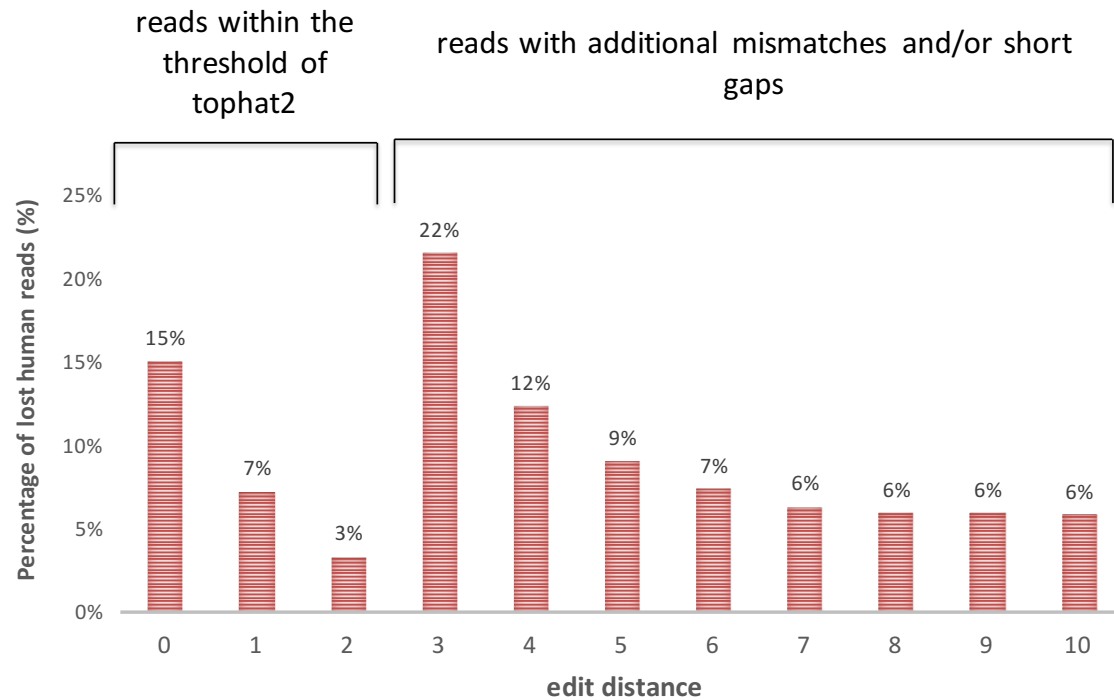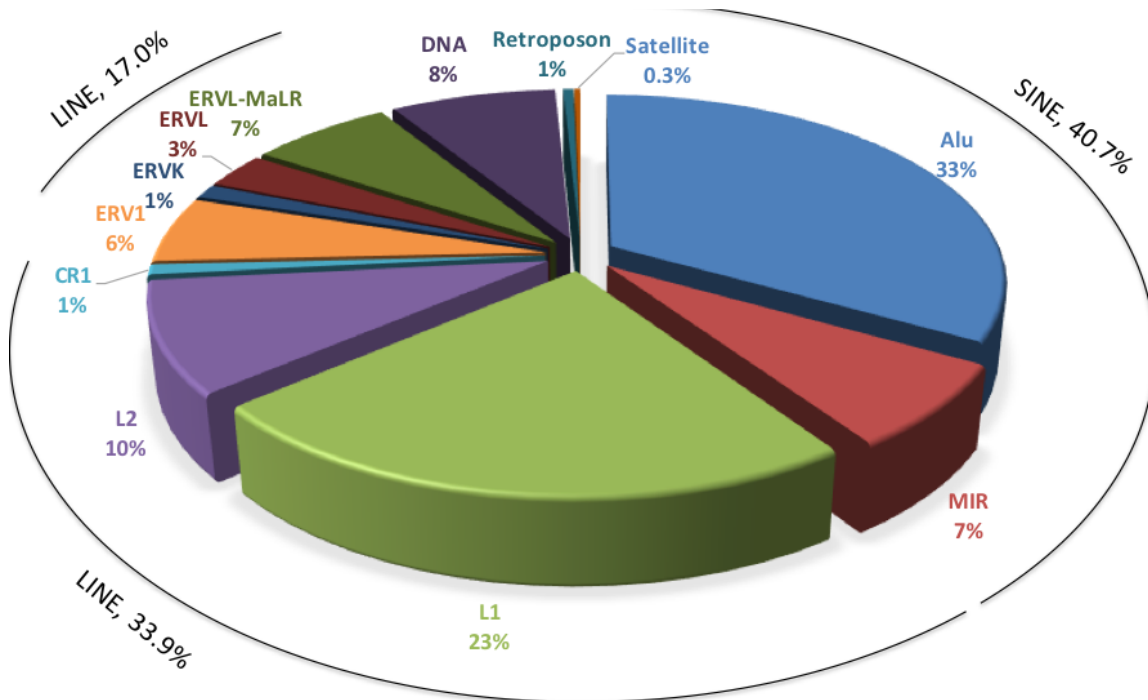1    *Supplementary Figures*

2

3

4



5

6    **Supplemental Figure S1. Edit distance of lost human reads.**

7    Unmapped reads were remapped to the human references using Megablast. Edit distance

8    was calculated as the minimum number of operations required to transform a read

9    sequence into the corresponding reference subsequence. Reads are grouped by edit

10    distance with the transcriptome or the genome reference. The percentages are the

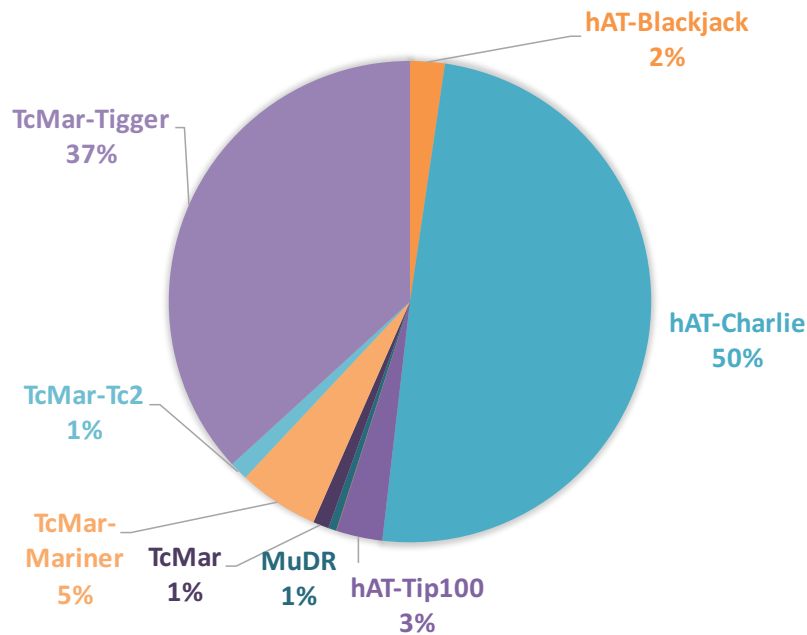11    averages across 10,641 samples.

12

On average 7% of RNA-Seq reads are categorized as repeats

**Supplemental Figure S2. Profile of repeat elements based on repeat sequences inferred from mapped and unmapped reads (lost repeat reads).**

ROP identifies and categorizes repetitive sequences among the mapped and unmapped reads. Mapped reads were categorized based on the overlap with the repeat instances prepared from RepeatMasker annotation (RepeatMasker v3.3, Repeat Library 20120124). Lost repeat reads are unmapped RNA-Seq reads aligned onto the reference repeat sequences (prepared from Repbase v20.07). The percentages are the averages across 10,641 samples.

25

# GTEx DNA repeats



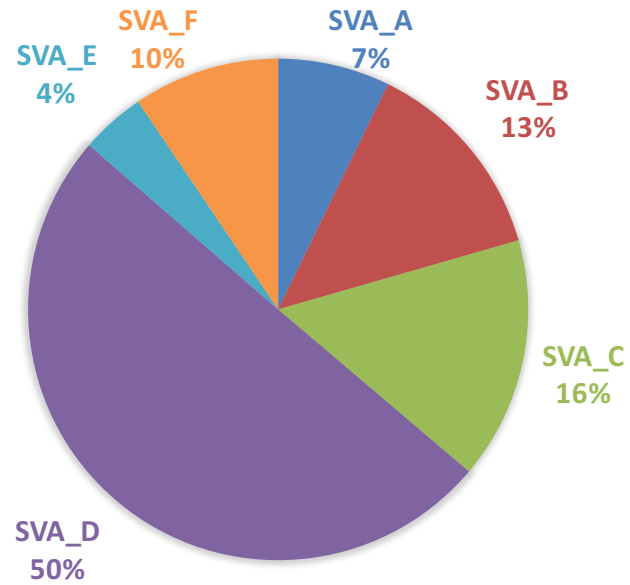*Percentages are calculated as a fraction from the reads matching DNA repeats

26

**Supplemental Figure S3. Profile of DNA repeats based on repeat sequences inferred from mapped and unmapped reads (lost repeat reads).**

