

# Landscape of X chromosome inactivation across human tissues

Taru Tukiainen<sup>1,2</sup>, Alexandra-Chloé Villani<sup>2,3</sup>, Angela Yen<sup>2,4</sup>, Manuel A. Rivas<sup>1,2,5</sup>, Jamie L. Marshall<sup>1,2</sup>, Rahul Satija<sup>2,6,7</sup>, Matt Aguirre<sup>1,2</sup>, Laura Gauthier<sup>1,2</sup>, Mark Fleharty<sup>2</sup>, Andrew Kirby<sup>1,2</sup>, Beryl B. Cummings<sup>1,2</sup>, Stephane E. Castel<sup>6,8</sup>, Konrad J. Karczewski<sup>1,2</sup>, François Aguet<sup>2</sup>, Andrea Byrnes<sup>1,2</sup>, GTEx Consortium†, Tuuli Lappalainen<sup>6,8</sup>, Aviv Regev<sup>2,9</sup>, Kristin G. Ardlie<sup>2</sup>, Nir Hacohen<sup>2,3</sup> & Daniel G. MacArthur<sup>1,2</sup>

**X chromosome inactivation (XCI) silences transcription from one of the two X chromosomes in female mammalian cells to balance expression dosage between XX females and XY males. XCI is, however, incomplete in humans: up to one-third of X-chromosomal genes are expressed from both the active and inactive X chromosomes (Xa and Xi, respectively) in female cells, with the degree of ‘escape’ from inactivation varying between genes and individuals<sup>1,2</sup>. The extent to which XCI is shared between cells and tissues remains poorly characterized<sup>3,4</sup>, as does the degree to which incomplete XCI manifests as detectable sex differences in gene expression<sup>5</sup> and phenotypic traits<sup>6</sup>. Here we describe a systematic survey of XCI, integrating over 5,500 transcriptomes from 449 individuals spanning 29 tissues from GTEx (v6p release) and 940 single-cell transcriptomes, combined with genomic sequence data. We show that XCI at 683 X-chromosomal genes is generally uniform across human tissues, but identify examples of heterogeneity between tissues, individuals and cells. We show that incomplete XCI affects at least 23% of X-chromosomal genes, identify seven genes that escape XCI with support from multiple lines of evidence and demonstrate that escape from XCI results in sex biases in gene expression, establishing incomplete XCI as a mechanism that is likely to introduce phenotypic diversity<sup>6,7</sup>. Overall, this updated catalogue of XCI across human tissues helps to increase our understanding of the extent and impact of the incompleteness in the maintenance of XCI.**

Mammalian female tissues consist of two mixed cell populations, each with either the maternally or paternally inherited X chromosome marked for inactivation. To overcome this heterogeneity, assessments of human XCI have often been confined to the use of artificial cell systems<sup>1</sup> or to samples that have skewed XCI<sup>1,2</sup>, that is, preferential inactivation of one of the two X chromosomes; this is common in clonal cell lines but rare in karyotypically normal, primary human tissues<sup>8</sup> (Extended Data Fig. 1 and Supplementary Note). Others have used bias in DNA methylation<sup>3,4,9</sup> or in gene expression<sup>5,10</sup> between males and females as a proxy for XCI status. Surveys of XCI are powerful in engineered model organisms, for example, mouse models with completely skewed XCI<sup>11</sup>, but the degree to which these discoveries are generalizable to human XCI remains unclear given marked differences in XCI initiation and the extent of escape across species<sup>7</sup>. Here we describe a systematic survey of the landscape of human XCI using three complementary RNA sequencing (RNA-seq)-based approaches (Fig. 1) that together enable the assessment of XCI from individual cells to population level across a diverse range of human tissues.

Given the limited accessibility of most human tissues, particularly in large sample sizes, no global investigation into the impact of incomplete XCI on X-chromosomal expression has been conducted in datasets spanning multiple tissue types. We used the Genotype-Tissue Expression (GTEx) project<sup>12,13</sup> dataset (v6p release), which includes high-coverage RNA-seq data from diverse human tissues, to investigate male–female differences in the expression of 681 X-chromosomal genes that encode proteins or long non-coding RNA in 29 adult tissues (Extended Data Table 1), hypothesizing that escape from XCI should typically result in higher female expression of these genes. Previous work<sup>5,10,14</sup> has indicated that some of the genes that escape XCI (hereafter referred to as escape genes) show female bias in expression, but our analysis benefits from a larger set of profiled tissues and individuals, as well as the high sensitivity of RNA-seq.

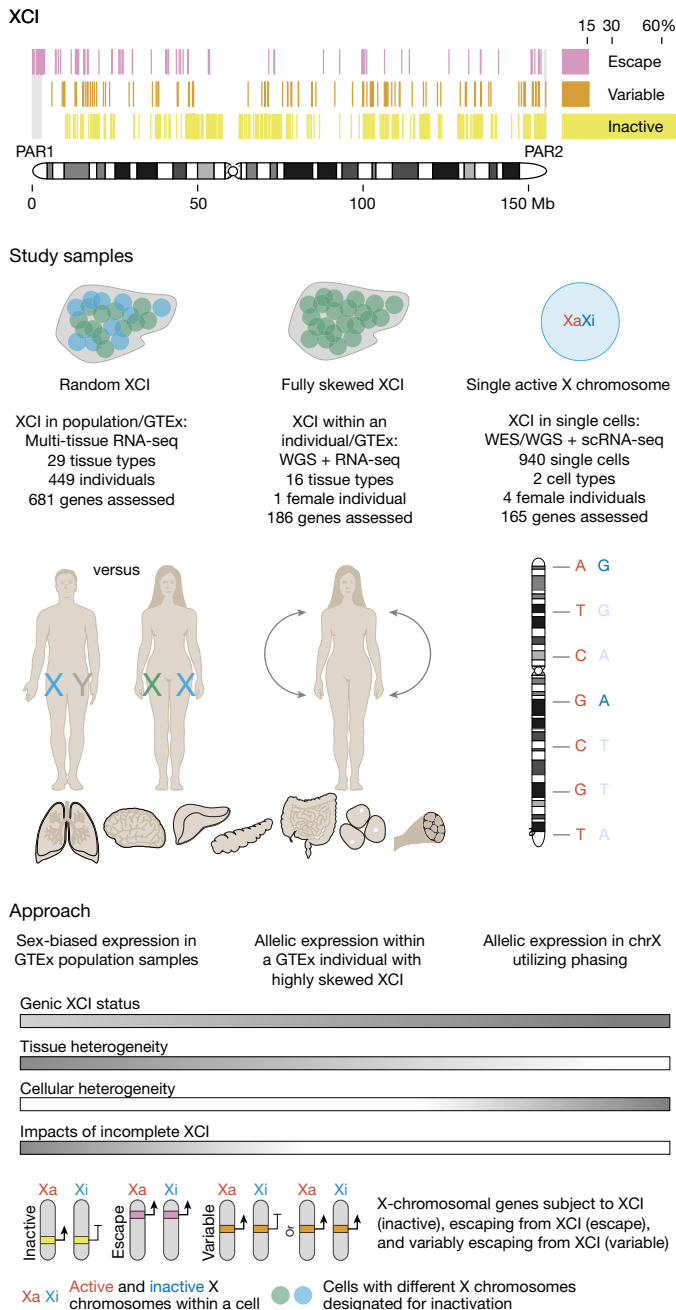
To confirm that male–female expression differences reflect incomplete XCI, we assessed the enrichment of sex-biased expression in known XCI categories using 561 genes with previously assigned XCI status, defined as escape ( $n = 82$ ), variable escape ( $n = 89$ ) or inactive ( $n = 390$ ) (Fig. 1 and Supplementary Table 1). Sex-biased expression is enriched in escape genes compared to both inactive genes (two-sided paired Wilcoxon rank-sum test,  $P = 3.73 \times 10^{-9}$ ) and variable escape genes ( $P = 3.73 \times 10^{-9}$ ) (Fig. 2b and Extended Data Fig. 2), with 74% of escape genes showing significant (false-discovery rate (FDR)  $q < 0.01$ ) male–female differences in at least one tissue (Fig. 2a, Extended Data Figs 3, 4 and Supplementary Table 2). In line with two active X-chromosomal copies in females, escape genes in the non-pseudoautosomal, that is, the X-specific, region (nonPAR) predominantly show female-biased expression across tissues (52 out of 67 assessed genes, binomial  $P = 6.46 \times 10^{-6}$ ). However, genes in the pseudoautosomal region PAR1, are expressed more highly in males (14 out of 15 genes, binomial  $P = 9.77 \times 10^{-6}$ ) (Fig. 2a), suggesting that combined Xa and Xi expression in females fails to reach the expression arising from X and Y chromosomes in males (discussed below).

Sex bias of escape genes is often shared across tissues; these genes show a higher number of tissues with sex-biased expression than genes in other XCI categories (Fig. 2a and Extended Data Fig. 2c), a result that is not driven by differences in the breadth of expression of escape and inactive genes (Extended Data Fig. 2e). Also, the direction of sex bias across tissues is consistent (Fig. 2a, c and Extended Data Fig. 2b). Together, these observations indicate that there is global and tight control of XCI, that potentially arises from early lockdown of the epigenetic marks regulating XCI. Previous reports have identified several epigenetic signatures associated with XCI escape in humans and mice<sup>15</sup>; in agreement with these discoveries we show that escape

<sup>1</sup>Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, Massachusetts 02114, USA. <sup>2</sup>Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02142, USA.

<sup>3</sup>Center for Immunology and Inflammatory Diseases, Massachusetts General Hospital, Charlestown, Massachusetts 02129, USA. <sup>4</sup>Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA. <sup>5</sup>Department of Biomedical Data Science, Stanford University, Stanford, California 94305, USA. <sup>6</sup>New York Genome Center, New York, New York 10013, USA. <sup>7</sup>Center for Genomics and Systems Biology, Department of Biology, New York University, New York, New York 10003, USA. <sup>8</sup>Department of Systems Biology, Columbia University, New York, New York 10032, USA. <sup>9</sup>Department of Biology, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA.

<sup>†</sup>Lists of participants and their affiliations are included in the online version of the paper.



**Figure 1 | Schematic overview of the study.** Previous expression-based surveys of XCI<sup>1,2</sup> have established the incomplete and variable nature of XCI, but these studies have been limited in the tissue types and samples assessed. To investigate the landscape of XCI across human tissues, we combined three approaches: (1) sex biases in expression using population-level GTEx data across 29 tissue types; (2) allelic expression in 16 tissue samples from a female GTEx donor with fully skewed XCI, and (3) validation using scRNA-seq by combining allelic expression and genotype phasing. WGS, whole-genome sequencing; WES, whole-exome sequencing.

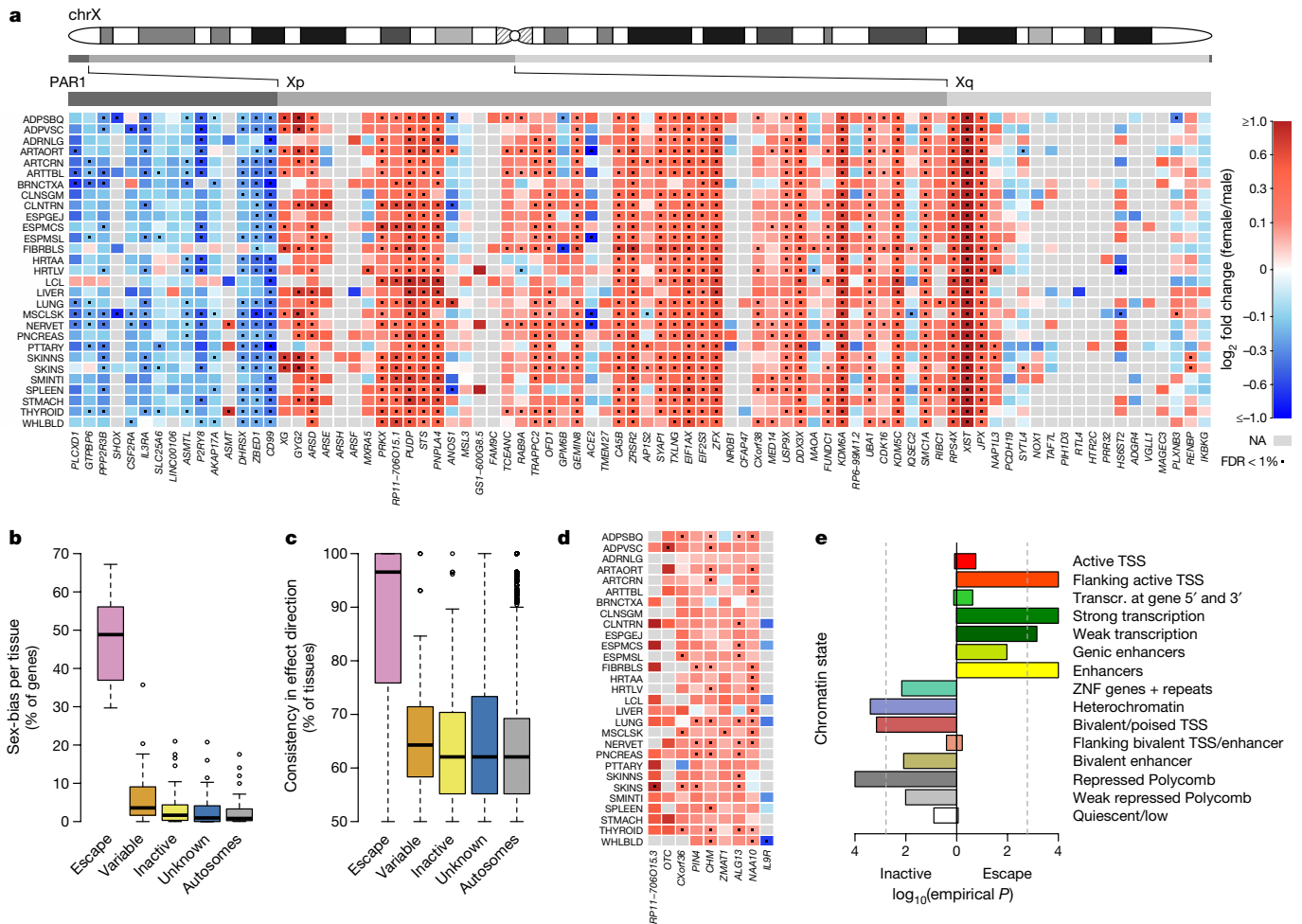
genes are enriched in chromatin states that are related to active transcription (Fig. 2e). Although sex bias on the X chromosome is broadly specific to escape genes, some genes show unexpected patterns. Eight genes with some previous evidence for inactivation show >90% concordance in effect direction and significant sex bias (Fig. 2d and Supplementary Table 3), suggesting that variable escape can also have considerable population-level effects. For example, *CHM* demonstrates such

concordance in sex bias and escape at this gene is confirmed when using single-cell RNA-seq (scRNA-seq; see below). One gene (*RP11-706O15.3*) without an assigned XCI status shows a similar sex bias pattern to escape genes. *RP11-706O15.3* resides between escape and variable escape genes *PRKX* and *NLGN4X* (Fig. 2d), consistent with known clustering of escape genes<sup>1,2</sup>. Some escape genes show more heterogeneous sex bias, for example, *ACE2* (Fig. 2a and Supplementary Discussion). Many such genes lie in the evolutionarily older region of the chromosome<sup>16</sup>, in Xq, where escape genes also show higher tissue-specificity and lower expression levels (Extended Data Fig. 5), characteristics that have been linked with higher protein evolutionary rates<sup>17,18</sup>.

Although sex bias serves as a proxy for XCI status, it provides only an indirect measurement of XCI. We identified a GTEx female donor with an unusual degree of skewing of XCI (Fig. 3a), in whom the same copy of chromosome X was silenced in approximately 100% of cells across all tissues, but with no X-chromosomal abnormality detected by whole-genome sequencing (Extended Data Fig. 6 and Supplementary Note), providing an opportunity to analyse allele-specific expression (ASE) across 16 tissues to investigate XCI. This approach is analogous to previous surveys in mouse<sup>11</sup> or in human cell lines with skewed XCI<sup>2</sup>, but extends the assessment to a larger number of tissues and avoids biases arising from genetic heterogeneity between tissue samples.

