1 **Widespread and tissue-specific expression of endogenous retroelements in human**

2 **somatic tissues**

3

4

5 Jean-David Larouche[1,2], Assya Trofimov[1,3], Leslie Hesnard[1,2], Gregory Ehx[1,2], Krystel

6 Vincent[1,2], Chantal Durette[1], Patrick Gendron[1], Jean-Philippe Laverdure[1], Éric Bonneil[1],

7 Caroline Côté[1], Sébastien Lemieux[1,3], Pierre Thibault[1,4] and Claude Perreault[1,2,5*].

8

9

10 1. Institute of Research in Immunology and Cancer, Université de Montréal, Montréal,

11 QC, Canada.

12 2. Department of Medicine, Université de Montréal, Montréal, QC, Canada.

13 3. Department of Informatics and Operational Research, Université de Montréal,

14 Montréal, QC, Canada.

15 4. Department of Chemistry, Université de Montréal, Montréal, QC, Canada.

16 5. Division of Hematology-Oncology, Hôpital Maisonneuve-Rosemont, Montréal, QC,

17 Canada.

18

19 *Correspondence:
20 Claude Perreault
21 IRIC - Université de Montréal,
22 P.O. Box 6128, Downtown Station
23 QC, Canada, H3C 3J7
24 claude.perreault@umontreal.ca
25
26
27
28

29

## Abstract

**Background:** Endogenous retroelements (EREs) constitute about 42% of the human genome and have been implicated in common human diseases such as autoimmunity and cancer. The dominant paradigm holds that EREs are expressed in embryonic stem cells (ESCs) and germline cells but are repressed in differentiated somatic cells. Despite evidence that some EREs can be expressed at the RNA and protein levels in specific contexts, a systems-level evaluation of their expression in human tissues is lacking.

**Methods:** Using RNA-sequencing data, we analyzed ERE expression in 32 human tissues, including medullary thymic epithelial cells (mTECs). A tissue-specificity index was computed to identify tissue-restricted ERE families. We also analyzed the transcriptome of mTECs in wild-type and AIRE-deficient mice. Finally, we developed a proteogenomic workflow combining RNA-sequencing and mass spectrometry (MS) in order to evaluate whether EREs might be translated and generate MHC I-associated peptides (MAP) in B-lymphoblastoid cell lines (B-LCL) from 16 individuals.

**Results:** We report that all human tissues express EREs but the breadth and magnitude of ERE expression are very heterogeneous from one tissue to another. ERE expression was particularly high in two MHC-I-deficient tissues (ESCs and testis) and one MHC-I-expressing tissue, mTECs. In mutant mice, we report that the exceptional expression of EREs in mTECs was AIRE-independent. MS sequencing identified 104 non-redundant MAPs in B-LCLs. These MAPs preferentially derived from sense translation of intronic

52 EREs. Notably, detailed analyses of their amino acid composition revealed that ERE-

53 derived MAPs presented homology to viral MAPs.

54

55 **Conclusions:** This study shows that ERE expression in somatic tissues is more pervasive

56 and heterogeneous than anticipated. The high and diversified expression of EREs in

57 mTECs and their ability to generate MAPs suggest that EREs may play an important role

58 in the establishment of self-tolerance. The viral-like properties of ERE-derived MAPs

59 suggest that those not expressed in mTECs can be highly immunogenic.

60

61 **Keywords:** Endogenous retroelements, immunopeptidome, major histocompatibility

62 complex, medullary thymic epithelial cells, somatic tissues, systems biology,

63 transcriptome.

64

65 **Background**

66 Endogenous retroelements (EREs) are remnants of transposable elements that successfully

67 integrated our germline DNA millions of years ago (1, 2). After initial integration in the

68 genome, EREs further increased their copy number via several successive waves of

69 retrotransposition (3, 4). Now, most ERE sequences contain mutated or truncated open

70 reading frames and have lost their capacity to transpose in the genome (2). Phylogenic

71 analyses have allowed the classification of EREs in families based on sequence homology

72 (5, 6). Most EREs are categorized in three groups, which altogether comprise ~40-50% of

73 the human genome: the long-terminal repeats (LTR) as well as the long and short

74 interspersed nuclear elements (LINE and SINE) (7-9).

75

Hosts repress ERE expression in order to protect their genomic integrity from deleterious

insertions of EREs in open reading frames (10, 11). Indeed, a strict epigenetic regulation

of ERE sequences is applied at both the DNA and histone levels (12). Growing evidence

suggests that KRAB zinc finger proteins (KZFPs) are involved in an evolutionary arms

race to repress the expression of novel ERE integrations (13). KZFPs recruit numerous

restriction factors to silence ERE sequences: the histone methyltransferase SETDB1, DNA

methyltransferase proteins, the nucleosome remodeling and deacetylase complex NuRD

and the heterochromatin protein HP1 (14). KZFP-independent mechanisms, such as the

HUSH complex (15) and the histone demethylase LSD1 (16), also apply non-redundant

epigenetic silencing on ERE sequences. Nevertheless, some "domesticated" EREs

contribute at many levels to human development and survival. Indeed, ERE sequences

provide promoters and enhancers to several human genes and thereby regulate the

expression of genes implicated in interferon responses, DNA damage response in the male

germline and maintenance of stem cell pluripotency (17-19). Additionally, a LINE-derived

transcript is essential to embryonic stem cells (ESCs) self-renewal via activation of rRNA

synthesis (20). Finally, syncytins are ERE-derived proteins that mediate cell-cell fusion to

allow formation of the placental syncytium (21, 22).

93

The dominant paradigm holds that EREs are expressed in ESCs as well as in germline cells,

but are repressed in other differentiated cells outside specific contexts in which they have

relevant functions (12). However, studies on ERE expression have been limited to subsets

of ERE families in one or few tissues. Additionally, to our knowledge, no study has

98  addressed ERE expression in the thymus where central T-cell immune tolerance is

99  established. Hence, we have no clue as to the ability of EREs to induce T-cell tolerance. In

100  the present study, we established an atlas of ERE expression in a panel of 30 healthy human

101  tissues and 2 cell types, including medullary thymic epithelial cells (mTECs). We first

102  demonstrate that ERE expression is widespread in human tissues, but with tissue-specific

103  profiles. Notably, three cell types showed particularly high and diversified expression of

104  EREs: ESCs, testis and mTECs. By analyzing the transcriptome of wild-type and AIRE-

105  deficient mice, we found that the impressive expression of EREs in mTECs was AIRE-

106  independent. In addition, our mass spectrometry (MS) analyses revealed that the three main

107  groups of EREs generate MHC I-associated peptides (MAPs) in healthy cells. Finally, we

108  demonstrate that ERE-derived MAPs (ereMAPs) retained strong homology to viruses.

109

110  **Methods**

111  **Transcriptomic data manifest**

112  RNA-seq data of 30 non-redundant human tissues were downloaded from the Genotype-

113  Tissue Expression (GTEx) on the dbGaP portal (accession number phs000424.v8.p2.c1)

114  (23). When possible, 50 samples were randomly selected per tissue, otherwise all available

115  samples were analyzed. Transcriptomic data of ESCs were downloaded from the sequence

116  read archive from Lister *et al* (24). RNA-seq data of purified hematopoietic cells were

117  obtained from the Gene Expression Omnibus (GEO) (projects PRJNA384650 and

118  PRJNA225999). Six human mTEC samples were analyzed: four from (25) and two

119  additional samples processed with the same protocol with minor modifications: i) after

120  transfer to our laboratory, thymic samples were frozen in cryovials containing a

121    cryoprotective medium composed of 5% DMSO and 95% Dextran-40 solution (5%

122    concentration), ii) CD45⁻ cells were magnetically enriched with the CD45 Microbeads

123    human kit from Miltenyi Biotec (no. 130-045-801) prior to mTEC sorting, iii) cDNA

124    libraries were prepared with the KAPA mRNAseq stranded kit (KAPA, Cat no. KK8421),

125    and iv) sequencing generated around $400 \times 10^6$ reads per sample. For the complete list of

126    human samples analyzed, see Table S1 of Additional File 2. Mature murine mTECs

127    (mTEC^hi) data were obtained from St-Pierre *et al* (26) on GEO (accession GSE65617).

