{ij}] y_i + b_j}{\sum_{i=1}^N y_i + \sum_{k=0}^K b_k}$$

10 where  $a_j$  and  $b_j$  are prior information for transcript  $t_j$ . Intuitively, these priors are  
11 equivalent to adding unique or ambiguous fragments to  $t_j$ ; providing non-zero values for  
12 these parameters prevents parameter estimates from converging to boundary values.  
13 Convergence of EM algorithms to local maxima has been shown by (Wu 1983), and is  
14 achieved when the absolute change in parameter estimates is less than a user defined  
15 level, typically  $\epsilon < 0.001$ .

# 17 **HERV Annotations**

18 A Telescope analysis requires an annotation that defines the transcriptional unit of each  
19 TE to be quantified. For HERV proviruses, the prototypical transcriptional unit contains  
20 an internal protein-coding region flanked by LTR regulatory regions. Existing  
21 annotations, such as those identified by RepeatMasker (Tarailo-Graovac and Chen 2009)  
22 (using the RepBase database (Jurka et al. 2005)) or Dfam (Wheeler et al. 2013) identify



1 sequence regions belonging to TE families but do not seek to annotate transcriptional  
2 units. Both databases represent the internal region and corresponding LTRs using  
3 separate models, and the regions identified are sometimes discontinuous. Thus, a HERV  
4 transcriptional unit is likely to appear as a collection of nearby annotations from the same  
5 HERV family.

6 Transcriptional units for HERV proviruses were defined by combining RepeatMasker  
7 annotations belonging to the same HERV family that are located in adjacent or nearby  
8 genomic regions. Briefly, repeat families belonging to the same HERV family (internal  
9 region plus flanking LTRs) were identified using the RepBase database (Jurka et al.  
10 2005). RepeatMasker annotations for each repeat family were downloaded using the  
11 UCSC table browser (Karolchik et al. 2004) and converted to GTF format, merging  
12 nearby annotations from the same repeat family. Next, LTR found flanking internal  
13 regions were identified and grouped using BEDtools (Quinlan and Hall 2010). HERV  
14 transcriptional units containing internal regions were assembled using custom python  
15 scripts. Each putative locus was categorized according to provirus organization; loci that  
16 did not conform to expected HERV organization or conflicted with other loci were  
17 visually inspected using IGV (Thorvaldsdóttir et al. 2013) and manually curated. As  
18 validation, we compared our annotations to the HERV-K(HML-2) annotations published  
19 by (Subramanian et al. 2011); the two annotations were concordant. Final annotations  
20 were output as GTF files and are available; all annotations, scripts, and supporting  
21 documentation are available at [https://github.com/mlbendall/telescope\\_annotation\\_db](https://github.com/mlbendall/telescope_annotation_db).

## 22 23 **Simulated HML-2 expression data**

1 We simulated 25 independent datasets, each consisted of randomly chosen 10  
2 HML-2 which were expressed at different level, ranging from 30 to 300 fragments per  
3 locus. Using the expression pattern and the chosen HML-2, we simulated sequencing  
4 fragments with the Bioconductor package for RNA-seq simulation, Polyester (Frazee et  
5 al. 2014). All simulations used the parameters of read length: 75 bp; average fragment  
6 size: 250; fragment size standard deviation: 25; and an Illumina error model with an error  
7 rate of 5e-3.

## 8

### 9 **Alignment to reference genome**

10 Sequenced fragments from each sample or simulation were aligned to human reference  
11 genome hg38 using bowtie2. Alignment options were specified to perform a sensitive  
12 local alignment search (--very-sensitive-local) with up to 100 alignments reported for  
13 each fragment pair (-k 100). The minimum alignment score threshold was chosen so that  
14 fragments with ~95% or greater sequence identity would be reported (--score-min  
15 L,0,1.6).

## 16

### 17 **Software Availability**

18 All scripts used for simulating and analyzing data are available at  
19 <https://github.com/mlbendall/TelescopeEncode>. The Telescope package is available at  
20 <https://github.com/mlbendall/telescope>.

21

22

# 1 **Author contributions**

2 M.L.B., L.P.I, K.A.C. A.L.S., M.P.-L. and C.E.O developed the mathematics and  
 3 statistics. L.P.I., D.G.H. and M.S. performed the experiments. M.M., and L.C.M designed  
 4 the experiments and performed the analysis. M.L.B. implemented the software. M.L.B.,  
 5 M.M.R., M.A.O, R.B.J., G.R-T, K.A.C. and D.F.N. conceived the research. All authors  
 6 wrote and approved the manuscript.

