# SUPPLEMENTAL MATERIAL

# **Quantitative analysis of Y-Chromosome gene expression across 36 human tissues**

Alexander K. Godfrey, Sahin Naqvi, Lukáš Chmátal, Joel M. Chick, Richard N. Mitchell, Steven P. Gygi, Helen Skaletsky, David C. Page.

SUPPLEMENTAL METHODS	2
Human transcriptome annotation and MSY genes	2
RNA-seq simulations	2
Quality control of GTEx RNA-seq analysis	4
Estimating gene mappability	6
Correlated expression of X and Y homologs	6
microRNA analyses	7
Analyses of EIF1AX/Y sequence and expression across species	8
Analysis of protein abundance in human heart tissue	8
Immunoblotting	13
SUPPLEMENT REFERENCES	14
SUPPLEMENTAL FIGURES	17

Note: Supplemental Tables and Supplemental Files uploaded individually, separate from this document.

# SUPPLEMENTAL METHODS

## Human transcriptome annotation and MSY genes

All human analyses use transcript/gene models defined in a custom subset of the comprehensive GENCODE version 24 transcript annotation, comprising the union of transcripts that (1) belong to the "GENCODE Basic" annotation and (2) are recognized by the Consensus Coding Sequence project (Pruitt et al. 2009). Filtering the comprehensive annotation in this manner enriches for full-length, manually curated transcripts defined by two distinct sources. The list of protein-coding human MSY genes analyzed in this study is based on our annotation of the male-specific region of the human Y Chromosome (Skaletsky et al. 2003) and is delineated in Supplementary Table S1. We note the following differences between the set of protein-coding genes analyzed here from the set of protein-coding described in the GENCODE v24 annotation. First, we excluded eight clone-based Ensembl genes (nomenclature: AC#######.#), which do not have official HGNC symbols and were either removed or reclassified as pseudogenes in subsequent versions of the GENCODE annotation. Second, we excluded PRORY from our analyses, which was not part of our original annotation (Skaletsky et al. 2003). Its cDNA sequence returns no matches by BLAST in the NCBI EST database, and we could not find compelling evidence of its transcription (e.g., RNA-seq reads that span exonexon junctions) in the GTEx samples. Third, we included PRKY and TXLNGY (CYorf15A/B) among the list of protein-coding genes. Both are currently listed as pseudogenes in the public annotation, as the result of significant structural differences with their X-linked homologs: PRKY lost an exon near the 3' end of its coding region, creating a premature termination codon and making it a candidate for nonsense mediated decay; TXLNGY comprises two transcription units (CYorf15A and CYorf15B), homologous to the 5' and 3' ends of X-linked TXLNG. However, both Y-linked genes retain significant open reading frames, and comparisons of their sequences on the human and rhesus macague Y Chromosomes suggest that the coding sequences of *PRKY* and *CYorf15A* remain under purifying selection (dN/dS < 1) (Hughes et al. 2012).

# **RNA-seq simulations**

Simulated RNA-seq libraries were generated using RSEM (v1.2.22) (Li and Dewey 2011). Sequencing parameters for the simulation were obtained by running rsem-calculate-expression with option --estimate-rpsd on GTEx testis sample GTEX-P4QS-2126-SM-3NMCF, supplying our modified transcriptome annotation as a reference. The output file of expression-level estimates ("isoforms.results") was then modified to set the expression levels of Y-Chromosome genes and their X-linked homologs to predetermined levels, following one of four scenarios:

- MSY genes/gene families = 0 TPM; X-linked homologs unmodified (i.e., kept at levels estimated in GTEX-P4QS-2126-SM-3NMCF
- (2) MSY genes/gene families = 1 TPM, X-linked homologs = 2 TPM
- (3) MSY genes/gene families = 5 TPM, X-linked homologs = 10 TPM
- (4) MSY genes/gene families = 5 TPM, X-linked homologs set to a random value between 0 and 10 TPM

For genes with multiple transcript isoforms, the relative abundance of each isoform was assigned in proportion to a random number drawn from a heavy-tailed distribution (Pareto with tail index  $\alpha = 0.5$ ). The relative abundance of individual members of multi-copy gene families were determined similarly, such that the summed expression level of genes in the family equaled 1, 2, 5, or 10 TPM as indicated. These relative isoform abundances were drawn anew in each simulation, to sample different configurations of alternative-isoform expression. 50 simulated RNA-seq libraries were generated for each of the four expression-level scenarios using rsem-simulate-reads, with 50 million 76bp paired-end reads in each library (median depth of samples in GTEx is ~78 million reads). Because of the random read-sampling process used in the simulations, the observed expression level for gene *g* in simulated library *j* (as given by the output of rsem-simulate-reads) deviated slightly from its idealized value (e.g., 1, 2, 5, or 10 TPM). To correct for this source of error when plotting results in Supplemental Fig. S2, the estimated expression levels obtained by various methods for gene *g* in simulated library *j* were multiplied by the ratio of idealized-to-observed values for that gene in that library.

Three methods were used to estimate the expression levels of Y-Chromosome genes and their X-linked homologs in these simulated libraries. First, simulated reads were aligned to the genome using TopHat2 (Kim et al. 2013) (v2.1.1; using parameters --no-mixed --nodiscordant) and the number of uniquely mapping reads overlapping the exons of each gene was counted with featureCounts (Liao et al. 2014) (v1.6.2), requiring both reads from each fragment to be mapped (-B) and each read to entirely overlap annotated exons (--fracOverlap 1). Read-count values for each gene were converted to TPM units, with the length of each gene given by the total length of the union of its exons. This procedure ("unique reads" in Supplemental Fig. S2) is similar to that followed by The GTEx Consortium. Second, after aligning simulated reads with TopHat2, Cufflinks (Trapnell et al. 2010) (v2.2.1) was used to estimate the expression levels of annotated transcripts (-G) in "multi-read-correct" mode (-u), and estimates in FPKM units were converted to TPM. Third, simulated reads were input to kallisto (Bray et al. 2016) (v0.42.5), and the expression levels of annotated transcripts were estimated with sequence-bias correction (--bias).

On average, the "unique reads" method over-estimated MSY gene expression levels and produced less precise (i.e., more variable) estimates than kallisto or Cufflinks in simulated datasets. The over-estimated expression levels are likely the result of discarding multi-mapping reads—with fewer mapped reads in each library, each gene receives a larger proportion of reads in the library overall, thus inflating the TPM value (TPM units convey a gene's expression level as a fraction of the total expression). The imprecision of the unique-reads method is likely due to the presence of alternative transcript isoforms. The "unique reads" method calculates the density of reads mapping to each gene using a fixed length for that gene, defined by the concatenation of all constitutive and alternative exons. When only a short alternative isoform is expressed, read density for the gene (and, correspondingly, its expression level) will be underestimated. The variability in the estimates produced by the unique-reads method thus likely reflects the random mixture of alternative isoforms present in any simulated library.