128

**Expression of transcripts derived from EREs and canonical genes**

130    RNA-seq reads of human samples were trimmed with Trimmomatic *0.35* (27) to remove

131    adapters and low quality sequences. Expression levels of transcripts and endogenous

132    retroelements were quantified in transcripts per million (TPM) with kallisto *0.43.1* (28)

133    with an index composed of i) GRCh38.88 transcripts and human ERE sequences from

134    RepeatMasker (downloaded on the UCSC Table Browser on July 19, 2018) or ii) GRCm38

135    transcripts and murine ERE sequences from RepeatMasker (downloaded on the UCSC

136    Table Browser on March 11, 2019) for human and murine samples, respectively. TPM

137    values of transcripts and ERE sequences were grouped in genes and ERE families based

138    on Ensembl and RepeatMasker annotations, respectively.

139

**ERE expression profiling in human tissues**

141    Expression levels of ERE families were computed for each tissue by calculating the median

142    expression across all samples for a given tissue. The numbers of standard deviations from

143    the mean (row Z-score) of ERE families for each tissue were determined using the scale

144    function in R. The Euclidean distance was then calculated between all tissues based on the

145    row Z-scores of ERE families, followed by an unsupervised hierarchical clustering. Finally,

146    the tree was manually separated in three clusters of tissues. Standard deviations of

147    expression of each ERE family between samples of a given tissue were also computed.

148

149    **Quintile ranking of ERE expression in somatic tissues**

150    Median expression of ERE families were calculated among all samples of a given tissue.

151    Tissues were then ranked based on their expression level of each ERE family individually

152    and assigned to quintiles of 6, 6, 8, 6 and 6 tissues, respectively. Finally, tissues were sorted

153    based on the number of times they were assigned to the fifth quintile.

154

155    **Identification and characterization of tissue-restricted EREs (TREs)**

156    The $\tau$-index of tissue specificity was calculated as per Yanai *et al* (29). Briefly, the $\tau$-index

157    is defined as:

158    $$\tau = \frac{\Sigma_{i=1}^{N}(1 - x_i)}{N - 1}$$

159    where $x_i$ is the level of expression of a gene or ERE family in tissue i normalized to its

160    maximal expression level among tissues and N is the number of tissues. Genes and ERE

161    families with $\tau \geq 0.8$ were considered as tissue-restricted. To determine in which tissue(s) a

162    tissue-restricted gene or ERE family was overexpressed, a binary pattern was computed as

163    reported by Yanai *et al* (29). Briefly, tissues were sorted based on their expression level for

164    each tissue-restricted gene (TRG) or ERE family (TRE). The distance between neighboring

165    tissues was calculated, and the maximal distance or 'gap' was used as threshold for the

166    binary pattern. Tissues with an expression level above the gap were considered as

167  overexpressing the TRG or TRE while other tissues were considered as underexpressing

168  them, and were given a value of 1 or 0, respectively. ERE groups were determined for all

169  identified TREs, and the proportions of LINE, LTR and SINE elements in TREs were

170  compared to their representation among ERE families. A chi-squared test was performed

171  to assess enrichment of discrete ERE groups among TREs. Using the above described

172  binary pattern, the number of overexpressing tissues was determined for each TRG or TRE.

173

174  **Impact of AIRE on ERE expression in mTECS**

175  Lists of AIRE-dependent, AIRE-independent and constitutively expressed genes were

176  generated as per St-Pierre *et al* (26). Expression levels of these three sets of genes as well

177  as ERE families were compared between wild-type (n=3) and AIRE knock-out (n=3)

178  murine mTEC$^{hi}$ using Wilcoxon tests. Expression levels of each individual ERE family

179  were also compared between wild-type and AIRE knock-out mice using Wilcoxon tests.

180

181  **MS analyses**

182  Peptidomic data of a cohort of 16 B-lymphoblastoid cell lines (B-LCL) samples from

183  Pearson *et al* (30) were downloaded from the Pride Archive (Project PXD004023). For the

184  detailed protocol of mild acid elution and peptide processing, see Granados *et al* (31).

185  Peptides were identified using Peaks X (Bioinformatics Solution Inc.) and peptide

186  sequences were searched against the personalized proteome of each sample. For peptide

187  identification, tolerance was set at 5 ppm and 0.02 Da for precursor and fragment ions,

188  respectively. Occurrence of oxidation (M) and deamination (NQ) were considered as post-

189  translational modifications.

190

**Identification of ereMAPs**

192 For individual B-LCL samples, RNA-seq reads were aligned to the reference genome

193 GRCh38.88 using STAR (32) with default parameters. Using the intersect mode of the

194 BEDTools suite (33), reads entirely mapping in RepeatMasker and Ensembl annotations

195 were separated in ERE and canonical datasets respectively, and any read seen in the

196 canonical dataset was discarded from the ERE dataset. Unmapped reads, secondary

197 alignments and low quality reads were then removed from the ERE dataset using Samtools

198 view (34) with the following parameters: -f "163", "147", "99" or "83" and -F "3852". In

199 order to keep a manageable database size, ambiguous nucleotides were trimmed from reads

200 of the ERE dataset, followed by translation in all possible reading frames. Finally, the

201 resulting ERE amino acid sequences were spliced to remove sequences following stop

202 codons. Only sequences of at least 8 amino acids were kept and given a unique ID to

203 generate a theoretical ERE proteome. In parallel, a canonical personalized proteome

204 containing the polymorphisms of the donor was generated as per (25) for each sample.

205 Briefly, single-nucleotide variants were detected using freebayes version 1.0.2 (35), and

206 variants with a minimal alternate count of 5 were inserted in transcript sequences using

207 pyGeno (36). Expression levels of transcripts were quantified with kallisto using

208 GRCh38.88 transcripts (downloaded from Ensembl) as index, and only transcripts with a

209 TPM>0 were translated into a canonical proteome, which was concatenated with the ERE

210 proteome to generate a Personalized Proteome unique to each sample.

211

**Peptide annotation and validation**

213    Following peptide identification, a list of unique peptides was extracted for each sample

214    and a false discovery rate (FDR) of 5% was applied on the peptide scores. Binding affinities

215    to the sample's HLA alleles were predicted with NetMHC4.0 (37) or with NetMHCpan-

216    4.0 (38) when an HLA allele was not included in NetMHC4.0, and only 8 to 11-amino-

217    acid-long peptides with a percentile rank $\leq 2\%$ were included for further annotation. For

218    each peptide, a binary code was generated based on the presence or absence of its amino

219    acid sequence in the ERE and canonical proteomes and an ERE status of "Yes", "Maybe"

220    or "No" was given to the peptide accordingly. Peptides that were seen only in the ERE

221    proteome or the canonical proteome were classified as "Yes" and "No" respectively. To

222    determine if candidates with a "Maybe" status were ereMAP candidates, we retrieved all

223    their possible nucleotide coding sequences from the sample's reads and split them in a set

224    of 24-nucleotide-long subsequences (k-mers). These k-mers were then queried in 24-

225    nucleotide-long k-mer databases generated from our ERE and canonical reads datasets

226    using Jellyfish version 2.2.3 (39) (with the -C argument to consider the read's sequence

227    and its reverse complement). Only peptides encoded by more than one read were kept for

228    further validation to reduce risks of sequencing errors. If at least one of the MAP-coding

229    sequences (MCS) was only seen in the canonical read dataset, the peptide was discarded.

