# Population variability in X-chromosome inactivation across 9 mammalian species

Jonathan M. Werner[1], John Hover[1], Jesse Gillis[1,2]

[1]Stanley Institute for Cognitive Genomics, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11724, USA

[2]Physiology Department and Donnelly Centre for Cellular and Biomolecular Research, University of Toronto, Toronto, ON, Canada

## Abstract:

One of the two X chromosomes in female mammals is epigenetically silenced in embryonic stem cells by X chromosome inactivation (XCI). This creates a mosaic of cells expressing either the maternal or the paternal X allele. The XCI ratio, the proportion of inactivated parental alleles, varies widely among individuals, representing the largest instance of epigenetic variability within mammalian populations. While various contributing factors to XCI variability are recognized, namely stochastic and/or genetic effects, their relative contributions are poorly understood. This is due in part to limited cross-species analysis, making it difficult to distinguish between generalizable or species-specific mechanisms for XCI ratio variability. To address this gap, we measured XCI ratios in nine mammalian species (9,143 individual samples), ranging from rodents to primates, and compared the strength of stochastic models or genetic factors for explaining XCI variability. Our results demonstrate the embryonic stochasticity of XCI is a general explanatory model for population XCI variability in mammals, while genetic factors play a minor role.

1

## Introduction

Every female mammalian embryo undergoes X-chromosome inactivation (XCI) as an essential step for successful development[1–3]. XCI evolved to balance the gene dosage between females with two X-chromosomes and males with one X-chromosome[4]. While the exact timing can vary across species[5], XCI usually occurs during preimplantation embryonic development[6]. During this process, one of the two X-alleles in each female cell is independently, randomly, and permanently chosen for transcriptional silencing to match the single X-allele in male embryos[1,7–9]. The choice of silenced X-allele is inherited through cell divisions, propagating the random choice of allelic inactivation down each cell's subsequent lineage. This produces whole-body mosaicism for allelic X-chromosome expression in each adult mammalian female, originating from very early embryonic development[10].

In humans, both X-alleles are equally likely to be inactivated, but XCI ratios vary widely among adult females, from balanced to highly skewed[11,12]. XCI ratios affect the phenotypes of X-linked diseases, as they can either protect or expose individuals to disease variants[10,13]. The factors that influence XCI variability are mostly studied in mice and humans, and include stochasticity[12] and genetics[14–16], but their relative roles are controversial[17]. Cross-species comparisons of XCI variability stand to reveal general or species-specific mechanisms of XCI. For instance, genetic determinants of XCI are well-established in lab mice[18–20], but not in humans[17,21,22], where they are harder to identify and measure. Exploring XCI variability in other mammals presents the opportunity to test models of stochasticity or genetics in the context of evolution.

Considering first a stochastic model for XCI variability, each cell within an embryo at the time of XCI independently selects an X-allele to inactivate, resulting in ratios of allelic-inactivation across embryos varying purely by chance (Fig. 1A). Closely following Mary Lyon's discovery of XCI in 1961[1], it was recognized that the inherent embryonic stochasticity and permanence of XCI is the simplest explanation for the observed variability in XCI among adults and positions this adult variability as a window into embryonic events[23–27]. For example, flipping 10 coins is more likely to result in 8 heads than flipping 100 coins is likely to result in 80 heads, meaning that the variability in heads-to-tails ratios depends on the number of coins flipped. Similarly, the variability of XCI ratios in a population of female mammalian embryos is determined by the number of cells at the time of XCI (Fig. 1A). Since each cell inherits its allelic-inactivation from its ancestor, measuring XCI variability in adults can approximate embryonic XCI variability and help infer cell counts at the time of XCI or other early lineage decisions[25,28] (Fig. 1D). Stochastic models have been used to estimate cell counts during embryonic events in human and mice populations for decades[20,23,25,27–29] – but their applicability has not been tested in other mammalian species.

In addition to stochasticity, genetic effects can influence the choice of allelic inactivation and contribute to population variability in XCI ratios. Allelic inactivation during XCI is mediated by the cis-acting long non-coding RNA XIST[30], which silences its corresponding X-allele through epigenetic modifications[31,32]. Heterozygous variants

93 affecting XIST expression can bias allelic inactivation[15]. For example, inbred mice show
94 preferential inactivation of specific X-alleles depending on the parental strains and their
95 corresponding X-chromosome controlling element (XCE) allele[18,20,33]. In humans,
96 genetic influence on XCI is mostly observed in small family studies or disease cases,
97 with no strong evidence for the broad allelic effects seen in mice[21,22]. Another genetic
98 influence on XCI is allelic selection, where natural or disease-causing variants favor
99 certain X-alleles[14,16,34–38]. However, evidence for allelic selection through natural
100 variation remains elusive in human populations. Thus, the relative contributions of
101 stochasticity and genetics to population XCI variability in mammals remain unclear with
102 currently limited data from mouse and human studies.
103
104     In this study, we assess population scale XCI variability and its determinants
105 across nine mammalian species.  We source female annotated bulk RNA-sequencing
106 samples from the Sequencing Read Archive (SRA), resulting in a total of 19,180 initial
107 samples (Fig. 1C), including human samples from the GTEx[39] dataset. Our approach
108 leverages natural genetic variation to sample X-linked heterozygosity and eliminates the
109 requirement for costly phased or strain specific genetic information to assess XCI ratios
110 across diverse mammals at population scale. We start by establishing the population-
111 level XCI ratio distributions for all nine mammalian species and use models of
112 embryonic stochasticity to predict the number of cells fated for embryonic lineages (Fig.
113 1D, Fig. 2). We then investigate how broad genetic diversity, as indicated by measures
114 of inbreeding (Fig. 3), as well as specific individual variants (Fig. 4), may impact
115 population XCI variability. Overall, our analyses explore how both models of
116 stochasticity and genetic factors can explain population XCI variability across diverse
117 mammalian species.
118
119 **Results**
120
121 **Reference aligned RNA-sequencing data enables scalable modeling of XCI ratios**
122
123     We use bulk RNA-sequencing (RNA-seq) data to measure the X-linked allelic
124 expression of a sampled tissue by computing allele-specific expression ratios of
125 heterozygous single nucleotide polymorphisms (SNPs). The parental proportion of X-
126 linked allelic reads are expected to follow a binomial distribution dependent on the
127 number of sampled reads and the XCI ratio of the tissue (see methods). The binomial
128 distribution is an appropriate model when the parental identity of sequencing reads is
129 known, which is not the case when aligning to a reference genome. A reference genome
130 will contain SNPs from both parents, making the parental identity of aligned reads
131 ambiguous and producing reference allelic expression ratios that represent expression
132 of both parental X-alleles (Fig. 1B).
133
134     We fold the distribution of reference allelic-expression ratios around 0.50 to
135 aggregate data across both alleles and enable a robust estimate of the XCI ratio
136 magnitude for the bulk RNA-seq sample (Fig. 1B). We fit folded-normal distributions to
137 the reference allelic expression ratios of multiple SNPs per sample, which serves as a
138 continuous approximation of the underlying depth-dependent mixture of folded-binomial