#### Quality control of GTEx RNA-seq analysis

## Initial screening of samples using sample- and donor-level metadata

All RNA-seq samples meeting the following criteria were downloaded for initial consideration: (1) RIN (SMRIN)  $\ge$  6.0; (2) not annotated as "severely" autolyzed (SMATSSCR != 3); (3) donor deemed eligible by GTEx (INCEXC == True); (4) not flagged for removal by GTEx (SMTORMVE != 'FLAGGED'); (5) was generated from a primary, solid tissue (i.e., whole-blood samples and cell-line samples were excluded); (6) was generated from a tissue where at least 10 samples from male donors were available.

#### Identification and removal of gene-expression outlier samples

Gene-expression outlier samples were detected and removed following a procedure similar to that used in Wright et al. (Wright et al. 2014). Separately for each tissue, the pairwise Pearson correlation coefficient,  $r_{ij}$ , was calculated between the log<sub>2</sub>(TPM + 0.1) expression levels of samples *i* and *j*, using only genes with a median expression level  $\geq$ 3 TPM among the samples of that tissue. The median similarity between sample *i* and other samples from that tissue was calculated as  $\overline{r_i} = \text{median}_i(r_{ij})$  and re-expressed in median absolute deviations,  $D_i =$ 

 $|\overline{r_i} - \overline{r}| / \text{median}_j(|\overline{r_j} - \overline{r}|)$ , where  $\overline{r} = \text{median}_j(\overline{r_j})$ . Samples with  $D_i > 6$  were marked as outliers and removed from subsequent analyses. This process was repeated iteratively until no such samples remained. After outlier removal, hierarchical clustering was performed on the remaining samples to confirm the efficacy of this approach. Even after outlier removal, we found that the breast samples clustered into two highly dissimilar clusters. Breast samples from one of these clusters were found to be indistinguishable from adipose samples and were also excluded as outliers. A list of samples passing all filtering and outlier detection steps is given in Supplemental File S1.

## Merging similar tissues

The filtered set of samples spanned 28 organs/body sites and 49 tissue types, meaning multiple tissue types were sometimes collected from the same organ/body site (e.g., three tissue types were taken from the esophagus: "gastroesophageal junction", "mucosa", and "muscularis"). In cases where we could not clearly separate samples of two or more tissue types by hierarchical clustering, we merged these tissue labels, treating them as single tissue types: "Brain – Cerebellum"  $\leftarrow$  (Brain – Cerebellum, Brain – Cerebellar Hemisphere); "Brain – Cortex"  $\leftarrow$  (Brain – Cortex, Brain – Anterior cingulate cortex (BA24), Brain – Frontal Cortex (BA9)); "Brain – Striatum"  $\leftarrow$  (Brain – Caudate (basal ganglia), Brain – Nucleus accumbens (basal ganglia), Brain – Putamen (basal ganglia)); "Esophagus – Muscularis"  $\leftarrow$  (Esophagus – Gastroesophageal Junction, Esophagus – Muscularis). Multiple samples from the same tissue type of a single donor were treated as technical replicates. The expression levels of each gene across these replicates were averaged to obtain a single sample representing this donor-tissue combination.

#### Adjusting expression levels for the effects of covariates

Expression levels were corrected for the effects of three covariates: the intronic read mapping rate (SMNTRNRT), sample ischemic time (SMTSISCH), and RNA integrity number (RIN) (SMRIN). The effects of ischemic time and RIN on gene expression have been previously noted (Gallego Romero et al. 2014; Ferreira et al. 2018), and all three variables were significantly correlated with the expression levels of many Y- and X-Chromosome genes across multiple tissues. For samples collected from the brain, donor ischemic time (TRISCHD) was used in place of sample ischemic time, because the latter information was not available. To perform this correction, for each tissue separately, the linear model  $y_g = b_0 + Xb + \varepsilon$  was fit, where  $y_g$ 

is a vector of gene g's normalized log<sub>2</sub>(TPM + 0.1) expression levels across *n* samples,  $b_0$  is an intercept term, *X* is the *n* × 4 matrix of covariates (the three sample-quality variables, plus a fourth variable for sex), *b* is a 4 × n matrix of fixed-effect coefficients for the covariates, and  $\varepsilon$  is a vector of n residuals. To increase the comparability of expression levels across tissues, the covariates were centered on common values in all tissues: 8.0 for RIN, 0.12 for intronic read mapping, 100 minutes for ischemic time (or 450 minutes for brain tissues, which were processed separately from other tissues and had uniformly longer ischemic time). Gene g's adjusted expression levels were calculated as  $y_g^* = y_g - Xb$ .

# Estimating gene mappability

To estimate the short-read mappability of each gene in GENCODE v19, its longest transcript isoform was selected, and all possible 76-nt reads were generated (n - 76 + 1 reads in total, where n is the length of the transcript in nucleotides) using a sliding window. These reads were aligned to the full transcriptome annotation with bowtie (v1.2) (Langmead et al. 2009), allowing 0 mismatches (-v 0) and reporting up to 200 alignments (-k 200). If a read aligned to the transcript isoforms of more than one gene or to more than one position within a single transcript isoform, it was classified as multi-mapping (and otherwise as uniquely mapping). The gene's mappability was then calculated as the fraction of reads from that gene that are uniquely mapping. A chromosome's mappability was calculated as the fraction of all reads generated from genes on that chromosome that mapped uniquely.

## Correlated expression of X and Y homologs

The significance of the correlation between gene *i* and gene *j* was assessed by calculating the proportion of genes in the genome as or more correlated in expression with gene *i* than gene *j*, and vice versa. Specifically, the correlation coefficients between gene *i* and all *N* expressed genes were calculated and ordered from largest to smallest; let  $r_{ij}$  be the rank of gene *j* in this list. (For example,  $r_{ij} = 2$  if only one gene in the genome shows a higher correlation with gene *i* than gene *j*.) The procedure was repeated for gene *j*, and the rank of gene *i*,  $r_{ji}$ , was obtained. The significance of the correlation between gene *i* and *j* was estimated by the average rank, normalized by the number of expressed genes in that tissue:  $(r_{ij} + r_{ji})/2N$ . For example, in skeletal muscle, where 9,445 genes are expressed above 5 TPM, eight genes show higher

correlation with *DDX3X* than *DDX3Y*, and two genes show higher correlation with *DDX3Y* than *DDX3X*; therefore, the average, normalized rank is (3 + 9)/(2 \* 9445) = 0.0006.

