1    Comprehensive analysis of RNA-sequencing to find the source of 1 trillion reads across

2    diverse adult human tissues

3    Supplementary Materials

4

## Table of Contents

83    *Supplementary Figures*

84

85

86



87

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

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

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

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

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

93    averages across 10641 samples.

94

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

**Supplemental Figure S2. Profile of repeat elements across 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 10641 samples.

# GTEx DNA repeats



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

**109** **Supplemental Figure S3. Profile of DNA repeats based on repeat sequences inferred from**

**110** **mapped and unmapped reads (lost repeat reads).**

**111** ROP identifies and categorizes DNA repetitive sequences among the mapped and

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

**113** instances prepared from RepeatMasker annotation (Repeatmasker v3.3, Repeat Library

**114** 20120124). Lost repeat reads are unmapped RNA-Seq reads aligned onto the reference

**115** repeat sequences (prepared from Repbase v20.07). The percentages are the averages

**116** across 10641 samples.

**117**

# SINE–VNTR–*Alu* (SVA)



SVA_F 10%
SVA_E 4%
SVA_A 7%
SVA_B 13%
SVA_C 16%
SVA_D 50%

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

118

119

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

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

122  *categorizes SVA retrotransposons sequences among the mapped and unmapped reads.*

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

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

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

126  *(prepared from Repbase v20.07). The percentages are the averages across 10641 samples.*

127

128

*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).

136

137

138



Average Number of SVA-F reads across Tissue

139

140

*Supplemental Figure 6.  Average number of SVA-F reads across GTEx tissues.* ROP identifies and categorizes *SVA* retrotransposons 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). Among the GTEx tissues, *testis* showed significantly higher expression of SVA F retrotransposons compared to other tissues ($\mathbf{p} = \mathbf{2.46 \times 10^{-33}}$)

**Alu and L1 co-expression in Individual Tissues**

$y = 0.7324x - 2118.5$
$R^2 = 0.76151$

149

150

151    *Supplemental Figure 7.* Co-expression of *Alu* and L1 elements across GTEx tissues. ROP

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

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

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

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

156    (prepared from Repbase v20.07).

157

158

159

**Hyper-edited reads** A.

**Hyper-editing reads** B.

161    Supplemental Figure S8. Distribution of hyper-edited reads.

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

163    A-to-G, indicating a high level of specificity for the hyper-editing screen. The 1,613,213

164    detected A-to-G reads contain 10,666,458 editing events (3,157,685 unique editing-sites).

165    **B.** Hyper-editing identified in the GTEx RNA-Seq data. Results showed that 80% of the reads

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

167    201,676,069 detected A-to-G reads contain 1,130,591,911 editing events (690,386,562

168    unique editing-sites).

169



170

B.



Sequence context of detected hyper-editing sites

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

176    **the detected hyper-edited A-to-G sites.**

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

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

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

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

181    RNA-Seq data (A) and GTEx RNA-Seq data (B) using the hyper-editing pipeline and human

182    editing-sites from the RADAR database.

183

184

185

186

Supplemental Figure S10. Distribution of non-co-linear (NCL) events across *across 10641*

*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.

Trans-splicing — Gene fusions — CircRNAs

P=8x10$^{-4}$ P=5x10$^{-8}$ P=3x10$^{-12}$ P=3x10$^{-12}$

* poly(A) enrichment
**ribo-depletion

194

**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.

195
196
197
198
199
200
201
202
203
204
205
206
207

208

209

210

211

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

213     are calculated from the total number of reads. Reads arising from trans-splicing, gene

214     fusion and circRNA events are captured by a TopHat-Fusion and CIRCexplorer2 tools and

215     reported a NCL reads.

216

217

218

219

220

221



222

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

224     *obtained by ROP protocol from GTEx RNA-Seq sample of EBV-transformed lymphoblastoid*

225     *cell lines (LCLs).*

226

227

228

richness__IGK_VJ

Artery - Tibial
Brain - Hypothalamus
Cells - EBV-transformed lymphocytes
Muscle - Skeletal
Skin - Sun Exposed (Lower leg)
Brain - Caudate (basal ganglia)
Brain - Substantia nigra
Nerve - Tibial
Brain - Amygdala
Heart - Left Ventricle
Brain - Nucleus accumbens (basal ganglia)
Brain - Anterior cingulate cortex (BA24)
Brain - Hippocampus
Brain - Putamen (basal ganglia)
Brain - Frontal Cortex (BA9)
Brain - Spinal cord (cervical c-1)
Brain - Cortex
Brain - Cerebellum
Brain - Cerebellar Hemisphere
Heart - Atrial Appendage
Thyroid
Adipose - Subcutaneous
Cells - Transformed fibroblasts
Pituitary
Esophagus - Muscularis
Skin - Not Sun Exposed (Suprapubic)
Esophagus - Gastroesophageal Junction
Pancreas
Colon - Sigmoid
Ovary
Whole Blood
Testis
Adrenal Gland
Prostate
Kidney - Cortex
Uterus
Breast - Mammary Tissue
Artery - Coronary
Artery - Aorta
Bladder
Adipose - Visceral (Omentum)
Liver
Vagina
Lung
Esophagus - Mucosa
Cervix - Ectocervix
Fallopian Tube
Stomach
Cervix - Endocervix
Small Intestine - Terminal Ileum
Colon - Transverse
Minor Salivary Gland
Spleen

0    50    100    150

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

231    *chain.*

232

233

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

236    *chain.*

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

* poly(A) enrichment

**ribo-depletion

257

**Supplemental Figure S16. Combinatorial diversity of immunoglobulin kappa locus (IGK) locus across in-house tissues.**

Samples were prepared by poly(A) selection (whole blood and nasal epithelium) and ribo-depletion (lung epithelium) protocols. The combinatorial diversity of IGK locus is determined based on the recombinations of the VJ gene segments. Shannon entropy measures the alpha diversity by incorporating the total number of VJ combinations and their relative proportions. Mean alpha diversity for blood samples was 4.2, for nasal samples, was 2.5, and for lung, was 1.0.

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, %

266

267    Supplemental Figure S17.  Association between microbial load and immune diversity.

268    (a) Scatterplot of the viral load and combinatorial immune diversity of IGK locus. Pearson

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

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

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

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

273

Supplemental Figure S18. 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=$5.9 \times 10^{-3}$) (c) Compositional similarities between the samples in terms of gain or loss of VJ combinations of IGK locus are measured using the Sørensen–

284    Dice index across pairs of samples from the same group (Asthma, Controls) and pairs of

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

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

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

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

289

290

291

292



293

294     **Supplemental Figure S19. Combinatorial diversity of T cell receptor beta (TCRB) locus**

295     **differentiates disease status.**

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

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

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

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

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

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

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

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

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

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

306     observed between nasal samples of asthmatic individuals compared to unaffected controls

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

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

309

310

311

312

313

314

Supplemental Figure S20. 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-value = $1.2 \times 10^{-2}$, ANOVA). (c) Compositional similarities between the samples in terms of gain or loss of VJ combinations of IGK locus are measured using the

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

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

327    observed between nasal samples of asthmatic individuals compared to unaffected controls

328    (p-value < 1.3 x $10^{-8}$,). Nasal samples of unaffected controls are more similar to each other

329    than to the asthmatic individuals (p-value < 8.2 x $10^{-6}$).


330


331


332


333


334


335


336


337


338


339


340

341     *Supplemental Tables*

342         ***Supplemental Table S1. RNA-Seq datasets  overview.*** in-house RNA-Seq data (n=86)

343     from the peripheral blood, nasal, and large airway epithelium of asthmatic and control

344     individuals (S1); (2) multi-tissue RNA-Seq data from Genotype-Tissue Expression (GTEx v6)

345     from 53 human body sites (Consortium & others, 2015) (n=8555) (S2); (3) randomly

346     selected RNA-Seq samples from the Sequence Read Archive (SRA) (n=2000) (S3). Unless

347     otherwise noted, we reported percentage of reads averaged across 3 datasets. For

348     counting purposes, the pairing information of the reads is disregarded, and each read from

349     a pair is counted separately.

350

351

352

| *Datasets* | S1 | S2 | S3 |
|---|---|---|---|
| *Number of samples* | 87 | 8555 | 1000 |
| *Read length* | 100bp | 76bp | 25-100bp |
| *Average number of  reads per sample, (million reads)* | 88.8 | 54.6 | 90.2 |
| *Percentage of mapped reads (%)* | 83.8% | 88.2% | 77.2% |

353

354

355

356

357

358    **Supplemental Table S2. Genomic profile of unmapped reads reported for each dataset (S1,**

359    **S2, S3).** Percentage of unmapped reads for each category is calculated as a fraction from

360    the total number of reads. Bars of the plot are not scaled. Human reads (black color)

361    mapped to reference genome and transcriptome via TopHat2. (a) Low quality/low-

362    complexity (light brown) and reads matching rRNA repeating unit (dark brown) were

363    excluded. (b) Hyper-edited reads are captured by hyper-editing pipeline proposed in

364    (Porath et al., 2014). (c) ROP identifies lost human reads (red color) from unmapped reads

365    using a more sensitive alignment. (d) ROP identifies lost repeat sequences (green color) by

366    mapping unmapped reads onto the reference repeat sequences. (e) Reads arising from

367    trans-spicing, gene fusion and circRNA events (orange color) are captured by a TopHat-

368    Fusion and CIRCexplorer2 tools. (f) IgBlast is used to identify reads spanning B and T cell

369    receptor gene rearrangement in the variable domain (V(D)J recombinations) (violet color).

370    (g) Microbial reads (blue color) are captured by mapping the reads onto the microbial

371    reference genomes.

372

373

374

375

376

377

378

|  | S1 | S2 | S3 | Averaged across 3 datasets |
|---|---|---|---|---|
| *Mapped* | 83.2% | 88.2% | 77.2% | 82.9% |
| *Unmapped* | 17% | 11.8% | 23% | 17.1% |
| *Low quality reads* | 4.8% | 7.0% | 9% | 7.0% |
| *rRNA repeat* | 3.8% | 0.1% | 3% | 2.4% |
| *Lost human reads* | 6.0% | 3.7% | 8% | 5.7% |
| *Hyper-edited reads* | 0.02% | 0.1% | 0.1% | 0.1% |
| *Lost repeat reads* | 0.3% | 0.1% | 0.1% | 0.2% |
| *NCL RNA* | 0.3% | 0.3% | 0.4% | 0.3% |
| *V(D)J recombinations* | 0.01% | 0.03% | 0.01% | 0.02% |
| *Microbial reads* | 1.5% | 0.5% | 2.3% | 1.4% |
| *Unaccounted reads* | 0.18% | 0.09% | 0.10% | 0.12% |

379

380

381

384    Supplemental Table S3. Relative genomic abundance of microbial taxa at different levels

385    of taxonomic classification after removal of reads with human origin (average over all

386    samples of tissues).