139  distributions per SNP. The mean of the fitted distribution is the estimate of the XCI ratio
140  (Fig. 1B). We also incorporate specific steps to address confounding factors that can
141  impact X-linked allelic expression, including reference bias and escape from XCI[40,41]
142  (Supp. Figs. 1-2, see methods). Interestingly, we find the strongest signals of escape
143  from XCI near chromosomal ends across all species (Supp. Fig. 2), suggesting escape
144  within pseudo-autosomal regions is conserved across mammals[40,42]. Previously, we
145  validated our SNP filtering and XCI modeling approach using phased RNA-seq data
146  (where haplotype information is known for each variant) from the EN-TEx consortium[43],
147  achieving nearly perfect agreement in XCI ratio estimates for samples with folded XCI
148  ratios of 0.60 or higher, demonstrating the accuracy of our approach.
149
150      By calling SNPs from RNA-seq reads and employing folded distributions to model
151  reference-aligned allelic expression, we can estimate the magnitude of XCI in any
152  female mammalian bulk RNA-seq sample. We source female annotated bulk RNA-seq
153  samples of 8 non-human mammalian species from the SRA database (Fig. 1C),
154  additionally including cross-tissue human samples from the GTEx dataset. After
155  processing, the number of samples with a minimum of 10 well-powered SNPs for
156  estimating XCI ratios are 130 macaca (mean of 28 SNPs +- 17 SD), 275 horse (mean of
157  54 SNPs +- 36 SD), 269 dog (mean of 29 SNPs +- 13 SD), 328 rat (mean of 26 SNPs
158  +- 13 SD), 383 goat (mean of 34 SNPs +- 14 SD), 624 pig (mean of 50 SNPs +- 28 SD),
159  731 sheep (mean of 79 SNPs +- 42 SD), 1328 cow (mean of 32 SNPs +- 19 SD), and
160  4877 human (mean of 56  SNPs +- 23 SD, 314 total individuals) samples (Fig. 1C,
161  Supp. Fig. 1). Aggregating reference allelic expression ratios for samples with similar
162  estimated XCI ratios (0.05 bins) clearly reveals the expected haplotype expression
163  distributions, demonstrating the applicability of folded models (Supp. Fig. 3). Following
164  XCI ratio modeling, we then generate population-level distributions by unfolding the
165  distribution of folded XCI ratio sample estimates per species (Fig. 1D).
166
167      To ensure the allelic variability we report from X-linked SNPs is specific to XCI, we
168  estimate autosomal allelic imbalances for all samples using the same pipeline and
169  approach as for the X-chromosome analysis (Supp. Fig. 4, see methods). Comparing
170  allelic imbalances across the two autosomes closest in size to the X-chromosome
171  reveals the vast majority of samples across all species are biallelically balanced for
172  autosomal expression, as expected (Supp. Fig. 4). Several species (Pig, Cow, Goat,
173  Rat, Sheep, and Dog) exhibit small subsets of samples that are consistently imbalanced
174  across the two autosomes and the X-chromosome, indicative of a global influence on
175  allelic-expression independent of XCI (Supp. Fig. 4). These samples with global allelic
176  imbalances are excluded from all downstream analysis, ensuring the population
177  distributions of XCI ratios reflect variability specific to XCI.
178
179  **Models of embryonic stochasticity explain adult population XCI variability**
180
181      After generating population distributions of XCI ratios for the 9 mammalian species,
182  we next explore how well models of embryonic stochasticity explain the observed adult
183  XCI ratio variability. The initial variability in XCI ratios among mammalian embryos is

184 dependent on the number of cells present during XCI (Fig. 1A), where adult variability
185 can be modeled to infer embryonic cell counts.
186
187     An important consideration when estimating embryonic cell counts from XCI
188 variability in adult tissues is the fact adult tissues only represent the embryonic lineage
189 of the blastocyst as opposed to extra-embryonic lineages. This positions XCI variability
190 of adult tissue samples as informative for the number of cells present within the last
191 common lineage decision for all adult cells, i.e. the number of cells present within the
192 epiblast of the mammalian blastocyst. If XCI occurs after epiblast specification, the
193 variability in XCI ratios is determined by the number of epiblast cells at the time of XCI.
194 On the other hand, if XCI occurs before epiblast specification, XCI variability within the
195 embryonic lineage is influenced by both the initial stochasticity of XCI and the
196 stochasticity associated with cell sampling during epiblast lineage specification. The
197 temporal ordering of XCI among these lineage events cannot be resolved without cross-
198 tissue sampling of both the extra-embryonic and embryonic tissues. As such, estimating
199 cell counts solely on XCI variability in adult tissues provides an estimate of the number
200 of cells present within the epiblast of the embryo.
201
202     Figure 2A presents the unfolded population distributions of XCI ratios in the 9
203 mammalian species we sampled, ranging from the least variable (macaca) to most
204 variable (dog). We fit normal distributions as continuous approximations to the
205 underlying binomial distribution that defines the relationship between cell counts and
206 XCI ratio variability (Fig. 1A,D, see methods). We focus on the tails of the distributions,
207 as our previous validation using phased data indicated increased uncertainty for folded
208 XCI ratio estimates between 0.5-0.6, which translates to unfolded estimates between
209 0.4-0.6.  At a broad level, population XCI ratio variability varies substantially across the
210 sampled mammalian species. Our estimates for the number of epiblast cells present at
211 the time of XCI include 65 (macaca), 31 (rat), 23 (pig), 16 (goat), 15 (horse), 14 (sheep),
212 14 (cow), 13 (human) and 8 (dog) cells, with associated 95% confidence intervals
213 presented in figure 2B. The error between the empirical XCI ratio distributions and the
214 normal fitted distributions is strikingly small, with a mean of 0.00538 (+- 0.0101 SD)
215 across the species (Supp. Fig. 5). This indicates models of embryonic stochasticity can
216 explain observed XCI ratio variability in adult populations exceptionally well.
217
218     For the least and most variable species (macaca and dog), the estimated autosomal
219 imbalances offer additional context for the reported XCI population variability. The
220 reported X-linked variability in macaca is in excess to the reported autosomal allelic
221 variability (Supp. Fig. 4). This demonstrates the X-linked population variability for
222 macaca, while strikingly small, is specific to XCI and informative for estimating cell
223 counts. On the other hand, the dog population is the only one that contains samples
224 with strong allelic imbalances on only one autosome, where autosomal imbalances in all
225 other species are global (Supp. Fig. 4). This is suggestive of broader genomic
226 incompatibilities within the dog population. The reported X-linked population variability in
227 dog is likely a combination of XCI and broader allelic incompatibilities, positioning our
228 estimate of 8 cells as a likely underestimate due to excess variability outside of XCI.
229

230      Modeling XCI ratio variability across numerous species allows comparisons in light
231   of evolution for determining generalizable or species-specific characteristics in XCI.
232   Broadly, we demonstrate XCI ratios are variable in each species we assess, revealing
233   variability in XCI ratios itself as a conserved characteristic of XCI. The exact variance in
234   XCI ratios varies across the species, with differences in the timing of XCI and/or
235   embryonic/extra-embryonic lineage specification (differences in cell counts) as one
236   putative explanation. We compare our estimated cell counts to the evolutionary
237   relationships among the species we assess (Fig. 2B), suggesting that variability in
238   timing for these early embryonic events can be recent evolutionary adaptations. This is
239   highlighted by the large differences in cell counts between macaca and humans. When
240   viewed through the lens of cell divisions (log2 of the estimated cell counts, Fig. 2B), the
241   differences in XCI ratio variability among the species can be explained by differences in
242   a range of only 3 cell divisions, a narrow developmental window. This demonstrates
243   even slight changes in the timing of XCI or embryonic/extra-embryonic lineage
244   specification across mammalian species can produce large differences in population
245   XCI ratio variability, as explained through the inherent stochasticity of XCI.
246
247   **XCI ratios are not associated with X-linked heterozygosity**
248
249      After determining stochastic models can explain population XCI ratio variability
250   across mammalian species, we turn to testing whether we can identify any genetic
251   correlates with XCI ratios. Our approach leveraging natural genetic variation to quantify
252   XCI ratios enables us to assess a large catalog of genetic variants for associations with
253   XCI ratios across mammalian species (10,735 macaca SNPs, 12,024 rat SNPs, 23,603
254   pig SNPs, 16,123 goat SNPs, 10,281 horse SNPs, 53,505 sheep SNPs, 18,509 cow
255   SNPs, 16,168 human SNPs, and 10,050 dog SNPs). One putative genetic contribution
256   to XCI ratio variability is allelic selection during development, where increased X-linked
257   heterozygosity (i.e., genetic distance), is more likely to produce selective pressures
258   between the two X-alleles. It follows that samples with higher X-linked heterozygosity
259   would be expected to exhibit more variability in XCI ratios.
260
261      We score X-linked heterozygosity per sample as the ratio of the detected SNPs
262   within a sample to the number of unique SNPs identified across all samples, relative for
263   each species (Fig. 3A). This quantification also serves as a measure of inbreeding, with
264   decreased heterozygosity associated with a higher degree of inbreeding[44]. The trend in
265   heterozygosity across species is as expected, with rats (likely laboratory strains) as the
266   most inbred (Fig. 3A). Next, we examine the correlations between sample
267   heterozygosity and the estimated XCI ratio, as well as the estimated XCI variability
268   across SNPs in each sample (mean and standard deviation of the fitted folded-normal
269   distribution per sample, Fig. 3B). Across all species, X-linked heterozygosity showed a
270   near-zero correlation with the estimated XCI ratio, indicating a lack of association
271   between X-linked genetic variability and XCI ratio variability (Fig. 3B). However, we
272   observe moderate correlations between sample heterozygosity and the estimated
273   variability in SNP allelic ratios in three species: rat (corr: 0.576), macaca (corr: 0.459),
274   and cow (corr: 0.364), notably the most inbred species (Fig. 3A, Supp. Fig. 6). The
275   increased variability in allelic expression present only within the most inbred species