Analysis of the X-chromosomal allelic counts (Supplementary Tables 4–6) from this GTEx donor highlights the incompleteness and consistency of XCI across tissues (Fig. 3b). Approximately 23% of the 186 X-chromosomal genes that were assessed show expression from both alleles, indicative of incomplete XCI, matching previous estimates of the extent of escape<sup>1,2</sup>. For 43% of the genes that were expressed from both alleles in this sample, Xi expression is of a similar magnitude between tissues, therefore supporting the observation of a general global and tight control of XCI. However, suggesting some tissue dependence in XCI, the rest of the genes that were expressed from both alleles show variability in Xi expression, including a subset of genes (5.8% of all genes) that appear biallelic in only one of the multiple tissues assayed. While tissue-specific escape is common in mouse<sup>11</sup>, limited evidence exists for such a pattern in human tissues other than for neurons<sup>3,4,9</sup>. In our data, one of the genes with the strongest evidence for tissue-specific escape is *KALI* (Fig. 3f and Supplementary Table 6), the causal gene for X-linked Kallmann syndrome. We show that *KALI* shows biallelic expression exclusively in the lung (Fig. 3f), in line with the strong female bias detected specifically in lung expression in the analysis described above and in Fig. 2a, suggesting that tissue differences in escape can directly translate into tissue-specific sex biases in gene expression. The predictions of XCI status in this sample not only align with previous assignments (Fig. 3c–f and Supplementary Table 7, for example, *TSR2*, *XIST* and *ZBED1*) but also suggest five new incompletely inactivated genes (Fig. 3g–k and Supplementary Table 5), three of which act in a tissue-specific manner. For instance, *CLIC2*, which in previous studies was shown to either be subject to<sup>2</sup> or variably escape from<sup>1</sup> XCI, shows considerable Xi expression only in skin tissue. Such specific patterns illustrate the need to assay multiple tissue types to fully uncover the diversity in XCI.

The emergence of scRNA-seq methods<sup>19</sup> presents an opportunity to directly assess XCI without the complication of cellular heterogeneity in bulk tissue samples (Fig. 1), as demonstrated recently in mouse studies<sup>20–23</sup> and in human fibroblasts<sup>24</sup> and preimplantation development<sup>25</sup>. To directly profile XCI in human samples, we examined scRNA-seq data in combination with deep genotype sequences from 940 immune-related cells from four females: 198 cells from lymphoblastoid cell lines (LCLs) sampled from three females of African (Yoruba) ancestry, and 742 blood dendritic cells from a female of Asian ancestry<sup>26</sup> (Fig. 1 and Extended Data Table 2). We used ASE to distinguish the expression coming from each of the two X-chromosomal haplotypes in a given cell (Supplementary Table 4). Because the inference of allele-specific phenomena in single cells is complicated by widespread



**Figure 2 | Assessment of tissue-sharing and population-level impacts of incomplete XCI in GTEx data.** **a**, Male–female expression differences in reported XCI-escaping genes ( $n = 82$ ) across 29 GTEx tissues. Definitions for the abbreviations can be found in Extended Data Table 1. **b**, Proportion of significantly biased ( $FDR < 1\%$ ) genes in each tissue by reported XCI status. **c**, Proportion of tissues where the bias direction is shared with the

reported XCI status. Genes expressed in at least five tissues are included. **d**, Sex bias pattern of nine genes not classified as full escape genes that follow a similar profile to established escape genes. **e**, Chromatin state enrichment between escape and inactive genes in the Roadmap Epigenomics<sup>31</sup> female samples.

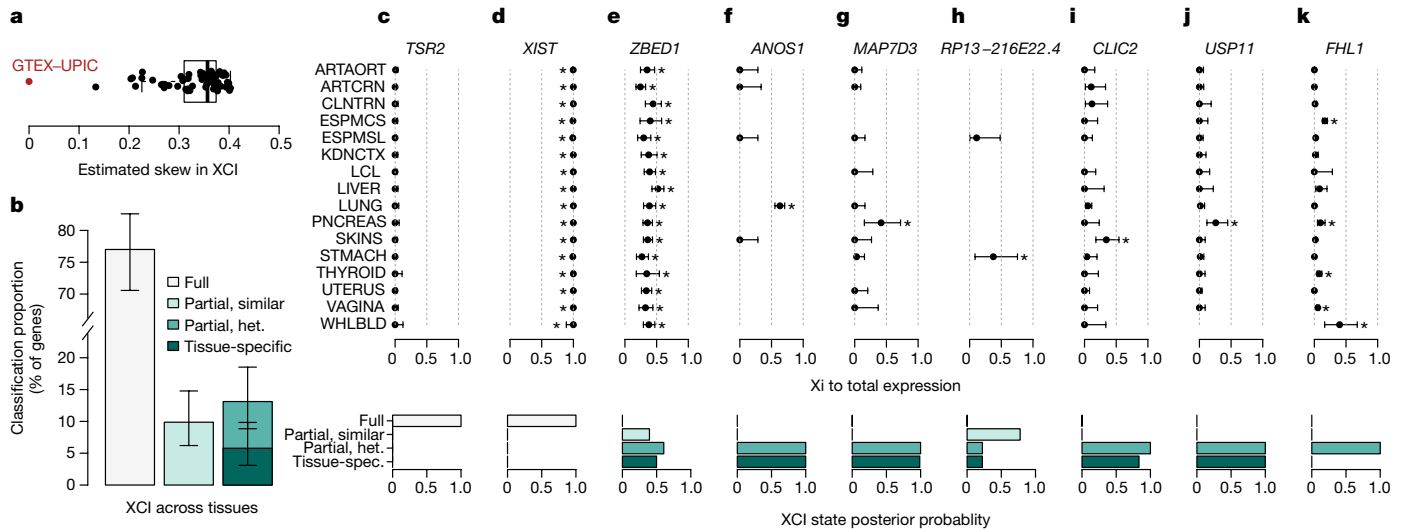
monoallelic expression<sup>21,27–29</sup>, besides searching for X-chromosomal sites with biallelic expression (Extended Data Fig. 7), we leveraged genotype phase information to detect sites for which the expressed allele was discordant with the active X chromosome in that cell.

Only 129 (78%) out of the 165 assayed genes (41–98 per sample) were fully inactivated in these data whereas the rest showed incomplete XCI in one or more samples (Fig. 4a, b and Supplementary Tables 8, 9); this is mostly consistent with previous assignments of XCI status to these genes (Fig. 4a and Supplementary Table 10). For instance, single-cell data reveal consistent expression from both X-chromosomal alleles for eleven genes in PAR1, in line with their known escape from XCI (for example, *ZBED1*, Fig. 4c), and replicate the known expression of *XIST* exclusively from Xi (Fig. 4d).

We next assessed whether our approach could extend the spectrum of escape from XCI. For seven genes that have previously been reported as inactivated, the data from single cells pointed to incomplete XCI (Fig. 4e–k and Supplementary Table 11), including *FHL1*, which was also highlighted as a candidate escape gene in the GTEx ASE analysis (Fig. 4e), and *ATP6AP2*, which displays predominantly female-biased expression across GTEx tissues (Fig. 4h). Both of these genes demonstrate significant Xi expression in only a subset of the scRNA-seq samples, a pattern that is consistent with variable escape<sup>1,2</sup>. Between-individual variability exists not only in the presence but also in the degree of expression from Xi (for example, *MSL3*, Fig. 4l). Highlighting

the capacity of scRNA-seq to provide information beyond bulk RNA-seq, we identify examples where Xi expression varies considerably between the two X-chromosomal haplotypes within an individual (for example, *ASMTL*; Supplementary Table 12), suggesting *cis*-acting variation as one of the determinants for the level of Xi expression<sup>3</sup>. As a further layer of heterogeneity in Xi expression, we find a unique pattern for *TIMP1*. For this gene, the level of Xi expression across cells is not significant, but exclusive to a subset of cells that express the gene biallelically (Extended Data Fig. 7), pointing to cell-to-cell variability in escape.

Using the ASE estimates from the scRNA-seq and GTEx analyses to infer the magnitude of the incompleteness of XCI, we find that expression from Xi at escape genes rarely reaches levels equal to expression from Xa, Xi expression remaining on average at 33% of Xa expression. However, there is a lot of variability along the chromosome (Extended Data Fig. 8a and Supplementary Discussion), as has previously been demonstrated in specific tissue types<sup>1,2</sup>. Balanced expression dosage between males and females in PAR1 requires full escape from XCI, however, Xi expression remains below Xa expression also in this region (mean Xi to Xa ratio is around 0.80), pointing to partial spreading of XCI beyond nonPAR. In further support that the consistent male bias in PAR1 expression (Fig. 2a) is due to the incompleteness of escape, we observe no systematic up- or downregulation of Y chromosome expression in PAR1 (Extended Data Fig. 8b and Supplementary Discussion).



**Figure 3 | Assessment of tissue-sharing of XCI in a GTEx donor with a highly skewed XCI.** **a**, Distribution of the skewness of XCI in GTEx female samples ( $n = 62$ , v3 release). Each data point shows the mean skew in XCI across tissue samples per individual. **b**, Classification of X-chromosomal genes ( $n = 186$ ) into full or incomplete and tissue-shared or heterogeneous XCI based on the analysis of ASE patterns across tissues. Error bars show the 95% credible interval. **c–e**, Examples of genes where the ASE-based assessment of XCI status match previously reported assignments (*TSR2*,

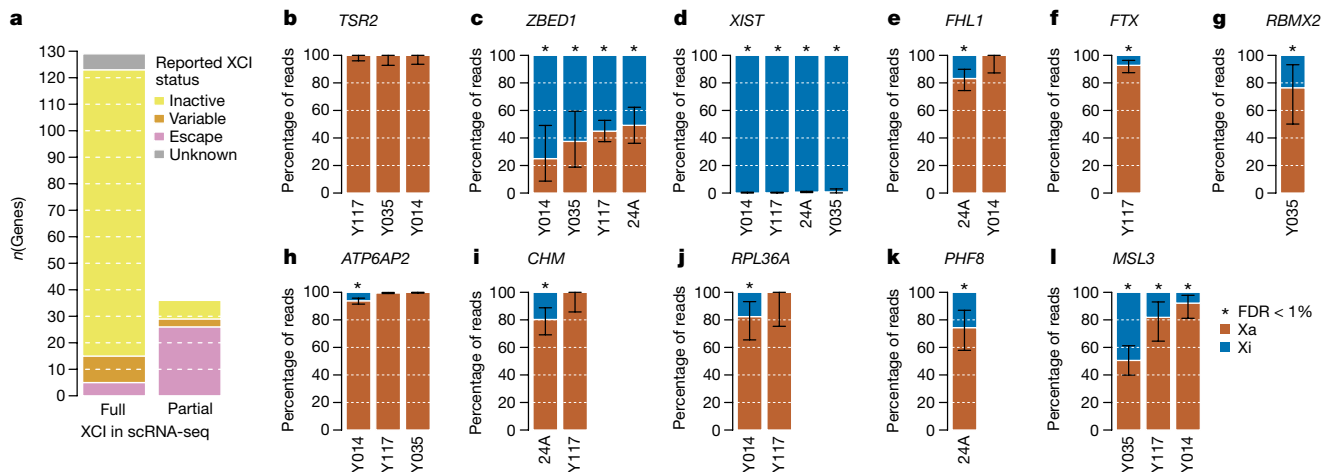
inactive; *XIST*, escape; *ZBED1*, escape). Note that *XIST* is only expressed monoallelically from Xi, which is unusual for an escape gene. **f**, *KAL1* shows strong evidence for tissue-specific escape. **g–k**, Genes without previous or conclusive evidence for escape from XCI that are classified as incompletely inactivated in this sample. In **c–k** asterisks indicate that the Xi expression in the given tissue was significant at  $FDR < 1\%$  (one-sided binomial test) and error bars show the 95% confidence interval.

As another consequence of the partial Xi expression, several of the X–Y homologous genes in nonPAR<sup>30</sup> become male-biased when expression from the Y chromosome counterpart is accounted for (Extended Data Fig. 8c).

By combining diverse types and analyses of high-throughput RNA-seq data, we have systematically assessed the incompleteness and heterogeneity in XCI across 29 human tissues (Supplementary Table 13). We establish that scRNA-seq is suitable for surveys of human XCI and present the first steps towards understanding the cellular-level variability in the maintenance of XCI. Our phasing-based approach enables the full use of low-coverage scRNA-seq, however, because any single

individual and cell type is only informative for restricted number of genes, larger datasets with more diverse cell types and conditions are required to fully profile XCI. We have therefore used the multi-tissue GTEx dataset to explore XCI in a larger number of X-chromosomal genes and to assess the tissue heterogeneity and impacts of XCI on gene expression differences between the sexes.

These analyses show that incomplete XCI is mostly shared between individuals and tissues, and extend previous surveys by pinpointing several examples of variability in the degree of XCI escape between cells, chromosomes, and tissues. In addition, our data demonstrate that escape from XCI results in sex-biased expression of at least 60 genes,



**Figure 4 | Analysis of XCI using scRNA-seq.** **a**, Proportion of genes demonstrating full and partial XCI in the ASE analysis in scRNA-seq data, and the concordance with previously reported XCI status. **b–l**, Examples of genes with different XCI patterns in scRNA-seq: previously reported inactive gene (**b**), known escape gene in PAR1 (**c**), escape gene with known exclusive expression from Xi (**d**), new candidates for escape genes that demonstrate incomplete XCI in only a subset of samples (**e–k**), and

a known escape gene that shows escape of varying degrees in the three samples (Pearson's  $\chi^2$  test for equal proportions,  $P = 3.80 \times 10^{-7}$ ) (**l**). **b–l**,  $x$  axis labels are sample identifiers. Asterisk above a bar indicates that the proportion of Xi expression, that is, blue bar, in a given sample is significantly greater than the expected baseline ( $FDR < 1\%$ , one-sided binomial test). Error bars show the 95% confidence interval.

potentially contributing to sex-specific differences in health and disease (Supplementary Discussion). As a whole, these results highlight the between-female and male–female diversity introduced by incomplete XCI, the biological implications of which remain to be fully explored.

**Online Content** Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

**Received 30 August 2016; accepted 8 September 2017.**

- Carrel, L. & Willard, H. F. X-inactivation profile reveals extensive variability in X-linked gene expression in females. *Nature* **434**, 400–404 (2005).
- Cotton, A. M. *et al.* Analysis of expressed SNPs identifies variable extents of expression from the human inactive X chromosome. *Genome Biol.* **14**, R122 (2013).
- Cotton, A. M. *et al.* Landscape of DNA methylation on the X chromosome reflects CpG density, functional chromatin state and X-chromosome inactivation. *Hum. Mol. Genet.* **24**, 1528–1539 (2015).
- Schultz, M. D. *et al.* Human body epigenome maps reveal noncanonical DNA methylation variation. *Nature* **523**, 212–216 (2015).
- Johnston, C. M. *et al.* Large-scale population study of human cell lines indicates that dosage compensation is virtually complete. *PLoS Genet.* **4**, e9 (2008).
- Tukialainen, T. *et al.* Chromosome X-wide association study identifies loci for fasting insulin and height and evidence for incomplete dosage compensation. *PLoS Genet.* **10**, e1004127 (2014).
- Deng, X., Berletch, J. B., Nguyen, D. K. & Distech, C. M. X chromosome regulation: diverse patterns in development, tissues and disease. *Nat. Rev. Genet.* **15**, 367–378 (2014).
- Amos-Landgraf, J. M. *et al.* X chromosome-inactivation patterns of 1,005 phenotypically unaffected females. *Am. J. Hum. Genet.* **79**, 493–499 (2006).
- Lister, R. *et al.* Global epigenomic reconfiguration during mammalian brain development. *Science* **341**, 1237905 (2013).
- Zhang, Y. *et al.* Transcriptional profiling of human liver identifies sex-biased genes associated with polygenic dyslipidemia and coronary artery disease. *PLoS ONE* **6**, e23506 (2011).
- Berletch, J. B. *et al.* Escape from X inactivation varies in mouse tissues. *PLoS Genet.* **11**, e1005079 (2015).
- The GTEx Consortium. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science* **348**, 648–660 (2015).
- The GTEx Consortium. Genetic effects on gene expression across tissues. <https://doi.org/10.1038/nature24277> (2017).
- Melé, M. *et al.* The human transcriptome across tissues and individuals. *Science* **348**, 660–665 (2015).
- Balaton, B. P. & Brown, C. J. Escape artists of the X chromosome. *Trends Genet.* **32**, 348–359 (2016).
- Ross, M. T. *et al.* The DNA sequence of the human X chromosome. *Nature* **434**, 325–337 (2005).
- Pál, C., Papp, B. & Hurst, L. D. Highly expressed genes in yeast evolve slowly. *Genetics* **158**, 927–931 (2001).
- Winter, E. E., Goodstadt, L. & Ponting, C. P. Elevated rates of protein secretion, evolution, and disease among tissue-specific genes. *Genome Res.* **14**, 54–61 (2004).
- Stegle, O., Teichmann, S. A. & Marioni, J. C. Computational and analytical challenges in single-cell transcriptomics. *Nat. Rev. Genet.* **16**, 133–145 (2015).
- Chen, G. *et al.* Single-cell analyses of X chromosome inactivation dynamics and pluripotency during differentiation. *Genome Res.* **26**, 1342–1354 (2016).
- Deng, Q., Ramsköld, D., Reinius, B. & Sandberg, R. Single-cell RNA-seq reveals dynamic, random monoallelic gene expression in mammalian cells. *Science* **343**, 193–196 (2014).
- Reinius, B. *et al.* Analysis of allelic expression patterns in clonal somatic cells by single-cell RNA-seq. *Nat. Genet.* **48**, 1430–1435 (2016).
- Wang, M., Lin, F., Xing, K. & Liu, L. Random X-chromosome inactivation dynamics *in vivo* by single-cell RNA sequencing. *BMC Genomics* **18**, 90 (2017).
- Wainer-Katsir, K. & Linnal, M. Single cell expression data reveal human genes that escape X-chromosome inactivation. Preprint at <http://www.biorxiv.org/content/early/2016/10/09/079830> (2016).

