

# 1                    **Non-invasive Preimplantation Genetic Testing of** 2                    **Embryonic Genome in Spent Culture Medium**

3    Lei Huang<sup>1,2#\*</sup>, Yangyun Zou<sup>3#</sup>, Ruiqi Zhang<sup>4#</sup>, Jin Huang<sup>5,6,7,8#</sup>, Guangjun Yin<sup>1,2,9</sup>,  
4    Quangui Wang<sup>9</sup>, Yingying Xia<sup>3</sup>, Jialin Jia<sup>5,6,7,8</sup>, Zeyu Wu<sup>10</sup>, Dandan Cao<sup>9</sup>, Weiliang  
5    Song<sup>1,2</sup>, Yaqiong Tang<sup>1,2,9</sup>, Kai Liu<sup>1,2,9</sup>, Xiaoran Chai<sup>1,2</sup>, Guo-Bo Chen<sup>11</sup>, Sijia Lu<sup>3\*</sup>, Hao  
6    Ge<sup>1,4\*</sup>, Jie Qiao<sup>1,2,5,6,7,8,12\*</sup> and Xiaoliang S. Xie<sup>1,2,9\*</sup>

- 7
- 8    1. Biomedical Pioneering Innovation Center (BIOPIC), Peking University, Beijing,  
9        100871, China
  - 10    2. Beijing Advanced Innovation Center for Genomics, Peking University, Beijing,  
11        100871, China
  - 12    3. Yikon Genomics Co., Ltd., Shanghai, 310120, China
  - 13    4. Beijing International Center for Mathematical Research (BICMR), Peking  
14        University, Beijing, 100871, China
  - 15    5. State Key Laboratory of Female Fertility Promotion, Centre for Reproductive  
16        Medicine, Department of Obstetrics and Gynecology, Peking University Third  
17        Hospital, Beijing 100191, China
  - 18    6. National Clinical Research Center for Obstetrics and Gynecology, Peking  
19        University Third Hospital, Beijing 100191, China
  - 20    7. Key Laboratory of Assisted Reproduction (Peking University), Ministry of  
21        Education, Beijing 100191, China
  - 22    8. Beijing Key Laboratory of Reproductive Endocrinology and Assisted Reproductive  
23        Technology, Beijing 100191, China
  - 24    9. Changping Laboratory, Beijing, China
  - 25    10. School of Mathematical Sciences, Peking University, Beijing, 100871, China
  - 26    11. Center for Reproductive Medicine, Department of Genetics and Genomic Medicine,  
27        and Clinical Research Institute, Zhejiang Provincial People's Hospital, People's  
28        Hospital of Hangzhou Medical College, Hangzhou, Zhejiang, 310024, China
  - 29    12. Peking-Tsinghua Center for Life Sciences, Peking University, Beijing, 100871,  
30        China

31    # These authors contribute equally to the work.

32    \*Corresponding authors: X. S. Xie ([sunneyxie@biopic.pku.edu.cn](mailto:sunneyxie@biopic.pku.edu.cn)), J. Qiao  
33    ([jie.qiao@263.net](mailto:jie.qiao@263.net)), L. Huang([leihuangmed@gmail.com](mailto:leihuangmed@gmail.com)), H. Ge  
34    ([haoge@pku.edu.cn](mailto:haoge@pku.edu.cn)), S. Lu ([lusijia@yikongenomics.com](mailto:lusijia@yikongenomics.com))

35

36 **NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.**

37

38

## ABSTRACT

39

40 Traditionally, preimplantation genetic testing (PGT) for in vitro  
41 fertilization (IVF) requires invasive trophoctoderm (TE) biopsy, which  
42 might be detrimental to the embryo. Recently proposed non-invasive PGT  
43 (ni-PGT) utilizing cell-free DNA from spent embryo culture medium  
44 (SCM) also faces serious challenges in accuracy, especially for monogenic  
45 diseases (niPGT-M), due to trace DNA content, maternal cell  
46 contamination, and high Allele Drop-Out (ADO) rates. In this study, an  
47 improved linear single-cell whole genome amplification method and a  
48 Bayesian linkage analysis model were used to enhance accuracy in niPGT-  
49 M. We achieved about 75% report rate across all samples and 100%  
50 accuracy in the reported samples. Meanwhile, we reconstructed the  
51 embryonic genome and calculated the risk of type II diabetes (T2D) via  
52 niPGT-P, consistent well with those from TE biopsy samples. Our approach  
53 alleviated the limitations of ni-PGT and offers a promising avenue for  
54 advancing noninvasive PGT with potential clinical applications.

55

56 Keywords: non-invasive preimplantation genetic testing, monogenic  
57 diseases, polygenic diseases, spent embryo culture medium, Bayesian  
58 linkage analysis

59

60

## INTRODUCTION

61 Preimplantation genetic testing (PGT) is widely used in clinical *in vitro*  
62 fertilization (IVF) to screen embryos for monogenic diseases, aneuploidy,  
63 or structural rearrangements. Next-generation sequencing (NGS)-based  
64 PGT allows simultaneous analysis of chromosome copy number, linkage,  
65 and targeted haplotyping for mutations [1-8]. However, traditional PGT  
66 requires trophoctoderm (TE) biopsy, which poses risks to embryo health  
67 and is prone to sampling bias due to embryo mosaicism [9]. Recently,

68 minimally invasive or non-invasive methods using cell-free DNA (cfDNA)  
69 from blastocoele fluid (BF) [10] or spent embryo culture medium (SCM)  
70 [11,12] have been reported as alternatives for PGT with less risk and lower  
71 price. Despite their potential, challenges like trace DNA amount, poor  
72 DNA integrality, and complicated origin led to highly variable success rates  
73 in the published references from different clinical centers [12-15].  
74 Although the success amplification rate in SCM was higher than that in BF  
75 [13,16] using the same amplification method, it still does not meet the yield  
76 required for clinical diagnosis application.

77

78 Even worse, genetic contents amplified from SCM can be potentially  
79 affected by maternal cell contamination (MCC), as maternal DNA is often  
80 longer than embryonic DNA [17,18] thus preferentially being amplified,  
81 and by the presence of high levels of human serum albumin (HSA) protein  
82 in culture medium. Current commercial single-cell amplification kits, such  
83 as DOP-PCR, MDA, or MALBAC are not specially optimized to address  
84 these challenges. In 2017, Linear Amplification via Transposon Insertion  
85 (LIANTI) was developed [19], significantly reducing the false positive rate  
86 of Single Nucleotide Polymorphisms (SNP) detection and much lower  
87 amplification bias caused by fragment lengths. However, the yield of  
88 LIANTI amplification for SCM is still insufficient. There is an urgent need  
89 to develop an efficient and accurate amplification method to obtain  
90 embryonic DNA information from SCM and/or BF.

91

92 Numerous reports have explored on non-invasive preimplantation genetic  
93 testing for aneuploidy (niPGT-A) [11,20,21], but studies on non-invasive  
94 preimplantation genetic testing for monogenic diseases (niPGT-M) remain  
95 limited [12,13,15,16, 22-24]. In niPGT-M studies using SCM, the genotype  
96 concordance with trophoctoderm biopsy results varied widely, from 21%  
97 to 88%, limiting the clinical utility of niPGT-M [12-15]. The key criterion  
98 for evaluating PGT-M is not only genotype concordance but also the  
99 misdiagnosis rate, which is very low (<0.1%) with traditional methods  
100 according to the ESHRE PGT consortium [25]. This emphasizes the need

101 to assess misdiagnosis risk in new niPGT-M approaches.

102

103 Accurate linkage analysis is crucial alongside the amplification method to  
104 pinpoint disease-carrying chromosomes in the embryos and minimize  
105 misdiagnosis, especially given the trace DNA content in SCM. Traditional  
106 linkage analysis typically requires at least two fully informative markers  
107 like STRs or SNPs, near the mutation site, with ADO rates below 5% [26].  
108 Ou et al. used 4 informative SNPs [15], and Liu et al. used 10 [12].  
109 However, since ADO rates in SCM or BF can exceed 5%, relying on a few  
110 informative sites for linkage analysis may lead to misdiagnosis.  
111 Informative SNP markers could be too far from the mutation site,  
112 increasing recombination risk. Meanwhile, maternal cell contamination  
113 can also raise the misdiagnosis rate. These challenges indicate that  
114 traditional linkage analysis methods for niPGT-M may not be suitable for  
115 clinical applications. Therefore, developing a highly accurate linkage  
116 analysis method tailored for high ADO rates and maternal cell  
117 contamination is crucial.

118

119 Furthermore, preimplantation genetic testing for polygenic disorders  
120 (PGT-P) using polygenic risk scores (PRS) holds promise in reducing the  
121 risk of polygenic diseases, such as coronary heart diseases, diabetes, and  
122 cancers at the pre-conception stage [27-29]. Non-invasive PGT-P (niPGT-  
123 P) utilizing genetic information from SCM may, help alleviate some of the  
124 ethical risk-benefit concerns by avoiding the harms of biopsy to the  
125 embryos [30]. However, the challenge lies in the low coverage and inferior  
126 genotype quality of embryo genomes due to limited genetic materials and  
127 poor quality of cell-free DNA in SCM, which hinders the availability of  
128 genome-wide genetic variants necessary for PRS evaluation. Therefore,  
129 accurate whole genome reconstruction is crucial for the successful  
130 implementation of niPGT-P.

131

132 In the present study, we improved the LIANTI method to amplify trace  
133 DNA from SCM, overcoming challenges posed by high-concentration



134 protein existence and amplification bias caused by DNA size difference.  
135 We introduced a Bayesian-based linkage analysis method (BASE-niPGT-  
136 M) that incorporates high ADO rates and maternal cell contamination. This  
137 method could accurately detect disease-carrying chromosomes and provide  
138 a confidence level for each diagnosis. Furthermore, using the fragmented  
139 and incomplete genomic data from SCM, we employed Pedigree-  
140 Population-based Imputation with Haploid Assumption (PPIHA) method  
141 to reconstruct the embryonic genome. After reconstruction, the high  
142 genome-wide coverage and accuracy of genetic information were obtained,  
143 enabling PRS analysis. In summary, our study demonstrates high detection  
144 accuracy for ni-PGT, offering a reliable and potentially applicable non-  
145 invasive method for preimplantation genetic testing intended for clinical  
146 use.

147

148

## RESULTS

### 149 **Characterization of Cell-Free DNA in Spent Embryo Culture Medium**

150 Understanding the properties of DNA in SCM, such as fragment size  
151 and DNA quantity, is crucial for selecting appropriate amplification  
152 methods. Unlike biopsy samples, which typically contain one to three  
153 complete cell nuclei, cfDNA in SCM or BF is present in low quantities and  
154 of poor quality. Our results showed the cfDNA in the culture medium  
155 exhibits highly variable DNA fragment size distribution and distinct DNA  
156 ladder patterns (Extended data Fig. 1 and Supplementary Note S1), posing  
157 a significant challenge for DNA amplification. Based on the above  
158 evidences, an effective amplification method should minimize  
159 amplification bias caused by DNA size variation.

160 Analysis of the sequencing results of cfDNA revealed low genome  
161 coverage in SCM, ranging from 10% to 50% (Extended data Fig. 2), in  
162 contrast to the genome coverage of biopsy samples, which typically  
163 exceeds 95%. This deficiency is not primarily attributed to the  
164 amplification-induced allelic dropout (ADO) but rather due to the large-  
165 scale DNA information loss of paternal or maternal alleles, we called it

166 large fragment loss of haplotype (LFLoH), as shown in Extended data Fig.  
167 3. This LFLoH phenomenon often results in the loss of informative SNPs  
168 both upstream and downstream adjacent to the pathogenic site, which are  
169 typically used in traditional linkage analysis. Informative SNPs are those  
170 where the SNP is heterozygous in the parent carrying the disease and  
171 homozygous in the other parent, and where the SNP remains heterozygous  
172 in the embryo (Extended data Fig. 4).

173 Maternal cell contamination in embryo culture medium has long been  
174 a challenge in diagnostic procedures, often leading to misdiagnosis. Using  
175 informative SNPs with homozygous but different parental genotypes at the  
176 whole-genome level, we observed MCC values ranging from -55.1% to  
177 99.7% (Extended data Fig. 5). Negative values may result from copy  
178 number gains in segments of paternal chromosomes or losses of maternal  
179 chromosomal segments. Employing an iterative strategy within a 10 Mb  
180 range upstream and downstream of pathogenic loci in our Bayesian linkage  
181 analysis method (See below and Methods), we observed the modified  
182 MCC values ranging from 0.1% to 99.7%. Aside from the absence of  
183 negative values, with this region, 58 samples have an MCC of less than 2%,  
184 indicating minimal maternal cell contamination, 43 samples have relatively  
185 low levels of maternal cell contamination between 2% and 30%, and only  
186 10 samples have maternal cell contamination levels between 30% and 80%,  
187 with only 6 samples exceeding 80% (Extended data Fig. 5, Supplementary  
188 Table S4).

189 It's important to note that the cfDNA in SCM may potentially  
190 originate from cell apoptotic processes. Using a ligation-based method  
191 (Supplementary Notes S1), we observed a clear DNA ladder in the culture  
192 medium, as shown in Extended data Fig. 1.

193

## 194 **Whole Genome Amplification Method for SCM**

195 Given the unique characteristics of cfDNA in SCM, a whole genome  
196 amplification method with minimal fragment size bias is essential. We  
197 compared the fragment size amplification biases of the current WGA  
198 methods for PGT-M and found that LIANTI exhibited the least bias (short:

199 long = 1:5), while MDA and MALBAC showed a stronger preference for  
200 amplifying longer fragments, with short: long amplification biases of 1:500  
201 and 1:6000, respectively (Supplementary Table S1). LIANTI is a linear  
202 amplification method known for its high SNP detection accuracy [19].

203 However, the original LIANTI method has low efficiency in  
204 amplifying such minute quantities of DNA, failing to meet clinical  
205 requirements. To address this, we improved the LIANTI method by  
206 modifying the lysis step, adding more transposons for insertion,  
207 introducing DNA primers during first-strand DNA synthesis, and  
208 implementing exponential amplification during second-strand synthesis, as  
209 depicted in Extended data Fig. 6 (see Methods for details).

210 The improved LIANTI method achieved a 100% amplification  
211 success rate for clinical samples and demonstrated high SNP detection  
212 accuracy, as shown in Extended data Fig. 7 and Supplementary Table S2.  
213 We compared the sequencing results of cfDNA from BJ cell line culture  
214 medium collected at different time points using the improved LIANTI  
215 method and other scWGA methods. The improved LIANTI method  
216 demonstrated superior amplification efficiency for cfDNA compared to the  
217 other scWGA methods.

218

### 219 **Linkage Analysis Method for Noninvasive PGT-M of Embryos**

220 Figure 1 illustrates the complete workflow employed in this study.  
221 Briefly, data collection was divided into two stages (Figure 1A): the Pre-  
222 phasing stage and the SCM genetics information acquisition stage. In the  
223 Pre-phasing stage, family trio sequencing was performed using DNA from  
224 the father, mother, and proband or alternatively from grandparents or  
225 discarded embryos. During the SCM genetics information acquisition stage,  
226 the DNA content in SCM was amplified using the improved LIANTI  
227 method, and then was followed by sequencing and SNP calling (See Fig.  
228 1A and Methods).