276 could potentially reflect gene-specific regulatory events between parental haplotypes[45]
277 rather than a direct genetic effect on XCI.
278
279 **Low frequency variants exhibit moderate associations with XCI ratios**
280
281 After investigating relationships between genetic variation and XCI ratios at a
282 broad level across the whole X-chromosome, we next asked if individual variants might
283 be associated with extreme XCI ratios. Variants that affect the expression and/or
284 function of the genetic elements that control XCI can result in highly skewed XCI ratios,
285 as documented in human studies[15]. This can also occur in other X-linked genes, if the
286 resulting differential in gene activity exerts a selective pressure across the X-alleles, as
287 documented in disease cases[14,16]. We test the association between XCI ratios and
288 individual variants for all variants detected in each species with a minimum of 10
289 samples, quantified through the area-under-the-receiver-operating-curve statistic
290 (AUROC). For each species, we rank the samples based on their estimated XCI ratio
291 and score the placement of samples carrying a given variant within the ordered list (Fig.
292 4A). If all the samples with that variant are at the top of the ordered list, the XCI ratio
293 can be said to have perfectly predicted the presence of that variant, quantified with an
294 AUROC of exactly 1. An AUROC of 0.50 indicates the XCI ratio performs no better than
295 random chance for predicting the presence of the variant.
296
297 The distribution of AUROCs for each species show striking similarities to a null
298 comparison (Fig. 4B, see methods), indicating a pervasive lack of association between
299 XCI ratios and individual variants. However, a small subset of variants in each species
300 exhibits moderate associations (AUROCs >= 0.75). By comparing each variant's
301 AUROC with its frequency in the species, we find that the variants with moderate
302 associations occur at low frequencies within the sampled populations (Fig. 4C, Supp.
303 Fig. 7). We investigate whether this relationship is simply due to a lack in power with
304 bootstrap simulations, demonstrating moderate AUROCs (>= 0.75) are robust to their
305 small sample sizes (Supp. Fig. 7). Figure 4D displays these variants along with their
306 gene annotations for each species. Notably, several genes in humans with moderate
307 AUROCs have prior evidence for associations with skewed XCI, namely MECP2[46],
308 IDS[47] (also identified in macaca), IRAK1[48], and FLNA[49]. This suggests our analysis is
309 able to recover putative examples of selection impacting XCI ratios via disease-variants,
310 though with small effect sizes and low frequencies in our sampled population. In
311 general, we are unable to identify strong associations between genetic variation and
312 XCI ratios across all 9 mammalian species, both along the whole X-chromosome and
313 for individual variants.
314
315 **Discussion**
316
317 We modeled tissue XCI ratios from bulk RNA-seq samples across 9 mammalian
318 species and found population-level variation in XCI ratios, reflecting differences in
319 developmental events such as XCI timing or lineage specification. We showed that
320 embryonic stochasticity models fit the XCI data well and estimated epiblast cell counts
321 at the time of XCI across species. We also searched for genetic factors influencing XCI

322  ratios and found a pervasive lack of strong genetic associations with XCI ratios,
323  indicating that XCI variability is mainly driven by stochasticity rather than genetic
324  variation in mammals.
325
326  The lack of cross-mammalian comparisons of population XCI variability has
327  previously limited our understanding on the sources of XCI variability in mammals. The
328  existence of XCE-alleles in laboratory mice[18–20,33] has supported the hypothesis that a
329  similar genetic mechanism can exist in humans and drive population XCI variability[21],
330  though evidence for XCE-alleles in human populations remains inconclusive[22] and data
331  from other mammalian species is historically absent. Although genetic influences on
332  XCI, particularly variants affecting XIST[15] or disease-associated variants[34–37], have been
333  identified, they do not constitute a general mechanism that can fully account for
334  observed population-level XCI variability. Comprehensive assessment of genetic
335  influence on XCI would require combined DNA and RNA sequencing data, which is
336  challenging to perform at a large scale across mammalian populations. Our approach
337  for extracting heterozygous variants from RNA-seq data[28], while providing a sample of
338  genetic variability, is still able to assess hundreds of X-linked genes per species for
339  associations with XCI and culminated in only weak evidence for limited genetic
340  influence on XCI ratios. In contrast, we demonstrated models of embryonic stochasticity
341  can explain population XCI variability with exceedingly small amounts of error
342  consistently across mammalian species, providing a much more general explanation for
343  population XCI variability.
344
345  Besides X-linked disorders and XIST-variants, other factors that may affect XCI
346  ratio variability are genomic incompatibilities[45] and stochastic allelic drift during
347  development[20]. We found a link between the variance in X-linked allelic expression and
348  the inbreeding level of some species (Fig. 2B), as well as autosome-specific allelic
349  imbalances in dog (Supp. Fig. 4). This implies that X-linked allelic expression variability
350  may result from both the bulk XCI ratio and the genomic incompatibilities between the
351  parental genomes[45], depending on the species. We controlled for global allelic
352  imbalances by excluding samples that showed them (Supp. Fig. 4), which confirms that
353  the allelic-expression variability on the X-chromosome is specific to XCI. Moreover,
354  developmental allelic drift may introduce XCI ratio variability beyond the initial random
355  choice of allelic inactivation[20]. While our previous cross-tissue analysis of XCI ratios in
356  humans[28] showed consistent XCI ratios across tissues, suggesting allelic drift is not a
357  major factor in XCI ratio variability, similar data for non-human mammals is missing.
358  These factors indicate that our cell count estimates are lower bound estimates for the
359  number of cells needed to produce the observed XCI ratio variability as purely due to
360  embryonic stochasticity.
361
362  A general model of X-linked genetic variability depletion (due to strong purifying
363  selection in males[50–53]) accounts for the lack of evidence for broad allelic-selection or
364  individual variants influencing XCI ratio variability in mammals, as both parental alleles
365  are mostly equivalent. This does not apply to disease variants, but they cannot explain
366  the widespread XCI ratio variability across mammalian species. We find genes
367  associated with increased XCI ratios that have prior evidence for causing highly skewed

368 XCI in disease cases, but their effect sizes and population frequencies are small in our
369 samples. Therefore, the inherent stochasticity of XCI during embryogenesis is the main
370 source of the observed XCI ratio variability in mammalian populations.
371
372 **Methods**
373
374 **Snakemake pipeline for RNA-seq alignment and variant identification**
375
376 All non-human mammalian fastq data was downloaded from the Sequencing
377 Read Archive (SRA, https://www.ncbi.nlm.nih.gov/sra ), where only samples annotated
378 as female were selected, using the metadata provided through SRA. Details for
379 download and processing of the GTEx[39] data can be found here[28]. The entire sample
380 processing pipeline uses a standard collection of bioinformatics software tools, all
381 available for installation via Conda (STAR[54] v2.7.9a, GATK[55] v4.2.2.0, samtools[56] v1.13,
382 igvtools[57] v2.5.3, and sra-tools 2.11.0). All Snakemake workflow rules, environment
383 setup procedure, analysis commands and options, and underlying libraries are available
384 on Github at https://github.com/gillislab/cross_mammal_xci , and
385 https://github.com/gillislab/xskew. Briefly, a .fastq file acts as input, for either single- or
386 pair-end sequencing experiments, and a .vcf and .wig file are produced as outputs for
387 subsequent compiling of allele-specific read counts in R v4.3.0. The R script used for
388 combining the .vcf and .wig information is also made available at
389 https://github.com/gillislab/cross_mammal_xci/tree/main/R. Genome generation and
390 alignment was performed with STAR, with the addition of the WASP[58] algorithm for
391 identifying and excluding reference biased reads. We extract chromosome-specific
392 alignments from the .bam file (X chromosome or specific autosomes) and use GATK
393 tools to identify heterozygous SNPs from that chromosome. The suite of GATK tools for
394 identifying heterozygous variants from RNA-sequencing data was used following the
395 GATK Best Practices recommendations. Specifically, the tools utilized include
396 AddOrReplaceReadGroups -> MarkDuplicates -> SplitNCigarReads -> HaplotypeCaller
397 -> SelectVariants -> VariantFiltration.
398
399 Reference genomes and gene annotations (.gtf files) for each species were
400 sourced from the NCBI Refseq database (https://www.ncbi.nlm.nih.gov/refseq/ ). In
401 each case the latest assembly version path was used, and the genomic.fna and
402 genomic.gtf was downloaded. Annotated and indexed genomes were generated with
403 STAR using --runMode genomeGenerate with default parameters.
404
405 **SNP filtering**
406
407 Only SNPs with exactly two identified genotypes were included for analysis and
408 indels were excluded. We required each SNP to have a minimum of 10 reads mapped
409 to both alleles for a minimum read depth of 20 reads per SNP. Gene annotations for all
410 SNPs were extracted from the species-specific .gtf files. For XCI ratio modeling, we only
411 used SNPs found within annotated genes. For any sample with multiple SNPs identified
412 in a gene, we took the SNP with the highest read count to be the max-powered
413 representative of that gene, so each individual SNP is representative of a single gene.