- Petropoulos, S. *et al.* Single-cell RNA-seq reveals lineage and X chromosome dynamics in human preimplantation embryos. *Cell* **165**, 1012–1026 (2016).
- Villani, A. C. *et al.* Single-cell RNA-seq reveals new types of human blood dendritic cells, monocytes, and progenitors. *Science* **356**, eaah4573 (2017).
- Borel, C. *et al.* Biased allelic expression in human primary fibroblast single cells. *Am. J. Hum. Genet.* **96**, 70–80 (2015).
- Kim, J. K., Kolodziejczyk, A. A., Ilicic, T., Teichmann, S. A. & Marioni, J. C. Characterizing noise structure in single-cell RNA-seq distinguishes genuine from technical stochastic allelic expression. *Nat. Commun.* **6**, 8687 (2015).
- Marinov, G. K. *et al.* From single-cell to cell-pool transcriptomes: stochasticity in gene expression and RNA splicing. *Genome Res.* **24**, 496–510 (2014).
- Bellott, D. W. *et al.* Mammalian Y chromosomes retain widely expressed dosage-sensitive regulators. *Nature* **508**, 494–499 (2014).
- Roadmap Epigenomics Consortium *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317–330 (2015).

**Supplementary Information** is available in the online version of the paper.

**Acknowledgements** We thank J. Maller, F. Zhao and M. Lek for technical assistance and P. J. Siponen for assistance with figure design. T.T. was supported by the Academy of Finland (285725), Finnish Cultural Foundation, Orion-Farmos Research Foundation and Emil Aaltonen Foundation. K.J.K. is supported by a NIGMS Fellowship (F32GM115208). This work was supported by NIH grants U54DK105566, R01MH101820 and R01GM104371 to D.G.M. The Genotype-Tissue Expression (GTEx) project was supported by the Common Fund of the Office of the Director of the National Institutes of Health. Additional funds were provided by the NCI, NHGRI, NHLBI, NIDA, NIMH and NINDS. Donors were enrolled at Biospecimen Source Sites funded by NCI\SAIC-Frederick, Inc. (SAIC-F) subcontracts to the National Disease Research Interchange (10XS170), Roswell Park Cancer Institute (10XS171) and Science Care, Inc. (X10S172). The Laboratory, Data Analysis, and Coordinating Center (LDACC) was funded through a contract (HHSN26820100029C) to The Broad Institute; this grant also provided funding to D.G.M. and T.T. Biorepository operations were funded through an SAIC-F subcontract to the Van Andel Institute (10ST1035). Additional data repository and project management were provided by SAIC-F (HHSN261200800001E). The Brain Bank was supported by supplements to University of Miami grants DA006227 and DA033684 and to contract N01MH000028. Statistical Methods development grants were made to the University of Geneva (MH090941 and MH101814), the University of Chicago (MH090951, MH090937, MH101820 and MH101825), the University of North Carolina, Chapel Hill (MH090936 and MH101819), Harvard University (MH090948), Stanford University (MH101782), Washington University St. Louis (MH101810) and the University of Pennsylvania (MH101822).

**Author Contributions** T.T. and D.G.M. designed the study. A.-C.V. designed and conducted the scRNA-seq experiments. T.T., A.Y., M.A.R., M.A., L.G., M.F. and B.B.C. analysed the data. J.L.M., R.S., S.E.C., A.K., K.J.K., F.A., A.B., T.L., A.R., K.G.A., N.H. and D.G.M. provided tools and reagents. T.T. and D.G.M. wrote the manuscript with input from other authors.

**Author Information** Reprints and permissions information is available at [www.nature.com/reprints](http://www.nature.com/reprints). The authors declare competing financial interests: details are available in the online version of the paper. Readers are welcome to comment on the online version of the paper. Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations. Correspondence and requests for materials should be addressed to T.T. ([ttuk@broadinstitute.org](mailto:ttuk@broadinstitute.org)) or D.G.M. ([danmac@broadinstitute.org](mailto:danmac@broadinstitute.org)).

**Reviewer Information** *Nature* thanks A. Clark and the other anonymous reviewer(s) for their contribution to the peer review of this work.



This work is licensed under a Creative Commons Attribution 4.0 International (CC BY 4.0) licence. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in the credit line; if the material is not included under the Creative Commons licence, users will need to obtain permission from the licence holder to reproduce the material. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

**GTEx Consortium****Laboratory, Data Analysis & Coordinating Center (LDACC)—Analysis Working Group**

François Aguet<sup>1</sup>, Kristin G. Ardlie<sup>1</sup>, Beryl B. Cummings<sup>1,2</sup>, Ellen T. Gelfand<sup>1</sup>, Gad Getz<sup>1,3</sup>, Kane Hadley<sup>1</sup>, Robert E. Handsaker<sup>1,4</sup>, Katherine H. Huang<sup>1</sup>, Seva Kashin<sup>1,4</sup>, Konrad J. Karczewski<sup>1,2</sup>, Monkol Lek<sup>1,2</sup>, Xiao Li<sup>1</sup>, Daniel G. MacArthur<sup>1,2</sup>, Jared L. Nedzel<sup>1</sup>, Duyen T. Nguyen<sup>1</sup>, Michael S. Noble<sup>1</sup>, Ayellet V. Segre<sup>1</sup>, Casandra A. Trowbridge<sup>1</sup>, Taru Tuikainen<sup>1,2</sup>

**Statistical Methods groups—Analysis Working Group**

Nathan S. Abel<sup>5,6</sup>, Brunilda Balliu<sup>5</sup>, Ruth Barshir<sup>7</sup>, Omer Basha<sup>7</sup>, Alexis Battle<sup>8</sup>, Gireesh K. Bogu<sup>9,10</sup>, Andrew Brown<sup>11,12,13</sup>, Christopher D. Brown<sup>14</sup>, Stephane E. Castel<sup>15,16</sup>, Lin S. Chen<sup>17</sup>, Colby Chiang<sup>18</sup>, Donald F. Conrad<sup>19,20</sup>, Nancy J. Cox<sup>21</sup>, Farhan N. Damani<sup>8</sup>, Joe R. Davis<sup>5,6</sup>, Olivier Delaneau<sup>11,12,13</sup>, Emmanouil T. Dermizakis<sup>11,12,13</sup>, Barbara E. Engelhardt<sup>22</sup>, Eleazar Eskin<sup>23,24</sup>, Pedro G. Ferreira<sup>25,26</sup>, Laure Frésard<sup>5,6</sup>, Eric R. Gamazon<sup>21,27,28</sup>, Diego Garrido-Martín<sup>9,10</sup>, Ariel D.H. Gewirtz<sup>29</sup>, Genna Gliner<sup>30</sup>, Michael J. Gloudemans<sup>5,6,31</sup>, Roderic Guigo<sup>9,10,32</sup>, Ira M. Hall<sup>18,19,33</sup>, Buhm Han<sup>34</sup>, Yuan He<sup>35</sup>, Farhad Hormozdizari<sup>23</sup>, Cedric Howald<sup>11,12,13</sup>, Hae Kyung Im<sup>36</sup>, Brian Jo<sup>29</sup>, Eun Yong Kang<sup>23</sup>, Yungil Kim<sup>8</sup>, Sarah Kim-Hellmuth<sup>15,16</sup>, Tuuli Lappalainen<sup>15,16</sup>, Gen Li<sup>37</sup>, Xin Li<sup>6</sup>, Boxiang Liu<sup>5,6,38</sup>, Serghei Mangul<sup>23</sup>, Mark I. McCarthy<sup>39,40,41</sup>, Ian C. McDowell<sup>42</sup>, Pejman Mohammadi<sup>15,16</sup>, Jean Monlong<sup>9,10,43</sup>, Stephen B. Montgomery<sup>5,6</sup>, Manuel Muñoz-Aguirre<sup>9,10,44</sup>, Anne W. Ndungu<sup>39</sup>, Dan L. Nicolae<sup>36,45,46</sup>, Andrew B. Nobel<sup>47,48</sup>, Meritxell Oliva<sup>36,49</sup>, Halit Ongen<sup>11,12,13</sup>, John J. Palowitch<sup>47</sup>, Nikolaos Panousis<sup>11,12,13</sup>, Panagiotis Papanikolaou<sup>9,10</sup>, YoSon Park<sup>14</sup>, Princy Parsana<sup>8</sup>, Anthony J. Payne<sup>39</sup>, Christine B. Peterson<sup>50</sup>, Jie Quan<sup>51</sup>, Ferran Reverter<sup>9,10,52</sup>, Chiara Sabbati<sup>53,54</sup>, Ashis Saha<sup>8</sup>, Michael Sammeth<sup>55</sup>, Alexandra J. Scott<sup>18</sup>, Andrew A. Shabalin<sup>56</sup>, Reza Sodaei<sup>9,10</sup>, Matthew Stephens<sup>45,46</sup>, Barbara E. Stranger<sup>36,49,57</sup>, Benjamin J. Strober<sup>35</sup>, Jae Hoon Sul<sup>58</sup>, Emily K. Tsang<sup>63,1</sup>, Sarah Urbut<sup>46</sup>, Martijn van de Bunt<sup>39,40</sup>, Gao Wang<sup>46</sup>, Xiaoquan Wen<sup>59</sup>, Fred A. Wright<sup>60</sup>, Hualin S. Xi<sup>51</sup>, Esti Yeger-Lotem<sup>7,61</sup>, Zachary Zappala<sup>56</sup>, Judith B. Zaugg<sup>62</sup>, Yi-Hui Zhou<sup>60</sup>

**Enhancing GTEx (eGTEx) groups** Joshua M. Akey<sup>29,63</sup>, Daniel Bates<sup>64</sup>, Joanne Chan<sup>5</sup>, Lin S. Chen<sup>17</sup>, Melina Claussnitzer<sup>1,65,66</sup>, Kathryn Demanelis<sup>17</sup>, Morgan Diegel<sup>64</sup>, Jennifer A. Doherty<sup>67</sup>, Andrew P. Feinberg<sup>35,68,69,70</sup>, Marian S. Fernando<sup>36,49</sup>, Jessica Halow<sup>64</sup>, Kasper D. Hansen<sup>68,71,72</sup>, Eric Haugen<sup>64</sup>, Peter F. Hickey<sup>72</sup>, Lei Hou<sup>1,73</sup>, Farzana Jasmine<sup>17</sup>, Ruiqi Jian<sup>5</sup>, Lihua Jiang<sup>5</sup>, Audra Johnson<sup>64</sup>, Rajinder Kaul<sup>64</sup>, Manolis Kellis<sup>1,73</sup>, Muhammad G. Kibriya<sup>17</sup>, Kristen Lee<sup>64</sup>, Jin Billy Li<sup>5</sup>, Qin Li<sup>5</sup>, Xiao Li<sup>5</sup>, Jessica Lin<sup>5,74</sup>, Shinn Lin<sup>5,75</sup>, Sandra Linder<sup>5,6</sup>, Caroline Linka<sup>36,49</sup>, Yaping Liu<sup>1,73</sup>, Matthew T. Maurano<sup>76</sup>, Benoit Molinie<sup>1</sup>, Stephen B. Montgomery<sup>5,6</sup>, Jemma Nelson<sup>64</sup>, Fidencio J. Neri<sup>64</sup>, Meritxell Oliva<sup>36,49</sup>, Yongjin Park<sup>1,73</sup>, Brandon L. Pierce<sup>17</sup>, Nicola J. Rinaldi<sup>1,73</sup>, Lindsay F. Rizzardi<sup>68</sup>, Richard Sandstrom<sup>64</sup>, Andrew Skol<sup>36,49,57</sup>, Kevin S. Smith<sup>5,6</sup>, Michael P. Snyder<sup>5</sup>, John Stamatoyannopoulos<sup>64,74,77</sup>, Barbara E. Stranger<sup>36,49,57</sup>, Hua Tang<sup>5</sup>, Emily K. Tsang<sup>63,1</sup>, Li Wang<sup>41</sup>, Meng Wang<sup>5</sup>, Nicholas Van Wittenbergh<sup>1</sup>, Fan Wu<sup>36,49</sup>, Rui Zhang<sup>5</sup>

**NIH Common Fund** Concepcion R. Nierras<sup>78</sup>

**NIH/NCI** Philip A. Branton<sup>79</sup>, Latarsha J. Carithers<sup>79,80</sup>, Ping Guan<sup>79</sup>, Helen M. Moore<sup>79</sup>, Abhi Rao<sup>79</sup>, Jimmie B. Vaught<sup>79</sup>

**NIH/NHGRI** Sarah E. Gould<sup>81</sup>, Nicole C. Lockart<sup>81</sup>, Casey Martin<sup>81</sup>, Jeffery P. Struewing<sup>81</sup>, Simona Volpi<sup>81</sup>

**NIH/NIMH** Anjene M. Addington<sup>82</sup>, Susan E. Koester<sup>82</sup>

**NIH/NIDA** A. Roger Little<sup>83</sup>

**Biospecimen Collection Source Site—NDRI** Lori E. Brigham<sup>84</sup>, Richard Hasz<sup>85</sup>, Marcus Hunter<sup>86</sup>, Christopher Johns<sup>87</sup>, Mark Johnson<sup>88</sup>, Gene Kopen<sup>89</sup>, William F. Levinweber<sup>89</sup>, John T. Lonsdale<sup>89</sup>, Alisa McDonald<sup>89</sup>, Bernadette Mesticelli<sup>89</sup>, Kein Myer<sup>86</sup>, Brian Roe<sup>86</sup>, Michael Salvatore<sup>89</sup>, Saboor Shad<sup>89</sup>, Jeffrey A. Thomas<sup>89</sup>, Gary Walters<sup>88</sup>, Michael Washington<sup>88</sup>, Joseph Wheeler<sup>87</sup>

**Biospecimen Collection Source Site—RPCI** Jason Bridge<sup>90</sup>, Barbara A. Foster<sup>91</sup>, Bryan M. Gillard<sup>91</sup>, Ellen Karasik<sup>91</sup>, Rachna Kumar<sup>91</sup>, Mark Miklos<sup>90</sup>, Michael T. Moser<sup>91</sup>

**Biospecimen Core Resource—VARI** Scott D. Jewell<sup>92</sup>, Robert G. Montroy<sup>92</sup>, Daniel C. Rohrer<sup>92</sup>, Dana R. Valley<sup>92</sup>

**Brain Bank Repository—University of Miami Brain Endowment Bank** David A. Davis<sup>93</sup>, Deborah C. Mash<sup>93</sup>

**Leidos Biomedical—Project Management** Anita H. Undale<sup>94</sup>, Anna M. Smith<sup>95</sup>, David E. Tabor<sup>95</sup>, Nancy V. Roche<sup>95</sup>, Jeffrey A. McLean<sup>95</sup>, Negin Vatanian<sup>95</sup>, Karna L. Robinson<sup>95</sup>, Leslie Sobin<sup>95</sup>, Mary E. Barcus<sup>96</sup>, Kimberly M. Valentino<sup>95</sup>, Liqun Qi<sup>95</sup>, Steven Hunter<sup>95</sup>, Pushpa Hariharan<sup>95</sup>, Shilpi Singh<sup>95</sup>, Ki Sung Um<sup>95</sup>, Takunda Matose<sup>95</sup>, Maria M. Tomaszewski<sup>95</sup>