ROP identifies and categorizes DNA repetitive sequences among the mapped and unmapped reads. Mapped reads were categorized based on the overlap with the repeat instances prepared from RepeatMasker annotation (RepeatMasker v3.3, Repeat Library 20120124). Lost repeat reads are unmapped RNA-Seq reads aligned onto the reference repeat sequences (prepared from Repbase v20.07). The percentages are the averages across 10,641 samples.
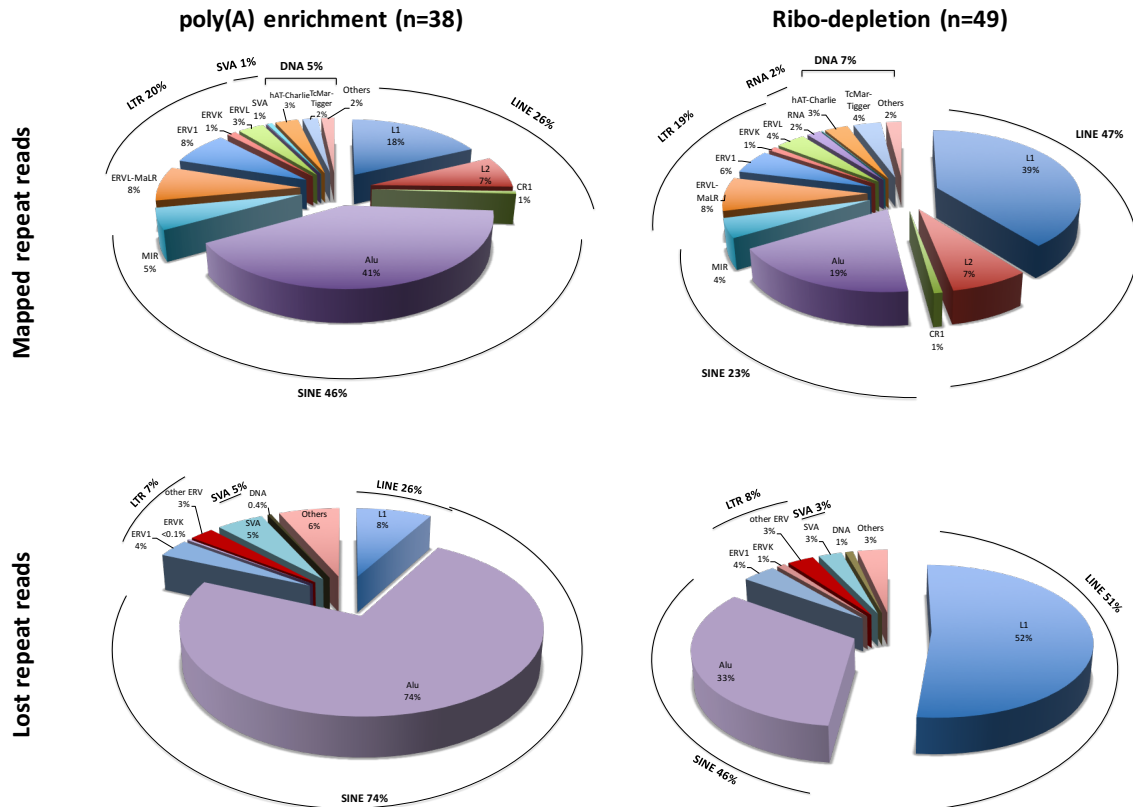
# SINE–VNTR–*Alu* (SVA)



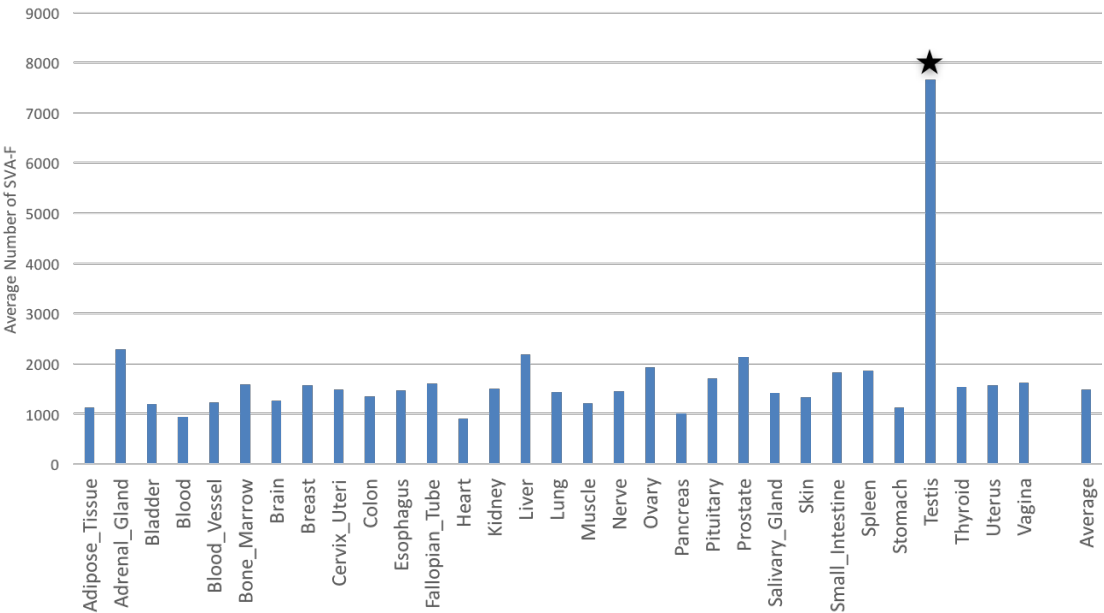*Percentages are calculated as a fraction from the reads matching SVA Retroposons

36

37

38 **Supplemental Figure S4. Profile of SVA retrotransposons based on repeat sequences**

39 **inferred from mapped and unmapped reads (lost repeat reads).** ROP identifies and

40 categorizes SVA retrotransposons sequences among the mapped and unmapped reads.

41 Mapped reads were categorized based on the overlap with the repeat instances prepared

42 from RepeatMasker annotation (Repeatmasker v3.3, Repeat Library 20120124). Lost

43 repeat reads are unmapped RNA-Seq reads aligned onto the reference repeat sequences

44 (prepared from Repbase v20.07). The percentages are the averages across 10,641

45 samples.

46

**Supplemental Figure S5.** Profile of repeat elements across poly(A) enrichment and ribo-depletion libraries. ROP identifies and categorizes repetitive sequences among the mapped and unmapped reads. RNA-Seq samples were prepared by poly(A) enrichment protocol (n=38) and ribo-depletion protocol (n=49). Mapped reads were categorized based on the overlap with the repeat instances prepared from RepeatMasker annotation (RepeatMasker v3.3, Repeat Library 20120124). Lost repeat reads are unmapped RNA-Seq reads aligned onto the reference repeat sequences (prepared from Repbase v20.07).

57

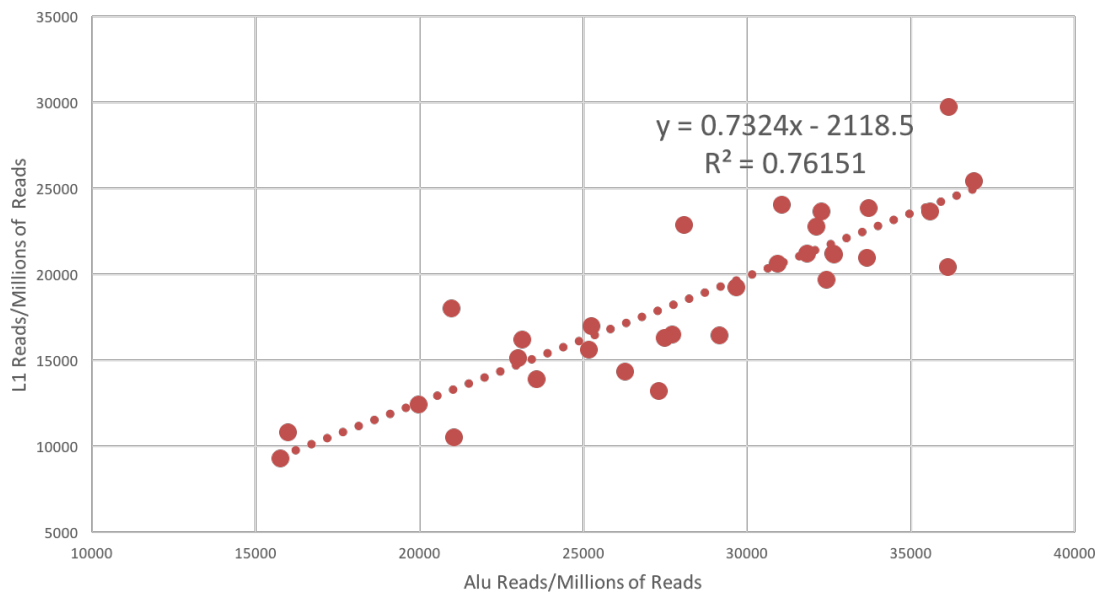## Average Number of SVA-F reads across Tissue



58

59

60  **Supplemental Figure S6.** Average number of SVA-F reads across GTEx tissues. ROP

61  identifies and categorizes SVA retrotransposons sequences among the mapped and

62  unmapped reads. Mapped reads were categorized based on the overlap with the repeat

63  instances prepared from RepeatMasker annotation (RepeatMasker v3.3, Repeat Library

64  20120124). Lost repeat reads are unmapped RNA-Seq reads aligned onto the reference

65  repeat sequences (prepared from Repbase v20.07). Among the GTEx tissues, *testis*

66  showed significantly higher expression of SVA F retrotransposons compared to other

