# seGMM: a new tool to infer sex from massively parallel sequencing data

1  **Sihan Liu[1], Yuanyuan Zeng[2], Meilin Chen[1], Qian Zhang[1], Lanchen Wang[1], Chao Wang[1], Yu**
2  **Lu[1,*], Hui Guo[3,*], Fengxiao Bu[1,*]**

3  [1] Institute of Rare Diseases, West China Hospital of Sichuan University, Chengdu 610041, China

4  [2] School of Medicine, National Institute for Data Science in Health and Medicine, Xiamen
5  University, Xiamen, China

6  [3] Center for Medical Genetics & Hunan Provincial Key Laboratory of Medical Genetics, School of
7  Life Sciences, Central South University, Changsha, Hunan 410078, China

8  **\* Correspondence:**
9  Fengxiao Bu: bufengxiao@wchscu.cn;

10  Hui Guo: guohui@sklmg.edu.cn;

11  Yu Lu: samuelluyu@163.com

14  **Abstract**

15  Inspecting concordance between self-reported sex and genotype-inferred sex from genomic data is a
16  significant quality control measure in clinical genetic testing. Numerous tools have been developed to
17  infer sex for genotyping array, whole-exome sequencing, and whole-genome sequencing data.
18  However, improvements in sex inference from targeted gene sequencing panels are warranted. Here,
19  we propose a new tool, seGMM, which applies unsupervised clustering (Gaussian Mixture Model) to
20  determine the gender of a sample from the called genotype data integrated aligned reads. seGMM
21  consistently demonstrated > 99% sex inference accuracy in publicly available (1000 Genomes) and
22  our in-house panel dataset, which achieved obviously better sex classification than existing popular
23  tools. Compared to including features only in the X chromosome, our results show that adding
24  additional features  from Y chromosomes (e.g. reads mapped to the Y chromosome) can increase sex
25  classification accuracy. Notably, for WES and WGS data, seGMM also has an extremely high degree
26  of accuracy. Finally, we proved the ability of seGMM to infer sex in single patient or trio samples by
27  combining with reference data and pinpointing potential sex chromosome abnormality samples. In
28  general, seGMM provides a reproducible framework to infer sex from massively parallel sequencing
29  data and has great promise in clinical genetics.

30

31

32

## 1    Introduction

Next-generation sequencing (NGS) has revolutionized the clinical field by transforming the landscape of clinical genetic testing and has been adopted as a standard for diagnosing hereditary disorders, expanding our understanding of clinical genetics, and offering new opportunities for personalized precision medicine over the last decade (Phillips and Douglas, 2018; Phillips et al., 2020). Clinical genetic testing usually refers to the analysis of DNA to identify pathogenic variants to aid in the diagnosis of disease (McPherson, 2006). It may focus on a single gene, multi-gene panels [targeted gene sequencing (TGS)], whole exome [whole exome sequencing (WES)], or whole genome [whole genome sequencing (WGS)](Di Resta et al., 2018). TGS is highly recommended in genetic testing because of its validity, utility, and cost-effectiveness, especially in hearing loss, cardiovascular disorders, and renal disorders (Lin et al., 2012; Saudi Mendeliome, 2015).

Parallelized TGS analysis of patients from the large cohort is commonplace in clinical genetic testing. Considerable efforts are required for quality control (QC) and preprocessing of these data before detecting pathogenic variants (Lee et al., 2017). Mismatched genders indicate potential sample swap, pollution, sex chromosome abnormalities, or sequencing error, which will substantially lead to erroneous conclusions and affect treatment decisions (Taylor et al., 2015; Webster et al., 2019). Thus, one essential QC step is verifying concordance between self-reported sex and genotype-inferred sex. Cytogenetic analyses, such as karyotyping, are gold standard methods of inferring sex but are time and effort consuming. Leveraging computational tools to infer genotypic sex from sequencing data of X and Y chromosomes is a convenient and powerful alternative strategy to verify sex concordance.

