

1 **Title:**

2

3 Identifying, understanding, and correcting technical biases on the sex chromosomes in
4 next-generation sequencing data

5

6 **Authors and Affiliations:**

7

8 Timothy H. Webster¹, Madeline Couse², Bruno M. Grande³, Eric Karlins⁴, Tanya N.
9 Phung⁵, Phillip A. Richmond^{6,7}, Whitney Whitford⁸, Melissa A. Wilson Sayres^{1,9}

10

11 ¹School of Life Sciences, Arizona State University

12 ²Child and Family Research Institute, University of British Columbia

13 ³Department of Molecular Biology and Biochemistry, Simon Fraser University

14 ⁴Division of Cancer Epidemiology and Genetics, National Cancer Institute, National
15 Institutes of Health

16 ⁵Interdepartmental Program in Bioinformatics, UCLA

17 ⁶Centre for Molecular Medicine and Therapeutics, University of British Columbia

18 ⁷BC Children's Hospital

19 ⁸School of Biological Sciences, The University of Auckland

20 ⁹Center for Evolution and Medicine, Arizona State University

21

22 **Corresponding Authors:**

23

24 Timothy H. Webster

25 School of Life Sciences

26 Arizona State University

27 Tempe, AZ, USA 85281

28 Timothy.h.webster@asu.edu

29

30 Melissa A. Wilson Sayres

31 School of Life Sciences

32 Arizona State University

33 Tempe, AZ, USA 85281

34 melissa.wilsonsayres@asu.edu

35

36

37 **Abstract**

38 Mammalian X and Y chromosomes share a common evolutionary origin and retain
39 regions of high sequence similarity. This sequence homology can cause the mismapping
40 of short sequencing reads derived from the sex chromosomes and affect variant calling
41 and other downstream analyses. Understanding and correcting this problem is critical for
42 medical genomics and population genomic inference. Here, we characterize how
43 sequence homology can affect analyses on the sex chromosomes and present XYalign, a
44 new tool that: (1) facilitates the inference of sex chromosome complement from next-
45 generation sequencing data; (2) corrects erroneous read mapping on the sex
46 chromosomes; and (3) tabulates and visualizes important metrics for quality control such
47 as mapping quality, sequencing depth, and allele balance. We show how these metrics
48 can be used to identify XX and XY individuals across diverse sequencing experiments,
49 including low and high coverage whole genome sequencing, and exome sequencing. We
50 also show that XYalign corrects mismapped reads on the sex chromosomes, resulting in
51 more accurate variant calling. Finally, we discuss how the flexibility of the XYalign
52 framework can be leveraged for other use cases including the identification of aneuploidy
53 on the autosomes. XYalign is available open source under the GNU General Public
54 License (version 3).

55

56 **Keywords**

57 X chromosome; Y chromosome; ploidy; aneuploidy; genomics; variant calling; mapping

58 Introduction

59 Accurate genotyping and variant calling are priorities in medical genetics,
60 including molecular diagnostics, and population genomics (Taylor *et al.*, 2015; Ashley,
61 2016). Despite the availability of numerous powerful tools developed to infer genotypes
62 from sequencing data, sequence homology among genomic regions still presents a major
63 challenge to genome assembly, short read mapping, and variant calling. Specifically,
64 similar sequence content can confound the mapping of short next-generation sequencing
65 reads to a reference genome and lead to technical artifacts in downstream analyses and
66 applications. Heteromorphic sex chromosomes, in particular, present a case of sequence
67 homology likely to affect all individuals in a given species.

68 Sex chromosomes in therians—the clade containing eutherian mammals and
69 marsupials—share a common evolutionary origin as a pair of homologous autosomes
70 (Glas *et al.*, 1999). Approximately 180 to 210 million years ago, they began
71 differentiating from each other through a series of recombination suppression events and
72 subsequent gene loss on the Y chromosome (Rens *et al.*, 2007; Lahn and Page, 1999;
73 Livernois *et al.*, 2012; Wilson Sayres and Makova, 2013). However, this pattern is not
74 unique to mammalian evolution or even XX/XY systems, and occurs often across taxa
75 with genetic sex determination (Bergero and Charlesworth, 2009; Wilson and Makova,
76 2009). This shared origin and complex history characteristic of sex chromosomes lead to
77 unique challenges for genome assembly and analysis, including large blocks of
78 homologous sequence between the sex chromosomes—called gametologous sequence—
79 that we hypothesize can lead to the mismapping of reads between the sex chromosomes.
80 Further, the sex chromosomes of many species contain pseudoautosomal regions (PARs;
81 of which humans have two: PAR1 and PAR2)—regions that have not differentiated
82 between the sex chromosomes and are identical in sequence between the two sex
83 chromosomes (Simmler *et al.*, 1985; Ross *et al.*, 2005). A reference genome that includes
84 the entire sequence content from both sex chromosomes will thus duplicate the PARs and
85 substantially reduce mapping quality in these regions because most reads will identically
86 map to two regions in the reference assembly. This stands in contrast to autosomal
87 sequence, for which each diploid autosome is represented just once in the reference
88 genome.

89 The technical challenges presented by the biological realities of the sex
90 chromosomes can lead to erroneous genotype calls, so the sex chromosomes are routinely
91 excluded from genome-wide analyses (e.g., Wise *et al.*, 2013). This is unfortunate
92 because the sex chromosomes contribute to phenotype and disease etiology (e.g., Chang
93 *et al.*, 2014) and are useful in population genetic inference of demography and patterns of
94 natural selection (Webster and Wilson Sayres, 2016; Wilson Sayres, 2018; Vicoso and
95 Charlesworth, 2006; Ellegren, 2009; Meisel and Connallon, 2013).

96 A number of tools, methods, and frameworks have been developed to aid in the
97 identification of sex-linked sequence (e.g., Muyle *et al.*, 2016), inference of an
98 individual's sex chromosome complement (e.g., Madel *et al.*, 2016), and handling of
99 some of the technical challenges sex chromosomes present in genome-wide association
100 studies (e.g., Gao *et al.*, 2015). However, to our knowledge, there is no tool that
101 simultaneously facilitates the identification of sex chromosome complement and corrects
102 for associated technical biases for the purposes of short read mapping and variant calling.

103 Out of the urgent need to understand the effects of sex chromosome homology on
104 next-generation sequencing analyses, in this paper we first test whether sequence
105 homology between sex chromosomes can confound aspects of read mapping and lead to
106 downstream errors in sequence analysis. We then present XYalign, a tool developed to
107 perform three major tasks: (1) aid in the characterization of an individual's sex
108 chromosome complement; (2) identify and correct for technical artifacts arising from sex
109 chromosome sequence homology; and (3) tabulate and visualize important metrics for
110 quality control such as mapping quality, sequencing depth, and allele balance. We show
111 how XYalign can be used to identify XX and XY individuals across sequencing depths
112 and capture techniques. We also show that the default steps taken by XYalign correct
113 many mismapped reads on the sex chromosomes, resulting in more accurate variant
114 calling. Finally, because XYalign is designed to be both scalable and customizable, we
115 discuss how it can be used in a variety of situations including genetic sex identification in
116 both XX/XY and ZZ/ZW systems, identification of sex-linked sequences and
117 pseudoautosomal regions in new draft genomes, correction of technical biases in genomic
118 and transcriptomic data, detection of aneuploidy, and investigation of mapping success
119 across arbitrary chromosomes.

120

121 **Methods**

122 *Implementation*

123 XYalign is implemented in Python and uses a number of third-party Python
124 packages including Matplotlib (Hunter, 2007), NumPy (Oliphant, 2006), Pandas
125 (McKinney, 2010), PyBedTools (Quinlan and Hall, 2010; Dale *et al.*, 2011), PySam
126 (<https://github.com/pysam-developers/pysam>), and SciPy (Jones *et al.*, 2001). It further
127 wraps the following external tools: repair.sh and shuffle.sh from BBTools
128 (<https://sourceforge.net/projects/bbmap/>), BWA (Li, 2013), Platypus (Rimmer *et al.*,
129 2014), Sambamba (Tarasov *et al.*, 2015), and SAMtools (Li *et al.*, 2009).