230    "Maybe" peptides were considered as ereMAP candidates if the minimal occurrence of

231    their most abundant MCS was at least 10 times higher in the ERE k-mer database than in

232    the canonical k-mer database. Because leucine and isoleucine variants are not

233    distinguishable by standard MS approaches, all possible I/L variants for each ereMAPs

234    candidates were searched in the personalized proteome. If one of the I/L variants had a

235    higher expression in the personalized proteome, the ereMAP candidate was discarded. The

236 genomic region generating each ereMAP candidate was determined by mapping the reads

237 coding for the peptide on the GRCh38.88 assembly of the reference genome with the BLAT

238 algorithm of the UCSC Genome Browser. If a clear genomic region could not be found,

239 the peptide was discarded. Genomic regions coding for ereMAPs candidates were then

240 inspected in IGV (40) to see if the MCS contained known germline polymorphisms (using

241 dbSNP v.149), and candidates were kept or discarded based on their orientation in ERE

242 and annotated sequences. Briefly, any ereMAP candidate whose MCS mapped in the sense

243 of a gene coding sequence was discarded, whereas candidates whose coding sequences

244 mapped in intergenic regions were considered as ereMAPs no matter their orientation.

245 Candidates were also discarded if they fulfilled these two conditions: i) their MCS mapped

246 in the sense of an intron and in antisense of the ERE, and ii) if their MCS did not map in

247 other ERE sequences (for the complete decision tree, see Figure S3). Finally, MS/MS

248 spectra of the ereMAPs candidates were manually validated to ensure the quality of the

249 identification. Peptides that passed all these validation steps were then considered as

250 ereMAPs.

251

252 **Characterization of ereMAPs**

253 During manual validation in IGV, characteristics regarding the family and group of the

254 ERE generating the peptides, the type of genomic region encoding the peptide (coding

255 sequence, intronic or intergenic) and the orientation of the peptide sequences (sense or

256 antisense) were retrieved for individual ereMAPs. When a peptide was identified in

257 multiple samples and had different characteristics depending upon the sample, all

258 possibilities were kept, otherwise they were aggregated to reduce redundancy. The

259    expression levels of ERE families that were source or non-source of ereMAPs were

260    averaged among B-LCL samples, and their distributions were compared with a Mann-

261    Whitney test. We next compared the proportions of the three main groups of EREs (LINE,

262    LTR and SINE) in the genome, transcriptome and immunopeptidome. Representation of

263    EREs in the transcriptome was assessed in our B-LCL samples: the expression levels of

264    LINE, LTR and SINE elements were summed in each sample and divided by the expression

265    level of all EREs. We then averaged these transcriptomic proportions across all B-LCL

266    samples. We used immunopeptidomic proportions of LINE, LTR and SINE elements from

267    the ereMAPs identified in this work, whereas the genomics proportions were taken from

268    Treangen *et al* (8). A chi-squared test was performed to compare the proportions of ERE

269    groups at the genomic, transcriptomic and immunopeptidomic levels. The proportions of

270    ERE sequences located in intergenic and intronic regions as well as in coding sequences

271    were determined by intersecting the genomic localization of ERE sequences with the

272    localization of introns and exons from the UCSC Table Browser (files downloaded on

273    August 21, 2019). A chi-squared test was used to determine the enrichment of a certain

274    genomic region for ereMAPs generation. Finally, Pearson correlation between the number

275    of ereMAPs generated by each ERE family and the number of copies of the family's

276    sequence in the human genome (determined from RepeatMasker annotations) was

277    computed with a confidence level of 95%.

278

279    **GTEx profiling of ereMAP expression**

280    To evaluate the expression of the ereMAP-coding sequences in peripheral tissues, we

281    downloaded RNA-seq data of 30 tissues from the GTEx consortium (phs000424.v7.p2).

282    For the complete protocol of this analysis, see Laumont *et al* (25). Briefly, we generated

283    24-nucleotide-long k-mer databases for each sample, in which we queried each ereMAP-

284    coding sequence's 24-nucleotide-long k-mer set. For each ereMAP, the minimal

285    occurrence in the k-mer set was used as the number of reads coding for the peptide in a

286    given sample ($r_{overlap}$). The number of reads coding for a peptide was normalized between

287    RNA-seq experiments by dividing $r_{overlap}$ by the total number of reads of the sample and

288    multiplying this number by $10^8$ to obtain the number of reads detected per hundred million

289    reads sequenced (rphm). We then averaged the log-transformed rphm values (*$log_{10}$(rphm*

290    *+ 1)*) for each tissue, and an average expression superior to 10 rphm in a tissue was

291    considered as significant.

292

293    **Amino acid composition of ereMAPs**

294    In addition to the list of ereMAPs identified on our B-LCL samples, two linear and MHC

295    I-restricted epitopes' sequences datasets were downloaded from the Immune Epitope

296    Database: a first dataset of 36 472 MAPs from any virus infecting human cells and a second

297    one of 282 069 human canonical MAPs (downloaded on August 7, 2019). Lists of 8 to 11-

298    amino-acid-long MAPs were extracted from these two datasets. Usage frequency of each

299    amino acid was calculated by dividing their occurrences by the total number of amino acids

300    in the ERE, viral and human canonical MAPs datasets. In parallel, datasets were separated

301    in subsets of 8, 9, 10 and 11-amino-acid-long MAPs, and frequencies of amino acids were

302    computed for each peptide position of each subset of MAPs. The 11-amino-acid-long MAP

303    subset was discarded because of an insufficient number of ereMAPs (n = 2).

304

**Viral homology**

To assess the similarity between ereMAPs and viral peptides, we used the same datasets of viral and human canonical MAPs from the Immune Epitope Database used for the amino acid composition analysis (see section "Amino acid composition of ereMAPs" of the Methods). We aligned ereMAP sequences to this database of viral peptides using version 2.2.28 of the Protein Basic Local Alignment Tool (BLASTp) (41) in the blastp-short mode with the following arguments: -word_size 2, -gapopen 5, -gapextend 2, -matrix PAM30, and -evalue 10 000 000. As a control, human canonical MAPs were aligned to the viral peptides dataset with BLASTp. For the viral homology analysis, we compared the 104 ERE MAPs to 10,000 groups of 104 randomly sampled canonical MAPs. We calculated the percentage of identity (%I) of ereMAPs and canonical MAPs with viral peptides as:

$$\%_I = \frac{M_{max} \times L_a}{L_p} \times 100\%$$

where $M_{max}$ is the maximal percentage of identical matches with the viral MAPs database, $L_a$ is the length of the alignment and $L_p$ is the length of the ereMAP or the canonical MAP. The average percentage of identity of ereMAPs and each subgroup of the bootstrap distribution was computed, and the p-value was determined as the number of times that the percentage of identity of the bootstrap distribution was higher than the percentage of identity of ereMAPs divided by the number of bootstrap iterations (10,000) as per Granados *et al* (42).