Several tools, such as PLINK, seXY, and XYalign, have been developed to infer sex using data from genotyping array, WES, or WGS. PLINK calculated the F coefficient with X chromosome heterozygosity to infer sex for genotype array data (Purcell et al., 2007). In contrast, seXY considered both X chromosome heterozygosity and Y chromosome missingness to infer sex in genotype array data by logistic regression (Qian et al., 2017). In particular, XYalign extract read count mapped to sex chromosomes and calculated the ratio of X and Y counts (Webster et al., 2019). Together with self-reported sex and calculated ratios, by plotting a scatter plot, users can infer sex with eyeballs for WES and WGS data. However, the accuracy of these methods in TGS panel data has not yet been fully evaluated, and improvements in sex inference from gene panel data are warranted.

In this study, we propose a new sex inference tool, seGMM, that determines the gender of a sample from called genotype data integrated aligned reads and jointly considers information on the X and Y chromosomes in diverse genomic data, including TGS panel data. seGMM applies Gaussian Mixture Model (GMM) clustering to classify the samples into different clusters. Compared to previous methods that use logistic regression and training data to infer sex, seGMM is more powerful for modeling data with different covariance structures and different numbers of mixture components for various genomic data.

## 2    Materials and Methods

### 2.1    Data

To evaluate the accuracy of existing methods and seGMM in inferring sex for TGS panel data, we used 2 datasets from publicly available sources (Supplementary Table S1) and 1 dataset from our in-house resource: (1) exon-targeted sequencing data for 1000 genes from 110 males and 98 females from the 1000 Genomes Project (Dataset 1, (Genomes Project et al., 2010)); (2) massively parallel sequencing of 785 deafness-related genes (Supplementary Table S2) from 8,950 males and 7,737

76  females (Dataset 2); and (3) targeted sequencing data for 189 autism risk genes across a cohort from
77  the Autism Clinical and Genetic Resources in China (ACGC), including 42 females and 205 males
78  (Dataset 3, (Guo et al., 2018)).

79  In addition, to assess the application of seGMM in inferring sex for WES and WGS data, we used 2
80  datasets from publicly available sources (Supplementary Table S3 and Supplementary Table S4) and
81  1 dataset from our in-house source: (1) exome sequencing data from 164 males and 118 females from
82  the 1000 Genomes Project (Dataset 4, (Genomes Project et al., 2015)); (2) 27 high-coverage whole
83  genomes from the 1000 Genomes Project including 11 males and 16 females (Dataset 5; (Genomes
84  Project et al., 2015)) and (5) exome sequencing data from 1,255 males and 1,138 females from our
85  in-house resource (Dataset 6).

86  The publicly available BAM files were previously mapped to the reference genome (GRCh37),
87  which we can directly use for downstream analyses. For our in-house datasets, Fastp was used to
88  remove adapters and low-quality reads, as well as evaluate the quality of sequencing data via several
89  measures, including Q20, sequence duplication levels, coverage, and GC content (Chen et al., 2018).
90  After evaluation, none of the samples were excluded. Clean DNA sequencing reads were mapped to
91  the human reference genome hg19 using the BWA-MEM algorithm (Li and Durbin, 2009).
92  Duplicated reads for public BAM files and our in-house BAM files were removed using PicardTools.
93  Genomic variants were called following the Genome Analysis Toolkit software best practices
94  (McKenna et al., 2010). Variants were filtered by VCFtools (Danecek et al., 2011) with (1) missing
95  in more than 50% of samples; (2) minor allele count < 3; (3) quality < 30 and (4) DP < 5.

## 2.2 Inferring genetic sex with seGMM

97  As expected, five features may be associated with sex, including X chromosome heterozygosity
98  (XH), reads mapped to the X chromosome (Xmap), reads mapped to the Y chromosome (Ymap), the
99  ratio of X/Y counts (XYratio), and the mean depth of exons in the sex-determining region of the Y
100 chromosome (SRY) gene (SRY_dep). Usually, seGMM computes XH as the fraction of all genotypes
101 on the X chromosome with two different allele calls, excluding missing genotypes. Xmap/Ymap was
102 computed as the fraction of high-quality reads (mapq>30) that mapped to the X/Y chromosome in all
103 high-quality reads that mapped to the genome with samtools (Li et al., 2009). XYratio was computed
104 as the ratio of Xmap divided by Ymap. SRY_dep was extracted by mosdepth with  high-quality reads
105 (mapq>30) (Pedersen and Quinlan, 2018). Considering that the panel data may only contain genes in
106 the X or Y chromosome, seGMM allows users to select features put into the GMM model.