387    Taxonomic classification is performed using MetaPhlAn2, which is able to assign the

388    filtered unmapped reads to the microbial marker genes.

| Tissue | Whole blood | Nasal epithelium | Lung epithelium |
|---|---|---|---|
| N | 19 | 19 | 49 |
| Library preparation method | poly(A) enrichment | poly(A) enrichment | ribo-depletion |
| **Phylum** | | | |
| Proteobacteria | 0.0% | 0.9% | 100.0% |
| Actinobacteria | 0.0% | 99.1% | 0.0% |
| **Class** | | | |
| Betaproteobacteria | 0.0% | 0.5% | 86.7% |
| Gammaproteobacteria | 0.0% | 0.5% | 13.3% |
| Actinobacteria | 0.0% | 98.9% | 0.0% |
| **Order** | | | |
| Burkholderiales | 0.0% | 0.0% | 87.0% |
| Enterobacteriales | 0.0% | 0.0% | 12.0% |
| Actinomycetales | 0.0% | 99.5% | 0.0% |
| Pseudomonadales | 0.0% | 0.5% | 1.0% |

389

390    **Supplementary Table S4. Parameters for each RNA-Seq aligner for default, sensitive, and**

391    **very sensitive settings.**

392    Sensitive setting has more relaxed parameters for filtering.

| | Default | Sensitive | Very Sensitive |
|---|---|---|---|
| Topha t | -D 10 -R 2 -N 0 -L 22 -i S,0,2.50 | -D 15 -R 2 -L 22 -i S,1,1.15 | -D 20 -R 3 -N 0 -L 20 -i S,1,0.50 |
| STAR | -- seedNoneLociPerWindo w 10 — outFilterMismatchNmax 10 — seedPerReadMax 1000 | -- seedNoneLociPerWindo w 15 -- outFilterMismatchNmax 15 --seedPerReadNmax 1500 | -- seedNoneLociPerWindo w 15 -- outFilterMismatchNmax 15 --seedPerReadNmax 1500 --twopassMode Basic |

393

394    **Supplementary Table S5. Average mapping rate for different aligners with different**

395    **mapping settings.**

396    The average rate is noted, and the standard deviation is noted within parenthesis.

| | Default/Fast | Sensitive | Very Sensitive |
|---|---|---|---|
| Tophat | 89.06% (3.84) | 89.22% (3.51) | 89.18% (3.62) |
| STAR | 80.86% (9.22) | 81.70% (9.25) | 81.74% (9.35) |

397

*Supplemental Methods*

<u>In-house RNA-Seq data</u>

400    *Subject Recruitment*

401    **Poly(A) selected RNA-Seq samples (n=38).** In this analysis, we used a subset of Puerto Rican

402    Islanders recruited as part of the on-going Genes-environments & Admixture in Latino

403    Americans study (GALA II) (Anders, Pyl, & Huber, 2014; Jin, Tam, Paniagua, & Hammell,

404    2015; Melé et al., 2015; Tarailo-Graovac & Chen, 2009). We classified asthma by physician

405    diagnosis and the presence of at least two symptoms (wheezing, coughing, or shortness of

406    breath) during 2 years prior to the enrollment. All study subjects had no history of smoking

407    or recent (within 4 weeks of recruitment) nasal steroid use. The study was approved by

408    local institutional review boards, and written assent/consent was received from all subjects

409    and, if applicable, parents of subjects under the age of legal consent.

410    **Ribo-Zero RNA-Seq samples (n=49).** Via community-based advertising, we recruited adults

411    aged 18-70 years to participate in a study, in which they underwent research

412    bronchoscopy. The study was approved by the University of California at San Francisco

413    Committee on Human Research. Written informed consent was obtained from all subjects,

414    and all studies were performed in accordance with the principles expressed in the

415    Declaration of Helsinki.

416

*Sample Collection*

418 **Poly(A) selected RNA-Seq samples (n=38).** Methods for nasal epithelial cell collection and

419 processing are described in Poole et al. (Tarailo-Graovac & Chen, 2009). Briefly, nasal

420 epithelial cells were collected from behind the inferior turbinate with a cytology brush

421 using a nasal illuminator. The collected brush was submerged in a mixture of RLT Plus lysis

422 buffer and beta-mercaptoethanol, and frozen at -80 C until extraction was performed with

423 a Qiagen Allprep RNA/DNA extraction kit (Qiagen, Valencia, CA). We collected 10ml of

424 whole blood using PAXgene RNA blood tubes (PreAnalytiX, Valencia, CA) and isolated RNA

425 using PAXgene RNA blood extraction kits, according to the manufacturers' protocol.

426 Portions of the nasal airway epithelial whole transcriptome data were published in a

427 previous manuscript (Tarailo-Graovac & Chen, 2009).

428 **Ribo-Zero RNA-Seq samples (n=49).** During bronchoscopy airway epithelial brushings,

429 samples were obtained from $3^{rd}$-$4^{th}$ generation bronchi. RNA was extracted from the

430 epithelial brushing samples using the Qiagen RNeasy mini-kit (Qiagen, Valencia, CA),

431 according to manufacturer's protocol.

432

433 *Whole Transcriptome Sequencing*

434 **Poly(A) selected RNA-Seq samples (n=38).** We constructed Poly-A RNA-seq libraries using

435 500 ng of blood and nasal airway epithelial total RNA from 9 atopic asthmatics and 10 non-

436 atopic controls. Libraries were constructed and barcoded with the Illumina TruSeq RNA

437 Sample Preparation v2 protocol. Barcoded nasal airway RNA-seq libraries from each of the

438 19 subjects were pooled and sequenced as 2 x 100bp paired-end reads across two flow

439    cells of an Illumina HiSeq 2000. Barcoded blood RNA-seq libraries from each of the 19

440    subjects were pooled and sequenced as 2 x 100bp paired end reads across 4 lanes of an

441    Illumina Hiseq 2000 flow cell.

442    **Ribo-Zero RNA-Seq samples (n=49).**   We used 100ng of isolated RNA from a total of 61

443    samples to construct ribo-depleted RNA-seq libraries using the TruSeq Stranded Total RNA

444    with Ribo-Zero Human/Mouse/Rat library preparation kit, per manufacturer's protocol.

445    Barcoded bronchial epithelial RNA-seq libraries were multiplexed and sequenced as 2 x

446    100bp paired end reads on an Illumina HiSeq 2500. On average, 37 million reads were

447    generated per sample. We excluded 12 samples from further analyses due to high

448    ribosomal RNA read counts (library preparation failure), leaving a total of 49 samples

449    suitable for further analyses.

450

451    <u>GTEx RNA-Seq data</u>

452    We used RNA-Sequencing data from Genotype-Tissue Expression study (GTEx Consortium

453    v.6) corresponding to 8,555 samples collected from 544 individuals from 53 tissues

454    obtained from Genotype-Tissue Expression study (GTEx v6). RNA-Seq data is from Illumina

455    HiSeq sequencing of 75 bp paired-end reads. The data was derived from 38 solid organ

456    tissues, 11 brain subregions, whole blood, and three cell lines of postmortem donors. The

457    collected samples are from adults matched for age across males and females.   We

458    downloaded the mapped and unmapped reads in BAM format from dbGap

459    (http://www.ncbi.nlm.nih.gov/gap).

460

461    SRA RNA-Seq data

462

463    Samples (n=2000) were randomly selected using SQLite database from R/Bioconductor

464    package SRAdb (https://bioconductor.org/packages/release/bioc/html/SRAdb.html). We

465    have                used               a               script               from

466    https://github.com/nellore/runs/blob/master/sra/define_and_get_fields_SRA.R to select

467    run_accessions from the sra table with platform = 'ILLUMINA', library_strategy = 'RNA-

468    Seq', and taxon_id = 9606 (human).

469

470    *Workflow to categorize the mapped reads*

471    *Map reads onto human genome and transcriptome*

472    We mapped reads onto the human transcriptome (Ensembl GRCh37) and genome

473    reference (Ensembl hg19) using tophat2 (v 2.0.13) with the default parameters. Tophat2

474    was supplied with a set of known transcripts (as a GTF formatted file, Ensembl GRCh37)

475    using –G option. The mapped reads of each sample are stored in a binary format (.bam).

476

477    *Categorize mapped reads into genomic categories*

478    ROP categorizes the reads into genomic categories based on the compatibility of each read

479    from the pair with the features defined by Ensembl (GRCh37) gene annotations. First, we

480    determined CDS, UTR3, UTR5 coordinates. We downloaded annotations for CDS, UTR3,

481    UTR5 from UCSC Genome Browser (http://genome.ucsc.edu/cgi-bin/hgTables) in BED

482 (browser extensible data) format. Next, we used gene annotations (a GTF formatted file,

483 Ensembl GRCh37) to determine intron coordinates and inter-genic regions. We defined

484 two types of inter-genic regions: '(proximate) inter-genic' region (1Kb from the gene

485 boundaries) and 'deep inter-genic' (beyond a proximity of 1Kb from the gene boundaries).

486

487 Next, we checked the compatibility of the mapped reads with the defined genomic

488 features, as follows:

489

490     a.  Read mapped to multiple locations on the reference genome is categorized

491         as a multi-mapped read.

492     b.  Read fully contained within the CDS, intron, UTR3, or UTR5 boundaries of a

493         least one transcript is classified as a CDS, intronic, UTR3, or UTR5,

494         respectively.

495     c.  Read simultaneously overlapping UTR3 and UTR5 regions is classified as a

496         UTR read.

497     d.  Read spanning exon-exon boundary is defined as a junction read.

498     e.  Read mapped outside of gene boundaries and within a proximity of 1Kb is

499         defined as a (proximal) inter-genic read.

500     f.  Read mapped outside of gene boundaries and beyond the proximity of 1Kb

501         is defined as a deep inter-genic read.

502     g.  Read mapped to mitochondrial DNA (MT tag in hg19) is classified as a

503         mitochondrial read.

504               h.   Reads from a pair mapped to different chromosomes are classified as a

505               fusion read.