7

# 8 **Disclosure Declaration**

9 The authors declare no conflict of interest.

10

# 1    **Acknowledgments**

2    The work was supported in part by US National Institutes of Health grants: CA206488  
3    (DFN), AI076059 (DFN), UL1TR001876 (KAC), GM113886 (LCFM), and GM113886-  
4    01S1 (LCFM). MP-L was partially supported by DC CFAR pilot and CFAR  
5    1P30AI117970 awards. MLB is a predoctoral student in the Systems Biology Program of  
6    the Institute for Biomedical Sciences at the George Washington University. This work is  
7    from a dissertation to be presented to the above program in partial fulfillment of the  
8    requirements for the Ph.D. degree.

9

10   We thank Timothy Powell, Rodrigo Duarte and Deepak Srivastava for constructive  
11   reading of the manuscript.

12

# References

- Ardeljan D, Taylor MS, Ting DT, Burns KH. 2017. The Human Long Interspersed Element-1 Retrotransposon: An Emerging Biomarker of Neoplasia. *Clin Chem* **63**: 816–822.
- Becker J, Pérot P, Cheynet V, Oriol G, Mugnier N, Mommert M, Tabone O, Textoris J, Veyrieras J, Mallet F. 2017. A comprehensive hybridization model allows whole HERV transcriptome profiling using high density microarray. *BMC Genomics* **18**: 286.
- Christensen T. 2016. Human endogenous retroviruses in neurologic disease. *APMIS* **124**: 116–126.
- Conesa A, Madrigal P, Tarazona S, Gomez-Cabrero D, Cervera A, McPherson A, Szczesniak MW, Gaffney DJ, Elo LL, Zhang X, et al. 2016. A survey of best practices for RNA-seq data analysis. *Genome Biol* **17**: 13.
- Contreras-Galindo R, Kaplan MH, Contreras-Galindo AC, Gonzalez-Hernandez MJ, Ferlenghi I, Giusti F, Lorenzo E, Gitlin SD, Dosik MH, Yamamura Y, et al. 2012. Characterization of Human Endogenous Retroviral Elements in the Blood of HIV-1-Infected Individuals. *J Virol* **86**: 262–276.
- Criscione SW, Zhang Y, Thompson W, Sedivy JM, Neretti N. 2014. Transcriptional landscape of repetitive elements in normal and cancer human cells. *BMC Genomics* **15**: 583.
- Dempster A, Laird N, Rubin D. 1977. Maximum likelihood from incomplete data via

1 the EM algorithm. *J R Stat Soc Ser B* **39**: 1–38.

2 Dunham I, Kundaje A, Aldred SF, Collins PJ, Davis CA, Doyle F, Epstein CB, Frietze S,  
3 Harrow J, Kaul R, et al. 2012. An integrated encyclopedia of DNA elements in  
4 the human genome. *Nature* **489**: 57–74.

5 Flockerzi A, Ruggieri A, Frank O, Sauter M, Maldener E, Kopper B, Wullich B, Seifarth  
6 W, Müller-Lantzsch N, Leib-Mösch C, et al. 2008. Expression patterns of  
7 transcribed human endogenous retrovirus HERV-K(HML-2) loci in human  
8 tissues and the need for a HERV Transcriptome Project. *BMC Genomics* **9**: 354.

9 Francis OE, Bendall M, Manimaran S, Hong C, Clement NLNL, Castro-Nallar E, Snell Q,  
10 Schaalje GBB, Clement MJMJ, Crandall KAKA, et al. 2013. Pathoscope: species  
11 identification and strain attribution with unassembled sequencing data.  
12 *Genome Res* **23**: 1721–9.

13 Frazee AC, Jaffe AE, Langmead B, Leek J. 2014. *Polyester: simulating RNA-seq datasets*  
14 *with differential transcript expression*. Cold Spring Harbor Labs Journals.

15 Gage FH, Muotri AR. 2012. What makes each brain unique. *Sci Am* **306**: 26–31.

16 Garber M, Grabherr MG, Guttman M, Trapnell C. 2011. Computational methods for  
17 transcriptome annotation and quantification using RNA-seq. *Nat Methods* **8**:  
18 469–77.

19 Garrison KE, Jones RB, Meiklejohn D a, Anwar N, Ndhlovu LC, Chapman JM, Erickson  
20 AL, Agrawal A, Spotts G, Hecht FM, et al. 2007. T cell responses to human  
21 endogenous retroviruses in HIV-1 infection. *PLoS Pathog* **3**: e165.

1 Gnanakkan VP, Jaffe AE, Dai L, Fu J, Wheelan SJ, Levitsky HI, Boeke JD, Burns KH.  
2 2013. TE-array--a high throughput tool to study transposon transcription. *BMC*  
3 *Genomics* **14**: 869.

4 Göke J, Lu X, Chan Y-S, Ng H-H, Ly L-H, Sachs F, Szczerbinska I. 2015. Dynamic  
5 Transcription of Distinct Classes of Endogenous Retroviral Elements Marks  
6 Specific Populations of Early Human Embryonic Cells. *Cell Stem Cell* **16**: 135–  
7 141.