107 After extracting features from BAM and VCF files, features were normalized to the same level using
108 the scale function in R. Then, the R package mclust was used to perform model-based clustering via
109 the EM algorithm to classify the samples into different clusters (Scrucca et al., 2016). Samples with
110 an uncertainty greater than 0.1 were considered outliers. In addition, seGMM can infer genetic sex
111 for a single patient or trio samples with additional reference data containing the same features
112 selected to put into the model (Figure 1).

## 2.3 Identifying potential sex abnormity samples

114 To pinpoint sex abnormity samples, we first define these values as mean_xmap_z, sd_xmap_z,
115 mean_ymap_z and sd_ymap_z, where $z \in \{m, f\}$ denotes whether the summaries were conditioned
116 on the genetically determined males or females. Then, we defined the following six gates to classify
117 individuals according to the values above:
118      ● XY gate:

119     ○   mean_xmap_m - 3 sd_xmap_m < x < mean_xmap_m + 3 sd_xmap_m
120     ○   mean_ymap_m - 3 sd_ymap_m < y < mean_ymap_m + 3 sd_ymap_m
121   ● XYY gate:
122     ○   mean_xmap_m - 3 sd_xmap_m < x < mean_xmap_m + 3 sd_xmap_m
123     ○   y > 2 mean_ymap_m
124   ● XX gate:
125     ○   mean_xmap_f - 3 sd_xmap_f < x < mean_xmap_f+ 3 sd_xmap_f
126     ○   mean_ymap_f - 3 sd_ymap_f < y < mean_ymap_f + 3 sd_ymap_f
127   ● XXY gate:
128     ○   x > 2 mean_xmap_f
129     ○   mean_ymap_m- 3 sd_ymap_m < y < mean_ymap_m + 3 sd_ymap_m
130   ● XXX gate:
131     ○   x > 3 mean_xmap_f
132     ○   mean_ymap_f - 3 sd_ymap_f < y < mean_ymap_f + 3 sd_ymap_f
133   ● X gate:
134     ○   x < 0.5 mean_xmap_f
135     ○   mean_ymap_f - 3 sd_ymap_f < y < mean_ymap_f + 3 sd_ymap_f

## 2.4    Comparing performance with existing methods

137 To compare the performance between seGMM and existing methods. First, we downloaded and
138 configured three tools for sex inference: PLINK 1.9, XYalign and seXY. For PLINK 1.9, X
139 chromosome pseudoautosomal region was first splited off with --split-x. Then, -- check-sex was
140 running once without parameters, eyeball the distribution of F estimates, and rerun with parameters
141 corresponding to the empirical gap. For XYalign, following the method described in their published
142 paper, we used the CHROM_STATS module to obtain the depth of the 19 chromosome, X
143 chromosome and Y chromosome. Then, the depth of the X and Y chromosomes was normalized by
144 dividing it by the depth of chromosome 19. Finally, we plotted a scatter plot to assess sex-
145 mismatched samples. For seXY, the first step was obtain the X.ped and Y.ped files with PLINK.
146 Next, sex inference was conducted with seXY using X.ped, Y.ped and the training data set (subjects
147 in prostate cancer and ovarian cancer GWAS) provided by seXY. PLINK was applied to all datasets.
148 Since the target gene panel data of Datasets 2 and 3 do not contain genes located on the Y
149 chromosome, XYalign and seXY were only applied to Dataset 1. In addition, XYalign was also
150 applied to WES and WGS data.

