1    **Full Title**

2    *"Reference genome and transcriptome informed by the sex chromosome complement of the*

3    *sample increases ability to detect sex differences in gene expression from RNA-Seq data"*

4

5    **Short Title**

6    "*Sex chromosome complement informed alignment*"

7

8    **Authors**

9    Kimberly C. Olney[1,2], Sarah M. Brotman[1,4], Jocelyn P. Andrews[1, 5], Valeria A. Valverde-

10   Vesling[1], and Melissa A. Wilson[1,2,3*]

11

12   **Affiliations**

13   1. School of Life Sciences,

14   2. Center for Evolution and Medicine,

15   3. Center for Mechanisms of Evolution, The Biodesign Institute

16   Arizona State University, Tempe AZ 85282 USA

17   4. Department of Genetics, University of North Carolina, Chapel Hill NC 27599 USA

18   5. College of Osteopathic Medicine of the Pacific, Western University of Health Sciences,

19   Pomona CA 91766 USA

20

21   *Corresponding author

22   Melissa A. Wilson

23   School of Life Sciences | Arizona State University | PO Box 874501 | Tempe, AZ 85287-4501

24      mwilsons@asu.edu

**Abstract**

**Background:** Human X and Y chromosomes share an evolutionary origin and, as a consequence, sequence similarity. We investigated whether sequence homology between the X and Y chromosomes affects alignment of RNA-Seq reads and estimates of differential expression. We tested the effects of using reference genomes and reference transcriptomes informed by the sex chromosome complement of the sample's genome on measurements of RNA-Seq abundance and sex differences in expression.

**Results:** The default genome includes the entire human reference genome (GRCh38), including the entire sequence of the X and Y chromosomes. We created two sex chromosome complement informed reference genomes. One sex chromosome complement informed reference genome was used for samples that lacked a Y chromosome; for this reference genome version, we hard-masked the entire Y chromosome. For the other sex chromosome complement informed reference genome, to be used for samples with a Y chromosome, we hard-masked only the pseudoautosomal regions of the Y chromosome, because these regions are duplicated identically in the reference genome on the X chromosome. We analyzed transcript abundance in the whole blood, brain cortex, breast, liver, and thyroid tissues from 20 genetic female (46, XX) and 20 genetic male (46, XY) samples. Each sample was aligned twice; once to the default reference genome and then independently aligned to a reference genome informed by the sex chromosome complement of the sample, repeated using two different read aligners, HISAT and STAR. We then quantified sex differences in gene expression using featureCounts to get the raw count estimates followed by Limma/Voom for normalization and differential expression. We additionally created sex chromosome complement informed transcriptome references for use in pseudo-alignment using Salmon. Transcript abundance was quantified twice for each sample; once to the default target transcripts

2

48    and then independently to target transcripts informed by the sex chromosome complement of the

49    sample.

50    **Conclusions:** We show that regardless of the choice of read aligner, using an alignment protocol

51    informed by the sex chromosome complement of the sample results in higher expression estimates

52    on the pseudoautosomal regions of the X chromosome in both genetic male and genetic female

53    samples, as well as an increased number of unique genes being called as differentially expressed

54    between the sexes. We additionally show that using a pseudo-alignment approach informed on the

55    sex chromosome complement of the sample eliminates Y-linked expression in female XX samples.

56    **Key words:** RNA-Seq, sex chromosomes, differential expression, transcriptome, mapping,

57    alignment, pseudo-alignment, quantification.

58    **Author summary**

59    The human X and Y chromosomes share an evolutionary origin and sequence homology, including

60    regions of 100% identity; this sequence homology can result in reads misaligning between the sex

61    chromosomes, X and Y. We hypothesized that misalignment of reads on the sex chromosomes

62    would confound estimates of transcript abundance if the sex chromosome complement of the

63    sample is not accounted for during the alignment step. For example, because of shared sequence

64    similarity, X-linked reads could misalign to the Y chromosome. This is expected to result in

65    reduced expression for regions between X and Y that share high levels of homology. For this

66    reason, we tested the effect of using a default reference genome versus a reference genome

67    informed by the sex chromosome complement of the sample on estimates of transcript abundance

68    in human RNA-Seq samples from whole blood, brain cortex, breast, liver, and thyroid tissues of

69    20 genetic female (46, XX) and 20 genetic male (46, XY) samples. We found that using a reference

70    genome with the sex chromosome complement of the sample resulted in higher measurements of

71    X-linked gene transcription for both male and female samples and more differentially expressed

72    genes on the X and Y chromosomes. We additionally investigated the use of a sex chromosome

73    complement informed transcriptome reference index for alignment free quantification protocols.

74    We observed no Y-linked expression in female XX samples only when the transcript quantification

75    was performed using a transcriptome reference index informed on the sex chromosome

76    complement of the sample. We recommend that future studies requiring aligning RNA-Seq reads

77    to a reference genome or pseudo-alignment with a transcriptome reference should consider the sex

78    chromosome complement of their samples prior to running default pipelines.

**Background**

Sex differences in aspects of human biology, such as development, physiology, metabolism, and disease susceptibility are partially driven by sex specific gene regulation (Arnold et al., 2012; Khramtsova et al., 2018; Raznahan et al., 2018; Traglia et al., 2017). There are reported sex differences in gene expression across human tissues(Gershoni and Pietrokovski, 2017; Goldstein et al., 2014; Shi et al., 2016) and while some may be attributed to hormones and environment, there are documented genome-wide sex differences in expression based solely on the sex chromosome complement (Arnold and Chen, 2009). However, accounting for the sex chromosome complement of the sample in quantifying gene expression has been limited due to shared sequence homology between the sex chromosomes, X and Y, that can confound gene expression estimates.

The X and Y chromosomes share an evolutionary origin: mammalian X and Y chromosomes originated from a pair of indistinguishable autosomes ~180-210 million years ago that acquired the sex-determining genes (Charlesworth, 1991; Lahn and Page, 1999; Ross et al., 2005). The human X and Y chromosomes formed in two different segments: a) one that is shared across all mammals called the X-conserved region (XCR) and b) the X-added region (XAR) that is shared across all eutherian animals (Ross et al., 2005). The sex chromosomes, X and Y, previously recombined along their entire lengths, but due to recombination suppression from Y chromosome-specific inversions (Lahn and Page, 1999; Pandey et al., 2013), now only recombine at the tips in the pseudoautosomal regions (PAR) PAR1 and PAR2 (Charlesworth, 1991; Lahn and Page, 1999; Ross et al., 2005). PAR1 is ~2.78 million bases (Mb) and PAR2 is ~0.33 Mb; these sequences are 100% identical between X and Y (Aken et al., 2017; Charchar et al., 2003; Ross et al., 2005) (Figure 1A). The PAR1 is a remnant of the XAR Ross et al. 2005) and shared among eutherians, while the PAR2 is recently added and human-specific (Charchar et al., 2003). Other

5

102 regions of high sequence similarity between X and Y include the X-transposed-region (XTR) with

103 98.78% homology (Veerappa et al., 2013) (Figure 1A). The XTR formed from an X chromosome

104 to Y chromosome duplication event following the human-chimpanzee divergence (Ross et al.,

105 2005; Skaletsky et al., 2003). Thus, the evolution of the X and Y chromosomes has resulted in a

106 pair of chromosomes that are diverged, but still share some regions of high sequence similarity.

107      To infer which genes or transcripts are expressed, RNA-Seq reads can be aligned to a

108 reference genome. The abundance of reads mapped to a transcript is reflective of the amount of

109 expression of that transcript. RNA-Seq methods rely on aligning reads to an available high quality

110 reference genome sequence, but this remains a challenge due to the intrinsic complexity in the

111 transcriptome of regions with a high level of homology (Piskol et al., 2013). By default, the

112 GRCh38 version of the human reference genome includes both the X and Y chromosomes, which

113 is used to align RNA-Seq reads from both male XY and female XX samples. It is known that

114 sequence reads from DNA will misalign along the sex chromosomes affecting downstream

115 analyses (Webster et al., 2019). However, this has not been tested using RNA-Seq data and the

116 effects on differential expression analysis are not known. Considering the increasing number of

117 human RNA-Seq consortium datasets (e.g., the Genotype-Tissue Expression project (GTEx)

118 (GTEx Consortium, 2015), The Cancer Genome Atlas (TCGA) (Cancer Genome Atlas Research

119 Network et al., 2013), Geuvadis project (Lappalainen et al., 2013), and Simons Genome Diversity

120 Project (Mallick et al., 2016)), there is an urgent need to understand how aligning to a default

121 reference genome that includes both X and Y may affect estimates of gene expression on the sex

122 chromosomes (Khramtsova et al., 2018; Tukiainen et al., 2016). We hypothesize that regions of

123 high sequence similarity will result in misaligning of RNA-Seq reads and reduced expression

124 estimates (Figure 1A & B).

6

125    Here, we tested the effect of sex chromosome complement informed read alignment on the

126    quantified levels of gene expression and the ability to detect sex-biased gene expression. We

127    utilized data from the GTEx project, focusing on five tissues, whole blood, brain cortex, breast,

128    liver, and thyroid, which are known to exhibit sex differences in gene expression  (Gershoni and

129    Pietrokovski 2017; R. Li and Singh 2014; de Perrot et al. 2000; Melé et al. 2015; Mayne et al.

130    2016). Many genes have been reported to be differentially expressed between male and female

131    brain samples (Gershoni and Pietrokovski, 2017; Goldstein et al., 2014; Shi et al., 2016) and

132    differential expression in blood samples between males and females has also been documented

133    (Gershoni and Pietrokovski, 2017; Goldstein et al., 2014). An analysis of all GTEx tissue samples

134    reported that breast mammary gland tissues are the most sex differentially expressed tissue

135    (Gershoni and Pietrokovski, 2017). It has also been reported that there are sex disparities in thyroid

136    cancer (Rahbari et al., 2010) and liver cancer (Natri et al., 2019; Naugler et al., 2007) suggesting

137    possible sex differences in gene expression. We used whole blood, brain cortex, breast, liver, and

138    thyroid tissues from 20 genetic male (46, XY) and 20 genetic female (46, XX) individuals for a

139    total of 200 samples evenly distributed among tissues. Male and female samples, for each tissue,

140    were age-matched between the sexes and only included samples of age 55 to 70. We aligned all

141    samples to a default reference genome that includes both the X and Y chromosomes and to a

142    reference genome that is informed on the sex chromosome complement of the genome: Male XY

143    samples were aligned to a reference genome that includes both the X and Y chromosome, where

144    the Y chromosome PAR1 and PAR2 are hard-masked with Ns (Figure 1C) so that reads will align

145    uniquely to the X PAR sequences. Conversely, female XX samples were aligned to a reference

146    genome where the entirety of the Y chromosome is hard-masked (Figure 1C). We tested two

147    different read aligners, HISAT (Kim et al., 2015) and STAR (Dobin et al., 2013), to account for

7

148    variation between alignment methods and measured differential expression using Limma/Voom

149    (Law et al., 2014). We found that using a sex chromosome complement informed reference

150    genome for aligning RNA-Seq reads increased expression estimates on the pseudoautosomal

151    regions of the X chromosome in both male XY and female XX samples and uniquely identified

152    differentially expressed genes.

153          We additionally investigated the effect of transcriptome references on pseudo-alignment

154    methods. We quantified abundance using Salmon (Patro et al., 2017) in male and female brain

155    cortex samples twice, once to a default reference transcriptome index that includes both the X

156    and Y chromosome linked transcripts and to a reference transcriptome index that is informed on

157    the sex chromosome complement of the sample. We found that using a sex chromosome

158    complement informed reference transcriptome index for RNA-Seq pseudo-alignment

159    quantification eliminated Y-linked expression estimates in female XX samples, that were

160    observed in the default approach.

161          Regardless of alignment or pseudo-alignment approach, we recommended carefully

162    considering the annotations of the sex chromosomes in the references used, as theses will affect

163    quantifications and differential expression estimates, especially of sex chromosome linked genes.

164

165    **Methods**

166    *Building sex chromosome complement informed reference genomes*

167    All GRCh38.p10 unmasked genomic DNA sequences, including autosomes 1-22, X, Y,

168    mitochondrial DNA (mtDNA), and contigs were downloaded from ensembl.org release 92 (Aken

169    et al., 2017). The default reference genome here includes all 22 autosomes, mtDNA, the X

170    chromosome, the Y chromosome, and contigs. For the two sex chromosome complement informed

8

171     reference assemblies, we included all 22 autosomes, mtDNA, and contigs from the default

172     reference and a) one with the Y chromosome either hard-masked for the "Y-masked reference

173     genome" or b) one with the pseudoautosomal regions, PAR1 and PAR2, hard-masked on the Y

174     chromosome for "YPARs-masked reference genome" (Figure 1C). Hard-masking with Ns will

175     force reads to not align to those masked regions in the genome. Masking the entire Y chromosome

176     for the sex chromosome complement informed reference genome, Y-masked, was accomplished

177     by changing all the Y chromosome nucleotides [ATGC] to N using a sed command in linux.