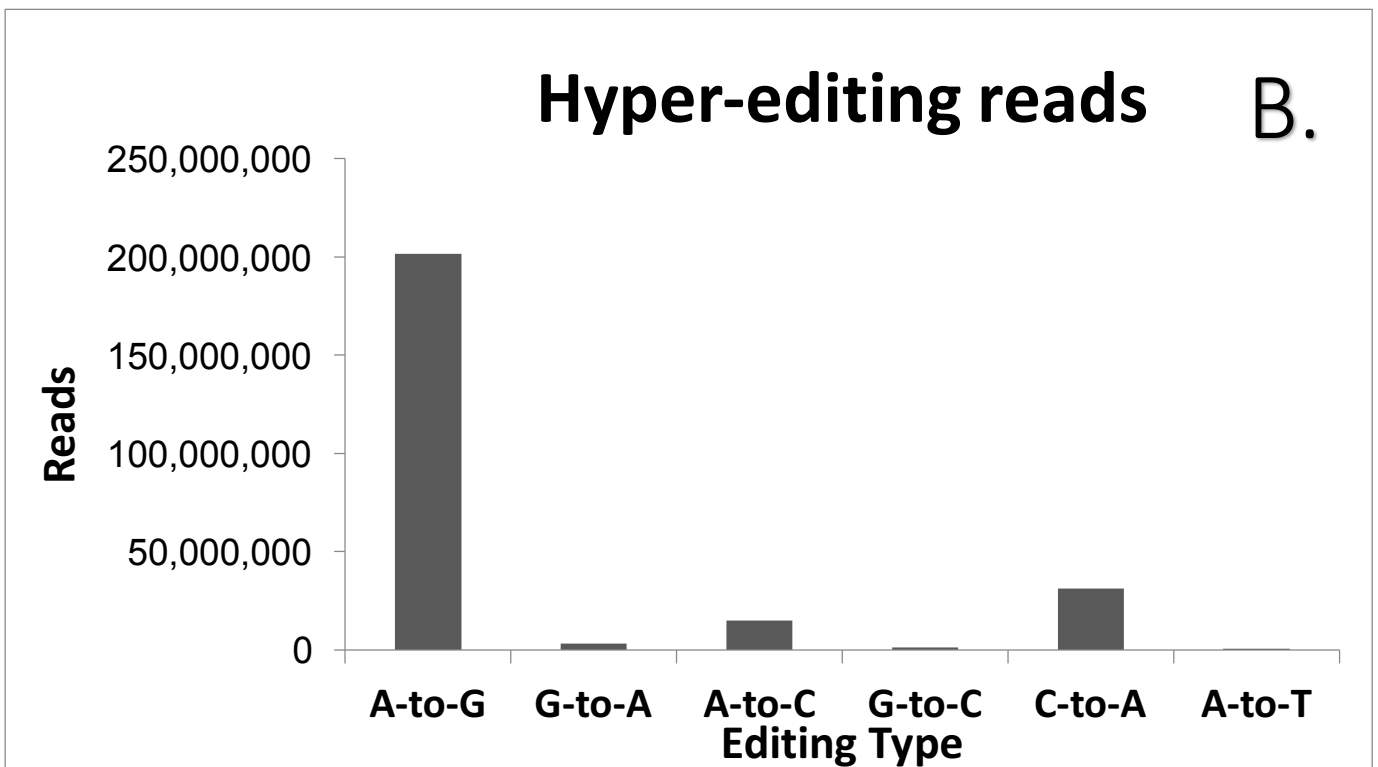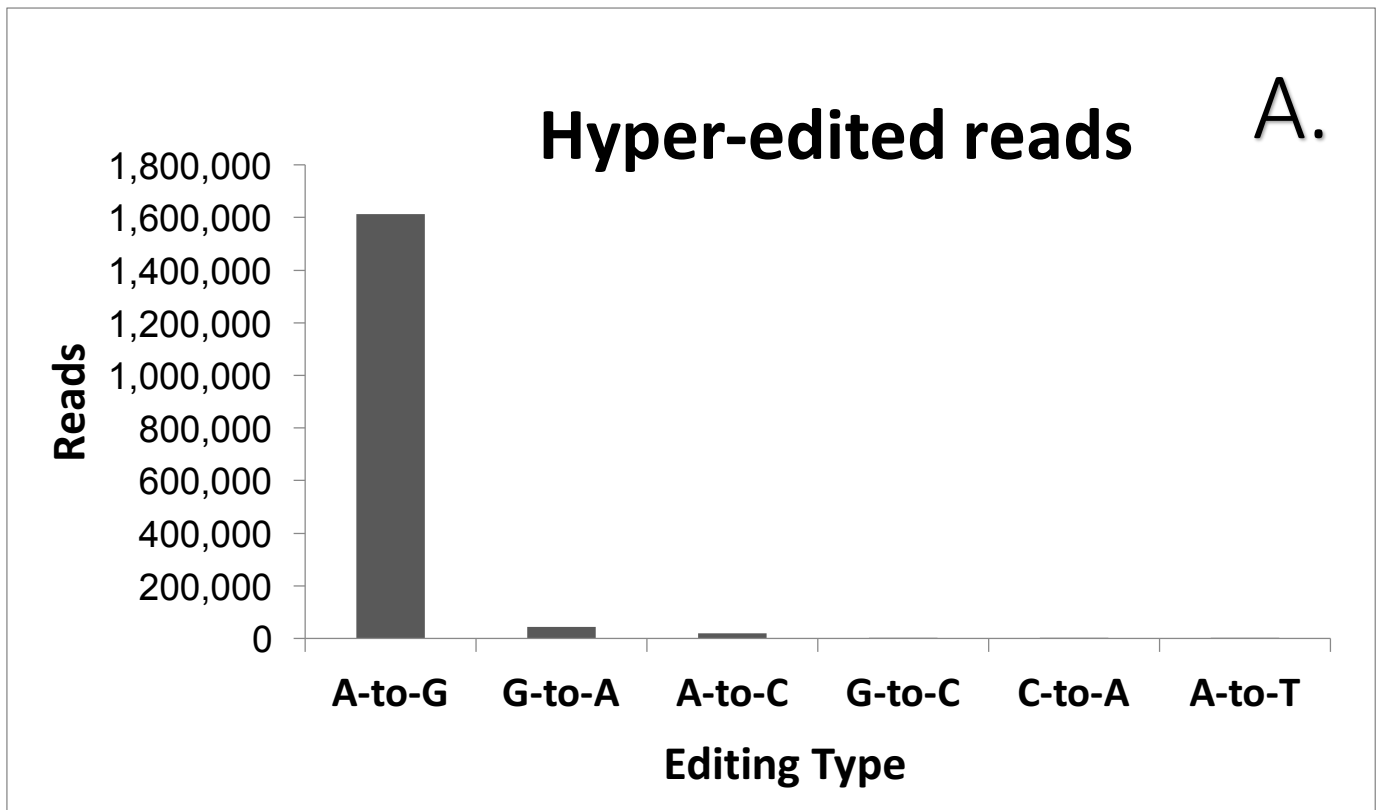67  tissues ($\mathbf{p = 2.46 \times 10^{-33}}$).

## Alu and L1 co-expression in Individual Tissues



68

69

70   **Supplemental Figure S7.** Co-expression of Alu and L1 elements across GTEx tissues. ROP

71   identifies and categorizes repetitive sequences among the mapped and unmapped reads.

72   Mapped reads were categorized based on the overlap with the repeat instances prepared

73   from RepeatMasker annotation (RepeatMasker v3.3, Repeat Library 20120124). Lost

74   repeat reads are unmapped RNA-Seq reads aligned onto the reference repeat sequences

75   (prepared from Repbase v20.07).

**Hyper-edited reads** A.

Reads / Editing Type

A-to-G, G-to-A, A-to-C, G-to-C, C-to-A, A-to-T

**Hyper-editing reads** B.

Reads / Editing Type

A-to-G, G-to-A, A-to-C, G-to-C, C-to-A, A-to-T

77   **Supplemental Figure S8. Distribution of hyper-edited reads.**

78   **A.** Hyper-editing identified in the in-house data. Results showed that 96% of the reads

79   were A-to-G, indicating a high level of specificity for the hyper-editing screen. The

80   1,613,213 detected A-to-G reads contain 10,666,458 editing events (3,157,685 unique

81   editing-sites). **B.** Hyper-editing identified in the GTEx RNA-Seq data. Results showed that

82   80% of the reads were A-to-G, indicating a high level of specificity for the hyper-editing

83   screen. The 201,676,069 detected A-to-G reads contain 1,130,591,911 editing events