229 With these data in hand, we present a computational method BASE-  
230 niPGT-M (BAYesian linkage analySis mEthod for non-invasive  
231 Preimplantation Genetic Testing of Monogenic disorders) (Fig. 1B,

232 Extended data Fig. 8 and Supplementary Note S2). This approach  
233 addresses the challenges posed by MCC and LFLoH in SCM samples. The  
234 Bayesian model takes the sequencing data of SCM and the phased  
235 haplotypes of the parents, as input. Samples with extremely low quality are  
236 filtered out. Unlike previous approaches for high-quality sequencing data,  
237 which often predetermined the number of SNPs to be considered, our  
238 model incrementally incorporates SNPs, starting from the disease-causing  
239 mutation site. This process allows us to calculate the log-likelihood ratio  
240 of inheriting disease-carrying chromosomes versus disease-free  
241 chromosomes for each SNP subset, ultimately resulting in a log-likelihood  
242 ratio curve (Fig. 1B, Extended data Fig. 8 and Supplementary Note S2).  
243 This curve, with its distinct characteristics, enables us to classify the results  
244 into four categories: Highly Confident, Moderately Confident, Possibly  
245 confident, and Undetermined. In the first three categories, we can pinpoint  
246 the inherited chromosome.

247

### 248 **Case Study of our niPGT-M Method**

249 We selected one case (Family 2) to illustrate the SNP distribution of  
250 cfDNA and the results of our niPGT-M method, which consists of WGA  
251 and NGS followed by BASE-niPGT-M. This case involved an autosomal  
252 dominant genetic disease, where the husband carried a genetic variant for  
253 Marfan syndrome, a heritable connective tissue disorder. This disorder is  
254 known for its characteristic features such as elongated limbs, aortic  
255 aneurysms, and mitral valve prolapse. Genetic diagnosis of the husband  
256 revealed a mutation c.643C>T at the FBN1 gene, known to cause Marfan  
257 syndrome. Eight embryos of the couple developed to the blastocyst stage,  
258 and several trophectoderm cells were biopsied from each embryo for  
259 traditional PGT-M. All eight SCM samples were collected, and cfDNA was  
260 amplified using the improved LIANTI method. The linkage analysis results  
261 of these eight SCMs are shown in Fig. 2A. The detected SNP distribution  
262 of 2T-SCM-01 and 2T-SCM-02 appeared relatively scattered, with longer  
263 intervals between informative SNPs (Fig. 2A), and the SNPs flanking the  
264 causal mutation site indicated different haplotypes from the father (Fig. 2B),

265 leading to their classification as 'undetermined'. In contrast, 2T-SCM-03,  
266 2T-SCM-04, 2T-SCM-05, 2T-SCM-06 and 2T-SCM-08 showed denser  
267 SNP distributions, and SNPs upstream and downstream of the causal  
268 mutation site indicated the same haplotype (Fig. 2B), resulting in a  
269 classification of 'Highly confident'. The detected SNPs of SCM-07 were  
270 dense, and the haplotype upstream and downstream was concordant.  
271 However, an upstream SNP located relatively close to the causal mutation  
272 site pointed to the opposite haplotype (Fig. 2B), leading to a classification  
273 of 'moderately confident' due to our conservative diagnosis strategy. All  
274 these results were consistent with those of trophoctoderm biopsies,  
275 demonstrating no misclassifications.

276

### 277 **Validation Using all Matched Biopsy/ Whole Embryo.**

278 To further evaluate the performance of our niPGT-M method, we  
279 compared the results obtained from all SCM samples to the gold standard,  
280 which was either TE-biopsy or discarded embryo PGT-M results (Table 1).  
281 We categorized the cases into three groups: autosomal dominant disease,  
282 autosomal recessive disease, and X-linked recessive disease. For  
283 autosomal recessive diseases, we compared each allele of embryos.

284 In the seven cases with autosomal dominant diseases, there were 45  
285 SCM samples. For embryos that the father carried the causal mutation, 23  
286 out of 28 SCMs (82.1% report rate) were successfully diagnosed. For  
287 embryos that the mother carried the causal mutation, 11 out of 17 SCMs  
288 (64.7% report rate) were successfully diagnosed. In the five cases with  
289 autosomal recessive disorders, 36 SCMs were analyzed. We successfully  
290 detected the paternal allele in 25 SCMs (69.4% report rate) and the  
291 maternal allele in 27 SCMs (75% report rate). For the two cases with X-  
292 linked recessive disorders, 8 out of 9 SCMs (88.9% report rate) were  
293 diagnosed. Overall, we successfully detected 94 out of 126 alleles (74.6%  
294 average report rate). The average report rate for paternal alleles was 75%,  
295 while the average report rate for maternal alleles was 74.2%.

296 Importantly, all detected results were concordant with the gold  
297 standard, demonstrating a 100% accuracy rate of BASE-niPGT-M (Table

298 1).

299

## 300 **Whole Genome Reconstruction of Embryos**

301 The directly detected genetic information of embryos obtained from  
302 the SCM samples is much limited due to low amount of cfDNA materials  
303 released from blastocyst cells. We therefore developed a Pedigree-  
304 Population-based Imputation with Haploid Assumption (PPIHA) to  
305 acquire the missing genetic information. This approach allows us to  
306 reconstruct the whole genome information of embryos from SCM samples.

307 The initial density of detected sites in the raw data was  $222\pm 127$  SNPs  
308 per megabase (Mb). However, after performing PPIHA (shown in Fig. 3A),  
309 the median coverage of SNPs significantly improved from 6.4% to 93.5%,  
310 and the density of imputed sites increased to  $1479\pm 554$  SNPs/Mb (Fig. 3B  
311 and 3D).

312 To assess the accuracy of the predicted embryo's haplotype obtained  
313 from SCM samples, we compared it with the haplotype of corresponding  
314 discarded embryos or TE biopsy samples. When the MCC of SCM was low,  
315 we observed a high similarity of haplotype configuration between SCM  
316 and the corresponding TE biopsy. Conversely, SCM samples with high  
317 MCC exhibited reduced haplotype similarity, especially in the maternal  
318 haplotypes (Fig. 3C). For a more comprehensive exploration, we obtained  
319 43 paired samples with whole genome sequencing data from both SCMs  
320 and their corresponding discarded embryos or TE biopsy samples for  
321 comparison. To avoid bias we removed paired samples with failed quality  
322 control or aneuploid embryos for further analysis (Supplementary Table  
323 S3). Paternal or Maternal haplotype concordance at each physical site  
324 between SCM and discarded embryos/TE biopsy samples was applied to  
325 measure the haplotype accuracy constructed in SCM (Fig. 3E). We found  
326 that the concordance of maternal and paternal haplotypes in SCM versus  
327 discarded embryos/TE was around 98.6% and 99.8% respectively, for  
328 SCM samples without significant maternal cell contamination (Fig. 3E).  
329 High levels of maternal cell contamination could largely impact the  
330 haplotype construction of SCM, especially for maternal haplotypes (Fig.



331 3E).

332 Furthermore, after filtering out SNPs with Mendelian errors, we  
333 compared the constructed genome-wide genotypes from SCM samples  
334 with those of corresponding discarded embryos/TE biopsy samples. The  
335 concordance of genotypes ranged from 79.2% to 99.9%, partially impacted  
336 by maternal cell contamination in SCM (Fig. 3E). For those SCM without  
337 significant maternal cell contamination, the average genotype accuracy  
338 was approximately 96.9%. Additionally, we compared the constructed  
339 genotypes from two SCM samples with the genotypes of corresponding  
340 amniotic fluid samples, which exhibited concordance rates of 94.5% and  
341 94.1%, respectively.

342

### 343 **Polygenic Risk Evaluation of Type II Diabetes for Embryos**

344 We used type II diabetes (T2D) as a complex disease model to assess  
345 the feasibility of polygenic risk evaluation based on constructed whole  
346 genome information of SCM. After applying various filters, including  
347 removing SCM samples with high maternal cell contamination,  
348 chromosomal abnormality or poor-quality, 25 SCM samples were left for  
349 polygenic risk score (PRS) calculation of T2D. We employed a PRS model  
350 of T2D from a previous study [31]. The PRS of T2D model involves a total  
351 of 114 loci, where only 28.7% of the loci can be directly detected, 52.5%  
352 being obtained through family-based genome reconstruction, and the  
353 remaining 18.8% being imputed by population-based reconstruction (Fig.  
354 3B, Fig. 4A).

355 When comparing the genotypes of SNPs in the T2D PRS model  
356 between the SCM samples and corresponding discarded embryos or TE  
357 biopsy samples, we found that median 97.4% of the SNPs were identical.  
358 Among these identical SNPs, 52.5% are homozygous REF genotypes, 24.9%  
359 are homozygous ALT genotypes, and 22.6% are heterozygous genotypes  
360 (Fig. 4B).

361 The PRS of T2D calculated from SCM showed a high level of  
362 concordance with those calculated from the corresponding discarded  
363 embryos, TE biopsy samples, or amniotic fluid, with an R-squared value

364 of 0.79 (Fig. 4C). In clinical practice, PRS percentiles are typically used to  
365 evaluate the relative genetic risk of common diseases instead of raw PRS  
366 scores. To determine an individual's PRS percentile for a particular disease,  
367 a reference population comprising both patients (e.g., with T2D) and  
368 healthy individuals is needed to establish the PRS distribution, which then  
369 allows each individual to be placed within a specific percentile of the  
370 distribution. For example, if an individual's PRS is in the 90th percentile,  
371 it means he/she has a higher PRS than 90% of the reference population.  
372 Generally, the higher the PRS percentile, the greater the risk of a genetic  
373 disease. We used a reference population from the UK Biobank  
374 (<https://www.ukbiobank.ac.uk/>) to construct the PRS distribution for T2D  
375 and calculated each embryo's PRS percentile. We observed a correlation of  
376 T2D PRS percentiles between SCM samples and the corresponding  
377 discarded embryos, TE biopsy samples, or amniotic fluid, indicated by an  
378 R-squared value of 0.68 (Fig. 4D).

379

380

## DISCUSSION

381 In this study, we introduced a method for non-invasive  
382 preimplantation genetic testing, addressing key challenges in the accuracy  
383 of genetic assessments using SCM. We improved a linear single-cell whole  
384 genome amplification method to successfully amplify DNA content in  
385 SCM and introduced a Bayesian linkage analysis model to improve the  
386 accuracy of ni-PGT for monogenic diseases, matching the reliability of  
387 invasive procedures. Our approach advances the potential clinical  
388 implementation of ni-PGT, reducing the risks associated with traditional  
389 PGT while maintaining high accuracy.

390 PGT-P holds promise for reducing the incidence of complex genetic  
391 diseases at the population level. However, its widespread use is hindered  
392 by controversies regarding ethical concerns, scientific validity, and  
393 practical issues, particularly the potential harm of the biopsy procedure to  
394 embryos if the benefits of PGT-P are not conclusive. Non-invasive PGT-P,  
395 which detects genetic materials from embryos through SCM, may

396 accelerate scientific research and clinical implementation by significantly  
397 reducing potential risks to embryos. We achieved comprehensive whole-  
398 genome construction of embryos, accurately assessing the risk of polygenic  
399 diseases (niPGT-P) in line with results from trophoctoderm (TE) biopsy  
400 samples. Our PPIHA approach to reconstruct whole genome information  
401 of embryos from SCM samples fully considers the genetic characteristics  
402 of SCM specimens, such as LFLoH, which violates the general biological  
403 principles of diploidy, and the low amount and quality of cfDNA released  
404 from blastocysts. The high genome-wide coverage and accuracy of genetic  
405 information in SCM samples constructed through PPIHA enable non-  
406 invasive assessment of polygenic disease risk in embryos, advancing the  
407 prevention and control of common diseases to the pre-implantation stage.

408 Maternal cell contamination from cumulus cells or polar bodies in  
409 SCM is an issue that requires attention, which would impact the accuracy  
410 of niPGT-M and niPGT-P. Special techniques for the removal of maternal  
411 cell contamination are necessary to be developed in the future. At current  
412 phase, we created a quantitative estimation of MCC rate which plays a  
413 warning role when it is high, and ensures accurate ni-PGT results when it  
414 is not high.

415 In contrast to TE samples that contain alleles from two haplotypes,  
416 SCM only provides alleles from a single haplotype in many genomic  
417 regions, resulting in a high rate of allelic dropout and random segments of  
418 single haplotypes in the SCM data. To address this issue, we developed a  
419 specialized approach to first estimate the haplotype status of each SNP,  
420 specifically whether only the parental or maternal haplotype detected.  
421 Assuming the independence of haplotype status among measured SNPs  
422 would introduce an excessive number of parameters, leading to the risk of  
423 overfitting. Therefore we have adopted a cautious approach: the default  
424 assumption for the haplotype status of each SNP is set to be “both parental  
425 haplotypes are detected”, unless the SNP data strongly contradicts this  
426 default.

427 The accuracy of haplotype pre-phasing of the parents is important for  
428 both niPGT-M and niPGT-P. In cases where no proband is available,

429 sequencing data from discarded embryos must be used for pre-phasing.  
430 However, the data quality, particularly the coverage, obtained from  
431 discarded embryos is often limited, which can hinder the accuracy of  
432 haplotype phasing results for the parents. Improved precision in haplotype  
433 phasing of the parents, possibly achievable through more advanced  
434 sequencing technologies, would enhance the accuracy of niPGT-M and  
435 niPGT-P.

436 In the context of clinical applications, our methodology extends  
437 beyond merely estimating inherited chromosomes in the embryo; it also  
438 provides a measure of confidence in this estimation. By utilizing a  
439 Bayesian model, the log-likelihood ratio curve enables a nuanced  
440 classification of our diagnoses. This approach is crucial for aiding both  
441 medical professionals and parents in determining the most suitable course  
442 of action for the embryo. The precision of our method, combined with the  
443 refined classification outputs, significantly enhances the reliability of  
444 identifying healthy embryos. This enhancement underscores the potential  
445 translational impact of our work, promising increased report rates in  
446 pregnancies and the delivery of healthy infants.

