

1

1 Restriction Enzyme Based Enriched L1Hs sequencing (REBELseq)

2

3 Benjamin C. Reiner, PhD<sup>1,†,\*</sup>; Glenn A. Doyle, PhD<sup>1,†</sup>; Andrew E. Weller, MD<sup>1</sup>; Rachel N.  
4 Levinson, BA<sup>1</sup>; Esin Namoglu, BS<sup>1</sup>; Alicia Pigeon, BS<sup>1</sup>; Gabriella Arauco-Shapiro, MS<sup>1</sup>; Emilie  
5 Dávila Perea<sup>1</sup>; Cyndi Shannon Weickert, PhD<sup>2</sup>; Gustavo Turecki, MD, PhD<sup>3</sup>; Deborah C. Mash,  
6 PhD<sup>4</sup>; Richard C. Crist, PhD<sup>1</sup> and Wade H. Berrettini, MD, PhD<sup>1</sup>

7

8 <sup>1</sup> Center for Neurobiology and Behavior, Department of Psychiatry, Perelman School of  
9 Medicine, University of Pennsylvania; <sup>2</sup> Department of Neuroscience and Physiology, Upstate  
10 Medical University; <sup>3</sup> McGill Group for Suicide Studies, Douglas Mental Health University  
11 Institute, McGill University; <sup>4</sup> Dr. Kiran C. Patel College of Allopathic Medicine, NOVA  
12 Southeastern University, <sup>†</sup> These authors contributed equally to this work, <sup>\*</sup> corresponding  
13 author: [bcreiner@penmedicine.upenn.edu](mailto:bcreiner@penmedicine.upenn.edu)

14

15 Corresponding author information

16 Benjamin C. Reiner, PhD

17 125 S. 31<sup>st</sup> St., room 2208-1

18 Philadelphia, PA 19104

19 Office: (215)-573-0278

20 Fax: (215)-573-2041

21 Email: [bcreiner@penmedicine.upenn.edu](mailto:bcreiner@penmedicine.upenn.edu)

22

23

## 24 **Abstract**

25 Long interspersed element-1 retrotransposons (LINE-1 or L1) are ~6 kb mobile DNA  
26 elements implicated in the origins of many Mendelian and complex diseases. The actively  
27 retrotransposing L1s are mostly limited to the L1 human specific Ta subfamily. In this  
28 manuscript, we present REBELseq as a method for the construction of differentially amplified  
29 next-generation sequencing libraries and bioinformatic identification of Ta subfamily long  
30 interspersed element-1 human specific elements. REBELseq was performed on DNA isolated  
31 from NeuN+ neuronal nuclei from postmortem brain samples of 177 individuals and empirically-  
32 driven bioinformatic and experimental cutoffs were established. REBELseq reliably identified  
33 both known and novel Ta subfamily L1 insertions distributed throughout the genome.  
34 Differences in the proportion of individuals possessing a given reference or non-reference  
35 retrotransposon insertion were identified. We conclude that REBELseq is an unbiased, whole  
36 genome approach to the amplification and detection of Ta subfamily L1 retrotransposons.

37

## 38 **Introduction**

39 Retrotransposons are a class of mobile DNA elements capable of replicating and  
40 inserting copies of themselves elsewhere in the genome using an RNA intermediate (1). Long  
41 interspersed element-1 (L1), a type of retrotransposon that is currently active in the human  
42 genome, is estimated to constitute approximately 17% of the human genome (2). Despite their  
43 abundance, most L1s are truncated or mutated to the point of no longer being able to  
44 retrotranspose (3). There are ~100 L1s in each human genome that are competent, meaning  
45 they are full length and capable of replicating, the vast majority of which belong to the L1 human  
46 specific (L1Hs) Ta subfamily (4). Alternatively, in agreement with Kazazian *et al.* (5), Beck *et al.*  
47 described two competent L1s of the pre-Ta subfamily (2). A competent L1Hs is ~6 kb in length  
48 and contains a promoter, 5' and 3' untranslated regions, and two open reading frames (ORF):

49 ORF1, encoding an RNA binding protein (6) , and ORF2, encoding a fusion protein that acts as  
50 both a reverse transcriptase (7) and endonuclease (8).

51 Whereas competent retrotransposons are mostly limited to full length L1Hs of the Ta  
52 subfamily, the remnants of vast numbers of retrotransposons, from both evolutionarily older and  
53 Ta subfamily L1Hs can be found throughout the human genome. Notably, Ta subfamily L1Hs  
54 contain identifying nucleotides in the 3' UTR that can be utilized for their differential amplification  
55 by polymerase chain reaction (PCR) (4). Retrotransposition of competent L1Hs elements  
56 frequently results in 5' truncation of the new L1Hs element insertions. Therefore, to reliably  
57 amplify and detect full length and truncated L1Hs elements, including germline polymorphic and  
58 individual somatic mutations, from the Ta subfamily, PCR primers for amplification must be  
59 targeted near the 3' end of the L1Hs Ta subfamily sequence.

60 In this manuscript, we present REBELseq (Restriction Enzyme Based Enriched L1Hs  
61 sequencing), a scalable technique for the differential amplification of both full length and 5'  
62 truncated L1s. As described below, we specifically target the Ta subfamily of L1Hs and  
63 summarize the results of the application of REBELseq to DNA samples from 177 individuals.

## 64 **Materials and Methods**

### 65 Brain samples and purification of genomic DNA from NeuN+ nuclei

66 177 fresh frozen human postmortem prefrontal cortex brain samples (Brodmann's Area  
67 9 or 10) were provided by the Douglas-Bell Canada Brain Bank at McGill University, the Human  
68 Brain and Spinal Fluid Resource Center at UCLA, the Sydney Brain Bank at Neuroscience  
69 Research Australia, or the University of Miami Brain Endowment Bank. All work was approved  
70 by the University of Pennsylvania Institutional Review Board as category four exempt human  
71 subject research. All reported age, sex and ethnicity data are based on associated medical  
72 records. Sample set contained 138 males (Age  $45.6 \pm 15.5$ ; 93 European, 25 African, 15

73 Hispanic, 2 Asian and 3 of unknown ethnic origin) and 39 females (Age  $55.5 \pm 19.2$ ; 33  
74 European, 1 African, 3 Hispanic and 2 of unknown ethnic origin). See Supplemental Methods for  
75 full details.