## 2.5    STR analysis for verifying sex

152 The STR analysis was conducted in our own-designed multiplex STR system (modified based on
153 PowerPlex® 16 System), which allows coamplification and four-color detection of sixteen loci
154 (fifteen STR loci and Amelogenin). The primers for Amelogenin were designed as 5'-
155 GTTCAGACGTGTGCTTCAACTTCAGCTATGAGGTAATTTTTC – 3' and 5'-
156 ATCCGACGGTAGTGTCCAACCATCAGAGCTTAAACTGG – 3'. All sixteen loci were amplified
157 simultaneously in a single tube and analyzed in a single injection or gel lane. One primer for each of
158 the vWA, amelogenin, FGA and TPOX loci was labeled with carboxyrhodamine (ROX). The
159 amplicons were separated on an ABI 3730XL Genetic Analyzer, and data were collected using
160 GeenMapper ID v3.2. Since females are XX, only a single peak is observed when testing female
161 DNA, whereas males, which possess both X and Y chromosomes, exhibit two peaks with a 6 bp
162 difference.

## 3    Results

## 3.1   seGMM achieved better sex classification in TGS data than existing methods

We calculated XH, Xmap, Ymap and XYratio for Dataset 1. We found that Ymap and XYratio plot as distinct clusters for the majority of males and females (Figure 2). The distribution of XH and Xmap for males and females has a lot of intersections, suggesting that the accuracy of sex inference is limited if we only include features extracted from the X chromosome. Our results proved this hypothesis: with features extracted only from the X chromosome, seGMM reported 59 samples as outliers, and the accuracy for the remaining samples was only 84.56%. Next, we evaluate the performance of seGMM in Dataset 1 using the 4 features we calculated before. We found that no samples were reported as outliers, and the accuracy increased to 99.52% after including features from the Y chromosome (Table 1). Moreover, looking into different genders, the accuracy of seGMM in females was 98.98%, and that in males was 100%. The only female sample (NA19054) misclassified by seGMM had an XY ratio that closely mirrored those of males.

To assess the performance of seGMM, we also applied PLINK, seXY and XYalign in Dataset 1 to assess the accuracy of these existing tools in inferring sex. We discovered that the distribution of the F coefficient is concentrated between 0-0.9 and without an empirical gap (Supplementary Figure S1). The accuracy of PLINK was 81.44 %, with 94 samples whose predicted sex was clear (Table 1). The accuracy of seXY is only 62.5%. For XYalign, which does not directly provide a predicted sex, we plot the normalized sequence depth of chromosome X and chromosome Y. A couple of females and males are mixed, indicating the loss of accuracy with XYalign compared to seGMM (Supplementary Figure S2). Moreover, we have tested the computation time of different methods using 1 core, 10 cores and 20 cores on a server with 64 Intel(R) Xeon(R) CPU E7-8895 v3 @ 2.60 GHz. We can see that, as expected, seGMM costs much more time than PLINK and seXY, which don't collect features of reads mapped to the X and Y chromosomes. Contrary to PLINK and seXY, with 1 core, seGMM costs fairly time compared to XYalign, while when using 20 cores, seGMM achieves 10 times faster than XYalign. (Supplementary Table S5)

To validate the performance of seGMM in inferring sex from gene panel data, we further applied seGMM to Dataset 2 and Dataset 3. Since the target gene panel data of these datasets do not contain genes located on the Y chromosome, we calculated XH and Xmap as features. In contrast to Dataset 1, the distribution of XH and Xmap for Dataset 2 plot as distinct clusters for the majority of males and females (Figure 3). We compared the performance of seGMM in inferring sex with PLINK. In total, the accuracy of seGMM was 99.10%, while the accuracy of PLINK was 86.58% (Table 2). Additionally, looking into different genders, the accuracy of seGMM in females was nearly 100% and 98.34% in males. The few males misclassified by seGMM have XH values that closely mirror those of females. Furthermore, the accuracy of seGMM in Dataset 3 is 92.31% (Table 2), while the accuracy of PLINK is relatively lower (38.87%).

## 3.2   seGMM has good performance in inferring sex for WES and WGS data

Recently, WES and WGS have shown promise in becoming a first-tier diagnostic test for patients with Mendelian disorders. Therefore, we evaluated the performance of seGMM in inferring sex from WES and WGS data. First, we applied seGMM to publicly available WES data (Dataset 4). The accuracy of seGMM was 100%, indicating that seGMM also has excellent performance in WES data (Table 3 and Supplementary Figure S3). Meanwhile, the accuracy of PLINK was 100%, while the accuracy of XYalign  was 99.65%. NA12413 was mixed with female samples (Supplementary Figure S4).