414 In addition to implementing the WASP algorithm for excluding reference biased reads,
415 we filter out SNPs within each species whose mean expression ratios across samples
416 deviate strongly from 0.50 (mean allelic ratio < 0.40 and > 0.60, Supp. Fig. 1). This SNP
417 filtering also excludes potential eQTL effects that may impact allelic-expression outside
418 of the underlying XCI ratio.
419
420 **Identifying and excluding chromosomal regions that escape XCI**
421
422 We reasoned robust escape from XCI would produce more balanced biallelic
423 expression in samples with skewed XCI. We performed an initial pass at XCI ratio
424 modeling including all well-powered SNPs in a sample to identify samples with skewed
425 XCI ratios (XCI ratios >= 0.70 for all species except rat and macaca, where a threshold
426 of 0.60 was used due to a reduced incidence of skewed XCI in these species). Using
427 the subset of skewed samples for each species, we averaged the folded allelic-
428 expression ratios for all SNPs present in 1 mega-base (MB) bins across the X-
429 chromosome (Supp. Fig. 2). Chromosomal-bins that displayed balanced allelic
430 expression in opposition to the clearly skewed allelic expression of the rest of the
431 chromosome were excluded from analysis. Specifically, chromosomal bins with an
432 average allelic-expression < 0.65 for pig, goat, horse, sheep, and cow, < 0.60 in rat and
433 macaca, and <0.675 in dog were excluded (Supp. Fig. 2) The ends of the X-
434 chromosome in all species, except rat, demonstrated strong balanced biallelic
435 expression, indicative of escape within putative pseudo-autosomal regions. We
436 excluded any bin within these putative pseudo-autosomal regions regardless of average
437 allelic expression. The escape threshold for dog was increased to exclude all bins within
438 the dog putative pseudo-autosomal region.
439
440 **Modeling XCI ratios with the folded-normal distribution**
441
442 Starting with a single parental allele, the sampled maternal allelic-expression of a
443 heterozygous X-linked SNP can be modeled with a binomial distribution, dependent on
444 the ratio of active maternal X-alleles in the sample and the read depth of the SNP.
445
446 $$\frac{X_{mat}}{n_{reads}} \sim \frac{Bin(n_{reads}, p_{mat})}{n_{reads}}; \; E\left[\frac{X_{mat}}{n_{reads}}\right] = p_{mat}; \; Var\left(\frac{X_{mat}}{n_{reads}}\right) = \frac{p_{mat}(1-p_{mat})}{n_{reads}},$$  eq. 1
447
448 where $X_{mat}$ is the number of maternal allelic reads, $n_{reads}$ is the read depth of the SNP,
449 and $p_{mat}$ is the ratio of active maternal X-alleles. When aligned to a reference genome,
450 the parental phasing information is lost and the allelic-expression of X-linked SNPs can
451 instead be modeled with the folded-binomial model[59,60]. Since SNPs vary in read-depth,
452 we use a folded-normal model as an approximation of the underlying mixture of depth-
453 dependent folded-binomial distributions. The probability of allelic-expression under the
454 folded-normal model is defined as:
455

$$\Pr(x_{ratio}; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_{ratio} - \mu)^2}{2\sigma^2}} + \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_{ratio} + \mu - 1)^2}{2\sigma^2}}, \text{for } \mu \in [0.50, 1],$$  eq. 2

456

10

457  where $x_{ratio}$ is the folded allelic-expression ratio of a SNP, $\mu$ is the folded XCI ratio of
458  the sample, and $\sigma$ is the standard deviation of the folded-normal distribution. We utilize
459  a maximum-likelihood approach (negative log-likelihood minimization of eq. 2) to fit
460  folded-normal distributions to the observed folded allelic-expression ratios of at least 10
461  filtered SNPs per sample, taking the $\mu$ parameter of the maximum-likelihood folded-
462  normal distribution as the folded XCI ratio estimate of the sample.
463
464  **Modeling autosomal imbalances**
465
466     The folded-normal model can also be applied to autosomal data to estimate
467  allelic-imbalances. For each species, we extract chromosome-specific alignments from
468  the .bam file for the two autosomes closest in size to the X-chromosome (Supp. Fig. 4).
469  We employ the exact same processing pipeline and thresholds as used for the X-
470  chromosome. Any sample that displayed an autosomal imbalance greater than or equal
471  to a folded estimate of 0.60 (dotted lines in Supp. Fig. 4A) on either autosome was
472  excluded from downstream analysis.
473
474  **Modeling population XCI variability with models of embryonic stochasticity**
475
476     XCI is a binomial sampling event, where the number of cells choosing to inactivate
477  the same X-allele follows a binomial distribution defined as:
478

$$X \sim Bin(n_{cells}, p_{inact}),$$

eq. 3

479
480  where $X$ is the number of cells inactivating the same X-allele, $n_{cells}$ is the number of
481  cells present at the time of XCI, and $p_{inact}$ is the probability of inactivation (0.50).
482
483  Embryonic XCI ratios can be modeled as:                                        eq. 4
484

$$\frac{X}{n_{cells}} \sim \frac{Bin(n_{cells}, p_{inact})}{n_{cells}}$$

485
486  We estimate $n_{cells}$ by fitting normal distributions to the unfolded population XCI ratio
487  distributions of each species, as a continuous approximation for the underlying binomial
488  distribution. The variance of the normal distribution is defined as:
489

$$var_{normal} = Var\left(\frac{Bin(n_{cells}, p_{inact})}{n_{cells}}\right) = \frac{p_{inact}(1 - p_{inact})}{n_{cells}} = \frac{.5(1 - .5)}{n_{cells}}$$    eq. 5

490
491  We model population XCI ratios as:
492

$$\frac{X}{n_{cells}} \sim Norm\left(\mu, \sqrt{var_{normal}}\right),$$

eq. 6

493
494  where $\mu = p_{inact} = 0.50$ and $var_{normal}$ is computed for $n_{cells} \in [2, 200]$.
495

11

496   We identify the normal distribution with minimum sum-squared error between its
497 CDF and the empirical population XCI ratio CDF, minimizing error over the tails of the
498 distributions with percentiles <= 0.40 or >= 0.60 (Supp. Fig. 5). We compute 95%
499 confidence intervals about the cell number estimate $n_{cells}$ through bootstrap simulations.
500 We sample with replacement from the empirical population XCI ratio distribution,
501 matching the sample size of the original empirical population distribution, and fit a
502 normal model to derive a bootstrap estimate of $n_{cells}$. We repeat this for 2000
503 simulations to generate a bootstrapped distribution of $n_{cells}$, from which we derive the
504 95% confidence intervals, defined as the interval where 2.5% of the bootstrapped
505 distribution lies outside either end.
506
507 **Measuring sample X-linked heterozygosity**
508
509  We compute sample heterozygosity as the ratio of SNPs detected in a sample (20
510 read minimum) to the total number of unique SNPs identified across all samples for a
511 given species. We quantify associations between X-linked heterozygosity and XCI ratios
512 as the spearman correlation coefficient between the sample X-linked heterozygosity
513 ratio and the fitted mean and variance of the maximum-likelihood folded-normal
514 distribution of the sample (Fig. 3B-C, Supp. Fig. 6). We only consider samples with at
515 least 10 detected SNPs.
516
517 **Quantifying variant associations with extreme XCI ratios**
518
519   We quantify the strength of XCI ratios as a predictor for the presence of a given
520 variant through the AUROC metric. Given a ranked list of data (XCI ratios) and an
521 indicator of true positives (samples with a given variant), the AUROC quantifies the
522 probability a true positive is ranked above a true negative. An AUROC of 1 indicates all
523 true positive samples were ranked above all true negative samples, demonstrating XCI
524 ratios were a perfect predictor for the presence of that variant. An AUROC of 0.50
525 indicates random placement of true positives and negatives in the ranked list,
526 demonstrating XCI ratios performed no better than random chance for predicting the
527 presence of that variant. We compute the AUROC through the Mann-Whitney U-test,
528 defined as:
529