447

## 448 REFERENCES

- 449 1. Yan, L. et al. Live births after simultaneous avoidance of monogenic diseases and  
450 chromosome abnormality by next-generation sequencing with linkage analyses.  
451 *Proc Natl Acad Sci U S A*. 112(52):15964-9 (2015)
- 452 2. Esteki, Z.M. et al. Concurrent whole-genome haplotyping and copy-number  
453 profiling of single cells. *Am J Hum Genet*. 96(6):894-912 (2015)
- 454 3. Xiong, L. et al. Bayesian model for accurate MARSALA (mutated allele revealed  
455 by sequencing with aneuploidy and linkage analyses). *J Assist Reprod Genet*. 36(6):  
456 1263–1271 (2019)
- 457 4. Backenroth, D. et al. Haploseek: a 24-hour all-in-one method for preimplantation  
458 genetic diagnosis (PGD) of monogenic disease and aneuploidy. *Genet Med*.  
459 21(6):1390-1399 (2019)
- 460 5. Masset, H. et al. Single-cell genome-wide concurrent haplotyping and copy-number  
461 profiling through genotyping-by-sequencing. *Nucleic Acids Res*. 50(11):e63 (2022)
- 462 6. Yan, L. et al. Chinese experts' consensus guideline on preimplantation genetic  
463 testing of monogenic disorders. *Hum Reprod*. 38(Supplement\_2):ii3-ii13 (2023)

- 464 7. Backenroth, D. et al. SHaploseek is a sequencing-only, high-resolution method for  
465 comprehensive preimplantation genetic testing. *Sci Rep.* 13(1):18036 (2023)
- 466 8. De Coster, T. et al. Genome-wide equine preimplantation genetic testing enabled by  
467 simultaneous haplotyping and copy number detection. *Sci Rep.* 14(1):2003 (2024)
- 468 9. McCoy, R.C. Mosaicism in preimplantation human embryos: when chromosomal  
469 abnormalities are the norm. *Trends Genet.* 33(7):448-463 (2017) .
- 470 10. Palini, S. et al. Genomic DNA in human blastocoele fluid. *Reprod Biomed Online*  
471 26, 603–610 (2013).
- 472 11. Xu, J. et al. Noninvasive chromosome screening of human embryos by genome  
473 sequencing of embryo culture medium for in vitro fertilization. *Proc Natl Acad Sci*  
474 113, 11907–11912 (2016).
- 475 12. Liu, W.Q. et al. Non-invasive preimplantation aneuploidy screening and diagnosis  
476 of beta thalassemia IVSII654 mutation using spent embryo culture medium. *Ann*  
477 *Med* 49, 319–328 (2017)
- 478 13. Capalbo, A. et al. Diagnostic efficacy of blastocoel fluid and spent media as sources  
479 of DNA for preimplantation genetic testing in standard clinical conditions. *Fertil*  
480 *Steril* 110, 870–879 (2018).
- 481 14. Cimadomo, D. et al. The dawn of the future: 30 years from the first biopsy of a  
482 human embryo. The detailed history of an ongoing revolution. *Human*  
483 *Reproduction Update.* 26, 453-73 (2020).
- 484 15. Ou, Z. et al. Improved non-invasive preimplantation genetic testing for beta-  
485 thalassemia using spent embryo culture medium containing blastocoelic fluid.  
486 *Frontiers in endocrinology* 12, 1941 (2022).
- 487 16. Galluzzi, L. et al. Extracellular embryo genomic DNA and its potential for  
488 genotyping applications. *Futur Sci OA* 1, FSO62 (2015).
- 489 17. Lo, Y.M.D. et al. Maternal plasma DNA sequencing reveals the genome-wide  
490 genetic and mutational profile of the fetus. *Sci. Transl. Med.* 2, 61ra91 (2010)
- 491 18. Yu, S.C.Y. et al. Size-based molecular diagnostics using plasma DNA for  
492 noninvasive prenatal testing. *Proc. Natl. Acad. Sci. USA* 111, 8583–8588 (2014)
- 493 19. Chen, C. et al. Single-cell whole-genome analyses by Linear Amplification via  
494 Transposon Insertion (LIANTI). *Science* 356, 189-94 (2017).
- 495 20. Huang, L. et al. Noninvasive preimplantation genetic testing for aneuploidy in spent  
496 medium may be more reliable than trophectoderm biopsy. *Proc. Natl. Acad. Sci.*  
497 *USA* 116, 14105-14112 (2019)
- 498 21. Leaver, M. & Wells, D. Non-invasive preimplantation genetic testing (niPGT): the  
499 next revolution in reproductive genetics? *Human reproduction update* 26, 16-42  
500 (2020).
- 501 22. Wu, H. et al. Medium-based noninvasive preimplantation genetic diagnosis for

- 502 human  $\alpha$ -thalassemias-SEA. *Medicine (Baltimore)* 94, e669 (2015).
- 503 23. Zhang, Y. et al. Molecular analysis of DNA in blastocoele fluid using next-  
504 generation sequencing. *J Assist Reprod Genet* 33, 637–645 (2016).
- 505 24. Shangguan, T. et al. Detection and analysis of DNAmaterial in human blastocoele  
506 fluid. *Biomed Genet Genomics* 2, 1–5 (2017).
- 507 25. De Rycke, M. et al. ESHRE PGD Consortium data collection XIV-XV: Cycles from  
508 January 2011 to December 2012 with pregnancy follow-up to October 2013. *Hum.*  
509 *Reprod.* 32, 1974–1994 (2017).
- 510 26. Carvalho, F. et al. ESHRE PGT Consortium good practice recommendations for the  
511 detection of monogenic disorders. *Human Reproduction Open* 2020(3), hoaa018  
512 (2020).
- 513 27. Treff, N.R. et al. Utility and first clinical application of screening embryos for  
514 polygenic disease risk reduction. *Frontiers in Endocrinology* 10, 845 (2019).
- 515 28. Kumar, A. et al. Whole-genome risk prediction of common diseases in human  
516 preimplantation embryos. *Nature Medicine* 28, 513–516 (2022).
- 517 29. Lennon, N.J. et al. Selection, optimization and validation of ten chronic disease  
518 polygenic risk scores for clinical implementation in diverse US populations. *Nature*  
519 *Medicine* 30, 480–487 (2024).
- 520 30. Turley, P. et al. Problems with Using Polygenic Scores to Select Embryos. *New*  
521 *England Journal of Medicine* 385, 78–86 (2021).
- 522 31. Scott, R.A. et al. An Expanded Genome-Wide Association Study of Type 2 Diabetes  
523 in Europeans. *Diabetes* 66, 2888–2902 (2017).
- 524

## 525 METHODS

### 526 Institutional Review Board Approval.

527 This study was approved by the Peking University Third Hospital  
528 Ethical Review Committee (Approval no. 2020SZ-005). The  
529 patients/participants provided written informed consent to participate in  
530 this study.

### 531 Study Subjects.

532 Patients (n=14 couples) using ICSI and PGT-M with TE biopsy for  
533 preventing their inherited genetic diseases at Peking University Third  
534 Hospital during the period from September 2019 to April 2021 were  
535 included. All of the 14 patients had genetic testing reports.

### 536 ART Laboratory Protocols and Blastocyst Biopsy for PGT-M.

537 Standard protocols were used for ICSI. Embryos were cultured



538 individually in 25- $\mu$ L microdrops of the equilibrated Vitrolife G-series  
539 media with human serum albumin (HSA) (LifeGlobal) overlain with  
540 mineral oil in incubators with 6% CO<sub>2</sub> and 5% O<sub>2</sub> at 37 °C. Embryos were  
541 evaluated on day 5 or 6 for trophoctoderm (TE) biopsy. A few TE cells  
542 from each hatching blastocyst were biopsied and sent to PGT lab for PGT-  
543 M testing.

#### 544 **Collection of Culture Medium Sample and Discarded Embryo**

545 Immediately after the blastocysts were transferred to another dish for  
546 biopsy, 20- $\mu$ L spent embryo culture media were removed from each of the  
547 residual spent medium drops. The discarded embryos and their  
548 corresponding spent embryo culture media samples were also collected for  
549 analysis. Each discarded embryo was gently moved with a pipette tip to the  
550 edge of its microdrop and then removed and transferred into a RNase-  
551 DNase-free PCR tube containing 3  $\mu$ L of lysis buffer. Pipette tips were  
552 changed between collections of each sample to avoid cross-sample  
553 contamination. Media drops incubated and collected under identical  
554 conditions to those used for blastocyst culture, but never containing  
555 embryos, served as negative controls. All samples were immediately frozen  
556 after collection and stored at -80 °C until analyzed.

#### 557 **Spent Embryo Culture Medium Lysis and Discarded Embryo Lysis**

558 The frozen 20- $\mu$ L spent embryo culture medium samples were thawed  
559 and gently mixed, and 10.8  $\mu$ L was removed to a PCR tube (Maximum  
560 Recovery, Axygen) containing 1.2  $\mu$ L 10x lysis buffer (1x lysis buffer: 60  
561 mM Tris-Ac pH 8.3, 2 mM EDTA pH 8.0, 15 mM DTT, 0.5  $\mu$ M carrier  
562 ssDNA (5'-TCAGGTTTTCCTGAA-3', Thermo Fisher Scientific oligo  
563 with PAGE purification)), spun down. It was heated at 75°C for 30 min,  
564 then incubated at 75°C for 30 min. 0.5  $\mu$ L 35mg/mL QIAGEN protease  
565 (dissolved in water and stored at 4°C) was added, spun down, and  
566 incubated at 55°C for 2 hrs followed by 80°C for 30 min. The resulting  
567 lysate in PCR tubes could be immediately subject to improved linear  
568 transposon-based amplification as described herein or stored in a freezer  
569 for later use.

570 The frozen discarded embryos were thawed and heated at 75°C for 30

571 min. Then the discarded embryo samples were lysed by adding 0.5  $\mu$ L  
572 5mg/mL Qiagen Protease and heating (2 h at 55  $^{\circ}$ C and 30 min at 80  $^{\circ}$ C)  
573 in 3  $\mu$ L of lysis buffer.

#### 574 **Whole Genome Amplification**

575 Transposon DNA  
576 (5'/Phos/CTGTCTCTTATACACATCTGAACAGAATTTAATACGACTC  
577 ACTATAGGGGAGATGTGTATAAGAGACAG-3', Thermo Fisher  
578 Scientific oligo with PAGE purification) was annealed by gradual cooling  
579 in annealing buffer (20 mM Tris-Ac pH 8.3, 50 mM NaCl, 2 mM EDTA  
580 pH 8.0) into a 1.5  $\mu$ M self-looping structure. The 1.5  $\mu$ M annealed  
581 transposon DNA was then mixed with an equal volume of  $\sim$ 1  $\mu$ M Tn5  
582 transposase (Lucigen, EZ-Tn5<sup>TM</sup> Transposase) and incubated at room  
583 temperature for 30 min, dimerizing into transposomes with a final  
584 concentration of  $\sim$ 0.25  $\mu$ M. The transposomes can be stored at -20 $^{\circ}$ C for a  
585 long time, and each single-cell transposonbased amplification needs  $\sim$ 0.5  
586  $\mu$ L transposome.

587 Starting with the culture medium lysate, a 20  $\mu$ L transposition mixture  
588 was assembled in the buffer containing 2.5 mM MgCl<sub>2</sub> and 18.75 nM  
589 transposome. The transposition reaction was carried out at 55 $^{\circ}$ C for 12 min.  
590 Then 0.84  $\mu$ L mixture of EDTA and ssDNA was added to each tube and  
591 was incubated at 68  $^{\circ}$ C for 30 min. After removing transposase reaction,  
592 0.2  $\mu$ L Q5 HF DNA Polymerase (New England Biolabs) was added and  
593 was heated at 73 $^{\circ}$ C for 45 seconds in the presence of 2 mM MgCl<sub>2</sub> and  
594 200  $\mu$ M dNTPs for fragments end filling and extension. 0.5  $\mu$ L 4mg/mL  
595 Qiagen protease was added and the PCR tube was heated at 50  $^{\circ}$ C for 1  
596 hour in the presence of 4 mM EDTA to inactivate DNA polymerase,  
597 followed by protease heat inactivation by 77 $^{\circ}$ C for 20 min in the presence  
598 of 450 mM NaCl in a total volume of  $\sim$ 25  $\mu$ L. DNA fragments were  
599 amplified to RNAs in a 90  $\mu$ L T7 in vitro transcription reaction mixture  
600 overnight at 37  $^{\circ}$ C as described in Cheng et al's paper [19].

601 The next day, 10  $\mu$ L 0.5M EDTA was added to each tube and RNAs  
602 were column purified (Zymo Research). 18  $\mu$ L RNAs were transferred to  
603 a PCR tube containing 3.1  $\mu$ L mixture of 6.5mM each dNTPs, 9.7  $\mu$ M

604 carrier ssDNA (5'-TCAGGTTTTTCCTGAA-3', Thermo Fisher Scientific  
605 oligo with PAGE purification) and 3.2 U/ $\mu$ L SUPERase In RNase Inhibitor  
606 (Invitrogen). The denaturation incubation was at 70 °C 1min, 90 °C 15s,  
607 and followed by ice quenching. The reverse transcription reaction for the  
608 first stand was carried out in a 30  $\mu$ L SuperScript IV reserve reaction  
609 system including 0.67 mM each dNTPs, 0.6 U/ $\mu$ L SUPERase In RNase  
610 Inhibitor (Invitrogen) and 6 U/ $\mu$ L SuperScript IV Reserve Transcriptase  
611 (Invitrogen) in SuperScript IV buffer, with the primer 5'-  
612 AGATGTGTATAAGAGACAG-3', Thermo Fisher Scientific oligo with  
613 PAGE purification, and the incubation program was 25°C 1 min, 37°C 1  
614 min, 42°C 1 min, 50°C 1 min, 55°C 15 min, 60°C 10 min, 65°C 12 min,  
615 70°C 8 min, 75°C 5 min, 80°C 10 min. RNA was then removed by  
616 incubation at 37°C for 30 min with 10 ng/ $\mu$ L affinity-purified RNase A  
617 (Invitrogen) and 0.08 U/ $\mu$ L RNase H (New England Biolabs). Second  
618 strand synthesis was carried out in a 100  $\mu$ L Q5 DNA polymerase system  
619 (New England Biolabs, 1X Q5 reaction buffer, 1X Q5 High GC enhancer,  
620 200  $\mu$ M dNTPs, 0.5  $\mu$ M primer, 0.02 U/ $\mu$ L Q5 DNA polymerase), with the  
621 primer 5'-NNNNNNNNGGGAGATGTGTATAAGAGACAG-3', Thermo  
622 Fisher Scientific oligo with PAGE purification. Each tube was heated at  
623 98°C for 30 s, then using 10 cycles of 10 s at 98°C, 30s at 58°C, 30s at  
624 60°C, 30s at 65°C, 2.5 min at 70°C, then 72°C 6 min for strand extension.  
625 The resulted amplicons were column purified into 23  $\mu$ L elution buffer and  
626 stored at -20°C. The DNA concentration of the product after amplification  
627 was measured using a Qubit 2.0 fluorometer (ThermoFisher Scientific)  
628 with the Qubit dsDNA HS Assay kit (Life Technologies).

629 Starting with the discarded embryo lysate, DNA was amplified  
630 following the LIANTI amplification method described in [19], except the  
631 DNA primer 19\_1 5'-AGATGTGTATAAGAGACAG-3', Thermo Fisher  
632 Scientific oligo with PAGE purification was added in the reaction of the  
633 first strand cDNA synthesis.

#### 634 **Sequencing and Data Analysis.**

635 Upon library preparation, the amplicons were subject to a  
636 conventional sonication process (Covaris S2) to the length required by

637 sequencing platforms. The final insert size was suitable for 2x150 bp pair-  
638 end Illumina sequencing. Library preparation was performed by NEBNext  
639 Ultra II DNA Library Prep Kit for Illumina (New England Biolabs),  
640 following the instructions of the manufacturer and skipping the optional  
641 size selection step. For the gDNA of blood samples, all the samples were  
642 subject to a conventional sonication process (Covaris S2). Library  
643 preparation was performed by a PCR-free Library Prep Kit, NEBNext  
644 Ultra II DNA Library Prep Kit for Illumina (NEB #E7645S/L).

645 The bulk samples, culture medium samples, and discarded embryo  
646 samples were sequenced on the Illumina NovaSeq 6000 System or Illumina  
647 HiSeq HiSeq 4000 platform with 2x150 bp pair-end sequencing.  
648 Sequencing lanes were shared by 24 samples with NEBNext Indexes. Each  
649 gDNA sample was sequenced to generate ~90 Gb raw data and each culture  
650 medium sample or discarded embryo was sequenced to generate ~30 Gb  
651 raw data.

## 652 **Mapping and SNP Calling**

653 Raw paired-end reads were trimmed adapters and low-quality ends by  
654 cutadapt. Clean reads were mapped to human reference genome hg19 with  
655 BWA mem. Each BAM file was marked duplicates and sorted by  
656 MarkDuplicatesSpark of GATK 4.2.0. Base quality score recalibration  
657 (BQSR) was performed using BaseRecalibrator to generate a recalibration  
658 table, followed by ApplyBQSR to adjust the base quality accordingly. For  
659 variant calling, GATK HaplotypeCaller produced a genomic variant call  
660 format (GVCF) file for each sample. GenotypeGVCFs was utilized to  
661 perform joint genotyping across all samples. SNPs were filtered using  
662 variant quality score recalibration (VQSR) with HapMap 3.3, Omni 2.5,  
663 1000 Genome phase I, and dbSNP 151 as SNP training sets. A 99%  
664 sensitivity threshold was chosen to filter SNPs accurately.