84   (690,386,562 unique editing-sites).
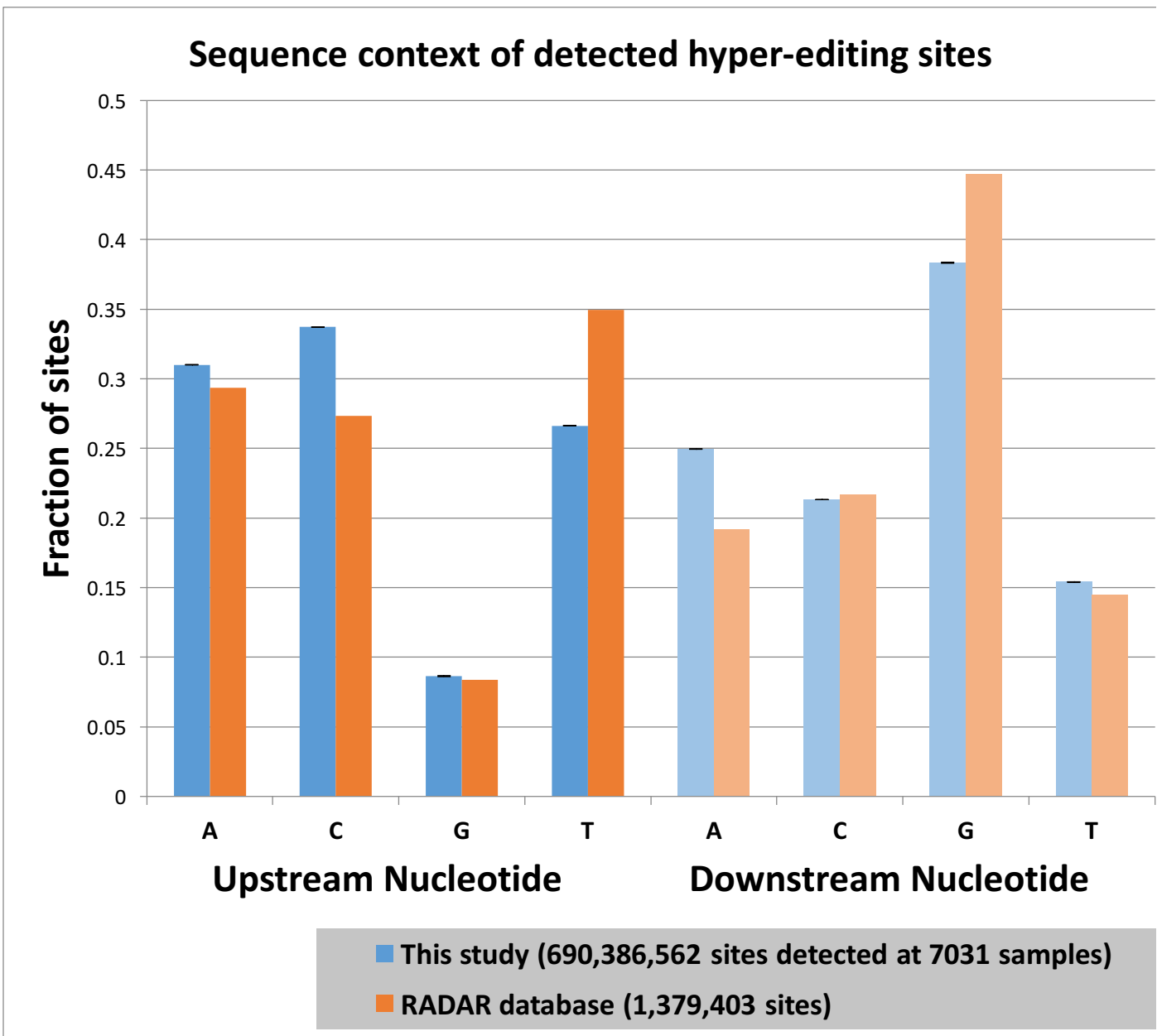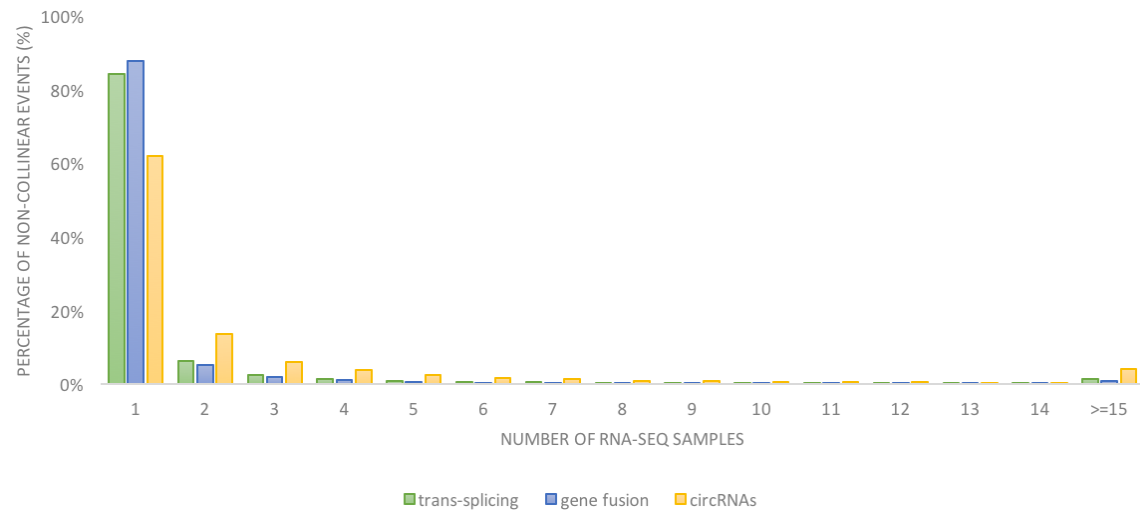
85

A.



**Sequence context of detected hyper-editing sites**

86

87

B.

88.

89

90



Sequence context of detected hyper-editing sites

This study (690,386,562 sites detected at 7031 samples)
RADAR database (1,379,403 sites)

91     **Supplemental Figure S9. The sequence context of the Figure S8. The sequence context**

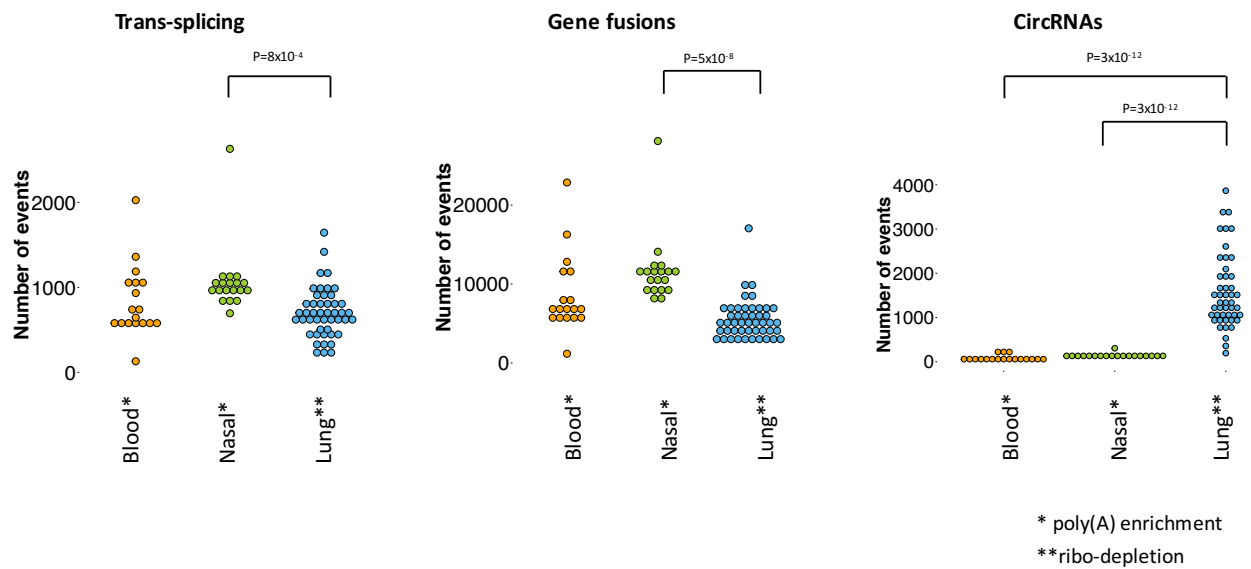92     **of the detected hyper-edited A-to-G sites.**

93     The sequence near the detected hyper-editing sites is depleted of Gs upstream and

94     enriched with Gs downstream, in agreement with previously known data about the ADAR

95     motif. The bars correspond to the fraction of editing sites with each type of nucleotide

96     one base upstream and downstream of the site. Results are shown for sites detected in-

97     house RNA-Seq data (A) and GTEx RNA-Seq data (B) using the hyper-editing pipeline and