530
$$AUROC = \frac{U}{n_{pos} + n_{neg}},$$
     eq. 7

531
532 where $U$ is the Mann-Whitney U-test test statistic, computed in R with
533 wilcox.test(alternative = 'two.sided'), $n_{pos}$ is the number of true positive samples and
534 $n_{neg}$ is the number of true negative samples. We generate a null AUROC per variant by
535 randomly shuffling the true positive and negative labels. The variant frequency is
536 defined as the number of samples that carry a given variant over the total number of
537 samples for a given species. The p-value for a given AUROC is the p-value associated
538 with the Mann-Whitney U-test test statistic ($U$), where we determine significance as an
539 FDR-corrected p-value <= 0.05. We perform FDR correction for all p-values computed

540    for all variants across the 9 species through the Benjamini-Hochberg method,
541    implemented in R via p.adjust(method = 'BH').
542
543        We estimate the power of each variant through bootstrap simulations. We
544    randomly sample with replacement the XCI ratios of the true positive and true negative
545    samples, those that either carry or do not carry a given variant. We match the sample
546    size of the original true positive and negative labels. We compute a bootstrapped
547    AUROC and p-value from the simulated data, repeating for 2000 simulations to
548    compute a bootstrapped distribution of AUROCs. The AUROC power (Supp. Fig. 7B) is
549    defined as the fraction of bootstrapped AUROCs that are significant, using a
550    significance threshold of p-value <= 0.05. The AUROC effect size power (Supp. Fig. 7C)
551    is defined as the fraction of bootstrapped AUROCs that are >= 0.75. We also report the
552    variance of the bootstrapped AUROC distribution per variant in Supp. Fig. 7D. We
553    exclude all variants classified as reference biased from Supp. Fig. 1, with the
554    distributions of AUROCs for the reference biased and non-reference biased SNPs
555    presented in Supp. Fig. 7E.
556
557    **Software**
558
559        All analysis was performed in R[61] v4.3.0. All plots were generated using ggplot2[62]
560    v3.4.2 functions. The phylogenetic tree in Fig. 2B was generated from TimeTree
561    http://www.timetree.org/.
562
563    **Data and Code availability**
564
565        All associated code can be found at
566    https://github.com/gillislab/cross_mammal_xci. This includes the snakemake pipeline
567    used for processing the non-human mammalian data as well as all R notebooks used
568    for data analysis and figure generation.
569
570    **Author Contributions**
571

581
582
583
584
585

## References

1. Lyon, M. F. Gene Action in the X -chromosome of the Mouse ( Mus musculus L.). *Nature* **190**, 372–373 (1961).

2. Migeon, B. R. An overview of X inactivation based on species differences. *Semin. Cell Dev. Biol.* **56**, 111–116 (2016).

3. Okamoto, I. *et al.* Eutherian mammals use diverse strategies to initiate X-chromosome inactivation during development. *Nature* **472**, 370–374 (2011).

4. Ohno, S. *Sex Chromosomes and Sex Linked Genes*. (Springer Berlin, Heidelberg, 1966).

5. Lyon, M. F. X-chromosome inactivation and developmental patterns in mammals. *Biol. Rev. Camb. Philos. Soc.* **47**, 1–35 (1972).

6. van den Berg, I. M. *et al.* X Chromosome Inactivation Is Initiated in Human Preimplantation Embryos. *Am. J. Hum. Genet.* **84**, 771–779 (2009).

7. Evans, H. J., Ford, C. E., Lyon, M. F. & Gray, J. DNA Replication and Genetic Expression in Female Mice with Morphologically Distinguishable X Chromosomes. *Nature* **206**, 900–903 (1965).

8. Wu, H. *et al.* Cellular resolution maps of X-chromosome inactivation: implications for neural development, function, and disease. *Neuron* **81**, 103–119 (2014).

9. Mutzel, V. *et al.* A symmetric toggle switch explains the onset of random X inactivation in different mammals. *Nat. Struct. Mol. Biol.* **26**, 350–360 (2019).

10. Migeon, B. *Females Are Mosaics: X Inactivation and Sex Differences in Disease. Females Are Mosaics* (Oxford University Press, 2013).

608    11.     Amos-Landgraf, J. M. *et al.* X Chromosome–Inactivation Patterns of 1,005 Phenotypically

609            Unaffected Females. *Am. J. Hum. Genet.* **79**, 493–499 (2006).

610    12.     Shvetsova, E. *et al.* Skewed X-inactivation is common in the general female population.

611            *Eur. J. Hum. Genet.* **27**, 455–465 (2019).

612    13.     Fang, H., Deng, X. & Disteche, C. M. X-factors in human disease: impact of gene content

613            and dosage regulation. *Hum. Mol. Genet.* **30**, R285–R295 (2021).

614    14.     Migeon, B. R. Non-random X chromosome inactivation in mammalian cells. *Cytogenet.*

615            *Cell Genet.* **80**, 142–148 (1998).

616    15.     Plenge, R. M. *et al.* A promoter mutation in the XIST gene in two unrelated families with

617            skewed X-chromosome inactivation. *Nat. Genet.* **17**, 353–356 (1997).

618    16.     Belmont, J. W. Genetic control of X inactivation and processes leading to X-inactivation

619            skewing. *Am. J. Hum. Genet.* **58**, 1101–1108 (1996).

620    17.     Brown, C. & Robinson, W. The causes and consequences of random and non-random X

621            chromosome inactivation in humans: X chromosome inactivation in humans. *Clin. Genet.* **58**,

622            353–363 (2000).

623    18.     Cattanach, B. M. & Isaacson, J. H. Genetic control over the inactivation of autosomal

624            genes attached to the X-chromosome. *Z Vererbungsl* **96**, 313–323 (1965).

625    19.     Simmler, M. C., Cattanach, B. M., Rasberry, C., Rougeulle, C. & Avner, P. Mapping the

626            murine Xce locus with (CA)n repeats. *Mamm. Genome Off. J. Int. Mamm. Genome Soc.* **4**,

627            523–530 (1993).

628    20.    Sun, K. Y. *et al.* Bayesian modeling of skewed X inactivation in genetically diverse mice

629        identifies a novel Xce allele associated with copy number changes. *Genetics* **218**, iyab034

630        (2021).

631    21.    Peeters, S. B., Yang, C. & Brown, C. J. Have humans lost control: The elusive X-controlling

632        element. *Semin. Cell Dev. Biol.* **56**, 71–77 (2016).

633    22.    Bolduc, V. *et al.* No evidence that skewing of X chromosome inactivation patterns is

634        transmitted to offspring in humans. https://www.jci.org/articles/view/33166/pdf (2008)

635        doi:10.1172/JCI33166.

636    23.    Gandini, E., Gartler, S. M., Angioni, G., Argiolas, N. & Dell'Acqua, G. Developmental

637        implications of multiple tissue studies in glucose-6-phosphate dehydrogenase-deficient

638        heterozygotes. *Proc. Natl. Acad. Sci.* **61**, 945–948 (1968).

639    24.    Gandini, E. & Gartler, S. M. Glucose-6-phosphate Dehydrogenase Mosaicism for studying

640        the Development of Blood Cell Precursors. *Nature* **224**, 599–600 (1969).

641    25.    Nesbitt, M. N. X chromosome inactivation mosaicism in the mouse. *Dev. Biol.* **26**, 252–

642        263 (1971).

643    26.    Fialkow, P. J. Primordial cell pool size and lineage relationships of five human cell types*.

644        *Ann. Hum. Genet.* **37**, 39–48 (1973).

645    27.    McMahon, A., Fosten, M. & Monk, M. X-chromosome inactivation mosaicism in the

646        three germ layers and the germ line of the mouse embryo. *J. Embryol. Exp. Morphol.* **74**,

647        207–220 (1983).

648   28.    Werner, J. M., Ballouz, S., Hover, J. & Gillis, J. Variability of cross-tissue X-chromosome

649        inactivation characterizes timing of human embryonic lineage specification events. *Dev. Cell*

650        **57**, 1995-2008.e5 (2022).

651   29.    Bittel, D. C. *et al.* Comparison of X-chromosome inactivation patterns in multiple tissues

652        from human females. *J. Med. Genet.* **45**, 309–313 (2008).

653   30.    Brown, C. J. *et al.* The human XIST gene: Analysis of a 17 kb inactive X-specific RNA that

654        contains conserved repeats and is highly localized within the nucleus. *Cell* **71**, 527–542

655        (1992).

656   31.    Dossin, F. & Heard, E. The Molecular and Nuclear Dynamics of X-Chromosome

657        Inactivation. *Cold Spring Harb. Perspect. Biol.* a040196 (2021)

658        doi:10.1101/cshperspect.a040196.