665 We exclusively retain biallelic SNPs for subsequent analysis. To  
666 deduce the two parental haplotypes, we employ Plink 1.9 with the option  
667 "--mendel" to identify sites with Mendelian errors.

668

## 669 **Parental haplotype phasing**

670 For parental haplotype phasing, genetic information of pedigrees such as  
671 sibling embryos or other family members are needed. Likelihood-based  
672 haplotyping approach using Hidden Markov Model (HMM) and dynamic  
673 programming Viterbi algorithm were used to determine the most likely  
674 haplotype configuration of parents [32,33]. The transition probability of  
675 haplotype state changes between consecutive loci in the HMM was  
676 calculated from the recombination fraction obtained from the genetic  
677 distance map of the 1000 Genomes Project Phase 3 references  
678 (<https://www.internationalgenome.org/>). To address the ADO issue  
679 stemming from scWGA bias during parental haplotyping with sibling  
680 embryos, we filtered out sites that violated Mendelian rules or phased  
681 variants that did not adhere to chromosome interference theory [34]. Each  
682 chromosome of parent was then independently processed until the parental  
683 haplotypes of all chromosomes were phased.

684

#### 685 **BASE-niPGT-M: a Bayesian linkage analysis method for niGPT-M**

686

687 **Bayesian model.** We present BASE-niPGT-M, a Bayesian linkage analysis  
688 method tailored for niPGT-M. Our approach incorporates the MCC rate  
689 and haplotype status of each SNP into the likelihood function. To enhance  
690 the robustness of our data, SNPs with extremely low quality are filtered  
691 out. Then the sequencing error rate, MCC rate and haplotype status were  
692 carefully estimated. Subsequently, we calculated the likelihood of  
693 observing SNP data from the SCM samples within regions proximal to the  
694 disease-causing mutation site. The calculation is based on the single-SNP  
695 likelihood, which was attained from a binomial model using the measured  
696 allelic depth (AD) values and estimated parameters, and the recombination  
697 probability that can be obtained from the DECODE dataset or other  
698 datasets containing the recombination fractions. See Extended data Fig. 8  
699 and Supplementary Notes S2 for more details.

700

701 **Recursively estimating the MCC rate and haplotype status.** MCC rate  
702 was defined as the ratio of DNA fragments originating from the maternal

703 chromosome. Haplotype status, representing the true parental origin of  
704 DNA at each measured SNP, is categorized into three distinct classes:  
705 'Paternal Chromosome Only,' 'Maternal Chromosome Only,' or 'Both  
706 Parental Chromosomes.' For instance, the designation 'Paternal  
707 Chromosome Only' indicates that the DNA identified at a specific SNP  
708 exclusively originates from the father. Initial MCC rate and haplotype  
709 status were estimated, and these parameters were then iteratively and  
710 jointly calibrated until convergence (See Supplementary Notes S2 for  
711 details).

712

713 **Determining the inherited chromosome.** Unlike previous publications,  
714 where a fixed number of SNPs were predetermined, our methodology  
715 involves a sequential addition of SNPs. Starting from the disease-causing  
716 mutation site, we incrementally incorporate SNPs one by one, calculating  
717 the log-likelihood ratio for each subset. This iterative process yields a curve  
718 of log-likelihood ratio plotted against the physical distance of the terminal  
719 SNP.

720

721 Typically, the curve initiates at a point where the ordinate is in close  
722 proximity to zero, gradually stabilizing into a plateau with the addition of  
723 sufficient SNPs. The distinctive properties of this curve enable the  
724 determination of the inherited chromosome. Subsequently, the confidence  
725 level of our disease-carrying analysis is categorized into four classes:  
726 Highly Confident, Moderately Confident, Possibly confident, and  
727 Undetermined, based on the characteristics of the curve. See  
728 Supplementary Notes S2 for more details.

729

### 730 **Haplotype phasing and whole genome reconstructing of SCM**

731 In many regions of the genome from the SCM samples, genetic  
732 information from paternal and maternal origins of the embryo was not both  
733 detected. This was due to allelic dropout resulting from scWGA bias or  
734 LFLoH. The conventional diploid assumption of the genome could not be  
735 applied to haplotyping and genotype imputation of SCM samples.



736 Therefore, we developed a Pedigree-Population-based Imputation with  
737 Haploid Assumption (PPIHA) to reconstruct the whole genome  
738 information of SCM. Specifically, we combined pre-phased parental  
739 haplotype information and selected informative SNPs where at least one  
740 parental origin haplotype could be confidently phased, to create the  
741 haplotype scaffold of SCM (Extended Data Fig. 4). For example, for an  
742 informative locus with paternal genotype being AB (allele A comes from  
743 haplotype P1 and allele B from P2, based on the pre-phased parental  
744 haplotype), the maternal genotype being AA, and the genotype called from  
745 SCM being BB, it could infer that the B allele of this variant in SCM is  
746 derived from haplotype P2, based on the Mendelian rule (see Extended data  
747 Fig. 4 and Fig. 9). Maternal origin haplotype was not determined at this  
748 site due to ADO bias or LFLoH feature in SCM. And vice versa, for the  
749 locus with maternal, paternal and SCM genotypes being AB, AA and BB,  
750 respectively, maternal haplotype of SCM could be inferred while paternally  
751 genetic information unavailable. After examining all such kind of  
752 informative variants from the sequenced SCM data, we constructed a  
753 sparse haplotype scaffold with single parent of SCM (Fig. 3A and Extended  
754 data Fig. 9).

755 Paternal or maternal haplotypes in the scaffold were first corrected if  
756 continuous variants of the opposite haplotype were observed on both  
757 nearby sides. This correction was applied when the probability of a double  
758 crossover between the nearest variants with the opposite haplotype,  
759 calculated as the square of their genetic distances multiplied by 0.0001,  
760 was less than 0.005 for the maternal haplotype or 0.001 for the paternal  
761 haplotype, according to chromosome interference theory ([34], Extended  
762 Data Fig. 9). Then, with parental haplotype support and HMM strategy,  
763 paternal and maternal haplotypes of other variants were independently  
764 phased based on pre-determined haplotype scaffold of SCM (Extended  
765 Data Fig. 9). We assumed that the state (P1 or P2 for paternal haplotypes  
766 and M1 or M2 for maternal haplotypes) at variant  $t$  could be inferred by  
767 SNPs in the scaffold from  $t-5$  to  $t+5$ . The transition probability between  
768 two SNPs was considered equivalent to their genetic distance multiplied

769 by 0.01. Joint probabilities of the two states of paternal or maternal  
770 haplotypes at variant  $t$  were calculated using a forward-backward  
771 algorithm based on the adjoining ten variants (from  $t-5$  to  $t+5$ ). The  
772 haplotype of variant  $t$  was then assigned to the state with the higher  
773 probability. If the joint probabilities of the two states (P1/P2 or M1/M2)  
774 for a variant were equal, phasing of that variant failed (Extended Data Fig.  
775 9). By organizing the haplotype states of all variants according to their  
776 chromosomal positions, the pedigree-based haplotypes of SCMs were  
777 reconstructed. Subsequently, the genotypes of SCM were imputed with the  
778 genome information of the parents (Fig. 3A, Extended Data Fig. 9).

779 Following pedigree-based genome reconstruction, population-based  
780 genotype imputation was conducted using the 1000 Genomes haplotype  
781 reference panel with Minimac4 [35], resulting in the reconstruction of the  
782 whole genome information of the embryo from the SCM samples (Fig. 3A).

### 783 **Polygenic risk score calculation**

784 We employed a PRS model of T2D including 128 SNPs from a  
785 previous study [31]. To validate the model and determine the PRS  
786 percentile (indicating relative disease risk) in the UK Biobank (UKB)  
787 cohort, SNPs whose genotype not available or failing to be imputed in  
788 UKB cohort were excluded, leaving a total of 114 SNPs. We then directly  
789 used the effect sizes from the original study to calculate the PRS,  
790 employing the following formula:

$$791 \quad \text{PRS} = \sum_{i=1}^{114} \beta_i \times G_i$$

792 where  $G_i$  is the allelic dosage and  $\beta_i$  is the effect size of SNP  $i$ .

793 Using the T2D model and control information in UKB cohort, PRS  
794 percentile for T2D was obtained to evaluate the polygenic risk of T2D for  
795 each individual.

796

797

798

## 798 **DATA AVAILABILITY**

799 Access to anonymized patient data is subject to a data-sharing

800 agreement and protocol approval from the institutional review board  
801 committee. The study-specific analyzable dataset is, therefore, not  
802 publicly available. All summary statistics are available in  
803 Supplementary Tables.

804

805

### CODE AVAILABILITY

806 Code for BASE-niPGT-M is available at [https://github.com/Ge-lab-](https://github.com/Ge-lab-pku/BASE-niPGT-M)  
807 [pku/BASE-niPGT-M](https://github.com/Ge-lab-pku/BASE-niPGT-M)

808

809

### REFERENCES

- 810 32. Abecasis, G.R. et al. Merlin---rapid analysis of dense genetic maps using sparse  
811 gene flow trees. *Nat Genet* 30, 97-101 (2002).
- 812 33. Lander, E.S. and Green, P. Construction of multilocus genetic linkage maps in  
813 humans. *Proc Natl Acad Sci U S A* 84, 2363-2367 (1987).
- 814 34. Housworth, E.A. and Stahl, F.W. Crossover interference in humans. *Am J Hum*  
815 *Genet* 73, 188-197 (2003).
- 816 35. Howie, B. et al. Fast and accurate genotype imputation in genome-wide association  
817 studies through pre-phasing. *Nat Genet* 44, 955-959 (2012).

818

819

### ACKNOWLEDGEMENTS

820 We thank Long Gao, Yuhang Huan, Ye Li, Cuiping Pan, Xiaodan Shi,  
821 Wenjie Sun, Zhongwei Wang, Liying Yan, Jingyi Yang, Jiakun Zhang  
822 and Nannan Zhang for their help. This work is funded by Beijing  
823 Advanced Innovation Center for Genomics at Peking University and  
824 Changping Laboratory. H.G. is supported by National Natural Science  
825 Foundation of China (No. 11971037 and T2225001). This study was  
826 also supported by the National Key R&D Program of China  
827 (2023YFC2705604), the National Science Foundation of China  
828 (82071721 and 82371706) and the special fund of the National  
829 Clinical Key Specialty Construction Program, P. R. China (2023).

830

831

## AUTHOR CONTRIBUTIONS

832 X.S.X., J.Q., L.H. and S.L. conceived and supervised the research. J.H.  
833 and J.J collected SCM samples and performed PGT-M experiments of  
834 biopsy samples. L.H. designed and performed the niPGT-M  
835 experiments with the help from G.Y., D.C., Y.T. and K.L.. R.Z. and  
836 H.G. performed the lineage analysis for niPGT-M with the help from  
837 L.H., Q.W., Z.W., W.S. and X.C. Y.Z. and Y.X. performed the data  
838 analysis for niPGT-P with the help from L.H. and G.C.. L.H., H.G.,  
839 Y.X. and Y.Z. wrote the initial manuscript with the help from Q.W.  
840 and G.Y. All authors approved the final version of the manuscript.

841

842

## COMPETING INTERESTS

843 L.H., H.G., R.Z. and X.S.X. are coinventors on patent application  
844 PCT/CN2023/125605 that includes the experimental discovery and  
845 BASE-niPGT-M in this manuscript. S.L., Y.Z. and Y.X. are current  
846 employees of Yikon Genomics.

847

848

## FIGURE CAPTIONS:

849

850 **Figure 1 niPGT-M workflow using SCM samples.**

851 **A.** Haplotype pre-phasing of parents is conducted using sequencing results  
852 from the proband, grandparents, or at least a discarded embryo within the  
853 family. Fertilized embryos are cultured, and immediately after transferring  
854 the blastocysts to another dish for biopsy, SCMs are removed from each of  
855 the residual drops. A few trophectoderm cells from each hatching  
856 blastocyst were biopsied and sent to an external PGT lab for PGT-M testing,  
857 serving as the gold standard. All the SCM samples underwent scWGA and  
858 NGS, followed by SNP calling.

859 **B.** The SNP data of SCM samples, together with the phased haplotypes of  
860 the parents, are input into the BASE-niPGT-M method. After iteratively

861 estimating the MCC rate and haplotype states of each SNP, the log-  
862 likelihood ratio of inheriting disease-carrying chromosomes versus  
863 disease-free chromosomes was calculated. niPGT-M diagnose is  
864 performed based on the characteristics of the log-likelihood ratio curve as  
865 a function of distance to the disease-causing mutation site.

866

## 867 **Figure 2 An illustrative case for our niPGT-M method.**

868 **A.** The SNP distribution and the linkage-analysis results of cfDNA in an  
869 illustrative case. This case involved an autosomal dominant genetic  
870 disease, where the husband carried a genetic variant for Marfan  
871 syndrome, a heritable connective tissue disorder. All eight SCM  
872 samples were collected, and cfDNA was amplified using the improved  
873 LIANTI method.

874 **B.** The curves of log-likelihood ratio and diagnosis results using BASE-  
875 niPGT-M. The diagnosis results of five SCM samples resulted in a  
876 classification of 'Highly confident', while one SCM sample was  
877 classified into 'Moderately confident' due to our conservative diagnosis  
878 strategy. The remaining two SCM samples were classified into  
879 'Undetermined' due to either a very low density of SNPs or an opposite  
880 inference of chromosome upstream and downstream of the disease-  
881 carrying mutation site. All these results were consistent with those of  
882 trophoctoderm biopsies, demonstrating no misclassifications.

883

## 884 **Figure 3 Whole genome reconstruction of the embryo from SCM data.**

885 **A.** Schematic illustration of the procedure for whole genome  
886 reconstruction. This involves four steps. Firstly, parental haplotypes were  
887 phased with sibling embryos or other family members. Secondly, using  
888 pre-phased parental haplotype information and selected informative SNPs,  
889 haplotype scaffold of SCM was created. Thirdly, pedigree-based  
890 imputation was conducted on SCM with the genome information of parents.  
891 Finally, whole genome information of SCM was reconstructed based on  
892 population-based genotype imputation.

893 **B.** Whole-genome SNP density and T2D related SNPs coverage of 10T-

894 SCM1. The blue triangles were SNPs in PRS model of T2D directly called  
895 from raw data; the orange triangles were SNPs imputed from pedigree  
896 information; the gray triangles represent SNPs imputed from population  
897 information.

898 **C.** Haplotypes of SCM samples with different MCC and corresponding TE  
899 biopsy samples. Different haplotypes were showed with different colors.  
900 The uncolored areas denoted regions where reconstruction failed.

901 **D.** Whole-genome coverage of SCM samples before and after our  
902 construction. See Supplementary Table S6 for the data.

903 **E.** Haplotype and genotype concordances between SCM samples and the  
904 corresponding discarded embryos or TE biopsy samples. See  
905 Supplementary Table S6 for the data.

906

907 **Figure 4 T2D associated genotype imputation and polygenic risk score**  
908 **calculation.**

909 **A.** Distribution of reconstructed genotypes of T2D associated SNPs in  
910 SCM samples. T2D related genotypes in SCM were obtained from direct  
911 detection by NGS in cfDNA from SCM (blue bar), pedigree-based  
912 genotype imputation (green bar) and population-based genotype  
913 imputation (yellow bar).

914 **B.** Genotype concordance of T2D associated SNPs in SCM versus  
915 corresponding TE or discarded embryo samples. Genotype 0/0, 1/1 and 0/1  
916 respectively denote homozygous reference allele, homozygous alternate  
917 allele and heterozygous genotypes.