178     YPARs-masked was created by hard-masking the Y PAR1: 6001-2699520 and the Y PAR2:

179     154931044-155260560 regions. The GRCh38.p10 Y PAR1 and Y PAR2 chromosome start and

180     end location was defined using Ensembl GRCh38 Y PAR definitions (Aken et al., 2017). After

181     creating the Y chromosome PAR1 and PAR2 masked fasta files, we concatenated all the Y

182     chromosome regions together to create a YPARs-masked reference genome. After creating the

183     GRCh38.p10 default reference genome and the two sex chromosome complement informed

184     reference genomes, we indexed the reference genomes and created a dictionary for each using

185     HISAT version 2.1.0 (Kim et al., 2015) hisat2-build -f option and STAR version 2.5.2 (Dobin et

186     al., 2013), using option --genomeDir and --sjdbGTFfile. Reference genome indexing was followed

187     by picard tools version 1.119 CreateSequenceDictionary (2020), which created a dictionary for

188     each         reference         genome         (Pipeline         available         on         GitHub,

189     https://github.com/SexChrLab/XY_RNAseq).

190

191     *Building sex chromosome complement informed transcriptome index*

192     Ensembl's GRCh38.p10 cDNA reference transcriptome fasta consists of transcript sequences

193     resulting from Ensemble gene predictions. Ensembl's cDNA was downloaded from ensembl.org

9

194    release 92 (Aken et al., 2017). The default transcriptome reference includes 199,234 transcripts

195    which includes autosomal, mtDNA, X chromosome, Y chromosome and contig transcripts. The

196    default Ensembl cDNA does not contain Y chromosome PAR linked transcript sequences, it only

197    contains the X chromosome PAR sequence transcripts. For the sex chromosome complement

198    informed reference transcriptome index, we included all 22 autosomes, mtDNA, X, and contigs

199    from the default cDNA transcriptome and we hard-masked all available Y chromosome linked

200    transcript sequences. Hard-masking the Y chromosome linked transcripts was accomplished by

201    changing all the Y chromosome nucleotides [ATGC] to N using a sed command in linux. After

202    downloading the GRCh38.p10 default reference transcriptome and creating the Y-masked sex

203    chromosome complement informed reference transcriptome fasta files, we then generated a

204    decoy-aware transcriptome for each transcriptome reference. For generating the default decoy-

205    aware reference transcriptome, we used the default genome as the decoy sequence. This was

206    accomplished by concatenating the default genome fasta to the end of the default transcriptome

207    fasta to populate the decoy file with the chromosome names, as suggested by Salmon (Patro et

208    al., 2017). The default transcriptome fasta and the default decoy file were then used to create the

209    mapping-based index using the Salmon version 1.2.0 index function (Patro et al., 2017). The Y-

210    masked decoy-aware transcriptome fasta was generated by concatenating the Y-masked genome

211    fasta to the end of the Y-masked transcriptome fasta to populate the decoy file with the

212    chromosome names. The Y-masked transcriptome fasta and the decoy file were then used as

213    inputs for generating the Y-masked mapping-based index using the salmon index function. For

214    both the default and the Y-masked mapping-based index, a k-mer of 31 was used as this was

215    suggested to work well for reads of 75bp.

216    In addition to the Ensembl reference, we investigated the effects of a sex chromosome

217    complement reference transcriptome index using the gencode transcript reference fasta

218    GRCh38.p12 that contains 206,694 transcripts which includes autosomal, mtDNA, X, Y and

219    contigs. The gencode transcriptome reference includes both the X and Y PAR transcripts (J et al.,

220    2012). Following the same parameters for the Ensembl decoy-aware transcriptome, we created

221    two gencode sex chromosome complement decoy-aware transcriptome references, in addition to

222    a default gencode decoy-aware transcriptome reference. The pipeline is available on GitHub,

223    https://github.com/SexChrLab/XY_RNAseq.

224

225    *RNA-Seq samples*

226    From the Genotyping-Tissue Expression (GTEx) Project data, we downloaded SRA files for whole

227    blood, brain cortex, breast, liver, and thyroid tissues from 20 separate genetic female (46, XX) and

228    20 separate genetic male (46, XY) individuals (Consortium, 2015; GTEx Consortium, 2015) that

229    were age matched between the sexes and ranged from age 55 to 70 (Additional file 1 & 2). Age

230    matching exactly was accomplished using the matchit function in the R package MatchIt  (Ho et

231    al. 2011). The GTEx data is described and available through dbGaP under accession

232    phs000424.v6.p1; we received approval to access this data under dbGaP accession #8834. GTEx

233    RNA-Seq samples were sequenced to 76bp reads and the median coverage was ~82 million (M)

234    reads (Consortium, 2015). Although information about the genetic sex of the samples was provided

235    in the GTEx summary downloads, it was additionally investigated by examining the gene

236    expression of select genes that are known to be differentially expressed between the sexes or are

237    known X-Y homologous genes: *DDX3X*, *DDX3Y*, *PCDH11X*, *PCDH11Y*, *USP9X*, *USP9Y*, *ZFX*,

238    *ZFY*, *UTX*, *UTY*, *XIST*, and *SRY* (Figure 2; Additional file 3 & 4).

239

240  *RNA-Seq trimming and quality filtering*

241  RNA-Seq sample data was converted from sequence read archive (sra) format to the paired-end

242  FASTQ format using the SRA toolkit (Leinonen et al., 2011). Quality of the samples' raw

243  sequencing reads was examined using FastQC (Andrews) and MultiQC . Subsequently, adapter

244  sequences were removed using Trimmomatic version 0.36 (Bolger et al., 2014). More specifically,

245  reads were trimmed to remove bases with a quality score less than 10 for the leading strand and

246  less than 25 for the trailing strand, applying a sliding window of 4 with a mean PHRED quality of

247  30 required in the window and a minimum read length of 40 bases.

248

249  *RNA-Seq read alignment*

250  Following trimming, paired RNA-Seq reads from all samples were aligned to the default reference

251  genome. Unpaired RNA-Seq reads were not used for alignment. Reads from the female (46, XX)

252  samples were aligned to the Y-masked reference genome and reads from male (46, XY) individuals

253  were aligned to the YPARs-masked reference genome. Read alignment was performed using

254  HISAT version 2.1.0 (Kim et al., 2015), keeping all parameters the same, only changing the

255  reference genome used, as described above. Read alignment was additionally performed using

256  STAR version 2.5.2 (Dobin et al., 2013), where all samples were aligned to a default reference

257  genome and to a reference genome informed on the sex chromosome complement, keeping all

258  parameters the same (Pipeline available on GitHub, https://github.com/SexChrLab/XY_RNAseq).

259

260  *Processing of RNA-Seq alignment files*

12

261     Aligned RNA-Seq samples from HISAT and STAR were output in Sequence Alignment Map

262     (SAM) format and converted to Binary Alignment Map (BAM) format using bamtools version

263     2.4.0 (Li et al., 2009). Summaries on the BAM files including the number of reads mapped were

264     computed using bamtools version 2.4.0 package (Barnett et al., 2011). RNA-Seq BAM files were

265     indexed, sorted, duplicates were marked, and read groups added using bamtools, samtools, and

266     Picard (Barnett et al., 2011; Li et al., 2009, 2020). All RNA-Seq BAM files were indexed using

267     the default reference genome using Picard ReorderSam (2020), this was done so that all samples

268     would include all chromosomes in the index files. Aligning XX samples to a Y-masked reference

269     genome using HISAT indexes would result in no Y chromosome information in the aligned BAM

270     and BAM index bai files. For downstream analysis, some tools require that all samples have the

271     same chromosomes, which is why we hard-masked rather than removed. Reindexing the BAM

272     files to the default reference genome does not alter the read alignment, and thus does not alter our

273     comparison between default and sex chromosome complement informed alignment.

274

275     *Gene expression level quantification*

276     Read counts for each gene across all autosomes, sex chromosomes, mtDNA, and contigs were

277     generated using featureCounts version 1.5.2 (Liao et al., 2014) for all aligned and processed RNA-

278     Seq BAM files. Female XX samples when aligned to a sex chromosome complement informed

279     reference genome will show zero counts for Y-linked genes, but will still include those genes in

280     the raw counts file. This is an essential step for downstream differential expression analysis

281     between males and females to keep the total genes the same between the sexes for comparison.

282     Only rows that matched gene feature type in Ensembl Homo_sapiens.GRCh38.89.gtf gene

283     annotation (Aken et al., 2017) were included for read counting. There are 2,283 genes annotated

284    on the X chromosome and a total of 56,571 genes across the entire genome for GRCh38 version

285    of the human reference genome (Aken et al., 2017). Only primary alignments were counted and

286    specified using the --primary option in featureCounts.

287

288    *RNA-seq quantification for transcriptome index*

289    Transcript quantification for trimmed paired RNA-seq brain cortex samples were estimated twice,

290    once to a default decoy-aware reference transcriptome index and once to a sex chromosome

291    complement informed decoy-aware reference transcriptome index using Salmon with the –

292    validateMappings flag. Salmon's –validateMappings adopts a scheme for finding protentional

293    mapping loci of a read using a chain algorithm introduced in minimap2 (Li, 2018).  Transcript

294    quantification for female (46, XX) samples was estimated using a Y-masked reference

295    transcriptome index and male (46, XY) transcript quantification was estimated using a Y PAR

296    masked reference transcriptome index when the Y PAR sequence information was available for

297    the transcriptome build. This was repeated for both the Ensembl and the gencode cDNA

298    transcriptome builds, keeping all parameters the same, only changing the reference transcriptome

299    index used, as described above.

300

301    *DGEList object*

302    Differential expression analysis was performed using the limma/voom pipeline (Law et al., 2014)

303    which has been shown to be a robust differential expression software package (Costa-Silva et al.,

304    2017; Seyednasrollah et al., 2015) for both reference-based and pseudo-alignment quantification.

305    Quantified read counts from each sample for the reference-based quantification were generated

306    from featureCounts were combined into a count matrix, each row representing a unique gene ID

307    and each column representing the gene counts for each unique sample. This was repeated for each

308    tissue type and read into R using the DGEList function in the R limma package  (Love et al., 2014).

309    A sample-level information file related to the genetic sex of the sample, male or female, and the

310    reference genome used for alignment, default or sex chromosome complement informed, was

311    created and corresponds to the columns of the count matrix described above.

312        Pseudo-aligned transcript read counts from each brain cortex sample quantified using

313    Salmon were combined into a count matrix using tximport (Soneson et al., 2015) with each row

314    representing a unique transcript ID and each column representing the transcript counts for each

315    unique sample. To create length scaled transcripts per million (TPM) values to pass into limma,

316    tximport function lengthScaledTPM was employed (Soneson et al., 2015). The reference assembly

317    annotation file was read into R using tximport function makeTxDbFromGFF. Following this, a

318    key of the transcript ID corresponding to the gene ID was created was created using the keys

319    function (Soneson et al., 2015). Gene level TPM values were then generated using the tx2gene

320    function. This was repeated for the Ensembl and the gencode default and sex chromosome

321    complement informed transcriptome quantification estimates.

322

323    *Multidimensional Scaling*

324    Multidimensional Scaling (MDS) was performed using the DGEList-object containing gene

325    expression count information for each sample. MDS plots were generated using the plotMDS

326    function in in the R limma package (Law et al., 2014). The distance between each pair of samples

327    is shown as the $\log_2$ fold change between the samples. The analysis was done for each tissue

328    separately using all shared common variable genes for dimensions (dim) 1 & 2 and dim 2 & 3.

329    Samples that did not cluster with reported sex or clustered in unexpected ways in either dim1, 2 or

330     3 were removed from all downstream analysis (Additional file 5). MDS plots for each tissue

331     containing the samples that were used for quality control are located in Additional file 6. Briefly,

332     one male XY whole blood did not cluster with any of the other samples and was removed. One

333     female XX breast sample clustered with the opposite sex and was thus removed. In brain cortex,

334     three male XY brain cortex samples didn't cluster neatly with the other male XY samples in dim

335     1 & 2 were thus removed. Another male brain cortex sample, although clustered with other male

336     samples, had the lowest number of sequencing remaining after trimming for quality, 23.9M, and

337     thus was also removed. To keep the number of samples in each sex roughly equal, four female XX

338     brain cortex samples were randomly selected for removal. For liver and thyroid tissue, no samples

339     appeared to cluster in any unexpected ways and thus no liver or thyroid tissue samples were

340     removed. For all aligners the first component of variation in the MDS plot is explained by the sex

341     of the sample (Figure 3).

342

343     *Differential expression*

344     Using edgeR (Robinson et al., 2010), raw counts were normalized to adjust for compositional

345     differences between the RNA-Seq libraries using the voom normalize quantile function, which

346     normalizes the reads by the method of trimmed mean of values (TMM) (Law et al., 2014). Counts

347     were then transformed to $\log_2(\text{CPM}+0.25/\text{L})$, where CPM is counts per million, L is library size,

348     and 0.25 is a prior count to avoid taking the log of zero (Robinson et al., 2010). For this dataset,

349     the average library size is about 79.76 million, therefore L is 79.76. Thus, the minimum

350     $\log_2(\text{CPM}+0.25/\text{L})$ value for each sample, representing zero transcripts, is $\log_2(0+0.25/15) = -8.32$.

351         A mean minimum of 1 CPM, or the equivalent of 0 in $\log_2(\text{CPM}+2/\text{L})$, in at least one sex

352     per tissue comparison was required for the gene to be kept for downstream analysis. A CPM value

16

353    of 1 was used in our analysis to separate expressed genes from unexpressed genes, meaning that

354    in a library size of ~79.76 million reads, there are at least an average of 79 counts in at least one

355    sex. After filtering for a minimum CPM, 53,804 out of the 56,571 quantified genes were retained

356    for the whole blood samples, 53,822 for brain cortex, 54,184 for breast, 53,830 for liver, and

357    53,848 for thyroid. A linear model was fitted to the DGEList-object, which contains the filtered

358    and normalized gene counts for each sample, using the limma lmfit function which will fit a

359    separate model to the expression values for each gene (Law et al., 2014).