8 Gonzalez-Hernandez MJ, Cavalcoli JD, Sartor M a, Contreras-Galindo R, Meng F, Dai  
9 M, Dube D, Saha AK, Gitlin SD, Omenn GS, et al. 2014. Regulation of the Human  
10 Endogenous Retrovirus K (HML-2) Transcriptome by the HIV-1 Tat Protein. *J*  
11 *Virol* **88**: 8924–35.

12 Grow EJ, Flynn R a., Chavez SL, Bayless NL, Wossidlo M, Wesche DJ, Martin L, Ware  
13 CB, Blish C a., Chang HY, et al. 2015. Intrinsic retroviral reactivation in human  
14 preimplantation embryos and pluripotent cells. *Nature* **522**: 221–225.

15 Haase K, Mösch A, Frishman D. 2015. Differential expression analysis of human  
16 endogenous retroviruses based on ENCODE RNA-seq data. *BMC Med Genomics*  
17 **8**: 71.

18 Hanke K, Hohn O, Bannert N. 2016. HERV-K(HML-2), a seemingly silent subtenant -  
19 but still waters run deep. *Apmis* **124**: 67–87.

20 Hohn O, Hanke K, Bannert N. 2013. HERV-K(HML-2), the Best Preserved Family of  
21 HERVs: Endogenization, Expression, and Implications in Health and Disease.  
22 *Front Oncol* **3**: 246.

1 Jin Y, Tam OH, Paniagua E, Hammell M. 2015. TEtranscripts: a package for including  
2 transposable elements in differential expression analysis of RNA-seq datasets.  
3 *Bioinformatics* **31**: 3593–3599.

4 Jones RB, John VM, Hunter D V, Martin E, Mujib S, Mihajlovic V, Burgers PC, Luider  
5 TM, Gyenes G, Sheppard NC, et al. 2012. Human endogenous retrovirus K(HML-  
6 2) Gag- and Env-specific T-cell responses are infrequently detected in HIV-1-  
7 infected subjects using standard peptide matrix-based screening. *Clin Vaccine*  
8 *Immunol* **19**: 288–92.

9 Jurka J, Kapitonov V V, Pavlicek A, Klonowski P, Kohany O, Walichiewicz J. 2005.  
10 Repbase Update, a database of eukaryotic repetitive elements. *Cytogenet*  
11 *Genome Res* **110**: 462–7.

12 Karolchik D, Hinrichs AS, Furey TS, Roskin KM, Sugnet CW, Haussler D, Kent WJ.  
13 2004. The UCSC Table Browser data retrieval tool. *Nucleic Acids Res* **32**: D493-6.

14 Kellis M, Wold B, Snyder MP, Bernstein BE, Kundaje A, Marinov GK, Ward LD, Birney  
15 E, Crawford GE, Dekker J, et al. 2014. Defining functional DNA elements in the  
16 human genome. *Proc Natl Acad Sci U S A* **111**: 6131–8.

17 Krzywinski M et al. 2009. Circos: an Information Aesthetic for Comparative  
18 Genomics. *Genome Res* **19**: 1639–1645.

19 Li W, Lee M-H, Henderson L, Tyagi R, Bachani M, Steiner J, Campanac E, Hoffman DA,  
20 von Geldern G, Johnson K, et al. 2015. Human endogenous retrovirus-K  
21 contributes to motor neuron disease. *Sci Transl Med* **7**: 307ra153-307ra153.



1 Magiorkinis G, Belshaw R, Katzourakis A. 2013. “There and back again”: revisiting  
2 the pathophysiological roles of human endogenous retroviruses in the post-  
3 genomic era. *Philos Trans R Soc B Biol Sci* **368**: 20120504–20120504.

4 Marioni JC, Mason CE, Mane SM, Stephens M, Gilad Y. 2008. RNA-seq: an assessment  
5 of technical reproducibility and comparison with gene expression arrays.  
6 *Genome Res* **18**: 1509–17.

7 Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B. 2008. Mapping and  
8 quantifying mammalian transcriptomes by RNA-Seq. *Nat Methods* **5**: 621–8.

9 Mortelmans K, Wang-Johanning F, Johanning GL. 2016. The role of human  
10 endogenous retroviruses in brain development and function. *Apmis* **124**: 105–  
11 115.

12 Muotri AR, Marchetto MCN, Coufal NG, Oefner R, Yeo G, Nakashima K, Gage FH. 2010.  
13 L1 retrotransposition in neurons is modulated by MeCP2. *Nature* **468**: 443–  
14 446.

15 Muradrasoli S, Forsman A, Hu L, Blikstad V, Blomberg J. 2006. Development of real-  
16 time PCRs for detection and quantitation of human MMTV-like (HML)  
17 sequences. *J Virol Methods* **136**: 83–92.

18 Nexø B a, Villesen P, Nissen KK, Lindegaard HM, Rossing P, Petersen T, Tarnow L,  
19 Hansen B, Lorenzen T, Hørslev-Petersen K, et al. 2015. Are human endogenous  
20 retroviruses triggers of autoimmune diseases? Unveiling associations of three  
21 diseases and viral loci. *Immunol Res* **64**: 55–63.