506  Scripts to categorize mapped reads into genomic categories are distributed with ROP

507  protocol.

508

509  *Categorize mapped reads overlapping repeat instances*

510  Mapped reads were categorized based on the overlap with the repeat instances defined

511  by RepeatMasker annotation (Repeatmasker v3.3, Repeat Library 20120124).

512  RepeatMasker masks the repeats using the RepBase library:

513  (http://www.girinst.org/repbase/update/index.html), which contains prototypic

514  sequences representing repetitive DNA from different eukaryotic species. We use GTF files

515  generated from the RepeatMasker annotations by Jin, Ying, et al. (Jin et al., 2015) and

516  downloaded from:

517  http://labshare.cshl.edu/shares/mhammelllab/www-

518  data/TEToolkit/TE_GTF/hg19_rmsk_TE.gtf.gz

519

520  Following Melé, Marta, et al. (Melé et al., 2015), repeat elements overlapping CDS regions

521  are excluded from the analysis. We filtered out 6,873 repeat elements overlapping CDS

522  regions. Prepared repeat annotations (bed formatted file) are available at

523  https://drive.google.com/file/d/0Bx1fyWeQo3cORi1UNWhxOW9kYUk/view?pref=2&pli=

524  1

525

526    The prepared repeat annotations contain 8 Classes and 43 Families.  Number of elements

527    per family and class represented below (Supplemental Methods Table SM1):

528

| classID | N |
|---|---|
| DNA | 458223 |
| LINE | 1478382 |
| LTR | 707384 |
| RC | 2226 |
| SVA | 3582 |
| RNA | 717 |
| Satellite | 8950 |
| SINE | 1765403 |

529

530    **Supplemental Methods Table SM1. Number of repeat elements per class.** Repeat instances

531    are defined by RepeatMasker (RepeatMasker v3.3, Repeat Library 20120124) based on

532    RepBase library. RepBase library contains prototypic sequences representing repetitive

533    DNA from different eukaryotic species.

534

| familyID | n |
|---|---|
| acro | 44 |
| Alu | 1173282 |

| | |
|---|---|
| centr | 2272 |
| CR1 | 60577 |
| Deu | 1262 |
| DNA | 4609 |
| Dong-R4 | 554 |
| ERV | 579 |
| ERV1 | 172612 |
| ERVK | 10446 |
| ERVL | 159606 |
| ERVL-MaLR | 343266 |
| Gypsy | 18553 |
| hAT | 15418 |
| hAT-Blackjack | 19578 |
| hAT-Charlie | 251618 |
| hAT-Tip100 | 30204 |
| Helitron | 2226 |
| L1 | 937636 |
| L2 | 461296 |
| LTR | 2322 |
| Merlin | 55 |
| MIR | 589496 |
| MuDR | 1978 |

| | |
|---|---|
| Penelope | 51 |
| PiggyBac | 2352 |
| RNA | 717 |
| RTE | 17617 |
| RTE-BovB | 651 |
| Satellite | 6247 |
| SINE | 1363 |
| SVA_A | 257 |
| SVA_B | 465 |
| SVA_C | 279 |
| SVA_D | 1358 |
| SVA_E | 232 |
| SVA_F | 991 |
| TcMar | 5354 |
| TcMar-Mariner | 16253 |
| TcMar-Tc2 | 8098 |
| TcMar-Tigger | 102706 |
| telo | 387 |

535

536    **Supplemental Methods Table SM2. Number of repeat elements per family.** Repeat

537    instances are defined by RepeatMasker (RepeatMasker v3.3, Repeat Library 20120124)

538    based on RepBase library.

539

540　We determined the coordinates of repeat elements (*class_id* and *family_id attributes from*

541　*the GTF file*) from the repeat annotations. Next, we checked the compatibility of the

542　mapped reads with the repeat instances. We disregarded the pairing information for the

543　unmapped reads and count each end as a separate read. Reads entirely mapped to the

544　corresponding repeat instance are counted. Scripts to categorize mapped reads based on

545　the overlap with the repeat instances are distribuited with ROP protocol.

546

547　*Categorize mapped reads overlapping B cell receptor (BCR) and T cell receptor (TCR) loci*

548　We used the gene annotations (Ensembl GRCh37) to extract BCR and TCR genes. We

549　extracted gene annotations of the 'constant' (labeled as IG_C_gene, Ensembl GRCh37),

550　'variable' (labeled as IG_V_gene, Ensembl GRCh37), 'diversity' (labeled as IG_D_gene,

551　Ensembl GRCh37), and 'joining' genes (labeled as IG_J_gene, Ensembl GRCh37) of BCR and

552　TCR loci.  We excluded the BCR and TCR pseudogenes (labeled as IG_C_pseudogene,

553　IG_V_pseudogene,　　IG_D_pseudogene,　　IG_J_pseudogene,　　TR_C_pseudogene,

554　TR_V_pseudogene, TR_D_pseudogene, and TR_J_pseudogene). In addition, we excluded

555　the patch contigs *HG1592_PATCH and HG7_PATCH,* as they are not part of the Ensembl

556　hg19 reference, and reads are not mapped on the patch contigs by high throughput

557　aligners.  After following the filtering steps described above, we extracted a total of 386

558　immune genes: 207 BCR genes and 179 TCR genes.  The gene annotations for antibody

559　genes　　　　(GTF　　　　formatted　　　　file)　　　are　　　available　　　at

560　https://drive.google.com/file/d/0Bx1fyWeQo3cObFZNT3kyQlZUS1E/view?pref=2&pli=1

561

562 The number of VDJ genes per locus is reported in the Table 3.

563

|  | C domain | V domain | D domain | J domain |
|---|---|---|---|---|
| *IGH* locus | 8 | 55 | 38 | 6 |
| *IGK* locus | 1 | 46 | - | 5 |
| *IGL* locus | 4 | 37 | - | 7 |
| TCRA locus | 1 | 46 | - | 57 |
| TCRB locus | 1 | 39 | 0 | 8 |
| TRG locus | 2 | 9 | - | 5 |
| TRD locus | 1 | 3 | 11 | 4 |

564

565 **Supplemental Methods Table SM3. The number of VDJ genes for each antibody chains.**

566 Antibody genes were extracted from the gene annotations (Ensembl GRCh37).

567

568 The list of the genes encoding the C region of the BCR and TCR chains is presented in

569 Supplemental Methods Table SM4.

570

| Name of the chain | Genes encoding for the C region of the chain |
|---|---|
| IG@ locus | |
| α heavy IG chain | IGHA1, IGHA2 |

| | |
|---|---|
| δ heavy IG chain | IGHD |
| γ heavy IG chain | IGHG1, IGHG2, IGHG3, IGHG4 |
| ε heavy IG chain | IGHE |
| μ heavy IG chain | IGHM |
| κ light IG chain | IGKC |
| λ light IG chain | IGLC1, IGLC2, IGLC3, IGLC7 |
| TCR@ locus | |
| α TCR chain | TRAC |
| B TCR chain | TRBC2 |
| γ TCR chain | TRGC1, TRGC2 |
| δ TCR chain | TRDC |

571

572 **Supplemental Methods Table SM4. List of the genes encoding the C region of the BCR and**

573 **TCR chains.** Genes were extracted from the gene annotations (Ensembl GRCh37).

574

575 The number of reads mapping to each C-V-D-J genes was *obtained by counting the number*

576 *of* sequencing reads that align, with high confidence, to each of the genes (HTSeq v0.6.1)

577 (Anders et al., 2014). Script "htseq-count" is supplied with the gene annotations for BCR

578 and TCR genes (genes_Ensembl_GRCh37_BCR_TCR.gtf) and a bam file. The bam file

579 contains reads mapped to the human genome and transcriptome using tophat2 (See

580 Section *"Map reads onto human genome and transcriptome"* for details). The script

581 generates individual gene counts by examining the read compatibility with BCR and TCR

582 genes. We chose a conservative setting (--mode=intersection-strict) to handle reads

583 overlapping more than one feature. Thus, a read overlapping several genes simultaneously

584 is marked as a read with no feature and is excluded from the consideration.

585

586 *Workflow for categorizing the unmapped reads*

587 We first converted the unmapped reads saved by tophat2 from a BAM file into a FASTQ

588 file (using bamtools). The FASTQ file of unmapped contain full read pairs (both ends of a

589 read pair were unmapped) and discordant read pairs (one read end was mapped while the

590 other end was unmapped). We disregarded the pairing information of the unmapped reads

591 and categorize unmapped reads using the following steps:

592

593 *A. Quality Control*

594 Low quality reads, defined as reads that have quality lower than 30 in at least 75% of their

595 base pairs, were identified by FASTX (v 0.0.13).  Low complexity reads, defined as reads

596 with sequences of consecutive repetitive nucleotides, are identified by SEQCLEAN.  As a

597 part of the quality control, we also excluded unmapped reads mapped onto the rRNA

598 repeat sequence (HSU13369 Human ribosomal DNA complete repeating unit) (BLAST+

599 2.2.30). We prepared the index from rRNA repeat sequence using makeblastdb and

600 makembindex from BLAST+.  We used the following command for makeblastdb:

601    ➢ makeblastdb -parse_seqids -dbtype nucl -in <fasta file>.

602 We used the following command for makembindex:

603    ➢ makembindex -input <fasta file> -output <index> -iformat blastdb

604

605    *B. Mapping unmapped reads onto the human references.*

606    We remapped the unmapped reads to the human reference sequences using Megablast

607    (BLAST+ 2.2.30). We mapped reads onto the following references:

608        • Reference transcriptome (known transcripts), Ensembl GRCh37

609        • Reference genome, hg19 Ensembl

610    We prepared the index from each reference sequence using makeblastdb and

611    makembindex. We mapped the reads separately onto each reference in the order listed

612    above. Reads mapped to the reference genome and transcriptome were merged into a

613    'lost human reads' category. The following options were used to map the reads using

614    Megablast: for each reference: task = megablast, use_index = true, perc_identity = 90,

615    outfmt = 6, max_target_seqs =1, e-value = $1e^{-05}$.

616

617    *C. Identification of hyper-edited reads*

618    We        have        used        hyper-editing        pipeline        (HE-pipeline

619    http://levanonlab.ls.biu.ac.il/resources/zip), which is capable of identifying hyper-edited

620    reads.  When running HE-pipeline, additional changes can be made to parallelize the scripts

621    for use with UCLA's Hoffman2 cluster. Before proceeding, follow the instructions in

622    the README that is included with the scripts to prepare the reference and provide the

623    necessary third-party tools. Ensure that the output directory is set correctly

624    in config_file.sh (it is acceptable to use a single output directory), and check that the list of

625    input files has been prepared correctly.