5

207  In addition, we applied seGMM to our in-house WES data (Dataset 6). The accuracy of seGMM in
208  our in-house WES data was 99.75%, 99.76% for males and 99.74% for females (Supplementary
209  Table S6). Six samples (3 males and 3 females; Table 3 and Figure 4a) were mismatched between
210  SNP-inferred sex and self-reported sex, indicating potential misregistration of clinical information.
211  For PLINK, the accuracy was 99.54%, and 11 samples were mismatched (Table 3 and Figure 4b).
212  The accuracy for XYalign was 99.66%, and 8 samples were mismatched (Figure 4c). Six samples
213  were detected with mismatched sex in the three methods; however, 5 mismatched samples detected
214  with PLINK and 2 mismatched samples detected with XYalign were correct in seGMM (Figure 4d).

215  To verify the real sex of these 6 samples, we performed STR analysis with a sex marker, the
216  amelogenin gene. PCR products generated from the amelogenin gene are widely accepted for use in
217  sex identification. The amelogenin gene is highly conserved and occurs on both the X- and Y-
218  chromosomes. With a 6 bp deletion of the amelogenin gene in the Y chromosome, amplicons
219  generated from the X and Y chromosomes were distinguished from one another when electrophoretic
220  separation was performed to separate STR alleles. The results showed that the real sex matched the
221  seGMM prediction results, proving the pinpoint accuracy of seGMM (Table 4 and Figure 5). Finally,
222  seGMM was conducted on the WGS data (Dataset 5), and the accuracy was 100% (Table 3). The
223  accuracy for PLINK and XYalign is also 100%.

### 3.3  The ability of seGMM in clinical application

225  In clinical practice, individual patient or trio samples are usually sequenced to obtain a molecular
226  diagnosis. However, the GMM model requires a sufficient sample size to ensure the accuracy of
227  classification. To address this problem, seGMM permits users to provide additional reference data.
228  By combining the features from reference data, seGMM can ensure accuracy for clinical application.
229  Taking WES data as an example, using features including Xmap, Ymap, XYratio and normalized
230  SRY_dep (divided by XYratio), we found that with 1000 Genomes data points as a reference, all
231  samples in our in-house WES data were predicted accurately and vice versa.

232  Additionally, approximately 0.25% male and 0.15% female live births demonstrated some form of
233  sex chromosome abnormality. A previous study examined the feasibility of defining five gates to
234  classify individuals according to the normalized X and Y chromosome ratio, calculated on
235  genetically determined males and females, respectively (Turro et al., 2020). Similarly, following this
236  strategy, seGMM can automatically classify samples into 6 sex chromosome karyotypes (XX, XY,
237  XYY, XXY, XXX and X) according to the Xmap and Y map. For publicly available WES and WGS
238  data, none of the samples had sex chromosome abnormalities. For our in-house WES data, 3 samples
239  (HBSY-012-ge, HBSY-012 and XiZ-086) were identified with XYY chromosome karyotypes.

### 4  Discussion

241  This article introduces a sex inference tool, seGMM, to infer genotype sex from NGS data, especially
242  TGS panel data. seGMM integrates several tools and algorithms into a single workflow. Using
243  features extracted from sex chromosomes, seGMM applied unsupervised clustering to classify the
244  samples. Compared to many existing supervised methods that attempt to infer sex by training a
245  logistic regression classifier based on limited available data, seGMM can be applied directly to
246  different types of genomic data. By comparing PLINK, seXY and XYalign, we proved that seGMM
247  outperforms existing tools in inferring sex with TGS panel data and has excellent performance in
248  WES and WGS data.

249  Additionally, compared with including features extracted from only the X chromosome, we
250  discovered that jointing valuable information on the Y chromosome improved the accuracy of
251  inferring sex. Our data suggest that adding probes targeting unique regions of the Y chromosome,
252  particularly the exon of the SRY, which is involved in male-typical sex development (Gubbay et al.,
253  1990; Parma and Radi, 2012), is helpful in inferring genders using TGS panel data.