360        For differential expression analysis a design matrix containing the genetic sex of the sample

361    (male or female) and which reference genome the sample was aligned to (default or sex

362    chromosome complement informed) was created for each tissue type for contrasts of pairwise

363    comparisons between the sexes. Pairwise contrasts were generated using limma makecontrasts

364    function (Law et al., 2014). We identified genes that exhibited significant expression differences

365    defined using an Benjamini-Hochberg adjusted p-value cutoff that is less than 0.01 (1%) to account

366    for multiple testing in pairwise comparisons between conditions using limma/voom decideTests

367    vebayesfit (Law et al., 2014). A conservative adjusted p-value cutoff of less than 0.01 was chosen

368    to be highly confident in the genes that were called as differentially expressed when comparing

369    between reference genomes used for alignment.    Pipeline available on GitHub,

370    https://github.com/SexChrLab/XY_RNAseq.

371

372    *GO analysis*

373    We examined differences and similarities in gene enrichment terms between the differentially

374    expressed genes obtained from the differential expression analyses of the samples aligned to the

375    default and sex chromosome complement informed reference genomes, to investigate if the

17

376    biological interpretation would change depending on the reference genome the samples were

377    aligned to. We investigated gene ontology enrichment for lists of genes that were identified as

378    showing overexpression in one sex versus the other sex for whole blood, brain cortex, breast, liver,

379    and thyroid samples (adjusted p-value < 0.01). We used the GOrilla webtool, which utilizes a

380    hypergeometric distribution to identify enriched GO terms (Eden et al., 2009). A modified Fisher

381    exact p-value cutoff < 0.001 was used to select significantly enriched terms (Eden et al., 2009).

382

383    **Results**

384    *RNA-Seq reads aligned to autosomes do not vary much between reference genomes*

385    We compared total mapped reads when reads were aligned to a default reference genome and to a

386    reference genome informed on the sex chromosome complement. Reads mapped across the whole

387    genome, including the sex chromosomes, decreased when samples were aligned to a reference

388    genome informed on the sex chromosome complement, paired t-test p-value < 0.05 (Additional

389    files 7 - 9). This was true regardless of the read aligner used, HISAT or STAR, or of the sex of the

390    sample, XY or XX. To test the effects of realignment on an autosome, we selected chromosome

391    8, because of its similar size to chromosome X. Overall, there is a slight mean increase in reads

392    mapping to chromosome 8 when samples are aligned to a sex chromosome complement informed

393    reference genome compared to aligning to a default reference genome (Additional file 9). For

394    female XX samples, the mean increase in reads mapping for chromosome 8 was 42.2 reads for

395    whole blood, 50.25 for brain cortex, 109.9 for breast, 68.5 for liver, and 98.2 for thyroid

396    (Additional file 9), which was significant using a paired t-test, p-value < 0.05 in all tissues

397    (Additional file 9). Male XY samples also showed a mean increase in reads mapping for

398    chromosome 8. The mean increase in reads mapping to chromosome 8 for male whole blood

399      samples was 0.84, 2.38 for brain cortex, 5.58 for breast, 3.2 for liver, and 5 for thyroid (Additional

400      file 9). There was a significant increase, p-value < 0.05 paired t-test, for reads mapping to

401      chromosome 8 for male brain cortex, breast, liver, and thyroid samples. There was no significant

402      increase in reads mapping for male whole blood for chromosome 8 (Additional file 9).

403

404      *Reads aligned to the X chromosome increase in both XX and XY samples when using a sex*

405      *chromosome complement informed reference genome*

406      We found that when reads were aligned to a reference genome informed by the sex chromosome

407      complement for both male XY and female XX tissue samples, reads on the X chromosome

408      increased by ~0.12% when aligned using HISAT. For all tissues and both sexes we observe an

409      average increase of 1,991 reads on chromosome X. We observe an increase in reads mapping to

410      the X chromosome for all tissues and for each sex, which was significant using a paired t-test, p-

411      value < 0.05 (Additional file 9). Reads on the Y chromosome decreased 100% (67,033 reads on

412      average) across all female XX samples and by ~57.32% (69,947 reads on average) across all male

413      XY samples when aligned using HISAT (Additional file 7 & 9). Similar increases in X

414      chromosome and decreases in Y chromosome reads when aligned to a sex chromosome

415      complement informed reference were observed when STAR was used as the read aligner for both

416      male XY and female XX samples (Additional file 8 & 9).

417

418      *Aligning to a sex chromosome complement informed reference genome increases the X*

419      *chromosome PAR1 and PAR2 expression*

420      We next explored the effect of changes in read alignment on gene expression. There was an

421      increase in pseudoautosomal regions, PAR1 and PAR2, expression when reads were aligned to a

422    reference genome informed on the sex chromosome complement for both male XY and female

423    XX samples (Additional file 10 & 11). We found an average of 2.73 $\log_2$ fold increase in

424    expression in PAR1 expression for female XX brain cortex samples and 2.75 $\log_2$ fold increase in

425    expression in PAR1 for male XY brain cortex samples using HISAT (Figure 4). The X-transposed

426    region, XTR, in female XX brain cortex samples showed a 1.22 $\log_2$ fold increase in expression

427    and no change in male XY brain cortex samples. PAR2 showed an average of 2.13 $\log_2$ fold

428    increase for female XX brain cortex samples and 2.19 $\log_2$ fold increase in PAR2 for male XY

429    brain cortex samples using HISAT, with similar results for STAR read aligner (Additional file 10

430    & 11). Complete lists of the $\log_2$(CPM+0.25/L) values for each X chromosomal gene and each

431    gene within the whole genome for male XY and female XX samples are in Additional file 12

432    available on Dryad for download under  https://doi.org/10.5061/dryad.xksn02vbv.

433

434    *Regions outside the PARs and XTR show little difference in expression between reference genomes*

435    Intriguingly, regions outside the PARs on the X chromosome were less affected by the choice of

436    reference genome. Across the entire X-conserved region, we observed practically no change in

437    estimates of gene expression between the default and sex chromosome complement informed

438    references (e.g., a 0.99 $\log_2$ fold in male thyroid samples, and 1.00  $\log_2$ fold change in female

439    brain cortex samples, essentially showing no difference (Additional file 10 & 11)). Additionally,

440    X and Y homologous genes (*AMELX*, *ARSD*, *ARSE*, *ARSF*, *CASK*, *GYG2*, *HSFX1*, *HSFX2*,

441    *NLGN4X*, *OFD1*, *PCDH11X*, *PRKX*, *RBMX*, *RPS4X*, *SOX3*, *STS*, *TBL1X*, *TGIF2LX*, *TMSB4X*,

442    *TSPYL2*, *USP9X*, *VCX*, *VCX2*, *VCX3A*, *VCX3B*, *ZFX*) showed little to no increase in expression

443    when aligned to a sex chromosome complement informed reference genome compared to aligning

444    to a default reference genome (Additional file 13). *PCDH11X* showed the highest increase in

445    expression for all tissues regardless of read aligner. The log$_2$ fold increase in expression for

446    *PCDH11X* for female samples when aligned using HISAT was 0.4, 0.28, 0.33, 0.16, and 0.16 for

447    whole blood, brain cortex, breast, liver, and thyroid, respectively. Other X and Y homologous

448    genes sometimes increased in expression depending on the tissue and sometimes there was no

449    change in expression (Additional file 13). Next to *PCDH11X*, the most increase in expression in

450    an X and Y homologous genes was *VCX3B, NLGN4X,* and *VCX3A. NLGN4X* in whole blood

451    showed a 0.14 log$_2$ fold increase when aligned using HISAT. *VCX3B* showed a 0.2 log$_2$ fold

452    increase in brain, *NLGN4X* showed a 0.04 log$_2$ fold increase in breast, *VCX3A* showed a 0.07 log$_2$

453    fold increase in liver, and *VCX3B* showed a 0.04 log$_2$ fold increase in thyroid, when aligned using

454    HISAT (Additional file 13).

455

456    *A sex chromosome complement informed reference genome increases the ability to detect sex*

457    *differences in gene expression*

458    We next investigated how this would affect gene differential expression between the sexes.

459    Generally, we find that more genes are differentially expressed on the sex chromosomes between

460    the sexes when the sex chromosome complements are taken into account. The number of

461    differentially expressed genes on the autosomes remained the same or increased. At a conservative

462    Benjamini-Hochberg adjusted p-value of $< 0.01$ and aligning with HISAT, we find 4 new genes

463    (3 Y-linked and 1 X-linked) that are only called as differentially expressed between the sexes in

464    the brain cortex when aligned to reference genomes informed on the sex chromosome complement

465    (Figure 5; Additional file 14). We observed similar trends in changes for differential expression

466    between male XY and female XX for whole blood, breast, liver, and thyroid samples using either

467    HISAT or STAR as the aligner (Additional file 14). For example, in whole blood, 3 additional

21

468    genes are called as being differentially expressed between the sexes using HISAT, while 1

469    additional gene is called differentially expressed when aligned using STAR. Additionally, when

470    taking sex chromosome complement into account, the number of genes called as differentially

471    expressed between the sexes for the breast samples increased by 13 genes (8 autosomal, 3 X-linked

472    and 2 Y-linked) using HISAT and by 8 genes using STAR (6 autosomal and 2 X-linked)

473    (Additional file 14 & 15). For all tissues, no genes were uniquely called as being differentially

474    expressed between the sexes when aligned to a default reference genome compared to a reference

475    genome informed on the sex chromosome complement (Additional file 14 & 15). Rather, only

476    when samples were aligned to a sex chromosome complement did we observe an increase in the

477    genes called as being differentially expressed (Figure 5; Additional file 14 & 15).

478

479    *Increase in gene enrichment pathways when samples are aligned to a sex chromosome complement*

480    *informed reference genome*

481    A sex chromosome complement informed reference genome increases the ability to detect genes

482    as differentially expressed between the sexes and thus alters gene enrichment results. When the

483    thyroid samples were aligned using a sex chromosome complement informed reference genome

484    using HISAT, genes up-regulated in male XY samples still show enrichment for positive

485    regulation of transcription from RNA polymerase II (found when aligning to a default reference

486    genome), but additionally find postsynaptic membrane assembly, postsynaptic membrane

487    organization, and vocalization behavior (Additional file 16). These additional GO enrichments in

488    the male XY thyroid samples involve *NRXN1* and *NLGN4Y* genes, both of these genes are located

489    on the Y chromosome. GO enrichment analysis of genes that are more highly expressed in female

490    liver compared to male liver samples, when samples were aligned to a default reference genome

491 using HISAT, were genes involved in modification histone lysine demethylation (Additional file

492 16). However, when these samples were aligned to a sex chromosome complement informed

493 reference genome, genes upregulated in females were enriched for histone lysine demethylation

494 as well as negative regulation of endopeptidase activity, negative regulation of peptidase activity,

495 cytoplasmic actin-based contraction involved in cell motility (Additional file 16). These additional

496 GO enrichments in the female XX liver samples include the involvement of *KDM6A, DDX3X,* and

497 *VIL1. KDM6A, DDX3X* are X-linked and *VIL1* is on chromosome 2. Whole blood, brain cortex,

498 male liver, and female thyroid samples showed no difference in GO enrichment pathways when

499 using a default reference genome compared to a sex chromosome complement reference genome

500 for alignment when using HISAT with similar results for STAR as the read aligner (Additional

501 file 17). Thus, while there won't always be a difference, aligning to a sex chromosome complement

502 informed reference genome can increase ability to detect enriched pathways.

503

504 *Using sex-linked genes alone is inefficient for determining the sex chromosome complement of a*

505 *sample*

506 The sex of each sample used in this analysis was provided in the GTEx manifest. We investigated

507 the expression of genes that could be used to infer the sex of the sample. We studied X and Y

508 homologous genes (*DDX3X/Y*, *PCDH11X/Y*, *USP9X/Y*, *ZFX/Y*, *UTX/Y*), *XIST*, and *SRY* gene

509 expression in male and female whole blood, brain cortex, breast, liver, and thyroid (Figure 2;

510 Additional file 3 & 4). Both males and females are expected to show expression for the X-linked

511 homologs, whereas only XY samples should show expression of the Y-linked homologs. Further,

512 *XIST* expression should only be observed in XX samples and *SRY* should only be expressed in

513 samples with a Y chromosome. Using the default reference genome for aligning samples, we

514    observed a small number of reads aligning to the Y-linked genes in female XX samples, but also

515    observed clustering by sex for *DDX3Y*, *USP9Y*, *ZFY*, and *UTY* gene expression (Figure 2). Male

516    XY samples showed expression for *DDX3X*, *DDX3Y*, *USP9X*, *ZFX*, and *UTX* (greater than 5

517    $\log_2$(CPM+2/L). Female XX samples showed expression for *XIST* (greater than 4.0

518    $\log_2$(CPM+2/L) and male XY samples showed little to no expression for *XIST* (less than 0

519    $\log_2$(CPM+2/L) with the exception of 2 male whole blood samples and 1 male liver sample, which

520    showed greater than 5 $\log_2$(CPM+2/L) expression). In contrast to the default reference genome,

521    when aligned to a sex chromosome complement informed reference genome, samples cluster more

522    distinctly by sex for *DDX3Y*, *USP9Y*, ZFY, and *UTY*, all showing at least a 4 $\log_2$(CPM+2/L)

523    difference between the sexes (Figure 2; Additional file 3 & 4). *SRY* is predominantly expressed in

524    the testis (Albrecht et al., 2003; Turner et al., 2011) and typically one would expect *SRY* to show

525    male-specific expression. In our set, we did not observe *SRY* expressed in any sample, and so it

526    could not be used to differentiate between XX and XY samples (Figure 2, Additional file 3 & 14).