## Results

**Expression of ERE transcripts in normal human tissues and cells**

327    To assess ERE expression in heathy human tissues, we quantified the expression levels of

328    the 809 ERE families contained in the RepeatMasker annotations in 1371 samples from 32

329    different healthy human tissues and cell types. We calculated the median expression of

330    each ERE family among samples of a given tissue or cell type (Table S2) and then

331    computed the row Z-score across tissues. Unsupervised hierarchical clustering of tissues

332    based on ERE expression allowed us to identify 3 clusters of high (cluster 1), intermediate

333    (cluster 2) and low (cluster 3) ERE expression (Fig. 1). High ERE expression (cluster 1) in

334    ESCs and testis was expected. The salient finding was the high ERE expression in mTECs

335    which, to the best of our knowledge, has never been reported before. Comparison with

336    hematopoietic cell types at several differentiation stages confirmed the high ERE

337    expression in mTECs and ESCs (Figure S1A). For brevity, mTECs and ESCs will be

338    referred to as tissues in the following paragraphs. Computing the standard deviation of ERE

339    expression among individual samples for each tissue also revealed that most ERE families

340    displayed low interindividual variability (Figure S1B). Finally, while quintile ranking

341    analysis showed that ERE expression was generally concordant among ERE families in

342    each tissue analyzed, almost all tissues expressed some ERE families at high level (Figure

343    S2), suggesting that some tissue-specific factors regulate ERE expression in human tissues.

344

345    **Most human tissues show a tissue-specific ERE expression.**

346    To ascertain if expression of discrete ERE families was restricted to specific tissues, we

347    computed the $\tau$-index of tissue-specificity as defined by Yanai *et al* (29). Briefly, the $\tau$-

348    index compares the expression of a gene in a set of tissues and has a value ≤0.4 for

349    housekeeping genes and ≥0.8 for tissue-restricted genes (43). We identified a total of 124

350 ERE families with a tissue-restricted expression. As control, we computed the **τ**-index for

351 annotated genes and known tissue-restricted genes (TRGs), such as *INS*, *CRP* and

352 *CHRNA1*. The majority (108/124) of the tissue-restricted ERE families (TREs) were

353 identified in ESCs, testis and mTECs, revealing that in addition to their high expression of

354 EREs, these tissues expressed a broader repertoire of EREs than other tissues (Fig. 1, Fig.

355 2A). Nonetheless, tissue-restricted expression of EREs is a widespread phenomenon across

356 human tissues because we identified TREs in 17 out of the 32 human tissues analyzed. For

357 a given tissue, the number of TREs is positively associated with the number of TRGs (Fig.

358 2A) suggesting some commonality between expression of TRGs and TREs. We also

359 identified in TREs a significant enrichment of LTRs relative to LINE and SINE families

360 (Fig. 2B). Finally, TREs' expression was typically restricted to fewer tissues than TRGs,

361 with 91.7% of TREs being tissue-specific (Fig. 2C, Table S3). Altogether, these results

362 show that ERE expression in healthy human tissues is widespread but not homogenous.

363 Indeed, 124 ERE families, most of which are LTR elements with low copy numbers,

364 showed tissue-specific expression.

365

366 **Impact of the *AIRE* gene on ERE expression in mTECS**

367 Out of the three tissues with high ERE expression (Fig. 1), two are known to express no or

368 barely detectable MHC-I molecules (testis and ESCs, respectively), whereas mTECs

369 express standard levels of MHC I (44-46). Promiscuous expression of genomic sequences

370 is a quintessential feature of mTECs that is driven in part by the *AIRE* gene and also by

371 other genes whose identity is still debated (47). Since the role of mTECs is to induce

372 tolerance to the MAPs that they display, EREs expressed in mTECs could be tolerogenic.

373    However, T cell-mediated responses towards EREs were previously observed, suggesting

374    that the establishment of central tolerance towards EREs in the thymus is incomplete (48,

375    49). Therefore, we next investigated the contribution of the AIRE transcription factor to

376    ERE expression in mTECs. To do so, we quantified the expression of ERE families as well

377    as canonical genes in mTECs extracted from wild-type and AIRE knock-out mice.

378    Canonical genes were sorted in three categories based on St-Pierre *et al* (26) : i)

379    constitutively expressed genes, ii) AIRE-independent TRGs and iii) AIRE-dependent

380    TRGs. As expected, expression of AIRE-dependent TRGs significantly decreased in the

381    absence of AIRE, whereas constitutively expressed genes and AIRE-independent TRGs

382    were minimally affected by AIRE absence (Fig. 3A) Strikingly, global ERE expression

383    was independent of AIRE since it was unchanged in AIRE knock-out relative to wild-type

384    mice (Fig. 3A). Furthermore, computing Mann-Whitney tests for each ERE family revealed

385    that the absence of AIRE did not affect the expression of any ERE family (Fig. 3B). Hence,

386    expression of all ERE families was independent of AIRE in mTECs.

387

388    **Translation of ERE transcripts by healthy cells**

389    We next sought to determine whether some ERE transcripts are translated in healthy cells.

390    When performed on whole cell extracts, MS is strongly biased for identification of

391    abundant and stable proteins at the proteome level. We therefore decided to investigate the

392    contribution of EREs to the immunopeptidome, which is mainly composed of peptides

393    derived from rapidly degraded proteins (50, 51). To do so, we reanalyzed previously

394    reported transcriptomic and peptidomic data from 16 B-lymphoblastoid cell lines (B-LCL)

395    (Table S4) (30). As conventional approaches do not include ERE sequences, precluding

396    identification of ereMAPs, we developed a proteogenomic workflow combining RNA-

397    sequencing and MS to enable ereMAP identification (Fig. 4A, Figure S3). Briefly, we

398    generated for each B-LCL a personalized proteome that contained only the sample's

399    expressed sequences as well as its polymorphisms. Canonical and ERE RNA sequences

400    were translated *in silico* and concatenated to generate a personalized proteome that was

401    used to identify MAPs in MS analyses (Fig. 4A). For each MAP identified, we retrieved

402    the peptide's coding sequence and proceeded to its annotation. Two categories of peptides

403    were kept as ereMAP candidates to be further manually validated: i) peptides that were

404    only seen in the ERE proteome, and ii) peptides seen in both the ERE and canonical

405    proteomes ("Maybe" candidates) and for which the occurrence of the coding sequences

406    was at least 10-fold higher in ERE reads compared to canonical reads. Our proteogenomic

407    approach enabled the identification of 130 ereMAPs in the 16 B-LCL samples analyzed,

408    revealing that ERE sequences are translated in non-neoplastic cells (Fig. 4B). Of those, 104

409    were non-redundant, confirming that ereMAPs can be shared by multiple individuals

410    (Table S5). Of course, the extent of interindividual sharing would be considerably greater

411    in cohorts of HLA-matched individuals since various HLA allotypes present different sets

412    of MAPs (50). Profiling of the ereMAPs' RNA expression in healthy human tissues showed

413    that 26% (27/104) of ereMAPs' coding sequences were expressed at high levels by multiple

414    tissues (Figure S4). Hence, since highly expressed transcripts are preferential sources of

415    MAPs (30), ereMAPs derived from abundant transcripts could be presented on the surface

416    of a wide range of tissues (Figure S4). We also observed that ereMAPs were generated by

417    the three main groups of ERE sequences (SINE, LINE, LTR), confirming that they all have

418    the potential to be translated in healthy cells (Fig. 4C). Together, these proteogenomic

419    analyses show that several EREs are translated and generate ereMAPs in B-LCLs, and

420    suggest that this is also the case for a wide range of human tissues.