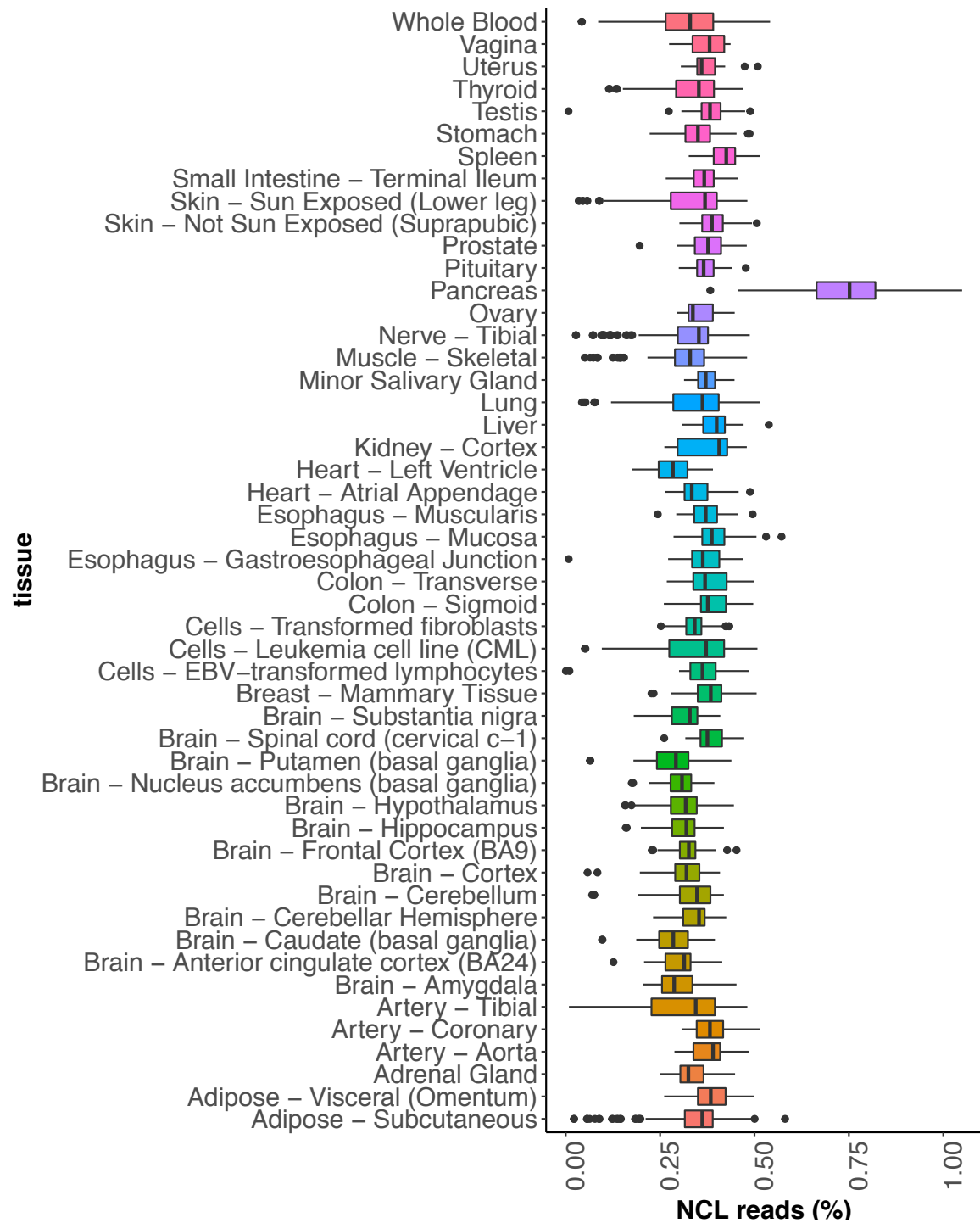98     human editing-sites from the RADAR database.

99

100

101

102

**Supplemental Figure S10. Distribution of non-co-linear (NCL) events across 10,641**

**samples. .**

Reads arising from trans-splicing, gene fusion and circRNA events are captured by a

TopHat-Fusion and CIRCexplorer2 tools. Trans-splicing events are identified from reads

that are spliced distantly on the same chromosome. Gene fusion events are identified

from reads spliced across different chromosomes. CircRNAs are identified from reads

spliced in a head-to-tail configuration.

110

**Supplemental Figure S11. Number of NCL events across in-house tissues and library preparation protocols.**

NCL events per sample are detected by TopHat-Fusion and CIRCexplorer tools. Samples were prepared with poly(A) selection (whole blood and nasal epithelium) and ribo-depletion (lung epithelium) protocols. Trans-splicing events are identified from reads spliced distantly on the same chromosome. Gene fusion events are identified from reads spliced across different chromosomes.
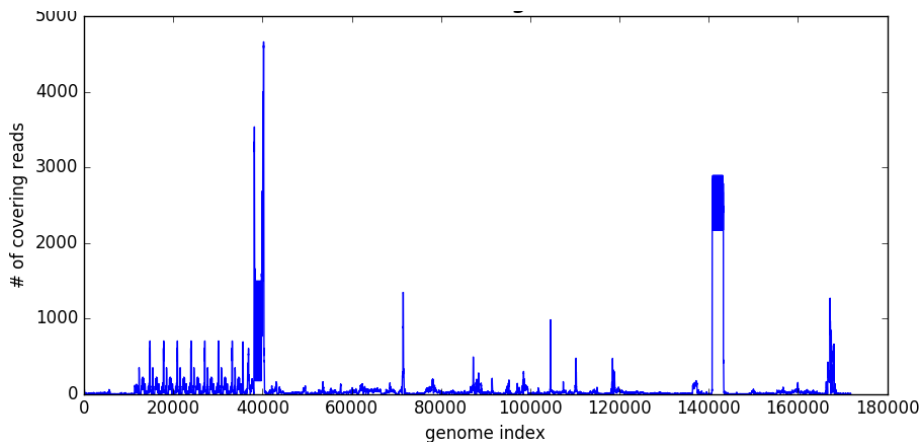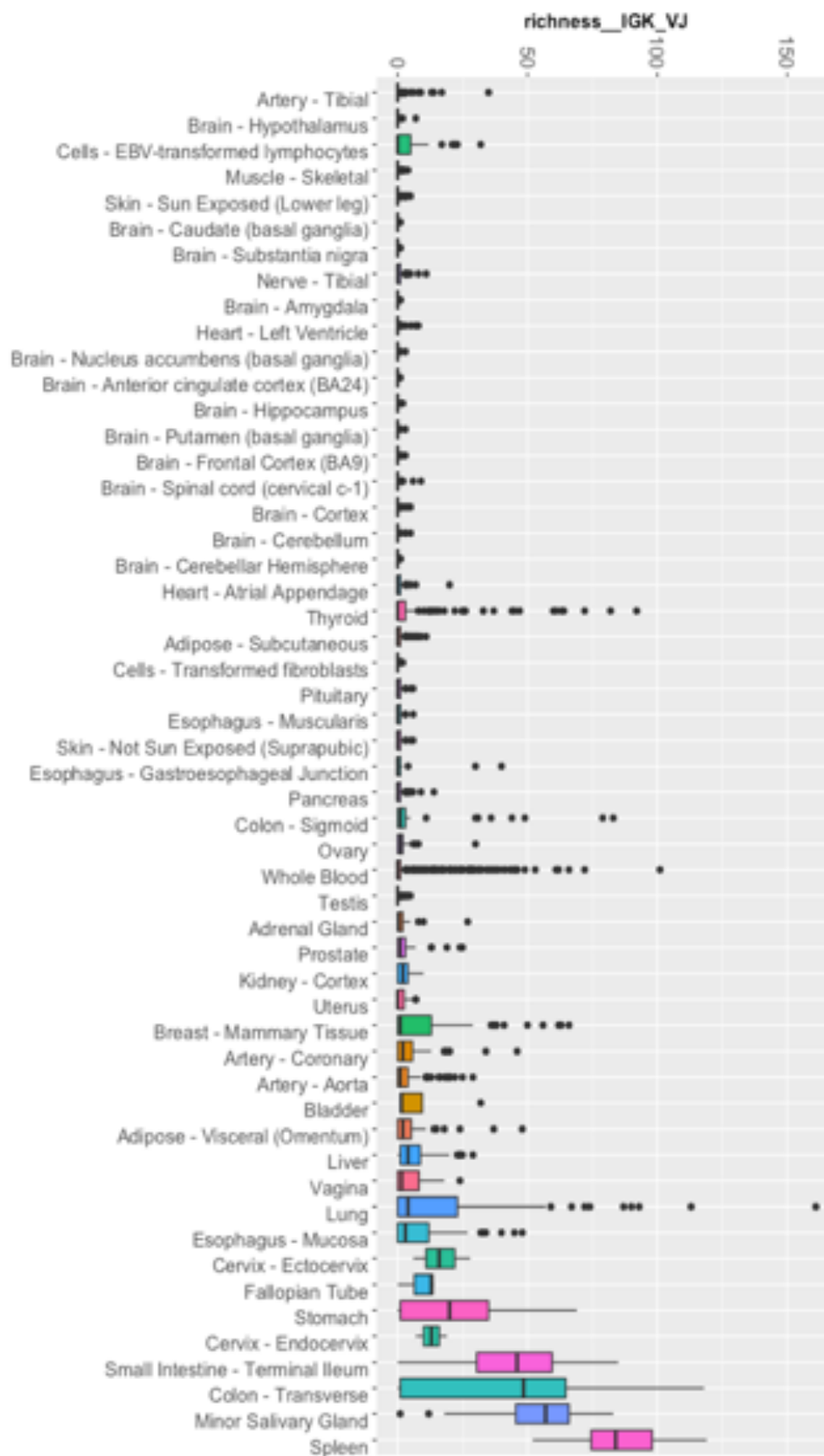
124

125

126

127

128 **Supplemental Figure S12. Percentage of NCL reads across GTEx tissues (n=54).**

129 Percentages are calculated from the total number of reads. Reads arising from trans-

130 splicing, gene fusion and circRNA events are captured by a TopHat-Fusion and

131 CIRCexplorer2 tools and reported a NCL reads.

132

133

134

135

136

137



138

139 **Supplemental Figure S13.** An example of coverage plot of EBV virus. Viral reads were

140 obtained by ROP protocol from GTEx RNA-Seq sample of EBV-transformed

141 lymphoblastoid cell lines (LCLs).