## microRNA analyses

For each X–Y gene pair, the 3' UTRs of the X homolog, the Y homolog, and their autosomal chicken ortholog were aligned using PRANK (Löytynoja and Goldman 2005) with default parameters. Scripts from TargetScan 6.0 (Friedman et al. 2009) were then used to identify all potential miRNA target sites in the aligned sequences (targetscan\_60.pl) and calculate their context+ scores and percentiles (targetscan 60 context scores.pl). Context+ scores, rather than the more recent context++ scores (Agarwal et al. 2015), were used because 3P-seq data needed to calculate context++ scores were not available for the human Y Chromosome or chicken. Sites identified in X homologs were validated in the context++ model (TargetScan 7.2 (Agarwal et al. 2015)) (Supplemental Table S9). miRNA-target-site presence/absence in X- and Y-homolog 3'-UTRs was then compared to miRNA expression patterns across human tissues to generate predictions about their differential effects on X- and Y-homolog expression. miRNA expression patterns were assessed using quantile normalized expression values from Ludwig et al. (2016) (Ludwig et al. 2016) (https://ccb-web.cs.uni-saarland.de/tissueatlas/). Expression levels from the two donors were averaged, and only tissues matching a tissue in the GTEx dataset were analyzed. Among target sites for tissue-specific, highly expressed miRNAs, the miR-1 target site in EIF1AX is the target site with the highest context+ score-percentile (i.e., greatest predicted efficacy) preserved in one homolog of an X-Y pair but not the other.

For luciferase assays, the entire *EIF1AY* 3'-UTR and the first 1015bp of the *EIF1AX* 3'-UTR were amplified from human genomic DNA and cloned into the psiCheck-2 vector backbone (Promega) by restriction digest (Pme1, Not1). Using the QuikChange II kit (Agilent), *EIF1AX*'s miR-1 site was changed to shuffled sequence; *EIF1AY*'s disrupted miR-1 site was changed to match that of *EIF1AX*. Each psiCheck plasmid, along miR-1 or miR-124 duplexes, was transfected into HEK293 cells with Lipofectamine 2000 (Thermo Fisher Scientific). *Renilla* and firefly absorbance were quantified 24h post-transfection using the Dual-Luciferase Reporter Assay System (Promega), and the ratio of *Renilla*-to-firefly absorbance was calculated. Primers for cloning and mutagenesis, and miRNA oligonucleotide sequences, are listed in Supplemental File S2.

#### Analyses of EIF1AX/Y sequence and expression across species

Non-human RNA-seq data are from Brawand et al. and Merkin et al. (Brawand et al. 2011; Merkin et al. 2012). Kallisto was used to estimate transcript abundances (with options --bias and, for Brawand et al. data, --single -l 275 -s 15), supplying Ensembl version 98 transcript annotations for chimpanzee (Pan\_tro\_3.0), rhesus macaque (Mmul\_10), and chicken (GRCg6a) and the GENCODE vM23 Basic annotation for mouse. For species with an intact, Y-linked *EIF1AY* ortholog (chimpanzee, rhesus), the cDNA sequences of *EIF1AX* and *EIF1AY* orthologs were aligned, and the well-aligned portion of each sequence was inserted into the annotation in place of the annotated sequence(s) listed by Ensembl. This was done to prevent differences in the completeness or correctness of the annotated X- and Y-linked sequences from skewing estimated Y/X expression ratios. After transcript abundance estimation, Y/X expression ratios were calculated in each tissue sample. For the expression patterns shown in Supplemental Figure S12, the expression levels from the samples of a given species were adjusted using the housekeeping normalization method described in Methods; the expression level of *EIF1AX* in each sample was then divided by *EIF1AX*'s median expression level observed among all samples from that species.

#### Analysis of protein abundance in human heart tissue

#### Selection of human heart tissue samples for protein quantification

GTEx heart (left ventricle) samples from 21 male donors and 12 female donors were selected for quantitative proteomic analysis after thoroughly screening all left ventricle samples by donor medical history and histopathological analysis. The two goals of this screening process were to identify samples with minimal pathology and to minimize differences between XX and XY samples (e.g., adiposity, fibrosis, hypertrophy) that might introduce biases. In the first round of screening, the pathology notes released by the The GTEx Consortium (SMPTHNTS) were reviewed, and samples were excluded if the notes indicated >5% adipose tissue, more than "minimal" fibrosis or hypertrophy, or evidence of infarction, ischemia, or myocarditis. Next, samples were excluded based on donor medical history and circumstances of death. Specifically, a sample was excluded if the first underlying cause (DTHFUCOD) or immediate cause (DTHCOD) of death primarily affected the heart or cardiovascular system (e.g., cardiac arrest, myocardial infarction, cardiovascular disease), or if the donor had a recorded history of myocardial infarction (MHHRTATT), heart disease (MHHRTDIS, MHHRTDISB). Finally, the 56 remaining samples underwent a further round of expert histological review using the histology

images available on the GTEx portal (<u>https://gtexportal.org/</u>). Each sample was scored for the content of adipose tissue and interstitial fibrosis, and for the degree of myocyte hypertrophy. Samples were excluded if they showed >3% adipose tissue; >2% fibrosis; or moderate myocyte hypertrophy in combination with borderline adipose and/or borderline fibrosis. Tissue from the remaining 33 remaining samples (21 male, 12 female) was obtained from the GTEx biobank (Supplemental File S3).

#### Human heart tissue proteomics: data generation

A total of 33 human heart samples (left ventricle sampled 1 cm above apex), 21 males and 12 females, stored in PAXgene at –80 °C, were obtained from Gene Expression Tissue (GTEx) Biobank. Samples were rinsed from the PAXgene buffer by ice-cold PBS, pulverized in 1.5ml RIPA lysis buffer with Roche complete protease inhibitors, and sonicated for 2 min using 0.5 pulses.