421

422    We next investigated the mechanisms leading to presentation of ereMAPs on the cell

423    surface. First, we noted that ereMAPs preferentially derived from highly expressed ERE

424    transcripts (Fig. 5A). For the majority of ereMAPs, this transcription was in the same sense

425    as the ERE sequence in the genome, but ~30% of ereMAPs (34/104) resulted from

426    antisense transcription (Fig. 5B), which is common for EREs (52-54). Even though

427    ereMAPs were generated by the three main groups of EREs (Fig. 4C), the relative

428    frequency of LTR translation was higher than that of LINEs and SINEs (Fig. 5C). Indeed,

429    the representation of LTRs in the immunopeptidome was superior to the space they occupy

430    in the genome or their abundance in the transcriptome (Fig. 5C). Additionally, intronic

431    EREs were a preferential source of ereMAPs: while 51% of EREs were intronic, 79% of

432    ereMAPs derived from intronic EREs (Fig. 5D). Finally, when we assigned a genomic

433    location to ereMAPs, we noted that some ERE families generated several distinct ereMAPs

434    (Table S5). This can be explained in part by variations in the genomic space occupied by

435    the various ERE families. Indeed, for the various ERE families, we observed a moderate,

436    yet significant, correlation between the number of genomic copies and the number of

437    ereMAPs (Fig. 5E). Altogether, these results demonstrate that i) ereMAPs are generated by

438    both sense and antisense transcripts that are preferentially located in introns and expressed

439    at high levels, and ii) generation of ereMAPs is enhanced when a family belongs to the

440    LTR group occupying a large genomic space.

441

442 **ereMAPs have a viral-like amino acid composition**

443 We next asked to what extent ereMAPs and their coding transcripts might retain some

444 traces of their phylogeny ("viral features"). We found conspicuous differences between

445 amino acid frequencies in ereMAPs relative to both viral MAPs and canonical human

446 MAPs listed in the Immune Epitope Database (Fig. 6A). Indeed, ereMAPs showed lower

447 abundance of multiple amino acids (aspartic and glutamic acids, phenylalanine,

448 methionine, asparagine and tryptophan) and higher frequencies of leucine (L) and proline

449 (P) residues. Overall, ereMAPs had therefore a less balanced (i.e., more skewed) amino

450 acid composition. Furthermore, analysis of amino acid usage at individual MAP positions

451 revealed that, relative to human MAPs, some residues were specifically enriched in ERE

452 and viral MAPs, such as arginine (R) in P5 of 8 amino acid-long MAPs (Figure S5). We

453 therefore aligned ereMAPs sequences to the viral MAPs dataset using BLAST and

454 calculated the average percentage of identity between ereMAPs and viral MAPs. We then

455 compared this result with a bootstrap distribution (10,000 iterations) of randomly selected

456 canonical MAPs that were also aligned to the viral MAPs dataset (Fig. 6B). This analysis

457 revealed that ereMAPs had a significantly higher percentage of identity with viral MAPs

458 than all 10,000 randomly selected sets of canonical MAPs. Hence, ereMAPs clearly retain

459 features that reflect their viral origin.

460

461 **Discussion**

462 Hundreds of scientific articles have alluded to the potential implication of EREs in various

463 human diseases, particularly cancer and autoimmunity (2, 55-60). We therefore felt

464 compelled to draw the global landscape of ERE expression in human somatic cells. We

465 hope that this atlas will serve as a reference in further studies on EREs in various

466 physiological and pathological conditions. One salient point emerging from this atlas is

467 that ERE expression in somatic tissues is more pervasive and heterogeneous than

468 anticipated. All tissues express EREs but the breadth and magnitude of ERE expression are

469 very heterogeneous from one tissue to another. Thus, we identified 124 ERE families

470 expressed in a tissue-restricted fashion, most of which were LTR elements. LTRs can act

471 as promoters and enhancers to stimulate gene expression (17, 19), and some LTR families

472 are tissue-specifically enriched in intronic enhancer regions containing transcription factor

473 binding sites (61). Our work therefore suggests that EREs, and more particularly LTRs,

474 may regulate gene expression in a wide range of somatic tissues. In future experiments,

475 single cell analyses might unveil a further level of heterogeneity that we could not capture

476 by global tissue expression profiling. It was previously reported that EREs were expressed

477 at high levels in two MHC I-deficient cell types: ESCs and testis (62, 63). That similar

478 levels of expression were found in mTECs for three major groups of EREs (LINE, SINE

479 and LTR) (Fig. 1) is remarkable and raises fundamental questions as to the mechanism and

480 role of ERE expression in mTECs. The key role of mTECs is to induce central immune

481 tolerance to a vast repertoire of self-peptides displayed by somatic tissues (47, 64). Given

482 the large-scale expression of EREs in peripheral tissues highlighted in the present report,

483 we speculate that it may be important for gnathostomes to be tolerant to a wide array of

484 ERE-derived antigens. As a corollary, when EREs are overexpressed, for instance in cancer

485 cells (65, 66), only those that are not expressed in mTECs may be immunogenic. Induction

486 of tolerance to the multitude of self-peptides depends on the unique ability of mTECs to

487 promiscuously express thousands of otherwise tissue-specific genes (67, 68). Promiscuous

488     gene expression in mTECs is driven in part by *AIRE* and in part by other genes whose

489     identity is unresolved and may include *FEZF2* as well as genes involved in DNA

490     methylation, histone modification and RNA splicing (26, 47, 69-71). Our data clearly show

491     that the overexpression of numerous ERE families in mTECs is entirely AIRE-independent

492     (Fig. 3). This observation underscores the relevance of further studies on the mechanisms

493     of AIRE-independent promiscuous gene expression in mTECs.

494

495     A notable finding was that our MS analyses identified ereMAPs derived from LINEs (n =

496     48), SINEs (n = 29) and LTRs (n= 27). This means that these EREs are translated and

497     produce peptides that are adequately processed for presentation by MHC-I molecules. A

498     few ereMAPs have previously been identified in cancer cells (25, 59, 66). The presence of

499     ereMAPs on normal cells means that the mere identification of ereMAPs on cancer cells is

500     not sufficient to infer that these MAPs are cancer-specific nor immunogenic. Nevertheless,

501     we have previously shown in mice that some ereMAPs are truly cancer-specific,

502     immunogenic and can elicit protective anti-tumor responses (25). Furthermore, compelling

503     evidence has been reported that some LTRs can generate immunogenic ereMAPs in clear

504     cell renal cell carcinoma in humans (56). These studies coupled to our finding that

505     ereMAPs retain viral like features (Fig. 6) suggest that ereMAPs may represent particularly

506     attractive targets for the development of cancer vaccines. In line with this, we must also

507     emphasize that the number of translated EREs is certainly superior to the number of

508     ereMAPs identified in our study: i) collectively our 16 B-LCLs expressed 39 MHC-I

509     allotypes out of the thousands that can be found in human populations (Table S5), and ii)

510     like canonical proteins (30), some translated EREs may not generate MAPs.

511

512  We anticipate that the biogenesis of ereMAPs in normal and neoplastic cells will be a fertile

513  field of investigation. First, several observations suggest that the landscape of ereMAPs is

514  highly diversified: i) the MAP repertoire is shaped by several cell type-specific variations

515  in gene expression (72), and ii) ERE transcription is highly heterogeneous among various

516  cell types (Fig. 1) and can be drastically affected by neoplastic transformation (73). The

517  processing of ereMAPs is also intriguing. Indeed, following their integration in human

518  genomes, EREs have undergone several rounds of mutation and truncation and very few

519  have previously been shown to be translated (2, 74). Because ERE sequences are

520  degenerate, they are not expected to yield stable polypeptides. However, MAPs

521  preferentially derive from rapidly degraded unstable peptides, commonly referred to as

522  defective ribosomal products (51). We therefore hypothesize that for most EREs,

523  translation may yield ereMAPs but not stable long-lived proteins. In other words, the

524  products of ERE translation may be detectable only in the immunopeptidome and not in

525  the proteome.

526

527  **Conclusions**

528  In summary, transcriptomic analysis demonstrated that ERE expression is heterogeneous

529  in healthy human tissues, with a higher expression in mTECs, ESCs and testis than in other

530  tissues. mTECs are the sole normal human cells that express high levels of both EREs and

531  MHC-I molecules. In mutant mice, we report that the exceptional expression of EREs in

532  mTECs is AIRE-independent. We also identified ERE families expressed in a tissue-

533  restricted manner, revealing that most healthy human tissues have a unique ERE signature.