76 Brain samples were homogenized and co-stained with DAPI and an  $\alpha$ -NeuN-AF488  
77 antibody (Millipore # MAB377X) using a modification of a previously described method (9), and  
78 NeuN+ and NeuN- nuclei were isolated using a FACSAria II (Beckman-Coulter). Isolated nuclei  
79 were lysed overnight in digestion buffer in a  $56^\circ\text{C}$  water bath, and the following morning  
80 genomic DNA (gDNA) was isolated using the Zymo Research genomic DNA clean and  
81 concentrator kit (#D4011). The concentration of gDNA was quantified using a Qubit 3  
82 fluorometer and high-sensitivity double-stranded DNA assay (ThermoFisher #Q32854). See  
83 Supplemental Methods for full details.

#### 84 Ta subfamily L1Hs-enriched library construction and sequencing

85 Ta subfamily L1Hs-enriched next generation sequencing libraries were constructed  
86 utilizing the REBELseq technique (Fig. 1). 33 ng of gDNA extracted from NeuN+ neuronal nuclei  
87 was digested with HaeIII, in the presence of shrimp alkaline phosphatase. A single primer  
88 extension utilizing a 3' diagnostic 'A' nucleotide (L1HsACA primer, see Supplemental Oligomers  
89 for all primer sequences; all oligos from IDT, Iowa, USA), a more stringent variation of a primer  
90 originally designed by Ewing and Kazazian (10), extends only Ta subfamily L1 3' UTR sequence  
91 (4) and the adjacent downstream genomic DNA, leaving a terminal 3' A-overhang. The A-  
92 overhang products were ligated to a double stranded T-linker molecule, originally designed for a  
93 different technique (11), and the ligated products were amplified using the L1HsACA primer and  
94 a T-linker specific primer to enrich the number of copies of each unique Ta subfamily L1Hs  
95 insertion (Primary PCR). The Primary PCR product was diluted and used as a template for a  
96 hemi-nested Secondary PCR reaction using the T-linker primer and the Seq2-L1HsG primer.  
97 The purpose of the hemi-nested secondary PCR reactions is three-fold: to reduce the length of

98 3' L1 sequence carried forward, to add an Illumina sequencing primer, and, most importantly, to  
99 use the L1Hs diagnostic 3' G nucleotide of the Seq2-L1HsG primer to further enrich the L1Hs  
100 Ta subfamily. Secondary PCR products were cleaned and size selected using KAPA Pure  
101 Beads (Roche #KK8002) according to the manufacturer's protocol for a 0.55X-0.75X double-  
102 sided size selection, generating a purified library of 200-1,000 bp fragments. A Tertiary PCR  
103 using 5' overhang primers then added Illumina flow cell adapters, containing a single  
104 multiplexing index, to both ends of the amplicon. Sequencing libraries then underwent a final  
105 cleanup and size selection using KAPA Pure Beads 0.55X-0.75X double-sided size selection  
106 and the average fragment size of the library was determined using the 2100 Bioanalyzer and  
107 high sensitivity DNA kit (Agilent #5067-4626). The quantity of sequenceable library for each  
108 sample was determined using the KAPA Library Quantification Kit (Roche #KK4835) on the  
109 7900HT real-time PCR (Applied Biosystems). Barcoded samples were pooled in groups of 6,  
110 and pooled libraries were sequenced by the University of Pennsylvania Next-Generation  
111 Sequencing Core, using one pool per lane, on an Illumina HiSeq 4000 utilizing 150 bp paired-  
112 end sequencing. See Supplemental Methods for full details.

### 113 Bioinformatics

114 Demultiplexed sequencing data was cleaned and quality trimmed to a Phred quality  
115 score of  $Q \geq 20$  using BBTools bbduk (<https://sourceforge.net/projects/bbmap>), and trimmed  
116 read 1 data was aligned to the hg19 build of the human genome with Bowtie2-2.1.0 (12) using  
117 end-to-end, very-sensitive alignment to generate a sequence alignment map (SAM) file for each  
118 individual. ~10% of unaligned reads contained a poly(T) sequence after a genomic sequence  
119 which made it unalignable. Assuming these poly(T) stretches likely corresponded to the 3'  
120 poly(A) tail of the L1Hs sequence, they were trimmed, and the trimmed reads were aligned to  
121 hg19 with Bowtie2-2.1.0 using end-to-end, very-sensitive alignment parameters to generate a  
122 second SAM alignment file for each individual. SAM files were converted to BAM format, with

123 reads having alignment MapQ scores < 30 being discarded, using samtools-0.1.19 (13). Using  
124 samtools-0.1.19, the two BAM files for each individual were combined. The resulting single BAM  
125 files from all individuals were then merged using RG tags, a tag that allows for the combination  
126 of reads from multiple BAM files, while still retaining the source (i.e. the individual) from which  
127 each aligned read was generated, and the merged file was sorted and indexed. The merged  
128 BAM file was stripped into BAM files for individual chromosomes using samtools-0.1.19, and the  
129 chromosome specific BAM files were used as input for the custom python script  
130 REBELseq\_v1.0. REBELseq\_v1.0 utilizes the sorted and indexed reads in the chromosome  
131 specific BAM files to generate peaks of overlapping reads within a 150 base pair sliding window,  
132 and does so with respect to DNA strand. Each peak corresponds to a putative L1  
133 retrotransposon insertion, for which the peak's genomic coordinates, number of unique reads  
134 per peak, number of reads per individual sample and average read alignment quality (mean  
135 MapQ) are determined. The REBELseq\_v1.0 output file was then further annotated using the  
136 custom python script REBELannotate\_v1.0 for L1Hs annotated in hg19 repeat masker and  
137 L1Hs identified in the 1000 genomes data (14). REBELannotate\_v1.0 utilizes a browser  
138 extensible data (BED) formatted file to annotate genomic features of interest that overlap with or  
139 occur within 500 base pairs downstream, with respect to strand, of the peak being annotated.  
140 The REBELseq\_v1.0 and REBELannotate\_v1.0 custom scripts are based on work originally  
141 described by Ewing and Kazazian (10), and all python scripts and necessary reference files  
142 discussed in this manuscript are available online (<https://github.com/BenReiner/REBELseq>).  
143 The raw sequencing data that was analyzed to support the findings of this manuscript are  
144 available from the corresponding author upon request. See Supplemental Methods for full  
145 details.