Following the procedure outlined in Chick et al. (Chick et al. 2016), heart samples (~40mg) were reduced with 5mM dithiothreitol (DTT) for 30min at 54°C followed by alkylation with 20mM iodoacetamide for 30min at room temperature in the dark. The alkylation reaction was guenched by adding 15mM DTT for 15min at room temperature in the dark. A 200µl sample aliquot was then methanol/chloroform precipitated. The samples were allowed to air dry before being resuspended in 300 µl of 8 M urea buffer supplemented with 50 mM Tris at pH 8.2. The urea concentration was diluted down to ~1.5M urea with 50mM Tris. Proteins were quantified using a BCA assay. Protein was then digested using a combination of Lys-C/trypsin at an enzyme-to-protein ratio of 1:100. First, protein was digested overnight with Lys-C followed by 6-h digestion with trypsin all at 37°C. Samples were then acidified using formic acid to approximately pH 3. Samples were desalted using a SepPak column, and eluents were dried using a vacuum centrifuge. Peptide pellets were resuspended in 110µl of 200mM HEPES buffer, pH 8, and peptides were quantified by a BCA assay. Approximately 70µg of peptides  $(100 \mu l of sample + 30 \mu l of 100\%$  acetonitrile) were then labeled with 15 \mu l of 20  $\mu g \mu l^{-1}$  of the corresponding TMT 11-plex reagent for 2h at room temperature. The reaction was quenched using 8µl of 5% hydroxylamine for 15min. Peptides were then acidified using 150µl of 1% formic acid, each set of 11 samples was mixed and desalted using a SepPak column. In total, 3 TMT 11-plex reactions were performed to analyze all 33 samples. The full labeling scheme for the heart samples is provided in Figure 6-source data.

Each of the 3 TMT experiments was separated by basic, reversed-phase chromatography. Samples were loaded onto an Agilent 300 Extend C18 column (5µm particles, 4.6mm ID and 220mm in length). Using an Agilent 1100 quaternary pump equipped with a degasser and a photodiode array detector (set at 220- and 280-nm wavelength), peptides were separated using a 50min linear gradient from 18% to 40% acetonitrile in 10mM ammonium bicarbonate, pH 8, at a flow rate of 0.8ml min–1. Peptides were separated into a total of 96 fractions that were consolidated into 24. Samples were subsequently acidified with 1% formic acid and vacuum centrifuged to near dryness. Each fraction was desalted via StageTip, dried via vacuum centrifugation, and reconstituted in 1% formic acid for liquid chromatography tandem mass spectrometry (LC–MS/MS) processing.

Peptides from every odd fraction (12 fractions total) from basic reverse-phase fractionation were analysed using an Orbitrap Fusion Tribrid mass spectrometer (Thermo Scientific) equipped with a Proxeon ultra high pressure liquid chromatography unit. Peptide mixtures were separated on a 100 $\mu$ m ID microcapillary column packed first with ~0.5cm of 5 $\mu$ m Magic C18 resin followed by 40cm of 1.8 $\mu$ m GP-C18 resin. Peptides were separated using a 3-h gradient of 6–30% acetonitrile gradient in 0.125% formic acid with a flow rate of ~400nl min–1. In each data collection cycle, one full MS scan (400–1,400 m/z) was acquired in the Orbitrap (1.2×105 resolution setting and an automatic gain control (AGC) setting of 2×105). The subsequent MS2–MS3 analysis was conducted with a top 10 setting or a top speed approach using a 2-s duration. The most abundant ions were selected for fragmentation by collision induced dissociation (CID). CID was performed with a collision energy of 35%, an AGC setting of 4×103, an isolation window of 0.5 Da, a maximum ion accumulation time of 150ms and the rapid ion trap setting. Previously analyzed precursor ions were dynamically excluded for 40s.

During the MS3 analyses for TMT quantification, precursors were isolated using a 2.5-Da m/z window and fragmented by 35% CID in the ion trap. Multiple fragment ions (SPS ions) were co-selected and further fragmented by HCD. Precursor ion selection was based on the previous MS2 scan and the MS2–MS3 was conducted using sequential precursor selection (SPS) methodology. HCD used for the MS3 was performed using 55% collision energy and reporter ions were detected using the Orbitrap with a resolution setting of 60,000, an AGC setting of 50,000 and a maximum ion accumulation time of 150 ms. Human heart tissue proteomics: peptide quantification and protein-abundance estimation Software tools were used to convert mass spectrometric data from raw file to the mzxml format (Huttlin et al. 2015). Erroneous charge state and monoisotopic m/z values were corrected as per previous publication (Huttlin et al. 2015). MS/MS spectra assignments were made with the Sequest algorithm (Eng et al. 1994) using an indexed Ensembl database (Ensembl version GRCh37.61). Databases were prepared with forward and reversed sequences concatenated according to the target-decoy strategy (Elias and Gygi 2007). All searches were performed using a static modification for cysteine alkylation (57.0215 Da) and TMT on the peptide N termini and lysines. Methionine oxidation (15.9949 Da) was considered a dynamic modification. Mass spectra were searched with trypsin specificity using a precursor ion tolerance of 10 p.p.m. and a fragment ion tolerance of 0.8 Da. Sequest matches were filtered by linear discriminant analysis as described previously, first to a data set level error of 1% at the peptide level based on matches to reversed sequences (Elias and Gygi 2007). Peptide probabilities were then multiplied to create protein rankings and the data set was again filtered to a final data set level error of 1% false discovery rate (FDR) at the protein level. The final protein-level FDR fell well below 1% (~0.22% peptide level).

Peptide quantitation using TMT reporter ions was accomplished as previously published (Ting et al. 2011; McAlister et al. 2014). In brief, a 0.003 Da m/z window centered on the theoretical m/z value of each reporter ion was monitored for each of the 11 reporter ions, and the intensity of the signal closest to the theoretical m/z value was recorded. TMT signals were also corrected for isotope impurities based on the manufacturer's instructions. Peptides were only considered quantifiable if the total signal-to-noise for all channels was >200 with an isolation specificity of >0.75.

For proteins not encoded by X–Y gene pairs, peptides were assigned to protein matches using a reductionist model, where all peptides were explained using the least number of proteins. The signal-to-noise values in each channel were then divided by the sum of all signal-to-noise values in that channel, such that each channel had the same summed value. Protein quantitation was then performed by summing the signal-to-noise values for all peptides for a given protein. Within each 11-plex TMT experiment, protein quantitative measurements were then scaled to 100, such that equal expression across all channels would be equal to  $100/11 \approx 9.1$ .