918 **C.** PRS consistency of T2D between SCM and paired TE/discarded  
919 embryo or amniotic fluid samples.

920 **D.** PRS percentile of T2D between SCM and paired TE/discarded embryo  
921 or amniotic fluid samples. 20 SCM samples matched with TE samples  
922 (green round dots), 5 SCM samples matched with discarded embryo  
923 samples (blue rhombic dots), and 2 of the 20 TE biopsied embryos had  
924 corresponding amniotic fluid samples, also be compared with the  
925 corresponding SCM samples (red triangular dots).

926 See Supplementary Table S7 for the data shown in this figure.

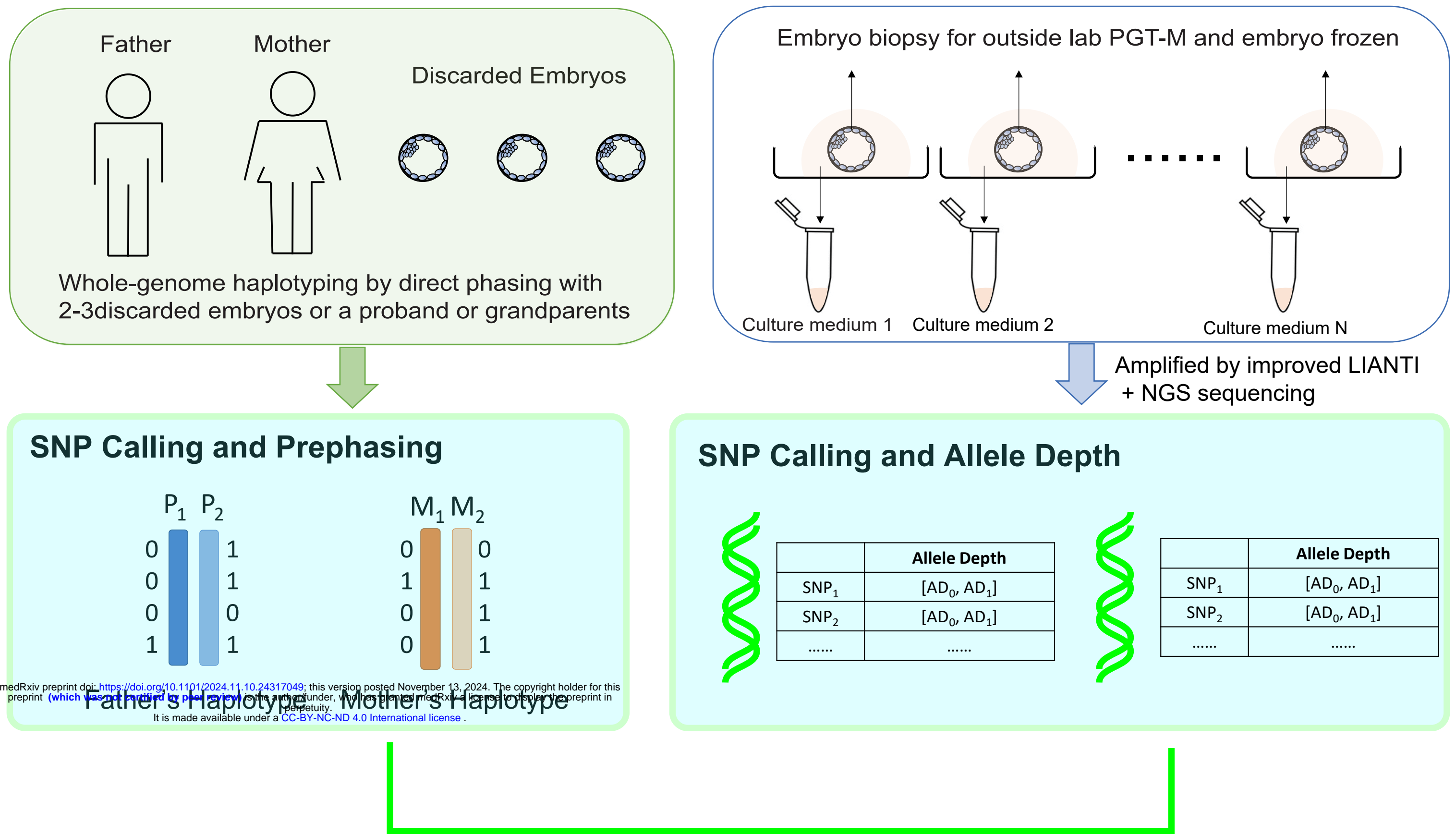


927

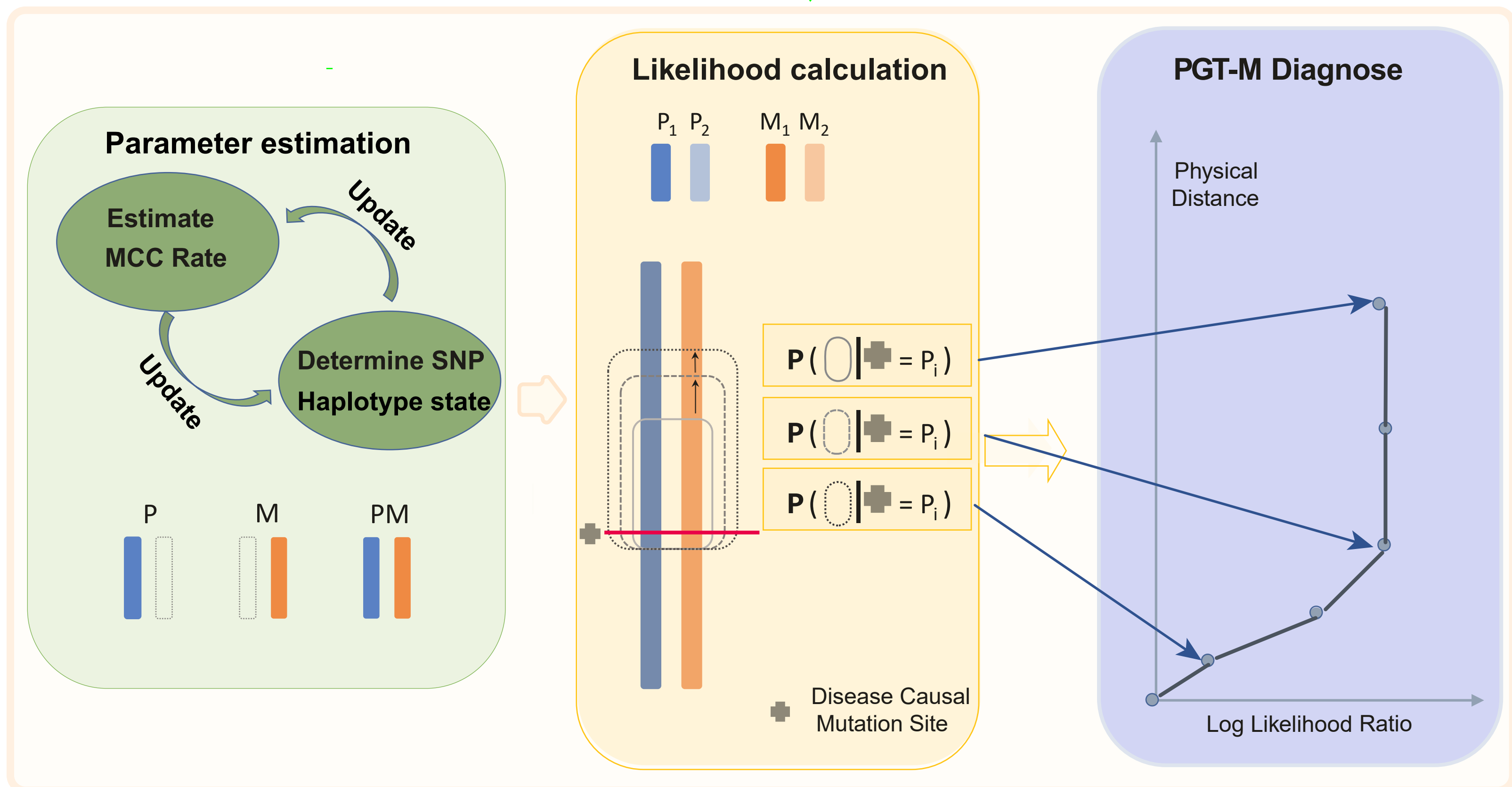
928 **Table 1. Summarized results of BASE-niPGT-M using all matched**  
929 **biopsy/discarded embryo.** We categorized the results into three groups:  
930 autosomal dominant disease, autosomal recessive disease, and X-linked  
931 recessive disease. For autosomal recessive diseases, we compared each  
932 allele of the embryos. The acronyms HC, MC, PO, UN, and RO stand for  
933 Highly Confident, Moderately Confident, Possibly Confident,  
934 Undetermined, and Ruled Out, respectively. Overall, we successfully  
935 detected 95 out of 126 alleles, achieving a 75.4% average report rate. All  
936 detected results were concordant with the gold standard, demonstrating a  
937 100% accuracy rate of BASE-niPGT-M. See Supplementary Table S5 for  
938 more detailed data.