659   32.    Dixon-McDougall, T. & Brown, C. J. Multiple distinct domains of human XIST are required

660        to coordinate gene silencing and subsequent heterochromatin formation. *Epigenetics*

661        *Chromatin* **15**, 6 (2022).

662   33.    Calaway, J. D. *et al.* Genetic Architecture of Skewed X Inactivation in the Laboratory

663        Mouse. *PLOS Genet.* **9**, e1003853 (2013).

664   34.    Migeon, B. R. Studies of skin fibroblasts from 10 families with HGPRT deficiency, with

665        reference in X-chromosomal inactivation. *Am. J. Hum. Genet.* **23**, 199–210 (1971).

666   35.    Migeon, B. R. *et al.* Adrenoleukodystrophy: evidence for X linkage, inactivation, and

667        selection favoring the mutant allele in heterozygous cells. *Proc. Natl. Acad. Sci. U. S. A.* **78**,

668        5066–5070 (1981).

669    36.    Devriendt, K. *et al.* Skewed X-chromosome inactivation in female carriers of dyskeratosis

670        congenita. *Am. J. Hum. Genet.* **60**, 581–587 (1997).

671    37.    Plenge, R. M., Stevenson, R. A., Lubs, H. A., Schwartz, C. E. & Willard, H. F. Skewed X-

672        chromosome inactivation is a common feature of X-linked mental retardation disorders. *Am.*

673        *J. Hum. Genet.* **71**, 168–173 (2002).

674    38.    Schmidt, M. & Du Sart, D. Functional disomies of the X chromosome influence the cell

675        selection and hence the X inactivation pattern in females with balanced X-autosome

676        translocations: a review of 122 cases. *Am. J. Med. Genet.* **42**, 161–169 (1992).

677    39.    Lonsdale, J. *et al.* The Genotype-Tissue Expression (GTEx) project. *Nat. Genet.* **45**, 580–

678        585 (2013).

679    40.    Bonora, G. & Disteche, C. M. Structural aspects of the inactive X chromosome. *Philos.*

680        *Trans. R. Soc. B Biol. Sci.* **372**, 20160357 (2017).

681    41.    Fang, H., Disteche, C. M. & Berletch, J. B. X Inactivation and Escape: Epigenetic and

682        Structural Features. *Front. Cell Dev. Biol.* **7**, (2019).

683    42.    Posynick, B. J. & Brown, C. J. Escape From X-Chromosome Inactivation: An Evolutionary

684        Perspective. *Front. Cell Dev. Biol.* **7**, (2019).

685    43.    Rozowsky, J. *et al.* The EN-TEx resource of multi-tissue personal epigenomes & variant-

686        impact models. *Cell* **186**, 1493-1511.e40 (2023).

687    44.    Miller, J. M. *et al.* Estimating genome-wide heterozygosity: effects of demographic

688        history and marker type. *Heredity* **112**, 240–247 (2014).

689    45.    Shorter, J. R. *et al.* Male Infertility Is Responsible for Nearly Half of the Extinction

690        Observed in the Mouse Collaborative Cross. *Genetics* **206**, 557–572 (2017).

691    46.    Knudsen, G. P. S. *et al.* Increased skewing of X chromosome inactivation in Rett

692        syndrome patients and their mothers. *Eur. J. Hum. Genet.* **14**, 1189–1194 (2006).

693    47.    Kloska, A., Jakóbkiewicz-Banecka, J., Tylki-Szymańska, A., Czartoryska, B. & Węgrzyn, G.

694        Female Hunter syndrome caused by a single mutation and familial XCI skewing: implications

695        for other X-linked disorders. *Clin. Genet.* **80**, 459–465 (2011).

696    48.    Morcillo, P. *et al.* Directional X Chromosome Skewing of White Blood Cells from Subjects

697        with Heterozygous Mosaicism for the Variant IRAK1 Haplotype. *Inflammation* **43**, 370–381

698        (2020).

699    49.    Robertson, S. P. *et al.* Localized mutations in the gene encoding the cytoskeletal protein

700        filamin A cause diverse malformations in humans. *Nat. Genet.* **33**, 487–491 (2003).

701    50.    Payseur, B. A., Cutter, A. D. & Nachman, M. W. Searching for Evidence of Positive

702        Selection in the Human Genome Using Patterns of Microsatellite Variability. *Mol. Biol. Evol.*

703        **19**, 1143–1153 (2002).

704    51.    Avery, P. J. The population genetics of haplo-diploids and X-linked genes. *Genet. Res.* **44**,

705        321–341 (1984).

706    52.    Casto, A. M. *et al.* Characterization of X-Linked SNP genotypic variation in globally

707        distributed human populations. *Genome Biol.* **11**, R10 (2010).

708    53.    Veeramah, K. R., Gutenkunst, R. N., Woerner, A. E., Watkins, J. C. & Hammer, M. F.

709        Evidence for Increased Levels of Positive and Negative Selection on the X Chromosome versus

710        Autosomes in Humans. *Mol. Biol. Evol.* **31**, 2267–2282 (2014).

711    54.    Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21

712        (2013).

713    55.     McKenna, A. *et al.* The Genome Analysis Toolkit: A MapReduce framework for analyzing

714       next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).

715    56.     Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinforma. Oxf. Engl.*

716       **25**, 2078–2079 (2009).

717    57.     Robinson, J. T. *et al.* Integrative genomics viewer. *Nat. Biotechnol.* **29**, 24–26 (2011).

718    58.     van de Geijn, B., McVicker, G., Gilad, Y. & Pritchard, J. K. WASP: allele-specific software

719       for robust molecular quantitative trait locus discovery. *Nat. Methods* **12**, 1061–1063 (2015).

720    59.     Urbakh, V. Yu. Statistical Testing of Differences in Causal Behaviour of Two

721       Morphologically Indistinguishable Objects. *Biometrics* **23**, 137–143 (1967).

722    60.     Gart, J. J. A Locally Most Powerful Test for the Symmetric Folded Binomial Distribution.

723       *Biometrics* **26**, 129–138 (1970).

724    61.     R Core Team. R: A Language and Environment for Statistical Computing. (2023).

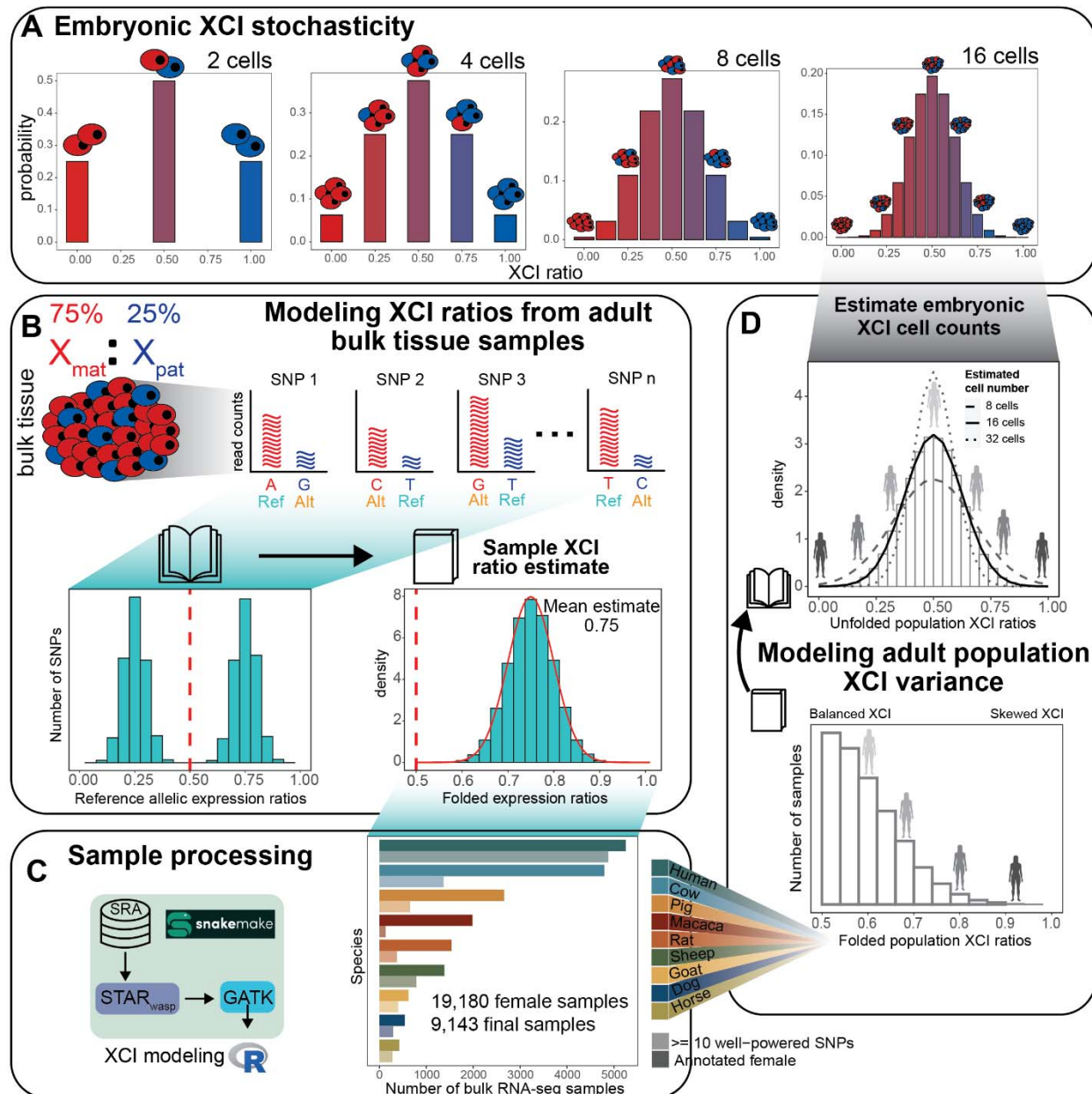725    62.     Wickham, H. ggplot2: Elegant Graphics for Data Analysis. *Springer-Verl. N. Y.* (2016).