130 XYalign is composed of six modules that can be called individually or serve as
131 steps in a full pipeline: PREPARE_REFERENCE, CHROM_STATS, ANALYZE_BAM,
132 CHARACTERIZE_SEX_CHROMS, STRIP_READS, and REMAPPING. Below, we
133 discuss each module as a step in the full XYalign pipeline using human samples (XX/XY
134 sex determination) as an example. Note, however, that XYalign will work with other sex
135 chromosome systems (e.g., ZZ/ZW) and on arbitrary chromosomes (e.g., detecting
136 autosomal aneuploidy). We describe examples of XYalign commands in the section titled
137 "Use Cases."

138 The PREPARE_REFERENCE module generates two versions of the same
139 reference genome: one for the homogametic sex (e.g., XX) and one for the heterogametic
140 sex (e.g., XY). In the simplest case, it will completely hard-mask the Y chromosome with
141 Ns in the XX version of the reference. Optionally, it will also accept one or more BED
142 files containing regions to hard mask in both reference versions. If pseudoautosomal
143 regions (PARs) are present on both sex chromosome sequences in the reference, we
144 strongly suggest masking the PARs on the Y chromosome, allowing reads from these
145 regions to map exclusively to the X chromosome in XY individuals. In XYalign, we use
146 hard masks, rather than omitting the Y chromosome in the XX reference version because
147 these hard masks allow files from both references to share the same sequence dictionaries

148 and indices, thus permitting seamless integration of files from both references into
149 downstream analyses (e.g., joint variant calling).

150 The CHROM_STATS module provides a relatively quick comparison of mapping
151 quality and sequencing depth across one or more chromosomes and over multiple BAM
152 files. While this provides a less detailed perspective than ANALYZE_BAM or
153 CHARACTERIZE_SEX_CHROMS (detailed below), we envision it to be especially
154 useful in at least two different cases. First, in well-characterized systems (e.g., human),
155 comparing chromosome-wide values of mean mapping quality and depth represent a
156 quick and often sufficient way to identify the sex chromosome complement (e.g., XX or
157 XY) of individuals across a population. Second, in uncharacterized systems, the
158 CHROM_STATS output provides information that can help with the identification of
159 sex-linked scaffolds. It is important to note, however, that results for both cases will vary
160 based on ploidy and with differences in the degree of sequence homology between the
161 sex chromosomes.

162 The ANALYZE_BAM module runs a series of analyses designed to aid in the
163 identification of sex-linked sequence and characterize the sex chromosome content of an
164 individual. In doing so, it provides more detailed metrics than CHROM_STATS. For
165 ANALYZE_BAM, XYalign runs Platypus (Rimmer *et al.*, 2014) across multiple threads,
166 if permitted, to identify variants. It then parses the output VCF file containing the
167 variants, applies filters for site quality, genotype quality, and read depth, and plots the
168 read balance at variant sites. Here, we define read balance at a given site as the number of
169 reads containing the alternate allele (i.e., nonreference allele) divided by the total number
170 of reads mapped to the position. XYalign produces plots and tables for read balance per
171 site, as well as mean read balance and variant count per genomic bin or window across a
172 chromosome. We anticipate these data will not only be useful for masking regions
173 containing incorrect genotypes but will also aid in the identification of PARs as well.
174 XYalign next traverses the BAM file, calculating mean mapping quality and an
175 approximation of mean depth in windows across the genome. During traversal, secondary
176 and supplementary read mappings are ignored, and depth is calculated as the total length
177 of all reads mapping to a genomic window divided by the total length of the window. We
178 have found that this heuristic approximation is very similar to calculations of exact depth,
179 particularly as window sizes increase, and is much faster to compute across entire
180 chromosomes. XYalign will output a table containing genomic coordinates, mean depth,
181 and mean mapping quality for each window. It will then filter windows based on user-
182 defined thresholds of mean depth and mapping quality and output two BED files
183 containing windows that passed and failed these thresholds, respectively, which can be
184 used for additional masking in downstream applications. Finally, XYalign will output
185 plots of mapping quality and depth in each window across each chromosome.

186 After running ANALYZE_BAM, the windows meeting thresholds can be used by
187 the CHARACTERIZE_SEX_CHROMS module to systematically compare mean depth
188 in pairs of chromosomes using three different approaches. The first is a bootstrap analysis
189 that provides 95% confidence intervals of mean window depth for each of the
190 chromosomes in a given pair to test for overlap. The second is a permutation analysis to
191 test for differences in depth between the two chromosomes. The third is a two-sample
192 Kolmogorov-Smirnov test (Massey Jr., 1951). Though all three tests are implemented in
193 XYalign, we only present results from the bootstrap analyses in this manuscript. Further,

194 while we present analyses pairing sex chromosomes with an autosome (here we use
195 chromosome 19), the chromosome pairs are arbitrary and can feature any scaffolds or
196 chromosomes in a reference genome, depending on a user's needs.

197 Finally, the REMAPPING module will infer the presence or absence of a Y
198 chromosome based on the results of CHARACTERIZE_SEX_CHROMS. If a Y
199 chromosome is not detected, the STRIP_READS module will iteratively remove reads
200 from the sex chromosomes by read group ID using SAMtools (Li *et al.*, 2009), writing
201 FASTQ files for each. XYalign will use repair.sh from BBTools to sort and re-pair
202 paired-end reads or shuffle.sh from BBTools to sort single-end reads for each read group.
203 The REMAPPING module then maps reads with BWA-MEM (Li, 2013) and sorts
204 alignments with SAMtools (Li *et al.*, 2009) by read group. If more than one read group is
205 present, the resulting BAM files are merged using SAMtools (Li *et al.*, 2009). Finally,
206 XYalign uses Sambamba (Tarasov *et al.*, 2015) to isolate all scaffolds not associated with
207 sex chromosomes from the original BAM file and then SAMtools (Li *et al.*, 2009) to
208 merge this file with the BAM file containing the new sex chromosome mappings.

209 When run as a full pipeline on a sample, XYalign will first call
210 PREPARE_REFERENCE to generate XX and XY reference genomes with appropriate
211 masks. Next, it will call ANALYZE_BAM and CHARACTERIZE_SEX_CHROMS to
212 preliminarily analyze the unprocessed input BAM file. Then, based on the results of
213 CHARACTERIZE_SEX_CHROMS, XYalign will call STRIP_READS to extract reads
214 from the sex chromosomes and REMAPPING to remap to the appropriate reference
215 genome output from PREPARE_REFERENCE. Finally, XYalign will re-run the
216 ANALYZE_BAM module to analyze the remapped BAM file and provide metrics to
217 allow a before-and-after comparison.

218 While we anticipate that this full pipeline will be useful in certain situations, it is
219 neither the only nor the best-suited option for most users. Rather, we expect that most
220 users will call modules individually. We give examples of other implementations below
221 and provide recommendations for incorporating XYalign into bioinformatic pipelines in
222 the discussion.

223

224 *Operation*

225 XYalign is available via PyPI (<https://pypi.python.org/pypi>), Bioconda (Grüning
226 *et al.*, 2017), and Github (<https://github.com/WilsonSayresLab/XYalign>), with
227 documentation hosted at Read the Docs (<http://xyalign.readthedocs.io/en/latest/>). A full
228 environment containing all dependencies can be most easily installed and managed using
229 Anaconda (<https://www.continuum.io/>) and Bioconda (Grüning *et al.*, 2017). It has been
230 tested on a variety of UNIX operating systems (including Linux and MacOS), but it is not
231 currently supported for the Windows operating system.

232 XYalign is typically invoked from the command line, but it can be imported into
233 Python scripts for more customized use cases. The next section lists a number of example
234 commands that illustrate how to call the full pipeline as well as individual modules.