Protein abundances of the X and Y isoforms were estimated separately as follows. Within each 11-plex experiment, raw signal-to-noise values in each of the 11 channels were first normalized (divided) by the summed signal/noise value for all peptides in that channel:  $\tilde{y}_{ij} = y_{ij} / \sum_{p \in P} y_{pj}$ , where  $y_{ij}$  is the raw signal/noise value for peptide *i* in channel *j*,  $\tilde{y}_{ij}$  is the channel-normalized signal/noise value, and *P* is the set of all detected peptides. Among all detected peptides, we then identified those that specifically matched the amino-acid sequence of an X homolog of an X–Y pair, of a Y homolog of an X–Y pair, or both X and Y homologs of an X–Y pair but not the sequence of any other protein. We detected all three classes (X-specific, Y-specific, X-Y-shared) of peptides for RPS4Y1/RPS4X, EIF1AY/EIF1AX, and DDX3Y/DDX3X, the three most highly expressed X–Y pair genes. (We further detected X-specific peptides for USP9X/USP9Y but no Y-specific or X-Y-shared peptides.) Y-specific peptides showed low but roughly constant signal in female channels, with mean signal/noise = 2.57: we used this value as an estimate of the non-specific background for all peptides and subtracted this value from all peptides in all channels, setting any negative values to 0. For each of the three X–Y pairs separately, we then estimated the relative abundance of the X isoform in channel *j*,  $a_j^{(X)}$ , as the percentage of signal from all X-specific peptides in channel j out of the total signal across all channels,

$$a_j^{(\mathrm{X})} = \frac{\sum_{p \in P_X} \tilde{y}_{pj}}{\sum_{k=1}^{11} \sum_{p \in P_X} \tilde{y}_{pk}} * 100,$$

where  $P_X$  is the set of all X-specific peptides for the given X–Y pair. We repeated this calculation using X–Y-shared peptides to obtain  $a_j^{(XY)}$ , the relative abundance of the sum of X and Y isoforms in channel *j*. To obtain the male-to-female expression ratio for the X homolog specifically,  $\phi_{MF}^{(X)}$ , we pooled the abundance estimates across all three 11-plex experiments and divided its average abundance in male channels by its average abundance in female channels, i.e.,

$$\phi_{\rm MF}^{(\rm X)} = \frac{\frac{1}{21} \sum_{j \in C_M} a_j^{(\rm X)}}{\frac{1}{12} \sum_{j \in C_F} a_j^{(\rm X)}},$$

where  $C_M$  and  $C_F$  are the sets of channels from male and female donors, respectively. The male-to-female expression ratio for the sum of X and Y homolog expression,  $\phi_{MF}^{(XY)}$ , was obtained similarly. *p*-values for the sex bias in expression were estimated by permuting the sample labels within each 11-plex experiment one million times and calculating the proportion of permutations that yielded more extreme male-to-female expression ratios. (This procedure was also used to estimate *p*-values for the sex bias of non-X–Y-pair proteins, shown in

Supplemental Fig. S18.) Finally, the Y-to-X expression ratio within males was estimated as  $(\phi_{MF}^{(XY)} - \phi_{MF}^{(X)})/\phi_{MF}^{(X)}$ .

## Immunoblotting

To prepare human heart lysates, 40 mg of human heart tissue (obtained originally for the mass spectrometry analysis) was rinsed twice with 2 ml of ice-cold PBS and pulverized in 1500 µl of RIPA buffer with protease inhibitor cocktail (Roche, Catalog number: 11836170001). Lysates were incubated 30 min on ice and spun at 10,000 g for 20 min at 4 °C, and the supernatant was collected. Aliquots of the heart lysates were pooled by sex (21 male samples, 12 female samples), mixed with NuPAGE Sample Reducing Agent 10x (Invitrogen, #NP0009) and NuPAGE LDS Sample Buffer 4x (Invitrogen, #NP0007), incubated for 10 min at 90 °C, and then chilled on ice for 2 min. Proteins were separated for 3.5 hr at 80 V on a NuPAGE 4-12% Bis-Tris gel (Invitrogen, #NP0322BOX) and transferred to a nitrocellulose membrane. The membrane was blocked for 1 hr in Tris-buffered saline containing 0.1% Tween-20 (TBST) and 5% non-fat milk at room temperature, and then incubated with primary antibodies in TBST overnight at 4 °C on a shaking platform (GAPDH: Ambion AM4300, anti-mouse, 1:106 dilution; eIF-1A: Abcam Ab177939, anti-rabbit, 1:5000 dilution). The monoclonal eIF-1A antibody was generated using a proprietary synthetic peptide within amino-acids 50 - 144 of human EIF1AX. After three washes with TBST, the membrane was incubated at room temperature for 1 hr with TBST and 1% milk containing anti-mouse and anti-rabbit secondary antibodies labeled with fluorescent dyes detectable at different wavelengths (LI-COR IRDye 680RD Goat anti-Mouse 925-68070, 1:20000 dilution; LI-COR IRDye 800CW Goat anti-Rabbit 926-32211, 1:20000 dilution). Following three further washes in Tris-buffered saline lacking Tween-20, fluorescent signal was recorded using an Odyssey CLx imager (LI-COR) with Image Studio software (version 5.2.5). Fluorescent signal corresponding to EIF1A was normalized to signal for GAPDH in each lane, with four technical replicates (i.e., lanes) per sex.

Lymphoblastoid cell lines (LCLs) from individuals with sex chromosome aneuploidies were derived in our lab, except the following that were previously reported: 47,XYY (Repping et al. 2003), 49,XYYYY (GM11419, from Coriell) and Sirota et al. (Sirota et al. 1981). Cells were pelleted, washed in ice-cold PBS, and snap frozen in liquid nitrogen. Approximately 5 million cells per line were lysed in 250  $\mu$ I M-PER lysis buffer (Thermo Scientific #78503) supplemented with protease inhibitor. Lysates were incubated on ice for 15 min and spun for at 14,000 rpm for 15 min at 4 °C, and the supernatant was collected. LCL lysates were mixed with sample

13

reducing agent and sample buffer, incubated for 10 min at 90 °C, and chilled on ice for 2 min. Proteins were separated for 3 hr at 80 V on a NuPAGE 4-12% Bis-Tris gel and transferred to a nitrocellulose membrane. Primary-antibody incubation was performed as described above. Membranes were subsequently incubated with TBST and 1% milk containing peroxidaseconjugated secondary antibodies at 1:5000 dilution (Jackson ImmunoResearch Peroxidase AffiniPure Donkey Anti-Mouse 715-035-151 & Anti-Rabbit 711-035-152) for 1 hr at room temperature. Following three washes with TBST, proteins on the membranes were detected by addition of Lumi-Light Western Blotting Substrate (Roche 12015200001).