146

147

## 148 PCR Validations

149 Primers used in PCR experiments for method validation were designed using Primer3-2.2.3 in a  
150 Perl script (makeprimers.pl), originally written by Adam Ewing (10), with an optimal T<sub>m</sub> setting of  
151 58°C (minimum 56°C and max 63°C), an optimal primer length of 24 nucleotides (minimum 21  
152 and max 27) and the GC-clamp option set to 1. 25 µL reactions were constructed using 1x Go-  
153 Taq colorless hotstart master mix (#M5133, Promega), 1 ng of gDNA for both an individual  
154 predicted to have a particular L1Hs insertion or an individual predicted not to have the insertion,  
155 and either 0.2 µM of the filled site and empty site (empty allele detection) or the filled site and  
156 L1HsACA primer (filled allele detection). Samples were thermally cycled as 95 °C 2:30 min, then  
157 10 cycles of touchdown PCR at 95 °C 0:30 min, 69-60 °C 0:30 min, 72 °C 1:30 min, then 25  
158 cycles of 95 °C 0:30 min, 60 °C 0:30 min, 72 °C 1:30 min, then 72 °C 5:00 min, 4 °C hold. Go-  
159 Taq green flexi buffer (#M8911, Promega) was added to all samples after amplification, and  
160 samples were separated by 1.2% gel electrophoresis and visualized on a GelDoc XR+ (Bio-  
161 Rad).

## 162 Statistics

163 Data are presented as mean ± standard deviation.

164

## 165 **Results**

### 166 REBELseq identifies L1Hs retrotransposon insertions

167 We performed REBELseq using DNA isolated from postmortem cortex NeuN+ nuclei  
168 samples of 177 individuals and identified a total of 157,178 independent putative Ta subfamily  
169 L1Hs insertions, with 1,050 annotating to hg19 known reference L1Hs elements (ref L1Hs) and  
170 156,128 not annotating to hg19 known reference L1Hs (non-ref L1Hs). The ref L1Hs we

171 detected represented ~68% of all L1Hs annotated in hg19 repeat masker (1,050/1,544; see  
172 Discussion) and were distributed across all chromosomes (Fig. 2a). Of the non-ref L1Hs, 432  
173 annotated to known L1s in the 1000 genomes data (known non-ref L1Hs), while the other  
174 155,696 are putative novel non-ref L1Hs. The ref L1Hs had an average mean MapQ of  $38.82 \pm$   
175 4.10. Having high confidence that the detected ref L1Hs were real, we used the ref L1Hs  
176 average mean MapQ minus two standard deviations ( $\text{MapQ} \geq 30.62$ ) as a bioinformatic cutoff of  
177 our data, after which we retained 974 ref L1Hs (92.76% of the 1,050 ref L1s), 430 known non-  
178 ref L1Hs (99.54% of 432) and 127,976 putative novel non-ref L1Hs (82.20% of 155,696).

179 In the remaining known (ref L1Hs and known non-ref L1Hs) and putative novel L1Hs, we  
180 next examined the average number of sequencing reads per person for a given L1Hs insertion  
181 (Fig. 2b), and a clear difference between the known and novel L1Hs emerged, with 92.6% of  
182 putative novel L1Hs having average reads  $< 10$ , and 87.7% of known L1Hs having average  
183 reads  $\geq 10$  and 74.2% having average reads  $\geq 100$ . Having confidence in the known L1Hs pool,  
184 an average of  $\geq 100$  sequencing reads per person was used as a bioinformatic cutoff, which  
185 retained 699 ref L1Hs (66.57%), 332 known non-ref L1Hs (76.85%) and 1,842 putative novel  
186 non-ref L1Hs (1.18%).

### 187 L1Hs insertion validation by PCR

188 Having established bioinformatic cutoffs for our data, we next sought to experimentally  
189 determine the proportion of bioinformatically-detected putative L1Hs insertions that were real  
190 using PCR of the 3' insertion junction of the allele containing the insertion (filled allele) and the  
191 corresponding genomic region on the 'empty' allele (see Fig. 3 for example image). Having  
192 confidence that the known L1Hs insertions (ref L1Hs and known non-ref L1Hs) are real because  
193 they were previously reported, we focused on the validation of the putative novel non-ref L1Hs  
194 insertions predicted by REBELseq. The percentages of putative novel non-ref L1Hs insertions  
195 that could be validated by average read number bin were experimentally determined, and these



196 data were used to calculate the number of detected putative novel non-ref L1Hs insertions that  
197 we would expect to be true positives in each read bin. Using the number of putative novel non-  
198 ref L1Hs insertions predicted to be true positives and the numbers of known L1Hs detected, we  
199 calculated the probability of validating any detected L1Hs insertion per read bin and the  
200 cumulative probability of validating any detected L1Hs insertions using the lower value of a read  
201 bin as a minimum cutoff for the data (Table 1).

### 202 Distribution of detected L1Hs

203 The genomic distribution of known and novel L1Hs, per 10 MB genomic window, of our  
204 highest quality L1Hs detection data was ascertained ( $\geq 1,000$  average reads, Fig. 4a), and the  
205 distribution of the number of individuals having a given known (Fig. 4b) or novel (Fig. 4c) L1Hs  
206 for these data was assessed.