1 Ormsby CE, Sengupta D, Tandon R, Deeks SG, Martin JN, Jones RB, Ostrowski M a,  
2 Garrison KE, Vázquez-Pérez J a, Reyes-Terán G, et al. 2012. Human endogenous  
3 retrovirus expression is inversely associated with chronic immune activation in  
4 HIV-1 infection. *PLoS One* **7**: e41021.

5 Pérot P, Mugnier N, Montgiraud C, Gimenez J, Jaillard M, Bonnaud B, Mallet F. 2012.  
6 Microarray-based sketches of the HERV transcriptome landscape. *PLoS One* **7**:  
7 e40194.

8 Perron H, Hamdani N, Faucard R, Lajnef M, Jamain S, Daban-Huard C, Sarrazin S,  
9 LeGuen E, Houenou J, Delavest M, et al. 2012. Molecular characteristics of  
10 Human Endogenous Retrovirus type-W in schizophrenia and bipolar disorder.  
11 *Transl Psychiatry* **2**: e201.

12 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing  
13 genomic features. *Bioinformatics* **26**: 841–842.

14 Rakoff-Nahoum S, J. Kuebler P, J. Heymann J, E. Sheehy M, M. Ortiz G, S. Ogg G,  
15 Barbour JD, Lenz J, Steinfeld AD, Nixon DF. 2006. Detection of T Lymphocytes  
16 Specific for Human Endogenous Retrovirus K (HERV-K) in Patients with  
17 Seminoma. *AIDS Res Hum Retroviruses* **22**: 52–56.

18 Rangwala SH, Zhang L, Kazazian HH. 2009. Many LINE1 elements contribute to the  
19 transcriptome of human somatic cells. *Genome Biol* **10**: R100.

20 Rodić N, Steranka JP, Makohon-Moore A, Moyer A, Shen P, Sharma R, Kohutek Z a,  
21 Huang CR, Ahn D, Mita P, et al. 2015. Retrotransposon insertions in the clonal  
22 evolution of pancreatic ductal adenocarcinoma. *Nat Med* **21**: 1060–4.

1     Santoni F a, Guerra J, Luban J. 2012. HERV-H RNA is abundant in human embryonic  
2     stem cells and a precise marker for pluripotency. *Retrovirology* **9**: 111.

3     Seifarth W, Spiess B, Zeilfelder U, Speth C, Hehlmann R, Leib-Mösch C. 2003.  
4     Assessment of retroviral activity using a universal retrovirus chip. *J Virol*  
5     *Methods* **112**: 79–91.

6     Subramanian RP, Wildschutte JH, Russo C, Coffin JM. 2011. Identification,  
7     characterization, and comparative genomic distribution of the HERV-K (HML-2)  
8     group of human endogenous retroviruses. *Retrovirology* **8**: 90.

9     Takahashi Y, Harashima N, Kajigaya S, Yokoyama H, Cherkasova E, McCoy JP,  
10     Hanada K, Mena O, Kurlander R, Tawab A, et al. 2008. Regression of human  
11     kidney cancer following allogeneic stem cell transplantation is associated with  
12     recognition of an HERV-E antigen by T cells. *J Clin Invest* **118**: 1099–109.

13    Tang Z, Steranka JP, Ma S, Grivainis M, Rodić N, Huang CRL, Shih I-M, Wang T-L,  
14    Boeke JD, Fenyö D, et al. 2017. Human transposon insertion profiling: Analysis,  
15    visualization and identification of somatic LINE-1 insertions in ovarian cancer.  
16    *Proc Natl Acad Sci* **114**: E733–E740.

17    Tarailo-Graovac M, Chen N. 2009. Using RepeatMasker to identify repetitive  
18    elements in genomic sequences. *Curr Protoc Bioinformatics* **Chapter 4**: Unit  
19    4.10.

20    Thorvaldsdóttir H, Robinson JT, Mesirov JP. 2013. Integrative Genomics Viewer  
21    (IGV): high-performance genomics data visualization and exploration. *Brief*  
22    *Bioinform* **14**: 178–92.