254  One important innovation of seGMM is that seGMM is adapted to clinical applications that can be
255  applied to individual patients and automatically report sex chromosome abnormality samples. When
256  applying seGMM with reference data, attention should be given to the consistency of the analytical
257  methods between reference data and testing data, since inconsistent analysis methods can introduce
258  bias.

259  There are several limitations for seGMM. First, as a method expected to find sex chromosome
260  abnormal samples, seGMM classifies the samples by calculating the Xmap and Ymap intervals in
261  which they are located. However, the interval is fixed and when the number of male or female
262  samples in the testing dataset is small, the distribution of the Xmap/Ymap is approximately
263  nonnormal and then the sex chromosome karyotypes classification of samples in this case may be
264  inaccurate. Hence, we recommended that if the sample size of the male or female is small, combining
265  with a reference data is an effective strategy to ensure the accuracy of results. Second, the
266  computational time of seGMM depends on many factors, such as the number of features, the number
267  of samples, and the number of threads used. seGMM costs much time in collecting features of reads
268  mapped to the X and Y chromosomes, causing it to run slower than PLINK and seXY but faster than
269  XYalign. The running speed of seGMM can increased by adding cores, and we also consider
270  rewriting it as a GPU program to speed up in the future.

271  In conclusion, we have developed a new tool to infer genetic sex based on a Gaussian Mixture Model
272  called seGMM, which combines stable predictive ability and clinical application. In addition, when
273  the genomic data are TGS, seGMM is one of the best choices for inferring sex, which could meet the
274  needs of clinical genetics.

## 5    Data Availability Statement

276  seGMM is publicly available at https://github.com/liusihan/seGMM and can be installed directly by
277  Conda. Users can also download the source code from GitHub or PyPI and install related software.
278  The data used in this study were retrieved from the 1000 Genomes database (https://ftp-
279  trace.ncbi.nih.gov/1000genomes/ftp/). Further inquiries can be directed to the corresponding authors.

## 6    Conflict of Interest

281  The authors declare that the research was conducted in the absence of any commercial or financial
282  relationships that could be construed as a potential conflict of interest.

## 7    Author Contributions

284  SL developed the tool and drafted the manuscript. YZ, MC and QZ participated in data preprocessing
285  and testing. LW and CW performed the experiment. YL and HG reviewed and revised the
286  manuscript. FB designed and supervised the study and reviewed the manuscript. All authors
287  contributed to the article and approved the submitted version.

## 8    Funding

289 This work was supported by the 1·3·5 project for disciplines of excellence, West China

290 Hospital, Sichuan University.

## 9    Reference

292 Chen, S., Zhou, Y., Chen, Y., and Gu, J. (2018). fastp: an ultrafast all-in-one FASTQ preprocessor.
293        *Bioinformatics* 34(17)**,** i884-i890. doi: 10.1093/bioinformatics/bty560.

294 Danecek, P., Auton, A., Abecasis, G., Albers, C.A., Banks, E., DePristo, M.A., et al. (2011). The
295        variant call format and VCFtools. *Bioinformatics* 27(15)**,** 2156-2158. doi:
296        10.1093/bioinformatics/btr330.

297 Di Resta, C., Galbiati, S., Carrera, P., and Ferrari, M. (2018). Next-generation sequencing approach
298        for the diagnosis of human diseases: open challenges and new opportunities. *EJIFCC* 29(1)**,**
299        4-14.

300 Genomes Project, C., Abecasis, G.R., Altshuler, D., Auton, A., Brooks, L.D., Durbin, R.M., et al.
301        (2010). A map of human genome variation from population-scale sequencing. *Nature*
302        467(7319)**,** 1061-1073. doi: 10.1038/nature09534.

303 Genomes Project, C., Auton, A., Brooks, L.D., Durbin, R.M., Garrison, E.P., Kang, H.M., et al.
304        (2015). A global reference for human genetic variation. *Nature* 526(7571)**,** 68-74. doi:
305        10.1038/nature15393.

