

Supporting Information

Ergun et al. 10.1073/pnas.1311839110

SI Materials and Methods

ImmGen Consortium. David A. Blair^a, Michael L. Dustin^a, Susan A. Shinton^b, Richard R. Hardy^b, Tal Shay^c, Aviv Regev^c, Nadia Cohen^d, Patrick Brennan^d, Michael Brenner^d, Francis Kim^e, Tata Nageswara Rao^e, Amy Wagers^e, Tracy Heng^f, Jeffrey Ericson^f, Katherine Rothamel^f, Adriana Ortiz-Lopez^f, Diane Mathis^f, Christophe Benoist^f, Taras Kreslavsky^g, Anne Fletcher^g, Kutlu Elpek^g, Angélique Bellemare-Pelletier^g, Deepali Malhotra^g, Shannon Turley^g, Jennifer Miller^h, Brian Brown^h, Miriam Merad^h, Emmanuel L. Gautier^{h,i}, Claudia Jakubczik^h, Gwendalyn J. Randolph^{h,i}, Paul Monach^j, Adam J. Best^k, Jamie Knell^k, Ananda Goldrath^k, Vladimir Jovic^l, Daphne Koller^l, David Laidlaw^m, Jim Collinsⁿ, Roi Gazit^o, Derrick J. Rossi^o, Nidhi Malhotra^p, Katelyn Sylvia^p, Joonsoo Kang^p, Natalie A. Bezman^q, Joseph C. Sun^q, Gundula Min-Oo^q, Charlie C. Kim^q, and Lewis L. Lanier^q.

^aSkirball Institute of Biomolecular Medicine, New York University School of Medicine, New York, NY 10016; ^bFox Chase Cancer Center, Philadelphia, PA 19111; ^cBroad Institute and Department of Biology, Massachusetts Institute of Technology, Cambridge, MA 02142; ^dDivision of Rheumatology, Immunology and Allergy, Brigham and Women's Hospital, Boston, MA 02115; ^eSection on Islet Cell and Regenerative Biology, Joslin Diabetes Center, Boston, MA 02215; ^fDivision of Immunology, Department of Microbiology and Immunobiology, Harvard Medical School, Boston, MA 02115; ^gDepartment of Cancer Immunology and AIDS, Dana-Farber Cancer Institute and Harvard Medical School, Boston, MA 02115; ^hIcahn Medical Institute, Mount Sinai Hospital, New York, NY 10029; ⁱDepartment of Pathology and Immunology, Washington University in St. Louis, St. Louis, MO 63110; ^jDepartment of Medicine, Boston University, Boston, MA 02118; ^kDivision of Biological Sciences, University of California, San Diego, La Jolla, CA 92093; ^lComputer Science Department, Stanford University, Stanford, CA 94305; ^mComputer Science Department, Brown University, Providence, RI 02912; ⁿDepartment of Biomedical Engineering, Howard Hughes Medical Institute, Boston University, Boston, MA 02115; ^oProgram in Molecular Medicine, Children's Hospital, Boston, MA 02115; ^pDepartment of Pathology, University of Massachusetts Medical School, Worcester, MA 01655; and ^qDepartment of Microbiology and Immunology, University of California, San Francisco, CA 94143.

Alternative Splicing Discovery from RNA-Sequencing Data. RNA from CD4+ T and CD19+ B cells was used for the preparation of nonstrand-specific paired-end cDNA libraries in a format compatible with Illumina sequencing technology, which was modified from the work by Parkhomchuk et al. (1). In brief, an mRNA fraction was enriched from total cellular RNA using oligo(dT) selection, treated with DNase, fragmented, and converted to cDNA with random hexamers. cDNA libraries were subjected to strand-independent Illumina library preparation and sequenced using an Illumina sequencing platform (2 × 76 bp).

We mapped the RNA-sequencing (RNA-seq) data to the mm9 (NCBIM37) genome assembly using TopHat (2), version 1.1.4, which required that reads map uniquely. Using the UCSC known gene annotations, we computed gene expression estimates using Cufflinks (3).