235

236 *Use Cases*

237 To highlight some features of XYalign and its flexibility, we used two datasets
238 from publicly available sources (Supplemental Table S1): (1) exome, low-coverage
239 whole-genome, and high-coverage whole-genome sequencing data from one male

240 (HG00512) and one female (HG00513) from the 1000 Genomes Project (Dataset 1; (The
241 1000 Genomes Project Consortium, 2015); and (2) 24 high-coverage whole genomes
242 from the 1000 Genomes Project (Dataset 2; (Sudmant *et al.*, 2015). For Dataset 1, we
243 mapped reads to the hg19 version of the human reference genome (International Human
244 Genome Sequencing Consortium, 2001) using BWA MEM (Li, 2013), marked duplicates
245 with SAMBLASTER (Faust and Hall, 2014), and used SAMtools (Li *et al.*, 2009) to sort,
246 index, and merge BAM files. The publicly available BAM files for Dataset 2 were
247 previously mapped using a different version of hg19 (from the Broad Institute's GATK
248 Resource Bundle; <https://software.broadinstitute.org/gatk/download/bundle>), which we
249 used for analyses involving this dataset.

250 With these datasets, we examined three potential uses of XYalign. First, to
251 explore the effects of simple corrections for technical biases arising from sequence
252 homology on the sex chromosomes, we ran the full XYalign pipeline on all six BAM
253 files from Dataset 1. In all cases below, the exact commands are included in the
254 Supplementary Material and in a Snakemake (Köster and Rahmann, 2012) pipeline
255 available with the XYalign software distribution (Webster *et al.*, 2018), and templates are
256 shown here for convenience. Because we were using the same output directory for these
257 analyses, we avoided conflicts by initially preparing separate XX and XY references
258 using the following command:

259

```
260 xyalign --PREPARE_REFERENCE --ref <hg19 reference genome> --xx_ref_out  
261 hg19.XXonly.fasta --xy_ref_out hg19.XY.fasta --output_dir <output_directory> --  
262 x_chromosome chrX --y_chromosome chrY --bwa_index True
```

263

264 where *<hg19 reference genome>* was the path to the FASTA file containing the hg19
265 reference, *<input bam file>* was a sorted BAM file, and *<output directory>* was the
266 directory where XYalign wrote output. We then ran the full pipeline on all six files using
267 the following command template:

268

```
269 xyalign --ref <hg19 reference genome> --bam <input bam file> --output_dir <output  
270 directory> --sample_id <sample ID> --cpus 4 --reference_mask  
271 hg19_PAR_Ymask_startEnd.bed --window_size 5000 --chromosomes chr19 chrX chrY --  
272 x_chromosome chrX --y_chromosome chrY --xmx 4g --fastq_compression 4 --  
273 min_depth_filter 0.2 --max_depth_filter 2 --xx_ref_in hg19.XXonly.fasta --xy_ref_in  
274 ref_out hg19.XY.fasta,
```

275

276 where *<sample ID>* was the identification code for a given sample,
277 *hg19_PAR_Ymask_startEnd.bed* was a BED file containing the genomic coordinates of
278 the PARs in the hg19 assembly, and *hg19.XXonly.fasta* and *hg19.XY.fasta* were the two
279 FASTA formatted reference genomes prepared in the previous step.

280 Next, we examined how the metrics generated by XYalign can be used to identify
281 the sex chromosome complement of individuals with both datasets. Here, we used the
282 CHARACTERIZE_SEX_CHROMS module of XYalign. This was automatically done
283 for Dataset 1 when running the full pipeline (see above). For Dataset 2, we used the
284 following command template for BAM files:

285

```
286 xyalign --CHARACTERIZE_SEX_CHROMS --ref <1000 genomes reference genome> --  
287 bam <input bam file> --output_dir <output directory> --sample_id <sample ID> --cpus  
288 4 --window_size 5000 --chromosomes 19 X Y --x_chromosome X --y_chromosome Y  
289
```

290 Finally, we explored the utility of the CHROM_STATS module for identifying
291 sex chromosome complement and potentially sex-linked scaffolds with both datasets
292 using the following command template for BAM files:

```
293  
294 xyalign --CHROM_STATS --chromosomes chr1 chr8 chr19 chrX chrY chrM --bam  
295 <input_bam_1> <input_bam_2> <input_bam_3> --ref null --sample_id  
296 <name_of_analysis> --output_dir <output_dir>  
297
```

298 We additionally ran CHROM_STATS using the above command with the
299 addition of the “--use_counts” flag to calculate metrics using only the number of reads
300 mapping to each chromosome.

301 We visualized all CHROM_STATS results using the plot_count_stats utility, with the
302 command template:

```
303  
304  
305 plot_count_stats --input <chrom_stats output file> --output_prefix <output prefix>--  
306 meta <metadata text file> --exclude_suffix <suffix> --first_chr chrX --second_chr chrY -  
307 -const_chr chr19 --var1_marker color --var1_marker_vals darkslateblue thistle --  
308 var2_marker shape --var2_marker_vals o s v --marker_size 1700 --legend_marker_scale  
309 0.4  
310
```

311 where *<chrom_stats_output_file>* was either the count, mapping quality, or depth output
312 of CHROM_STATS, *<metadata text file>* was the appropriate metadata text file, and
313 *<suffix>* was the string to remove from filenames.

314
315 *Sex chromosome coordinates*

316 To better understand the genomic context of technical artifacts, we explored
317 variation in mapping quality and depth in association with genomic features on the X and
318 Y chromosomes. On the Y chromosome, we used coordinates from Poznik et al. (2013)
319 based on Skaletsky et al. (2003) (provided by D. Poznik, personal communication). On
320 the X chromosome, we obtained coordinates for ampliconic regions from Cotter et al.
321 (2016) and all other regions (PARs, telomeres, centromere, and XTR) from the UCSC
322 Table Browser (Karolchik et al., 2004). We define the X-transposed region (XTR) on the
323 X chromosome as beginning at the start of DXS1217 and ending at the end of DXS3
324 (Mumm et al., 1997).

325 To count variants falling in major genomic regions, we intersected a BED file
326 containing coordinates with VCF files using BEDTools (Quinlan and Hall, 2010). We
327 first filtered VCF files using BCFtools (Li et al., 2009) with the following command
328 template:

```
329  
330 bcftools filter --include 'INFO/MQ>=30 && %QUAL>=30' <input_vcf>  
331
```


332 We then identified variants unique to each file through iterations of the “subtract”
333 command in BEDtools (Quinlan and Hall, 2010):

334
335 *bedtools subtract -header -a <first_vcf> -b <second_vcf>*
336

337 Finally, in each region, we counted variants present in a given filtered VCF file using the
338 BEDtools (Quinlan and Hall, 2010) “intersect” command:

339
340 *bedtools intersect -c -a <BED file> -b <vcf_file>*
341

342 where *<BED_file>* is the BED file containing genomic coordinates (Supplemental Table
343 S2).