306 Gubbay, J., Collignon, J., Koopman, P., Capel, B., Economou, A., Munsterberg, A., et al. (1990). A
307        gene mapping to the sex-determining region of the mouse Y chromosome is a member of a
308        novel family of embryonically expressed genes. *Nature* 346(6281)**,** 245-250. doi:
309        10.1038/346245a0.

310 Guo, H., Wang, T., Wu, H., Long, M., Coe, B.P., Li, H., et al. (2018). Inherited and multiple de novo
311        mutations in autism/developmental delay risk genes suggest a multifactorial model. *Mol*
312        *Autism* 9**,** 64. doi: 10.1186/s13229-018-0247-z.

313 Lee, C., Bae, J.S., Ryu, G.H., Kim, N.K.D., Park, D., Chung, J., et al. (2017). A Method to Evaluate
314        the Quality of Clinical Gene-Panel Sequencing Data for Single-Nucleotide Variant Detection.
315        *J Mol Diagn* 19(5)**,** 651-658. doi: 10.1016/j.jmoldx.2017.06.001.

316 Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows–Wheeler
317        transform. *Bioinformatics* 25(14)**,** 1754-1760. doi: 10.1093/bioinformatics/btp324.

318 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., et al. (2009). The Sequence
319        Alignment/Map format and SAMtools. *Bioinformatics* 25(16)**,** 2078-2079. doi:
320        10.1093/bioinformatics/btp352.

321 Lin, X., Tang, W., Ahmad, S., Lu, J., Colby, C.C., Zhu, J., et al. (2012). Applications of targeted
322        gene capture and next-generation sequencing technologies in studies of human deafness and
323        other genetic disabilities. *Hear Res* 288(1-2)**,** 67-76. doi: 10.1016/j.heares.2012.01.004.

324 McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., et al. (2010). The
325        Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA
326        sequencing data. *Genome Res* 20(9)**,** 1297-1303. doi: 10.1101/gr.107524.110.

327 McPherson, E. (2006). Genetic diagnosis and testing in clinical practice. *Clin Med Res* 4(2)**,** 123-129.
328        doi: 10.3121/cmr.4.2.123.

329 Parma, P., and Radi, O. (2012). Molecular mechanisms of sexual development. *Sex Dev* 6(1-3), 7-17.
330     doi: 10.1159/000332209.

331 Pedersen, B.S., and Quinlan, A.R. (2018). Mosdepth: quick coverage calculation for genomes and
332     exomes. *Bioinformatics* 34(5), 867-868. doi: 10.1093/bioinformatics/btx699.

333 Phillips, K.A., and Douglas, M.P. (2018). The Global Market for Next-Generation Sequencing Tests
334     Continues Its Torrid Pace. *J Precis Med* 4.

335 Phillips, K.A., Douglas, M.P., and Marshall, D.A. (2020). Expanding Use of Clinical Genome
336     Sequencing and the Need for More Data on Implementation. *JAMA* 324(20), 2029-2030. doi:
337     10.1001/jama.2020.19933.

338 Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M.A., Bender, D., et al. (2007).
339     PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J*
340     *Hum Genet* 81(3), 559-575. doi: 10.1086/519795.

341 Qian, D.C., Busam, J.A., Xiao, X., O'Mara, T.A., Eeles, R.A., Schumacher, F.R., et al. (2017). seXY:
342     a tool for sex inference from genotype arrays. *Bioinformatics* 33(4), 561-563. doi:
343     10.1093/bioinformatics/btw696.

344 Saudi Mendeliome, G. (2015). Comprehensive gene panels provide advantages over clinical exome
345     sequencing for Mendelian diseases. *Genome Biol* 16, 134. doi: 10.1186/s13059-015-0693-2.

346 Scrucca, L., Fop, M., Murphy, T.B., and Raftery, A.E. (2016). mclust 5: Clustering, Classification
347     and Density Estimation Using Gaussian Finite Mixture Models. *R J* 8(1), 289-317.