**ELSI Study** Laura K. Barker<sup>97</sup>, Maghboeba Mosavel<sup>98</sup>, Laura A. Siminoff<sup>97</sup>, Heather M. Traino<sup>97</sup>

**Genome Browser Data Integration & Visualization—EBI** Paul Flicek<sup>99</sup>, Thomas Juettemann<sup>99</sup>, Magali Ruffier<sup>99</sup>, Dan Sheppard<sup>99</sup>, Kieron Taylor<sup>99</sup>, Stephen J. Trevanion<sup>99</sup>, Daniel R. Zerbino<sup>99</sup>

**Genome Browser Data Integration & Visualization—UCSC Genomics Institute, University of California Santa Cruz** Brian Craft<sup>100</sup>, Mary Goldman<sup>100</sup>, Maximilian Haussler<sup>100</sup>, W. James Kent<sup>100</sup>, Christopher M. Lee<sup>100</sup>, Benedict Paten<sup>100</sup>, Kate R. Rosenbloom<sup>100</sup>, John Vivian<sup>100</sup>, Jingchun Zhu<sup>100</sup>

<sup>1</sup>The Broad Institute of Massachusetts Institute of Technology and Harvard University, Cambridge, Massachusetts 02142, USA. <sup>2</sup>Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, Massachusetts 02114, USA. <sup>3</sup>Massachusetts General Hospital Cancer Center and Department of Pathology, Massachusetts General Hospital, Boston, Massachusetts 02114, USA. <sup>4</sup>Department of Genetics, Harvard Medical School, Boston, Massachusetts 02114, USA. <sup>5</sup>Department of Genetics, Stanford University, Stanford, California 94305, USA. <sup>6</sup>Department of Pathology, Stanford University, Stanford, California 94305, USA. <sup>7</sup>Department of Clinical Biochemistry and Pharmacology, Faculty of Health Sciences, Ben-Gurion University of the Negev, Beer-Sheva 84105, Israel. <sup>8</sup>Department of Computer Science, Johns Hopkins University, Baltimore, Maryland 21218, USA. <sup>9</sup>Centre for Genomic Regulation (CRG), The Barcelona Institute for Science and Technology, 08003 Barcelona, Spain. <sup>10</sup>Universitat Pompeu Fabra (UPF), 08002 Barcelona, Spain. <sup>11</sup>Department of Genetic Medicine and Development, University of Geneva Medical School, 1211 Geneva, Switzerland. <sup>12</sup>Institute for Genetics and Genomics in Geneva (iG3), University of Geneva, 1211 Geneva, Switzerland. <sup>13</sup>Swiss Institute of Bioinformatics, 1211 Geneva, Switzerland. <sup>14</sup>Department of Genetics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, Pennsylvania 19104, USA. <sup>15</sup>New York Genome Center, New York, New York 10013, USA. <sup>16</sup>Department of Systems Biology, Columbia University Medical Center, New York, New York 10032, USA. <sup>17</sup>Department of Public Health Sciences, The University of Chicago, Chicago, Illinois 60637, USA. <sup>18</sup>McDonnell Genome Institute, Washington University School of Medicine, St. Louis, Missouri 63108, USA. <sup>19</sup>Department of Genetics, Washington University School of Medicine, St. Louis, Missouri 63108, USA. <sup>20</sup>Department of Pathology & Immunology, Washington University School of Medicine, St. Louis, Missouri 63108, USA. <sup>21</sup>Division of Genetic Medicine, Department of Medicine, Vanderbilt University Medical Center, Nashville, Tennessee 37232, USA. <sup>22</sup>Department of Computer Science, Center for Statistics and Machine Learning, Princeton University, Princeton, New Jersey 08540, USA. <sup>23</sup>Department of Computer Science, University of California, Los Angeles, California 90095, USA. <sup>24</sup>Department of Human Genetics, University of California, Los Angeles, California 90095, USA. <sup>25</sup>Instituto de Investigação e Inovação em Saúde (i3S), Universidade do Porto, 4200-135 Porto, Portugal. <sup>26</sup>Institute of Molecular Pathology and Immunology (IPATIMUP), University of Porto, 4200-625 Porto, Portugal. <sup>27</sup>Department of Clinical Epidemiology, Biostatistics and Bioinformatics, Academic Medical Center, University of Amsterdam, 1105 AZ Amsterdam, The Netherlands. <sup>28</sup>Department of Psychiatry, Academic Medical Center, University of Amsterdam, 1105 AZ Amsterdam, The Netherlands. <sup>29</sup>Lewis Sigler Institute, Princeton University, Princeton, New Jersey 08540, USA. <sup>30</sup>Department of Operations Research and Financial Engineering, Princeton University, Princeton, New Jersey 08540, USA. <sup>31</sup>Biomedical Informatics Program, Stanford University, Stanford, California 94305, USA. <sup>32</sup>Institut Hospital del Mar d'Investigacions Mèdiques (IMIM), 08003 Barcelona, Spain. <sup>33</sup>Department of Medicine, Washington University School of Medicine, St. Louis, Missouri 63108, USA. <sup>34</sup>Department of Convergence Medicine, University of Ulsan College of Medicine, Asan Medical Center, Seoul 138-736, South Korea. <sup>35</sup>Department of Biomedical Engineering, Johns Hopkins University, Baltimore, Maryland 21218, USA. <sup>36</sup>Section of Genetic Medicine, Department of Medicine, The University of Chicago, Chicago, Illinois 60637, USA. <sup>37</sup>Department of Biostatistics, Mailman School of Public Health, Columbia University, New York, New York 10032, USA. <sup>38</sup>Department of Human Genetics, Stanford University, Stanford, California 94305, USA. <sup>39</sup>Wellcome Trust Centre for Human Genetics, Nuffield Department of Medicine, University of Oxford, Oxford OX3 7BN, UK. <sup>40</sup>Oxford Centre for Diabetes, Endocrinology and Metabolism, University of Oxford, Churchill Hospital, Oxford OX3 7LE, UK. <sup>41</sup>Oxford NIHR Biomedical Research Centre, Churchill Hospital, Oxford OX3 7LJ, UK. <sup>42</sup>Computational Biology & Bioinformatics Graduate Program, Duke University, Durham, North Carolina 27708, USA. <sup>43</sup>Human Genetics Department, McGill University, Montreal, Quebec H3A 0G1, Canada. <sup>44</sup>Departament d'Estadística i Investigació Operativa, Universitat Politècnica de Catalunya, 08034 Barcelona, Spain. <sup>45</sup>Department of Statistics, The University of Chicago, Chicago, Illinois 60637, USA. <sup>46</sup>Department of Human Genetics, The University of Chicago, Chicago, Illinois 60637, USA. <sup>47</sup>Department of Statistics and Operations Research, University of North Carolina, Chapel Hill, North Carolina 27599, USA. <sup>48</sup>Department of Biostatistics, University of North Carolina, Chapel Hill, North Carolina 27599, USA. <sup>49</sup>Institute for Genomics and Systems Biology, The University of Chicago, Chicago, Illinois 60637, USA. <sup>50</sup>Department of Biostatistics, The University of Texas MD Anderson Cancer Center, Houston, Texas 77030, USA. <sup>51</sup>Computational Sciences, Pfizer Inc, Cambridge, Massachusetts 02139, USA. <sup>52</sup>Universitat de Barcelona, 08028 Barcelona, Spain. <sup>53</sup>Department of Biomedical Data Science, Stanford University, Stanford, California 94305, USA. <sup>54</sup>Department of Statistics, Stanford University, Stanford, California 94305, USA. <sup>55</sup>Institute of Biophysics Carlos Chagas Filho (IBCCF), Federal University of Rio de Janeiro (UFRJ), 21941902 Rio de Janeiro, Brazil. <sup>56</sup>Department of Psychiatry, University of Utah, Salt Lake City, Utah 84108, USA. <sup>57</sup>Center for Data Intensive Science, The University of Chicago, Chicago, Illinois 60637, USA. <sup>58</sup>Department of Psychiatry and Biobehavioral Sciences, University of California, Los Angeles, California 90095, USA. <sup>59</sup>Department of Biostatistics, University of Michigan, Ann Arbor, Michigan 48109, USA. <sup>60</sup>Biostatistics Research Center and Departments of Statistics and Biological Sciences, North Carolina State University, Raleigh, North Carolina 27695, USA. <sup>61</sup>National Institute for Biotechnology in the Negev, Beer-Sheva, 84105, Israel. <sup>62</sup>European Molecular Biology Laboratory, 69117 Heidelberg, Germany. <sup>63</sup>Department of Ecology and Evolutionary Biology, Princeton University, Princeton, New Jersey 08540, USA. <sup>64</sup>Altius Institute for Biomedical Sciences, Seattle, Washington 98121, USA. <sup>65</sup>Beth Israel Deaconess Medical Center, Harvard Medical School, Boston, Massachusetts 02215, USA. <sup>66</sup>University of Hohenheim, 70599 Stuttgart, Germany. <sup>67</sup>Huntsman Cancer Institute, Department of Population Health Sciences, University of Utah, Salt Lake City, Utah 84112, USA. <sup>68</sup>Center for Epigenetics, Johns Hopkins University School of Medicine, Baltimore, Maryland 21205, USA. <sup>69</sup>Department of Medicine, Johns Hopkins University School of Medicine, Baltimore, Maryland 21205, USA. <sup>70</sup>Department of Mental Health, Johns Hopkins University School of Public Health, Baltimore, Maryland 21205, USA. <sup>71</sup>McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins School of Medicine, Baltimore, Maryland 21205, USA. <sup>72</sup>Department of Biostatistics, Johns Hopkins

University, Baltimore, Maryland 21205, USA. <sup>73</sup>Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA.

<sup>74</sup>Department of Medicine, University of Washington, Seattle, Washington 98195, USA.

<sup>75</sup>Division of Cardiology, University of Washington, Seattle, Washington 98195, USA. <sup>76</sup>Institute for Systems Genetics, New York University Langone Medical Center, New York, New York 10016, USA.

<sup>77</sup>Department of Genome Sciences, University of Washington, Seattle, Washington 98195, USA. <sup>78</sup>Office of Strategic Coordination, Division of Program Coordination, Planning and Strategic Initiatives, Office of the Director, NIH, Rockville, Maryland 20852, USA.

<sup>79</sup>Biorepositories and Biospecimen Research Branch, Division of Cancer Treatment and Diagnosis, National Cancer Institute, Bethesda, Maryland 20892, USA. <sup>80</sup>National Institute of Dental and Craniofacial Research, Bethesda, Maryland 20892, USA.

<sup>81</sup>Division of Genomic Medicine, National Human Genome Research Institute, Rockville, Maryland 20852, USA. <sup>82</sup>Division of Neuroscience and Basic Behavioral Science, National Institute of Mental Health, NIH, Bethesda, Maryland 20892, USA.

<sup>83</sup>Division of Neuroscience and Behavior, National Institute on Drug Abuse, NIH, Bethesda, Maryland 20892, USA. <sup>84</sup>Washington Regional

Transplant Community, Falls Church, Virginia 22003, USA. <sup>85</sup>Gift of Life Donor Program, Philadelphia, Pennsylvania 19103, USA. <sup>86</sup>LifeGift, Houston, Texas 77055, USA. <sup>87</sup>Center for Organ Recovery and Education, Pittsburgh, Pennsylvania 15238, USA. <sup>88</sup>LifeNet Health, Virginia Beach, Virginia 23453, USA. <sup>89</sup>National Disease Research Interchange, Philadelphia, Pennsylvania 19103, USA. <sup>90</sup>Unyts, Buffalo, New York 14203, USA. <sup>91</sup>Pharmacology and Therapeutics, Roswell Park Cancer Institute, Buffalo, New York 14263, USA. <sup>92</sup>Van Andel Research Institute, Grand Rapids, Michigan 49503, USA. <sup>93</sup>Brain Endowment Bank, Miller School of Medicine, University of Miami, Miami, Florida 33136, USA. <sup>94</sup>National Institute of Allergy and Infectious Diseases, NIH, Rockville, Maryland 20852, USA. <sup>95</sup>Biospecimen Research Group, Clinical Research Directorate, Leidos Biomedical Research, Inc., Rockville, Maryland 20852, USA. <sup>96</sup>Leidos Biomedical Research, Inc., Frederick, Maryland 21701, USA. <sup>97</sup>Temple University, Philadelphia, Pennsylvania 19122, USA. <sup>98</sup>Department of Health Behavior and Policy, School of Medicine, Virginia Commonwealth University, Richmond, Virginia 23298, USA. <sup>99</sup>European Molecular Biology Laboratory, European Bioinformatics Institute, Hinxton CB10 1SD, UK. <sup>100</sup>UCSC Genomics Institute, University of California Santa Cruz, Santa Cruz, California 95064, USA.

## METHODS

**GTEX data.** The GTEX project<sup>12</sup> collected tissue samples from 554 postmortem donors (187 females, 357 males; age range, 20–70), carried out RNA-seq on 8,555 tissue samples and generated genotyping data for up to 449 donors (GTEX analysis v6p release). More detailed methods can be found in ref. 13. All GTEX data, including RNA, genome and exome sequencing data, used in the analyses described here are available through dbGaP under accession number phs000424.v6.p1, unless otherwise stated. Summary data and details on data production and processing are also available from the GTEX Portal (<http://gtexportal.org>).

**Single-cell samples.** For the human dendritic cells samples profiled, the healthy donor (ID: 24A) was recruited from the Boston-based PhenoGenetic project, a resource of healthy subjects that are re-contactable by genotype<sup>32</sup>. The donor was a female Asian individual from China, 25 years of age at the time of blood collection. She was a non-smoker, had a normal BMI (height: 168.7 cm; weight: 56.45 kg; BMI: 19.8), and normal blood pressure (108/74). The donor had no family history of cancer, allergies, inflammatory disease, autoimmune disease, chronic metabolic disorders or infectious disorders. She provided written informed consent for the genetic research studies and molecular testing, as previously reported<sup>26</sup>.

Daughters of three parent–child Yoruba trios from Ibadan, Nigeria (that is, YRI trios), collected as part of the International HapMap Project, were chosen for single-cell profiling, both to maximize heterozygosity and due to availability of parental genotypes enabling phasing. DNA and LCLs were ordered from the NHGRI Sample Repository for Human Genetic Research (Coriell Institute for Medical Research): LCLs from B lymphocytes for the three daughters (catalogue numbers: GM19240, GM19199 and GM18518) and DNA extracted from LCLs for all members of the three trios (catalogue numbers for DNA: NA19240, NA19238, NA19239, NA19199, NA19197, NA19198, NA18518, NA18519 and NA18520). These YRI samples are referred to by their family IDs: Y014, Y035 and Y117.

**Clinical muscle samples.** To assess whether PARI genes are equally expressed from X and Y chromosomes, a combination of skeletal muscle RNA-seq data and trio genotyping data from eight male patients with muscular dystrophy, sequenced as part of an unrelated study, was used. Patient cases with available muscle biopsies were referred from clinicians starting April 2013 until June 2016. All patients included for RNA-seq had previously available trio whole-exome sequencing (WES) data, with one sample having additional trio whole-genome sequencing (WGS). Muscle biopsies were shipped frozen from clinical centres by liquid nitrogen dry shipping and, where possible, frozen muscle was sectioned on a cryostat and stained with haematoxylin and eosin to assess muscle quality as well as the presence of overt freeze–thaw artefacts.

**Genotyping.** The GTEX v6p release includes WGS data for 148 donors, including GTEX-UPIC. WGS libraries were sequenced on the Illumina HiSeqX or Illumina HiSeq2000. WGS data was processed through a Picard-based pipeline, using base quality score recalibration and local realignment at known indels. BWA-MEM aligner was used for mapping reads to the human genome build 37 (hg19). Single-nucleotide polymorphisms (SNPs) and indels (insertions and deletions) were jointly called across all 148 samples and additional reference genomes using HaplotypeCaller v.3.1 of GATK. Default filters were applied to SNP and indel calls using the variant quality score recalibration (VQSR) approach of GATK. An additional hard filter  $\text{InbreedingCoeff} \leq -0.3$  was applied to remove sites that VQSR failed to filter.

WGS for one of the clinical muscle samples was performed on 500 ng to 1.5  $\mu\text{g}$  of genomic DNA using a PCR-Free protocol that substantially increases the uniformity of genome coverage. These libraries were sequenced on the Illumina HiSeq X10 with 151-bp paired-end reads and a target mean coverage of  $>30\times$ , and were processed similarly to the above description.