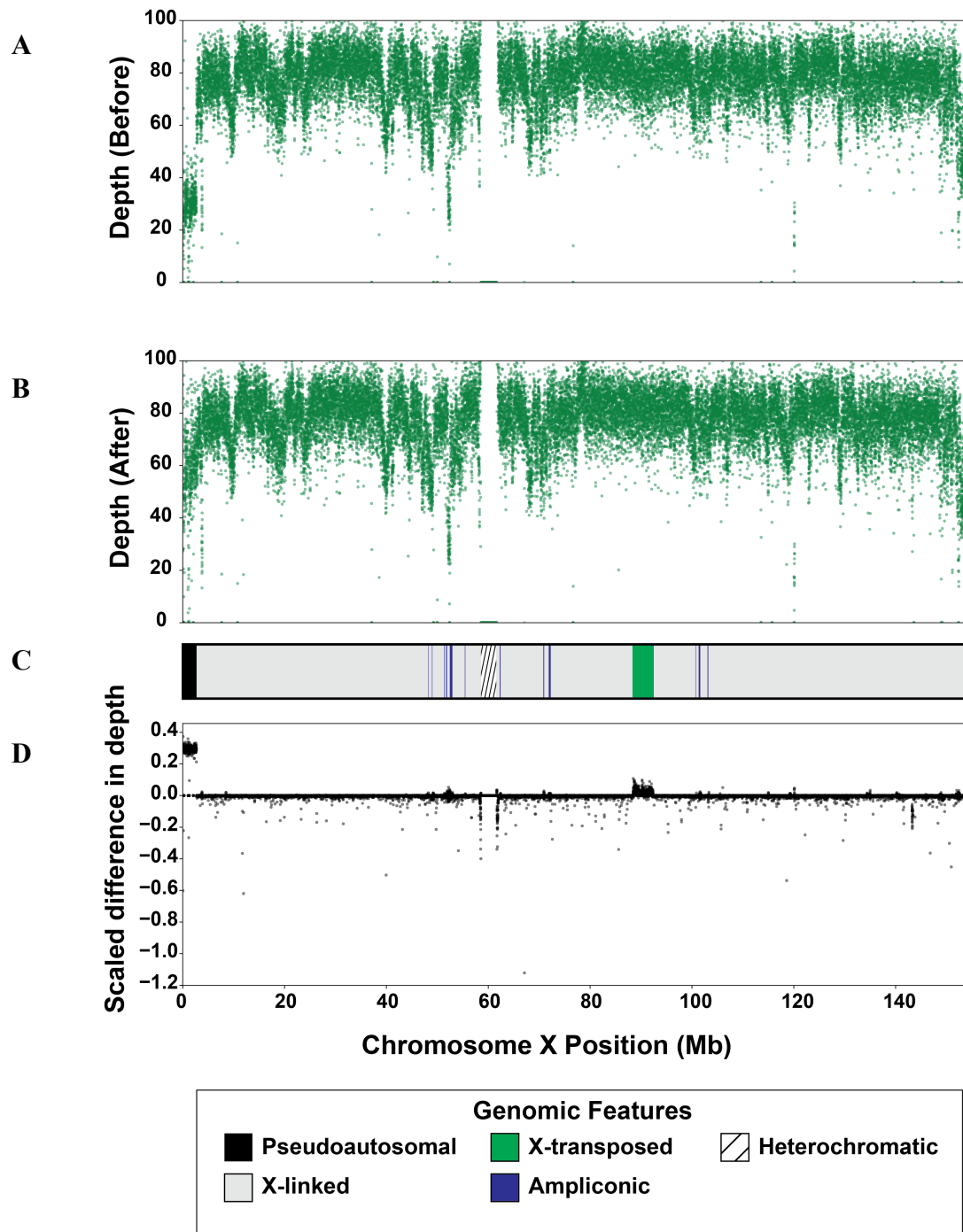
344
345 *Specific commands*

346 We provide Snakemake (Köster and Rahmann, 2012) workflows for all assembly
347 and analysis steps on Github (<https://github.com/WilsonSayresLab/XYalign>), Zenodo
348 (Webster *et al.*, 2018), and in the Supplemental Material.

349 350 **Results and Discussion**

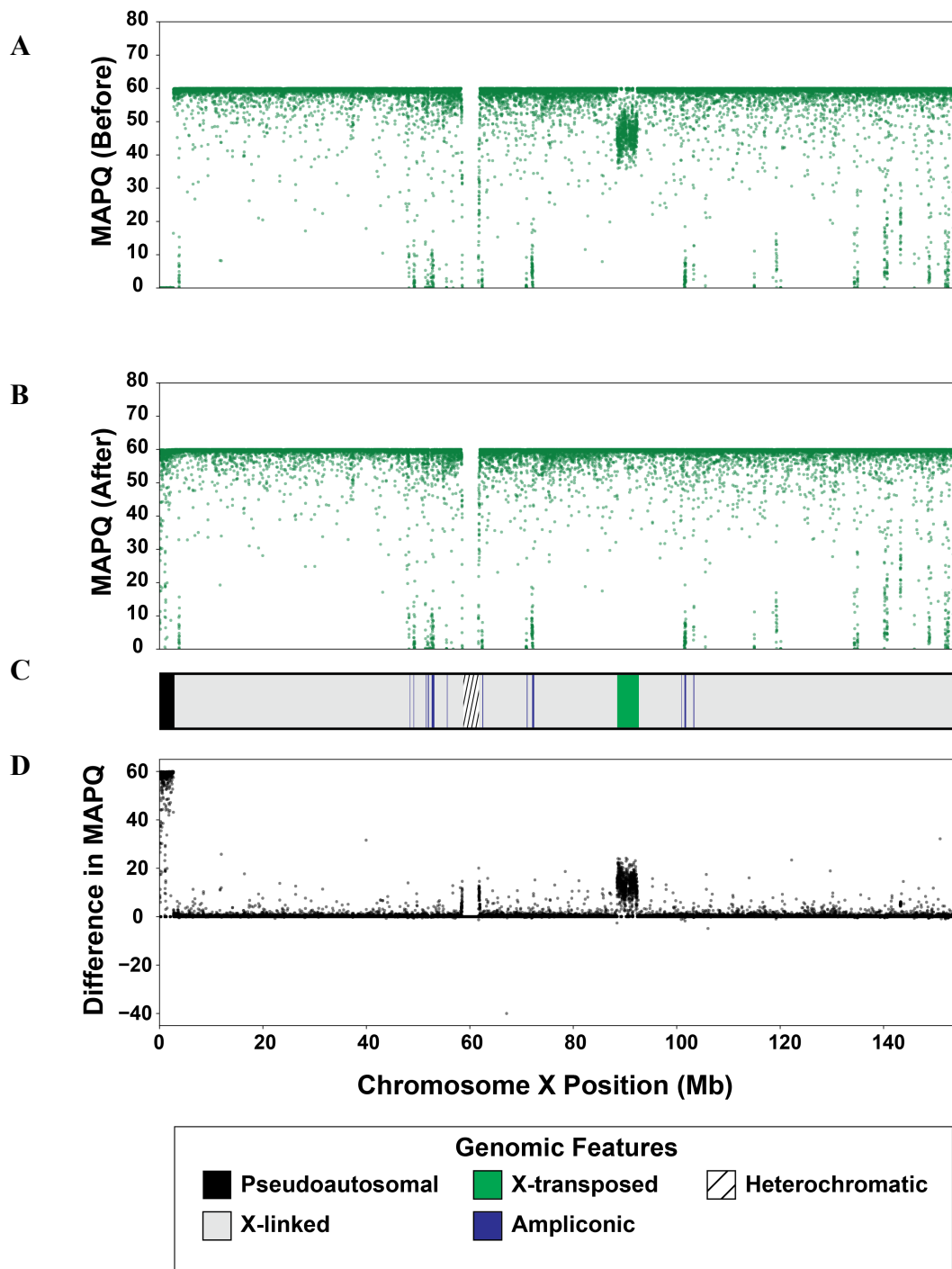
351 We observed a number of artifacts stemming from several methodological
352 challenges presented by the human sex chromosomes. First, PAR1 and PAR2 on both sex
353 chromosomes are clearly identifiable in genomic scatter plots of mapping quality and
354 depth in all datasets (Figures 1-3). While these results are not surprising given the
355 sequence homology in these regions (Ross *et al.*, 2005), they highlight the fact that these
356 measures can help identify other similarly problematic areas. For example, there is a
357 region of reduced mapping quality on the X chromosome beginning near 88.4 Mb and
358 ending near 92.3 Mb (Figure 2). This corresponds to the X-transposed region (XTR),
359 which arose by a duplication from the X to the Y chromosome in the human lineage since
360 its divergence with the chimpanzee-bonobo lineage (Page *et al.*, 1984; Ross *et al.*, 2005).
361 This region retains more than 98% sequence similarity between the X and Y chromosome
362 (Ross *et al.*, 2005), likely leading to the reduction in mapping quality. Interestingly, we
363 observe a similar decrease in mapping quality on the Y chromosome beginning near 2.9
364 Mb and ending near 6.6 Mb, corresponding with known coordinates of the XTR on the Y
365 chromosome (Figure 3). In fact, integrating mapping quality and depth recapitulates
366 established genomic features of both sex chromosomes (e.g., ampliconic regions, PARs,
367 and XTRs) described in previous studies (Figures 1-3; Poznik *et al.*, 2013; Mueller *et al.*,
368 2013). This suggests that, in at least some cases, the output of XYalign can be used to
369 quickly explore broad patterns of genomic architecture and mask regions likely to
370 introduce technical difficulties in genomic analyses.