TopHat outputs candidate splice junction predictions and a read count for the number of sequencing reads supporting each prediction. We took these reads to represent RNA splice junctions within our samples. Separately we computed alignment and coverage files to visualize the distribution of sequencing reads in

a genome browser. Using the summary read counts at each read position within the genome based on the coverages, we estimated the exon-level expression from the number of reads falling within exon boundaries by computing the mean of the piled-up reads with no subsequent normalization. (In total, we obtained 161,948 and 167,067 splice junctions, where 95,699 and 96,372 of these junctions were observed more than 10 times in CD4+ T and CD19+ B cells, respectively.) We further eliminated junctions ending in multiple genes or located entirely in intronic regions. This process resulted in a total of 92,562 and 91,882 junctions in CD4+ T and CD19+ B cells, respectively.

On the basis of the UCSC mm9 mouse exon annotation boundaries, we further categorized these splice junctions as canonical or alternative. If a junction end perfectly matched to two neighboring exon boundaries, we characterized it as canonical. Any junction with skipped exons or acceptor or donor sites previously not in existing databases was characterized as alternative and binned into one of the categories shown in Fig. 1.

We defined the skipping ratio of an exon as the total read count of junctions skipping that exon divided by the total number of reads skipping it and junctions including it. A junction was considered to include an exon if it fell within ±4 bp around exon-annotated boundaries (Fig. S24).

Alternative Splicing Discovery from Microarray Data. Transcriptional profiling within the Immunological Genome Project (ImmGen) cell populations has been previously described (4) (www.immgen.org). All ImmGen microarray data were generated using Affymetrix MoGene1.0 ST chips, which contain probes across the gene body. With two to three probes falling within an exon on average, we could characterize exon-level expression. To enable cross-comparison between microarray and RNA-seq data, we selected a set of probes with expression that was well-correlated (Pearson coefficient ≥ 0.7) between the microarray and RNA-seq datasets. Our final list was composed of 171,000 probes located in 75,673 exons and 7,261 genes. We estimated exon- and gene-level expression of Affymetrix probes after probe-level robust multiarray average (RMA) normalization (5) and summarized at the probe, exon, or gene level. For analysis of differential use, we selected 172 populations that excluded nonhematopoietic cell types and samples with lesser dynamic range.

Exon-Centric Discovery by Linear Regression. For each of 171,000 probes, we performed a linear regression between the probe's expression and the integrated expression of the corresponding gene across 802 ImmGen samples (three replicates per cell population on average) and computed a residual of each sample to a regression line estimated by a first-degree polynomial fit. We only retained those exons that had a consistent slope for each probe-gene regression for each of the probes and averaged the residual values for each exon and cell type. For each cell type, we selected the top 50 exons with the highest residual value, with exon/gene expression ratio ≤ 1 (Fig. S2C) and gene-level expression ≥ 100. We further required that at least five cell types have high residual value and that the difference in exon/gene expression ratio was at least 0.4 (75th percentile) among the cell types (Fig. S2B).

Population-Centric Analysis of Exon/Gene Ratios. For each cell type, we computed exon-gene ratios and selected exons that had a difference of exon- and gene-level expression > 0.4 (log10 scale) and a gene-level expression ≥ 100. To eliminate exons that consistently showed lower exon expression compared with gene-level