142

143

144

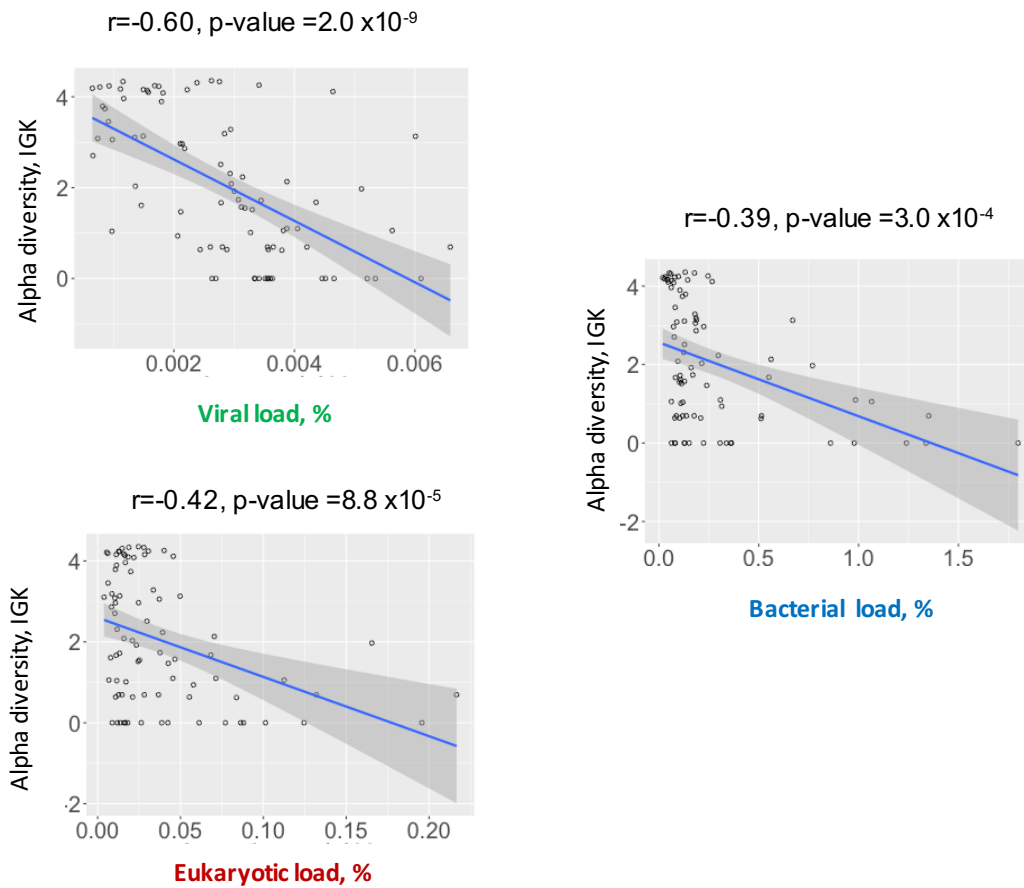richness__IGK_VJ

146     **Supplemental Figure S14. Number of VJ recombinations across GTEx human tissues for**

147     **IGK chain.**

148

149

150

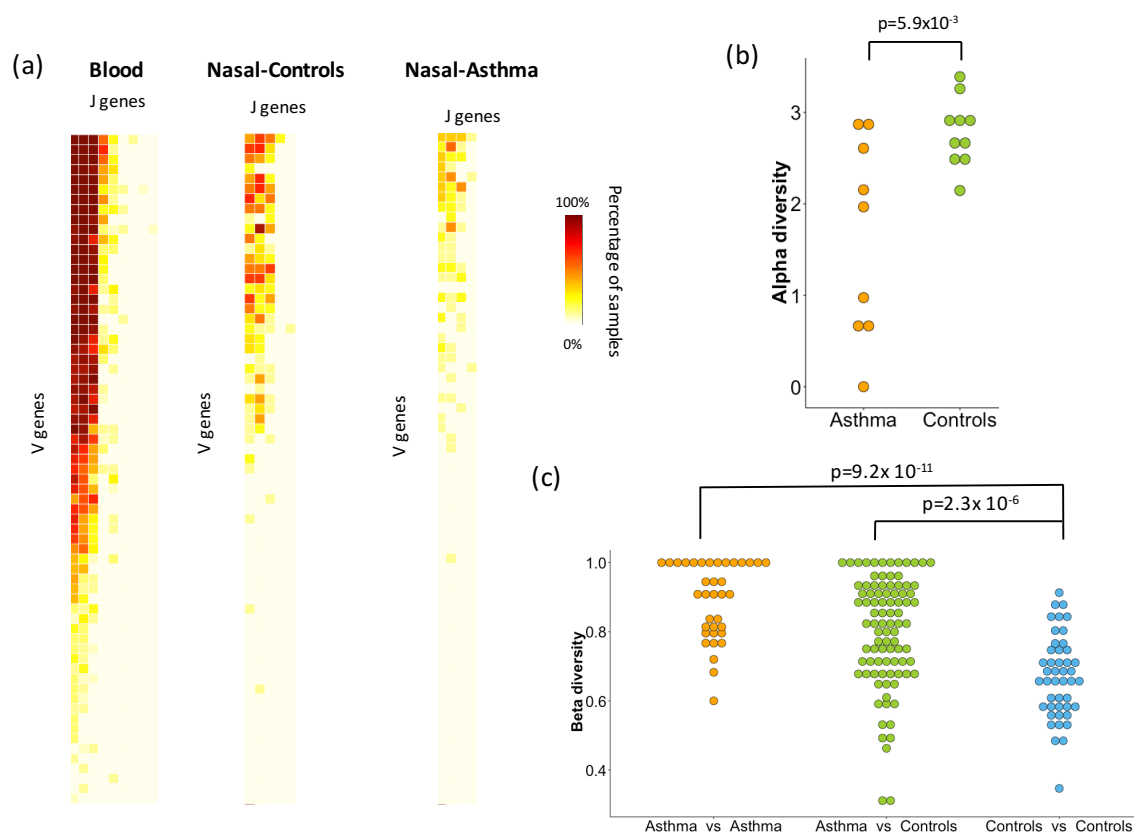151    **Supplemental Figure S15. Number of VJ recombinations across GTEx human tissues for**

152    **IGL chain.**

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

$r=-0.60$, p-value $=2.0 \times 10^{-9}$

Alpha diversity, IGK

**Viral load, %**

$r=-0.39$, p-value $=3.0 \times 10^{-4}$

Alpha diversity, IGK

**Bacterial load, %**

$r=-0.42$, p-value $=8.8 \times 10^{-5}$

Alpha diversity, IGK

**Eukaryotic load, %**

173

174  **Supplemental Figure S16. Association between microbial load and immune diversity.**
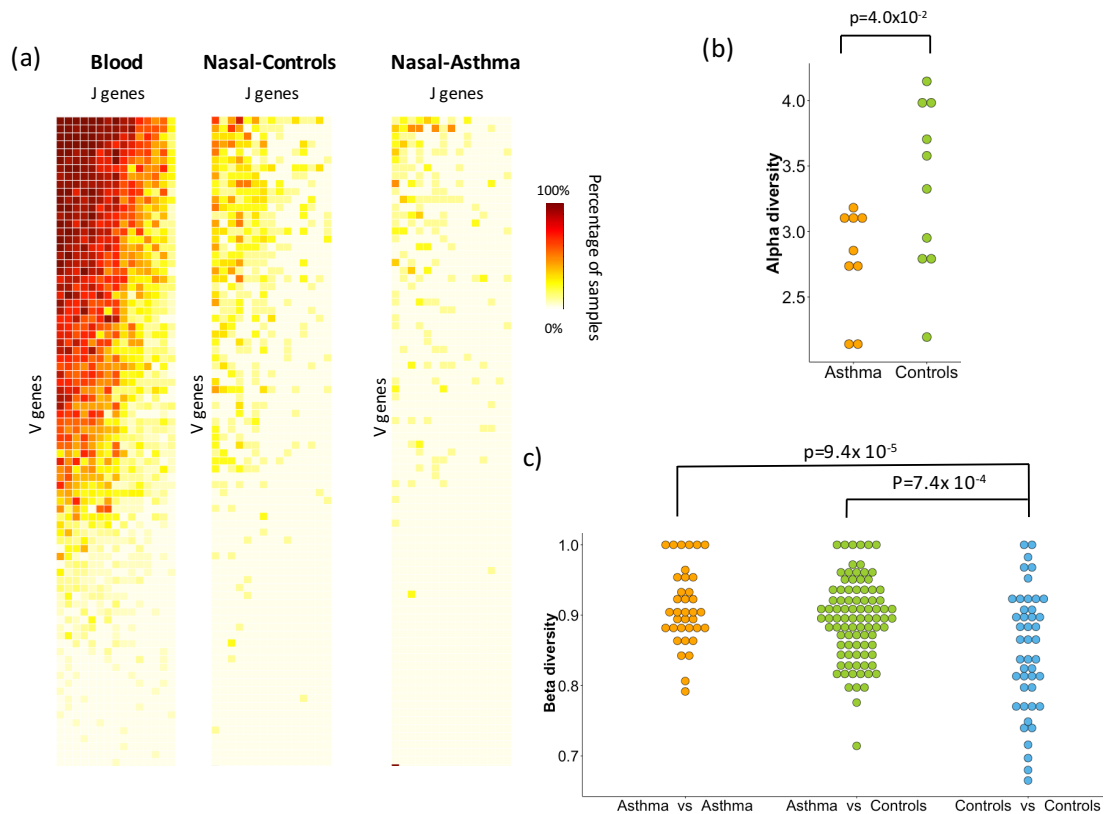
175  (a) Scatterplot of the viral load and combinatorial immune diversity of IGK locus. Pearson

176  correlation coefficient (r) and p -value are reported. (b) Scatterplot of the eukaryotic load

177  and combinatorial immune diversity of IGK locus. Pearson correlation coefficient (r) and

178  p -value are reported. (c) Scatterplot of the bacterial load and combinatorial immune

179  diversity of IGK locus. Pearson correlation coefficient (r) and p -value are reported.