371
372



373
374 **Figure 1. Sequencing depth on chromosome X before and after XYalign.** Mean
375 sequencing depth in 5 kb windows across the X chromosome before (A) and after (B)
376 XYalign processing. Changes in depth (D) are presented as the sign of the difference
377 times the absolute value of the \log_{10} difference, where the difference is depth after
378 XYalign minus depth before XYalign. The chromosome map (C) presents the location of
379 X chromosome genomic features depicted in the legend. X chromosome coordinates are
380 identical in all plots.

381

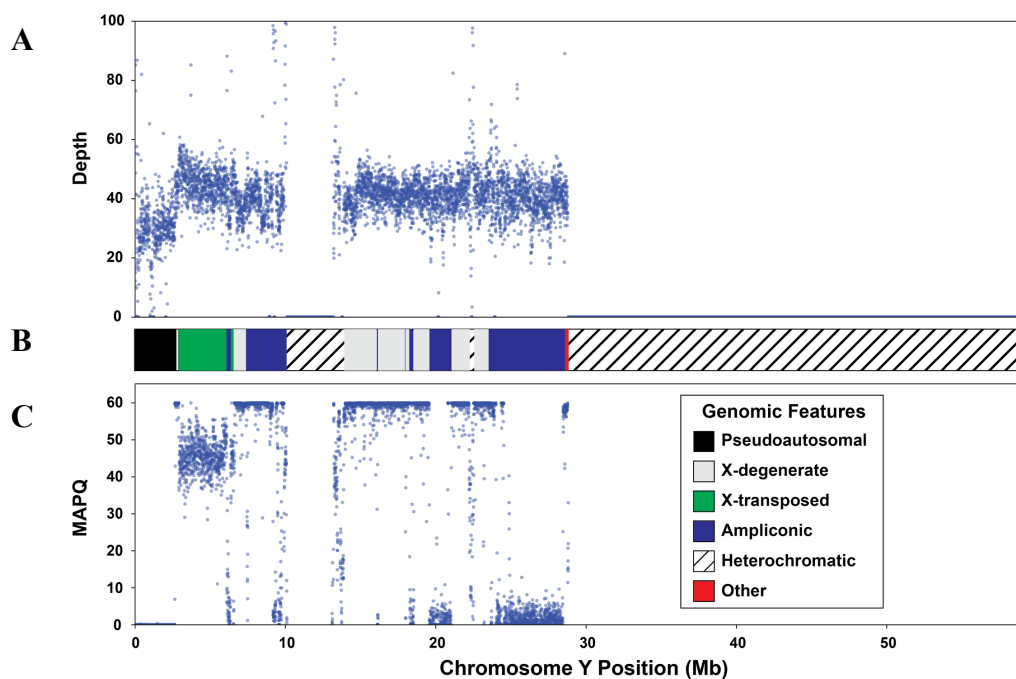


382

383

384 **Figure 2. Mapping quality on chromosome X before and after XYalign.** Mean
385 mapping quality (MAPQ) in 5 kb windows across the X chromosome before (A) and
386 after (B) XYalign processing. Changes in MAPQ (D) are presented as the difference is
387 MAPQ after XYalign minus MAPQ before XYalign. The chromosome map (C) presents
388 the location of X chromosome genomic features depicted in the legend. X chromosome
coordinates are identical in all plots.

389



390
391
392
393
394
395

Figure 3. Y chromosome sequencing depth and quality. Mean sequencing depth (A) and mapping quality (MAPQ; C) in 5 kb windows across the Y chromosome. The chromosome map (B) presents the location of Y chromosome genomic features depicted in the legend. Y chromosome coordinates are identical in all plots.

396 By hard-masking the Y chromosome in the XX reference genome, and the
397 pseudoautosomal regions (PAR1 and PAR2) in the reference genome for the XY
398 reference genome, we observed clear improvements in read mapping (Figures 1-2). On
399 the X chromosome, all metrics exhibited striking improvements in PAR1, PAR2, and
400 XTR (Figures 1 and 2). Furthermore, the Y chromosome of the XX individual no longer
401 exhibited any variant calls or mapped reads, though many passed filters before processing
402 (variants before: 4266; variants after: 0; mapped reads before: 5,729,007; reads mapped
403 after: 0). While this is expected given the hard masking of the Y chromosome, it is worth
404 emphasizing that this is consistent with the biological state of the individual.

405 We found that these improvements in mapping on the X chromosome after
406 masking the Y chromosome substantially impacted downstream variant calling (Table 1).
407 Unsurprisingly, the effect was most pronounced in the PARs, in which thousands of
408 variants were callable after masking the identical sequences present on the Y
409 chromosome in the reference assembly. The XTR also had a large increase in the number
410 of variants detected after Y masking—an average of 85.4 variants per megabase of
411 sequence (Table 1). However, effects were not limited to these regions of well-
412 documented homology: both the X-added region (XAR) and X-conserved region (XCR)
413 contained hundreds of affected variants, suggesting effects of more extensive homology
414 across the sex chromosomes.

415
416 **Table 1. The effect of sex chromosome homology on variant calling on the X**
417 **chromosome.**

Region^a	Length^b	False positives (per Mb)^c	False negatives (per Mb)^d
PAR1	2,589,520	0 (0)	7563 (2920.6)
PAR2	329,516	0 (0)	633 (1921)
XTR	4,287,237	40 (9.3)	366 (85.4)
XAR	55,982,492	299 (5.3)	400 (7.2)
XCR	89,011,795	610 (6.9)	523 (5.9)
<i>Total</i>	<i>152,250,560</i>	<i>949 (6.2)</i>	<i>9485 (62.3)</i>

418

419 ^aPAR1: pseudoautosomal region 1; PAR2: pseudoautosomal region 2; XTR: X-
420 transposed region; XAR: X-added region; XCR: X-conserved region.