The Y117 trio (sample IDs NA19240 (daughter), NA19238 (mother), and NA19239 (father)) was whole-genome-sequenced as part of the 1000 Genomes Project as described previously<sup>33</sup>. The VCF file containing the WGS-based genotypes for SNPs (YRI.trio.2010\_09.genotypes.vcf.gz) was downloaded from the FTP site of the project. The genotype coordinates (in human genome build 36) in the original VCF were converted to hg19 using the liftover script (liftOverVCF.pl) and chain files provided as part of the GATK package.

WES was performed using Illumina's capture Exome (ICE) technology (Y035, Y014, 24A) or Agilent SureSelect Human All Exon Kit v.2 exome capture (clinical muscle samples) with a mean target coverage of  $>80\times$ . WES data was aligned with BWA, processed with Picard, and SNPs and indels were jointly called with other samples using GATK HaplotypeCaller package v.3.1 (24A, clinical muscle samples) or v.3.4 (Y035, Y014). Default filters were applied to SNP and indel calls using the VQSR approach. A modified version of the Ensembl variant effect predictor was used for variant annotation for all WES and WGS data. For trio WES or WGS data the genotypes of the proband were phased using the PhaseByTransmission tool of the GATK toolkit.

**Single-cell data preparation and sequencing.** For profiling of healthy dendritic cells (DCs), peripheral blood mononuclear cells (PBMCs) were first isolated from fresh blood within 2 h of collection, using Ficoll–Paque density gradient centrifugation as previously described<sup>34</sup>. Single-cell suspensions were stained as per the manufacturer's recommendations with an antibody panel designed to enrich for all known blood DC population for single-cell sorting and scRNA-seq profiling<sup>26</sup>. A total of 24 single cells from four loosely gated populations were sorted per 96-well plate, with each well containing 10  $\mu\text{l}$  of lysis buffer. A total of eight plates were analysed by scRNA-seq.

All LCL cell lines were cultured according to Coriell's recommendations (medium: RPMI 1640, 2 mM L-glutamine, 15% fetal bovine serum (all three from ThermoFisher Scientific)) in T25 tissue culture flask with 10–20 ml medium at 37°C in 5% carbon dioxide. Cells were split upon reaching a cell density of approximately 300,000–400,000 viable cells per ml. All three lymphoblast cultures were split once before single-cell sorting. Cells were washed with  $1\times$  PBS, the pellet was resuspended and stained with DAPI (Biolegend) for viability according to the manufacturer's recommendations.

All single live cells (for both DCs and LCL cell lines) were sorted into a 96-well full-skirted Eppendorf plate chilled to 4°C, that were pre-prepared with 10  $\mu\text{l}$  TCL buffer (Qiagen) supplemented with 1%  $\beta$ -mercaptoethanol (lysis buffer), using a BD FACS Fusion instrument. Single-cell lysates were sealed, vortexed, spun down at 300g at 4°C for 1 min, immediately placed on dry ice and transferred for storage at  $-80^\circ\text{C}$ .

The Smart-Seq2 protocol was performed on single-sorted cells as described<sup>35,36</sup>, with some modifications as described in ref. 26 (Supplementary Methods). A total of 768 single DCs isolated from a healthy Asian female individual, along with 96 single cells from GM19240, 48 single cells from GM19199 and 48 single cells from GM18518 were profiled. In brief, single-cell lysates were thawed on ice, purified and reverse-transcribed using Maxima H Minus Reverse Transcriptase. PCR was performed with KAPA HiFi HotStart ReadyMix (KAPA Biosystems) and purified with Agencourt AMPureXP SPRI beads (Beckman-Coulter). The concentration of amplified cDNA was measured on the Synergy H1 Hybrid Microplate Reader (BioTek) using High-Sensitivity Qubit reagent (Life Technologies) and the size distribution of select wells was checked on a High-Sensitivity Bioanalyzer Chip (Agilent). The expected concentration was around 0.5–2  $\text{ng}\mu\text{l}^{-1}$  with a size distribution that sharply peaked around 2 kb.

Library preparation was carried out using the Nextera XT DNA Sample Kit (Illumina) with custom indexing adapters, allowing up to 384 libraries to be simultaneously generated in a 384-well PCR plate (note that DCs were processed in a 384-well plate whereas LCLs were processed in 96-well plate format). The concentration of the final pooled libraries was measured using the High-Sensitivity DNA Qubit (Life Technologies) and the size distribution was measured on a High-Sensitivity Bioanalyzer Chip (Agilent). The expected concentration of the pooled libraries was 10–30  $\text{ng}\mu\text{l}^{-1}$  with a size distribution of 300–700 bp. For the DCs, we created pools of 384 cells, whereas 96 LCL samples were pooled at the time. We sequenced one library pool per lane as paired-end 25-bp reads on a HiSeq2500 (Illumina). Barcodes used for indexing are listed in the Supplementary Methods.

**RNA-seq in GTEX.** RNA sequencing was performed using a non-strand-specific RNA-seq protocol with polyA selection of RNA using the Illumina TruSeq protocol with sequence coverage goal of 50 million 76-bp paired-end reads as has been previously described in detail<sup>12</sup>. The RNA-seq data, except for GTEX-UPIC, was aligned with TopHat v.1.4.1 to the UCSC human genome release version hg19 using the Gencode v.19 annotations as the transcriptome reference. Gene level read counts and reads per kilobase per million reads (RPKM) were derived using the RNA-SeQC tool<sup>37</sup> using the Gencode v.19 transcriptome annotation. The transcript model was collapsed into a gene model as described previously<sup>12</sup>. Read count and RPKM quantification include only uniquely mapped and properly paired reads contained within exon boundaries.

**RNA-seq alignment to personalized genomes.** For the four single-cell samples and for GTEX-UPIC RNA-seq, data were processed using a modification of the AlleleSeq pipeline<sup>38,39</sup> to minimize reference allele bias in alignment. A diploid personal reference genome for each of the samples was generated with the vcf2diploid tool<sup>38</sup> including all heterozygous biallelic single-nucleotide variants identified in WES or WGS either together with (YRI samples) or without (GTEX-UPIC, 24A) maternal and paternal genotype information. The RNA-seq reads were then aligned to both parental references using STAR<sup>40</sup> v.2.4.1a in a per-sample two-pass mode (GTEX-UPIC and YRI samples) or v.2.3.0e (24A) using hg19 as the reference. The alignments were combined by comparing the quality of alignment between the two references: for reads aligning uniquely to both references the alignment with the higher alignment score was chosen and reads aligning uniquely to only one reference were kept as such.



**RNA-seq of clinical muscle samples.** Patient RNA samples derived from primary muscle were sequenced using the GTEx sequencing protocol<sup>12</sup> with sequence coverage of 50 million or 100 million 76-bp paired-end reads. RNA-seq reads were aligned using STAR<sup>40</sup> 2-pass version v.2.4.2a using hg19 as the reference genome. Junctions were filtered after first pass alignment to exclude junctions with less than 5 uniquely mapped reads supporting the event and junctions found on the mitochondrial genome. The value for unique mapping quality was assigned to 60 and duplicate reads were marked with Picard MarkDuplicates (v.1.1099).

**Catalogue of X-inactivation status.** To compare results from the ASE and GTEx analyses with previous observations on genic XCI status we collated findings from two earlier studies<sup>1,2</sup> that represent systematic expression-based surveys into XCI. Each study catalogues hundreds of X-linked genes and together the data span two tissue types.

Carrel and Willard<sup>1</sup> surveyed in total 624 X-chromosomal transcripts expressed in primary fibroblasts in nine cell hybrids each containing a different human Xi. In order to find the gene corresponding to each transcript, the primer sequences designed to test the expression of the transcripts in the original study were aligned to reference databases based on the Gencode v.19 transcriptome and hg19 using in-house software (unpublished) (Supplementary Methods). In total 553 transcripts primer pairs were successfully matched to X-chromosomal Gencode v.19 reference mapping together with 470 unique X-chromosomal genes (Supplementary Methods). These 470 genes were split into three XCI status categories (escape, variable, inactive) based on the level of Xi expression (that is, the number of cell lines expressing the gene from Xi) resulting in 75 escape, 51 variable escape and 344 inactive genes.

Cotton *et al.*<sup>2</sup> surveyed XCI using allelic imbalance in clonal or near-clonal female LCL and fibroblast cell lines and provided XCI statuses for 508 genes (68 escape, 146 variable escape, 294 subject genes). The data were mapped to Gencode v.19 using the reported gene names and their known aliases (Supplementary Methods), resulting in a list of XCI statuses for 506 X-chromosomal genes.

The results were combined by retaining the XCI status in the original study where possible (that is, same status in both studies or gene unique to one study) and for genes where the reported XCI statuses were in conflict the following rules were applied: (1) a gene was considered 'escape' if it was called escape in one study and variable in the other; (2) 'variable escape' if classified as escape and inactive; and (3) 'inactive' if classified as inactive in one study and variable escape in the other. The final combined list of XCI statuses consisted of 631 X-chromosomal genes including 99 escape, 101 variable escape and 431 inactive genes.

**Analysis of sex-biased expression.** Differential expression analyses were conducted to identify genes that are expressed at significantly different levels between male and female samples using 29 GTEx v6p tissues with RNA-seq and genotype data available from more than 70 individuals after excluding samples flagged in QC and sex-specific, outlier (that is, breast tissue) and highly correlated tissues<sup>14</sup>. Only autosomal and X-chromosomal protein-coding or long non-coding RNA genes in Gencode v.19 were included, and all lowly expressed genes were removed (Extended Data Table 1 and Supplementary Methods).

Differential expression analysis between male and female samples was conducted following the voom-limma pipeline<sup>41–43</sup> available as an R package through Bioconductor (<https://bioconductor.org/packages/release/bioc/html/limma.html>) using the gene-level read counts as input. The analyses were adjusted for age, three principal components inferred from genotype data using EIGENSTRAT<sup>44</sup>, sample ischaemic time, surrogate variables<sup>45,46</sup> built using the sva R package<sup>47</sup>, and the cause of death classified into five categories based on the four-point Hardy scale (Supplementary Methods).

To control the FDR, the qvalue R package was used to obtain *q* values applying the adjustment separately for the differential expression results from each tissue. The null hypothesis was rejected for tests with *q* values below 0.01.

**XY homologue analysis.** A list of Y-chromosomal genes with functional counterparts in the X chromosome, that is, X–Y gene pairs, was obtained from ref. 30, which lists 19 ancestral Y chromosome genes that have been retained in the human Y chromosome. After excluding two of the genes (*MXRA5Y* and *OFDIY*), which were annotated as pseudogenes in ref 30, and a further four genes (*SRY*, *RBMV*, *TSPY* and *HSFY*) that according to ref. 30 have clearly diverged in function from their X-chromosomal homologues, the remaining 13 Y-chromosomal genes were matched with their X-chromosome counterparts using gene-pair annotations given in ref. 30 or by searching for known paralogues of the Y-chromosomal genes. To test for completeness of dosage compensation of the X–Y homologue genes, the sex-bias analysis in GTEx data was repeated replacing the expression of the X-chromosomal counterpart with the combined expression of the X and Y homologues.

**Chromatin state analysis.** To study the relationship between chromatin states and XCI, we used chromatin state calls from the Roadmap Epigenomics Consortium<sup>31</sup>. Specifically, we used the chromatin state annotations from the core 15-state model, publicly available at [http://egg2.wustl.edu/roadmap/web\\_portal/chr\\_state\\_learning.html#core\\_15state](http://egg2.wustl.edu/roadmap/web_portal/chr_state_learning.html#core_15state). We followed our previously published method<sup>48</sup> to calculate the covariate-corrected percentage of each gene body assigned to each chromatin state. After pre-processing, we filtered down to the 399 inactive and 86 escape genes on the X chromosome and down to 38 female epigenomes.

To compare the chromatin state profiles of the escape and inactive genes in female samples, we used the one-sided Wilcoxon rank-sum test. Specifically, for each chromatin state, we averaged the chromatin state coverage across the 38 female samples for each gene, and compared that average chromatin state coverage for all 86 escape genes to the average chromatin state coverage for all 399 inactive genes. We performed both one-sided tests, to test for enrichment in escape genes, as well as for enrichment in inactive genes.

Next, we performed simulations to account for possible chromatin state biases, such as the fact that the escape and inactive genes are all from the X chromosome. Specifically, we generated 10,000 randomized simulations where we randomly shuffled the escape or inactive labels on the combined set of 485 genes, while retaining the sizes of each gene set. For each of these simulated escape and inactive gene sets, we calculated both one-sided Wilcoxon rank-sum test *P* values as described above, and then, we calculated a permutation *P* value for the real gene sets based on these 10,000 random simulations (Supplementary Methods). Finally, we used Bonferroni multiple hypothesis corrections for our significance thresholds to correct for our 30 tests, one for each of 15 chromatin states, and both possible test directions.

**Allele-specific expression.** For ASE analysis the allele counts for biallelic heterozygous variants were retrieved from RNA-seq data using GATK ASEReadCounter (v.3.6)<sup>39</sup>. Heterozygous variants that passed VQSR filtering were first extracted for each sample from WES or WGS VCFs using GATK SelectVariants. The analysis was restricted to biallelic SNPs owing to known issues in mapping bias in RNA-seq against indels<sup>39</sup>. Sample-specific VCFs and RNA-seq BAMs were inputted to ASEReadCounter requiring minimum base quality of 13 in the RNA-seq data (scRNA-seq samples, GTEx-UPIC) or requiring coverage in the RNA-seq data of each variant to be at least 10 reads, with a minimum base quality of 10 and counting only reads with unique mapping quality (MQ = 60) (clinical muscle samples).

For downstream processing of the scRNA-seq and GTEx-UPIC ASE data, we applied further filters to the data to focus on exonic variation only and to conservatively remove potentially spurious sites (Supplementary Methods), for example, sites with non-unique mappability were removed, and furthermore, after an initial analysis of the ASE data, we subjected 22 of the X-chromosomal ASE sites to manual investigation. For GTEx-UPIC the X-chromosomal ASE data was limited to only one site per gene in case of multiple ASE sites, by selecting the site with coverage >7 reads in the largest number of tissues, to have equal representation of each gene for downstream analyses.

**Assessing ASE across tissues.** For the GTEx-UPIC individual, for whom ASE data from up to 16 tissues per each ASE site was available, we applied the two-sided hierarchical grouped tissue model (GTM\*) implemented in MAMBA v.1.0.0 (refs 49, 50) to ASE data. The Gibbs sampler was run for 200 iterations with a burn-in of 50 iterations.

GTM\* is a Bayesian hierarchical model that borrows information across tissues and across variants, and provides parameter estimates that are useful for interpreting global properties of variants. It classifies the sites into ASE states according to their tissue-wide ASE profiles and provides an estimate of the proportion of variants in each of the five different ASE states (strong ASE across all tissues (SNGASE), moderate ASE across all tissues (MODASE), no ASE across all tissues (NOASE) and heterogeneous ASE across tissues (HET1 and HET0)).

To summarize the GTM\* output in the context of XCI, SNGASE was considered to reflect full XCI, MODASE and NOASE were taken together to represent partial XCI with similar effects across tissues, and HET1 and HET0 were considered to reflect partial yet heterogeneous patterns of XCI across tissues. To combine estimates from two ASE states, we summed the estimated proportions in each class and subsequently calculated the 95% confidence intervals for each remaining ASE state using Jeffreys' prior.

**Determining XCI status in GTEx-UPIC.** In addition to the ASE states provided by the above MAMBA analysis, genic XCI status was assessed by comparing the allelic ratios at each X-chromosomal ASE site in each tissue individually. For each ASE site, the alleles were first mapped to Xa and Xi; the allele with lower combined relative expression across tissues was assumed to be the Xi allele. As an exception, at *XIST* the higher expressing allele was assumed to be the Xi allele. The significance of Xi expression at each ASE observation was tested using a one-sided binomial test, where the hypothesized probability of success was set at 0.025, that is, the

fraction of Xi expression from total expression was expected to be significantly greater than 0.025. To account for multiple testing, a FDR correction was applied, using the qvalue R package, to the *P* values from the binomial test for each of the 16 tissues separately. Observations with *q* values <0.01 were considered significant, that is, indicative of incomplete XCI at the given ASE site and tissue.

**Biallelic expression in single cells.** Biallelic expression in individual cells in the X chromosome was assessed only at ASE sites covered by the minimum of eight reads. A site was considered biallelically expressed when (1) allelic expression >0.05 and (2) the one-sided binomial test indicated allelic expression to be at least nominally significantly greater than 0.025. Only genes with at least two observations of biallelic expression across all cells within a sample were counted as biallelic.