180

**Supplemental Figure S17. Combinatorial diversity of immunoglobulin lambda locus (IGL) locus differentiates disease status.**

(a) Heat map depicting the percentage of RNA-Seq samples supporting particular VJ combination for whole blood, nasal epithelium of healthy controls and asthmatic individuals. Each row corresponds to a V gene and each column corresponds to a J gene. (b) Alpha diversity is measured using the Shannon entropy incorporating the total number of VJ combinations and their relative proportions. Nasal epithelium of asthmatic individuals exhibits decreased combinatorial diversity of IGK locus compared to that of healthy controls (p-value=$5.9x10^{-3}$). (c) Compositional similarities between the samples in terms of gain or loss of VJ combinations of IGK locus are measured using the Sørensen–
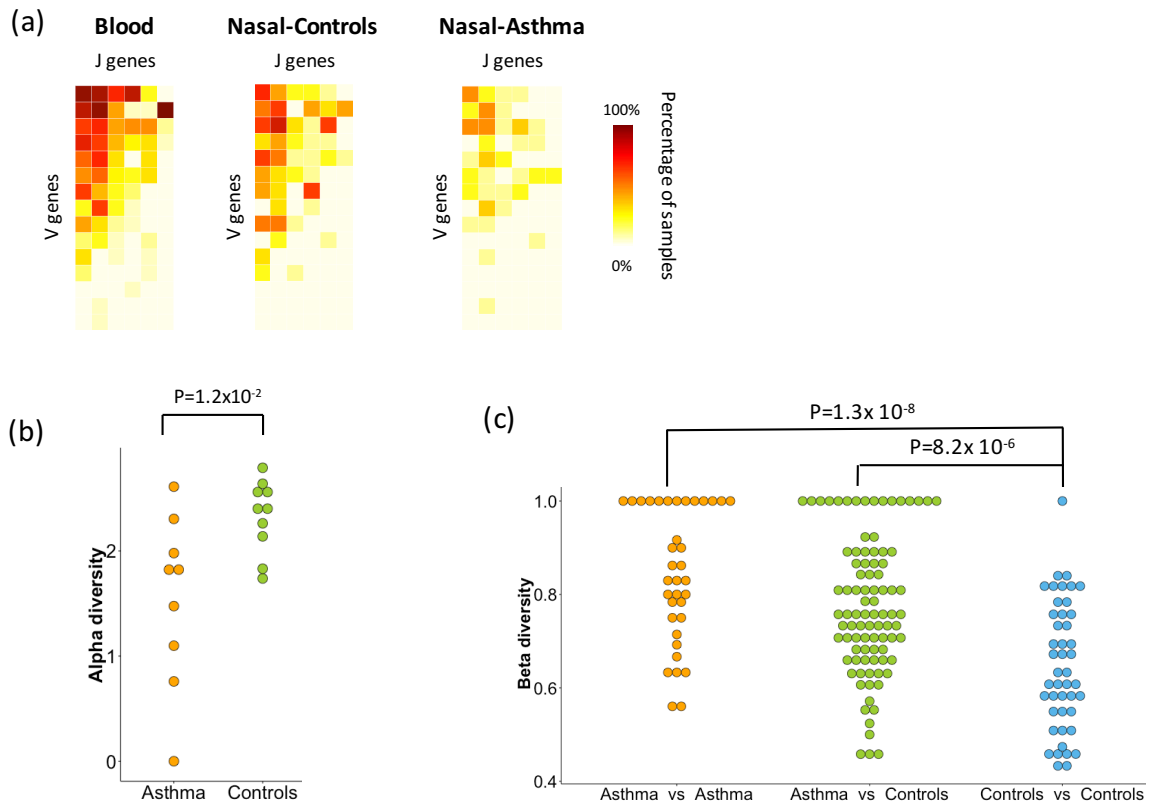
191    Dice index across pairs of samples from the same group (Asthma, Controls) and pairs of

192    sample from different groups (Asthma versus Controls). Lower level of similarity is

193    observed between nasal samples of the asthmatic individuals compared to the unaffected

194    controls (p-value<9.2 x $10^{-11}$). Nasal samples of the unaffected controls are more similar

195    to each other than to the asthmatic individuals (p-value<2.3 x $10^{-6}$).

196

197

198

199

200

201 **Supplemental Figure S18. Combinatorial diversity of T cell receptor beta (TCRB) locus**

202 **differentiates disease status.**

203 (a) Heat map depicting the percentage of RNA-Seq samples supporting of particular VJ

204 combination for whole blood, nasal epithelium of healthy controls and of asthmatic

205 individuals. Each row corresponds to a V gene and each column corresponds to a J gene.

206 (b) Alpha diversity is measured using the Shannon entropy incorporating the total number

207 of VJ combinations and their relative proportions. The nasal epithelium of asthmatic

208 individuals exhibits a decrease in combinatorial diversity of IGK locus compared to that of

209 healthy controls (p-value = $4.0 \times 10^{-2}$). (c) Compositional similarities between the samples

210 in terms of gain or loss of VJ combinations of IGK locus are measured using the Sørensen–

211 Dice index across pairs of sample from the same group (Asthma, Controls) and pairs of

212 sample from different groups (Asthma versus Controls). Lower level of similarity is

213 observed between nasal samples of asthmatic individuals compared to unaffected

214 controls (p-value < $9.4 \times 10^{-5}$). Nasal samples of unaffected controls are more similar to

215 each other than to the asthmatic individuals (p-value < $7.4 \times 10^{-4}$).

216

217

218

219

220

(a) **Blood** **Nasal-Controls** **Nasal-Asthma**

**Supplemental Figure S19. Combinatorial diversity of T cell receptor gamma (TCRG) locus differentiates disease status.**
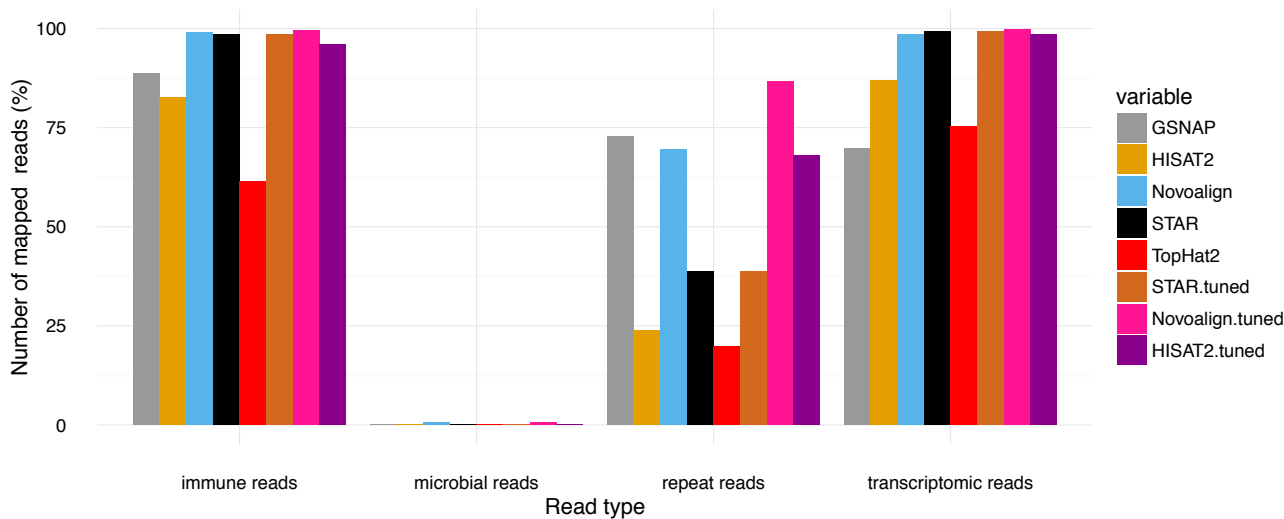
(a) Heat map depicting the percentage of RNA-Seq samples supporting of a particular VJ combination for whole blood, nasal epithelium of healthy controls and asthmatic individuals. Each row corresponds to a V gene and each column corresponds to a J gene. (b) Alpha diversity is measured using the Shannon entropy incorporating the total number of VJ combinations and their relative proportions. Nasal epithelium of asthmatic individuals exhibits decreased combinatorial diversity of IGK locus compared to that of healthy controls (p-valu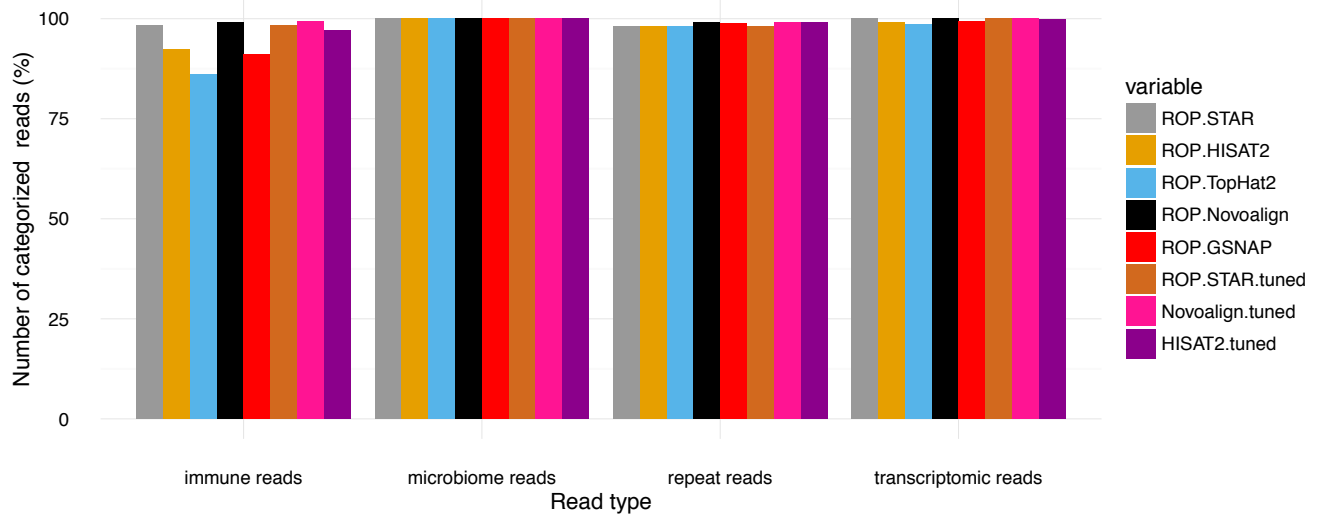e = $1.2 \times 10^{-2}$, ANOVA). (c) Compositional similarities between the samples in terms of gain or loss of VJ combinations of IGK locus are measured using the