527    In contrast, the X-linked gene *XIST* was differentially expressed between genetic males and genetic

528    females in both genome alignments (default and sex chromosome complement informed) for the

529    whole blood, brain cortex, breast, liver, and thyroid samples with the exception of 3 male XY

530    samples. *XIST* expression is important in the X chromosome inactivation process (Carrel and

531    Willard, 2005) and serves to distinguish samples with one X chromosome from those with more

532    than one X chromosome (Tukiainen et al., 2016). However, this does not inform about whether

533    the sample has a Y chromosome or not. For X-Y homologous genes, we do not find sex differences

534    in read alignment with either default or sex chromosome complement informed for the X-linked

535    homolog. When aligned to a default reference genome, female XX samples showed some

24

536  expression for homologous Y-linked genes, but only presence/absence of Y-linked reads alone is

537  insufficient to determine sex chromosome complement of the sample (Figure 2, Additional file 3).

538

539  *No Y-linked transcript expression in female XX samples when quantification was estimated using*

540  *a transcriptome index informed on the sex chromosome complement*

541  A pseudo-alignment shows similar effects of the reference to that of an alignment approach (Figure

542  5, Additional files 18 & 19). We observed no Y-linked expression in female XX samples when

543  transcript quantification was estimated using a Y-masked sex chromosome complement reference

544  transcriptome index. This was true for both the Ensembl and gencode pseudo-alignment with a sex

545  chromosome complement reference transcriptome index (Additional files 18 & 19). Interestingly,

546  there was a large difference between the Ensembl and gencode reference files. The transcript IDs

547  in the transcriptome cDNA fasta and the transcript IDs in the annotation file are not one-to-one for

548  the Ensembl assembly (Zhao and Zhang, 2015). There are 190,432 transcript sequences in the

549  Ensembl cDNA fasta file but there are 199,234 transcripts in the Ensembl annotation file. Notably,

550  Ensembl's cDNA reference transcriptome fastas does not contain known transcripts such as the

551  XIST transcripts (Eyras et al., 2004). The Ensembl reference transcriptome fasta also does not

552  contain the Y PARs transcript sequences, it only contains the X PAR transcript sequences. In

553  contrast, the gencode cDNA reference transcriptome fasta and annotation file both contain 206,694

554  sequences, including the Y PARs. Regardless of using an Ensembl or gencode transcriptome,

555  female XX sample show Y-linked expression when using a default refence transcriptome index

556  for pseudo-alignment, however the changes necessary for making a sex chromosome complement

557  informed reference are different for the two builds.

558

559 **Discussion**

560 For accuracy, the sex chromosome complement of the sample should be taken into account when

561 aligning RNA-Seq reads to reduce misaligning sequences. Neither Ensembl or Gencode human

562 reference genomes are correct for aligning both XX and XY samples. The Ensembl GRCh38

563 human reference genome includes all 22 autosomes, mtDNA, the X chromosome, the Y

564 chromosome with the Y PARs masked, and contigs (Aken et al., 2017). The Gencode hg19 human

565 reference genome includes everything with no sequences masked (Harrow et al., 2012).

566 Measurements of X chromosome expression increase for both male XY and female XX

567 whole blood, brain cortex, breast, liver, and thyroid samples when aligned to a sex chromosome

568 complement informed reference genome versus aligning to a default reference genome (Figure 4).

569 While we see increases in measured expression for PAR1 and PAR2 genes in both males and

570 females, we only observe a difference in measured XTR expression in females. This is because

571 while the PARs are 100% identical between the X and Y and so one copy (here we mask the Y-

572 linked copy) should be masked, the XTR is not hard-masked in the YPARs-masked reference

573 genome. The XTR is not identical between the X and Y; it shares 98.78% homology between X

574 and Y but no longer recombines between X and Y (Veerappa et al., 2013) (Figure 1A) and because

575 of this divergence, is therefore not hard-masked when aligning male XY samples. Tukiainen et al.,

576 (2016) and others have shown that PAR1 genes have a male bias in expression (Tukiainen et al.,

577 2016). Our findings here support this regardless if the samples were aligned to a default or a sex

578 chromosome complement reference genome (Additional file 11 & 12). Differential expression

579 results changed when using a sex chromosome complement informed alignment compared to using

580 a default alignment. When aligned to a default reference genome, due to sequence similarity, some

581 reads from female XX samples aligned to the Y chromosome (Figure 2; Figure 5). However, when

26

582  aligned to a reference genome informed by the sex chromosome complement, female XX samples

583  no longer showed Y-linked gene expression, and more Y-linked genes were called as being

584  differentially expressed between the sexes (Figure 2; Figure 5; Additional file 12 & 15). This

585  suggests that if using a default reference genome for aligning RNA-Seq reads, one would miss

586  some Y-linked genes as differentially expressed between the sexes (Figure 5). Furthermore, these

587  Y-linked genes serve in various important biological processes, thus altering the functional

588  interpretation of the sex differences (Additional file 16 & 17). Only when samples were aligned to

589  a sex chromosome complement reference genome did we observe more genes called as

590  differentially expressed between the sexes (Additional file 14). An increase in genes called

591  differentially expressed additionally alters the GO analysis results (Additional file 16 & 17). When

592  samples were aligned to a default reference genome we sometimes missed GO pathways or

593  misinterpreted which were the top pathways.

594      The choice of read aligner has long been known to give slightly differing results of

595  differential expression due to the differences in the alignment algorithms (Conesa et al., 2016;

596  Costa-Silva et al., 2017). Differences between HISAT and STAR could be contributed to

597  differences in default parameters for handling multi-aligning reads (Kim et al., 2015). We show

598  that regardless of choice of read aligner, HISAT or STAR, we observe similar results. Sample size

599  has also long been known to alter differential expression analysis (Ching et al., 2014; Lamarre et

600  al., 2018; Zhao et al., 2018). We therefore additionally replicated our findings in a smaller sample

601  size of 3 male XY compared to 3 female XX samples for whole blood and brain cortex tissue and

602  where the samples were randomly selected and confirmed the results from the larger sample size

603  (Additional file 20).

27

604    In addition to reference-based quantification, we tested whether quantifying sex-linked

605    reads with a pseudo-aligner would be affected by using a sex chromosome complement reference.

606    Previous studies have shown that reference-based alignment is not necessary for high-quality

607    estimation of transcript levels (Zielezinski et al., 2017). However, we observed expression

608    estimates for Y-linked transcripts in female XX samples when using a default reference

609    transcriptome index for pseudo-alignment quantification estimates. In contrast, when a sex

610    chromosome complement informed reference transcriptome index was used, we observed no Y-

611    linked expression in female XX samples. Salmon, and other alignment-free tools such as Kallisto

612    (Bray et al., 2015) and Sailfish (R et al., 2014), build an index of k-mers from a reference

613    transcriptome. The k-mer transcriptome index is used to group pseudoalignments belonging to the

614    same set of transcripts to directly estimate the expression of each transcript. A k-mer alignment

615    free approach is faster and less demanding than alignment protocols (Zielezinski et al., 2017);

616    however, a sex chromosome complement informed transcriptome index should be carefully

617    considered because even a k-mer approach is not sensitive to regions that are 100% identical in

618    sequence. Additionally, alignment-free methods are not as robust in quantifying expression

619    estimates for small RNAs and lowly-expressed genes (Wu et al., 2018).

620    The choice of reference transcriptome or reference genome can also give slightly differing

621    results of differential expression due to the difference in which transcripts are included in the

622    transcriptome (Zhao and Zhang, 2015). The Ensembl cDNA does not include the Y PAR linked

623    transcripts whereas the gencode transcriptome fasta includes both the X and Y PARs. The Ensembl

624    transcriptome does not include non-coding RNAs, such as *XIST* transcripts. The *XIST* gene is

625    called as being up-regulated in the female XX samples for all tissues and all comparisons except

626    for when transcript expression was estimated using the Ensembl reference transcriptome

627    (Additional file 15, 18, & 19). Given the current builds, for RNA-seq projects interested in sex

628    chromosome linked transcript expression, we suggest that researchers use a gencode sex

629    chromosome complement informed reference transcriptome index.

630        Ideally, one would use DNA to confirm presence or absence of the Y chromosome, but if

631    DNA sequence was not generated, one would need to confirm the genetic sex of the sample by

632    assessing expression estimates for X-linked and Y-linked genes. To more carefully investigate the

633    ability to use gene expression to infer sex chromosome complement of the sample, we examined

634    the gene expression for a select set of X-Y homologous genes, as well as *XIST* and *SRY* that are

635    known to be differentially expressed between the sexes (Figure 2, Additional file 13). The samples

636    broadly segregated by sex for Y-linked gene expression using default alignment. However, the

637    pattern was messy for each individual Y-linked gene. Thus, if inferring sex from RNA-Seq data,

638    we recommend using the estimated expression of multiple X-Y homologous genes and *XIST* to

639    infer the genetic sex of the sample. Samples should be aligned to a default reference genome first

640    to look at the expression for several Y-specific genes to determine if the sample is XY or XX.

641    Then samples should be realigned to the appropriate sex chromosome complement informed

642    reference genome. Independently assessing sex chromosome complement of samples becomes

643    increasingly important as karyotypically XY individuals are known to have lost the Y chromosome

644    in particular tissues sampled, as shown in Alzheimer Disease (Dumanski et al., 2016), age-related

645    macular degeneration (Grassmann et al., 2019), and in the blood of aging individuals (Forsberg,

646    2017), but should not have *XIST* expression. However, *XIST* may not be a sufficient marker alone

647    to infer sex chromosome complement, especially in cancer in samples from XX individuals, where

648    the inactive X can become reactivated (Chaligné et al., 2015). Self-reported sex may not match the

649    sex chromosome complement of the samples, even in karyotypic individuals.

29

650

## Conclusion

652 Here we show that aligning RNA-Seq reads to a sex chromosome complement informed reference

653 genome will change the results of the analysis compared to aligning reads to a default reference

654 genome. We previously observed that a sex chromosome complement informed alignment is

655 important for DNA as well (Webster et al., 2019). A sex chromosome complement informed

656 approach is needed for a sensitive and specific analysis of gene expression on the sex chromosomes

657 (Khramtsova et al., 2018). A sex chromosome complement informed reference alignment resulted

658 in increased expression of the PARs of the X chromosome for both male XY and female XX

659 samples. We further found different genes called as differentially expressed between the sexes and

660 identified sex differences in gene pathways that were missed when samples were aligned to a

661 default reference genome.

## Perspectives and Significance

663 The accurate alignment and pseudo-alignment of the short RNA-Seq reads to the reference genome

664 or reference transcriptome is essential for drawing reliable conclusions from differential

665 expression data analysis on the sex chromosomes. We strongly urge studies using RNA-Seq to

666 carefully consider the genetic sex of the sample when quantifying reads, and provide a framework

667 for doing so in the future (https://github.com/SexChrLab/XY_RNAseq).

**Funding**

**Author contributions**

KCO: Supervision, Formal Analysis, Investigation, Visualization, Writing - Original Draft Preparation, Writing - Review and Editing

SMB: Formal Analysis, Investigation, Writing - Original Draft Preparation, Writing - Review and Editing

JPA: Formal Analysis, Investigation, Writing - Review and Editing

VAVV: Investigation, Writing - Review and Editing

MAW: Conceptualization, Supervision, Visualization, Resources, Project Administration, Writing - Original Draft Preparation, Writing - Review and Editing, Funding Acquisition

**Acknowledgements**

**Competing Interests**

691     The authors declare no competing interests.

692

693     **Availability of Data and Material**

694     The RNA-Seq datasets analyzed during the current study are available from the GTEx project

695     through dbGaP under accession phs000424.v6.p1; we received approval to access this data under

696     dbGaP accession #8834. All codes used are available on GitHub:

697     https://github.com/SexChrLab/XY_RNAseq.

698

699     **Ethics Approval and Consent to Participate**

700     Not applicable.

701

702     **Consent for Publication**

703     Not applicable.

704 **Figure Legends**

705

706 **Figure 1. Homology between the human X and Y chromosomes where misaligning could**

707 **occur. A)** High sequence homology exists between the human X and Y chromosomes in three

708 regions: 100% sequence identity for the pseudoautosomal regions (PARs), PAR1 and PAR2, and

709 ~99% sequence homology in the X-transposed region (XTR). The X chromosome PAR1 is ~2.78

710 million bases (Mb) extending from X:10,001 to 2,781,479 and the X chromosome PAR2 is ~0.33

711 Mb extending from X:155,701,383 to 156,030,895. The X chromosome PAR1 and PAR2 are

712 identical in sequence to the Y chromosome PAR1 Y:10,001 - 2,781,479 and PAR2 Y:56,887,903

713 - 57,217,415. **B)** Using a standard alignment approach will result in reads misaligning between

714 regions of high sequence homology on the sex chromosomes. **C)** Using a reference genome that

715 is informed by the genetic sex of the sample may help to reduce misaligning between the X and Y

716 chromosomes. In humans, samples without evidence of a Y chromosome should be aligned to a

717 Y-masked reference genome and samples with evidence of a Y should be aligned to a YPARs-

718 masked reference genome.