# SUPPLEMENT REFERENCES

- Agarwal V, Bell GW, Nam J-W, Bartel DP. 2015. Predicting effective microRNA target sites in mammalian mRNAs. *eLife* **4**: e05005. doi:10.7554/eLife.05005
- Brawand D, Soumillon M, Necsulea A, Julien P, Csárdi G, Harrigan P, Weier M, Liechti A, Aximu-Petri A, Kircher M, et al. 2011. The evolution of gene expression levels in mammalian organs. *Nature* **478**: 343–348.
- Bray NL, Pimentel H, Melsted P, Pachter L. 2016. Near-optimal probabilistic RNA-seq quantification. *Nat Biotechnol* **34**: 525–527.
- Chick JM, Munger SC, Simecek P, Huttlin EL, Choi K, Gatti DM, Raghupathy N, Svenson KL, Churchill GA, Gygi SP. 2016. Defining the consequences of genetic variation on a proteome-wide scale. *Nature* **534**: 500–505.
- Elias JE, Gygi SP. 2007. Target-decoy search strategy for increased confidence in large-scale protein identifications by mass spectrometry. *Nat Methods* **4**: 207–214.
- Eng JK, McCormack AL, Yates JR. 1994. An approach to correlate tandem mass spectral data of peptides with amino acid sequences in a protein database. *J Am Soc Mass Spectrom* **5**: 976–989.
- Ferreira PG, Muñoz-Aguirre M, Reverter F, Godinho CPS, Sousa A, Amadoz A, Sodaei R, Hidalgo MR, Pervouchine D, Carbonell-Caballero J, et al. 2018. The effects of death and post-mortem cold ischemia on human tissue transcriptomes. *Nat Commun* **9**: 490. doi:10.1038/s41467-017-02772-x
- Friedman RC, Farh KK, Burge CB, Bartel DP. 2009. Most mammalian mRNAs are conserved targets of microRNAs. *Genome Res* **19**: 92–105.
- Gallego Romero I, Pai AA, Tung J, Gilad Y. 2014. RNA-seq: impact of RNA degradation on transcript quantification. *BMC Biol* **12**: 42. doi:10.1186/1741-7007-12-42
- Hughes JF, Skaletsky H, Brown LG, Pyntikova T, Graves T, Fulton RS, Dugan S, Ding Y, Buhay CJ, Kremitzki C, et al. 2012. Strict evolutionary conservation followed rapid gene loss on human and rhesus Y chromosomes. *Nature* **483**: 82–86.

- Huttlin EL, Ting L, Bruckner RJ, Gebreab F, Gygi MP, Szpyt J, Tam S, Zarraga G, Colby G, Baltier K, et al. 2015. The BioPlex Network: A Systematic Exploration of the Human Interactome. *Cell* **162**: 425–440.
- Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL. 2013. TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions. *Genome Biol* 14: R36. doi:10.1186/gb-2013-14-4-r36
- Langmead B, Trapnell C, Pop M, Salzberg SL. 2009. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol* **10**: R25. doi:10.1186/gb-2009-10-3-r25
- Li B, Dewey CN. 2011. RSEM: accurate transcript quantification from RNA-seq data with or without a reference genome. *BMC Bioinformatics* **12**: 323. doi:10.1186/1471-2105-12-323
- Liao Y, Smyth GK, Shi W. 2014. featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics* **30**: 923–930.
- Löytynoja A, Goldman N. 2005. An algorithm for progressive multiple alignment of sequences with insertions. *Proc Natl Acad Sci* **102**: 10557–10562.
- Ludwig N, Leidinger P, Becker K, Backes C, Fehlmann T, Pallasch C, Rheinheimer S, Meder B, Stähler C, Meese E, et al. 2016. Distribution of miRNA expression across human tissues. *Nucleic Acids Res* **44**: 3865–3877.
- McAlister GC, Nusinow DP, Jedrychowski MP, Wühr M, Huttlin EL, Erickson BK, Rad R, Haas W, Gygi SP. 2014. MultiNotch MS3 enables accurate, sensitive, and multiplexed detection of differential expression across cancer cell line proteomes. *Anal Chem* **86**: 7150–7158.
- Merkin J, Russell C, Chen P, Burge CB. 2012. Evolutionary Dynamics of Gene and Isoform Regulation in Mammalian Tissues. *Science* **338**: 1593–1599.
- Pruitt KD, Harrow J, Harte RA, Wallin C, Diekhans M, Maglott DR, Searle S, Farrell CM, Loveland JE, Ruef BJ, et al. 2009. The consensus coding sequence (CCDS) project: Identifying a common protein-coding gene set for the human and mouse genomes. *Genome Res* **19**: 1316–1323.
- Repping S, Skaletsky H, Brown L, Daalen SKM van, Korver CM, Pyntikova T, Kuroda-Kawaguchi T, Vries JWA de, Oates RD, Silber S, et al. 2003. Polymorphism for a 1.6-Mb deletion of the human Y chromosome persists through balance between recurrent mutation and haploid selection. *Nat Genet* **35**: 247–251.
- Sirota L, Zlotogora Y, Shabtai F, Halbrecht I, Elian E. 1981. 49, XYYYY. A case report. *Clin Genet* **19**: 87–93.
- Skaletsky H, Kuroda-Kawaguchi T, Minx PJ, Cordum HS, Hillier L, Brown LG, Repping S, Pyntikova T, Ali J, Bieri T, et al. 2003. The male-specific region of the human Y chromosome is a mosaic of discrete sequence classes. *Nature* **423**: 825–837.
- Ting L, Rad R, Gygi SP, Haas W. 2011. MS3 eliminates ratio distortion in isobaric multiplexed quantitative proteomics. *Nat Methods* **8**: 937–940.

- Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, Baren MJ van, Salzberg SL, Wold BJ, Pachter L. 2010. Transcript assembly and quantification by RNA-seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol* **28**: 511–515.
- Wright FA, Sullivan PF, Brooks AI, Zou F, Sun W, Xia K, Madar V, Jansen R, Chung W, Zhou Y-H, et al. 2014. Heritability and genomics of gene expression in peripheral blood. *Nat Genet* **46**: 430–437.

# SUPPLEMENTAL FIGURES

Supplemental Fig. S1



Supplemental Fig. S1. Discarding multi-mapping reads disproportionately underestimates MSY gene expression. (A) Bars show the average proportion of reads from genes on each chromosome that can be aligned uniquely. Chromosome "Y" refers to genes in the MSY; chromosome "X" similarly excludes genes in the pseudoautosomal region. (B) Each point shows the expression level of one gene in the testis as estimated without multi-mapping reads (GTEx Consortium, via RNA-SeQC) and with multi-mapping reads (our study, via kallisto). Each point is colored according to the proportion of reads from that gene that can be aligned uniquely (gene mappability).