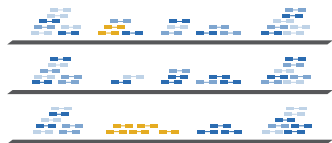
- 1 Trapnell C, Hendrickson DG, Sauvageau M, Goff L, Rinn JL, Pachter L. 2013.  
2 Differential analysis of gene regulation at transcript resolution with RNA-seq.  
3 *Nat Biotechnol* **31**: 46–53.
- 4 Trapnell C, Roberts A, Goff L, Pertea G, Kim D, Kelley DR, Pimentel H, Salzberg SL,  
5 Rinn JL, Pachter L. 2012. Differential gene and transcript expression analysis of  
6 RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc* **7**: 562–78.
- 7 Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Salzberg SL,  
8 Wold BJ, Pachter L. 2010. Transcript assembly and quantification by RNA-Seq  
9 reveals unannotated transcripts and isoform switching during cell  
10 differentiation. *Nat Biotechnol* **28**: 511–5.
- 11 Wang-Johanning F, Frost AR, Jian B, Epp L, Lu DW, Johanning GL. 2003. Quantitation  
12 of HERV-K env gene expression and splicing in human breast cancer. *Oncogene*  
13 **22**: 1528–1535.
- 14 Weiss RA. 2016. Human endogenous retroviruses: friend or foe? *APMIS* **124**: 4–10.
- 15 Wheeler TJ, Clements J, Eddy SR, Hubley R, Jones TA, Jurka J, Smit AFA, Finn RD.  
16 2013. Dfam: a database of repetitive DNA based on profile hidden Markov  
17 models. *Nucleic Acids Res* **41**: D70-82.
- 18 Wu CFJ. 1983. On the Convergence Properties of the EM Algorithm. *Ann Stat* **11**: 95–  
19 103.
- 20 Young GR, Mavrommatis B, Kassiotis G. 2014. Microarray analysis reveals global  
21 modulation of endogenous retroelement transcription by microbes.

1        *Retrovirology* **11**: 59.

2

3

## A Alignment



## B Annotation

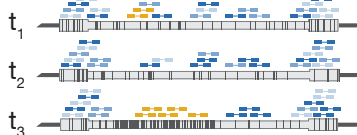


# Telescope

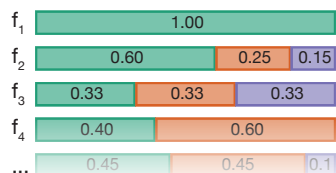
## C Initial proportions



## Initial fragment assignments

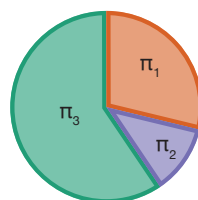


## D Expected assignment weights



## M-step

## E Transcript proportions



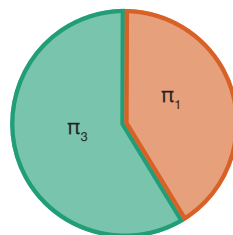
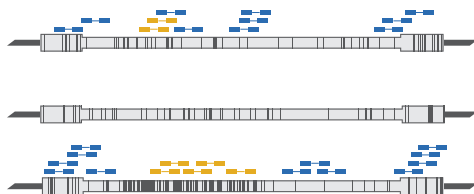
Expectation-Maximization (EM)