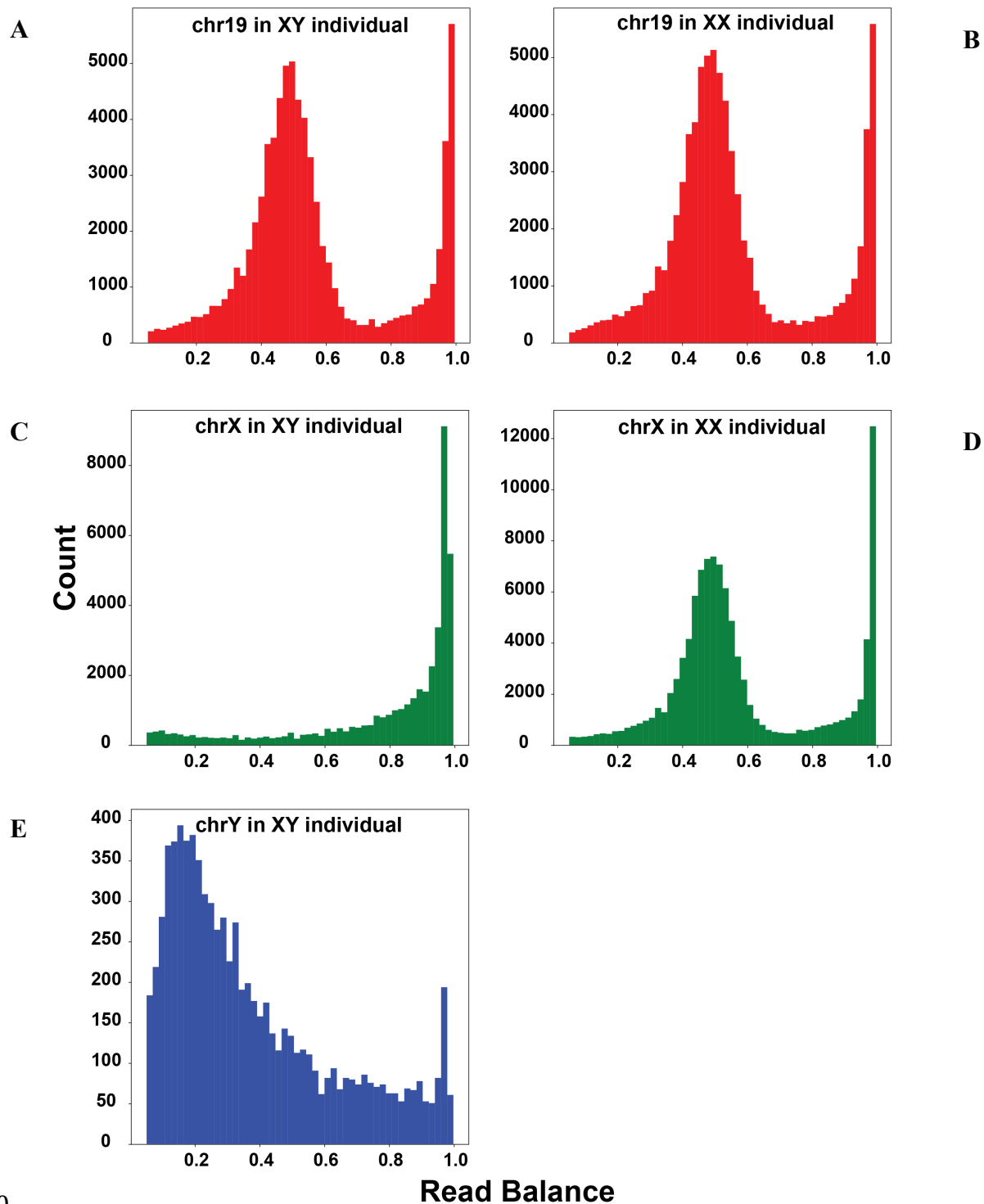
421 ^bTotal sequence length of region in base pairs.

422 ^cTotal number of false positive variants after filtering, defined as being present before but
423 not after Y chromosome masking. Variants per Mb of sequence are presented in
424 parentheses.

425 ^dTotal number of false negative variants after filtering, defined as being present after but
426 not before Y chromosome masking. Variants per Mb of sequence are presented in
427 parentheses.

428

429



430
431
432
433
434

Figure 4. Read balance in XY and XX samples. Histograms of read balance for an XY sample (Left Column; A, C, and E) and XX sample (Right Column; B and D) across chromosome 19 (Top; A and B), chromosome X (Middle; C and D), and chromosome Y (Bottom; E). Read balance at a given site is defined as the number of reads containing a

435 non-reference allele divided by the total number of reads mapped to a site. Read balances
436 between 0.05 and 1.0 are presented.

437

438 *Inferring Genetic Sex*

439 In our analyses, the most striking measure for assessing an individual's sex
440 chromosome complement was the distribution of read balances across a chromosome
441 (Figure 4). Specifically, when we plotted the distribution of the fraction of reads
442 containing a nonreference allele at a given variant site, we observed that diploid
443 chromosomes (e.g., autosomes, and chromosome X in XX individuals) exhibited peaks
444 both around 0.5 and 1.0, consistent with the presence of heterozygous sites and sites
445 homozygous for a nonreference allele, respectively (Figure 4). In the case of the X
446 chromosome in XY individuals, we observed a single peak near 1.0, consistent with an
447 expected haploid state (i.e., no heterozygous sites; Figure 4). We observed one exception
448 to this pattern: the Y chromosome exhibited a peak around 0.2 in addition to the one near
449 1.0 (Figure 4). All variants included in analyses met thresholds for depth, site quality, and
450 genotype quality, so quality does not appear to be a driving factor of this pattern. This
451 pattern also remained after genomic windows of low mapping quality and irregular depth
452 were removed. We are currently unable to explain these results and more work is thus
453 required to understand the factors responsible for this pattern and whether similar results
454 are obtained on the W chromosome in ZW systems.

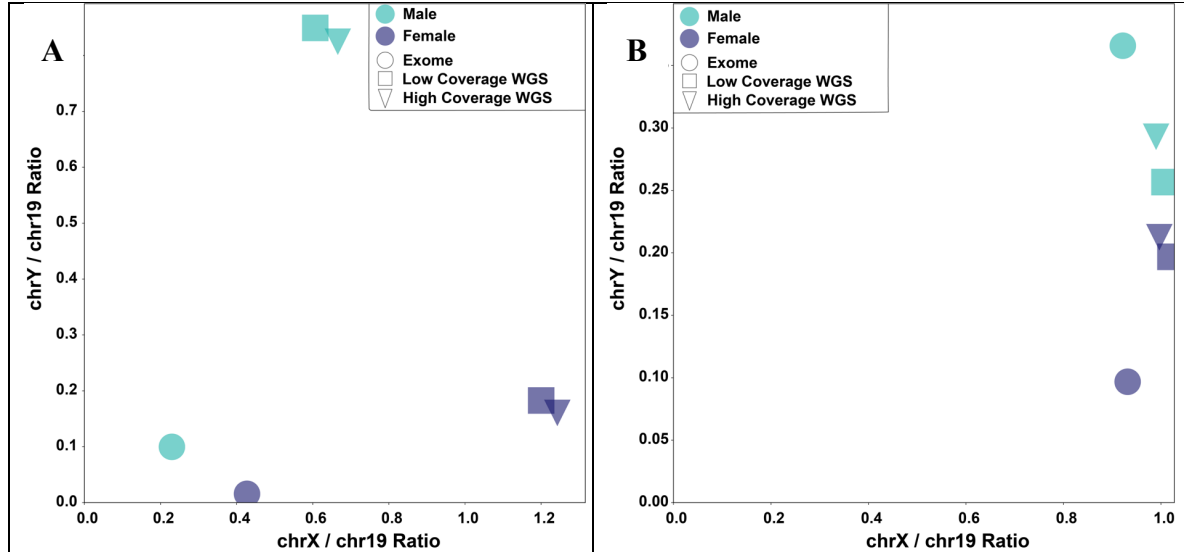
455 Across datasets, we observed variation in relative depth of the X and Y
456 chromosomes in XX and XY individuals, particularly among different sequencing
457 strategies: exome, low-coverage whole-genome, and high-coverage whole-genome
458 sequencing (Figure 5A). However, within datasets, XX and XY individuals were clearly
459 differentiated (Figure 5; Supplemental Figure S1). This pattern suggests that a general
460 threshold for assigning different genetic sexes across a range of organisms and
461 sequencing experiments might be difficult to implement. That being said, within species,
462 some combination of depth, mapping quality, and read balance is likely to be informative.
463 For example, in humans, relative mapping quality appears to be informative in some
464 sequencing strategies, particularly exome sequencing (Figure 5B). This should be
465 explored in each experiment, however, as we did not observe this differentiation in the
466 uncorrected 1000 Genomes high-coverage samples (Supplemental Figure S2).

467 Generating these results for all individuals in a study is easy to do with XYalign:
468 one can iteratively run the CHARACTERIZE_SEX_CHROMS module on preliminarily
469 mapped BAM files. Then, the results from all individuals can be analyzed together (see
470 the Supplementary Material for an example of such analysis). At least with human
471 samples, for which X and Y chromosomes are very differentiated, this process can be
472 sped up significantly with the CHROM_STATS module. In our data, read counts on the
473 X and Y chromosomes quickly and clearly clustered male and female samples within
474 sequencing strategies (i.e., exome, low-coverage whole-genome, and high-coverage
475 whole-genome; Supplemental Figures S3-S4). However, the success of this procedure
476 likely depends on the degree of differentiation between sex chromosomes; other
477 organisms might require the statistics output as part of the
478 CHARACTERIZE_SEX_CHROMS module.