719

720 **Figure 2. Genetic sex of RNA-Seq samples.** We investigated gene expression,

721 $\log_2$(CPM+0.25/L), of XY homologous genes (*DDX3X/Y, PCDH11X/Y, USP9X/Y, ZFX/Y,*

722 *UTX/Y*), and *XIST*, and *SRY* in all samples from all tissues analyzed here from genetic males (blue

723 squares) and genetic females (orange circles) **A)** when aligned to a default reference genome, and

724 **B)** when aligned to a sex chromosome complement informed reference genome, using HISAT as

725 the read aligner.

726

33

727 **Figure 3. Multidimensional scaling for the top 100 most variable genes.** We investigated

728 multidimensional scaling for the top 100 most common variable genes in brain cortex samples. **A)**

729 Salmon pseudo-alignment with Ensembl transcriptome reference **B)** HISAT read aligner and **C)**

730 STAR read aligner when quantifying using both the default and the sex chromosome complement

731 informed reference. The most variation in the data is explained by the sex of the sample.

732

733 **Figure 4. X chromosome RNA-Seq alignment differences in brain cortex.** We plot $\log_2$ fold

734 change (FC) across **A)** the entire X chromosome and **B)** the first 5 million bases (Mb) and show

735 **C)** average fold change in large genomic regions on the X chromosome between aligning brain

736 cortex using HISAT to the default genome and aligning to a sex chromosome complement

737 informed reference genome. For $\log_2$ FC, a value less than zero indicates that the gene showed

738 higher expression when aligned to a default reference genome, while values above zero indicate

739 that the gene shows higher expression when aligned to a reference genome informed by the sex

740 chromosome complement of the sample. Samples from genetic females are plotted in orange

741 circles, while samples from males are plotted in blue squares. Darker shades indicate which gene

742 points are in PAR1, XTR, and PAR2 while lighter shades are used for genes outside of those

743 regions.

744

745 **Figure 5. Sex chromosome complement informed alignment calls more sex-linked genes as**

746 **being differentially expressed. A)** Sex differences in gene expression, $\log_2(\text{CPM}+0.25/\text{L})$,

747 between the twenty samples from genetic males and females are shown when aligning all samples

748 to the default reference genome (left) and a reference genome informed on the sex chromosome

749 complement (right) for brain cortex. Each point represents a gene. Genes that are differentially

34

750    expressed, adjusted p-value < 0.01 are indicated in black for autosomal genes, blue for Y-linked

751    genes, and red for X-linked genes. **B)** We show overlap between genes that are called as

752    differentially expressed when all samples are aligned to the default genome, and genes that are

753    called as differentially expressed when aligned to a sex chromosome complement informed

754    genome. When samples were aligned to a reference genome informed on the sex chromosome

755    complement, 27 genes were called as differentially expressed between the sexes, of which 4 were

756    uniquely called in the sex chromosome complement informed alignment. There were no genes that

757    were uniquely called as differentially expressed when aligned to a default reference genome.

758

759    **Additional files**

760

761    **Additional file 1. Sample IDs.** RNA-Seq whole blood, brain cortex, breast, liver, and thyroid

762    tissue samples from 20 genetic female (46, XX) and 20 genetic male (46, XY) individuals were

763    downloaded from the Genotype-Tissue Expression (GTEx) project (GTEx Consortium, 2015) for

764    a total of 200 RNA-Seq tissue samples.

765

766    **Additional file 2. Histogram of sample reported age.** For each tissue, whole blood, brain cortex,

767    breast, liver, and thyroid, male XY and female XX samples were age matched perfectly between

768    age 55 to 70. Females are shown in blue and males are shown in lime green. Since the samples

769    were aged perfectly the histogram bars show only the overlap of female and male samples is a mix

770    color of the blue and lime green.

771

772    **Additional file 3. Genetic sex of RNA-Seq samples when aligned using STAR.** Gene expression

35

773    log$_2$(CPM+0.25/L) for select XY homologous genes (*DDX3X/Y, PCDH11X/Y, USP9X/Y, ZFX/Y,*

774    *UTX/Y*) and *XIST* and *SRY* when reads were aligned to a default reference genome **A)**, and for **B)**

775    when reads were aligned to a sex chromosome complement informed reference using STAR. Male

776    XY whole blood, brain cortex, breast, liver, and thyroid samples are shown in blue squares and

777    female XX in orange circles.

778

779    **Additional file 4. Genetic sex of RNA-Seq samples per tissue.** Gene expression

780    log$_2$(CPM+0.25/L) for select XY homologous genes (*DDX3X/Y, PCDH11X/Y, USP9X/Y, ZFX/Y,*

781    *UTX/Y*) and *XIST* and *SRY* when reads were aligned to a default reference genome **A)**, and for **B)**

782    when reads were aligned to a sex chromosome complement informed reference using HISAT and

783    **C)** and **D)**, for when the reads were aligned using STAR. Male XY whole blood, brain cortex,

784    breast, liver, and thyroid samples are shown in blue squares and female XX in orange circles.

785

786    **Additional file 5. List of samples that were removed from downstream analysis.** Samples

787    that did not cluster with the reported sex or clustered in unexpected ways were removed from the

788    differential expression analysis. One male XY whole blood, 4 female XX and 4 male XY brain

789    cortex, and one female XX breast sample were removed.

790

791    **Additional file 6. Multidimensional Scaling plots.** We investigated multidimensional scaling for

792    all shared common variable genes for dimensions 1 and 2, and for dimensions 2 and 3 in each

793    tissue. The most variation in each tissue is explained by the aligner **C.aligner**. The second most

794    variation in each tissue is explained by the sex of the sample **A.sex**.

795

36

796    **Additional file 7. HISAT mapped reads bar plot.**  Mean difference in expression for average

797    total reads mapped for each tissue and each sex when aligned to a sex chromosome informed

798    versus a default reference genome. Paired t-test to test for significant difference in total reads

799    mapped for the whole transcriptome, chromosome 8, and chromosome X. Nonparametric Wilcox

800    single rank sum test was used to test for significant difference in total reads mapped on the Y

801    chromosome for male samples in each tissue separately. Red  * indicate a significant, p-value <

802    0.05,  difference in average mapped reads, NS is no significant differences.

803

804    **Additional file 8. STAR mapped reads bar plot.**  Mean difference in expression for average total

805    reads mapped for each tissue and each sex when aligned to a sex chromosome informed versus a

806    default reference genome. Paired t-test to test for significant difference in total reads mapped for

807    the whole transcriptome, chromosome 8, and chromosome X. Nonparametric Wilcox single rank

808    sum test was used to test for significant difference in total reads mapped on the Y chromosome for

809    male samples in each tissue separately. Red  * indicate a significant, p-value < 0.05,  difference in

810    average mapped reads, NS is no significant differences.

811

812    **Additional file 9. Paired t-test for mapped reads in default compared to sex chromosome**

813    **complement reference genome.** Mean difference in expression for average total reads mapped

814    for each tissue and each sex when aligned to a sex chromosome informed versus a default reference

815    genome. Paired t-test to test for significant difference in total reads mapped for the whole

816    transcriptome (WT), chromosome 8, and chromosome X. Nonparametric Wilcox single rank sum

817    test was used to test for significant difference in total reads mapped on the Y chromosome for male

818    samples in each tissue separately.

819

820 **Additional file 10. X chromosome expression differences between default and sex**

821 **chromosome complement informed alignment.** X chromosome gene expression differences

822 between default and sex chromosome complement informed alignment. Increase in expression

823 when aligned to a sex chromosome complement informed reference genome is a $\log_2$ fold change

824 (FC) > 0. A decrease in expression when aligned to a sex chromosome complement informed

825 reference genome is $\log_2$ FC < 0. Female XX samples are indicated by red and pink circles for

826 PAR1, XTR, PAR2 genes, and for all other X chromosome genes respectively. Blue and light blue

827 squares represent male XY samples. Blue squares indicate which gene points are in PAR1, XTR,

828 and PAR2, and light blue squares are for genes outside of those regions. Differences in X

829 chromosome expression between reference genomes default and sex chromosome complement for

830 male XY and female XX samples aligned using HISAT for the whole X chromosome and the first

831 5Mb are shown for the whole blood (**A** and **B**, respectively), brain cortex (**E** and **F**, respectively),

832 breast (**I** and **J**, respectively), liver (**M** and **N**, respectively), and thyroid (**Q** and **R**, respectively).

833 Differences in X chromosome expression between reference genomes for male XY and female

834 XX samples aligned using STAR for the whole X chromosome and the first 5Mb are shown for

835 the whole blood (**C** and **D**, respectively), brain cortex (**G** and **H**, respectively), breast (**K** and **L**,

836 respectively), liver (**O** and **P**, respectively), and thyroid (**S** and **T**, respectively).

837

838 **Additional file 11. X chromosome regions mean and median expression values.** X

839 chromosome regions PAR1, PAR2, XTR, XDG, XAR, XCR mean and median CPM expression

840 for male XY and female XX samples for each tissue separately when aligned to a default or sex

841 chromosome complement informed reference genome using either HISAT and STAR. Paired t-

842    test was used to test for significant differences in expression. XTR and XAR show a significant

843    increase, p-value < 0.05, in female expression for each tissue type. XTR and XAR additionally

844    show a significant increase, p-value < 0.05, in male expression for liver and thyroid. PAR2 shows

845    a significant increase, p-value < 0.05, in female liver expression. Additionally reported fold change

846    in mean expression when using a sex chromosome complement informed compared to a default

847    reference genome. The mean fold change in expression either increased or stayed the same ranging

848    from 2.8 to 0.999 fold increase in expression. Finally, mean male over mean female expression

849    was reported for each X chromosome region for each tissue. Mean male over mean female

850    expression decreases for XTR when using a sex chromosome complement reference genome for

851    each tissue.

852

853    **Additional file 12. Whole genome gene expression values per sample, aligner and reference**

854    **genome used for alignment.** CPM values for male XY and female XX whole blood, brain cortex,

855    breast, liver and thyroid samples when aligned to a default and sex chromosome complement

856    informed reference genome for the whole genome (1-22, mtDNA, X, Y and non-chromosomal).

857

858    **Additional file 13. Gene expression for XY homologous genes.** X chromosome expression for

859    26 X and Y homologous genes (*AMELX*, *ARSD*, *ARSE*, *ARSF*, *CASK*, *GYG2*, *HSFX1*, *HSFX2*,

860    *NLGN4X*, *OFD1*, *PCDH11X*, *PRKX*, *RBMX*, *RPS4X*, *SOX3*, *STS*, *TBL1X*, *TGIF2LX*, *TMSB4X*,

861    *TSPYL2*, *USP9X*, *VCX*, *VCX2*, *VCX3A*, *VCX3B*, *ZFX)*. Difference in gene expression for when

862    male XY and female XX samples were aligned to a default and sex chromosome complement

863    informed reference genome for each tissue. Little to no difference in gene expression between

864    default and sex chromosome complement informed reference genome alignment was observed for

865　25 of the 26 X and Y homologous genes for both male XY and female XX samples using either

866　HISAT or STAR. The log2 fold increase in expression for *PCDH11X* when aligned using HISAT

867　was 0.4, 0.28, 0.33, 0.16, and 0.16 for whole blood, brain cortex, breast, liver, and thyroid,

868　respectively. The greatest increase in expression was observed for  *PCDH11X* in female whole

869　blood at a log2 fold increase of 0.4.

870

871　**Additional file 14. Differentially expressed genes between the sexes that were uniquely and**

872　**jointly called between reference genomes.** Genes that are differentially expressed between the

873　sexes, male XY and female XX, for whole blood, brain cortex, breast, liver, and thyroid samples.

874　Differentially expressed genes that are uniquely called when using either the default or sex

875　chromosome complement informed reference genome and differentially expressed genes that were

876　jointly called between the reference genomes.

877

878　**Additional file 15. Gene expression differences between male XY and female XX samples.**

879　Sex differences in gene expression for whole blood, brain cortex, breast, liver, and thyroid samples

880　for when samples were aligned to a default reference genome and to a reference genome informed

881　on the sex chromosome complement. Showing sex differences in gene expression between

882　reference genomes used for alignment and for when samples were aligned using HISAT and

883　STAR.

884

885　**Additional file 16. GO analysis of differentially expressed genes in female and male samples**

886　**with HISAT aligner.** Gene enrichment analysis of genes that are more highly expressed in one

887　sex versus the other sex for each tissue, whole blood, brain cortex, breast, liver and thyroid, when

888    samples were aligned to a default or sex chromosome complement informed reference genome

889    using HISAT.

890

891    **Additional file 17. GO analysis of differentially expressed genes in female and male samples**

892    **with STAR aligner.** Gene enrichment analysis of genes that are more highly expressed in one sex

893    versus the other sex for each tissue, whole blood, brain cortex, breast, liver and thyroid, when

894    samples were aligned to a default or sex chromosome complement informed reference genome

895    using STAR.