232   Sørensen–Dice index across pairs of sample from the same group (Asthma, Controls) and

233   pairs of sample from different groups (Asthma versus Controls). Lower level of similarity

234   is observed between nasal samples of asthmatic individuals compared to unaffected

235   controls (p-value $< 1.3 \times 10^{-8}$). Nasal samples of unaffected controls are more similar to

236   each other than to the asthmatic individuals (p-value $< 8.2 \times 10^{-6}$).
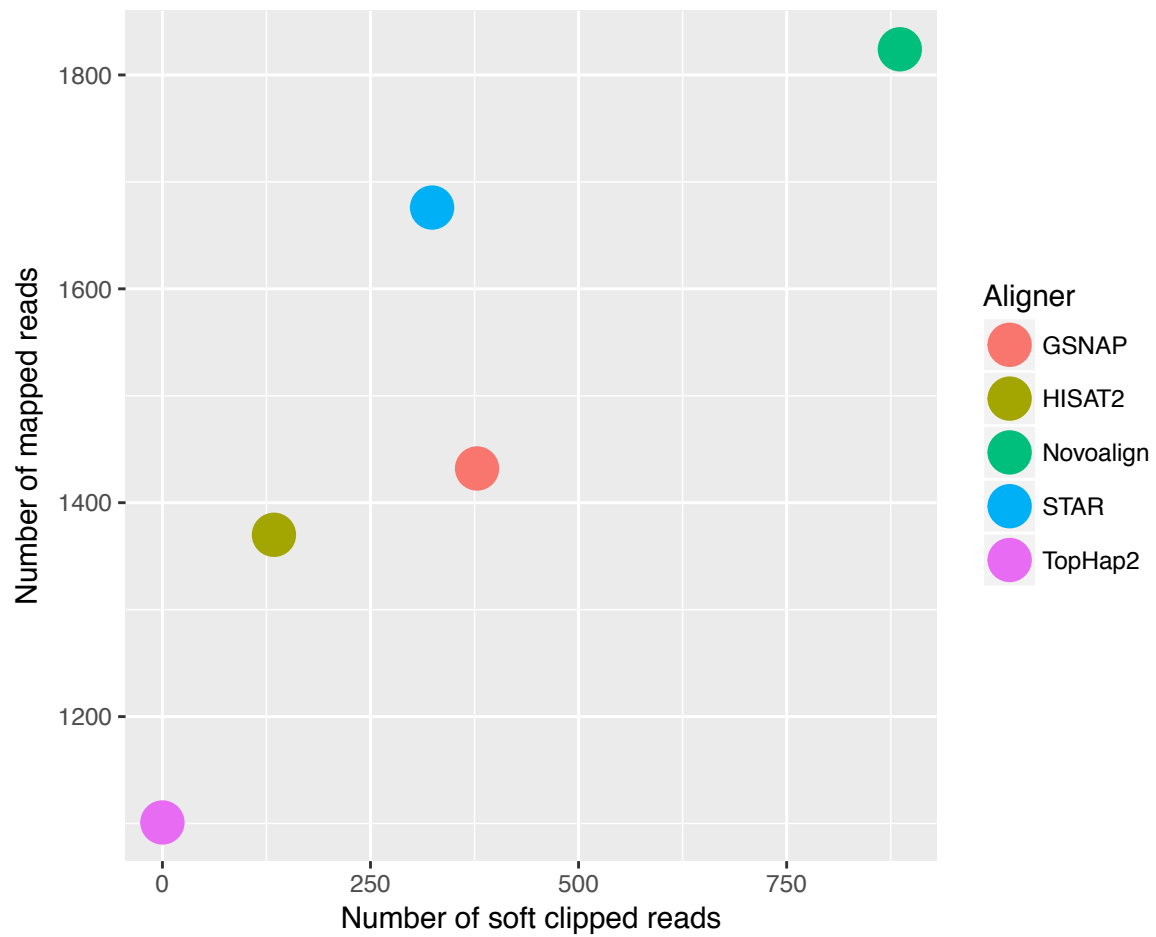
237

238

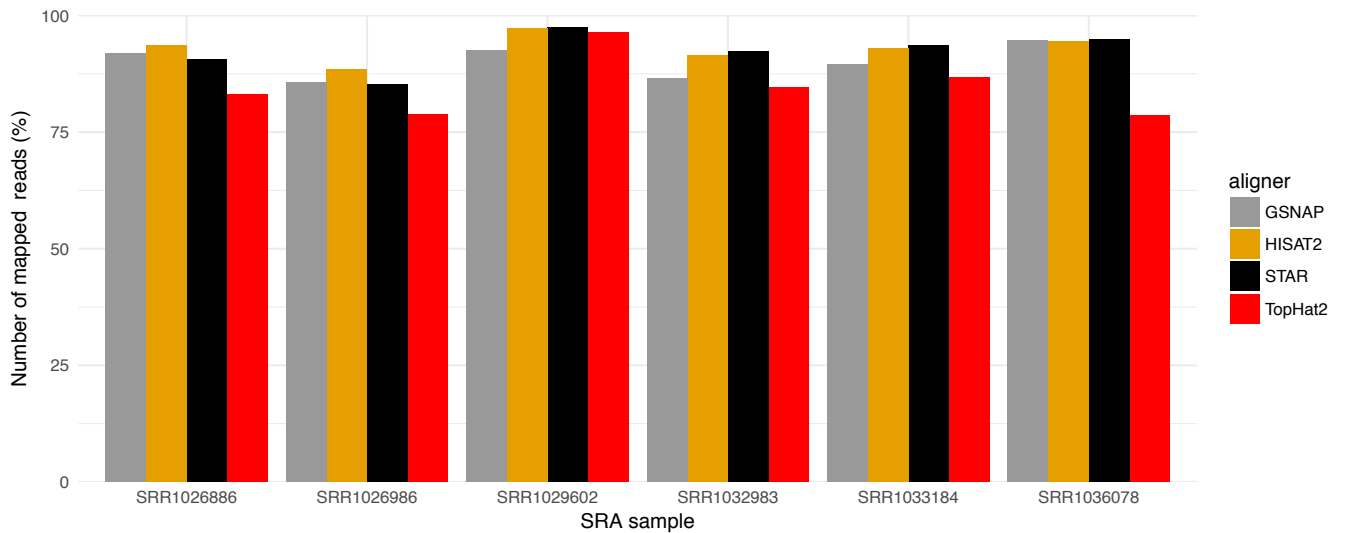239   (a) Number of reads categorized by RNA-Seq aligners



240

241   (b) Number of reads categorized by ROP using different RNA-Seq aligners

242

**Supplemental Figure S20.** The effect of RNA-Seq aligner on the fraction of reads

accounted by ROP. Percentages are calculated from the total number of reads in each

category. RNA-Seq aligners were run with default and optimized (tuned) parameters.

We use tuned setting recommended by Baruzzo et al. (2017). TopHat2 and GSNAP were

only run with default settings. Results are presented for simulated RNA-Seq data

composed of transcriptomic, repeat, immune, and microbial reads. (a) Percentages of

reads accounted by RNA-Seq aligners. (b) Percentages of reads categorized by ROP

across five state of the art aligners.

251

252

253

**Supplemental Figure S21.** Relationship between the number of soft clipped RNA-Seq

reads (partially mapped reads) and the total number of reads. Results are presented for

simulated RNA-Seq data composed of transcriptomic, repeat, immune, and microbial

reads.

254

255

256

257

258

259

260

261

262  **Supplemental Figure S22.** Number of the RNA-Seq reads mapped to the human

263  reference genome across five state-of-the-art RNA-Seq aligners. Number of mapped

264  reads is separately reported for each SRA sample. Percentages are calculated from the

265  total number of reads. Results are presented for 10 randomly selected SRA RNA-Seq

266  samples. Tools were run with default parameters. Novoalign was excluded from this

267  analysis because none of the experiments finished running within 24 hours.

268