### 207 **Discussion**

208 While other methods for preparation of L1-targeted next generation sequencing libraries  
209 exist (10, 11, 15-22), some are labor-intensive when scaled to a large number of individual  
210 genomes, because of reliance on gel purification (10, 11) or because they require preparation of  
211 multiple amplicon libraries per individual sample, due to the use of multiple hemi-degenerate  
212 primers (10, 19). Some methods rely on the random shearing of DNA (18, 22, 23), a process  
213 that requires specialized equipment and somewhat abundant DNA quantity, while others may  
214 have reduced specificity for the Ta subfamily of L1Hs, which constitutes the majority of actively  
215 replicating L1 retrotransposons in the human genome (16, 20, 21, 23). REBELseq was  
216 designed as a high throughput alternative method to reliably detect the 3' flanking region of Ta  
217 subfamily L1Hs elements with limited gDNA input. It should be noted that REBELseq was not  
218 intended to determine whether an L1Hs insertion is full length or truncated or if ORF1 and ORF2  
219 are intact, but merely to determine the presence or absence of the Ta subfamily L1Hs. After

220 identification of insertions of interest, these other aspects can be determined by long range PCR  
221 and Sanger sequencing.

222 Other methods to produce L1-enriched libraries employing restriction enzyme digests of  
223 gDNA as a starting point have been described (15). There is no set requirement for which  
224 restriction enzyme to use in REBELseq other than that it generates blunt ends and cuts 5' of the  
225 L1HsACA primer that identifies the 3' end of the Ta subfamily L1Hs elements. REBELseq  
226 begins with restriction enzyme digestion of gDNA with HaeIII. This restriction enzyme was  
227 chosen for the following reasons: 1) it creates blunt ends at 5'-GG|CC-3' sequences such that  
228 polishing of cleaved ends is unnecessary; 2) it cuts the human genome to an average size of  
229  $342 \pm 478$  base pairs (NEBcutter) which is similar to the range of fragments generated by  
230 sonication; 3) it does not cut within the 3' end of the L1Hs sequence targeted by our L1Hs Ta  
231 subfamily-specific primers. One possible concern of using a single restriction enzyme digestion  
232 in REBELseq library construction is the possibility that the restriction site for the enzyme occurs  
233 immediately downstream from a L1 insertion, thus preventing the single primer extension into  
234 the downstream gDNA (see Fig. 1). While this raises the possibility that using a single restriction  
235 enzyme digestion during REBELseq library construction may not detect all L1 insertions, it does  
236 not reduce the validity of L1 insertion that are detected. One possible solution to this concern  
237 would be to perform restriction enzyme digestions with two or more enzymes that meet the  
238 above criteria and pool the fragments before the single primer extension step.

239 Our method is specifically designed to leverage diagnostic nucleotides specific to the Ta  
240 subfamily of L1Hs elements for differential amplification. It can easily be adapted to detect the  
241 5' flanking region of full length L1Hs elements in the human genome, although diagnostic  
242 nucleotides are limited within the 5' end of L1Hs for the Ta subfamily. Reliable primers targeting  
243 the 5' end of Ta subfamily L1Hs and an enzyme other than HaeIII would be required for analysis  
244 of full length L1Hs elements, as HaeIII makes a single digestion in the L1Hs sequence near the  
245 5' end. REBELseq could also be utilized to differentially amplify and identify evolutionarily older

246 pre-Ta subfamily L1Hs insertions (by changing the last three nucleotides of the L1HsACA  
247 primer to 'ACG') or any type of L1Hs element (by eliminating the last nucleotide of the L1HsACA  
248 primer to end in 'AC'). Additionally, this method could be utilized to amplify and identify almost  
249 any other genomic element with a specific pair of nested primers (to replace L1HsACA and  
250 L1HsG).

251 The REBELseq technique is compatible with human gDNA purified from any tissue  
252 source, including gDNA purified from leukocytes in saliva and blood. In the application of  
253 REBELseq presented in this manuscript, 177 fresh frozen human postmortem prefrontal cortex  
254 brain samples were utilized for purification of gDNA from NeuN+ neuronal nuclei. Utilizing gDNA  
255 from NeuN+ neuronal nuclei allows for the identification of germline L1s, similar to using gDNA  
256 from a peripheral source, and individual neuronally relevant somatic mutations occurring: only in  
257 the brain (an L1 that retrotransposed early in neuronal development), in a limited number of  
258 neurons (an L1 that retrotransposed in a neuronal progenitor cell, producing a cluster of  
259 neurons harboring the insertion), or in a single neuron (an L1 that retrotransposed in the neuron  
260 post-mitotically), results which could be verified by comparing the REBELseq results from  
261 NeuN+ nuclei to those of a peripheral tissue. Using REBELseq to compare Ta subfamily L1Hs  
262 from multiple brain regions would also allow for the identification of regionally specific L1Hs  
263 insertions, possibly relevant in the context of neurological and psychiatric disease.

264 Using REBELseq, we identified 157,178 independent putative Ta subfamily L1Hs  
265 insertions, with 1,050 annotating to reference L1Hs in hg19 repeat masker. The ref L1Hs  
266 detected represent ~68% of all hg19 repeat masker L1Hs (Fig. 2a). The L1Hs elements  
267 annotated in hg19 repeat masker represent both Ta subfamily L1Hs, which were specifically  
268 targeted, and pre-Ta subfamily L1Hs. Thus, we did not expect to detect all L1Hs annotated in  
269 hg19 repeat masker. Previous work using the genomes of 15 individuals aligned to hg18  
270 showed that the pre-Ta subfamily of L1Hs represented ~36% of L1Hs insertions, with Ta

271 subfamily L1Hs elements composing the remainder (~64%) which closely approximates the  
272 number we detected (10).