348 Taylor, J.C., Martin, H.C., Lise, S., Broxholme, J., Cazier, J.B., Rimmer, A., et al. (2015). Factors
349     influencing success of clinical genome sequencing across a broad spectrum of disorders. *Nat*
350     *Genet* 47(7), 717-726. doi: 10.1038/ng.3304.

351 Turro, E., Astle, W.J., Megy, K., Graf, S., Greene, D., Shamardina, O., et al. (2020). Whole-genome
352     sequencing of patients with rare diseases in a national health system. *Nature* 583(7814), 96-
353     102. doi: 10.1038/s41586-020-2434-2.

354 Webster, T.H., Couse, M., Grande, B.M., Karlins, E., Phung, T.N., Richmond, P.A., et al. (2019).
355     Identifying, understanding, and correcting technical artifacts on the sex chromosomes in next-
356     generation sequencing data. *Gigascience* 8(7). doi: 10.1093/gigascience/giz074.

357

358

359

360

361

362

363

364

365    **Figure legends**

366    **Fig. 1. Schematic diagram of seGMM**. The input file to seGMM is VCF and BAM file, seGMM
367    will automatically collect features and build the GMM model.

368    **Fig. 2. Distribution of features collected from Dataset 1**. a. Distribution of X chromosome
369    heterozygosity rate (%) between males and females. b. Distribution of reads mapped to the X
370    chromosome (%) between males and females. c. Distribution of reads mapped to the Y chromosome
371    (%) between males and females. d. Distribution of XYratio between male and female.

372    **Fig. 3. Distribution of features collected from Dataset 2**. a. Distribution of X chromosome
373    heterozygosity rate (%) between males and females. b. Distribution of reads mapped to the X
374    chromosome (%) between males and females.

375    **Fig. 4. Accuracy of different methods in inferring sex with our in-house WES data**. a. Sample
376    clustering results of seGMM. b. Distribution of F coefficient. c. Scatter plot of normalized X and Y
377    ratio using XYalign. d. Venn plot of mismatched gender samples detected by the three methods.

378    **Fig. 5. Experimentally verified gender results**.

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395 **Tables**

396 Table 1. Accuracy of different methods in inferring sex with Dataset 1.

|  | Accuracy for all samples (%) | Accuracy for Male (%) | Accuracy for Female (%) |
|---|---|---|---|
| PLINK | 81.44 | 48.28 | 100 |
| seXY | 62.5 | 45.45 | 81.63 |
| seGMM | 99.52 | 100 | 98.98 |

397
398
399
400
401
402 Table 2. Accuracy of different methods in inferring sex with Dataset 2 and Dataset 3.

|  | Dataset 2 | Dataset 3 |
|---|---|---|
| PLINK | 86.58 | 38.87 |
| seGMM | 99.10 | 92.31 |

403
404
405
406
407 Table 3. Accuracy of different methods in inferring sex with WES and WGS data.

|  | PLINK | XYalign | seGMM |
|---|---|---|---|
| 1000G phase3 WES data | 100 | 99.65 | 100 |
| 1000G phase3 high quality WGS data | 100 | 100 | 100 |
| In-house WES data | 99.54 | 99.66 | 99.75 |

408
409
410
411
412
413
414
415
416
417
418
419
420

11

421    Table 4. Accuracy of different methods in inferring sex with WES and WGS data.

| Sample ID | Amelogenin | Self-reported gender | seGMM inferred gender | Experiment validate gender |
|---|---|---|---|---|
| GX-0524 | 209.15 - | Male | Female | Female |
| GX-0946 | 209.06 - | Male | Female | Female |
| GYF-0602212 | 209.04 214.8 | Female | Male | Male |
| GYF-0804464 | 209.06 214.85 | Female | Male | Male |
| GYF-0905712-ge | 209.11 - | Male | Female | Female |
| JL-102 | 209.18 214.92 | Female | Male | Male |

422

Input data

VCF data

WGS

WES

Panel

Genotype array

+

Bam file

Reference data

seGMM

Sample clustering: GMM model

Density

60

40

20

0

2

1

−0.06  −0.04  −0.02  0.00  0.02  0.04

Dir 1

Output

Cluster 2

Cluster 1

Abnormality

♀♂