**Phasing scRNA-seq data.** We assigned each cell to either of two cell populations distinguished by the parental X-chromosome designated for inactivation using genotype phasing. For the YRI samples, where parental genotype data was available, the assignment to the two parental cell populations was unambiguous for all cells where X-chromosomal sites outside PAR1 or frequently biallelic sites were expressed. For 24A, no parental genotype data were available, and we therefore used the correlation structure of the expressed X-chromosomal alleles across the 948 cells to infer the two parental haplotypes using the fact that in individual cells the expressed alleles at the chrX sites subject to full inactivation (that is, the majority chrX ASE sites), are from the X chromosome active in each cell (Supplementary Methods). In other words, while monoallelic expression in scRNA-seq in the autosomes is largely stochastic in origin, in the X chromosome the pattern of monoallelic expression is consistent across cells with the same parental X chromosome active<sup>22</sup>, unless a gene is expressed also from the inactive X. As such, for the phase inference calculations, we excluded all PAR1 sites and all additional sites that were frequently biallelic, to minimize the contribution of escape genes to the phase estimation. After assigning each informative cell to either of the parental cell populations, the reference and alternate allele reads for each ASE site were mapped to active and inactive allele reads within each sample using the actual or inferred parental haplotypes. The data were first combined per variant by taking the sum of active and inactive counts separately across cells, and further similarly combined per gene, if multiple SNPs per gene were available. For 24A the allele expressed at *XIST* was assumed the Xi allele, in line with the exclusive Xi expression in the Yoruba samples confirmed using the information on parental haplotypes.

**Determining XCI status from scRNA-seq ASE.** Before calling XCI status using the Xa and Xi read counts from the phased data aggregated across cells, we excluded all sites without fewer than five cells contributing ASE data at each gene and also all sites with coverage lower than eight reads across cells within each sample. To determine whether the observed Xi expression is significantly different from zero, and therefore indicative of incomplete XCI at the site or gene, we required the Xi to total expression ratio to be significantly (*q* value <0.01) greater than the hypothesized upper bound for error, 0.025. This threshold was determined using the proportion of miscalled alleles at *XIST* ASE sites (by definition, *XIST* should express only alleles from the inactive chrX) in the two YRI samples, which presented with fully skewed XCI, that is, the same active X chromosome across all assessed cells. The median proportion of miscalled *XIST* alleles was 0, yet one site in one of the samples showed up to 2.5% of other allele calls, and therefore this was chosen as the error margin. FDR correction, conducted using the qvalue R package, was applied to each sample individually. Genes where at least one of the samples showed significant Xi expression were considered partially inactivated, while the remaining were classified as subject to full XCI. Allelic dropout, which is extensive in scRNA-seq<sup>19,28</sup>, can lead to biases in allelic ratios in individual cells, that is, in our case resulting in false negatives where true escape genes are classified as inactivated, the used approach is based on using aggregate data across several cells and therefore the XCI status estimates are robust to such errors.

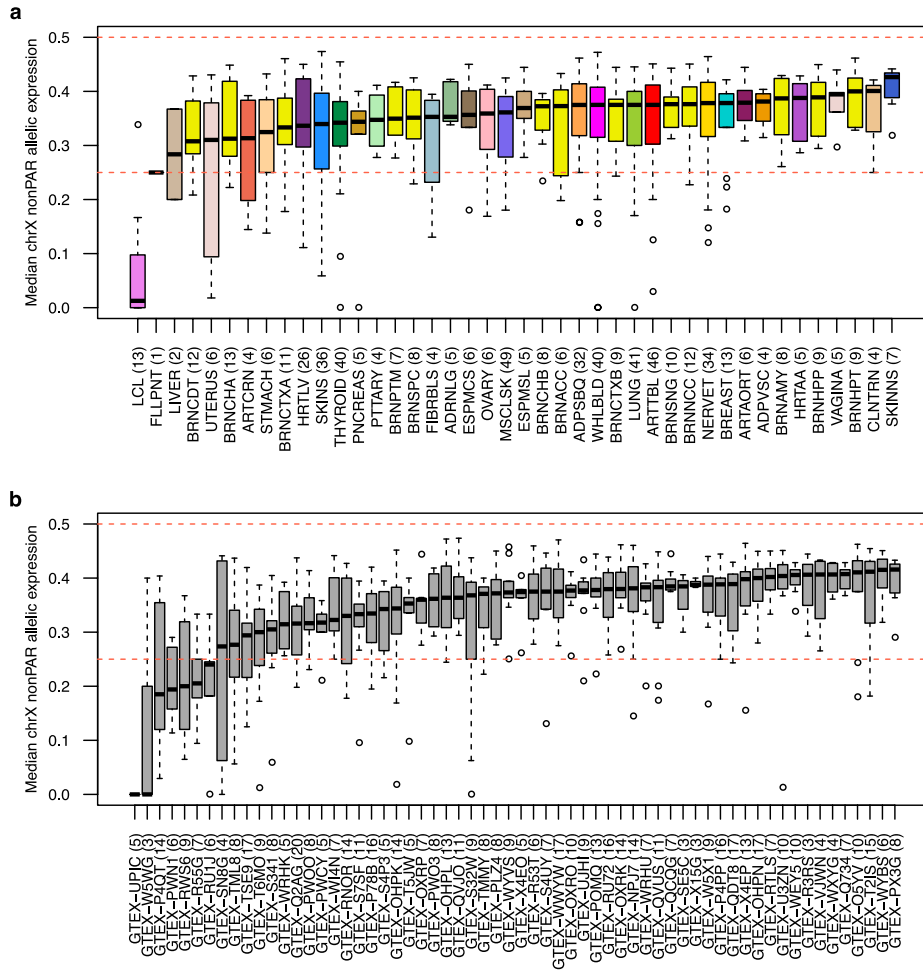
**ChrX and chrY expression in PAR1.** Using the parental origin of each allele reference and alternate allele read counts at PAR1 ASE sites were assigned to X and Y chromosomes (that is, maternally and paternally inherited alleles, respectively). For each sample, the PAR1 ASE data was summarized by gene by taking the sum of X and Y chromosome reads across all informative ASE sites within each gene.

Significance of deviation from equal expression was assessed using a two-sided binomial test.

**Manual curation of heterozygous variants from ASE analyses.** Twenty-two heterozygous variants assessed in chrX ASE analysis were subjected to manual curation because of results in the XCI analysis that were in conflict with previous assignment of the underlying gene to be subject to full XCI. For each sample, BWA-aligned germline BAM files were viewed in IGV using either WGS or WES data. The presence of a number of characteristics called into question the confidence of the variant read alignments and thus the variant itself (Supplementary Methods). Allele balance that deviated significantly from 50:50 was considered suspect and often coincided with the existence of homology between the reference sequence in the region surrounding the variant and another area of the genome, as ascertained using the UCSC browser self-chain track and/or BLAT alignment of variant reads from within IGV. Other sequence-based annotations added to the VCF by HaplotypeCaller were also evaluated in the interests of examining other signatures of ambiguous mapping. The phasing of nearby variants was also considered. If phased variants occurred in the DNA sequencing data that were not assessed in the ASE analysis, those variants were considered suspect.

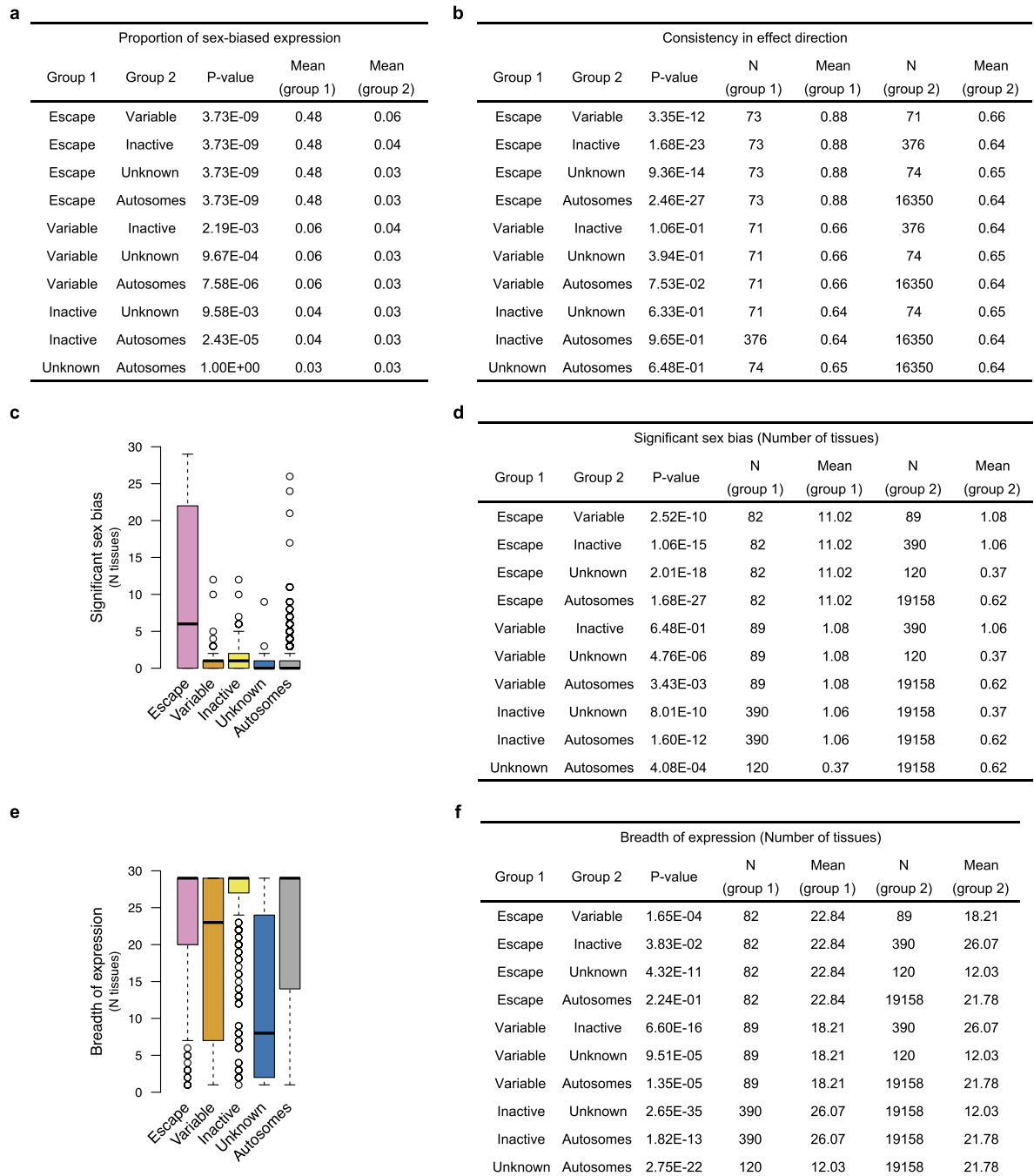
**Data availability.** Gene expression and genotype data from the GTEx v6p release are available in dbGaP (study accession phs000424.v6.p1; [http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study\\_id=phs000424.v6.p1](http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000424.v6.p1)). Raw RNA-seq data for 24A is available through dbGaP accession number phs001294.v1.p1 (<https://www.ncbi.nlm.nih.gov/bioproject/?term=phs001294.v1.p1>). The authors declare that all data supporting the findings of this study are available within the paper and its Supplementary Information. Source Data for Figs 2–4 are provided with the paper.

32. Xia, Z. *et al.* A 17q12 allele is associated with altered NK cell subsets and function. *J. Immunol.* **188**, 3315–3322 (2012).
33. The 1000 Genomes Project Consortium. A map of human genome variation from population-scale sequencing. *Nature* **467**, 1061–1073 (2010).
34. Lee, M. N. *et al.* Common genetic variants modulate pathogen-sensing responses in human dendritic cells. *Science* **343**, 1246980 (2014).
35. Picelli, S. *et al.* Smart-seq2 for sensitive full-length transcriptome profiling in single cells. *Nat. Methods* **10**, 1096–1098 (2013).
36. Trombetta, J. J. *et al.* Preparation of single-cell RNA-seq libraries for next generation sequencing. *Curr. Protoc. Mol. Biol.* **107**, 4.22.1–4.22.17 (2014).
37. DeLuca, D. S. *et al.* RNA-SeQC: RNA-seq metrics for quality control and process optimization. *Bioinformatics* **28**, 1530–1532 (2012).
38. Rozowsky, J. *et al.* AlleleSeq: analysis of allele-specific expression and binding in a network framework. *Mol. Syst. Biol.* **7**, 522 (2011).
39. Castel, S. E., Levy-Moonshine, A., Mohammadi, P., Banks, E. & Lappalainen, T. Tools and best practices for data processing in allelic expression analysis. *Genome Biol.* **16**, 195 (2015).
40. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
41. Law, C. W., Chen, Y., Shi, W. & Smyth, G. K. voom: precision weights unlock linear model analysis tools for RNA-seq read counts. *Genome Biol.* **15**, R29 (2014).
42. Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res.* **43**, e47 (2015).
43. Smyth, G. K. Linear models and empirical Bayes methods for assessing differential expression in microarray experiments. *Stat. Appl. Genet. Mol. Biol.* **3**, 1–25 (2004).
44. Price, A. L. *et al.* Principal components analysis corrects for stratification in genome-wide association studies. *Nat. Genet.* **38**, 904–909 (2006).
45. Leek, J. T. & Storey, J. D. Capturing heterogeneity in gene expression studies by surrogate variable analysis. *PLoS Genet.* **3**, 1724–1735 (2007).
46. Leek, J. T. & Storey, J. D. A general framework for multiple testing dependence. *Proc. Natl Acad. Sci. USA* **105**, 18718–18723 (2008).
47. Leek, J. T., Johnson, W. E., Parker, H. S., Jaffe, A. E. & Storey, J. D. The sva package for removing batch effects and other unwanted variation in high-throughput experiments. *Bioinformatics* **28**, 882–883 (2012).
48. Yen, A. & Kellis, M. Systematic chromatin state comparison of epigenomes associated with diverse properties including sex and tissue type. *Nat. Commun.* **6**, 7973 (2015).
49. Pirinen, M. *et al.* Assessing allele-specific expression across multiple tissues from RNA-seq read data. *Bioinformatics* **31**, 2497–2504 (2015).
50. Rivas, M. A. *et al.* Effect of predicted protein-truncating genetic variants on the human transcriptome. *Science* **348**, 666–669 (2015).



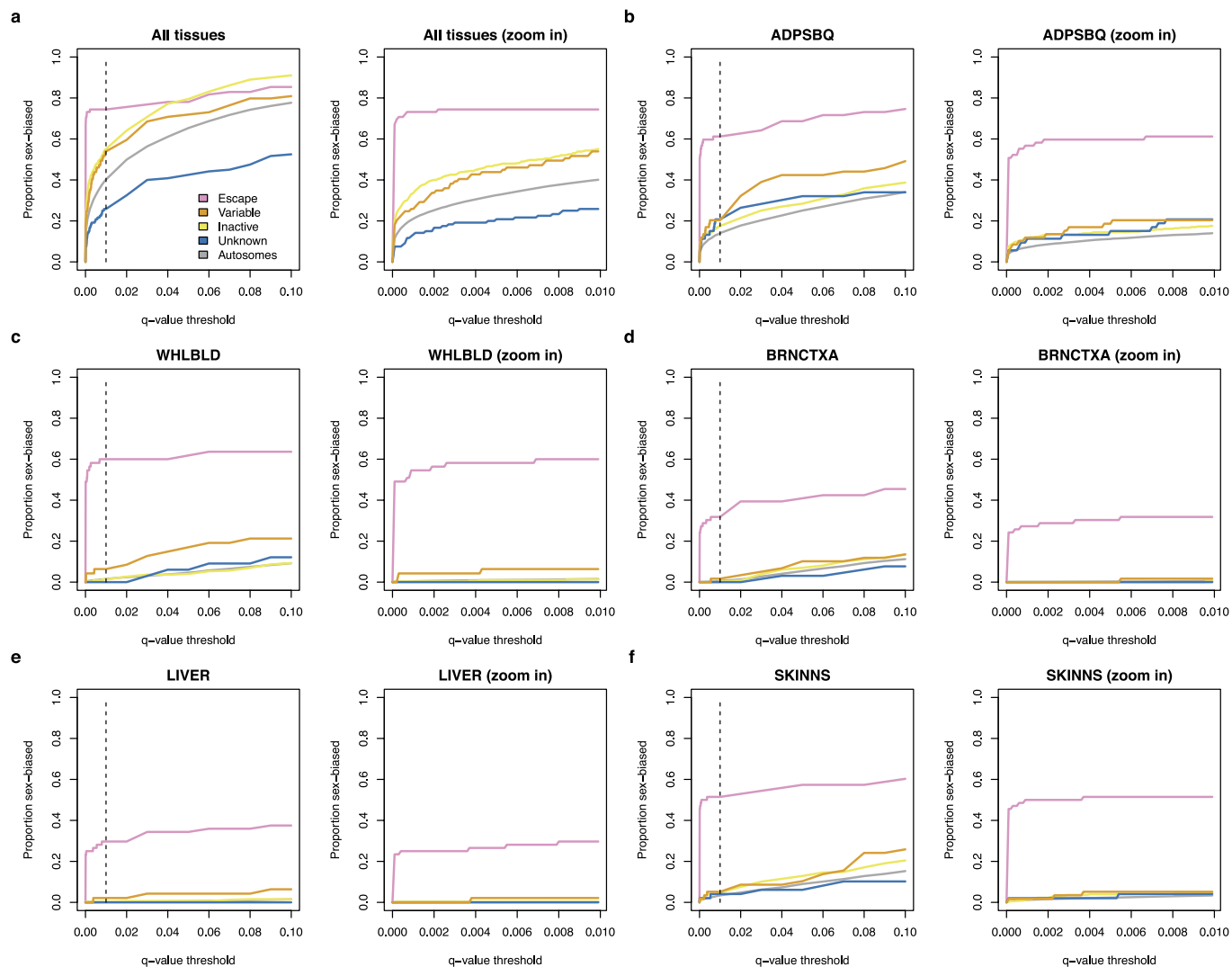
Extended Data Figure 1 | Assessment of skew in XCI in GTEx female samples (v3 analysis release). **a**, The estimated skew in XCI by tissue across individuals. **b**, The skew in XCI by individual across available tissue

samples. The number in brackets after the tissue or sample name indicates the number of individuals or tissues, respectively, contributing to each box plot. Details of the analysis can be found in the Supplementary Note.