273         The bioinformatic cutoffs described in this manuscript were empirically derived by the  
274 average mean MapQ score for the ref L1Hs, the average number of sequencing reads per  
275 individual per insertion (Fig. 2b), and experimental confirmation of putative novel non-ref L1Hs  
276 insertions (Table 1). We were unable to experimentally validate all putative novel L1Hs in any of  
277 our average read bins. This is possibly due to the detected L1Hs insertion being sequencing  
278 reads of a chimeric PCR product, or possibly because Ta subfamily L1Hs frequently  
279 retrotranspose into repetitive sequences in the genome (e.g. the remnants of older repetitive  
280 elements), making both their genomic alignment and PCR validation difficult. We believe this  
281 suggests the proportion we were able to validate likely represents a minimum value, with  
282 additional detected insertions being confirmable with more complex PCR methods. ~70% of all  
283 known L1Hs averaged at least 100 sequencing reads per person, with our experimentally  
284 determined probability of validating any detected L1Hs above 100 reads being ~59%.  
285 Examining the data with  $\geq 1,000$  reads per person, we still see ~47% of all detected known L1Hs  
286 and a probability of validating any detected L1Hs of almost 90%. These data suggest that novel  
287 non-ref L1Hs having  $\geq 1,000$  reads per person include what are likely real polymorphic and  
288 somatic L1Hs insertions.

289         When examining the genomic distribution of known and novel L1Hs having an average  
290 of  $\geq 1,000$  sequencing reads per person (Fig. 4a), we see that both the known and novel L1Hs  
291 are dispersed throughout the genome, suggesting that REBELseq is an unbiased whole  
292 genome approach. Assessing the distribution of the number of individuals having either a known  
293 (Fig. 4b) or novel (Fig. 4c) L1Hs insertion, we observe that both groups have L1 insertions  
294 occurring throughout the range of possible values (i.e. number of individuals). With respect to  
295 the known L1Hs, there appears to be a trimodal distribution with ~37% of L1Hs occurring in 45-

296 85 individuals and ~9% occurring in  $\leq 3$  and  $\geq 170$ , meaning ~55% of insertions occur in only  
297 ~29% of the possible numbers of individuals. Our finding that a large proportion of known L1Hs  
298 insertions occur in a small fraction of individuals disagrees with a previous report (10), but we  
299 believe this may be due to their small sample size, thereby diminishing the likelihood that an  
300 individual would not have a given insertion. ~70% of novel L1Hs occur in  $\leq 7$  people, while ~8%  
301 of novel L1Hs occur in more than two thirds of individuals, suggesting that these L1Hs insertions  
302 likely represent uncatalogued ref L1Hs occurring at common minor allele frequencies.

303       Taken together, we believe these data demonstrate that REBELseq reliably detects both  
304 known and novel Ta subfamily L1Hs insertions, and that this technique could be a powerful  
305 method for examining the association of the prevalence of L1Hs insertions in a broad range of  
306 human disease.

307

### 308 **Acknowledgements**

309 The authors wish to thank those providing support for this work. BR is supported, in part, by a  
310 2017 NARSAD Young Investigator Grant (#26634) from the Brain and Behavior Research  
311 Foundation as the Patrick A. Coffey Investigator, funding for which was generously provided by  
312 Ronald and Kathy Chandonais. BR was supported during part of this work by T32 MH014654  
313 (PI is WB). This work was funded by National Institute of Health grants to WB (R21NS095756,  
314 R01DA040972 & R01MH109260). RC was supported by K01DA036751. Additionally, the  
315 authors wish to thank all the participants, and their families, who donated their brain tissue to  
316 research for making these studies possible.