896

897    **Additional file 18. Sex chromosome complement informed transcriptome reference**

898    **eliminates Y-linked expression in female XX samples. A)** Sex differences in gene expression,

899    $\log_2(\text{CPM}+0.25/\text{L})$, between the sixteen samples from genetic males and females are shown when

900    aligning all samples to the default Ensembl reference transcriptome (left) and a reference

901    transcriptome informed on the sex chromosome complement (right) for brain cortex. Each point

902    represents a gene. Genes that are differentially expressed, adjusted p-value < 0.01 are indicated in

903    black for autosomal genes, blue for Y-linked genes, and red for X-linked genes. **B)** We show

904    overlap between genes that are called as differentially expressed when all samples are pseudo-

905    aligned to the default transcriptome, and genes that are called as differentially expressed when

906    pseudo-aligned to a sex chromosome complement informed transcriptome reference. When

907    samples were aligned to a reference transcriptome informed on the sex chromosome complement,

908    14 genes were called as differentially expressed between the sexes. *PLCXD1* was uniquely called

909    as differentially expressed when aligned to a default reference genome.

910

41

911

912 **Additional file 18. Ensembl sex chromosome complement informed transcriptome reference**

913 **eliminates Y-linked expression in female XX samples. A)** Sex differences in gene expression,

914 $\log_2$(CPM+0.25/L), between the sixteen samples from genetic males and females are shown when

915 aligning all samples to the default Ensembl reference transcriptome (left) and a reference

916 transcriptome informed on the sex chromosome complement (right) for brain cortex. Each point

917 represents a gene. Genes that are differentially expressed, adjusted p-value < 0.01 are indicated in

918 black for autosomal genes, blue for Y-linked genes, and red for X-linked genes. **B)** We show

919 overlap between genes that are called as differentially expressed when all samples are pseudo-

920 aligned to the default transcriptome, and genes that are called as differentially expressed when

921 pseudo-aligned to a sex chromosome complement informed transcriptome reference. When

922 samples were aligned to a reference transcriptome informed on the sex chromosome complement,

923 14 genes were called as differentially expressed between the sexes. *PLCXD1* was uniquely called

924 as differentially expressed when aligned to a default reference genome.

925

926 **Additional file 19. Gencode sex chromosome complement informed transcriptome reference**

927 **eliminates Y-linked expression in female XX samples. A)** Sex differences in gene expression,

928 $\log_2$(CPM+0.25/L), between the sixteen samples from genetic males and females are shown when

929 aligning all samples to the default gencode reference transcriptome (left) and a reference

930 transcriptome informed on the sex chromosome complement (right) for brain cortex. Each point

931 represents a gene. Genes that are differentially expressed, adjusted p-value < 0.01 are indicated in

932 black for autosomal genes, blue for Y-linked genes, and red for X-linked genes. **B)** We show

933 overlap between genes that are called as differentially expressed when all samples are pseudo-

42

934    aligned to the default transcriptome, and genes that are called as differentially expressed when

935    pseudo-aligned to a sex chromosome complement informed transcriptome reference. When

936    samples were aligned to a reference transcriptome informed on the sex chromosome complement,

937    17 genes were called as differentially expressed between the sexes. *ZBED1* was uniquely called as

938    differentially expressed when aligned to a default reference genome.

939

940    **Additional file 20. 3 male XY and 3 female XX brain cortex and whole blood differential**

941    **expression analysis.** Replicated analysis in a smaller sample size of 3 male XY compared to 3

942    female XX samples for whole blood and brain cortex tissue. Samples were randomly selected,

943    and confirm the results from the larger sample size.

944 **References**

945 Aken, B.L., Achuthan, P., Akanni, W., Amode, M.R., Bernsdorff, F., Bhai, J., Billis, K.,

946 Carvalho-Silva, D., Cummins, C., Clapham, P., et al. (2017). Ensembl 2017. Nucleic Acids Res.

947 *45*, D635–D642.

948 Albrecht, K.H., Young, M., Washburn, L.L., and Eicher, E.M. (2003). Sry expression level and

949 protein isoform differences play a role in abnormal testis development in C57BL/6J mice

950 carrying certain Sry alleles. Genetics *164*, 277–288.

951 Andrews, S. Babraham Bioinformatics - FastQC A Quality Control tool for High Throughput

952 Sequence Data.

953 Arnold, A.P., and Chen, X. (2009). What does the "four core genotypes" mouse model tell us

954 about sex differences in the brain and other tissues? Front. Neuroendocrinol. *30*, 1–9.

955 Arnold, A.P., Chen, X., and Itoh, Y. (2012). What a difference an X or Y makes: sex

956 chromosomes, gene dose, and epigenetics in sexual differentiation. Handb. Exp. Pharmacol. 67–

957 88.

958 Barnett, D.W., Garrison, E.K., Quinlan, A.R., Strömberg, M.P., and Marth, G.T. (2011).

959 BamTools: a C++ API and toolkit for analyzing and managing BAM files. Bioinformatics *27*,

960 1691–1692.

961 Bolger, A.M., Lohse, M., and Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina

962 sequence data. Bioinformatics *30*, 2114–2120.

963   Bray, N., Pimentel, H., Melsted, P., and Pachter, L. (2015). Near-optimal RNA-Seq

964   quantification. ArXiv150502710 Cs Q-Bio.

965   Cancer Genome Atlas Research Network, Weinstein, J.N., Collisson, E.A., Mills, G.B., Shaw,

966   K.R.M., Ozenberger, B.A., Ellrott, K., Shmulevich, I., Sander, C., and Stuart, J.M. (2013). The

967   Cancer Genome Atlas Pan-Cancer analysis project. Nat. Genet. *45*, 1113–1120.

968   Carrel, L., and Willard, H.F. (2005). X-inactivation profile reveals extensive variability in X-

969   linked gene expression in females. Nature *434*, 400–404.

970   Chaligné, R., Popova, T., Mendoza-Parra, M.-A., Saleem, M.-A.M., Gentien, D., Ban, K., Piolot,

971   T., Leroy, O., Mariani, O., Gronemeyer, H., et al. (2015). The inactive X chromosome is

972   epigenetically unstable and transcriptionally labile in breast cancer. Genome Res. *25*, 488–503.

973   Charchar, F.J., Svartman, M., El-Mogharbel, N., Ventura, M., Kirby, P., Matarazzo, M.R.,

974   Ciccodicola, A., Rocchi, M., D'Esposito, M., and Graves, J.A.M. (2003). Complex Events in the

975   Evolution of the Human Pseudoautosomal Region 2 (PAR2). Genome Res. *13*, 281–286.

976   Charlesworth, B. (1991). The evolution of sex chromosomes. Science *251*, 1030–1033.

977   Ching, T., Huang, S., and Garmire, L.X. (2014). Power analysis and sample size estimation for

978   RNA-Seq differential expression. RNA N. Y. N *20*, 1684–1696.

979   Conesa, A., Madrigal, P., Tarazona, S., Gomez-Cabrero, D., Cervera, A., McPherson, A.,

980   Szcześniak, M.W., Gaffney, D.J., Elo, L.L., Zhang, X., et al. (2016). A survey of best practices

981   for RNA-seq data analysis. Genome Biol. *17*, 13.

982    Consortium, T.Gte. (2015). The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue

983    gene regulation in humans. Science *348*, 648–660.

984    Costa-Silva, J., Domingues, D., and Lopes, F.M. (2017). RNA-Seq differential expression

985    analysis: An extended review and a software tool. PloS One *12*, e0190152.

986    Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson,

987    M., and Gingeras, T.R. (2013). STAR: ultrafast universal RNA-seq aligner. Bioinforma. Oxf.

988    Engl. *29*, 15–21.

989    Dumanski, J.P., Lambert, J.-C., Rasi, C., Giedraitis, V., Davies, H., Grenier-Boley, B., Lindgren,

990    C.M., Campion, D., Dufouil, C., European Alzheimer's Disease Initiative Investigators, et al.

991    (2016). Mosaic Loss of Chromosome Y in Blood Is Associated with Alzheimer Disease. Am. J.

992    Hum. Genet. *98*, 1208–1219.

993    Eden, E., Navon, R., Steinfeld, I., Lipson, D., and Yakhini, Z. (2009). GOrilla: a tool for

994    discovery and visualization of enriched GO terms in ranked gene lists. BMC Bioinformatics *10*,

995    48.

996    Eyras, E., Caccamo, M., Curwen, V., and Clamp, M. (2004). ESTGenes: Alternative Splicing

997    From ESTs in Ensembl. Genome Res. *14*, 976–987.

998    Forsberg, L.A. (2017). Loss of chromosome Y (LOY) in blood cells is associated with increased

999    risk for disease and mortality in aging men. Hum. Genet. *136*, 657–663.

1000    Gershoni, M., and Pietrokovski, S. (2017). The landscape of sex-differential transcriptome and

1001    its consequent selection in human adults. BMC Biol. *15*, 7.

46

1002  Goldstein, J.M., Holsen, L., Handa, R., and Tobet, S. (2014). Fetal hormonal programming of

1003  sex differences in depression: linking women's mental health with sex differences in the brain

1004  across the lifespan. Front. Neurosci. *8*.

1005  Grassmann, F., Kiel, C., den Hollander, A.I., Weeks, D.E., Lotery, A., Cipriani, V., Weber,

1006  B.H.F., and International Age-related Macular Degeneration Genomics Consortium (IAMDGC)

1007  (2019). Y chromosome mosaicism is associated with age-related macular degeneration. Eur. J.

1008  Hum. Genet. EJHG *27*, 36–41.

1009  GTEx Consortium (2015). Human genomics. The Genotype-Tissue Expression (GTEx) pilot

1010  analysis: multitissue gene regulation in humans. Science *348*, 648–660.

1011  Harrow, J., Frankish, A., Gonzalez, J.M., Tapanari, E., Diekhans, M., Kokocinski, F., Aken,

1012  B.L., Barrell, D., Zadissa, A., Searle, S., et al. (2012). GENCODE: the reference human genome

1013  annotation for The ENCODE Project. Genome Res. *22*, 1760–1774.

1014  J, H., A, F., Jm, G., E, T., M, D., F, K., Bl, A., D, B., A, Z., S, S., et al. (2012). GENCODE: the

1015  reference human genome annotation for The ENCODE Project. Genome Res. *22*, 1760–1774.

1016  Khramtsova, E., Davis, L., and Stranger, B. (2018). The role of sex in the genomics of human

1017  complex traits. Nat. Rev. Genet. *20*.

1018  Kim, D., Langmead, B., and Salzberg, S.L. (2015). HISAT: a fast spliced aligner with low

1019  memory requirements. Nat. Methods *12*, 357–360.

1020  Lahn, B.T., and Page, D.C. (1999). Four evolutionary strata on the human X chromosome.

1021  Science *286*, 964–967.

47

1022　Lamarre, S., Frasse, P., Zouine, M., Labourdette, D., Sainderichin, E., Hu, G., Le Berre-Anton,

1023　V., Bouzayen, M., and Maza, E. (2018). Optimization of an RNA-Seq Differential Gene

1024　Expression Analysis Depending on Biological Replicate Number and Library Size. Front. Plant

1025　Sci. *9*, 108.

1026　Lappalainen, T., Sammeth, M., Friedländer, M.R., 't Hoen, P.A.C., Monlong, J., Rivas, M.A.,

1027　Gonzàlez-Porta, M., Kurbatova, N., Griebel, T., Ferreira, P.G., et al. (2013). Transcriptome and

1028　genome sequencing uncovers functional variation in humans. Nature *501*, 506–511.

1029　Law, C.W., Chen, Y., Shi, W., and Smyth, G.K. (2014). voom: precision weights unlock linear

1030　model analysis tools for RNA-seq read counts. Genome Biol. *15*, R29.

1031　Leinonen, R., Sugawara, H., and Shumway, M. (2011). The Sequence Read Archive. Nucleic

1032　Acids Res. *39*, D19–D21.

1033　Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics *34*, 3094–

1034　3100.

1035　Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G.,

1036　Durbin, R., and 1000 Genome Project Data Processing Subgroup (2009). The Sequence

1037　Alignment/Map format and SAMtools. Bioinforma. Oxf. Engl. *25*, 2078–2079.

1038　Liao, Y., Smyth, G.K., and Shi, W. (2014). featureCounts: an efficient general purpose program

1039　for assigning sequence reads to genomic features. Bioinforma. Oxf. Engl. *30*, 923–930.

1040    Mallick, S., Li, H., Lipson, M., Mathieson, I., Gymrek, M., Racimo, F., Zhao, M., Chennagiri,

1041    N., Nordenfelt, S., Tandon, A., et al. (2016). The Simons Genome Diversity Project: 300

1042    genomes from 142 diverse populations. Nature *538*, 201–206.

1043    Natri, H.M., Wilson, M.A., and Buetow, K.H. (2019). Distinct molecular etiologies of male and

1044    female hepatocellular carcinoma. BMC Cancer *19*, 951.

1045    Naugler, W.E., Sakurai, T., Kim, S., Maeda, S., Kim, K., Elsharkawy, A.M., and Karin, M.

1046    (2007). Gender disparity in liver cancer due to sex differences in MyD88-dependent IL-6

1047    production. Science *317*, 121–124.

1048    Pandey, R.S., Wilson Sayres, M.A., and Azad, R.K. (2013). Detecting evolutionary strata on the

1049    human x chromosome in the absence of gametologous y-linked sequences. Genome Biol. Evol.