939

**A**

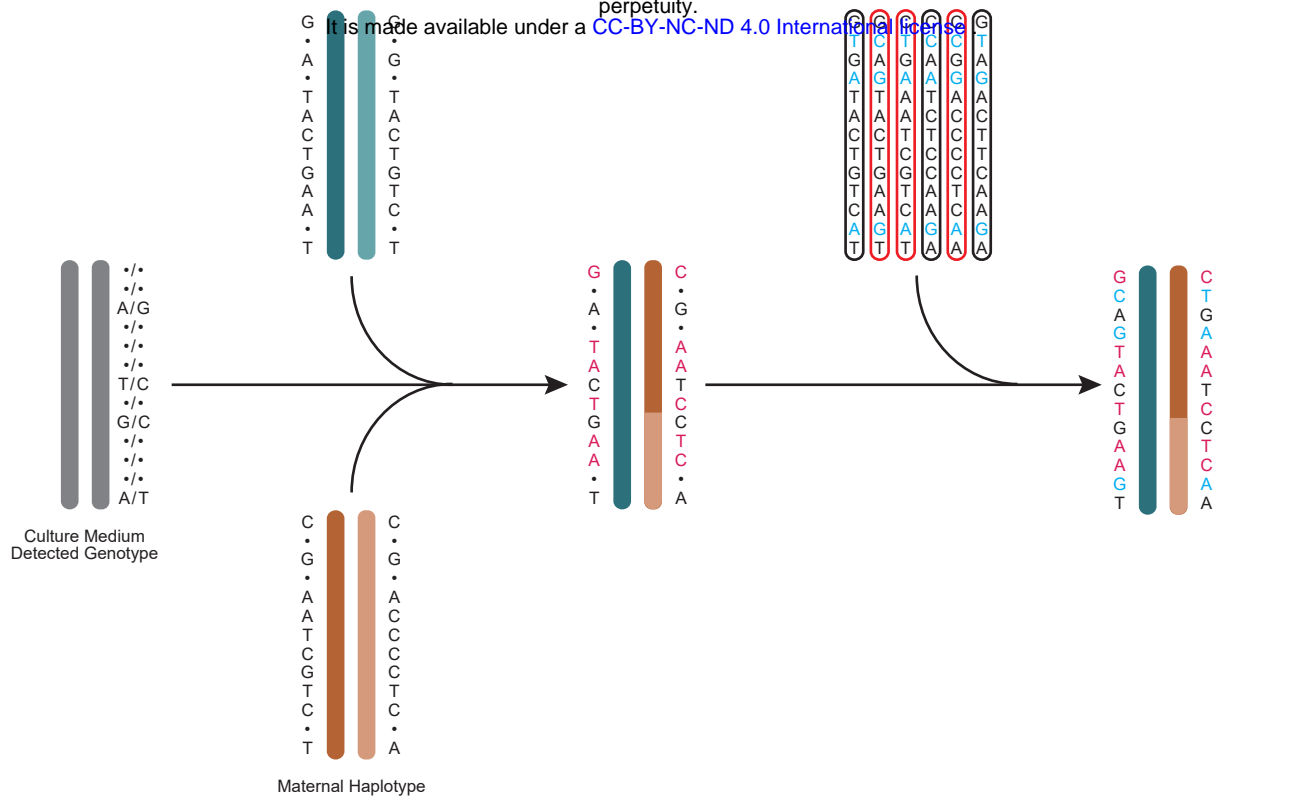


**B**

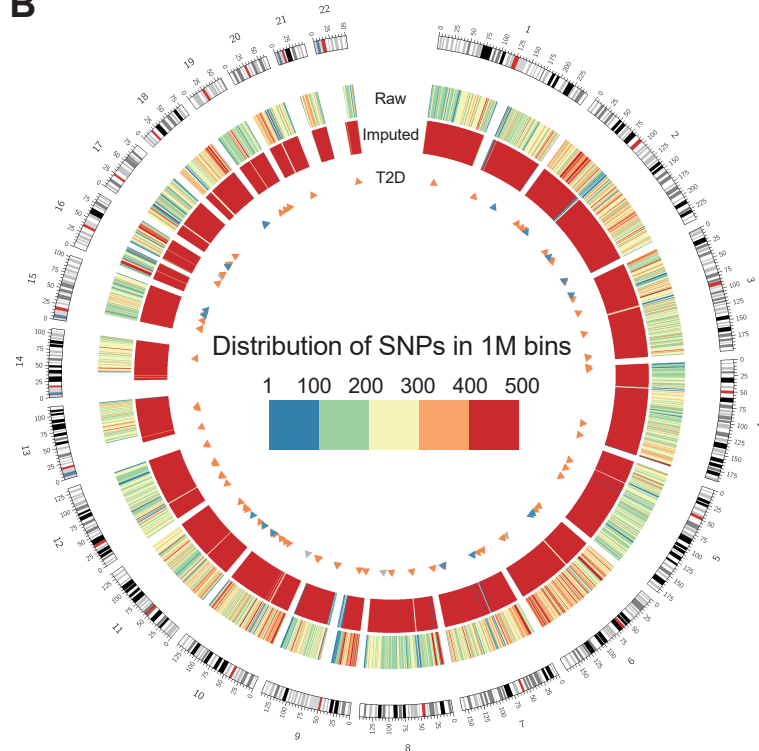




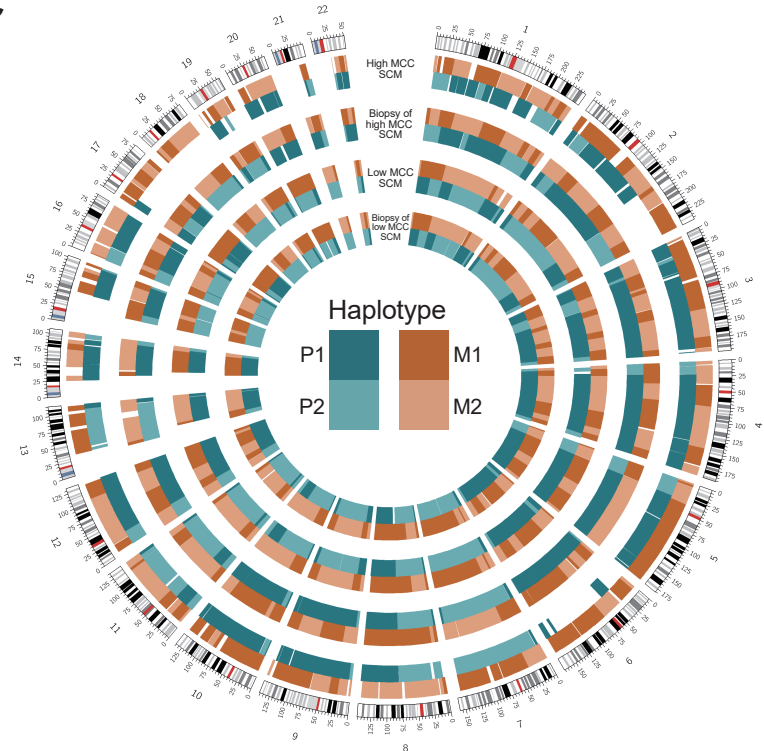
**A**



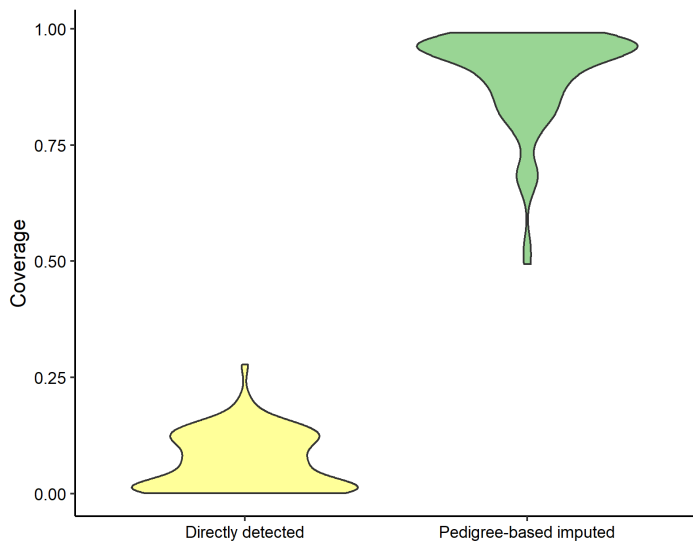
**B**



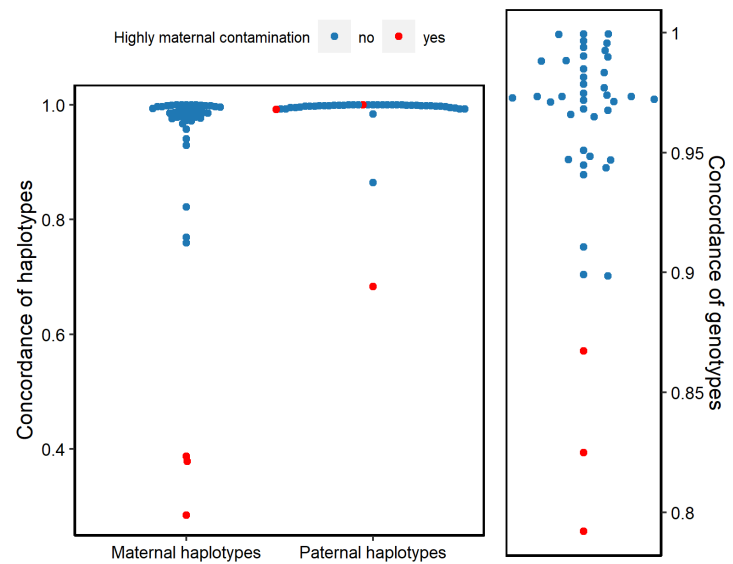
**C**



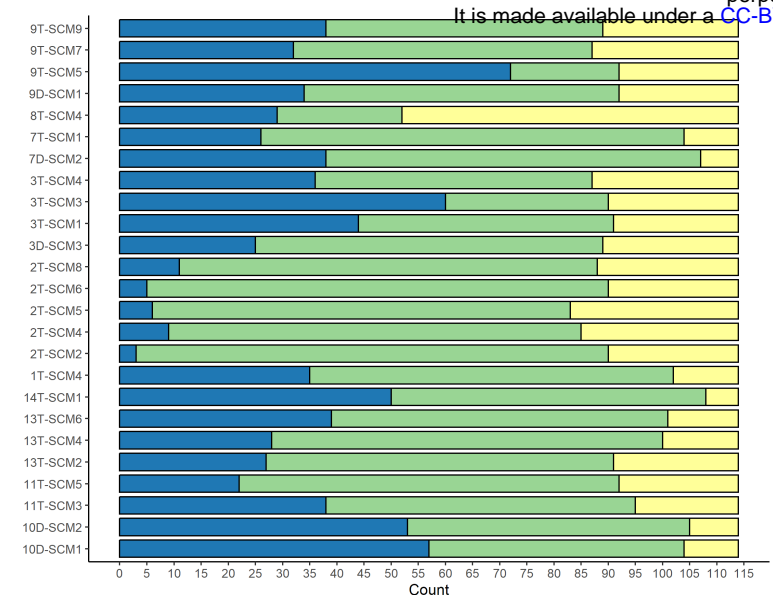
**D**



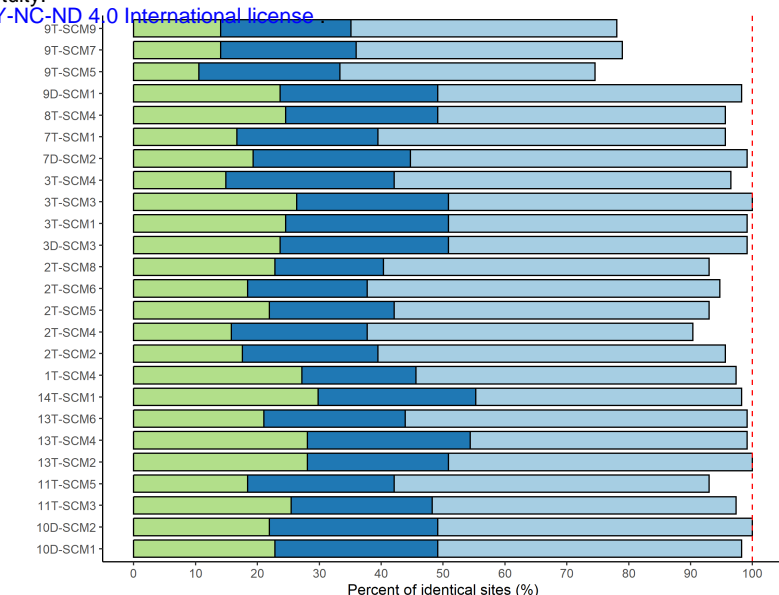
**E**



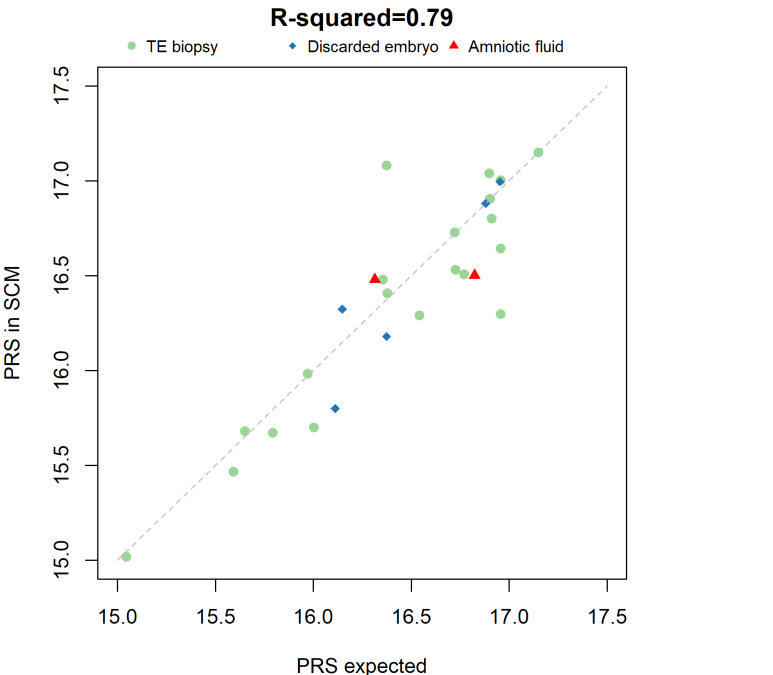
**A**



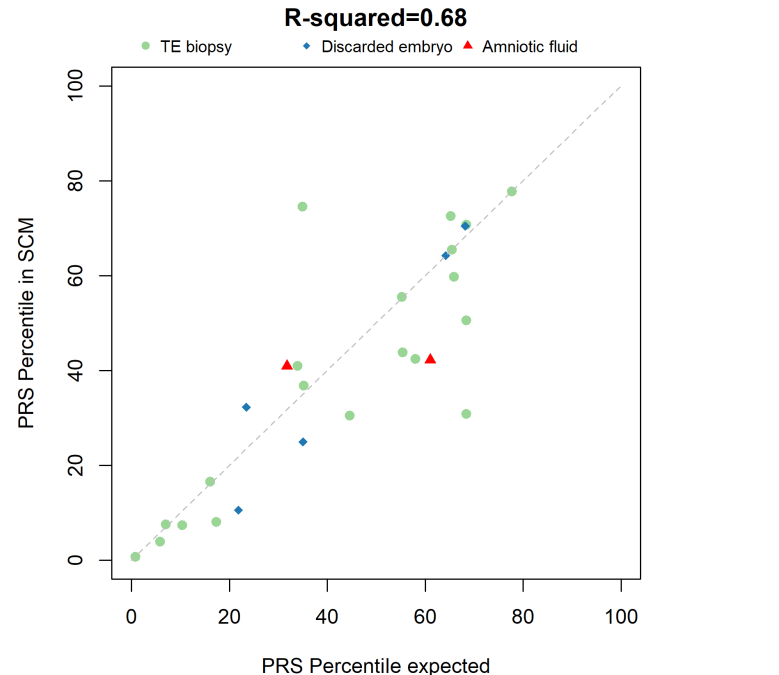
**B**



**C**



**D**



Type of Disease	Family Number	Parent	# of Samples	# of Reportable Samples	# of Correct Samples	HC	MC	PO	UN	RO	Report Rate(%)
Autosomal Dominant Disorder	1	M	6	2	2	2	0	0	0	4	33.3%
	2	P	8	6	6	5	1	0	2	0	75.0%
	3	P	8	7	7	5	0	2	0	1	87.5%
	4	M	11	9	9	3	6	0	0	2	81.8%
	5	P	4	4	3	1	2	0	0	1	75.0%
	6	P	4	4	4	2	0	2	0	0	100.0%
	7	P	4	3	3	1	0	2	1	0	75.0%
Total	\	\	45	34	34	19	9	6	3	8	75.6%
<b>P total</b>		\	28	23	23	14	3	6	3	2	<b>82.1%</b>
<b>M total</b>		\	17	11	11	5	6	0	3	6	<b>64.7%</b>
Autosomal Recessive Disorder	8	P	5	3	3	2	1	0	0	2	60.0%
		M	5	2	2	1	1	0	1	2	40.0%
	9	P	13	9	9	5	2	2	2	2	69.2%
		M	13	10	10	6	2	2	1	2	76.9%
	10	P	5	3	3	2	1	0	1	1	60.0%
		M	5	4	4	3	0	1	0	1	80.0%
	11	P	8	7	7	2	4	1	1	0	87.5%
		M	8	7	7	4	2	1	1	0	87.5%
	12	P	5	3	3	3	0	0	1	1	60.0%
		M	5	4	4	1	1	2	0	1	80.0%
Total alleles	\	\	72	52	52	29	14	9	8	12	72.2%
<b>P total</b>		\	36	25	25	14	8	3	5	6	<b>69.4%</b>
<b>M total</b>		\	36	27	27	15	6	6	3	6	<b>75.0%</b>
X-Recessive Disorder	13	M	6	5	5	5	0	0	1	0	83.3%
	14	M	3	3	3	3	0	0	0	0	100.0%
Total/ <b>M total</b>	\	\	9	8	8	8	0	0	1	0	<b>88.9%</b>
Total alleles			126	94	94	56	23	15	12	20	<b>74.6%</b>
P total			64	48	48	28	11	9	8	8	75.0%
M total			62	46	46	28	12	6	7	12	74.2%



940

941

## **ADDITIONAL INFORMATION**

### **Extended data**

943

#### **Extended Data Fig. 1 DNA size distribution in Spent Embryo Culture Medium (SCM)**

946 Size distributions of DNA molecules in six spent culture medium samples  
947 with different DNA fractions (Fragment Analyzer System). The horizontal  
948 axis represents the range of DNA fragments detected. 1bp is the lower  
949 marker, and 6000bp is the upper marker. The vertical axis represents a unit  
950 of relative fluorescence units. The first peak of insert fragments was  
951 approximately 180bp, with subsequent peaks spaced at 180bp intervals,  
952 corresponding to the length of DNA wrapped around a nucleosome. SCM  
953 2, SCM4, SCM5 and SCM6 showed varying levels of apoptosis,  
954 suggesting different degrees of DNA degradation. Additionally, SCM3  
955 failed to produce an effective library due to its low cfDNA content, and  
956 SCM1 did not exhibit the characteristic apoptosis signals. See  
957 Supplementary Notes S1 for more experimental details.

958

#### **Extended Data Fig. 2 Genome Coverage of Spent Embryo Culture Media**

961 Genome coverage of all SCM samples.

962

#### **Extended Data Fig. 3 Illustration of large fragment loss of haplotype (LFLoH) in one sample of SCM.**

965 Blue and orange bars denote haplotypes respectively from paternal and  
966 maternal chromosome 1. Black region represents no genetic information  
967 observed in specific haplotypes.

968

#### **Extended Data Fig. 4 Strategy for informative SNPs selection.**

970 0 and 1 respectively denote reference and alternative allele of variant. All  
971 the informative SNPs selected are listed.

972

973 **Extended Data Fig. 5 Estimated MCC rates.**

974 The top and bottom panels illustrate the MCC rates for autosomal dominant  
975 and autosomal recessive family samples, respectively. Circles represent  
976 pathogenic loci inherited paternally, while triangles denote maternal  
977 inheritance. Hollow symbols indicate the original MCC rates, and solid  
978 symbols represent the final MCC rates. Blue symbols correspond to  
979 diagnosable samples, whereas red symbols denote undiagnosable samples.  
980 See Supplementary Table S4 for the data.

981

982 **Extended Data Fig. 6 Improved protocol from LIANTI Method.**

983 In this study, we retained the LIANTI single-cell amplification method,  
984 which involves (a) : 1) Transposon-assisted DNA fragmentation; 2) Gap-  
985 filling and extension; 3) Reverse transcription; 4) Second-strand synthesis;  
986 5) Library construction and sequencing. To adapt this method for culture  
987 medium samples, we made the following improvements (b): 1) Adjusted  
988 the lysis system with a 10.8 µl sample volume and increased QIAGEN  
989 Protease (Qp) concentration to 35 mg/ml; 2) Doubled the transposon  
990 concentration; 3) Added a 19bp transposon ME sequence complementary  
991 primer to enhance reverse transcription efficiency; 4) Included 10  
992 additional PCR cycles in second-strand synthesis to improve the  
993 conversion of single-stranded cDNA to double-stranded DNA.

994

995 **Extended Data Fig. 7 Improved LIANTI amplification is more suitable  
996 for amplifying culture medium samples.**

997 Clinical collection of SCM (spent culture medium) yields limited volumes,  
998 and samples are precious. In this study, we utilized the BJ cell line to  
999 simulate the embryo culture process. a) Size distributions of DNA  
1000 molecules in BJ15 and BJ17 with different DNA fractions (Fragment  
1001 Analyzer System). b) Genome coverage analysis of BJ15 and BJ17, which  
1002 were amplified using Improved LIANTI, MALBAC, and MDA. c) SNP  
1003 number analysis of BJ15 and BJ17, which were amplified using Improved  
1004 LIANTI, MALBAC, and MDA. Genome coverage and SNP number

1005 normalization were carried out with Bulk sample. The sequencing depth  
1006 for bulk genomic DNA (gDNA) exceeds 30X, while for amplified products,  
1007 the sequencing depth exceeds 10X. BJ15: BJ cell lines medium collected  
1008 on Day5; BJ17: BJ cell lines medium collected on Day7 (see  
1009 Supplementary Table S2).