726
727
728
729
730
731
732
733
734
735
736
737

**Figure 1: Reference aligned RNA-sequencing data enables scalable modeling of XCI ratios**

**A** Schematic demonstrating the relationship between the number of cells present at the time of XCI and the probability of all possible XCI ratios. Increased cell numbers result in decreased XCI ratio variance.

**B** Schematic for modeling XCI ratios from bulk reference-aligned RNA-seq data. The reference SNPs will contain both maternal and paternal SNPs, representing allelic expression from both parental haplotypes. Folded normal models are fit to the folded reference allelic expression ratios (like folding a book closed), with the mean of the maximum-likelihood distribution as the sample XCI ratio estimate.

21

750   **C** Schematic for sample processing (genome alignment and variant identification) and a
751   bar graph depicting the number of annotated female samples initially downloaded for
752   each species (bold color), with the number of samples per species with at least 10 well-
753   powered SNPs for XCI ratio modeling after processing (faded color).
754   **D** Schematic demonstrating the population modeling of XCI variability. Folded
755   population distributions are first produced per species and then are unfolded. Normal
756   distributions are fit to the unfolded population distribution to estimate the number of
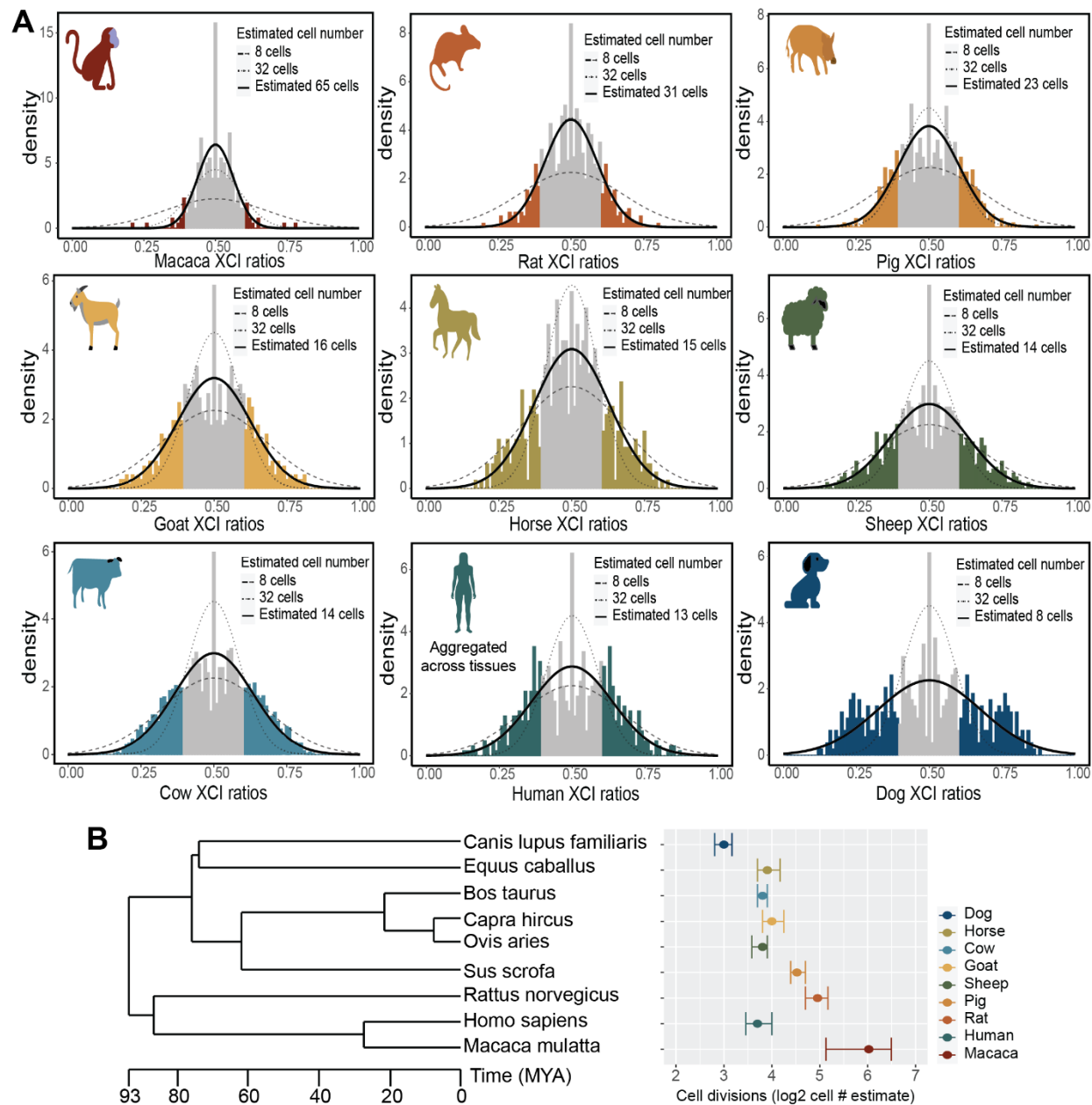757   embryonic cells required to produce the observed variance.
758
759
760
761
762

**Figure 2: Models of embryonic stochasticity explain adult population XCI variability**

**A** Unfolded distributions of XCI ratios per species, with the maximum-likelihood normal distribution depicted in bold, fitted to the tails of the distributions (shaded in sections of the distributions).