626

627     Details on how to run HE-pipeline are available here:

628     https://github.com/smangul1/rop/wiki/How-to-run-hyper-editing-pipeline

629

630     *D. Mapping unmapped reads onto the repeat sequences*

631     We filtered out the reads that failed QC and lost human reads. The remaining reads were

632     mapped to the reference repeat sequences.  The reference repeat sequences were

633     downloaded from Repbase v20.07 (http://www.girinst.org/repbase/). Human repeat

634     elements (humrep.ref and humsub.ref) were merged into a single reference. We prepared

635     the index from the merged repeat reference using makeblastdb and makembindex from

636     BLAST+. In total, we obtained sequences for 1,117 repeat elements. The following options

637     were used to map the reads using the Megablast: task = megablast, use_index = true,

638     perc_identity = 90, outfmt = 6, max_target_seqs = 1, e-value = $1e^{-05}$. Blast hits with

639     alignment length shorter than 80% of the read length were discarded (corresponding to

640     80bp of the 100bp read).

641

642     The repeat elements from humrep.ref and humsub.ref were classified into families and

643     classes using RepeatMasker annotations (hg19_rmsk_TE_prepared_noCDS.bed).

644     Repetitive reads identified from the unmapped reads were confirmed by directly applying

645     Repeatmasker (Tarailo-Graovac & Chen, 2009).

646

647 *E. Workflow to* detect 'non-co-linear' reads (trans-splicing, gene fusions, and circRNAs)

648

649 We divide non-co-linear reads into three categories:

650

651     1) gene fusion characterized by reads that map on different chromosomes

652     2) trans-splicing events characterized by reads that map on the same chromosome,

653        but are at least 1 Mb apart from each other

654     3) circRNAs characterized by reads that map in a head-to-tail configuration on the

655        same chromosome

656

657 To distinguish between these three categories, we make use of circExplorer2 (Zhang et al.,

658 2016), which was recently identified as one of the best tools to detect circRNAs (Hansen

659 et al., 2015). CircExplorer2 relies on Tophat-Fusion and thus allows also the monitoring

660 NCL events in the same run. TopHat-Fusion (v2.0.13, bowtie1 v0.12.9) and circExplorer2

661 (v2.2.4) were invoked with the following commands:

662

663 $ tophat2 -o tophat-output-directory -p 4 --fusion-search --keep-fasta-order --bowtie1 --

664 no-coverage-search bowtie1-index fastq-file

665

666 $ python CIRCexplorer2 parse -t TopHat-Fusion -o circrna-output-folder  tophat-output-

667 directory/accepted_hits.bam

668

669    $ python CIRCexplorer2 annotate -r ensemble-reference.txt -g genome.fa circrna-output-

670    folder

671

672    To separate potential gene and trans-fusions from the TopHat-Fusion output, we ran a

673    ruby custom script, which is part of the ROP pipeline.

674    *F. Mapping unmapped reads onto the V(D)J recombinations of B and T cell receptors*

675    Gene segments of B cell receptors (BCR) and T cell receptors (TCR) were imported from

676    IMGT      (International      ImMunoGeneTics      information      system):

677    (http://www.imgt.org/vquest/refseqh.html#V-D-J-C-sets).

678    IMGT database contains:

679    • Variable (V) gene segments

680    • Diversity (D) gene segments

681    • Joining (J) gene segments

682    Unmapped reads categorized by step (A)-(D) were filtered out. We used IgBLAST (v. 1.4.0)

683    with stringent e-value threshold (e-value < $10^{-20}$) to map the remaining high-quality

684    unmapped reads onto the V(D)J regions of the of the BCR and TCR loci.  Reference files

685    with BCR and TCR VDJ gene segments are distributed with ROP protocol and available at

686    https://drive.google.com/folderview?id=0Bx1fyWeQo3cOTkhKdHFDb3c5MjA&usp=shari
687    ng
688

689    The complete list of the references is presented in Supplemental Methods Table SM5.

| Name of the reference file | Description of the gene |
| --- | --- |

| BCR heavy chain | |
|---|---|
| IGHV.fa | V genes of BCR heavy chain |
| IGHD.fa | D genes of BCR heavy chain |
| IGHJ.fa | J genes of BCR heavy chain |
| **BCR light chains** | |
| IGLV.fa | V genes of BCR lambda chain |
| IGLJ.fa | J genes of BCR lambda chain |
| IGKV.fa | V genes of BCR kappa chain |
| IGKJ.fa | J genes of BCR kappa chain |
| **TCR chains** | |
| TCRAV.fa | V genes of TCR alpha chain |
| TCRAJ.fa | J genes of TCR alpha chain |
| TCRBV.fa | V genes of TCR beta chain |
| TCRBD.fa | D genes of TCR beta chain |
| TCRBJ.fa | J genes of TCR beta chain |
| TCRGV.fa | V genes of TCR gamma chain |
| TCRGJ.fa | J genes of TCR gamma chain |
| TCRDV.fa | V genes of TCR delta chain |
| TCRDD.fa | D genes of TCR delta chain |
| TCRDJ.fa | J genes of TCR delta chain |

690

691     Supplemental Methods Table SM5.  List of the references files prepare for V-D-J from BCR

692     and TCR loci.

693

694     We prepared the index from each reference sequence using makeblastdb and

695     makembindex from BLAST+. The following options were used to map the reads using

696     IgBLAST: -germline_db_V; germline_db_D; -germline_db_J; -organism=human; -outfmt =

697     7; –evalue = 1e-20.

698

699     The number of genes and gene alleles per antibody locus is presented in Supplemental

700     Methods Table SM6.

701

|  | V domain | D domain | J domain |
|---|---|---|---|
| *IGH* locus | **136**(370) | **27**(34) | **9**(16) |
| *IGK* locus | **100**(124) | - | **5**(9) |
| *IGL* locus | **70**(111) | - | **7**(10) |
| TCRA locus | **54**(112) | - | **61**(68) |
| TCRB locus | **77**(160) | **2**(3) | **14**(16) |
| TRG locus | **14**(26) | - | **5**(6) |
| TRD locus | **8**(22) | **0**(0) | **1**(4) |

702

703     Supplemental Methods Table SM6. The number of V-D-J genes and gene alleles per

704     antibody locus.  Number of genes is presented in bold and number of gene alleles is

705      presented in parenthesis. Gene and gene alleles of B cell receptors (BCR/IG) and T cell

706      receptors (TCR) were imported from IMGT.

707

708      We assessed combinatorial diversity of the antibody repertoire by looking at the

709      recombinations of the VJ gene segments of BCR and TCR loci. We extracted the reads

710      spanning the V-J gene boundaries.

711

712      *G. Identification of microbial reads*

713      Unmapped reads mapping in step (A -(E) were filtered out. The remaining reads were high-

714      quality non-human reads used to profile the taxonomic composition of the microbial

715      communities. We used MetaPhlAn2 (Metagenomic Phylogenetic Analysis, v 2.0) to assign

716      reads on microbial genes and to obtain a taxonomic profile. The database of the microbial

717      marker genes is provided by MetaPhlAn. We run MetaPhlAn in two stages as follow: the

718      first stage identifies the candidate microbial reads (i.e. reads hitting a marker), while the

719      second stage profiles metagenomes in terms of relative abundances – the commands used

720      are as follow:

721      ➢ metaphlan.py    &lt;fastq&gt;    &lt;map&gt;    --input_type    multifastq    --bowtie2db

722          bowtie2db/mpa -t reads_map --nproc 8 --bowtie2out

723      ➢ metaphlan.py --input_type blastout &lt;bowtie2out.txt&gt; -t rel_ab &lt;tsv&gt;

724

725      The output of the first stage is a file containing a list of candidate microbial reads with the

726      microbial taxa assigned (.map file). The second stage outputs the taxonomic profile (taxa

727 detected and its relative abundance, in tab separated format (.tsv file). We used taxa

728 detected from stage 2 to extract the reads associated with it in stage 1.

729 In addition to MetaPhlAn2 we used to create the curated database of taxa-specific genes,

730 we mapped the reads onto the entire reference genomes of microbial organisms. We used

731 Megablast (BLAST+ 2.2.30) to align reads onto the collection of bacterial, viral, and

732 eukaryotic pathogens reference genomes. Bacterial and viral genomes were downloaded

733 from NCBI [ftp://ftp.ncbi.nih.gov/](ftp://ftp.ncbi.nih.gov/) on February 1, 2015. Genomes of eukaryotic pathogens

734 were downloaded from EuPathDB database, which is available at:

735 [http://eupathdb.org/eupathdb/](http://eupathdb.org/eupathdb/).

736 The following parameters were used for the megablast alignment: e-value = $10^{-5}$,

737 perc_identity = 90. The Megablast hits shorter than 80% of the input read sequence were

738 removed (corresponding to 80bp of the 100bp read).

739

740 ___Comparing diversity across groups___

741 First, we sub-sampled unmapped reads to the number of reads corresponding to a sample

742 with the smallest number of unmapped reads. Diversity within a sample was assessed

743 using the richness and alpha diversity indices. Richness was defined as a total number of

744 distinct *events* in a sample. We used Shannon Index (SI), incorporating richness and

745 evenness components, to compute alpha diversity, which is calculated as follows:

746
$$\text{SI} = -\sum (p \times \log_2(p))$$

747    We used beta diversity (Sørensen–Dice index) to measure compositional similarities

748    between the samples in terms of gain or loss in the events.  We calculated the beta

749    diversity for each combination of the samples, and we produced a matrix of all pairwise

750    sample dissimilarities. The Sørensen–Dice beta diversity index is measured as $1 - \frac{2J}{A+B}$,

751    where J is the number of shared events, while A and B are the total number of events for

752    each sample, respectively.

753

754    *Percentage of unmapped reads calculation*

755    We calculated the percentage of unmapped reads using the following formula:

756
$$P_{unmapped} = \frac{(N_{ud} + (N_{uc} \times 2))}{(N_{total} \times 2)}$$

757    where,

758    $N_{ud}$ – number of discordant unmapped reads (one end is mapped, while the other end is

759    unmapped);

760    $N_{uc}$ – number of unmapped read pairs (both ends are unmapped);

761    $N_{total}$ – total number of read pairs (fragments).