Supplemental Fig. S2. Kallisto accurately estimates MSY gene expression levels in simulated RNA-seq datasets. Each point shows the expression level of an MSY gene/gene family (A – C) or the Y/X expression ratio (D) in a simulated RNA-seq dataset, estimated by a method that discards multi-mapping reads (gray), Cufflinks in multi-correct "-u" mode (blue), or kallisto (red). 50 simulated RNA-seq datasets were generated for each of three scenarios: (A) MSY genes not expressed (0 TPM) and X homologs of MSY genes kept at levels in sample GTEX-P4QS-2126-SM-3NMCF; (B, D) MSY genes set to 1 TPM and X homologs set to 2 TPM; (C) MSY genes set to 5 TPM and X homologs set to 10 TPM.



**Supplemental Fig. S3 (cont.) Expression patterns of individual MSY genes and gene families.** For each gene, blue bars show the gene's median expression level across the samples of each tissue; error bars: 5<sup>th</sup> and 95<sup>th</sup> percentiles.





**Supplemental Fig. S4. Expression of the DAZ gene family in the stomach.** (**A**) Each point shows the expression level of the DAZ gene family (sum of DAZ1, DAZ2, DAZ3, DAZ4) in a single tissue sample from an XY (blue) or XX (gray) donor. For skin, stomach, and testis, lines show the 25<sup>th</sup>, 50<sup>th</sup> (median), and 75<sup>th</sup> percentiles. The absence of DAZ expression in XX donors suggests that DAZ's stomach expression is not the result of mis-mapped reads from a gene on another chromosome. (**B**) The DAZ gene family arose from transposition of DAZL on Chr 3 to the Y Chromosome. Each point shows DAZL's expression in a single sample from XY or XX donors. The absence of DAZL expression in stomach samples suggests that DAZ acquired stomach expression after transposition to the Y Chromosome.





Supplemental Fig. S5. *NLGN4Y* and *SRY*, unlike other MSY genes, show broader expression than their X-linked homologs, leading to male-specific expression of *NLGN4X/Y* and *SOX3/SRY* in some tissues. (A) Expression levels of *NLGN4X* in XX donors (left) or the summed expression of *NLGN4X* and *NLGN4Y* in XY donors (right). (B) Expression levels of *SOX3* in XX donors (left) or the summed expression of *SOX3* and *SRY* in XY donors (right). In both panels, tissues are ordered by the fraction of expression coming from the Y-linked homolog. Tissues in gray text are those where the total expression of the X–Y pair is less than 1 TPM in both sexes.



**Supplemental Fig. S6. Comparison of X- and Y homolog expression breadth in tissues from XY donors.** Each cell shows the expression levels of the X-linked (upper left triangle) and Y-linked (lower right triangle) homologs of one X–Y gene pair in one tissue in XY donors.





**Supplemental Fig. S7. Replication of estimated Y/X expression ratios in RNA-seq data from the Human Protein Atlas (HPA).** (**A** – **B**) Mean Y/X ratios for 9 widely expressed X–Y pairs in GTEx (*x*-axis; error bars: 5<sup>th</sup> – 95<sup>th</sup> percentiles) and HPA samples (*y*-axis; error bars: min, max among samples) from (**A**) colon or (**B**) prostate. (**C**) Mean Y/X ratios for all homologous X–Y pairs in GTEx and HPA testis samples.



**Supplemental Fig. S8. Y/X expression ratios do not differ substantially between tissues.** Each point shows the estimated Y/X expression ratio of one X–Y pair in one tissue. Tissues are ordered by the average Y/X expression ratio across all pairs.



# Supplemental Fig. S9 (cont.). Co-expression of the X- and Y-linked members of X-Y gene pairs in

**individual tissues.** Each point shows the expression levels of the X- and Y-linked members of an X–Y gene pair in a single tissue sample. Each plot corresponds to a cell of the heatmap shown in Figure 3C. This subset of tissues (8/36) was selected to showcase a diversity of tissue types and cases where X–Y pairs are correlated and uncorrelated in expression.



**Supplemental Fig. S10. The expression levels of MSY genes and their corresponding X-linked homologs are independently estimated.** 50 simulated RNA-seq libraries were generated. In each, the expression level of the MSY gene/gene family was set to 5 TPM and its X-linked homolog was set to a random value between 0 and 10 TPM. The estimated expression levels of the X and Y homologs in each library are shown as a pair of red (X) and blue (Y) points. MSY genes were correctly estimated at ~5 TPM irrespective of the expression levels of their X-linked homologs.





**Supplemental Fig. S11. Similarity of EIF1AX and EIF1AY proteins.** Amino-acid sequence of human EIF1AX (positions 1 – 75 of 144) aligned with sequences of human EIF1AY and EIF1AX/EIF1AY homologs around the single position at which human EIF1AX and EIF1AY differ (position 50, yellow). "." indicates identity to human EIF1AX, shown at top. Position 50 is the only position in this region at which amino-acid substitutions are found among vertebrates.

Supplemental Fig. S12



Supplemental Fig. S12. Replication of EIF1AX/EIF1AY expression pattern across human tissues using Human Protein Atlas (HPA) RNA-seq data. (A) Y/X expression ratio for *EIF1AX/Y* in individual XY tissue samples from the HPA dataset. Highlighted box (gray) shows tissues with notable X–Y expression divergence (e.g., as shown in Fig. 3). (B) Y/X expression ratio for *EIF1AX/Y* in each of 16 tissues estimated using data from GTEx (*x*-axis) and HPA (*y*-axis). (C) Expression of *EIF1AX* (tan) in XX samples (left) vs. summed expression of *EIF1AX* (tan) and *EIF1AY* (blue) in XY samples (right) in HPA tissue samples collected from both sexes. Expression is normalized to mean expression level in XX samples from each tissue; error bars show min and max observed values. For each tissue, the number of samples from female and male samples, respectively, are given in parentheses. For comparison with *EIF1AX/Y*, similar plots are shown for *DDX3X/Y* (D – **F**) and *ZFX/Y* (**G** – **I**); figure legends as in **A** – **C**.



Supplemental Fig. S13. Expression patterns of *EIF1AX* orthologs, which retain intact miR-1 target sites.