1050    *5*, 1863–1871.

1051    Patro, R., Duggal, G., Love, M.I., Irizarry, R.A., and Kingsford, C. (2017). Salmon provides fast

1052    and bias-aware quantification of transcript expression. Nat. Methods *14*, 417–419.

1053    Piskol, R., Ramaswami, G., and Li, J.B. (2013). Reliable identification of genomic variants from

1054    RNA-seq data. Am. J. Hum. Genet. *93*, 641–651.

1055    R, P., Sm, M., and C, K. (2014). Sailfish Enables Alignment-Free Isoform Quantification From

1056    RNA-seq Reads Using Lightweight Algorithms (Nat Biotechnol).

1057    Rahbari, R., Zhang, L., and Kebebew, E. (2010). Thyroid cancer gender disparity. Future Oncol.

1058    Lond. Engl. *6*, 1771–1779.

1059    Raznahan, A., Parikshak, N.N., Chandran, V., Blumenthal, J.D., Clasen, L.S., Alexander-Bloch,

1060    A.F., Zinn, A.R., Wangsa, D., Wise, J., Murphy, D.G.M., et al. (2018). Sex-chromosome dosage

1061    effects on gene expression in humans. Proc. Natl. Acad. Sci. U. S. A. *115*, 7398–7403.

1062    Robinson, M.D., McCarthy, D.J., and Smyth, G.K. (2010). edgeR: a Bioconductor package for

1063    differential expression analysis of digital gene expression data. Bioinforma. Oxf. Engl. *26*, 139–

1064    140.

1065    Ross, M.T., Grafham, D.V., Coffey, A.J., Scherer, S., McLay, K., Muzny, D., Platzer, M.,

1066    Howell, G.R., Burrows, C., Bird, C.P., et al. (2005). The DNA sequence of the human X

1067    chromosome. Nature *434*, 325–337.

1068    Seyednasrollah, F., Laiho, A., and Elo, L.L. (2015). Comparison of software packages for

1069    detecting differential expression in RNA-seq studies. Brief. Bioinform. *16*, 59–70.

1070    Shi, L., Zhang, Z., and Su, B. (2016). Sex Biased Gene Expression Profiling of Human Brains at

1071    Major Developmental Stages. Sci. Rep. *6*, 21181.

1072    Skaletsky, H., Kuroda-Kawaguchi, T., Minx, P.J., Cordum, H.S., Hillier, L., Brown, L.G.,

1073    Repping, S., Pyntikova, T., Ali, J., Bieri, T., et al. (2003). The male-specific region of the human

1074    Y chromosome is a mosaic of discrete sequence classes. Nature *423*, 825–837.

1075    Soneson, C., Love, M.I., and Robinson, M.D. (2015). Differential analyses for RNA-seq:

1076    transcript-level estimates improve gene-level inferences. F1000Research *4*, 1521.

1077     Traglia, M., Bseiso, D., Gusev, A., Adviento, B., Park, D.S., Mefford, J.A., Zaitlen, N., and

1078     Weiss, L.A. (2017). Genetic Mechanisms Leading to Sex Differences Across Common Diseases

1079     and Anthropometric Traits. Genetics *205*, 979–992.

1080     Tukiainen, T., Villani, A.-C., Yen, A., Rivas, M.A., Marshall, J.L., Satija, R., Aguirre, M.,

1081     Gauthier, L., Fleharty, M., Kirby, A., et al. (2016). Landscape of X chromosome inactivation

1082     across human tissues. BioRxiv 073957.

1083     Turner, M.E., Ely, D., Prokop, J., and Milsted, A. (2011). Sry, more than testis determination?

1084     Am. J. Physiol.-Regul. Integr. Comp. Physiol. *301*, R561–R571.

1085     Veerappa, A.M., Padakannaya, P., and Ramachandra, N.B. (2013). Copy number variation-based

1086     polymorphism in a new pseudoautosomal region 3 (PAR3) of a human X-chromosome-

1087     transposed region (XTR) in the Y chromosome. Funct. Integr. Genomics *13*, 285–293.

1088     Webster, T.H., Couse, M., Grande, B.M., Karlins, E., Phung, T.N., Richmond, P.A., Whitford,

1089     W., and Wilson, M.A. (2019). Identifying, understanding, and correcting technical artifacts on

1090     the sex chromosomes in next-generation sequencing data. GigaScience *8*.

1091     Wu, D.C., Yao, J., Ho, K.S., Lambowitz, A.M., and Wilke, C.O. (2018). Limitations of

1092     alignment-free tools in total RNA-seq quantification. BMC Genomics *19*, 510.

1093     Zhao, S., and Zhang, B. (2015). A comprehensive evaluation of ensembl, RefSeq, and UCSC

1094     annotations in the context of RNA-seq read mapping and gene quantification. BMC Genomics

1095     *16*, 97.

1096    Zhao, S., Li, C.-I., Guo, Y., Sheng, Q., and Shyr, Y. (2018). RnaSeqSampleSize: real data based

1097    sample size estimation for RNA sequencing. BMC Bioinformatics *19*, 191.

1098    Zielezinski, A., Vinga, S., Almeida, J., and Karlowski, W.M. (2017). Alignment-free sequence

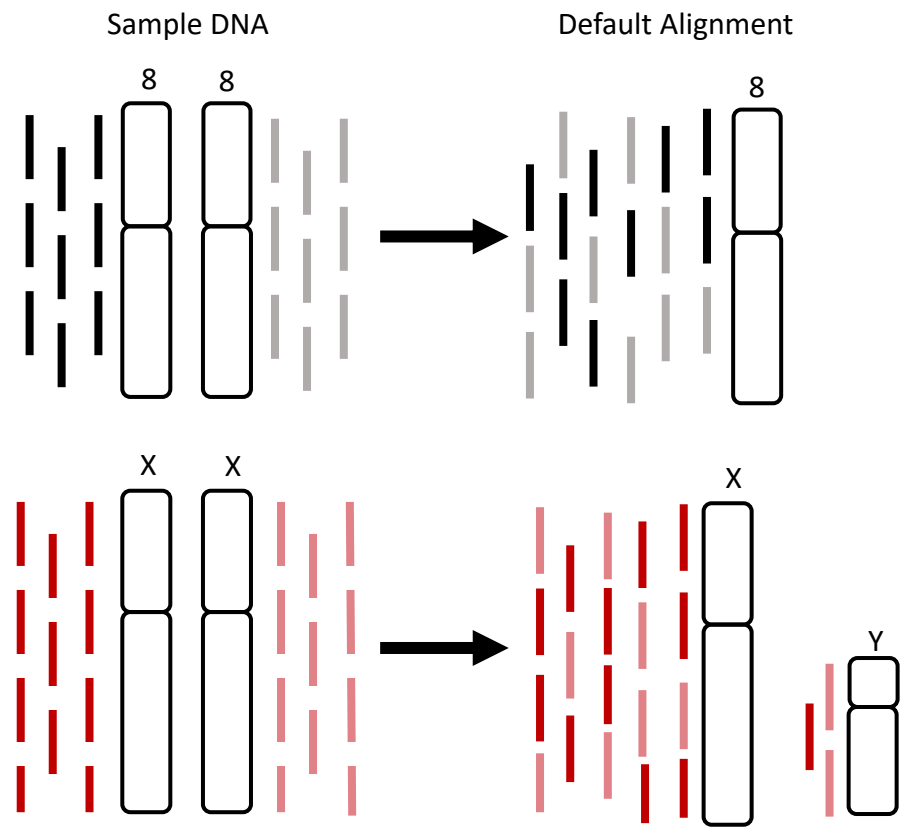1099    comparison: benefits, applications, and tools. Genome Biol. *18*, 186.

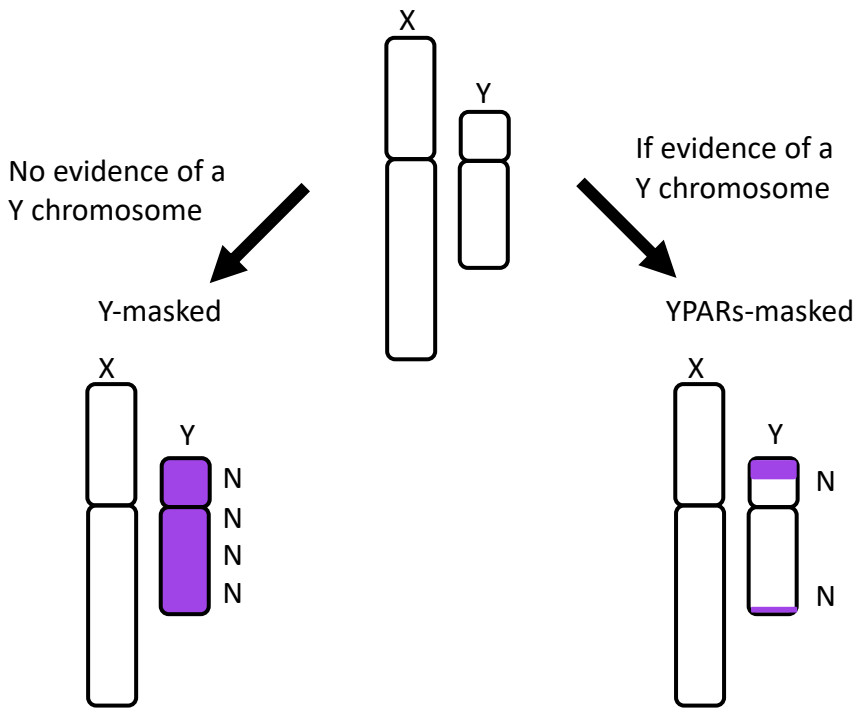1100    (2020). broadinstitute/picard (Broad Institute).
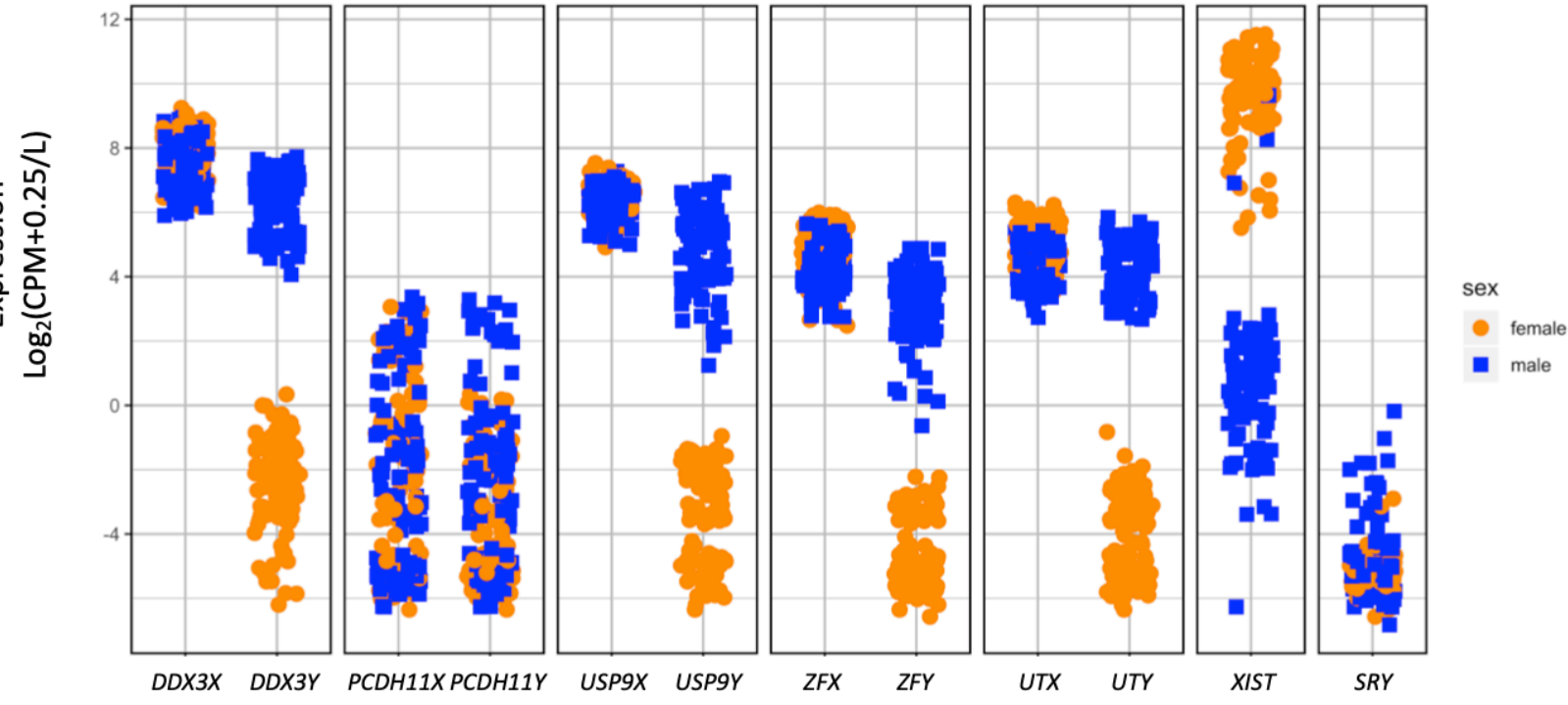
1101