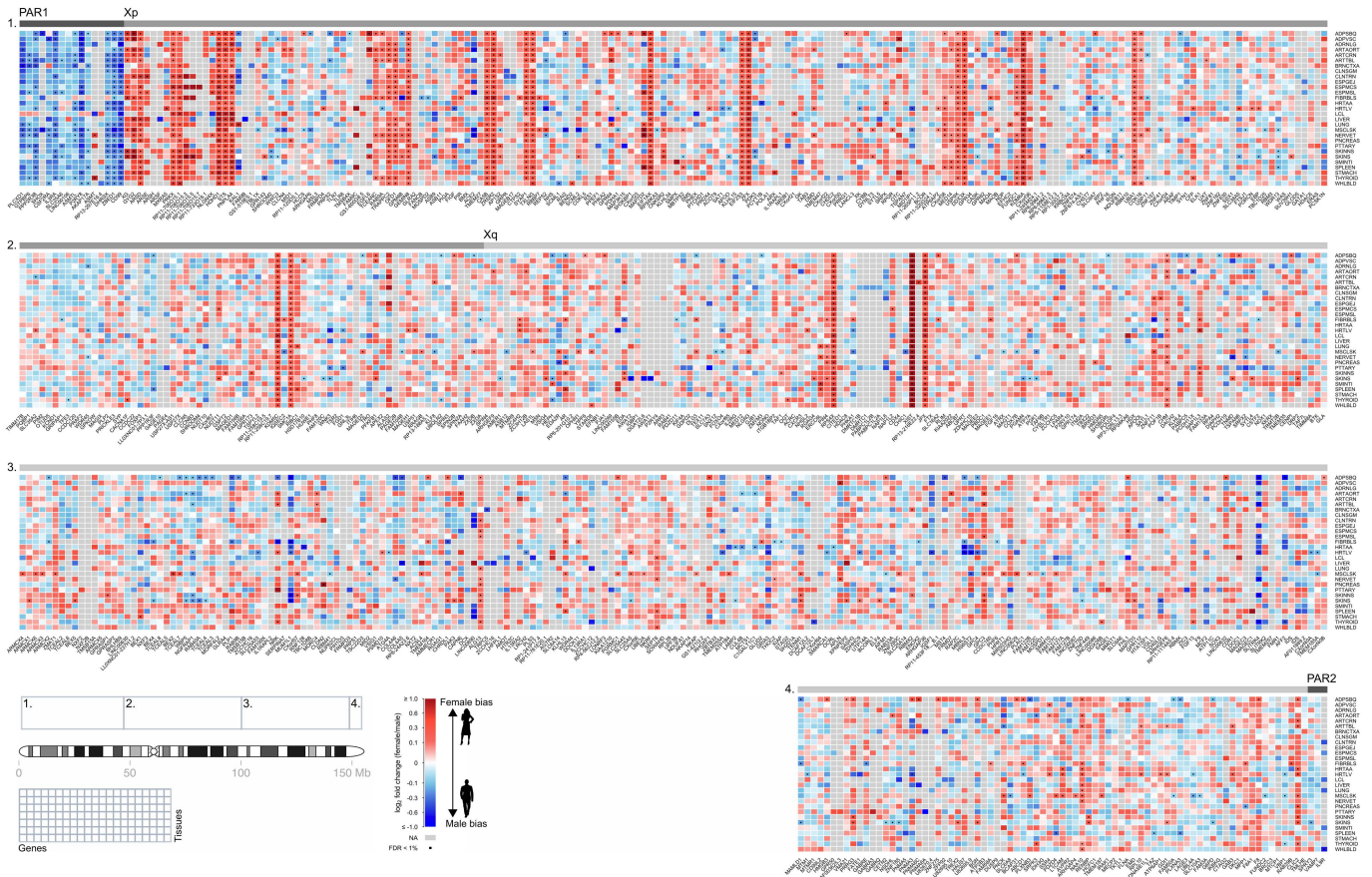
**Extended Data Figure 2 | Comparison of expression characteristics between reported genic XCI categories in the GTEx data. a,** The statistics for the comparison of the proportion of significantly biased (FDR <1%) genes by reported XCI status. Distributions are illustrated in Fig. 2b.  $n = 29$  for all comparisons. **b,** The statistics for the comparison of the consistency in effect sizes across tissues. Distributions are illustrated in Fig. 2c. Only genes expressed in at least five of the 29 tissues are included.

**c,** Number of tissues showing significant sex bias (FDR <1%) per gene by reported XCI status. **d,** Statistics for the comparison illustrated in **c**. **e,** Number of tissues in which genes are expressed by reported XCI status. **f,** Statistics for the comparison illustrated in **e**. All  $P$  values are from two-sided Wilcoxon rank-sum tests, except for **a**, where a paired, two-sided Wilcoxon rank-sum test was used. Only genes assessed for sex bias in at least one tissue are included, unless otherwise stated.



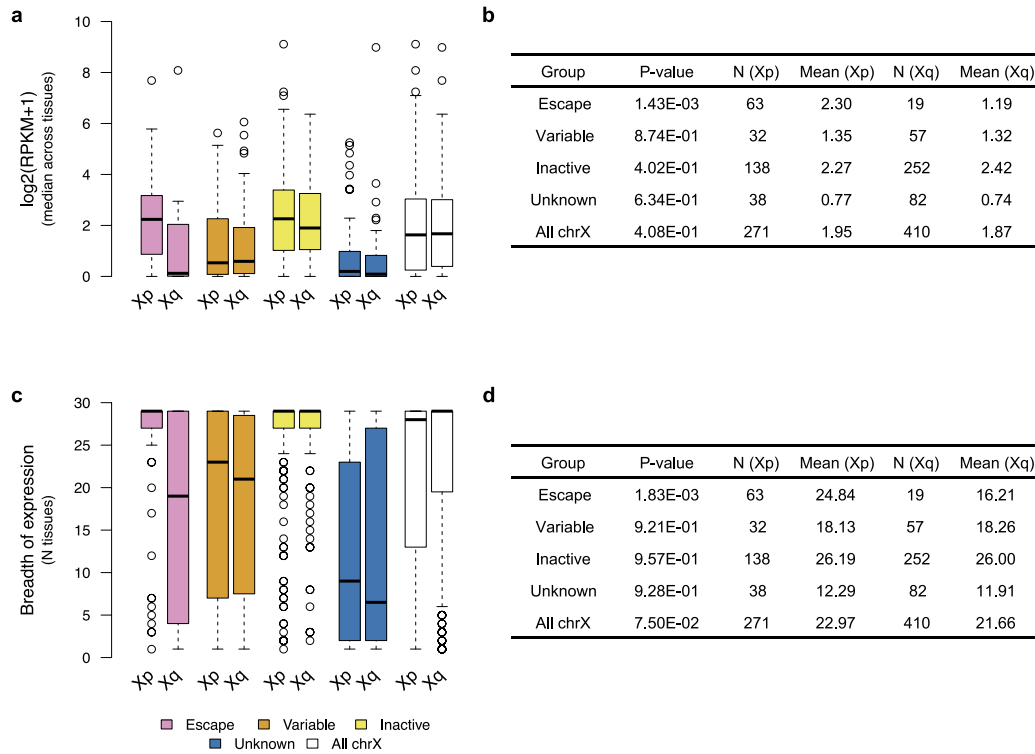
**Extended Data Figure 3 | Change in the proportion of discovered sex-biased genes by XCI category with varying  $q$  value cut-offs.** **a**, The proportion of sex-biased genes across tissues. Here a gene is classified as sex-biased if the  $q$  value for association falls below the given threshold in at least one tissue. **b–f**, Examples of the change in the proportion of

sex-biased expression in individual tissues. The dashed black line indicates the FDR < 1% cut-off applied in the analyses to determine sex-biased expression. ADPSBQ, adipose, subcutaneous; WHLBLD, whole blood; BRNCTXA, brain, cortex; SKINNS, skin not sun exposed (suprapubic).



**Extended Data Figure 4 | Heat map representation of male–female expression differences in all assessed X-chromosomal genes ( $n = 681$ ) across 29 GTEx tissues.** The colour scale displays the direction of sex bias, with red colour indicating higher female expression. Genes that were

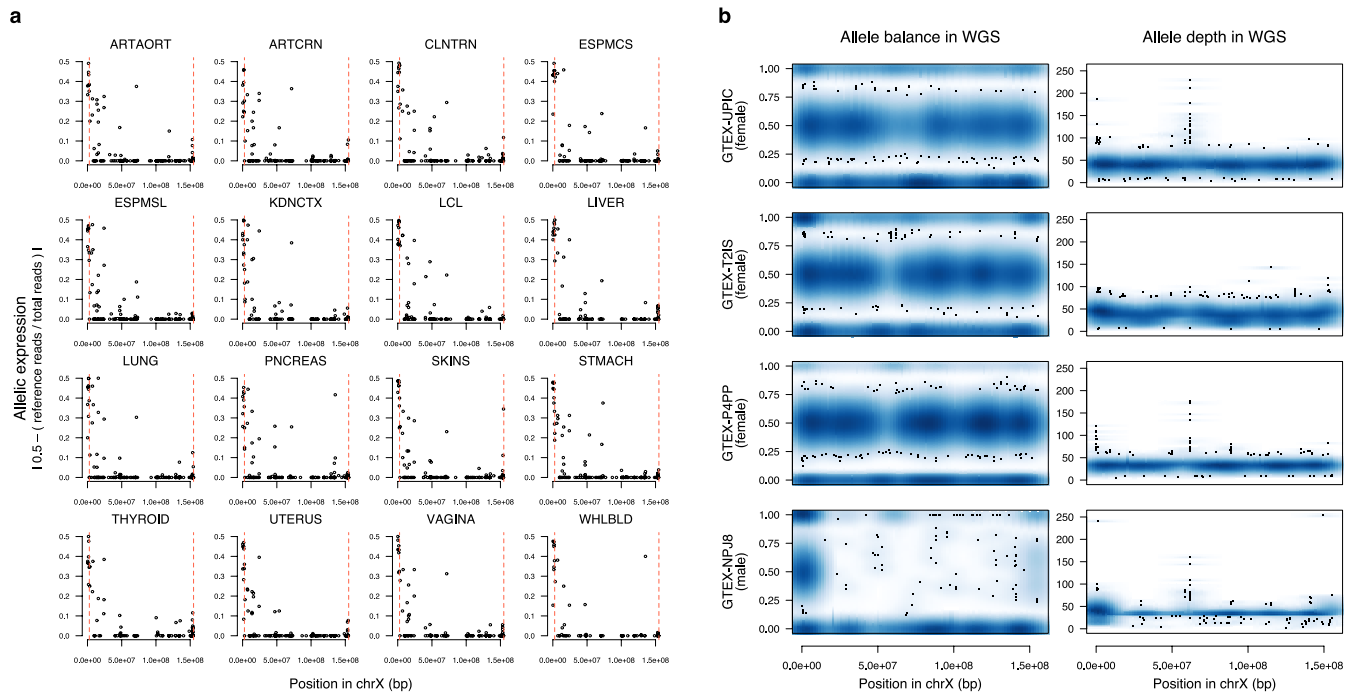
too weakly expressed to be assessed in a given tissue type in the sex bias analysis are coloured grey. Dots mark the observations where sex bias was significant at FDR < 1%.



**Extended Data Figure 5 | Comparison of expression characteristics between Xp and Xq, the evolutionary newer and older regions of chrX, respectively, by XCI status and for the whole chromosome.**

**a, b,** The level of median expression across GTEx tissues in  $\log_2$  RPKM units. **c, d,** The breadth of expression measured as the number of tissues

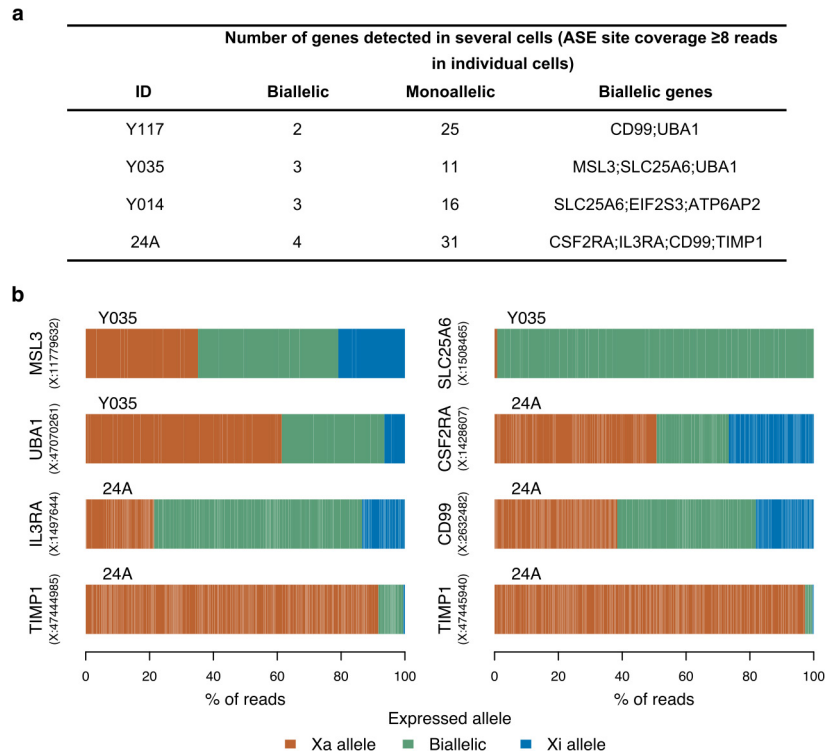
(max = 29) in which genes are expressed (median expression across samples  $>0.1$  RPKM and expressed in more than 10 individuals at  $>1$  counts per million). *P* values are calculated using the Wilcoxon rank-sum test. All genes expressed in at least one tissue are included in the comparisons.