762

763    The robustness of the ROP results against changing the thresholds for each of the ROP

764    steps

765

766    We have performed the robustness analysis to investigate the impact of the thresholds

767    used in each step of the ROP approach.  For each ROP step, we have reported number of

768    reads identified under different thresholds.  The results are presented as cumulative

769    frequency plots.

770

**a**



771

**b**



772

**c**



773

774

**d**



775

**e**



776

777

778 **Supplemental Methods Figure SM1.  Percentage of reads identified under different**

779 **threshold values**. Results are presented as cumulative frequency plots for each step of ROP.

780 ROP threshold is highlighted with red line.

781 The percentages are the averages across 87 samples.  (a)  Step 2 (Remap to human

782 references).  Cumulative frequency plot reporting the percentage of lost human reads

783 averaged across all samples (y-axis) identified under different threshold (edit distance) (x-

784 axis). Edit distance was calculated as the minimum number of operations required to

785 transform a read sequence into the corresponding reference subsequence. Reads are

786 grouped by edit distance with the transcriptome or the genome reference.  (b) Step 3 (Map

787 to repeat sequences). Cumulative frequency plot reporting the percentage of lost repeat

788 reads (y-axis) identified under different threshold averaged across (percentage identity) (x-

789 axis).  (c) Step 4 (NCl RNA profiling). Cumulative frequency plot of the percentage of NCL

790 reads averaged across all samples (y-axis) identified under different thresholds (number of

791     reads supporting NCL event) (x-axis). Results are reported separately for circRNAs, gene

792     fusions and trans-splicing events. (d) Step 5 (B and T cell receptors profiling). Cumulative

793     frequency plot reporting the percentage of immune reads averaged across all samples (y-

794     axis) identified under different threshold (e-value) (x-axis). (e) Step 6 (Microbiome

795     profiling). Cumulative frequency plot reporting the percentage of microbial reads averaged

796     across all samples (y-axis) identified under different threshold (percentage identity) (x-

797     axis). Results are reported separately for viral, bacterial and eukaryotic reads.

798

799

800     ***The impact of ROP step ordering on the read classification***

801     We have investigated the effect of the ordering on read classification. Ordering of ROP

802     steps will have an effect only when references of each step share homologous sequences.

803     For each ROP step, we have swapped its order with another ROP step. For example, we

804     considered swapping 'Remapping to human references' reads and 'QC' steps. Before

805     swapping, 'Remapping to human references' was number 2 in the queue. After swapping,

806     it became number 1.

807

808     We observed a major effect of swapping 'Remapping to human references' with all other

809     steps. For example, swapping 'Remapping to human references' and 'QC' steps results in

810     classifying 79.6% of rRNA reads as lost human reads. Similarly, swapping 'Remapping to

811     human references' and 'Microbiome profiling' steps results in classifying 0.2% of the lost

812     human reads as microbiome reads. In other words, this swap produces a 27.8% increase

813    of microbiome reads. Similarly, considering 'B and T lymphocytes profiling' prior to

814    'Remapping to human references' produces a 50.8% increase of identified immune reads.

815    Considering partial mapping of BCR and TCR reads prior to the 'Remapping to human

816    references' step may produce many false positives. Swapping other steps of ROP resulted

817    in minor effects (i.e. <1% of reads from each category were effected).

818

819

820

821    *The impact of mapping parameters and RNA-Seq aligners on the number of unmapped*

822    *reads*

823    Five samples were randomly selected among each library preparation protocol. In total,

824    we obtained ten samples for the mapping rate comparison. All selected samples were

825    aligned to the human genome (hg19) using two tools, Tophat2 and STAR, and three

826    different sensitivities for each tool – default, sensitive setting, and very sensitive setting –

827    as noted below in Supplemental Table S5. The average runtime for Tophat per million reads

828    was 2.5 hours; STAR, 0.13 hours; and Novoalign, 9.1 hours. Novoalign was not considered

829    in the analysis due to its substantially longer running time that made it infeasible for the

830    protocol.

831    The mapping rate for each tool and each setting is shown in Supplemental Table S6. The

832    mapping rate was significantly higher in Tophat when compared with STAR and using the

833    default option for each tool ($p < 0.03$). However, there is no significant difference in

834    mapping rate when comparing different mapping settings ($p > 0.92$ under two-tailed t-

835    tests for Tophat, p > 0.86 for STAR).

836

837

838

839

840    ### _Complexity analysis using Capture Recapture Model_

841    Given a sequencing experiment, the Read Origin Protocol (ROP) attempts to classify every

842    sequenced read in the experiment to an "origin" class. These origins can be considered to

843    be features of interest (e.g. exons, retroviral, immune, or bacterial). Since every read is

844    assigned to only one class, we can consider the reads assigned to a specific class to be a

845    random sample from the population of possibilities within that class. This leads us to

846    consider statistical models for population sampling, which are known as "capture-

847    recapture" models (Bunge & Fitzpatrick, 1993).

848    Using capture-recapture models allows us to make statistical inferences on several

849    quantities of interest. Of primary interest is the total number of possibilities in the feature.

850    We shall refer to this as the feature size but is commonly known in the statistics literature

851    as species richness (Bunge & Fitzpatrick, 1993; Deng, Daley, & Smith, n.d.). We also

852    consider the number of identified possibilities within a feature as a function of the number

853    of reads. We call this the complexity of the feature, in line with the notation of Daley and

854    Smith (T. Daley & Smith, 2013). The rate of change in the complexity curve is proportional

855    to the probability the next read in a previously unobserved class (T. P. Daley, 2014). This

856    quantity is commonly known in statistics literature as the mathematical coverage (Good,

857    1953), but to avoid confusion with sequencing coverage, we call this the discovery

858    probability (Favaro, Lijoi, & Prünster, 2012). One minus the discovery probability will be

859    called the saturation of the feature.

860    ***Statistical Model***

861    Suppose we sequence N reads from an experiment. There are C feature classes,

862    represented in the sequencing library with proportions $\pi_1, \ldots, \pi_C$. Features may overlap,

863    so it is not necessary that the proportions sum to one. The features are all known and

864    defined beforehand. This trait is in contrast to the number of classes within each feature.

865    Within each feature c, there are a fixed but unknown number of classes; Sc represented in

866    the experiment. Within the feature, these are represented with relative proportions

867
$$p_1, \ldots, p_{S_C}, \sum_{i=1}^{S_C} p_i = 1$$

868    If we are interested in the relative proportions within the experiment, we multiply the

869    relative proportion within the feature by the relative abundance of the feature within the

870    experiment.

871    The problem is that we only have information on the classes that were sequenced in the

872    experiment. We observed $D_C \leq S_C$ classes with observed frequencies $x_i$ = # reads from

873    class i with $\sum_{i=1}^{S_C} x_i = N_C$ and $\sum_{c=1}^{C} N_c = N$.

874    The problem of estimating the complexity is to estimate the number of expected distinct

875    classes observed as a function of reads sequenced. We use the non-parametric empirical

876    Bayesian? approach of Daley and Smith (T. Daley & Smith, 2013) to estimate the feature

877    complexity curve. The limit of the feature complexity curve can be regarded as an estimate

878    of the feature size (Colwell & Coddington, 1994).

879    The discovery probability of the observed experiment is the sum of the relative proportions

880    of the unobserved classes,

881
$$\sum_{i=1}^{S_c} p_i \mathbf{1}(x_i = 0).$$

882    The non-parametric empirical Bayes estimator for this quantity is given by the Good Turing

883    formula, $(\sum_{i=1}^{S_c} \frac{\mathbf{1}(x_i=1)}{N_C})$.

### Read Complexity Analysis

885    We first examine the read complexity as determined by the mapped start position of the

886    first end in the read pair. We observe little difference between the two libraries for the

887    single end complexity (Supplemental Methods Figure SM3). We observe only an average

888    of 20% and 29% of the mappable reads at the sequenced read depth. We estimate that all

889    libraries are an average of 58% saturated; that is, we observed 58% of the abundance. This

890    is natural since one would naturally sequence the most abundant reads first.

891

Supplemental Methods Figure SM3. Single end read complexity medians and interquartile

ranges across the two library preparations.

894

*Annotated Feature Complexity Analysis*

The mapped reads can be assigned to features within the genome. These include exons,

introns, coding sequences (CDS), and untranslated regions (UTR). In this section we shall

investigate the complexity of these features, which can be interpreted as estimating the

899   transcriptional diversity within these libraries.

900   As expected, more exons, CDSs, and UTRs were observed per sequenced fragment for the

901   polyA libraries than for the totalRNA libraries. Yet all libraries are very saturated. Most of

902   the abundant classes within these features have already been observed, and the

903   unobserved features are extremely rare. This is in line with the common practice of

904   sequencing a few tens of millions of reads for inferring differential expression.

905

906   To compare the saturation across libraries, we extrapolated the saturation to a common

907   value. The saturation is asymptotically normal (Mao, 2004), and the sequencing depth is

908   sufficiently high that we can use a standard t-test to investigate differences. The polyA

909   libraries are more saturated when all the features for all libraries are extrapolated out to

910   100 million observations (exons: p = 3.764E-16; CDS: p = 1.036E-14; UTR: p = 5.183E-14;

911   more significant differences were observed at lower depths, indicating that the differences

912   are not artifacts of the sampling depth).

913

914   Despite the large saturation for all features across libraries, a multitude of unobserved

915   classes remain (Supplemental Methods Table SM7). This means that most of the

916   unobserved classes are exceedingly rare. For example, we estimate that there are an

917   average of 41,990 unobserved exons in the polyA libraries. There is an average remaining

918   abundance of $1 - 0.9988 = 0.0012$, implying that the average abundance of the

919   unobserved exons is $\frac{0.0012}{41990} = 2.86\,E - 8$. Since, on average, a read has $2 \cdot 0.176 = 0.352$

920   probability of overlapping an exon, the average abundance of the unobserved exons is 1E-

921 8 and the total abundance, 0.00042, gives the marginal probability that the next sequenced

922 read is a new exon. For the totalRNA libraries, the average abundance of the unobserved

923 exons is 3.2E-8. Similarly, we calculated the average abundance of the unobserved CDS for

924 polyA and totalRNA libraries as 1.84E-8 and 7.78E-8, respectively, and for UTRs it was 1.1E-

925 8 and 6.48E-8.