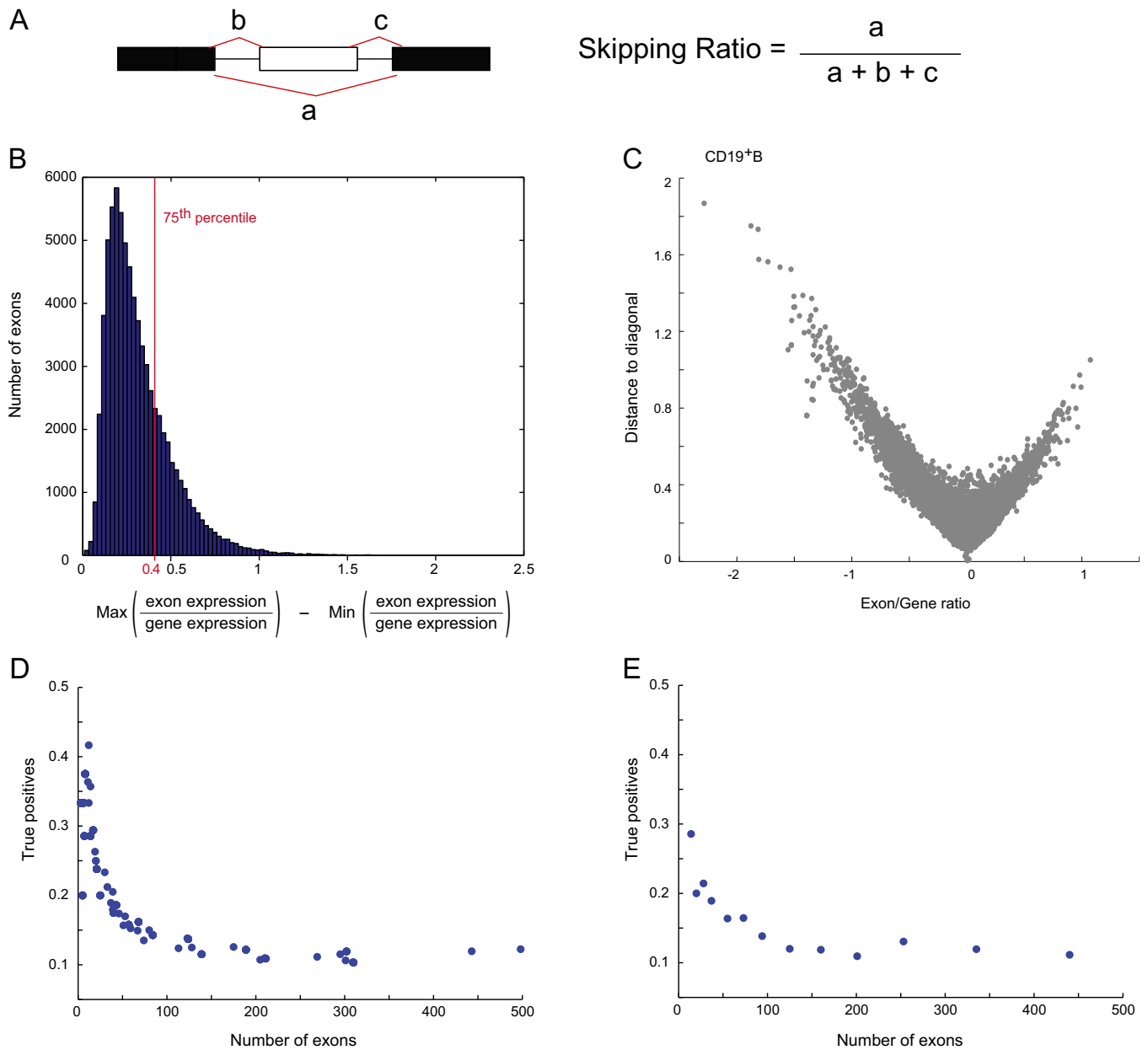


Fig. S2. (A) Skipping ratio for an exon computed as the ratio of junctions skipping that exon normalized by the total number of reads skipping it and junctions including it. (B) Difference between maximum and minimum exon inclusion/exclusion ratio for a given exon; 75th percentile (0.4) was chosen as a cutoff to ensure variability between inclusion/exclusion in flagged exons for the population-centric method. (C) Combined residual values for each population vs. exon/gene ratio from the exon-centric method. (D and E) Comparison between flagged exons in Affymetrix and skipped exons obtained from exon-exon junctions from RNA-seq data in CD19⁺ B cells. Percent of true positives flagged by the exon-centric method as a function of exon number and percent of true positives flagged by the population-centric method as a function of exon number, respectively, are shown.

Other Supporting Information Files

- [Dataset S1 \(XLSX\)](#)
- [Dataset S2 \(XLSX\)](#)
- [Dataset S3 \(XLSX\)](#)
- [Dataset S4 \(XLSX\)](#)
- [Dataset S5 \(XLSX\)](#)
- [Dataset S6 \(XLSX\)](#)
- [Dataset S7 \(XLSX\)](#)