534 MS analyses of 16 B-LCL samples enabled the identification of 104 non-redundant

535 ereMAPs, showing that EREs contribute to the immunopeptidome of healthy cells.

536 Interestingly, sharing of ereMAPs by multiple B-LCL samples was observed, and

537 ereMAPs' coding sequences are expressed at similar levels in other somatic tissues,

538 suggesting that ereMAPs could also be presented by other cell types. Finally, we found that

539 ereMAPs bear strong homology to viral MAPs and therefore have the potential to be

540 particularly immunogenic.

541

542 **Abbreviations**

543 B-LCL: B-lymphoblastoid cell line; ERE: Endogenous Retroelements; ereMAP: ERE-

544 derived MAP; ESC: Embryonic stem cells; FDR: False discovery rate; GTEx: Genotype-

545 Tissue Expression project; LINE: Long interspersed nuclear element; LTR: Long terminal

546 repeat; MCS : MAP-coding sequence; MAP: MHC I-associated peptide; mTEC: medullary

547 thymic epithelial cells; MS: Mass spectrometry; SINE: Short interspersed nuclear element;

548 TPM: transcripts per million; TRE: Tissue-restricted ERE; TRG: Tissue-restricted gene;

549 WT: Wild-type; KZFP: KRAB Zinc Finger Protein

550

551 # Declarations

552 **Ethics approval and consent to participate**

553 The study of MHC-associated peptides on human lymphoid cells was approved by the

554 Comité d'Éthique de la Recherche de l'Hôpital Maisonneuve-Rosemont (Permit Number

555 CÉR 2018-1396).

556

**Consent for publication**

Not applicable.


**Availability of data and materials**

<mark>XXXXXXXX</mark>


**Competing interests**

The authors declare that they have no competing interests.

**Authors' contributions**

JDL, KV and CP designed the study. LH and CC digested the thymic samples, isolated the mTECs and did the RNA extraction. JDL, AT, GE, PG and JPL contributed to the bioinformatic analyses. CD and EB did the PEAKS database searches and the MS/MS spectra validation. JDL and CP wrote the manuscript. All authors read and approved the final manuscript.

587

## References

589    1.      Dewannieux M, Heidmann T. Endogenous retroviruses: acquisition, amplification
590    and taming of genome invaders. Curr Opin Virol. 2013;3(6):646-56.
591    2.      Kassiotis G, Stoye JP. Immune responses to endogenous retroelements: taking the
592    bad with the good. Nat Rev Immunol. 2016;16(4):207-19.
593    3.      Sverdlov ED. Perpetually mobile footprints of ancient infections in human genome.
594    FEBS Lett. 1998;428(1-2):1-6.
595    4.      de Koning AP, Gu W, Castoe TA, Batzer MA, Pollock DD. Repetitive elements may
596    comprise over two-thirds of the human genome. PLoS Genet. 2011;7(12):e1002384.
597    5.      Tristem M. Identification and characterization of novel human endogenous
598    retrovirus families by phylogenetic screening of the human genome mapping project
599    database. J Virol. 2000;74(8):3715-30.
600    6.      Vargiu L, Rodriguez-Tome P, Sperber GO, Cadeddu M, Grandi N, Blikstad V, et al.
601    Classification and characterization of human endogenous retroviruses; mosaic forms are
602    common. Retrovirology. 2016;13:7.
603    7.      Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, et al. Initial
604    sequencing and analysis of the human genome. Nature. 2001;409(6822):860-921.
605    8.      Treangen TJ, Salzberg SL. Repetitive DNA and next-generation sequencing:
606    computational challenges and solutions. Nat Rev Genet. 2011;13(1):36-46.
607    9.      Bourque G, Burns KH, Gehring M, Gorbunova V, Seluanov A, Hammell M, et al. Ten
608    things you should know about transposable elements. Genome Biol. 2018;19(1):199.
609    10.     Argueso JL, Westmoreland J, Mieczkowski PA, Gawel M, Petes TD, Resnick MA.
610    Double-strand breaks associated with repetitive DNA can reshape the genome. Proc Natl
611    Acad Sci U S A. 2008;105(33):11845-50.
612    11.     Mills RE, Bennett EA, Iskow RC, Devine SE. Which transposable elements are active
613    in the human genome? Trends Genet. 2007;23(4):183-91.

614  12.     Deniz O, Frost JM, Branco MR. Regulation of transposable elements by DNA
615  modifications. Nat Rev Genet. 2019;20(7):417-31.
616  13.     Imbeault M, Helleboid PY, Trono D. KRAB zinc-finger proteins contribute to the
617  evolution of gene regulatory networks. Nature. 2017;543(7646):550-4.
618  14.     Bruno M, Mahgoub M, Macfarlan TS. The Arms Race Between KRAB-Zinc Finger
619  Proteins and Endogenous Retroelements and Its Impact on Mammals. Annu Rev Genet.
620  2019;53:393-416.
621  15.     Robbez-Masson L, Tie CHC, Conde L, Tunbak H, Husovsky C, Tchasovnikarova IA,
622  et al. The HUSH complex cooperates with TRIM28 to repress young retrotransposons and
623  new genes. Genome Res. 2018;28(6):836-45.
624  16.     Sheng W, LaFleur MW, Nguyen TH, Chen S, Chakravarthy A, Conway JR, et al. LSD1
625  Ablation Stimulates Anti-tumor Immunity and Enables Checkpoint Blockade. Cell.
626  2018;174(3):549-63 e19.
627  17.     Chuong EB, Elde NC, Feschotte C. Regulatory evolution of innate immunity through
628  co-option of endogenous retroviruses. Science. 2016;351(6277):1083-7.
629  18.     Beyer U, Moll-Rocek J, Moll UM, Dobbelstein M. Endogenous retrovirus drives
630  hitherto unknown proapoptotic p63 isoforms in the male germ line of humans and great
631  apes. Proc Natl Acad Sci U S A. 2011;108(9):3624-9.
632  19.     Fort A, Hashimoto K, Yamada D, Salimullah M, Keya CA, Saxena A, et al. Deep
633  transcriptome profiling of mammalian stem cells supports a regulatory role for
634  retrotransposons in pluripotency maintenance. Nat Genet. 2014;46(6):558-66.
635  20.     Percharde M, Lin CJ, Yin Y, Guan J, Peixoto GA, Bulut-Karslioglu A, et al. A LINE1-
636  Nucleolin Partnership Regulates Early Development and ESC Identity. Cell.
637  2018;174(2):391-405 e19.
638  21.     Mi S, Lee X, Li X, Veldman GM, Finnerty H, Racie L, et al. Syncytin is a captive
639  retroviral envelope protein involved in human placental morphogenesis. Nature.
640  2000;403(6771):785-9.
641  22.     Blaise S, de Parseval N, Benit L, Heidmann T. Genomewide screening for fusogenic
642  human endogenous retrovirus envelopes identifies syncytin 2, a gene conserved on
643  primate evolution. Proc Natl Acad Sci U S A. 2003;100(22):13013-8.
644  23.     Consortium GT. The Genotype-Tissue Expression (GTEx) project. Nat Genet.
645  2013;45(6):580-5.
646  24.     Lister R, Pelizzola M, Dowen RH, Hawkins RD, Hon G, Tonti-Filippini J, et al. Human
647  DNA methylomes at base resolution show widespread epigenomic differences. Nature.
648  2009;462(7271):315-22.
649  25.     Laumont CM, Vincent K, Hesnard L, Audemard E, Bonneil E, Laverdure JP, et al.
650  Noncoding regions are the main source of targetable tumor-specific antigens. Sci Transl
651  Med. 2018;10(470).
652  26.     St-Pierre C, Trofimov A, Brochu S, Lemieux S, Perreault C. Differential Features of
653  AIRE-Induced and AIRE-Independent Promiscuous Gene Expression in Thymic Epithelial
654  Cells. J Immunol. 2015;195(2):498-506.
655  27.     Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina
656  sequence data. Bioinformatics. 2014;30(15):2114-20.