479

480

481
482



483

484

485 **Figure 5. Relative sequencing depth and mapping quality on the X and Y**
486 **chromosomes across different sequencing strategies.** Values of relative (A) sequencing
487 depth and (B) mapping quality come from exome (circles), low-coverage whole-genome
488 sequencing (squares), and high-coverage whole-genome sequencing (triangles) for a
489 single male (green) and female (blue) individual. Mean depth and MAPQ on
489 chromosome 19 was used to normalize the sex chromosomes.

490

491

Recommendations for researchers

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

Based on these results, we can make the following recommendations for researchers. For organisms with multiple sex chromosomes assembled (e.g., both X and Y or both Z and W) and included in reference assemblies (e.g., human, chimpanzee, rhesus macaque, gorilla, mouse, rat, chicken, *Drosophila*), *if the genetic sex of every individual is known*, the user may: (1) prepare separate assemblies for the different sexes using the PREPARE_REFERENCE module; (2) map and process reads according to user's typical pipeline (mapping individuals by sex to their corresponding reference); (3) confirm genetic sex using the CHROM_STATS module; (4) remap any incorrectly assigned individuals; and (5) proceed with downstream analyses. *If genetic sexes of individuals are unknown*, the user should then: (1) prepare separate assemblies for the different sexes using the PREPARE_REFERENCE module; (2) map and process a suitable number of reads (e.g., whole dataset for exome or a single lane of WGS) according to user's typical pipeline using the reference genome of the heterogametic sex (i.e., XY or ZW); (3) infer the sex chromosome complement using either CHROM_STATS (for well-characterized and highly divergent sex chromosomes), CHARACTERIZE_SEX_CHROMS, or both; (4) map and process all reads using the prepared reference genome corresponding to the inferred sex of each individual; and (5) run downstream analyses.

For individuals of the homogametic sex (i.e., XX or ZZ), the above recommendations will likely completely remove artifacts stemming from sex

512 chromosome homology, assuming only a single unmasked sex chromosome is left after
513 XYalign processing. However, homology is unavoidable for individuals of the
514 heterogametic sex (i.e., XY or ZW) because both sex chromosomes are required in the
515 reference assembly for mapping. In this case, a more local masking or filtering approach
516 is likely the most promising option. For studies investigating specific variants, for which
517 false negatives are preferable to false positives, we suggest strict variant filtering that
518 includes high thresholds for mapping quality (e.g., thresholds of 55 or higher are required
519 to eliminate the effects of homology in the X-transposed region). However, for studies
520 investigating invariant sites as well (e.g., measures of genetic diversity require
521 information from all monomorphic and polymorphic sites), we recommend filtering
522 entire regions based on, at the very least, mapping and depth metrics. These masks are
523 output by the BAM_ANALYSIS module in XYalign, and for this use, we recommend
524 using small windows (e.g, 1 kb to 5 kb) and exploring a variety of depths. Finally, in all
525 cases, if pseudoautosomal regions are present in the reference genome, they should be
526 masked in the heterogametic sex's assembly output by the PREPARE_REFERENCE
527 module.

528

529 *Additional uses for XYalign*

530 While the development of XYalign was motivated by challenges surrounding
531 erroneous read mapping and variant calling due to sex chromosome homology in human
532 sequencing experiments, the software can be utilized in a number of additional scenarios.
533 First, it can be applied to any species with heteromorphic sex chromosomes to identify
534 relative quality and depth. The results output by CHROM_STATS, ANALYZE_BAM,
535 and CHARACTERIZE_SEX_CHROMS can be used to identify sex-linked scaffolds,
536 characterize sex chromosome complements, and determine the most appropriate
537 remapping strategy. Second, XYalign can be used to detect relative sequencing depth,
538 mapping quality, and read balance on any chromosome, not just the sex chromosomes. In
539 addition to exploring mapping artifacts, we anticipate that this will aid in detection of
540 aneuploidy in the autosomes. However, we note that many programs exist to calculate
541 depth of coverage (e.g., Quinlan and Hall, 2010; Pedersen and Quinlan, 2018; McKenna
542 *et al.*, 2010) and identify structural variants within statistical frameworks (e.g., Chen *et*
543 *al.*, 2016; Layer *et al.*, 2014; Abyzov *et al.*, 2011; Roller *et al.*, 2016). Accordingly,
544 XYalign might not be the most appropriate option for detecting local phenomena such as
545 copy number variants. Finally, XYalign may also be extended to other types of data,
546 including RNA sequencing data, where the same fundamental challenge (gametologous
547 sequence between the X and Y) can affect mapping and variant calling. In particular, we
548 expect biases to manifest in differential expression and biased-allelic expression, and
549 suggest that the PREPARE_REFERENCE module be considered for all RNA sequencing
550 experiments in systems with sex chromosomes.

551

552 **Conclusion**

553 We showed that the complex evolutionary history of the sex chromosomes creates
554 mapping biases in next-generation sequencing data that have downstream effects on
555 variant calling and other analyses. These technical artifacts are likely present in most
556 genomic datasets of species with chromosomal sex determination. However, many of
557 these biases can be corrected through the strategic use of masks during read mapping and

558 the filtering of variants. We developed XYalign, a tool that facilitates the characterization
559 of an individual's sex chromosome complement and implements this masking strategy to
560 correct these technical biases. We illustrated how XYalign can be used to identify the
561 presence or absence of a Y chromosome, characterize mapping biases across the genome,
562 and correct for these mapping biases. XYalign provides a framework to generate more
563 robust short read mapping and improve variant calling on the sex chromosomes.

564

565 **Software Availability**

566 XYalign is available on Github (<https://github.com/WilsonSayresLab/XYalign>) under a
567 GNU General Public License (version 3). We have also deposited a static version of the
568 source code used for analyses in this paper at Zenodo (Webster *et al.*, 2018).

569

570 **Author Contributions**

571 MAWS and THW conceived the research. All authors participated in the initial design of
572 the software. THW was responsible for subsequent design, development, and
573 implementation of the software. BG, EK, TNP, WW, and THW tested the software.
574 THW analyzed the data. THW and MAWS wrote the manuscript. All authors were
575 involved in the revision of the manuscript and have agreed to the final content.

576

577 **Competing Interests**

578

579 No competing interests were disclosed.

580

581 **Grant Information**

582 This study was supported by startup funds from the School of Life Sciences and the
583 Biodesign Institute at Arizona State University to MAWS. Furthermore, this study was
584 supported by the National Institute of General Medical Sciences of the National Institutes
585 of Health under Award Number R35GM124827 to MAWS. The content is solely the
586 responsibility of the authors and does not necessarily represent the official views of the
587 National Institutes of Health.

588 **Acknowledgements**

589

590 We thank the organizers of Hackseq 2016 for facilitating this project and supporting this
591 collaboration; members of the Wilson Sayres lab for helpful comments; and ASU
592 Research Computing for computational resources.

593

594 **References**

- 595 Abyzov, A. *et al.* (2011) CNVnator: an approach to discover, genotype, and characterize
596 typical and atypical CNVs from family and population genome sequencing.
597 *Genome Res.*, **21**, 974–984.
- 598 Ashley, E. A. (2016) Towards precision medicine. *Nat. Rev. Genet.*, **17**, 507–522.
- 599 Bergero, R. and Charlesworth, D. (2009) The evolution of restricted recombination in sex
600 chromosomes. *Trends Ecol. Evol.*, **24**, 94–102.
- 601 Chang, D. *et al.* (2014) Accounting for eXentricities: analysis of the X chromosome in