**A** X and Y sequence homology

**B** RNA-seq alignment to a default reference genome

Sample DNA

Default Alignment

**C** Sex chromosome complement informed alignment

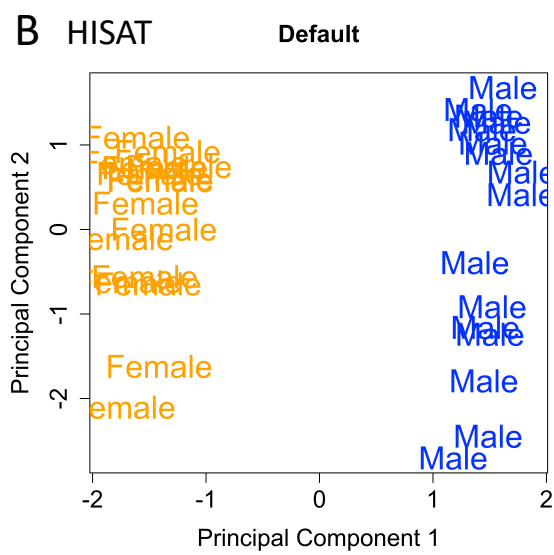No evidence of a Y chromosome
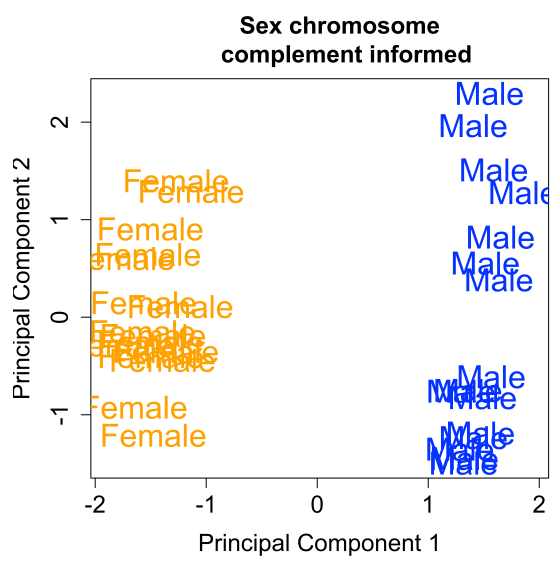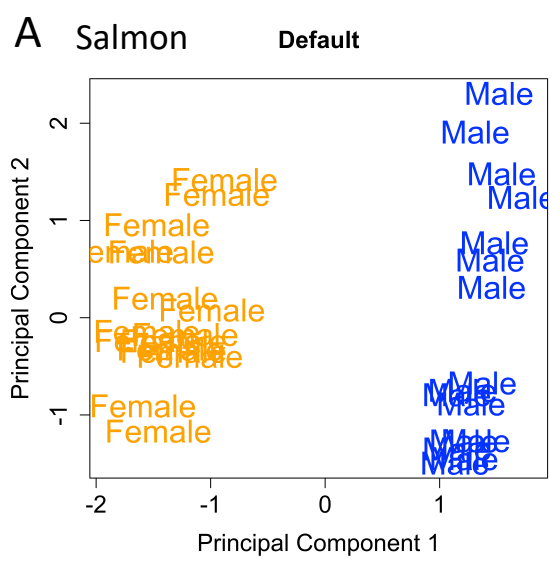
If evidence of a Y chromosome

Y-masked

YPARs-masked

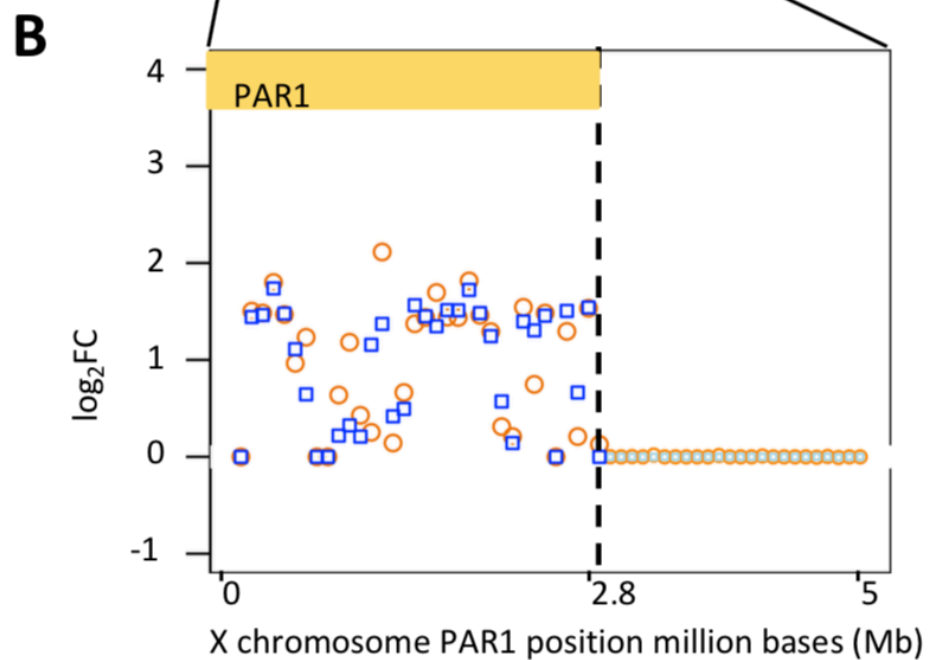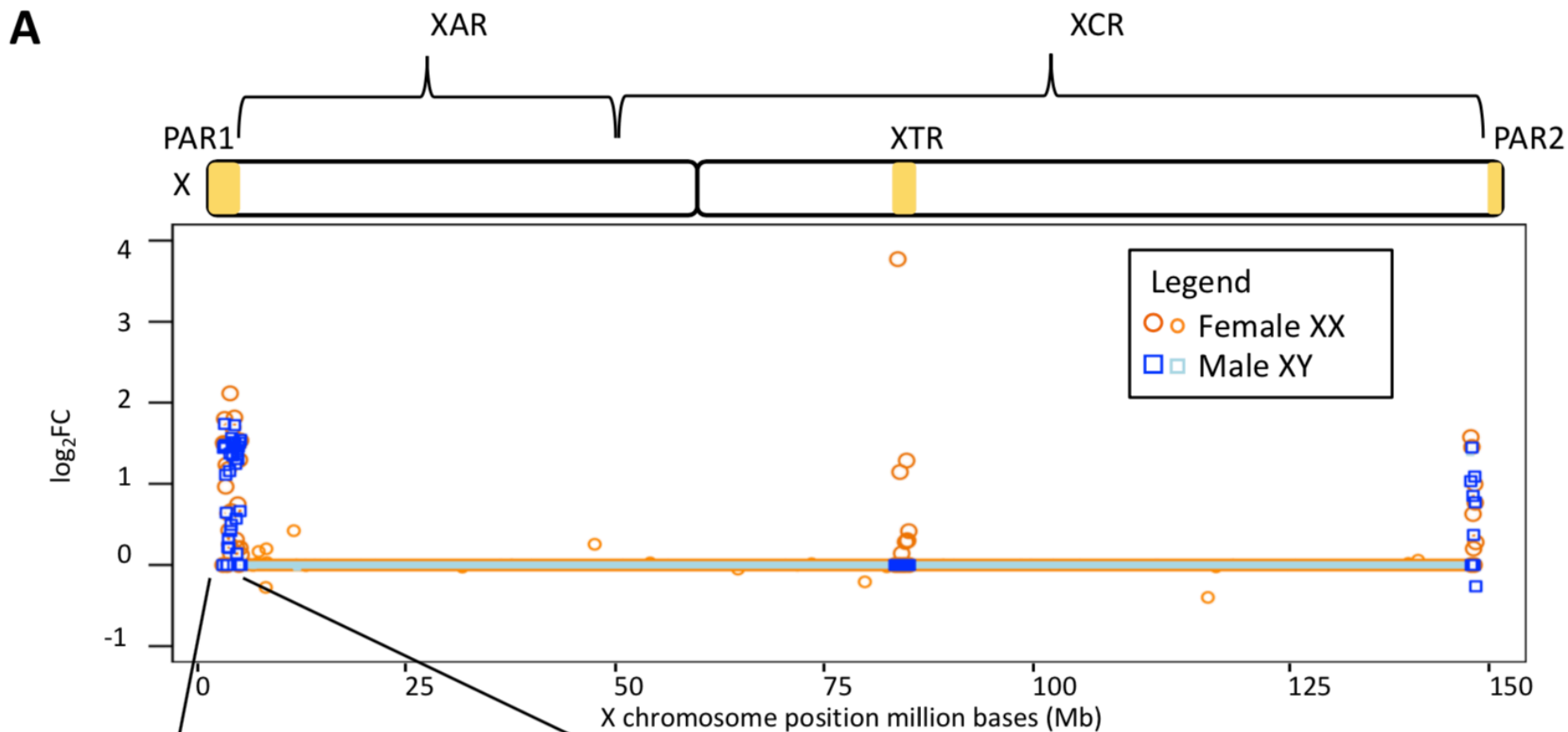**A** All TISSUES aligned to HISAT and default reference genome

**B** All TISSUES aligned to HISAT and sex chromosome complement reference genome

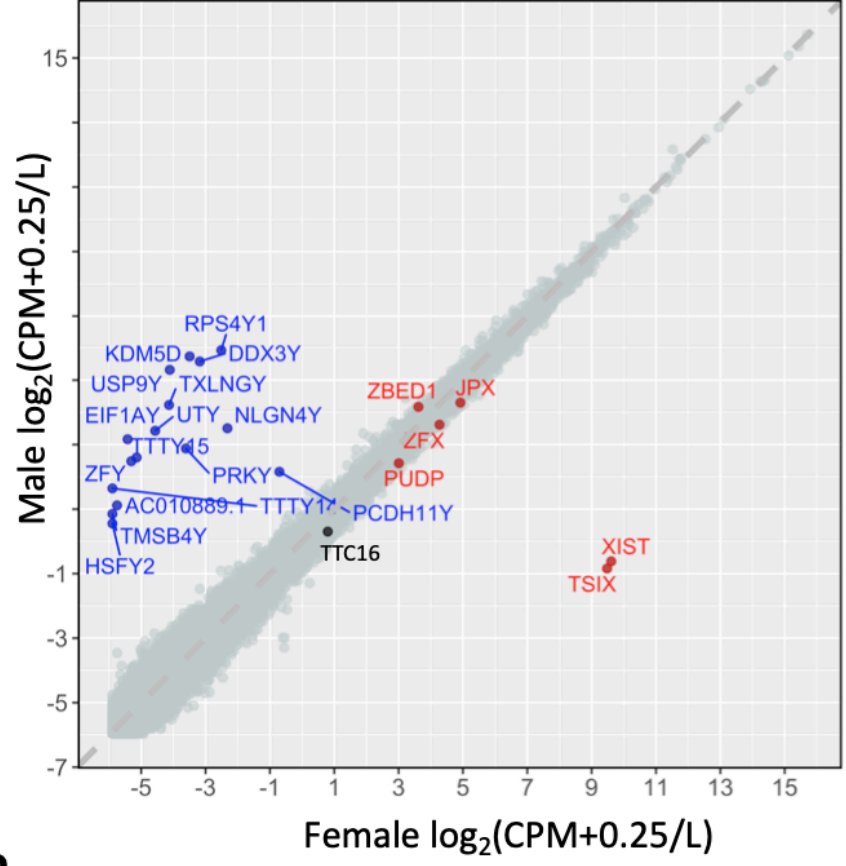**A** Salmon    **Default**        **Sex chromosome complement informed**

**B** HISAT    **Default**        **Sex chromosome complement informed**

**C** STAR    **Default**        **Sex chromosome complement informed**

**A**

**B**

**C** Fold change in mean expression for X chromosome regions

| X chromosome regions | Female 46, XX | Male 46, XY |
|---|---|---|
| PAR1 | 2.73 | 2.75 |
| XCR | 1.00 | 1.00 |
| XAR | 1.00 | 1.00 |
| XTR | 1.22 | 1.00 |
| XDG | 1.00 | 1.00 |
| PAR2 | 2.13 | 2.19 |

**A**

Default alignment

Sex chromosome complement informed alignment

Male log$_2$(CPM+0.25/L)

Female log$_2$(CPM+0.25/L)

Legend
FDR < 0.01
autosomal
X-linked
Y-linked

Left plot labels: RPS4Y1, KDM5D, DDX3Y, USP9Y, TXLNGY, EIF1AY, UTY, NLGN4Y, TTTY15, ZFY, PRKY, AC010889.1, TTTY14, PCDH11Y, TMSB4Y, HSFY2, ZBED1, JPX, ZFX, PUDP, TTC16, XIST, TSIX

Right plot labels: RPS4Y1, KDM5D, DDX3Y, USP9Y, TXLNGY, UTY, NLGN4Y, EIF1AY, PRKY, PSMA6P1, ZFY, PCDH11Y, TTTY15, EIF4A1P2, TTTY14, AC010889.1, HSFY2, TMSB4Y, VDAC1P6, TTC16, ZBED1, PLCXD1, JPX, ZFX, PUDP, XIST, TSIX

**B**

| Genes | Total | Unique |
|---|---|---|
| Autosomal | 1 | 0 |
| X-linked | 6 | 0 |
| Y-linked | 16 | 0 |

23    4

| Genes | Total | Unique |
|---|---|---|
| Autosomal | 1 | 0 |
| X-linked | 7 | 1 |
| Y-linked | 19 | 3 |