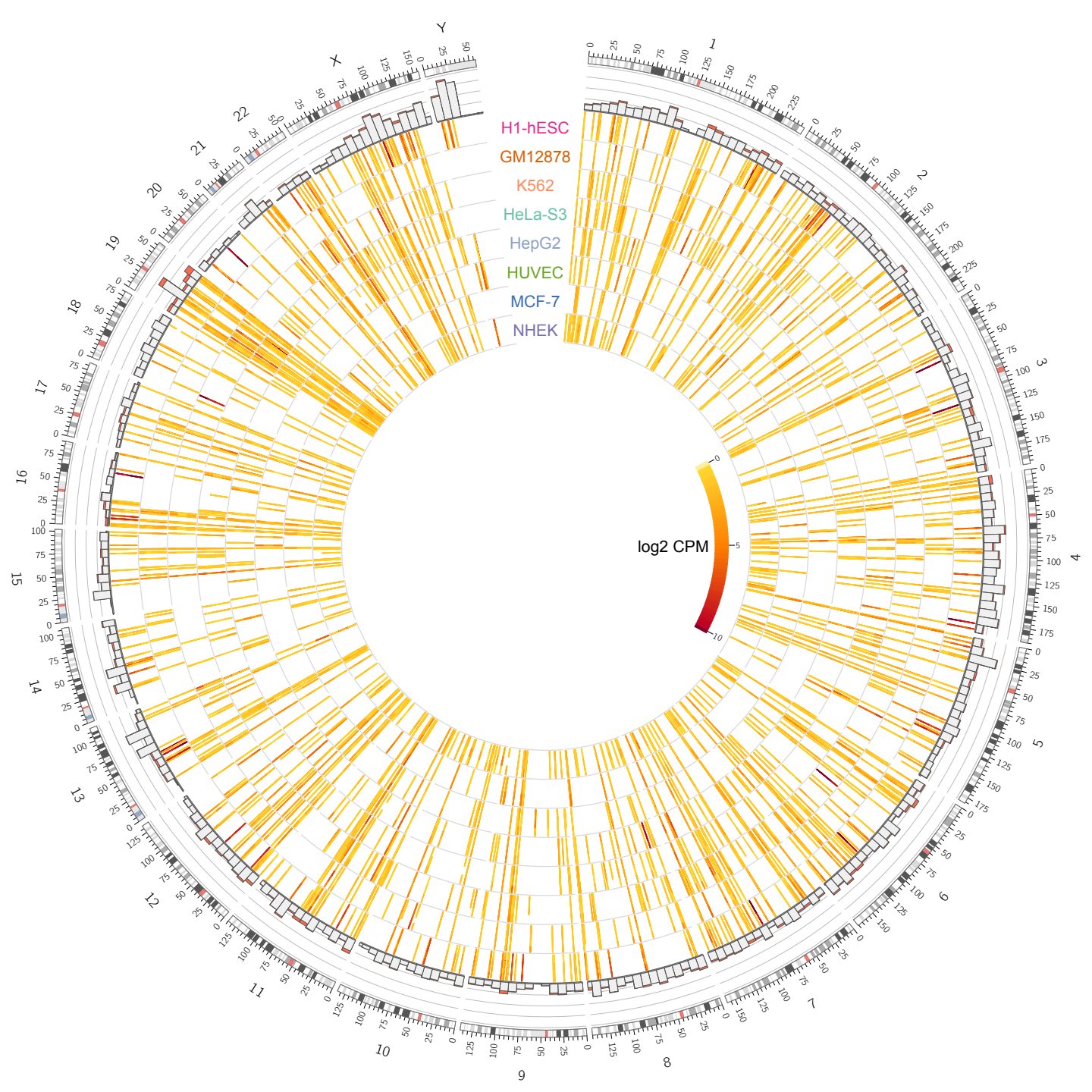
## E-step

## F

## Final fragment assignments

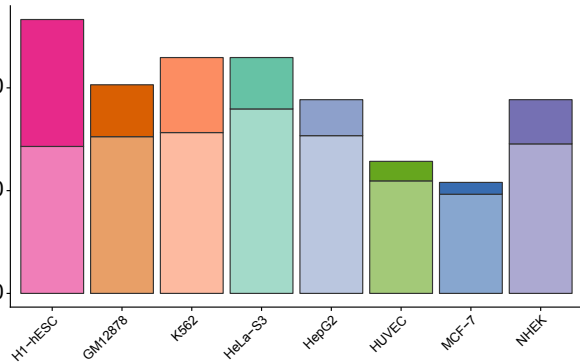
## MAP estimates of transcript proportions



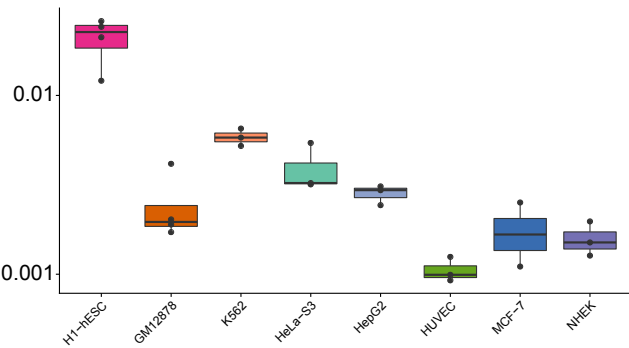
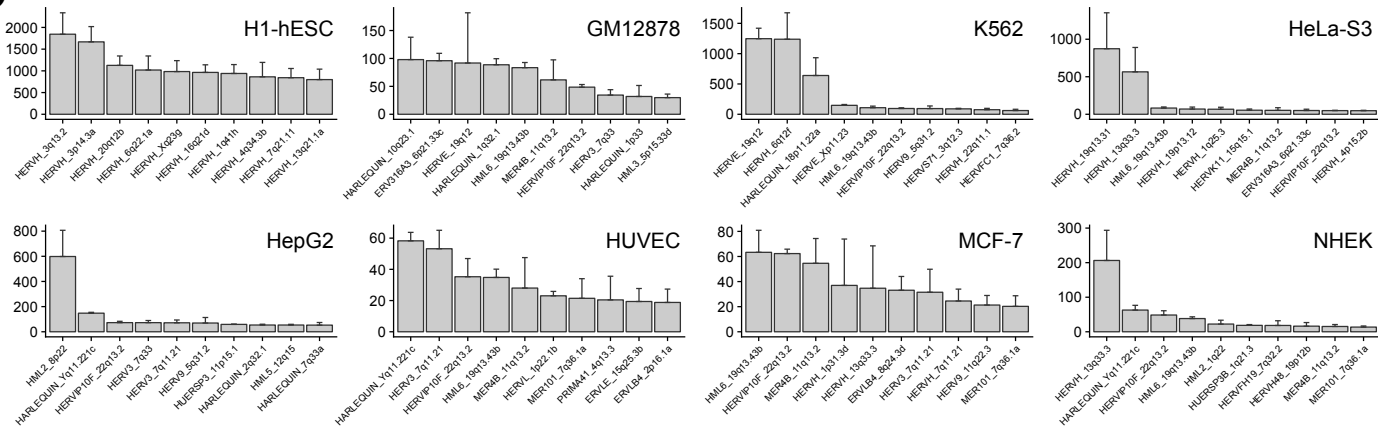