**B** Phylogenetic tree of the sampled mammalian species with their estimated embryonic cell counts on a log-2 scale, depicting the number of cell divisions that separate the estimated cell counts between the species. Error bars are 95% confidence intervals around the cell number estimate.
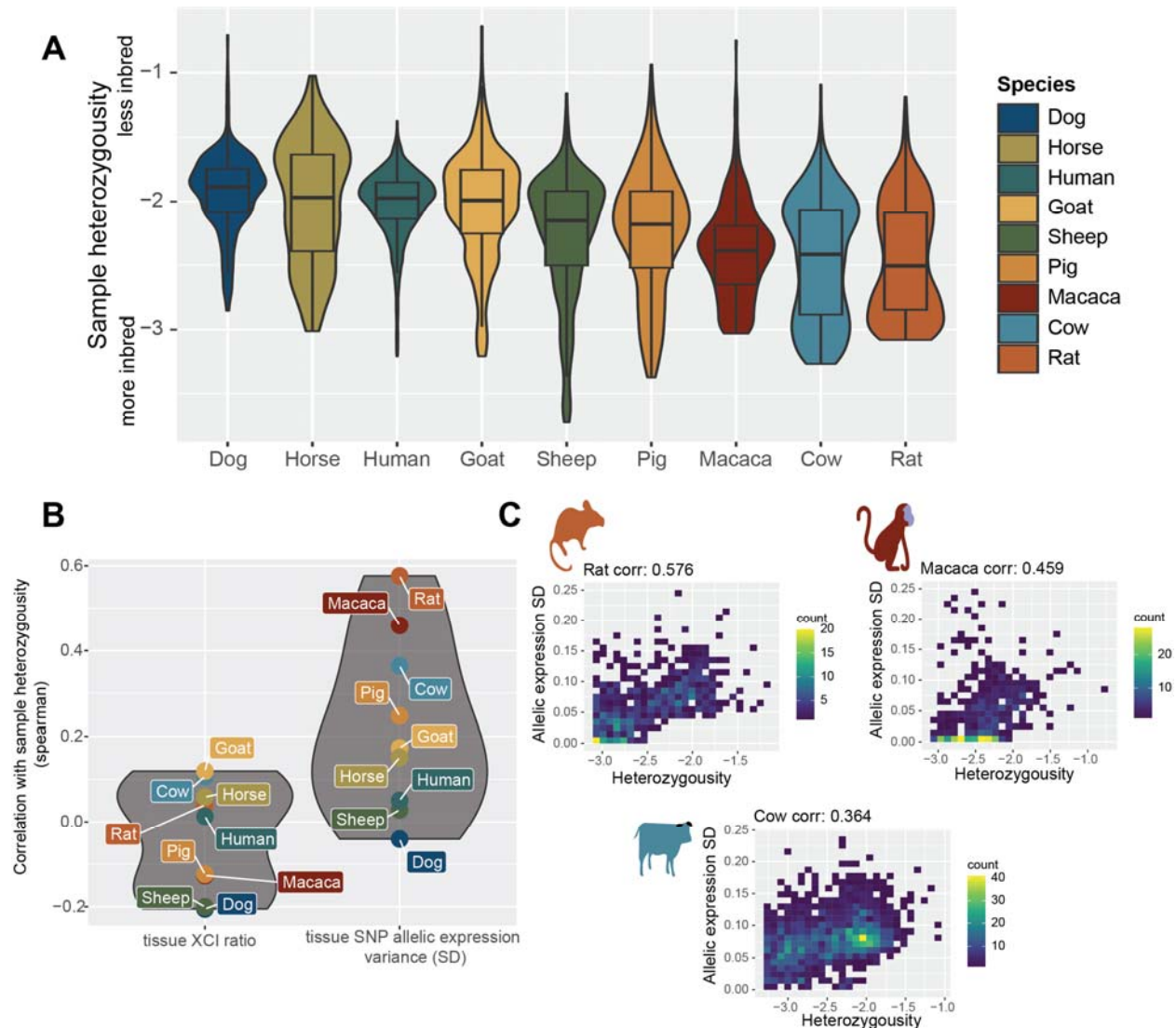
23

**Figure 3: XCI ratios are not associated with X-linked heterozygosity**

**A** Distributions of sample X-linked heterozygosity per species ordered by the median value. The y-axis is in log-10 scale, depicting the ratio of SNPs per sample to all unique identified SNPs per species. Boxplots depict the distributions' quartiles.

**B** The spearman correlation coefficients between sample X-linked heterozygosity and either the estimated standard deviation (SD) in X-linked allelic expression or the estimated XCI ratio of the sample (the SD and mean of the maximum-likelihood folded-normal model per sample).

**C** 2D Scatter plots of sample heterozygosity compared to the sample estimated X-linked allelic expression SD for the three species with moderate correlation coefficients. Color bars represent the number of samples in each 2D bin. Plots for the other species are in Supp. Fig. 6.
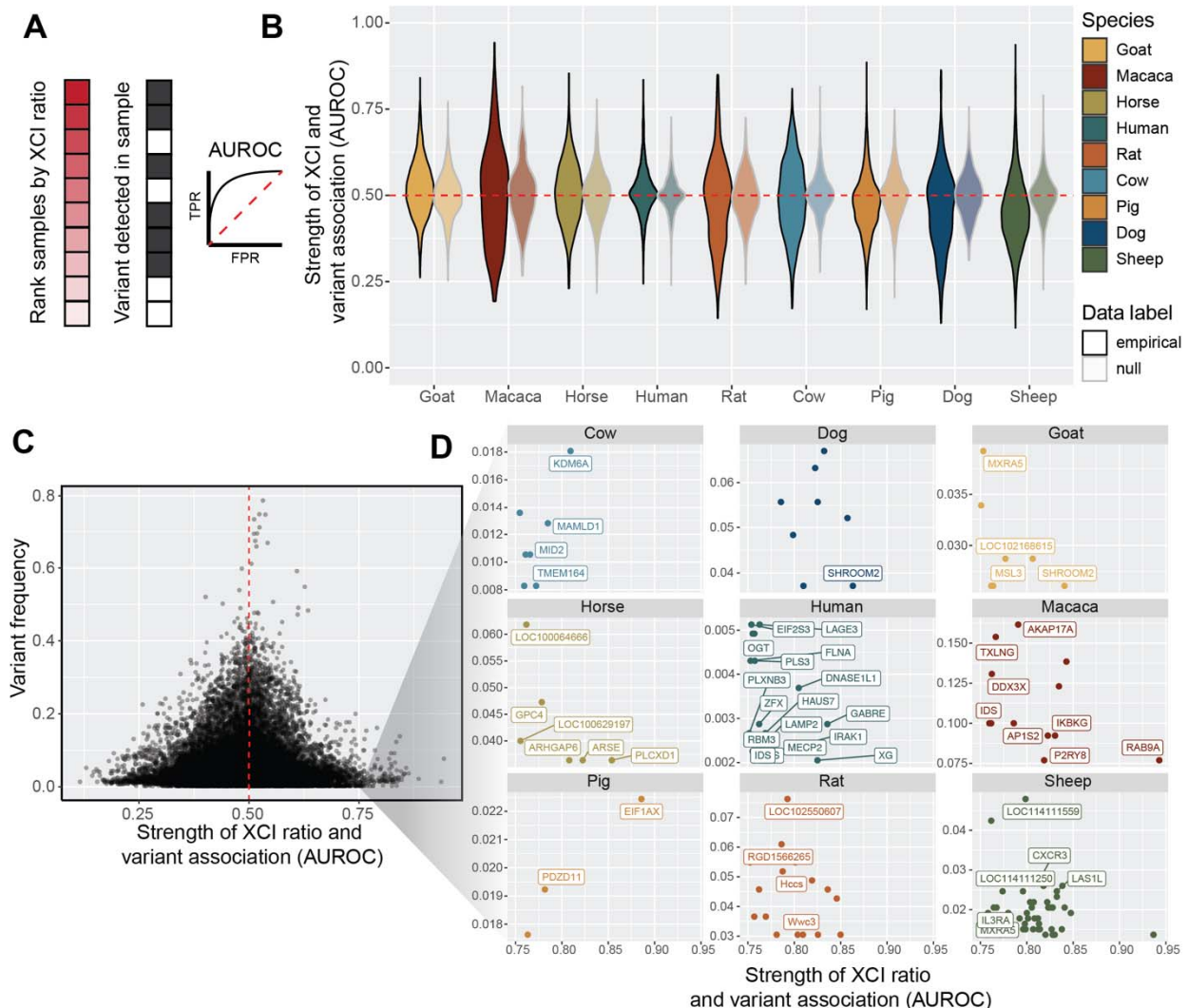
794



795

**Figure 4 Low frequency variants exhibit moderate associations with XCI ratios**
**A** Schematic depicting the AUROC quantification for testing the association between individual variants and extreme XCI ratios. Samples are ranked by their estimated XCI ratio, with the dark shaded red squares representing samples with more extreme XCI ratios. The position of samples with a given individual variant (grey squares) within the ranked list is used to compute the AUROC statistic. A variant with an AUROC value of 1 means all samples with that variant were at the top of the ranked list, whereas an AUROC value of 0.5 represents a random ordering of samples within the ranked list.
**B** Distributions of variant AUROCs for each species compared to a species-specific null distribution of AUROC values (faded distributions, see methods), ordered by the mean value of the empirical distributions. The red dotted line depicts an AUROC of 0.50, performance due to random chance.
**C** Scatter plot of variant AUROCs compared to each variant's prevalence (percent of samples with that variant, relative for each species) for all variants across all species. The red dotted line depicts an AUROC of 0.50, performance due to random chance. A threshold of AUROC >= 0.75 was used to identify SNPs with moderate associations with XCI ratios.

25

813 **D** Scatter plots depicting the same information as in C for the variants with moderate
814 associations with XCI ratios, but split by each species and including gene annotations.
815 SNPs not within annotated genes are unlabeled. Gene labels not present due to
816 overlapping labels are Macaca: ZBED1, Sheep: LOC101108113, LOC101115509,
817 LOC101117055, LOC105605313, LOC121818231, PPP2R3B, PRKX)
818
819
820