926

| Feature | Mean hits | | Mean observed | | Mean saturation | | Mean estimated total | |
|---------|-----------|----------|---------------|----------|-----------------|----------|----------------------|----------|
| | polyA | totalRNA | polyA | totalRNA | polyA | totalRNA | polyA | totalRNA |
| Exons | 10310521 | | 110553 | | 0.9969 | | 145950 | |
| | 17713362 | 5745436 | 115507 | 107498 | 0.9988 | 0.9956 | 157497 | 138829 |
| CDS | 4791394 | | 105820 | | 0.984 | | 131521 | |
| | 8804113 | 2316884 | 116068 | 99500 | 0.9977 | 0.9756 | 144062 | 123788 |
| UTR | 4359596 | | 33165 | | 0.9948 | | 43136 | |
| | 8035082 | 2093047 | 37448 | 30524 | 0.99913 | 0.99209 | 49849 | 38997 |

927

928 Supplemental Methods Table SM7. Mean number of observations, distinct observed

929 classes, observed saturation, and estimated total number of classes for exons, CDS, and

930    **UTR Features.**

931

932    Finally, we examined differences of diversity between case and controls for a fixed tissue

933    type and library type. The results are quite anticlimactic, as we found little differences

934    between cases and controls for extrapolated saturation and feature diversity. This

935    indicates that there are little differences in transcriptome diversity between the two

936    groups of case and controls. Alternateively, it may indicate that the differences between

937    the two groups are so small that a much larger cohort is required to accurately infer the

938    disparity.

939

940

941    <u>Genomic profiles across library preparation protocols</u>

942    Similar to Li, S. et al we observed that library preparation has a strong effect on the fraction

943    of both mapped and lost human reads mapping to CDS and intronic regions. Genomic

944    profile of mapped and unmapped reads across library preparation protocols is presented

945    in **Supplemental Methods Figure SM4.**

**A. poly(A) enrichment (n=38)**    **B. Ribo-depletion (n=49)**



946

947

948    Supplemental Methods Figure SM4. Genomic profile of mapped and lost human reads

949    across poly(A) enrichment and ribo-depletion libraries.

950    (A) RNA-Seq samples were prepared by poly(A) enrichment protocol (n=38). (B) RNA-Seq

951    samples were prepared by ribo-depletion protocol (n=49). Mapped human reads are

952    identified as RNA-Seq reads that mapped to the human reference genome and

953    transcriptome (ENSEMBL hg19 build, ENSEMBL GRCh37 transcripome) via tophat2. Lost

954    human reads are unmapped RNA-Seq reads that aligned to the human reference genome

955    and transcriptome (ENSEMBL hg19 build, ENSEMBL GRCh37 transcripome) via more

956    sensitive Megablast alignment. Single alignment is reported for each read by Megablast.

957    ROP categorizes the reads into genomic categories based on the compatibility of each read

958    from the pair with the features defined by the Ensembl gene annotations. Percentages are

959    calculated as a fraction of reads from a category from the total number of mapped or lost

960    human reads. Junction read is defined as a read spanning exon-exon boundary; CDS, UTR3,

961    UTR5: reads overlapping CDS, UTR3 or UTR5 region; UTR: reads simultaneously

962    overlapping UTR3 and UTR5 regions; intronic: reads overlapping intronic regions;

963    intergenic: reads mapped within the proximity of 1Kb from the gene boundaries; deep

964    intergenic: reads mapped beyond the proximity of 1Kb from the gene boundaries; MT:

965    mitochondrial reads; multi-mapped: reads mapped to multiple locations of the human

966    genome; fusion: reads from the read pair mapped to different chromosomes.

967

968     Genomic profile across tissue types and library preparation methods in S1. Genomic Profile

969     is obtained based on both mapped and lost human RNA-Seq reads.

**A. Genomic profile obtained based on mapped RNA-Seq reads. Mapped human reads are identified as the RNA-Seq reads mapped to the reference genome and transcriptome (ENSEMBL hg19 build, ENSEMBL GRCh37 transcripome) via tophat2.**

| Tissue | Whole blood | Nasal epithelium | Lung epithelium |
|---|---|---|---|
| N | 19 | 19 | 49 |
| Library preparation method | poly(A) enrichment | poly(A) enrichment | ribo-depletion |
| Splice junction reads, %*, mean (std) | 23.3% (3.3%) | 29.8% (2.2%) | 10.0% (3.3%) |
| CDS reads %, mean (std) | 18.0% (3.1%) | 16.9% ( 1.3%) | 6.9% (2.0%) |
| UTR3 reads %, mean (std) | 15.6% (3.1%) | 22.5% (1.7%) | 11.4% (2.5) |
| UTR5 reads  %, mean (std) | 3.2% (0.7%) | 2.2% (0.3%) | 2.6% (0.7%) |
| UTR** reads %, mean (std) | 4.3% (0.8%) | 5.9% (0.5%) | 1.9% (0.6%) |
| Intronic reads %, mean (std) | 5.6% (1.6%) | 4.4% (0.8%) | 39.4% (6.5%) |
| Proximate inter-genic*** reads %, mean (std) | 1.2% (0.6%) | 1.5% (0.6%) | 3.3% (0.4%) |
| Deep inter-genic reads**** %, mean (std) | 0.3% (0.1%) | 0.3% (0.1%) | 2.8% (0.9%) |
| Mitochondrial (MT) reads %*, mean (std) | 2.3% (1.0%) | 4.3% (1.3%) | 1.5% (1.8%) |
| Milti-mapped reads %, mean (std) | 10.6% (2.4%) | 1.9% (0.2%) | 1.9% (0.5%) |
| Fusion reads %, mean (std) | 0.2% (0.1%) | 0.4 % (0.1%) | 0.7% (0.2%) |

970

**B. Genomic profile obtained based on lost human reads.  Lost human reads are the unmapped RNA-Seq reads that aligned to the human reference genome and transcriptome (ENSEMBL hg19 build, ENSEMBL GRCh37 transcripome) via more sensitive Megablast alignment.**

| Tissue | Whole blood | Nasal epithelium | Lung epithelium |
|---|---|---|---|
| N | 19 | 19 | 49 |
| Library preparation method | poly(A) enrichment | poly(A) enrichment | ribo-depletion |
| Splice junction reads, %*, mean (std) | 1.5% (0.5%) | 0.7% (0.1%) | 0.6% (0.2%) |
| CDS reads %, mean (std) | 1.9% (0.7%) | 0.7% (0.1%) | 0.7% (0.2%) |
| UTR3 reads %, mean (std) | 1.3% (0.3%) | 0.9% (0.1%) | 1.1% (0.2%) |
| UTR5 reads  %, mean (std) | 0.4% (0.1%) | 0.2% (0.03%) | 0.3% (0.1%) |
| UTR** reads %, mean (std) | 0.4% (0.1%) | 0.2% (0.1%) | 0.2% (0.1%) |
| Intronic reads %, mean (std) | 1.0% (0.4%) | 1.3% ( 1.1%) | 5.9% (3.1%) |
| Proximate inter-genic*** reads %, mean (std) | 0.6% (0.4%) | 1.0% (1.1%) | 2.1% (2.5%) |
| Deep inter-genic reads**** %, mean (std) | 0.2% (0.1%) | 0.3% (0.3%) | 0.7% (0.4%) |
| Mitochondrial (MT) reads %*, mean (std) | 0.0% (0.0%) | 0.0% (0.0%) | 0.0% (0.0%) |

Notes :
* percentage from the total number of reads are reported
** reads simultaneously overlapping UTR3 and UTR5 regions
*** mapped with the 1K proximity from gene boundaries
**** mapped further than 1K from the gene boundaries

971

972

973 ## Repeat profile across tissues types and library preparation methods.

974 Repeat profile is based on both mapped and lost repeat reads.

**A. Repeat profile obtained based on mapped RNA-Seq reads. Mapped reads were categorized based on the overlap with the repeat instances prepared from RepeatMasker annotation (Repeatmasker v3.3, Repeat Library 20120124).**

| Tissue | Whole blood | Nasal epithelium | Lung epithelium |
|---|---|---|---|
| N | 19 | 19 | 49 |
| Library preparation method | poly(A) enrichment | poly(A) enrichment | ribo-depletion |
| L1, %*, mean | 0.4% | 0.5% | 5.5% |
| L2, %, mean | 0.2% | 0.2% | 1.0% |
| CR1, %, mean | 0.02% | 0.01% | 0.1% |
| Alu, %, mean | 1.0% | 1.0% | 2.5% |
| MIR, %, mean | 0.1% | 0.1% | 0.6% |
| ERVL-MaLR, %, mean | 0.2% | 0.2% | 1.1% |
| ERV1, %, mean | 0.2% | 0.2% | 0.8% |
| ERVK, %, mean | 0.0% | 0.0% | 0.1% |
| ERVL, %, mean | 0.1% | 0.1% | 0.5% |
| RNA, %, mean | 0.0% | 0.0% | 0.2% |
| hAT-Charlie, %, mean | 0.1% | 0.1% | 0.4% |
| TcMar-Tigger, %, mean | 0.04% | 0.1% | 0.5% |
| Others, %, mean | 0.05% | 0.1% | 0.3% |

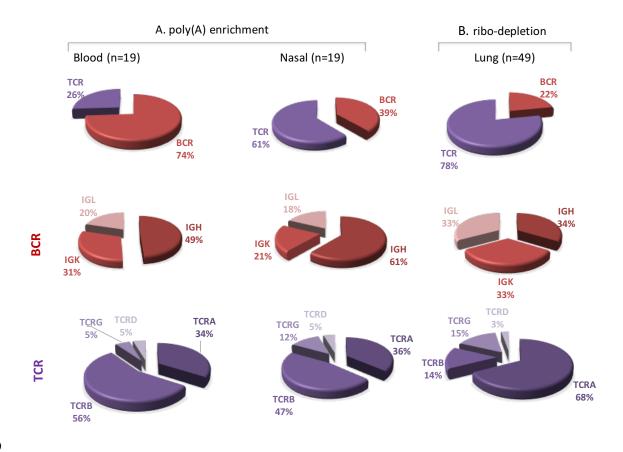\* Percentage from the total number of reads

975

**B. Repeat profile obtained based on lost repeat reads. Lost human reads are the unmapped RNA-Seq reads that aligned to human reference genome and transcriptome (ENSEMBL hg19 build, ENSEMBL GRCh37 transcripome) via more sensitive Megablast alignment.**

| Tissue | Whole blood | Nasal epithelium | Lung epithelium |
|---|---|---|---|
| N | 19 | 19 | 49 |
| Library preparation method | poly(A) enrichment | poly(A) enrichment | ribo-depletion |
| %, mean* | | | |
| hAT, mean | 0.0001% | 0.0004% | 0.0000% |
| TcMar-Mariner, mean | 0.0001% | 0.0005% | 0.0001% |
| TcMar-Tigger, mean | 0.0001% | 0.0015% | 0.0001% |
| L1, mean | 0.0045% | 0.1409% | 0.0048% |
| ERVK, mean | 0.0002% | 0.0026% | 0.0001% |
| ERV, mean | 0.0017% | 0.0082% | 0.0014% |
| ERV1, mean | 0.0025% | 0.0106% | 0.0016% |
| ERVL, mean | 0.0000% | 0.0014% | 0.0000% |
| Satellite, mean | 0.0001% | 0.0006% | 0.0000% |
| Alu, mean | 0.0495% | 0.0896% | 0.0382% |
| Deu, mean | 0.0001% | 0.0024% | 0.0001% |
| Others, mean | 0.0051% | 0.0072% | 0.0025% |

*Percentage from the total number of reads

976

977

## A. poly(A) enrichment

Blood (n=19)      Nasal (n=19)

## B. ribo-depletion

Lung (n=49)

979

980    *Supplemental Methods Figure SM5.*. Percentage of immune reads mapped to B-cell

981    receptor (BCR) and T-cell receptor (TCR) loci.