**A**

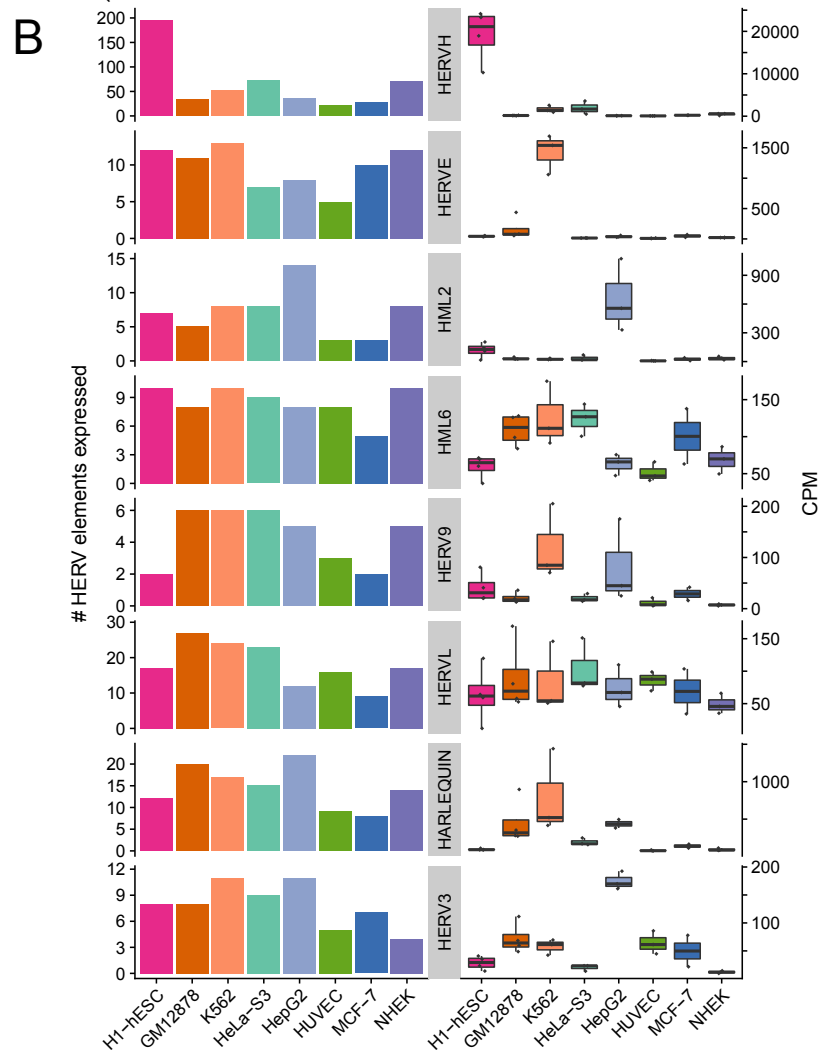
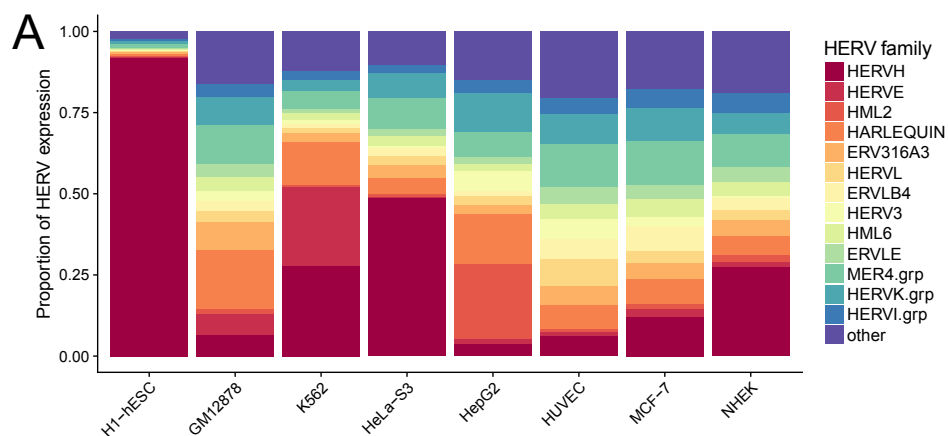
# HERV elements expressed