Each point shows the expression level of an *EIF1AX* ortholog in a single tissue sample collected from a male donor. Within each species, expression levels are normalized to the median expression level observed across the samples. Fewer points are shown for chimpanzee because this species was not analyzed in Merkin et al. 2012.





**Supplemental Fig. S14. Highly co-expressed X–Y gene pairs show little sex-biased expression.** Each point compares the degree to which the X and Y homologs of an X–Y pair show correlated expression in one tissue (*x*-axis) vs. the sex-biased expression of the X–Y pair in the same tissue (*y*-axis). The degree of correlated expression between gene A and gene B is measured as the proportion of genes in the genome more correlated in expression with gene A than gene B (Methods). Highly co-expressed X–Y pairs—in which the X and Y homologs are likely tightly co-regulated—show weak sex-biased expression. By contrast, when the expression of a pair of X and Y homologs is uncorrelated—indicating their regulation has likely diverged—more prominent sex-biased expression is observed.





# Supplemental Fig. S15. Strategy for estimating X and Y isoform expression from multiplexed

**proteomics data.** 12 XX and 21 XY heart (left ventricle) tissue samples were analyzed by mass spectrometry in three 11-plex experiments, consisting of 4 XX and 7 XY samples each. Within each 11-plex experiment, isobaric tandem mass tags (TMTs) are used as barcodes to quantify the relative abundance of a given peptide in each sample. Peptides from the X and Y protein isoforms of a given X–Y pair can match the sequences of both proteins (dark gray; X–Y-shared) or can be specific to the X (red) or Y (blue) isoform. Y-specific peptides are used to confirm the presence of the Y isoform in the sample. X-specific and X–Y-shared peptides are used to assess sex biases in expression; information from these two peptide classes is then integrated to infer the relative contribution of X and Y isoforms to overall expression in XY individuals.

A DDX3X/Y amino-acid sequence number of sets in which peptide was detected MSHV<sub>V</sub>V<sub>K</sub>N<sub>DPE</sub>LDQQ<sub>L</sub>A<sub>N</sub>GLDLNS<sub>EKQSG</sub>G<sub>A</sub>STASKGRYIPPHLRNREA<sub>S</sub>KGF<sub>H</sub>DKDSSGWS<sub>C</sub>SKDKDAYSS FGSR<sub>DS</sub>RGK<sub>PGY</sub>FS<sub>E</sub>RGSGSRGRFDDRG<sub>RSDYDGIG<sub>N</sub>R<sub>D</sub>PR<sub>D</sub>GFG<sub>K</sub>FER<sub>S</sub>G<sub>h</sub>SRwcD<sub>x</sub>S<sub>v</sub>EDDWSKPL<sub>P</sub> serLEQELFSGGNTGINFEKYDDIPVEATG<sub>N</sub>NCPPHIE<sub>N</sub>FSD<sup>VE</sup>MGEIIMGNIELTRYTRPTPVQKHAIPT IK<sub>G</sub>KRDLMACAQTGSGKTAAFLLPILSQIY<sub>T</sub>DGPGEAL<sub>K</sub>AMKENGRYGRRKQYPISLVLAPTRELAVQIYE EARKFSYRS<sup>1</sup> GFEPQIR<sup>2</sup>RIVEQDTMPPKGVRHTMMFSATFPKEIQMLARDFLDEYIFLAVGRVGSTSENITQKVWVVE<sub>DL</sub>D KRSFLLD<sub>L</sub>C<sub>0</sub><sup>8</sup>ATG<sub>S</sub>DSLTLVFVETKKGADSLEDFLYHEGYACTSIHGDRSQ<sup>1</sup>RDREEALHQFRSGKSPILVA TAVAARGLDISNV<sub>R</sub><sup>4</sup>HYINFDLPSDIEEYVHRIGRTG<sup>2</sup>RVGNLGLATSFFNE<sub>K</sub>N<sub>M</sub>NITKDLLDLLVEAKQEVP SWLENMAYEHHYKG<sup>6</sup>GSRGRSKS<sup>3</sup>RFSGGFGARD</mark>YRQSSG<sup>6</sup>SSS<sup>5</sup>G<sup>5</sup>SSR<sup>6</sup>GSSRSGGGG<sup>4</sup>G<sup>6</sup>SSRGFGGGGGY GGFYNSDGYGGNYNSOGVDWWGN</sub> B RPS4X/Y1 amino-acid sequence

**Supplemental Fig. S16. Recovery of peptides from DDX3X, DDX3Y, RPS4X, and RPS4Y1.** Amino-acid sequence of DDX3X/Y (**A**) and RPS4X/Y1 (**B**): X- and Y-specific amino acids are superscripted and subscripted, respectively. X-specific (gold), Y-specific (blue), and X–Y shared (purple) peptides detected by mass spectrometry are shown, along with the number of 11-plex experiments in which each peptide was detected.



Supplemental Fig. S17. Peptides from highly transcribed genes were more often detected by mass spectrometry. (A) The median transcript expression level (TPM) in heart (left ventricle) samples from XY donors is given for the X- and Y-linked homologs of the 9 most widely expressed X–Y gene pairs. Peptides were detected for the seven genes that had the highest transcript expression levels. (B) All genes genomewide were grouped into eight bins based on their transcript expression level in the heart. The percentage of genes in each bin with  $\geq$ 1 peptide detected by mass spectrometry is shown in pink.



**Supplemental Fig. S18. XX and XY expression of non-X–Y pair proteins.** Relative expression of five non-X–Y-pair proteins in XX (gray) and XY (green) samples, used as negative controls. Each protein corresponds to an X-Chromosome gene that is subject to X inactivation in XX cells and is not sex-biased at the transcript level.



**Supplemental Fig. S19. eIF-1A antibody recognizes both EIF1AX and EIF1AY proteins.** Western blots showing signal detected by an eIF-1A antibody in protein lysates prepared from human lymphoblastoid cell lines with various numbers of X (**A**) or Y (**B**) chromosomes. GAPDH was used as a loading control. Because *EIF1AX* escapes X-Chromosome inactivation (XCI), EIF1AX expression should increase in cells with additional numbers of X Chromosomes. Because the Y Chromosome is not subject to a program of epigenetic silencing analogous to XCI, expression of EIF1AY should also increase in cells with additional Y Chromosomes. As shown, signal detected by the eIF-1A antibody increases in cells with additional X Chromosomes (**A**) and in cells with additional Y Chromosomes (**B**), implying the eIF-1A antibody recognizes both EIF1AX and EIF1AY.

37