1010

1011 **Extended Data Fig. 8 The workflow of BASE-niPGT-M.**

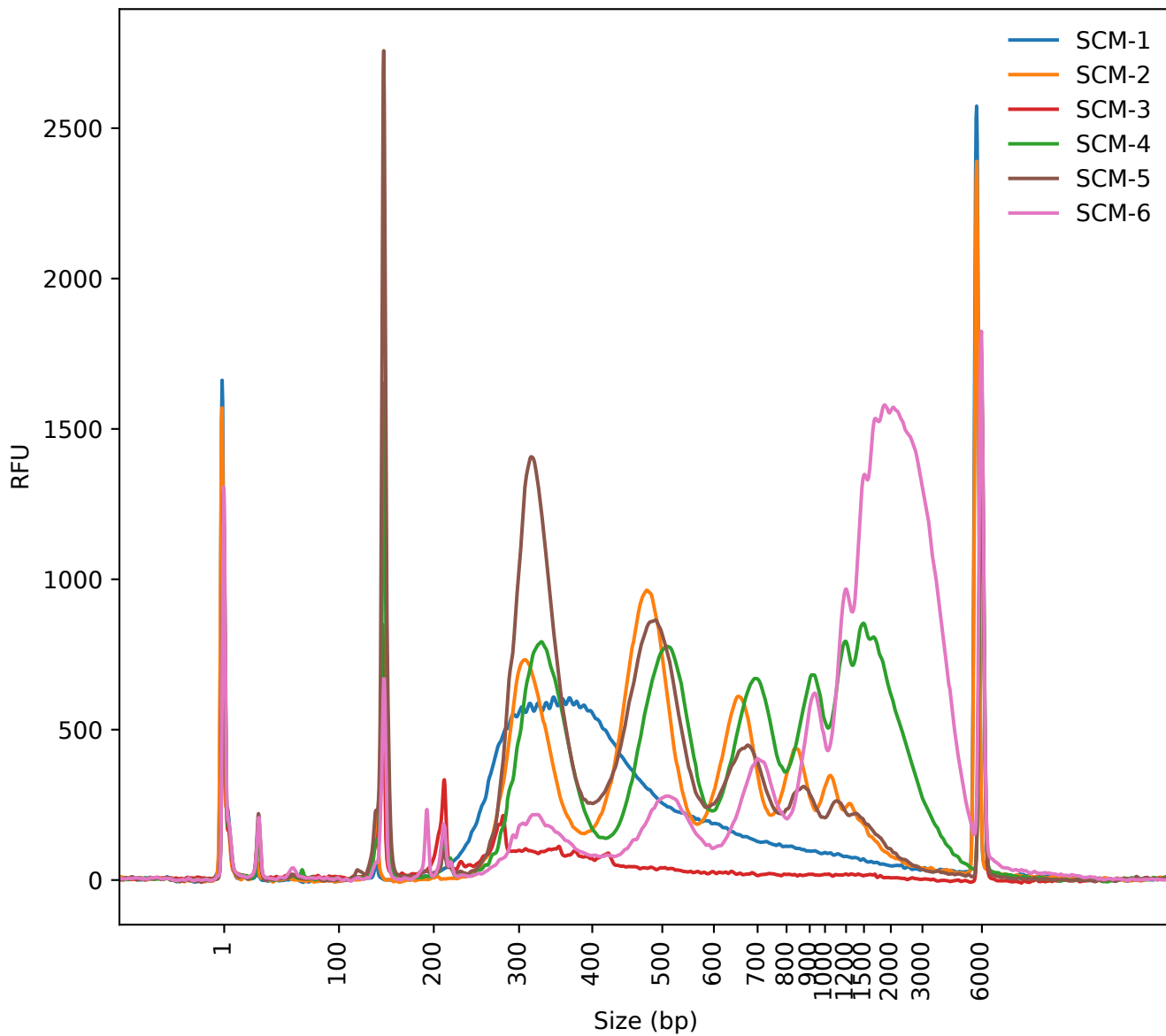
1012 The Bayesian model uses the sequencing data of SCM and the phased  
1013 haplotypes of the parents as input. By iteratively estimating the MCC rate  
1014 and haplotype states of each SNP, it calculates the log-likelihood ratio of  
1015 inheriting disease-carrying chromosomes versus disease-free  
1016 chromosomes as a function of distance to the disease-causing mutation site.  
1017 This curve, with its distinct characteristics, allows classification of the  
1018 results into four categories: Highly Confident, Moderately Confident,  
1019 Possibly confident, and Undetermined.

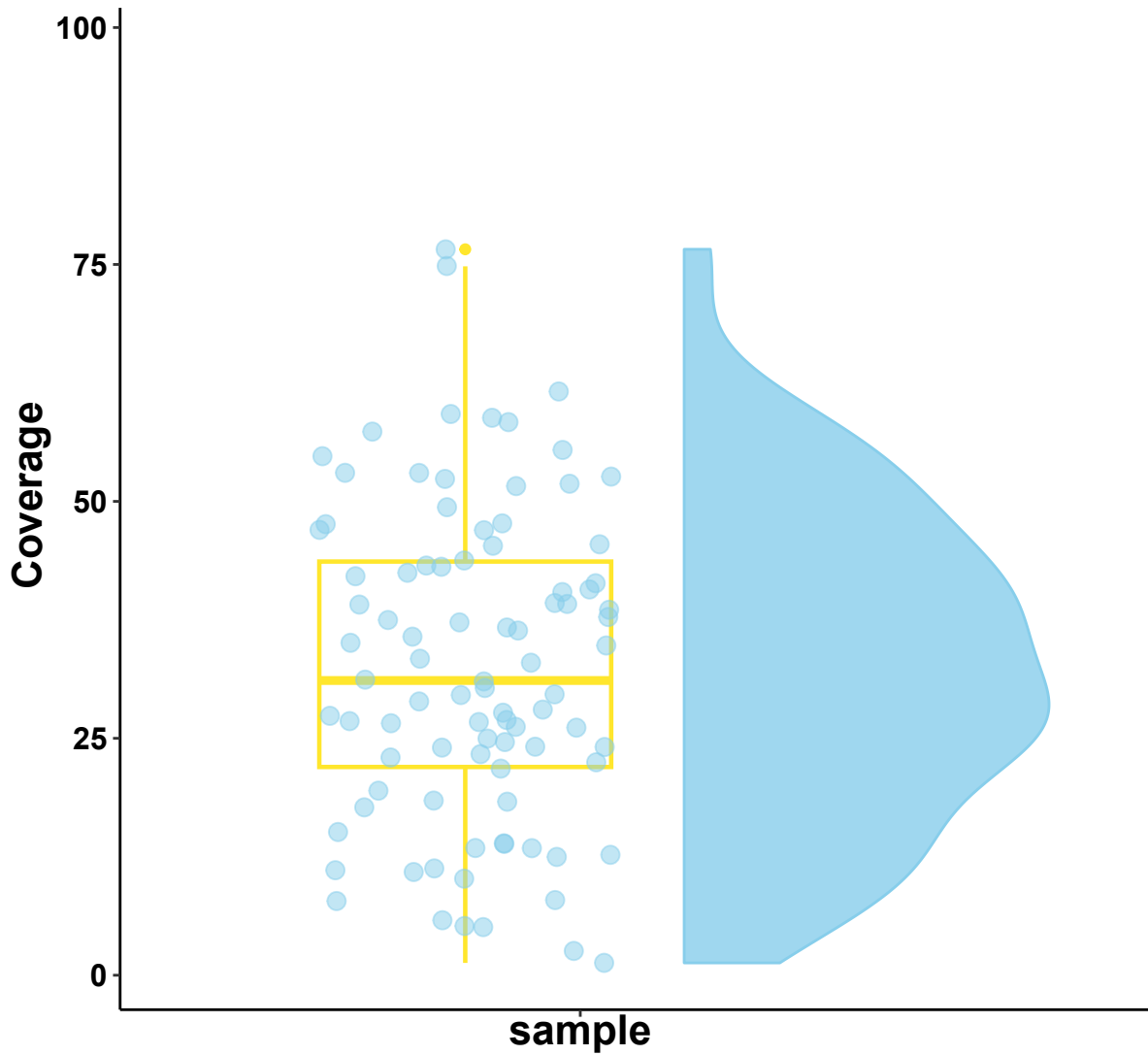
1020

1021 **Extended Data Fig. 9 An example illustrating the workflow of  
1022 pedigree-based whole genome reconstruction.**

1023 In the initial called genotype of the SCM, red numbers indicate the  
1024 informative alleles, which can identify the inherited allele from one parent.  
1025 Black regions represent the LFLoH pattern, characterized by consecutive  
1026 SNPs missing the same haplotype. Blue letters denote the imputed  
1027 haplotypes on non-informative SNPs, determined using the HMM method  
1028 with genetic distance as the transition probability. Green letters indicate  
1029 corrected haplotypes from the raw phasing results by the HMM method,  
1030 with incorrect ones shown in red. Green numbers show the imputed allele  
1031 or genotype of the genome.

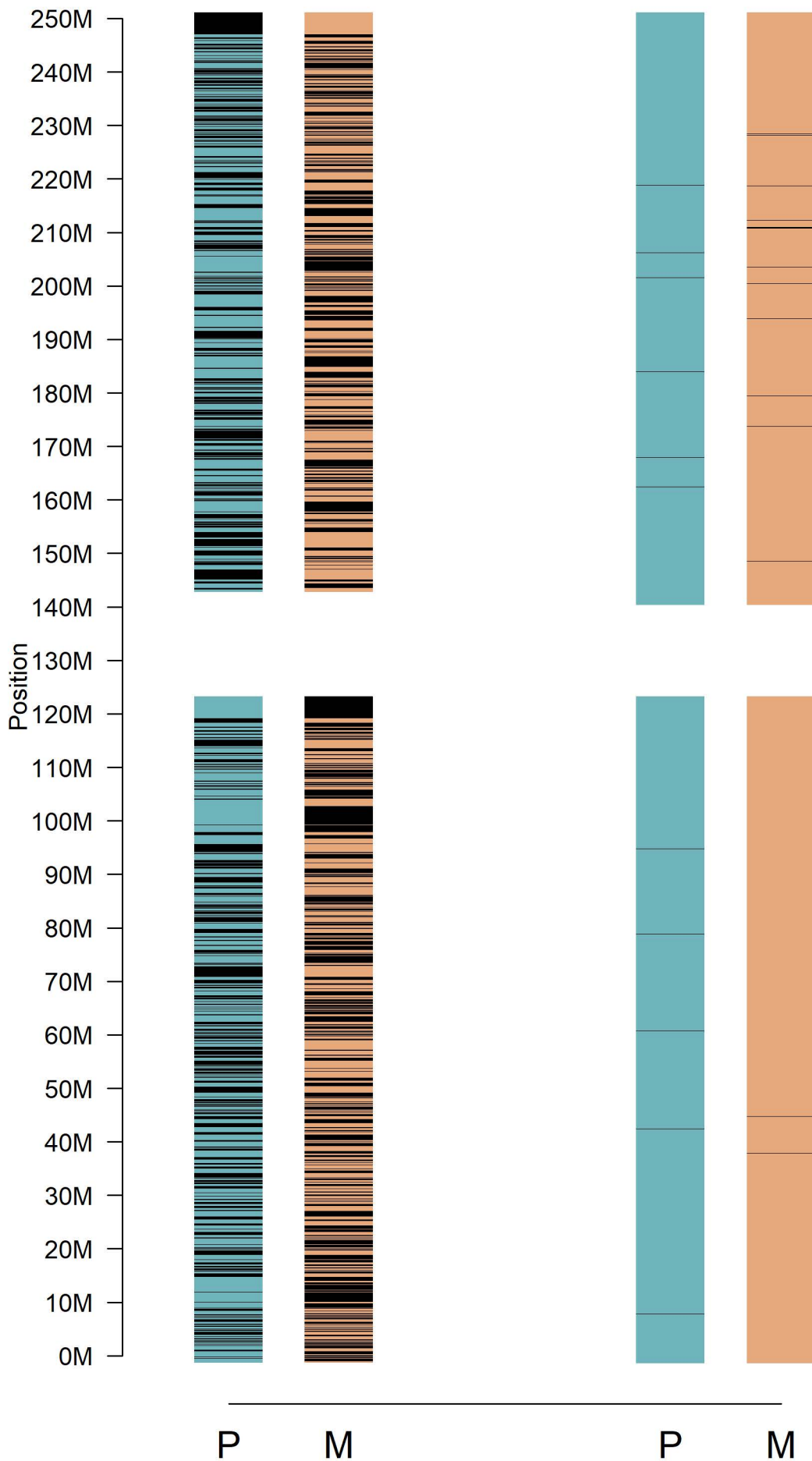
1032





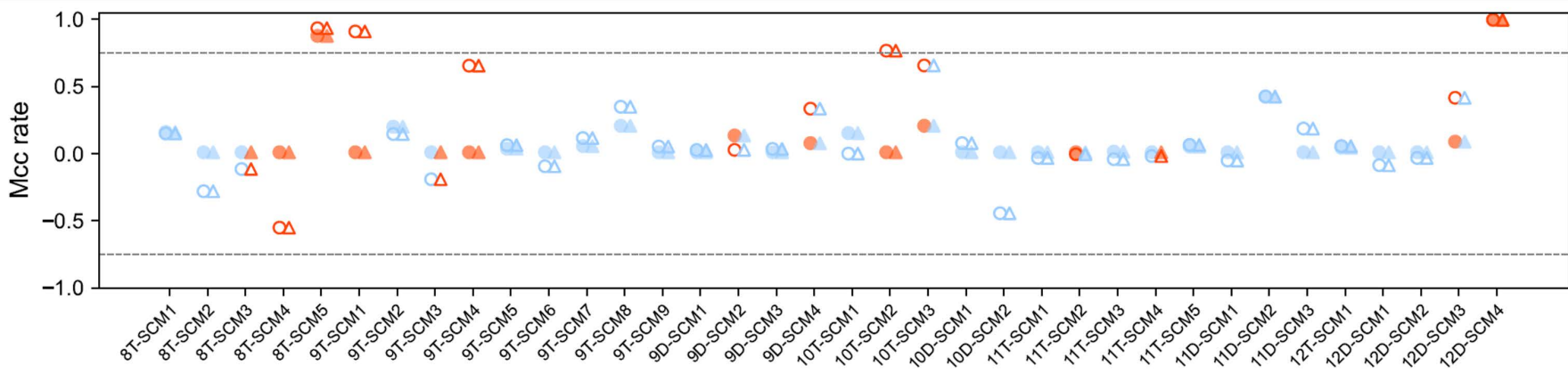
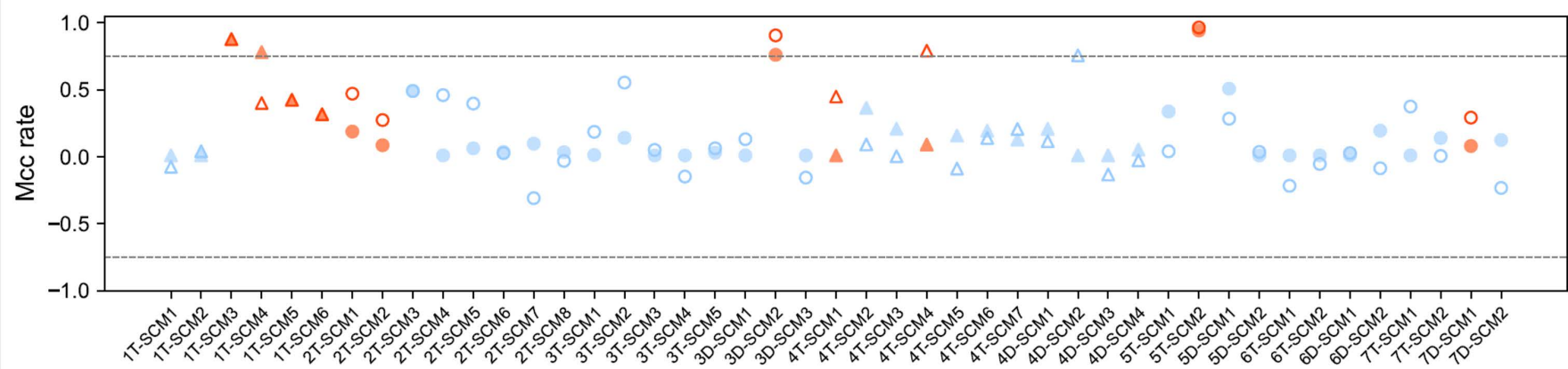
CM