317

### 318 **Conflict of Interests**

319 The authors report no conflicts of interest.

- 320 1. Richardson SR, Morell S, Faulkner GJ. L1 retrotransposons and somatic mosaicism in the brain.  
321 Annual review of genetics. 2014;48:1-27.
- 322 2. Beck CR, Collier P, Macfarlane C, Malig M, Kidd JM, Eichler EE, et al. LINE-1 retrotransposition  
323 activity in human genomes. Cell. 2010;141(7):1159-70.
- 324 3. Ostertag EM, Kazazian HH, Jr. Biology of mammalian L1 retrotransposons. Annual review of  
325 genetics. 2001;35:501-38.
- 326 4. Boissinot S, Chevret P, Furano AV. L1 (LINE-1) retrotransposon evolution and amplification in  
327 recent human history. Molecular biology and evolution. 2000;17(6):915-28.
- 328 5. Kazazian HH, Jr., Wong C, Youssoufian H, Scott AF, Phillips DG, Antonarakis SE. Haemophilia A  
329 resulting from de novo insertion of L1 sequences represents a novel mechanism for mutation in man.  
330 Nature. 1988;332(6160):164-6.
- 331 6. Naufer MN, Callahan KE, Cook PR, Perez-Gonzalez CE, Williams MC, Furano AV. L1  
332 retrotransposition requires rapid ORF1p oligomerization, a novel coiled coil-dependent property  
333 conserved despite extensive remodeling. Nucleic acids research. 2016;44(1):281-93.
- 334 7. Mathias SL, Scott AF, Kazazian HH, Jr., Boeke JD, Gabriel A. Reverse transcriptase encoded by a  
335 human transposable element. Science. 1991;254(5039):1808-10.
- 336 8. Feng Q, Moran JV, Kazazian HH, Jr., Boeke JD. Human L1 retrotransposon encodes a conserved  
337 endonuclease required for retrotransposition. Cell. 1996;87(5):905-16.
- 338 9. Jiang Y, Matevosian A, Huang HS, Straubhaar J, Akbarian S. Isolation of neuronal chromatin  
339 from brain tissue. BMC Neurosci. 2008;9:42.
- 340 10. Ewing AD, Kazazian HH, Jr. Whole-genome resequencing allows detection of many rare LINE-1  
341 insertion alleles in humans. Genome research. 2011;21(6):985-90.
- 342 11. Strevva VA, Jordan VE, Linker S, Hedges DJ, Batzer MA, Deininger PL. Sequencing, identification  
343 and mapping of primed L1 elements (SIMPLE) reveals significant variation in full length L1 elements  
344 between individuals. BMC Genomics. 2015;16:220.
- 345 12. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nat Methods.  
346 2012;9(4):357-9.
- 347 13. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence Alignment/Map  
348 format and SAMtools. Bioinformatics. 2009;25(16):2078-9.
- 349 14. Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, et al. An integrated  
350 map of structural variation in 2,504 human genomes. Nature. 2015;526(7571):75-81.
- 351 15. Badge RM, Alisch RS, Moran JV. ATLAS: a system to selectively identify human-specific L1  
352 insertions. Am J Hum Genet. 2003;72(4):823-38.
- 353 16. Baillie JK, Barnett MW, Upton KR, Gerhardt DJ, Richmond TA, De Sapio F, et al. Somatic  
354 retrotransposition alters the genetic landscape of the human brain. Nature. 2011;479(7374):534-7.
- 355 17. Costello M, Fleharty M, Abreu J, Farjoun Y, Ferriera S, Holmes L, et al. Characterization and  
356 remediation of sample index swaps by non-redundant dual indexing on massively parallel sequencing  
357 platforms. BMC Genomics. 2018;19(1):332.
- 358 18. Erwin JA, Paquola AC, Singer T, Gallina I, Novotny M, Quayle C, et al. L1-associated genomic  
359 regions are deleted in somatic cells of the healthy human brain. Nat Neurosci. 2016;19(12):1583-91.
- 360 19. Evrony GD, Cai X, Lee E, Hills LB, Elhosary PC, Lehmann HS, et al. Single-neuron sequencing  
361 analysis of L1 retrotransposition and somatic mutation in the human brain. Cell. 2012;151(3):483-96.
- 362 20. Iskow RC, McCabe MT, Mills RE, Torene S, Pittard WS, Neuwald AF, et al. Natural mutagenesis of  
363 human genomes by endogenous retrotransposons. Cell. 2010;141(7):1253-61.
- 364 21. Kvikstad EM, Piazza P, Taylor JC, Lunter G. A high throughput screen for active human  
365 transposable elements. BMC Genomics. 2018;19(1):115.
- 366 22. Upton KR, Gerhardt DJ, Jesuadian JS, Richardson SR, Sanchez-Luque FJ, Bodea GO, et al.  
367 Ubiquitous L1 mosaicism in hippocampal neurons. Cell. 2015;161(2):228-39.

368 23. Zhao B, Wu Q, Ye AY, Guo J, Zheng X, Yang X, et al. Somatic LINE-1 retrotransposition in cortical  
369 neurons and non-brain tissues of Rett patients and healthy individuals. PLoS Genet.  
370 2019;15(4):e1008043.

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389 **Figure legends**

390 Figure 1: Schematic of the construction of Ta subfamily enriched L1Hs sequencing libraries.

391 gDNA isolated from NeuN+ nuclei was enzymatically digested with HaeIII to fragment the  
392 genome. A single primer extension using Ta subfamily specific L1HsACA primer extends the 3'  
393 end of the L1 sequence into the downstream gDNA. The 3' 'A' overhang from the single primer  
394 extension is ligated to a custom T-linker, and primary PCR amplifies the construct using  
395 L1HsACA and T-linker specific primers. Hemi-nested secondary PCR using the L1Hs specific  
396 L1HsG primer and T-linker primer reduces the length of the L1 sequence carried forward and  
397 adds a sequencing adapter to the L1 end. Tertiary PCR uses 5' overhang primers to add  
398 barcodes to the L1 end and Illumina flow cell adapters to both ends of the amplicon.

399 Figure 2: Detection levels of Ta subfamily L1Hs. **a** The number of ref L1Hs detected per  
400 chromosome versus the number of L1Hs annotated in hg19 repeat masker. Overall, we  
401 detected ~68% of L1Hs annotated in hg19 repeat masker, which aligns with our expectations  
402 based on the literature (see discussion). **b** Average number of sequencing reads per person for  
403 a given L1 insertion. Data was binned to show contrast between the average number of  
404 sequencing reads per person seen for known and putative novel L1Hs insertions. An average  
405 number of sequencing reads per person  $\geq 100$  was used as a bioinformatic cutoff.

406 Figure 3: Example image of a successful and unsuccessful confirmatory PCR. PCR  
407 experiments were conducted to determine the proportion of putative novel non-ref L1Hs  
408 insertions detected by REBELseq that could be independently validated. For each insertion, a  
409 person predicted to have the L1Hs insertion (+) and a person not predicted to have the insertion  
410 (-) were used to amplify the genomic region purported to contain the L1Hs insertion (Filled site,  
411 F) and the same genomic region if it did not contain the insertion (Empty site, E). Insertion #1 is  
412 a positive confirmation of the results predicted by REBELseq, while Insertion #2 is a negative  
413 confirmation.



414 Figure 4: Distributions of high confidence L1Hs insertions. **a** The genomic distribution of known  
415 L1Hs (above the chromosome numbers) and novel L1Hs (below the chromosome numbers) per  
416 10 MB window of each chromosome. Alternating color pattern and labeled central blocks  
417 represent the different chromosomes. Known and novel L1Hs are distributed throughout the  
418 genome, suggesting REBELseq is an unbiased whole genome approach. **b** Number of  
419 individuals sharing a known L1Hs insertion. The number of known L1Hs insertions shared by  
420 different numbers of individuals shows a trimodal distribution. While some ref L1Hs occur with a  
421 rate of affected individuals approaching 1.0, the other local maxima are focused at few or less  
422 than half of surveyed individuals. This demonstrates that known L1Hs should be considered  
423 polymorphic in nature, rather than ubiquitous, in the human genome. **c** Number of individuals  
424 sharing a novel L1Hs insertion. The number of novel L1Hs insertions shared by different  
425 numbers of individuals shows a right skewed distribution. While some novel L1Hs were detected  
426 in most individuals, most novel L1Hs were detected in one or a few individuals.