- 602 GWAS reveals X-linked genes implicated in autoimmune diseases. *PloS One*, **9**,
603 e113684.
- 604 Chen, X. *et al.* (2016) Manta: rapid detection of structural variants and indels for
605 germline and cancer sequencing applications. *Bioinformatics*, **32**, 1220–1222.
- 606 Cotter, D. J. *et al.* (2016) Genetic Diversity on the human X chromosome does not
607 support a strict pseudoautosomal boundary. *Genetics*, **203**, 485–492.
- 608 Dale, R. K. *et al.* (2011) Pybedtools: a flexible Python library for manipulating genomic
609 datasets and annotations. *Bioinformatics*, **27**, 3423–3424.
- 610 Ellegren, H. (2009) The different levels of genetic diversity in sex chromosomes and
611 autosomes. *Trends Genet.*, **25**, 278–284.
- 612 Faust, G. G. and Hall, I. M. (2014) SAMBLASTER: fast duplicate marking and structural
613 variant read extraction. *Bioinformatics*, **30**, 2503–2505.
- 614 Gao, F. *et al.* (2015) XWAS: a software toolset for genetic data analysis and association
615 studies of the X chromosome. *J. Hered.*, **106**, 666–671.
- 616 Glas, R. *et al.* (1999) Cross-species chromosome painting between human and marsupial
617 directly demonstrates the ancient region of the mammalian X. *Mamm. Genome*
618 *Off. J. Int. Mamm. Genome Soc.*, **10**, 1115–1116.
- 619 Grüning, B. *et al.* (2017) Bioconda: a sustainable and comprehensive software
620 distribution for the life sciences. *bioRxiv*.
- 621 Hunter, J. D. (2007) Matplotlib: a 2D graphics environment. *Comput. Sci. Eng.*, **9**, 90–95.
- 622 International Human Genome Sequencing Consortium (2001) Initial sequencing and
623 analysis of the human genome. *Nature*, **409**, 860–921.
- 624 Jones, E. *et al.* (2001) SciPy: open source scientific tools for Python.
- 625 Karolchik, D. *et al.* (2004) The UCSC Table Browser data retrieval tool. *Nucleic Acids*
626 *Res.*, **32**, D493–D496.
- 627 Köster, J. and Rahmann, S. (2012) Snakemake--a scalable bioinformatics workflow
628 engine. *Bioinformatics*, **28**, 2520–2522.
- 629 Lahn, B. T. and Page, D. C. (1999) Four evolutionary strata on the human X
630 chromosome. *Science*, **286**, 964–967.
- 631 Layer, R. M. *et al.* (2014) LUMPY: a probabilistic framework for structural variant
632 discovery. *Genome Biol.*, **15**, R84.
- 633 Li, H. (2013) Aligning sequence reads, clone sequences and assembly contigs with
634 BWA-MEM. *arXiv*, **1303.3997**.
- 635 Li, H. *et al.* (2009) The Sequence Alignment/Map format and SAMtools. *Bioinformatics*,
636 **25**, 2078–2079.
- 637 Livernois, A. M. *et al.* (2012) The origin and evolution of vertebrate sex chromosomes
638 and dosage compensation. *Heredity*, **108**, 50–58.
- 639 Madel, M.-B. *et al.* (2016) TriXY-Homogeneous genetic sexing of highly degraded
640 forensic samples including hair shafts. *Forensic Sci. Int. Genet.*, **25**, 166–174.
- 641 Massey Jr., F. J. (1951) The Kolmogorov-Smirnov test for goodness of fit. *J. Am. Stat.*
642 *Assoc.*, **46**, 68–78.
- 643 McKenna, A. *et al.* (2010) The Genome Analysis Toolkit: a MapReduce framework for
644 analyzing next-generation DNA sequencing data. *Genome Res.*, **20**, 1297–1303.
- 645 McKinney, W. (2010) Data structures for statistical computing in Python., *Proceedings*
646 *of the 9th Python in Science Conference*, 51–56.
- 647 Meisel, R. P. and Connallon, T. (2013) The faster-X effect: integrating theory and data.

- 648 *Trends Genet.*, **29**, 537–544.
- 649 Mueller, J. L. *et al.* (2013) Independent specialization of the human and mouse X
650 chromosomes for the male germ line. *Nat. Genet.*, **45**, 1083.
- 651 Mumm, S. *et al.* (1997) Evolutionary Features of the 4-Mb Xq21.3 XY Homology
652 Region Revealed by a Map at 60-kb Resolution. *Genome Res.*, **7**, 307–314.
- 653 Muyle, A. *et al.* (2016) SEX-DETECTOR: a probabilistic approach to study sex
654 chromosomes in non-model organisms. *Genome Biol. Evol.*, **8**, 2530–2543.
- 655 Oliphant, T. E. (2006) A Guide to NumPy. Trelgol Publishing, USA.
- 656 Page, D. C. *et al.* (1984) Occurrence of a transposition from the X-chromosome long arm
657 to the Y-chromosome short arm during human evolution. *Nature*, **311**, 119–123.
- 658 Pedersen, B. S. and Quinlan, A. R. (2018) Mosdepth: quick coverage calculation for
659 genomes and exomes. *Bioinformatics*, **34**, 867–868.
- 660 Poznik, G. D. *et al.* (2013) Sequencing Y chromosomes resolves discrepancy in time to
661 common ancestor of males versus females. *Science*, **341**, 562–565.
- 662 Quinlan, A. R. and Hall, I. M. (2010) BEDTools: a flexible suite of utilities for
663 comparing genomic features. *Bioinformatics*, **26**, 841–842.
- 664 Rens, W. *et al.* (2007) The multiple sex chromosomes of platypus and echidna are not
665 completely identical and several share homology with the avian Z. *Genome Biol.*,
666 **8**, R243.
- 667 Rimmer, A. *et al.* (2014) Integrating mapping-, assembly- and haplotype-based
668 approaches for calling variants in clinical sequencing applications. *Nat. Genet.*,
669 **46**, 912.
- 670 Roller, E. *et al.* (2016) Canvas: versatile and scalable detection of copy number variants.
671 *Bioinformatics*, **32**, 2375–2377.
- 672 Ross, M. T. *et al.* (2005) The DNA sequence of the human X chromosome. *Nature*, **434**,
673 325–337.
- 674 Simmler, M. C. *et al.* (1985) Pseudoautosomal DNA sequences in the pairing region of
675 the human sex chromosomes. *Nature*, **317**, 692–697.
- 676 Skaletsky, H. *et al.* (2003) The male-specific region of the human Y chromosome is a
677 mosaic of discrete sequence classes. *Nature*, **423**, 825–837.
- 678 Sudmant, P. H. *et al.* (2015) An integrated map of structural variation in 2,504 human
679 genomes. *Nature*, **526**, 75.
- 680 Tarasov, A. *et al.* (2015) Sambamba: fast processing of NGS alignment formats.
681 *Bioinformatics*, **31**, 2032–2034.
- 682 Taylor, J. C. *et al.* (2015) Factors influencing success of clinical genome sequencing
683 across a broad spectrum of disorders. *Nat. Genet.*, **47**, 717–726.
- 684 The 1000 Genomes Project Consortium (2015) A global reference for human genetic
685 variation. *Nature*, **526**, 68–74.
- 686 Vicoso, B. and Charlesworth, B. (2006) Evolution on the X chromosome: unusual
687 patterns and processes. *Nat. Rev. Genet.*, **7**, 645–653.
- 688 Webster, T. H. *et al.* (2018) XYalign: Version 1.1.4. Zenodo.
689 <http://doi.org/10.5281/zenodo.1313870>
- 690 Webster, T. H. and Wilson Sayres, M. A. (2016) Genomic signatures of sex-biased
691 demography: progress and prospects. *Curr. Opin. Genet. Dev.*, **41**, 62–71.
- 692 Wilson, M. A. and Makova, K. D. (2009) Genomic analyses of sex chromosome
693 evolution. *Annu. Rev. Genomics Hum. Genet.*, **10**, 333–354.

- 694 Wilson Sayres, M. A. (2018) Genetic Diversity on the Sex Chromosomes. *Genome Biol.*
695 *Evol.*, **10**, 1064–1078.
- 696 Wilson Sayres, M. A. and Makova, K. D. (2013) Gene Survival and Death on the Human
697 Y Chromosome. *Mol. Biol. Evol.*, **30**, 781–787.
- 698 Wise, A. L. *et al.* (2013) eXclusion: toward integrating the X chromosome in genome-
699 wide association analyses. *Am. J. Hum. Genet.*, **92**, 643–647.
- 700