982    (A) RNA-Seq samples were prepared by poly(A) enrichment protocol (whole blood and

983    nasal epithelium). (B) RNA-Seq samples were prepared by ribo-depletion protocol (lung

984    epithelium). Immune reads that are entirely mapped to BCR and TCR genes are identified

985    by tophat2. Immune reads with extensive somatic hyper mutations (SHM) and reads arising

986    from V(D)J recombination are identified by IgBlast. Blood samples show a larger fraction of

987    reads mapped to BCR locus, while nasal and lung epithelium samples show a larger fraction

988    of reads mapped to TCR locus. BCR are composed of heavy (IGH) and light chains. Among

989    the reads mapped to BCR locus, the number of reads mapped to immunoglobulin heavy

990    locus (IGH), immunoglobulin kappa locus (IGK), and immunoglobulin lambda locus (IGL) is

991    determined. Among the reads mapped to TCR locus, the number of reads mapped to T cell

992    receptor alpha locus (TCRA), T cell receptor beta locus (TCRB), T cell receptor gamma locus

993    (TCRG), and T cell receptor delta locus (TCRD) is determined.

A. poly(A) enrichment             B. ribo-depletion

Blood (n=19)        Nasal (n=19)        Lung (n=49)

Blood: IGHM 43%, IGHA 30%, IGHD 7%, IGHG 20%, IGHE 0%

Nasal: IGHE 1%, IGHM 6%, IGHA 29%, IGHD 2%, IGHG 62%

Lung: IGHM 13%, IGHE 1%, IGHA 15%, IGHD 7%, IGHG 64%

994

995   *Supplemental Methods Figure SM6.*.  **Percentage of immune reads mapped to genes**

996   **encoding the constant region of immunoglobulin heavy locus (IGH).**

997   (A) RNA-Seq samples were prepared by poly(A) enrichment protocol (whole blood and

998   nasal epithelium).  (B)  RNA-Seq samples were prepared by ribo-depletion protocol (lung

999   epithelium).  Immune reads that are entirely mapped to IGHA (Immunoglobulin Heavy

1000   Constant Alpha), IGHD (Immunoglobulin Heavy Constant Delta), IGHG (Immunoglobulin

1001   Heavy Constant Gamma), IGHE (Immunoglobulin Heavy Constant Epsilon), and IGHM

1002   (Immunoglobulin Heavy Constant Mu) are identified by tophat2.

1003

1004

1005

1006

1007 **Number of RNA-Seq reads mapped to BCR and TCR genes (immune reads).**

1008 Reads entirely mapped to BCR and TCR genes are identified by Tophat2. Reads with

1009 extensive somatic hyper mutations (SHM) and reads arising from V(D)J recombination are

1010 identified by IgBLAST.

| Tissue | Whole blood | Nasal epithelium | Lung epithelium |
|---|---|---|---|
| N | 19 | 19 | 49 |
| | | | |
| Library preparation method | poly(A) enrichment | poly(A) enrichment | ribo-depletion |
| Number of immune reads (tophat2), RPM, mean | 4805 | 107 | 16 |
| Number of immune reads (IgBlast), RPM, mean | 270 | 7 | 1 |
| Total number of immune reads , RPM, mean | 5075 | 114 | 17 |

RPM : reads per million

1011

1012

1013

1014    **List of software tools used:**

1015    Tophat2 v.2.0.13 - http://ccb.jhu.edu/software/tophat/index.shtml

1016    STAR  v2.5.2b - https://github.com/alexdobin/STAR

1017    Bowtie v.0.12.9 -  http://bowtie-bio.sourceforge.net/index.shtml

1018    Bowtie2  v.2.2.9 - http://bowtie-bio.sourceforge.net/bowtie2/index.shtml

1019    Samtools v.0.1.18 - http://www.htslib.org/

1020    Bamtools v.2.3.0 -  https://github.com/pezmaster31/bamtools

1021    FASTX-Toolkit v.0.0.13 - http://hannonlab.cshl.edu/fastx_toolkit/

1022    SEQLEAN v(seqclean-x86_64) - http://sourceforge.net/projects/seqclean/files/

1023    BLAST+ v.2.2.30 - ftp://ftp.ncbi.nlm.nih.gov/blast/executables/blast+/LATEST/

1024    IgBlast v.1.4.0- http://www.ncbi.nlm.nih.gov/igblast/

1025    TopHat-Fusion v.2.0.13- http://ccb.jhu.edu/software/tophat/fusion_index.shtml

1026    circExplorer2 v.2.2.4 - http://circexplorer2.readthedocs.io/

1027    MetaPhlAn2 v.2.0 - http://huttenhower.sph.harvard.edu/metaphlan

1028    HTSeq v.0.6.1 - http://www-huber.embl.de/users/anders/HTSeq/

1029    Preseq v 2.0- http://smithlabresearch.org/software/preseq/

1030    Quicksect v.0.0.2  - https://github.com/brentp/quicksect

1031

1032

1033   **Databases**

1034   Ensembl hg19 - http://www.ensembl.org/Homo_sapiens/Info/Index

1035   Human          ribosomal          DNA          complete          repeating          unit          -

1036   http://www.ncbi.nlm.nih.gov/nuccore/U13369

1037   GTF          formatted          file          for          repeat          annotations-

1038   http://labshare.cshl.edu/shares/mhammelllab/www-

1039   data/TEToolkit/TE_GTF/hg19_rmsk_TE.gtf.gz

1040   Repeat elements (*RepBase20.07*) –  http://www.girinst.org/repbase/

1041   V(D)J genes of  *B and T cell receptor* - http://www.imgt.org/vquest/refseqh.html#V-D-J-C-

1042   sets

1043   Database of viral genomes: http://ftp.ncbi.nlm.nih.gov/genomes/Viruses

1044   Database of bacterial genomes:  http://ftp.ncbi.nlm.nih.gov/genomes/Bacteria/

1045   Database of eukaryotic pathogens - http://eupathdb.org/eupathdb/

1046

1047 **References:**

1048 Adiconis, X., Borges-Rivera, D., Satija, R., DeLuca, D. S., Busby, M. A., Berlin, A. M., … others.

1049     (2013). Comparative analysis of RNA sequencing methods for degraded or low-input

1050     samples. *Nature Methods*, *10*(7), 623–629.

1051 Anders, S., Pyl, P. T., & Huber, W. (2014). HTSeq--A Python framework to work with high-

1052     throughput sequencing data. *Bioinformatics*, btu638.

1053 Andrews, S. (2010). FastQC: A quality control tool for high throughput sequence data.

1054     Retrieved from http://www.bioinformatics.babraham.ac.uk/projects/fastqc/

1055 Ardlie, K. G., Deluca, D. S., Segrè, A. V, Sullivan, T. J., Young, T. R., Gelfand, E. T., … others.

1056     (2015). The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue gene

1057     regulation in humans. *Science*, *348*(6235), 648–660.

1058 Beck, J. M., Young, V. B., & Huffnagle, G. B. (2012). The microbiome of the lung.

1059     *Translational Research : The Journal of Laboratory and Clinical Medicine*, *160*(4), 258–

1060     66. https://doi.org/10.1016/j.trsl.2012.02.005

1061 Blachly, J. S., Ruppert, A. S., Zhao, W., Long, S., Flynn, J., Flinn, I., … others. (2015).

1062     Immunoglobulin transcript sequence and somatic hypermutation computation from

1063     unselected RNA-seq reads in chronic lymphocytic leukemia. *Proceedings of the*

1064     *National Academy of Sciences*, *112*(14), 4322–4327.

1065 Bunge, J., & Fitzpatrick, M. (1993). Estimating the number of species: a review. *Journal of*

1066     *the American Statistical Association*, *88*(421), 364–373.

1067 Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., & Madden, T.

1068     L. (2009). BLAST+: architecture and applications. *BMC Bioinformatics*, *10*(1), 421.

1069     Carrara, M., Beccuti, M., Cavallo, F., Donatelli, S., Lazzarato, F., Cordero, F., & Calogero, R.

1070         A. (2013). State of art fusion-finder algorithms are suitable to detect transcription-

1071         induced chimeras in normal tissues? *BMC Bioinformatics*, *14*(7), 1.

1072     Chuang, T.-J., Wu, C.-S., Chen, C.-Y., Hung, L.-Y., Chiang, T.-W., & Yang, M.-Y. (2015).

1073         NCLscan: accurate identification of non-co-linear transcripts (fusion, trans-splicing

1074         and circular RNA) with a good balance between sensitivity and precision. *Nucleic Acids*

1075         *Research*, gkv1013.

1076     Cloonan, N., Forrest, A. R. R., Kolle, G., Gardiner, B. B. A., Faulkner, G. J., Brown, M. K., …

1077         others. (2008). Stem cell transcriptome profiling via massive-scale mRNA sequencing.

1078         *Nature Methods*, *5*(7), 613–619.

1079     Colwell, R. K., & Coddington, J. A. (1994). Estimating terrestrial biodiversity through

1080         extrapolation. *Philosophical Transactions of the Royal Society B: Biological Sciences*,

1081         *345*(1311), 101–118.