427

428

429

430

431

432

433

434

435

436

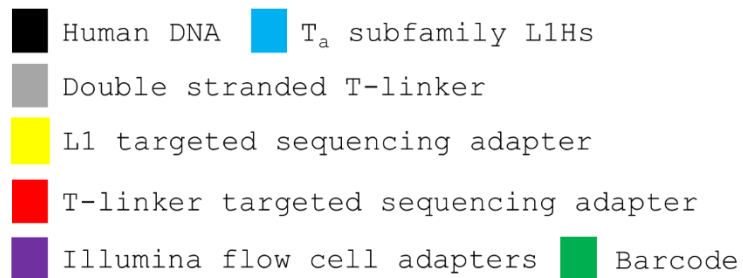
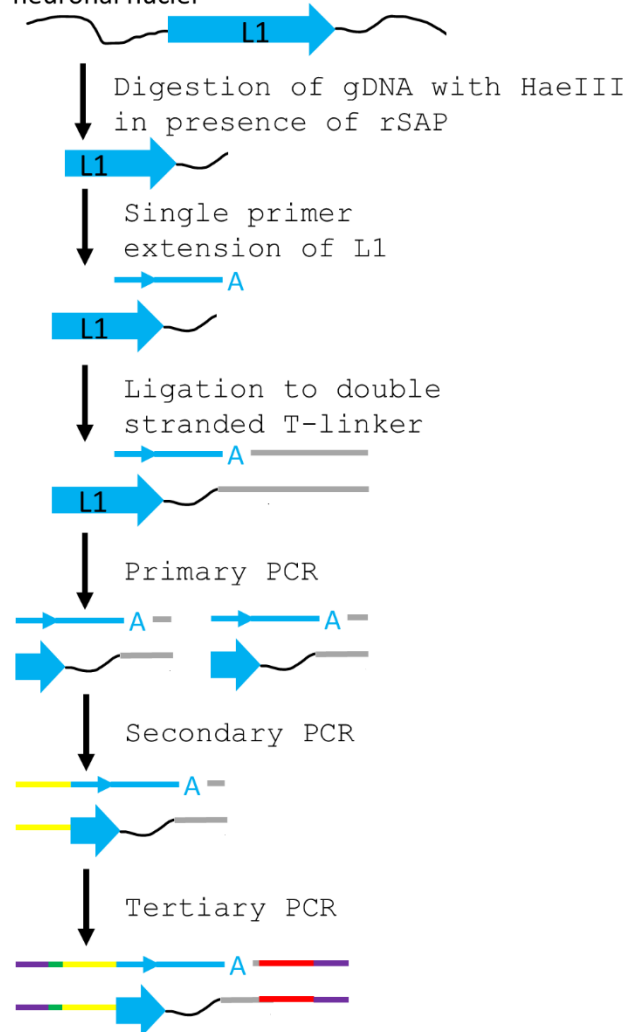
437