**B**

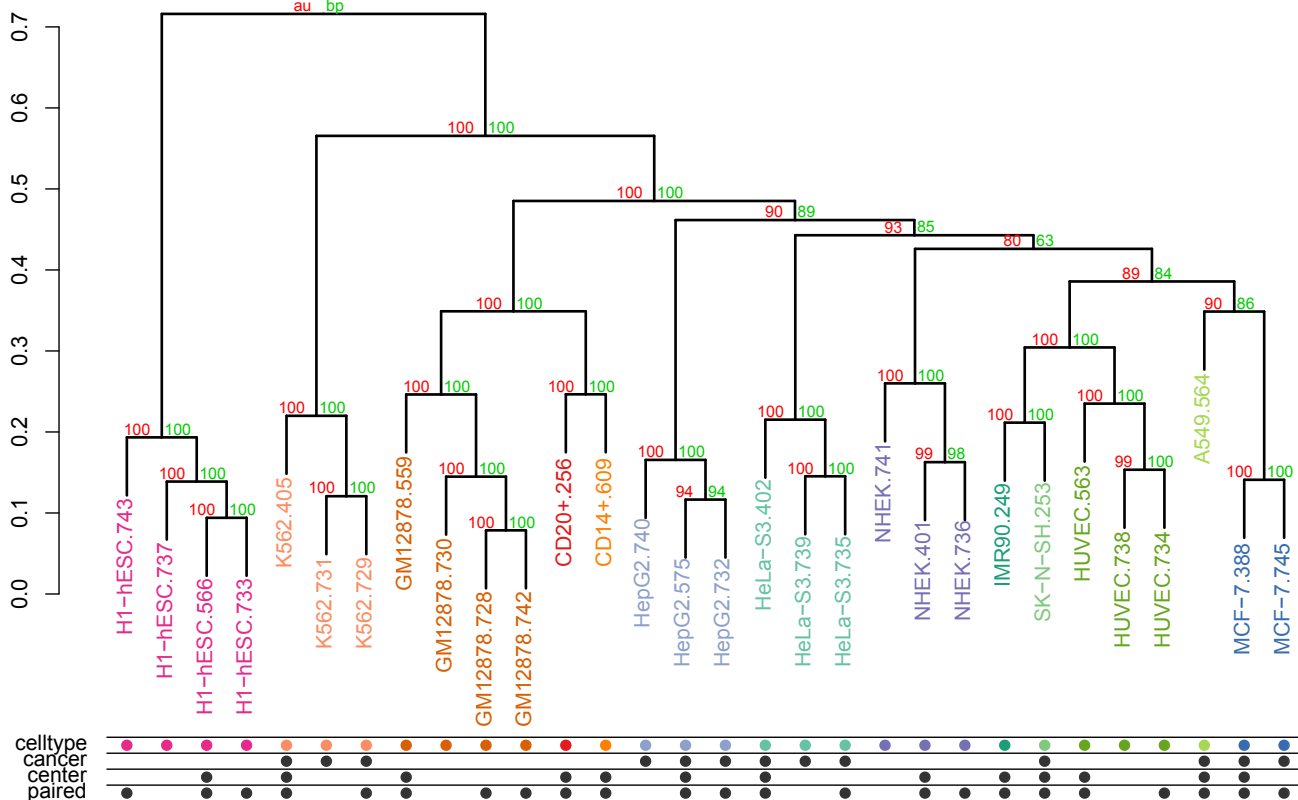
Proportion mapping to HERV

**C**

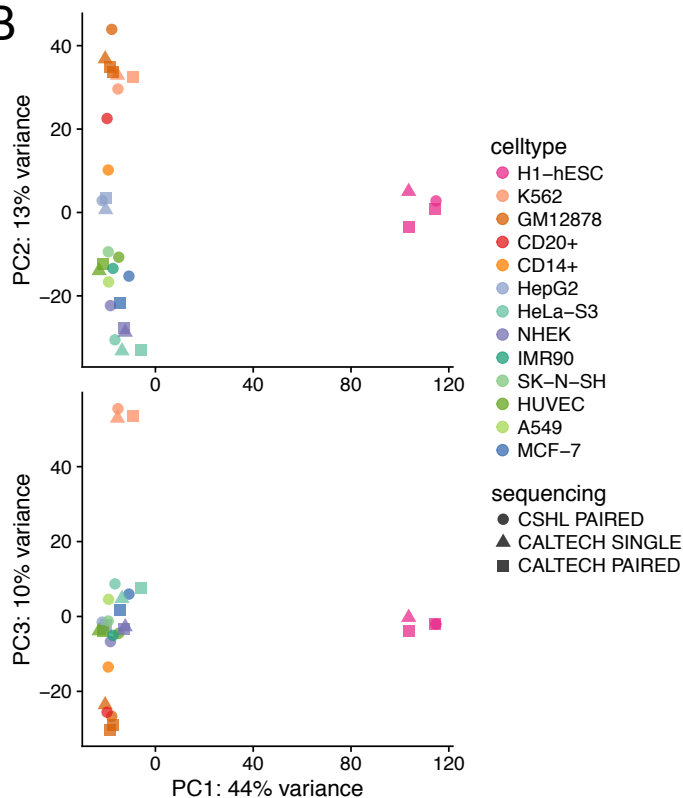




A



B



C