657    28.    Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq
658    quantification. Nat Biotechnol. 2016;34(5):525-7.
659    29.    Yanai I, Benjamin H, Shmoish M, Chalifa-Caspi V, Shklar M, Ophir R, et al. Genome-
660    wide midrange transcription profiles reveal expression level relationships in human tissue
661    specification. Bioinformatics. 2005;21(5):650-9.
662    30.    Pearson H, Daouda T, Granados DP, Durette C, Bonneil E, Courcelles M, et al. MHC
663    class I-associated peptides derive from selective regions of the human genome. J Clin
664    Invest. 2016;126(12):4690-701.
665    31.    Granados DP, Sriranganadane D, Daouda T, Zieger A, Laumont CM, Caron-Lizotte
666    O, et al. Impact of genomic polymorphisms on the repertoire of human MHC class I-
667    associated peptides. Nat Commun. 2014;5:3600.
668    32.    Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast
669    universal RNA-seq aligner. Bioinformatics. 2013;29(1):15-21.
670    33.    Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic
671    features. Bioinformatics. 2010;26(6):841-2.
672    34.    Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence
673    Alignment/Map format and SAMtools. Bioinformatics. 2009;25(16):2078-9.
674    35.    Garrison E, Marth G. Haplotype-based variant detection from short-read
675    sequencing. arXiv e-prints [Internet]. 2012 July 01, 2012. Available from:
676    https://ui.adsabs.harvard.edu/abs/2012arXiv1207.3907G.
677    36.    Daouda T, Perreault C, Lemieux S. pyGeno: A Python package for precision
678    medicine and proteogenomics. F1000Res. 2016;5:381.
679    37.    Andreatta M, Nielsen M. Gapped sequence alignment using artificial neural
680    networks: application to the MHC class I system. Bioinformatics. 2016;32(4):511-7.
681    38.    Jurtz V, Paul S, Andreatta M, Marcatili P, Peters B, Nielsen M. NetMHCpan-4.0:
682    Improved Peptide-MHC Class I Interaction Predictions Integrating Eluted Ligand and
683    Peptide Binding Affinity Data. J Immunol. 2017;199(9):3360-8.
684    39.    Marcais G, Kingsford C. A fast, lock-free approach for efficient parallel counting of
685    occurrences of k-mers. Bioinformatics. 2011;27(6):764-70.
686    40.    Robinson JT, Thorvaldsdottir H, Winckler W, Guttman M, Lander ES, Getz G, et al.
687    Integrative genomics viewer. Nat Biotechnol. 2011;29(1):24-6.
688    41.    Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search
689    tool. J Mol Biol. 1990;215(3):403-10.
690    42.    Granados DP, Yahyaoui W, Laumont CM, Daouda T, Muratore-Schroeder TL, Cote
691    C, et al. MHC I-associated peptides preferentially derive from transcripts bearing miRNA
692    response elements. Blood. 2012;119(26):e181-91.
693    43.    Fergusson JR, Morgan MD, Bruchard M, Huitema L, Heesters BA, van Unen V, et
694    al. Maturing Human CD127+ CCR7+ PDL1+ Dendritic Cells Express AIRE in the Absence of
695    Tissue Restricted Antigens. Front Immunol. 2018;9:2902.
696    44.    Boegel S, Lower M, Bukur T, Sorn P, Castle JC, Sahin U. HLA and proteasome
697    expression body map. BMC Med Genomics. 2018;11(1):36.
698    45.    Drukker M, Katz G, Urbach A, Schuldiner M, Markel G, Itskovitz-Eldor J, et al.
699    Characterization of the expression of MHC proteins in human embryonic stem cells. Proc
700    Natl Acad Sci U S A. 2002;99(15):9864-9.

701      46.      Klein L, Hinterberger M, Wirnsberger G, Kyewski B. Antigen presentation in the
702 thymus for positive selection and central tolerance induction. Nat Rev Immunol.
703 2009;9(12):833-44.

704      47.      Inglesfield S, Cosway EJ, Jenkinson WE, Anderson G. Rethinking Thymic Tolerance:
705 Lessons from Mice. Trends Immunol. 2019;40(4):279-91.

706      48.      Sacha JB, Kim IJ, Chen L, Ullah JH, Goodwin DA, Simmons HA, et al. Vaccination
707 with cancer- and HIV infection-associated endogenous retrotransposable elements is safe
708 and immunogenic. J Immunol. 2012;189(3):1467-79.

709      49.      Young GR, Ploquin MJ, Eksmond U, Wadwa M, Stoye JP, Kassiotis G. Negative
710 selection by an endogenous retrovirus promotes a higher-avidity CD4+ T cell response to
711 retroviral infection. PLoS Pathog. 2012;8(5):e1002709.

712      50.      Granados DP, Laumont CM, Thibault P, Perreault C. The nature of self for T cells-a
713 systems-level perspective. Curr Opin Immunol. 2015;34:1-8.

714      51.      Yewdell JW, Dersh D, Fahraeus R. Peptide Channeling: The Key to MHC Class I
715 Immunosurveillance? Trends Cell Biol. 2019;29(12):929-39.

716      52.      Chiappinelli KB, Strissel PL, Desrichard A, Li H, Henke C, Akman B, et al. Inhibiting
717 DNA Methylation Causes an Interferon Response in Cancer via dsRNA Including
718 Endogenous Retroviruses. Cell. 2015;162(5):974-86.

719      53.      Roulois D, Loo Yau H, Singhania R, Wang Y, Danesh A, Shen SY, et al. DNA-
720 Demethylating Agents Target Colorectal Cancer Cells by Inducing Viral Mimicry by
721 Endogenous Transcripts. Cell. 2015;162(5):961-73.

722      54.      Jung J, Lee S, Cho HS, Park K, Ryu JW, Jung M, et al. Bioinformatic analysis of
723 regulation of natural antisense transcripts by transposable elements in human mRNA.
724 Genomics. 2019;111(2):159-66.

725      55.      Attig J, Young GR, Stoye JP, Kassiotis G. Physiological and Pathological
726 Transcriptional Activation of Endogenous Retroelements Assessed by RNA-Sequencing of
727 B Lymphocytes. Front Microbiol. 2017;8:2489.

728      56.      Smith CC, Beckermann KE, Bortone DS, De Cubas AA, Bixby LM, Lee SJ, et al.
729 Endogenous retroviral signatures predict immunotherapy response in clear cell renal cell
730 carcinoma. J Clin Invest. 2018;128(11):4804-20.

731      57.      Treger RS, Pope SD, Kong Y, Tokuyama M, Taura M, Iwasaki A. The Lupus
732 Susceptibility Locus Sgp3 Encodes the Suppressor of Endogenous Retrovirus Expression
733 SNERV. Immunity. 2019;50(2):334-47 e9.

734      58.      De Cecco M, Ito T, Petrashen AP, Elias AE, Skvir NJ, Criscione SW, et al. L1 drives
735 IFN in senescent cells and promotes age-associated inflammation. Nature.
736 2019;566(7742):73-8.

737      59.      Attig J, Young GR, Hosie L, Perkins D, Encheva-Yokoya V, Stoye JP, et al. LTR
738 retroelement expansion of the human cancer transcriptome and immunopeptidome
739 revealed by de novo transcript assembly. Genome Res. 2019.

740      60.      Smith CC, Selitsky SR, Chai S, Armistead PM, Vincent BG, Serody JS. Alternative
741 tumour-specific antigens. Nat Rev Cancer. 2019.