438

439

440 **Figure 1**

Human gDNA extracted from NeuN+ neuronal nuclei



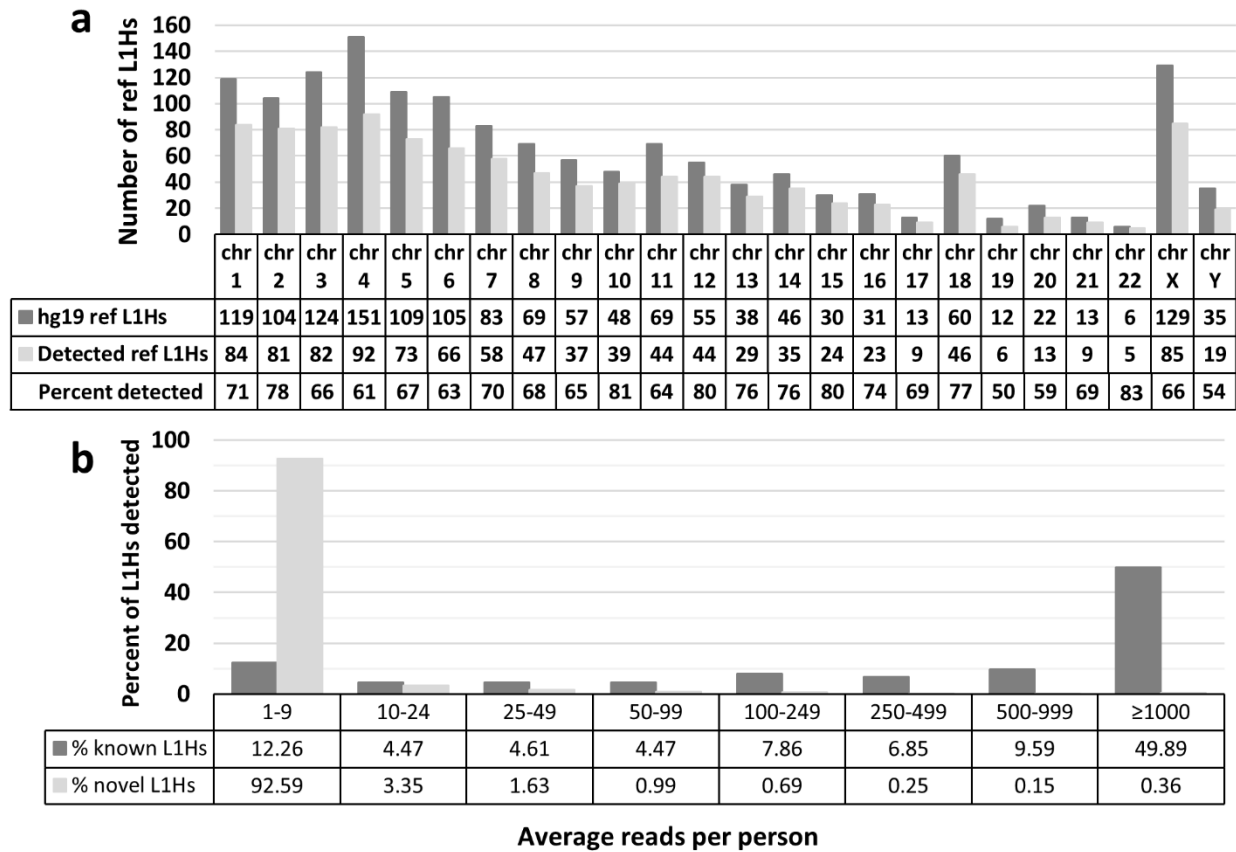
441

442

443

444

445 **Figure 2**



446

447

448

449

450

451

452

453

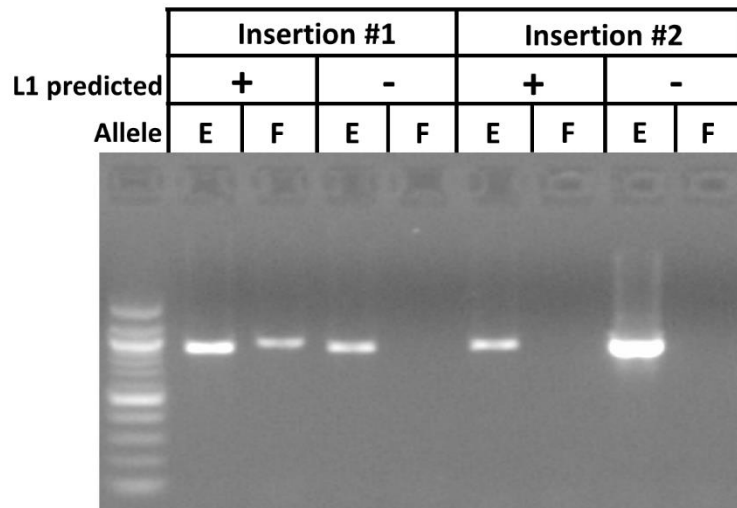
454

455

456

457

458 **Figure 3**



459

460

461

462

463

464

465

466

467

468

469

470

471

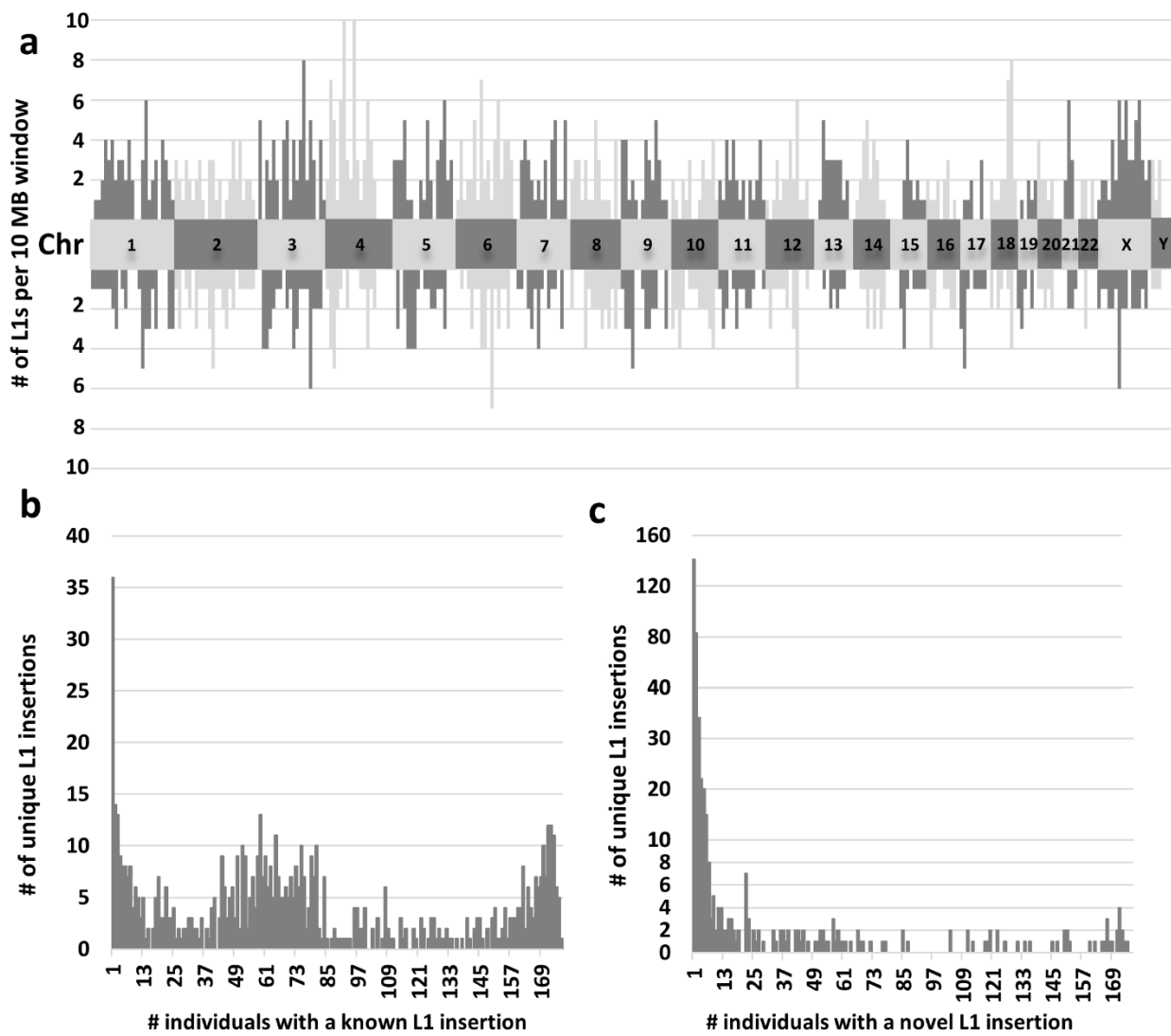
472

473

474

475

476 **Figure 4**



477

478

479