1082     Criscione, S. W., Zhang, Y., Thompson, W., Sedivy, J. M., & Neretti, N. (2014).

1083         Transcriptional landscape of repetitive elements in normal and cancer human cells.

1084         *BMC Genomics*, *15*(1), 583. https://doi.org/10.1186/1471-2164-15-583

1085     Daley, T. P. (2014). *Non-parametric Models for Large Capture-recapture Experiments with*

1086         *Applications to DNA Sequencing*. University of Southern California.

1087     Daley, T., & Smith, A. D. (2013). Predicting the molecular complexity of sequencing

1088         libraries. *Nature Methods*, *10*(4), 325–327.

1089     Deng, C., Daley, T., & Smith, A. D. (n.d.). Applications of species accumulation curves in

1090         large-scale biological data analysis. *Journal of Quantitative Biology*.

1091    Engström, P. G., Steijger, T., Sipos, B., Grant, G. R., Kahles, A., Rätsch, G., … others. (2013).

1092    Systematic evaluation of spliced alignment programs for RNA-seq data. *Nature*

1093    *Methods*, *10*(12), 1185–1191.

1094    Favaro, S., Lijoi, A., & Prünster, I. (2012). A new estimator of the discovery probability.

1095    *Biometrics*, *68*(4), 1188–1196.

1096    Good, I. J. (1953). The population frequencies of species and the estimation of population

1097    parameters. *Biometrika*, *40*(3–4), 237–264.

1098    Grabherr, M. G., Haas, B. J., Yassour, M., Levin, J. Z., Thompson, D. A., Amit, I., … Regev, A.

1099    (2011). Full-length transcriptome assembly from RNA-Seq data without a reference

1100    genome. *Nature Biotechnology*, *29*(7), 644–52. https://doi.org/10.1038/nbt.1883

1101    Hach, F., Hormozdiari, F., Alkan, C., Hormozdiari, F., Birol, I., Eichler, E. E., & Sahinalp, S. C.

1102    (2010). mrsFAST: a cache-oblivious algorithm for short-read mapping. *Nature*

1103    *Methods*, *7*(8), 576–577.

1104    Hansen, T. B., Venø, M. T., Damgaard, C. K., & Kjems, J. (2015). Comparison of circular RNA

1105    prediction tools. *Nucleic Acids Research* . https://doi.org/10.1093/nar/gkv1458

1106    Inman, C. F., Murray, T. Z., Bailey, M., & Cose, S. (2012). Most B cells in non-lymphoid

1107    tissues are naïve. *Immunology and Cell Biology*, *90*(2), 235–242.

1108    https://doi.org/10.1038/icb.2011.35

1109    Jeck, W. R., & Sharpless, N. E. (2014). Detecting and characterizing circular RNAs. *Nature*

1110    *Biotechnology*, *32*(5), 453–61. https://doi.org/10.1038/nbt.2890

1111    Jin, Y., Tam, O. H., Paniagua, E., & Hammell, M. (2015). TEtranscripts: a package for

1112        including transposable elements in differential expression analysis of RNA-seq

1113        datasets. *Bioinformatics*, btv422.

1114    Kim, D., Pertea, G., Trapnell, C., Pimentel, H., Kelley, R., & Salzberg, S. L. (2013). TopHat2:

1115        accurate alignment of transcriptomes in the presence of insertions, deletions and

1116        gene fusions. *Genome Biology*, *14*(4), R36. https://doi.org/10.1186/gb-2013-14-4-r36

1117    Kostic, A. D., Ojesina, A. I., Pedamallu, C. S., Jung, J., Verhaak, R. G. W., Getz, G., &

1118        Meyerson, M. (2011). PathSeq: software to identify or discover microbes by deep

1119        sequencing of human tissue. *Nature Biotechnology*, *29*(5), 393–396.

1120    Li, S., Tighe, S. W., Nicolet, C. M., Grove, D., Levy, S., Farmerie, W., … others. (2014). Multi-

1121        platform assessment of transcriptome profiling using RNA-seq in the ABRF next-

1122        generation sequencing study, *32*(9), 915–925. https://doi.org/10.1038/nbt.2972

1123    Mao, C. X. (2004). Predicting the conditional probability of discovering a new class. *Journal*

1124        *of the American Statistical Association*, *99*(468).

1125    Melé, M., Ferreira, P. G., Reverter, F., DeLuca, D. S., Monlong, J., Sammeth, M., … others.

1126        (2015). The human transcriptome across tissues and individuals. *Science*, *348*(6235),

1127        660–665.

1128    Mihaela Pertea, J. T. M. S. L. S. (2015). StringTie enables improved reconstruction of a

1129        transcriptome from RNA-seq reads. *Nature Biotechnology*, *33*, 290–295.

1130        https://doi.org/10.1038/nbt.3122

1131    Nicolae, M., Mangul, S., Mandoiu, I. I., & Zelikovsky, A. (2011). Estimation of alternative

1132        splicing isoform frequencies from RNA-Seq data. *Algorithms for Molecular Biology*,

1133        *6*(1), 9.

1134    Ozsolak, F., & Milos, P. M. (2011). RNA sequencing: advances, challenges and

1135        opportunities. *Nature Reviews. Genetics*, *12*(2), 87–98.

1136        https://doi.org/10.1038/nrg2934

1137    Perucheon, S., Chaoul, N., Burelout, C., Delache, B., Brochard, P., Laurent, P., … Richard, Y.

1138        (2009). Tissue-specific B-cell dysfunction and generalized memory B-cell loss during

1139        acute SIV infection. *PLoS ONE*, *4*(6), e5966.

1140        https://doi.org/10.1371/journal.pone.0005966

1141    Porath, H. T., Carmi, S., & Levanon, E. Y. (2014). A genome-wide map of hyper-edited RNA

1142        reveals numerous new sites. *Nature Communications*, *5*, 4726.

1143        https://doi.org/10.1038/ncomms5726

1144    Salter, S. J., Cox, M. J., Turek, E. M., Calus, S. T., Cookson, W. O., Moffatt, M. F., … Walker,

1145        A. W. (2014). Reagent and laboratory contamination can critically impact sequence-

1146        based microbiome analyses. *BMC Biology*, *12*(1), 87.

1147    Seqc/Maqc-Iii Consortium. (2014). A comprehensive assessment of RNA-seq accuracy,

1148        reproducibility and information content by the Sequencing Quality Control

1149        Consortium. *Nature Biotechnology*, *32*(9), 903–914.

1150        https://doi.org/10.1038/nbt.2957

1151 Siragusa, E., Weese, D., & Reinert, K. (2013). Fast and accurate read mapping with

1152    approximate seeds and multiple backtracking. *Nucleic Acids Research*, *41*(7), e78--

1153    e78.

1154 Spreafico, R., Rossetti, M., van Loosdregt, J., Wallace, C. A., Massa, M., Magni-Manzoni, S.,

1155    ... Albani, S. (2016). A circulating reservoir of pathogenic-like CD4+ T cells shares a

1156    genetic and phenotypic signature with the inflamed synovial micro-environment.

1157    *Annals of the Rheumatic Diseases*, *75*(2), 459–465.

1158 Strauli, N., & Hernandez, R. (2015). Statistical Inference of a Convergent Antibody

1159    Repertoire Response to Influenza Vaccine. *bioRxiv*, 25098.

1160 Sultan, M., Schulz, M. H., Richard, H., Magen, A., Klingenhoff, A., Scherf, M., ... others.

1161    (2008). A global view of gene activity and alternative splicing by deep sequencing of

1162    the human transcriptome. *Science*, *321*(5891), 956–960.

1163 Tang, F., Barbacioru, C., Wang, Y., Nordman, E., Lee, C., Xu, N., ... others. (2009). mRNA-

1164    Seq whole-transcriptome analysis of a single cell. *Nature Methods*, *6*(5), 377–382.

1165 Tarailo-Graovac, M., & Chen, N. (2009). Using RepeatMasker to identify repetitive

1166    elements in genomic sequences. *Current Protocols in Bioinformatics*, 4–10.

1167 Trapnell, C., Williams, B. a, Pertea, G., Mortazavi, A., Kwan, G., van Baren, M. J., ... Pachter,

1168    L. (2010). Transcript assembly and quantification by RNA-Seq reveals unannotated

1169    transcripts and isoform switching during cell differentiation. *Nature Biotechnology*,

1170    *28*(5), 511–515. https://doi.org/10.1038/nbt.1621

1171 Truong, D. T., Franzosa, E. A., Tickle, T. L., Scholz, M., Weingart, G., Pasolli, E., … Segata, N.

1172 (2015). MetaPhlAn2 for enhanced metagenomic taxonomic profiling. *Nature*

1173 *Methods*, *12*(10), 902–903.

1174 Wang, X.-S., Prensner, J. R., Chen, G., Cao, Q., Han, B., Dhanasekaran, S. M., … Chinnaiyan,

1175 A. M. (2009). An integrative approach to reveal driver gene fusions from paired-end

1176 sequencing data in cancer. *Nature Biotechnology*, *27*(11), 1005–11.

1177 https://doi.org/10.1038/nbt.1584

1178 Wang, Z., Gerstein, M., & Snyder, M. (2009). RNA-Seq: a revolutionary tool for

1179 transcriptomics. *Nature Reviews Genetics*, *10*(1), 57–63.

1180 Wu, C.-S., Yu, C.-Y., Chuang, C.-Y., Hsiao, M., Kao, C.-F., Kuo, H.-C., & Chuang, T.-J. (2014).

1181 Integrative transcriptome sequencing identifies trans-splicing events with important

1182 roles in human embryonic stem cell pluripotency. *Genome Research*, *24*(1), 25–36.

1183 Yan, M., Pamp, S. J., Fukuyama, J., Hwang, P. H., Cho, D. Y., Holmes, S., & Relman, D. a.

1184 (2013). Nasal microenvironments and interspecific interactions influence nasal

1185 microbiota complexity and S. aureus carriage. *Cell Host and Microbe*, *14*(6), 631–640.

1186 https://doi.org/10.1016/j.chom.2013.11.005

1187 Ye, J., Ma, N., Madden, T. L., & Ostell, J. M. (2013). IgBLAST: an immunoglobulin variable

1188 domain sequence analysis tool. *Nucleic Acids Research*, gkt382.

1189 Zhang, X.-O., Dong, R., Zhang, Y., Zhang, J.-L., Luo, Z., Zhang, J., … Yang, L. (2016). Diverse

1190 alternative back-splicing and alternative splicing landscape of circular RNAs. *Genome*

1191 *Research*. https://doi.org/10.1101/gr.202895.115

1192