742      61.      Trizzino M, Kapusta A, Brown CD. Transposable elements generate regulatory
743 novelty in a tissue-specific fashion. BMC Genomics. 2018;19(1):468.

744     62.     Gainetdinov I, Skvortsova Y, Kondratieva S, Funikov S, Azhikina T. Two modes of
745     targeting transposable elements by piRNA pathway in human testis. RNA.
746     2017;23(11):1614-25.
747     63.     Grow EJ, Flynn RA, Chavez SL, Bayless NL, Wossidlo M, Wesche DJ, et al. Intrinsic
748     retroviral reactivation in human preimplantation embryos and pluripotent cells. Nature.
749     2015;522(7555):221-5.
750     64.     Abramson J, Anderson G. Thymic Epithelial Cells. Annu Rev Immunol. 2017;35:85-
751     118.
752     65.     Rooney MS, Shukla SA, Wu CJ, Getz G, Hacohen N. Molecular and genetic
753     properties of tumors associated with local immune cytolytic activity. Cell. 2015;160(1-
754     2):48-61.
755     66.     Kong Y, Rose CM, Cass AA, Williams AG, Darwish M, Lianoglou S, et al.
756     Transposable element expression in tumors is associated with immune infiltration and
757     increased antigenicity. Nat Commun. 2019;10(1):5228.
758     67.     Klein L, Kyewski B, Allen PM, Hogquist KA. Positive and negative selection of the T
759     cell repertoire: what thymocytes see (and don't see). Nat Rev Immunol. 2014;14(6):377-
760     91.
761     68.     Sansom SN, Shikama-Dorn N, Zhanybekova S, Nusspaumer G, Macaulay IC,
762     Deadman ME, et al. Population and single-cell genomics reveal the Aire dependency,
763     relief from Polycomb silencing, and distribution of self-antigen expression in thymic
764     epithelia. Genome Res. 2014;24(12):1918-31.
765     69.     Ucar O, Rattay K. Promiscuous Gene Expression in the Thymus: A Matter of
766     Epigenetics, miRNA, and More? Front Immunol. 2015;6:93.
767     70.     Takaba H, Morishita Y, Tomofuji Y, Danks L, Nitta T, Komatsu N, et al. Fezf2
768     Orchestrates a Thymic Program of Self-Antigen Expression for Immune Tolerance. Cell.
769     2015;163(4):975-87.
770     71.     Danan-Gotthold M, Guyon C, Giraud M, Levanon EY, Abramson J. Extensive RNA
771     editing and splicing increase immune self-representation diversity in medullary thymic
772     epithelial cells. Genome Biol. 2016;17(1):219.
773     72.     Caron E, Vincent K, Fortier MH, Laverdure JP, Bramoulle A, Hardy MP, et al. The
774     MHC I immunopeptidome conveys to the cell surface an integrative view of cellular
775     regulation. Mol Syst Biol. 2011;7:533.
776     73.     Chong C, Müller M, Pak H, Harnett D, Huber F, Grun D, et al. Integrated
777     proteogenomic deep sequencing and analytics accurately identify non-canonical peptides
778     in tumor immunopeptidomes. bioRxiv. 2019.
779     74.     Bonnaud B, Bouton O, Oriol G, Cheynet V, Duret L, Mallet F. Evidence of selection
780     on the domesticated ERVWE1 env retroviral element involved in placentation. Mol Biol
781     Evol. 2004;21(10):1895-901.
782

783     **Figure legends**

784

785    **Fig. 1**. Expression profiling of endogenous retroelements in 30 healthy human tissues and

786    2 cell types. Hierarchical clustering of tissues based on the expression levels of the 809

787    ERE families sorted in LINE, LTR and SINE elements. For each tissue, mean expression

788    of ERE families was computed among available samples. Row Z-scores were then

789    determined for each ERE family across tissues.

790

791    **Fig. 2**. Tissue specificity of ERE expression in healthy human tissues. Tissue-specificity

792    indexes were computed for ERE families as well as annotated genes. (A) Barplots showing

793    the number of TRGs and TREs for each of the 32 healthy human tissues analyzed. (B) Pie

794    charts depicting the proportions of the 809 ERE families (left panel) or TREs (right panel)

795    belonging to the LINE, LTR and SINE groups (Chi-squared test, *$P{\leq}0.05$). (C) Histogram

796    showing the number of tissues in which each identified TRGs and TREs are overexpressed.

797

798    **Fig. 3**. ERE expression is independent of AIRE in mouse mTECs. (A) Boxplot showing

799    the expression levels of constitutively expressed genes, AIRE-dependent TRGs, AIRE-

800    independent TRGs (lists of genes based on St-Pierre *et al* (26)) as well as ERE families in

801    wild-type (n=3) and AIRE knock-out (n=3) mice. (B) Heatmap depicting the expression

802    levels of ERE families in each replicate of wild-type and AIRE knock-out murine mTECs.

803    A Mann-Whitney test was used for statistical analysis in both panels, n.s. not significant

804    ($P{>}0.05$), ***$P{\leq}0.001$.

805

806    **Fig. 4**. ERE sequences are translated and contribute to the immunopeptidome of B-LCLs.

807    (A) Schematic depicting how the personalized proteome of each B-LCL sample was

808    generated. The personalized proteome was generated by combining the ERE and the

809    canonical proteomes and then used to identify MAPs by MS. MAPs were annotated to keep

810    only ereMAPs. (B, C) Barplots showing the number of ereMAPs identified in B-LCL

811    samples separated by (B) individual samples analyzed and (C) according to the three main

812    groups of EREs.

813

814    **Fig. 5**. Sense transcription of intronic EREs is the main source of ereMAPs. (A) Boxplot

815    showing the mean expression levels ($\log_{10}(\text{TPM} + 1)$) of ERE families that are source or

816    non-source of ereMAPs in B-LCLs (Mann-Whitney test, ***$P{\leq}0.001$). (B) Barplot

817    showing the number of ereMAPs generated by sense or antisense transcription of ERE

818    sequences. (C) Stacked barplot depicting the proportions of LINE, LTR and SINE groups

819    in the genome, transcriptome and immunopeptidome. Statistical significance was

820    computed with a chi-squared test (**$P{\leq}0.01$). (D) Pie charts depicting the percentages of

821    all ERE sequences (left) and of ereMAPs-coding sequences (right) that are localized in

822    intergenic regions, introns or coding sequences (Chi-squared test, ***$P{\leq}0.001$). (E)

823    Scatterplot showing the Spearman correlation between the number of ereMAPs generated

824    by each ERE family and the number of copies of the ERE family's sequence in the human

825    genome based on RepeatMasker annotations.

826

827    **Fig. 6**. Endogenous retroelements retained sequence homology with viruses. (A) Barplot

828    showing the frequencies of each amino acid in ereMAPs, viral MAPs and human canonical

829    MAPs. Abbreviations for amino acids: Y, Tyrosine; W, Tryptophan; V, Valine; T,

830    Threonine; S, Serine; R, Arginine; Q, Glutamine; P, Proline; N, Asparagine; M,

831    Methionine; L, Leucine; K, Lysine; I, Isoleucine; H, Histidine; G, Glycine; F,

832    Phenylalanine; E, Glutamic Acid; D, Aspartic Acid; C, Cysteine; A, Alanine. (B) Human

833    canonical MAPs and ereMAPs were aligned to a database of viral peptides using BLAST,

834    and the percentage of identity of their sequences with viral peptides was computed. The

835    red line represents the average percentage of identity of ereMAPs with viral MAPs. A

836    bootstrap procedure was used to calculate the percentage of identity of 10,000 sets of 104

837    randomly selected human canonical MAPs with viral MAPs. P-value was calculated as the

838    number of times the bootstrap distribution had a higher percentage of identity with viral

839    MAPs than ereMAPs ($P<0.0001$).