**Extended Data Figure 6 | X-chromosomal RNA-seq and WGS data in the GTEX donor with fully skewed XCI (GTEX-UPIC). a,** Allelic expression in chrX in 16 RNA-sequenced tissue samples available from

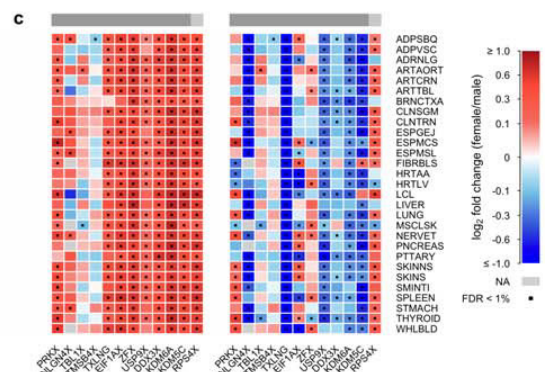
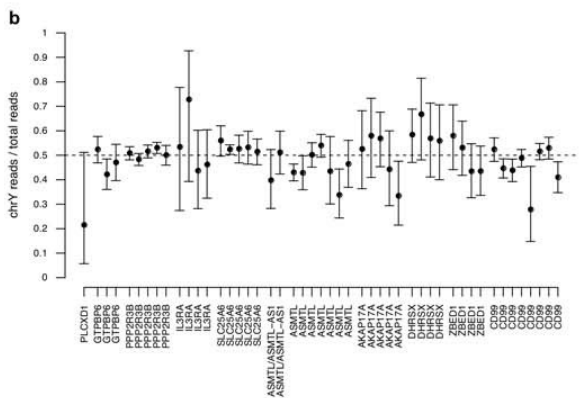
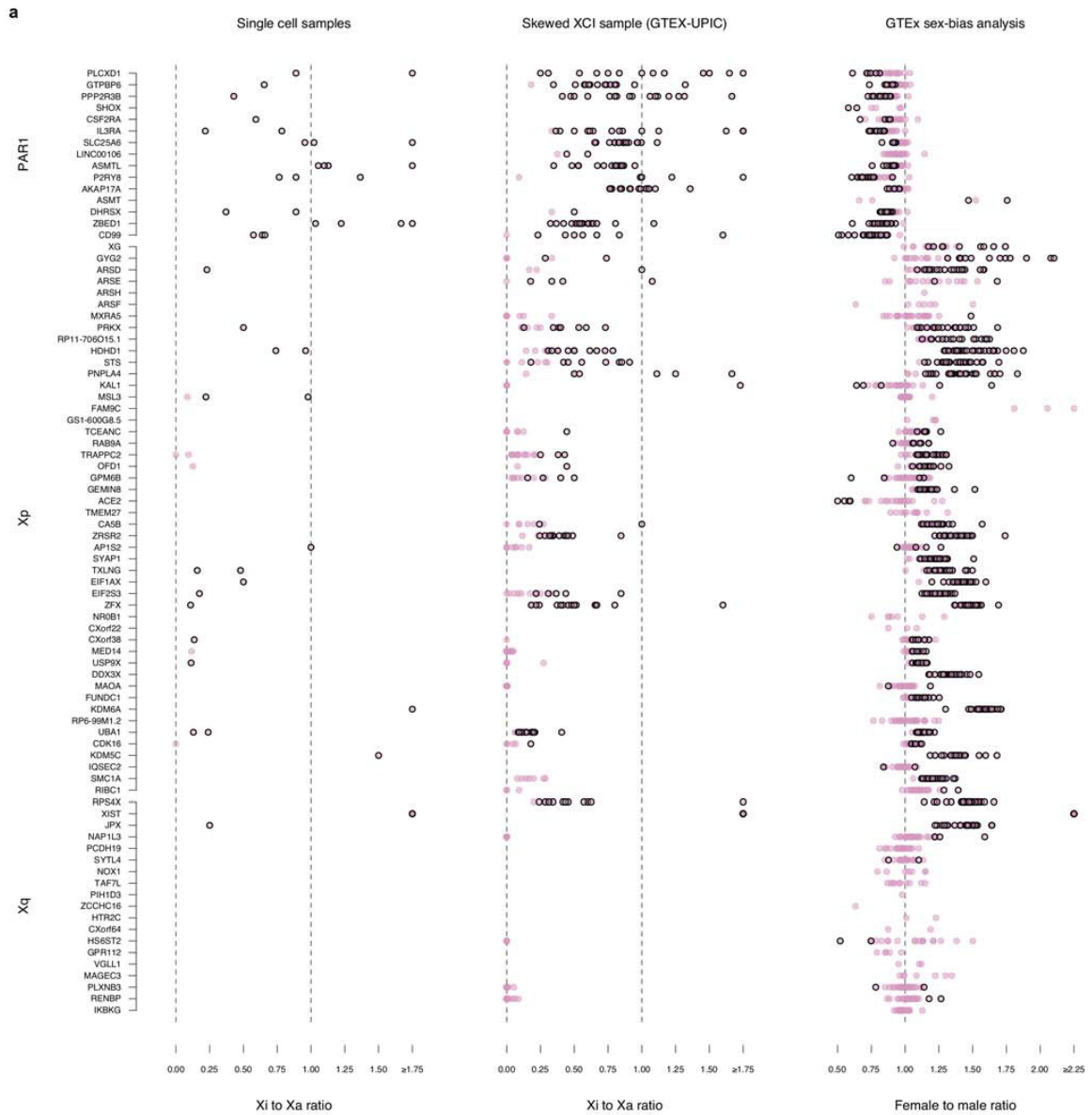
the donor. Dashed red lines indicate PAR1 and PAR2 boundaries. **b,** Allele balance and allele depth across chrX in WGS for GTEX-UPIC and two female and one male GTEX WGS samples that were randomly chosen.





**Extended Data Figure 7 | Expressed alleles at biallelically expressed ASE sites in scRNA-seq. a**, X-chromosomal genes repeatedly biallelic in scRNA-seq (see Methods for details). **b**, Illustration of the relative expression from the two alleles at all X-chromosomal ASE sites that were

repeatedly biallelically expressed across cells in either of the two scRNA-seq samples that showed random XCI (Y035 and 24A). Narrow white lines separate observations from individual cells.



Extended Data Figure 8 | See next page for caption.

**Extended Data Figure 8 | Assessment of the level of Xi expression at escape genes and in different regions of the X chromosome.** **a**, The ratio of Xi-to-Xa expression in the single-cell samples (left; each circle represents a sample), in the skewed XCI donor from GTEx (middle; each circle represents a tissue), and the female-to-male ratio in expression (right; each circle represents a tissue) at reported escape genes. Genes are ordered according to their location in the X chromosome with genes in the pseudoautosomal region residing in the top part. A dark border around a circle indicates that there was significant evidence for Xi expression greater than the baseline in the given sample or tissue (left and middle) or significant sex-bias in the given tissue (right). Given some outliers, for example, *XIST*, the Xi-to-Xa ratio is capped at 1.75 and female-to-male ratio at 2.25. **b**, The relative expression arising from the X and Y chromosome at PAR1 genes in skeletal muscle in eight males. The allelic expression at these genes was assigned to the two chromosomes using

parental genotypes available for these samples (see Methods for details). The dashed line at 0.5 indicates the point where expression from X and Y chromosomes is equal. The error bars give the 95% confidence intervals for the observed read ratio. **c**, Heat map representation of the change in pattern of sex-bias at 13 X–Y homologous gene pairs (see Methods for details) in nonPAR from only including the X-chromosomal expression (heat map on the left) to accounting for the Y-chromosomal expression (heat map on the right). The colour scale displays the direction of sex-bias with red colour indicating higher female expression. Genes that were too lowly expressed in the given tissue type to be assessed in the sex-bias analysis are coloured grey. Dots mark the observations where sex-bias was significant at  $FDR < 1\%$ . The grey bars on top of the heat maps indicate the location of the gene in the X chromosome: dark grey indicating Xp and lighter grey Xq.

Extended Data Table 1 | Tissues, individuals and genes in the GTEx sex-bias analysis

Abbreviation	Tissues		Individuals				Genes analyzed		
	Full name	All	Females	Males	Mean age	All	Autosomes	ChrX	
ADPSBQ	Adipose - Subcutaneous	297	186	111	52.15	15,273	14,735	538	
ADPVSC	Adipose - Visceral (Omentum)	184	117	67	51.97	15,301	14,765	536	
ADRNLG	Adrenal Gland	126	70	56	50.51	14,956	14,435	521	
ARTAORT	Artery - Aorta	197	126	71	51.11	14,675	14,137	538	
ARTCRN	Artery - Coronary	118	70	48	51.7	14,881	14,350	531	
ARTTBL	Artery - Tibial	284	183	101	50.26	14,501	13,981	520	
BRNCTXA	Brain - Cortex	92	66	26	57.67	15,339	14,791	548	
CLNSGM	Colon - Sigmoid	114	72	42	48.28	15,045	14,524	521	
CLNTRN	Colon - Transverse	255	159	96	50.93	15,732	15,181	551	
ESPG EJ	Esophagus - Gastroesophageal Junction	124	74	50	53.52	14,770	14,245	525	
ESPMCS	Esophagus - Mucosa	169	97	72	48.89	15,137	14,617	520	
ESPM SL	Esophagus - Muscularis	126	78	48	50.74	14,879	14,356	523	
FIBRBL S	Cells - Transformed fibroblasts	240	150	90	50.2	13,635	13,158	477	
HR TAA	Heart - Atrial Appendage	218	137	81	48.62	14,662	14,145	517	
HR TLV	Heart - Left Ventricle	159	105	54	53.64	14,075	13,586	489	
LCL	Cells - EBV-transformed lymphocytes	190	123	67	50.75	13,067	12,621	446	
LIVER	Liver	96	63	33	53.52	14,031	13,556	475	
LUNG	Lung	277	181	96	52.06	16,154	15,590	564	
MSCLSK	Muscle - Skeletal	361	228	133	51.85	13,623	13,153	470	
NERVET	Nerve - Tibial	256	163	93	51.65	15,563	15,020	543	
PNCREAS	Pancreas	149	87	62	50.09	14,355	13,861	494	
PTTARY	Pituitary	86	64	22	56.37	16,068	15,489	579	
SKINNS	Skin - Not Sun Exposed (Suprapubic)	195	128	67	53.06	15,601	15,069	532	
SKINS	Skin - Sun Exposed (Lower leg)	300	188	112	52.22	15,746	15,211	535	
SMINTI	Small Intestine - Terminal Ileum	77	43	34	47.62	15,594	15,046	548	
SPLEEN	Spleen	89	50	39	48.26	14,993	14,469	524	
STMACH	Stomach	169	97	72	48.2	15,604	15,057	547	
THYROID	Thyroid	278	179	99	52.14	15,974	15,417	557	
WHLBLD	Whole Blood	338	213	125	51.64	13,187	12,751	436	
Total		449	290	159	52.27	19,839	19,158	681	

Extended Data Table 2 | scRNA-seq samples

ID	24A	Y117	Y035	Y014
Ancestry	China, Asia	Yoruba / Nigeria, Africa	Yoruba / Nigeria, Africa	Yoruba / Nigeria, Africa
Design	Singleton	Trio	Trio	Trio
Genotype data	WES	WGS	WES	WES
Number of cells	742	96	48	48
Cell type	Dendritic cells	LCL	LCL	LCL
Sequenced read pairs (mean (range))	1,187,000 (335-7,403,000)	2,547,000 (38,190-5,126,000)	2,571,000 (46,940-5,038,000)	2,436,000 (69,130-5,457,000)
Aligned read pairs* (mean (range))	808,600 (197-5,727,000)	1,471,000 (14,910-3,309,000)	1,459,000 (16,400-2,893,000)	1,391,000 (14,920-3,067,000)
Alignment rate (mean (range))	0.667 (0.271-0.799)	0.545 (0.251-0.645)	0.551 (0.266-0.615)	0.526 (0.175-0.606)
Skew in XCI (% maternal active : % paternal active)	54:46 (373 cells where one parental chromosome active, 315 cells where the other parental chromosome active, 54 cells uninformative for X- chromosomal phasing)	100:0 (90 cells where maternal X chromosome active, 6 cells uninformative for X- chromosomal phasing)	79:21 (37 cells where maternal X chromosome active, 8 cells where paternal X chromosome active, 2 cells uninformative for X- chromosomal phasing)	100:0 (43 cells where maternal X chromosome active, 2 cells uninformative for X- chromosomal phasing)
Notes	Due to the unavailability of parental genotype information, the parental origin of the inferred X- chromosomal haplotypes is unknown			

\*Uniquely aligned, properly paired, quality-control passed reads.

## Life Sciences Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form is intended for publication with all accepted life science papers and provides structure for consistency and transparency in reporting. Every life science submission will use this form; some list items might not apply to an individual manuscript, but all fields must be completed for clarity.

For further information on the points included in this form, see [Reporting Life Sciences Research](#). For further information on Nature Research policies, including our [data availability policy](#), see [Authors & Referees](#) and the [Editorial Policy Checklist](#).

## ▶ Experimental design

## 1. Sample size

Describe how sample size was determined.

No statistical methods were used to predetermine sample size.

## 2. Data exclusions

Describe any data exclusions.

The GTEx samples were curated according to pre-established QC criteria as detailed in the accompanying manuscript by Aguet et al. scRNA-seq data was limited to those cells that were informative for chromosome X allelic expression.

## 3. Replication

Describe whether the experimental findings were reliably reproduced.

The analyses conducted were exploratory and the results were not replicated in independent data sets. However each analysis included multiple data points (individuals and/or tissues) thus providing further support for the conclusions drawn.

## 4. Randomization

Describe how samples/organisms/participants were allocated into experimental groups.

The experiments were not randomized. The study included no allocation into experimental groups.

## 5. Blinding

Describe whether the investigators were blinded to group allocation during data collection and/or analysis.

The investigators were not blinded to allocation during experiments and outcome assessment. The study included no allocation into experimental groups.

Note: all studies involving animals and/or human research participants must disclose whether blinding and randomization were used.

## 6. Statistical parameters

For all figures and tables that use statistical methods, confirm that the following items are present in relevant figure legends (or in the Methods section if additional space is needed).

- n/a | Confirmed
- The exact sample size ( $n$ ) for each experimental group/condition, given as a discrete number and unit of measurement (animals, litters, cultures, etc.)
  - A description of how samples were collected, noting whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
  - A statement indicating how many times each experiment was replicated
  - The statistical test(s) used and whether they are one- or two-sided (note: only common tests should be described solely by name; more complex techniques should be described in the Methods section)
  - A description of any assumptions or corrections, such as an adjustment for multiple comparisons
  - The test results (e.g.  $P$  values) given as exact values whenever possible and with confidence intervals noted
  - A clear description of statistics including central tendency (e.g. median, mean) and variation (e.g. standard deviation, interquartile range)
  - Clearly defined error bars

See the web collection on [statistics for biologists](#) for further resources and guidance.

## ► Software

Policy information about [availability of computer code](#)

### 7. Software

Describe the software used to analyze the data in this study.

RNA-seq alignment: Tophat version v1.4.1, STAR versions 2.4.2a, 2.4.1a or 2.3.0e; RNA-seq QC and quantification: RNA-SeQC; Allelic expression and variant calling: GATK version 3.1 or 3.4. Data processing: R version 3.4.0.

For manuscripts utilizing custom algorithms or software that are central to the paper but not yet described in the published literature, software must be made available to editors and reviewers upon request. We strongly encourage code deposition in a community repository (e.g. GitHub). *Nature Methods* [guidance for providing algorithms and software for publication](#) provides further information on this topic.

## ► Materials and reagents

Policy information about [availability of materials](#)

### 8. Materials availability

Indicate whether there are restrictions on availability of unique materials or if these materials are only available for distribution by a for-profit company.

All unique materials are readily available from the authors or from commercial sources as described in the Online Methods.

### 9. Antibodies

Describe the antibodies used and how they were validated for use in the system under study (i.e. assay and species).

The antibody panels used to enrich for all known blood DC population for single cell sorting and single cell RNA-sequencing (scRNA-seq) are described in Villani et al (Science 2017). All antibodies are commercially available as described in Supplementary Table 14.

### 10. Eukaryotic cell lines

a. State the source of each eukaryotic cell line used.

YRI LCLs were obtained from NHGRI Sample Repository for Human Genetic Research (Coriell Institute for Medical Research).

b. Describe the method of cell line authentication used.

None of the cell lines used were authenticated.

c. Report whether the cell lines were tested for mycoplasma contamination.

Coriell Biorepositories declares that their lymphoblastoid cell lines are free of bacterial, fungal or mycoplasma contamination. No other tests were run to test for mycoplasma contamination.

d. If any of the cell lines used are listed in the database of commonly misidentified cell lines maintained by [ICLAC](#), provide a scientific rationale for their use.

No commonly misidentified cell lines were used.

## ► Animals and human research participants

Policy information about [studies involving animals](#); when reporting animal research, follow the [ARRIVE guidelines](#)

### 11. Description of research animals

Provide details on animals and/or animal-derived materials used in the study.

No animals were used

Policy information about [studies involving human research participants](#)

### 12. Description of human research participants

Describe the covariate-relevant population characteristics of the human research participants.

24A: Female, Asian ancestry, 25 yo, healthy  
Y117, Y035 and Y014: Female, African ancestry, age and health status unknown  
GTEx-UPIC: Female, European ancestry, 21 yo, cause of death asphyxiation  
See Extended Data Table for information on other GTEx donors