TE



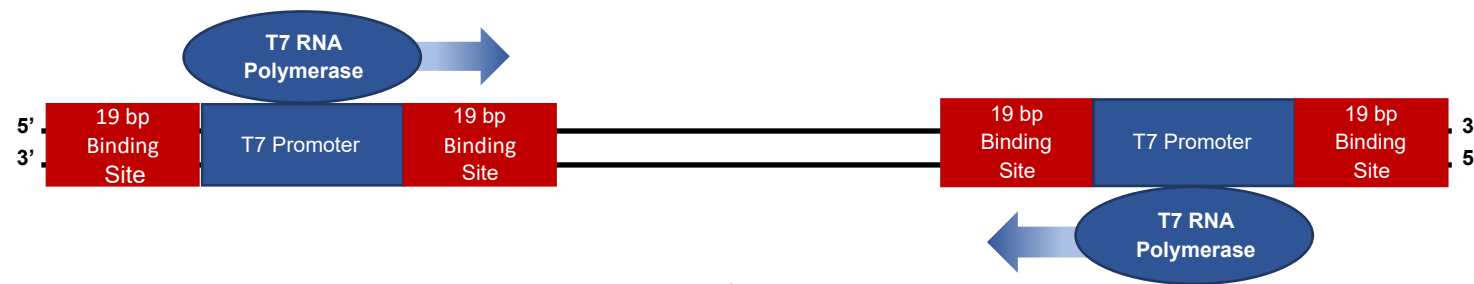


Paternal genotype	Maternal genotype	SCM genotype	Key allele	SNP type
0/1	0/0	0/1	1	Paternal informative SNP
0/1	0/0	1/1	1	Paternal informative SNP
0/1	1/1	0/1	0	Paternal informative SNP
0/1	1/1	0/0	0	Paternal informative SNP
0/0	0/1	0/1	1	Maternal informative SNP
0/0	0/1	1/1	1	Maternal informative SNP
1/1	0/1	0/1	0	Maternal informative SNP
1/1	0/1	0/0	0	Maternal informative SNP



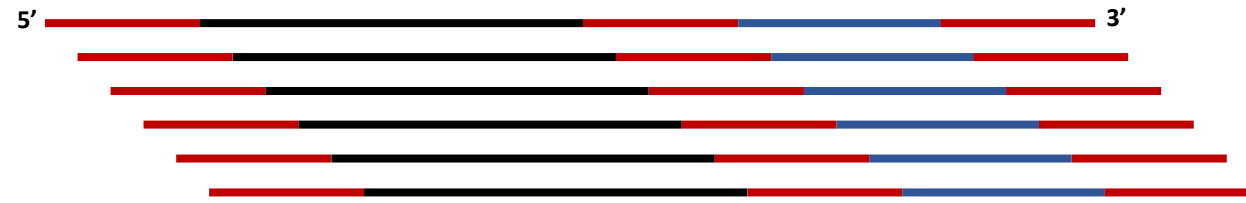
a

5'-CTGTCTCTTATACACATCTGAACAGAATTTAATACGACTCACTATAGGGAGATGTGTATAAGAGACAG-3' 5'-CTGTCTCTTATACACATCTGAACAGAATTTAATACGACTCACTATAGGGAGATGTGTATAAGAGACAG-3'



In vitro transcription amplification

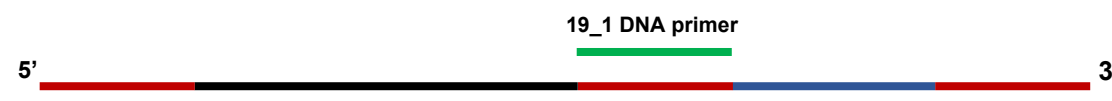
Linear Amplification into RNAs



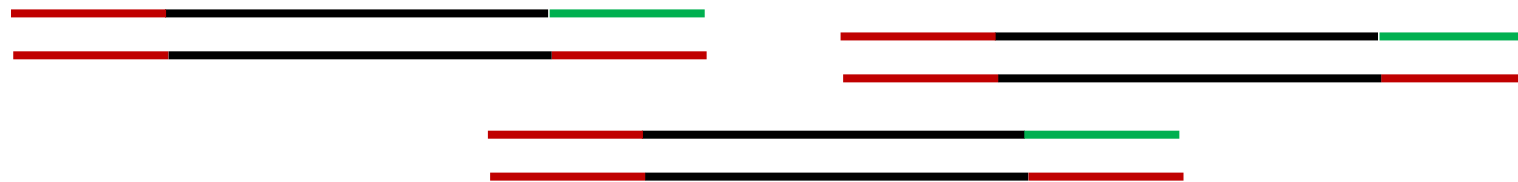
medRxiv preprint doi: <https://doi.org/10.1101/2024.11.10.24317049>; this version posted November 13, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity.

It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).

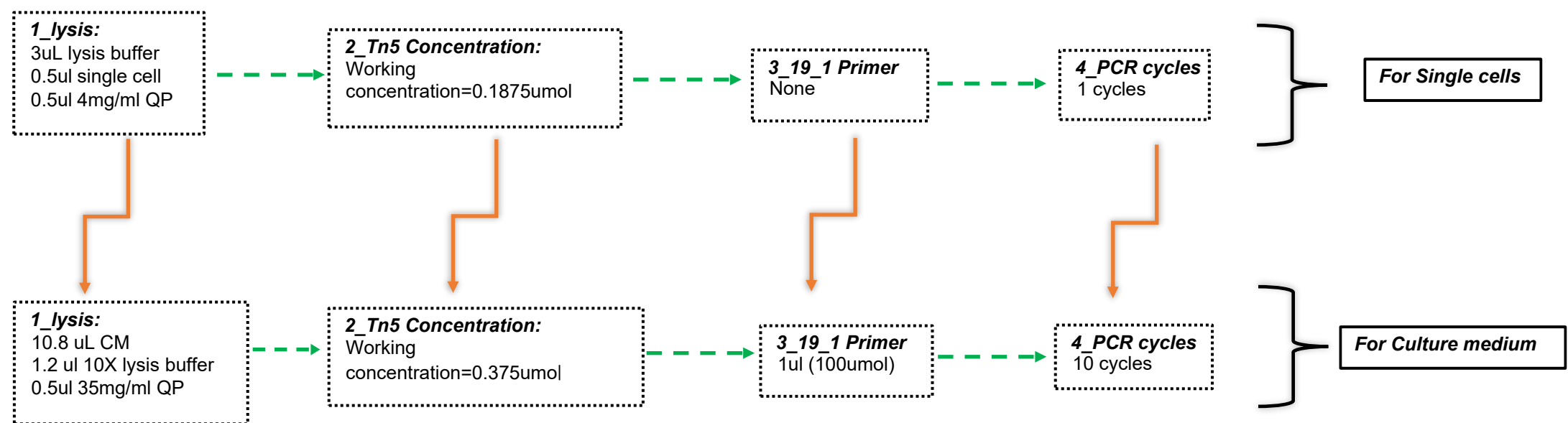
Add DNA primers

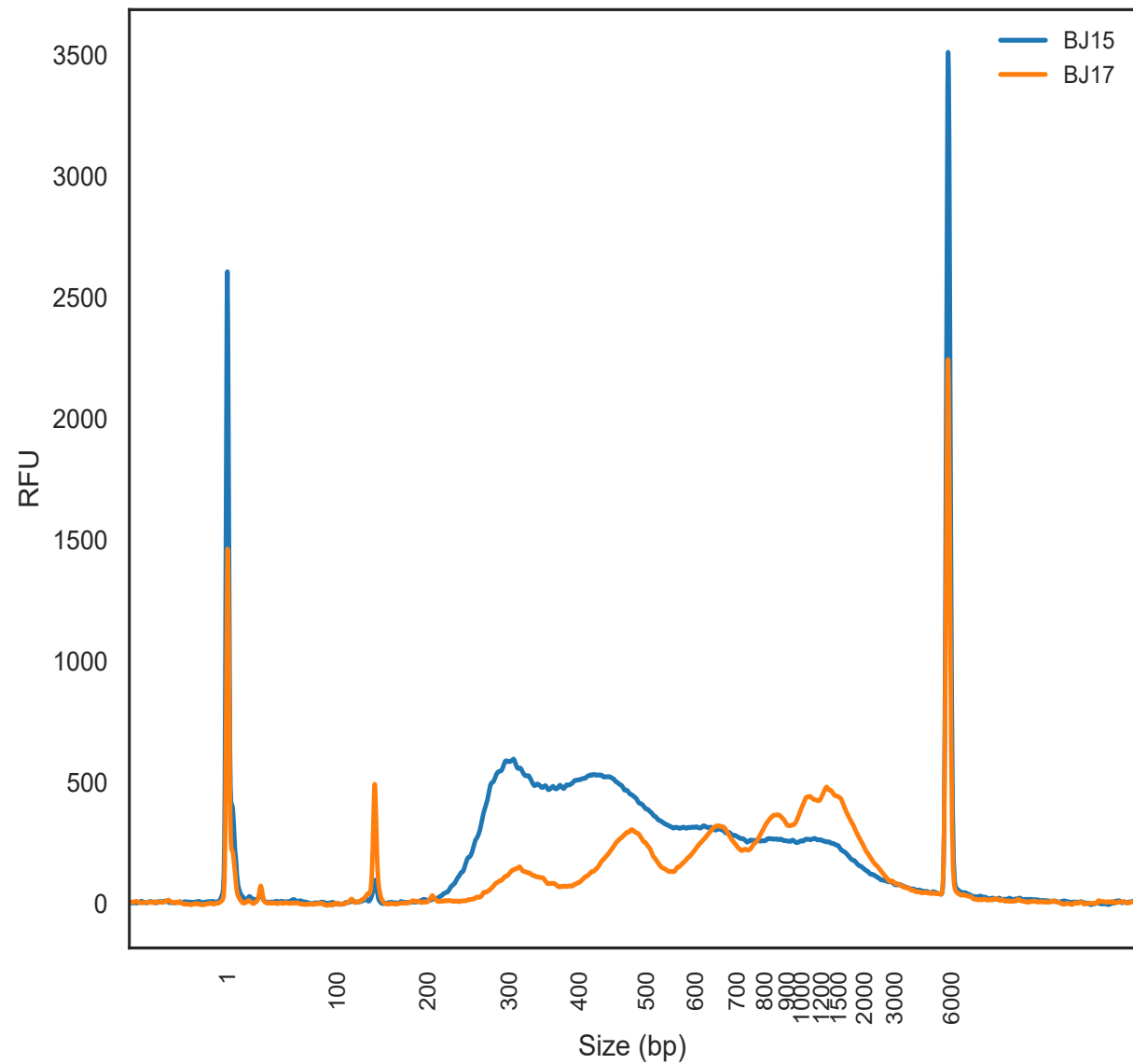
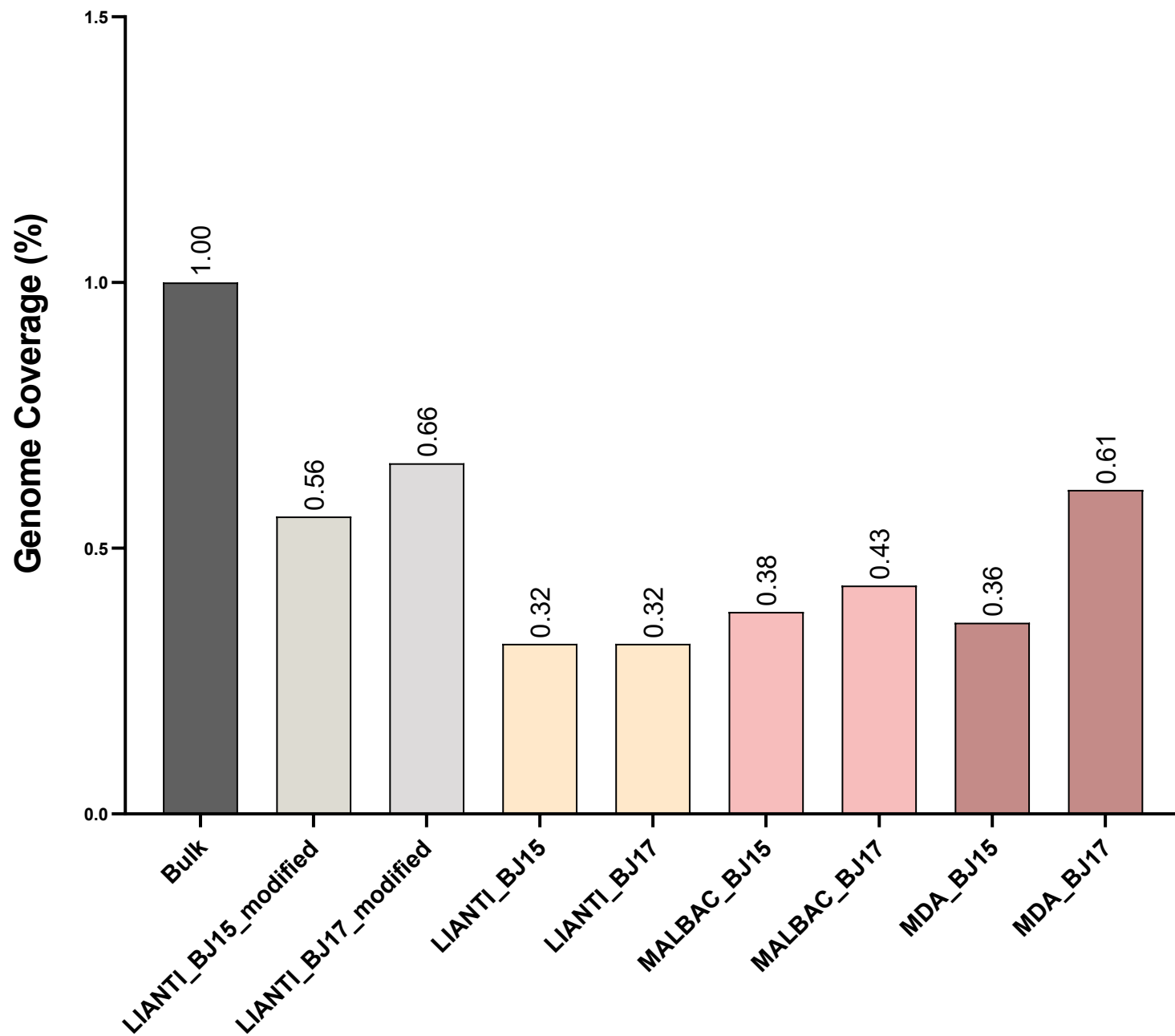
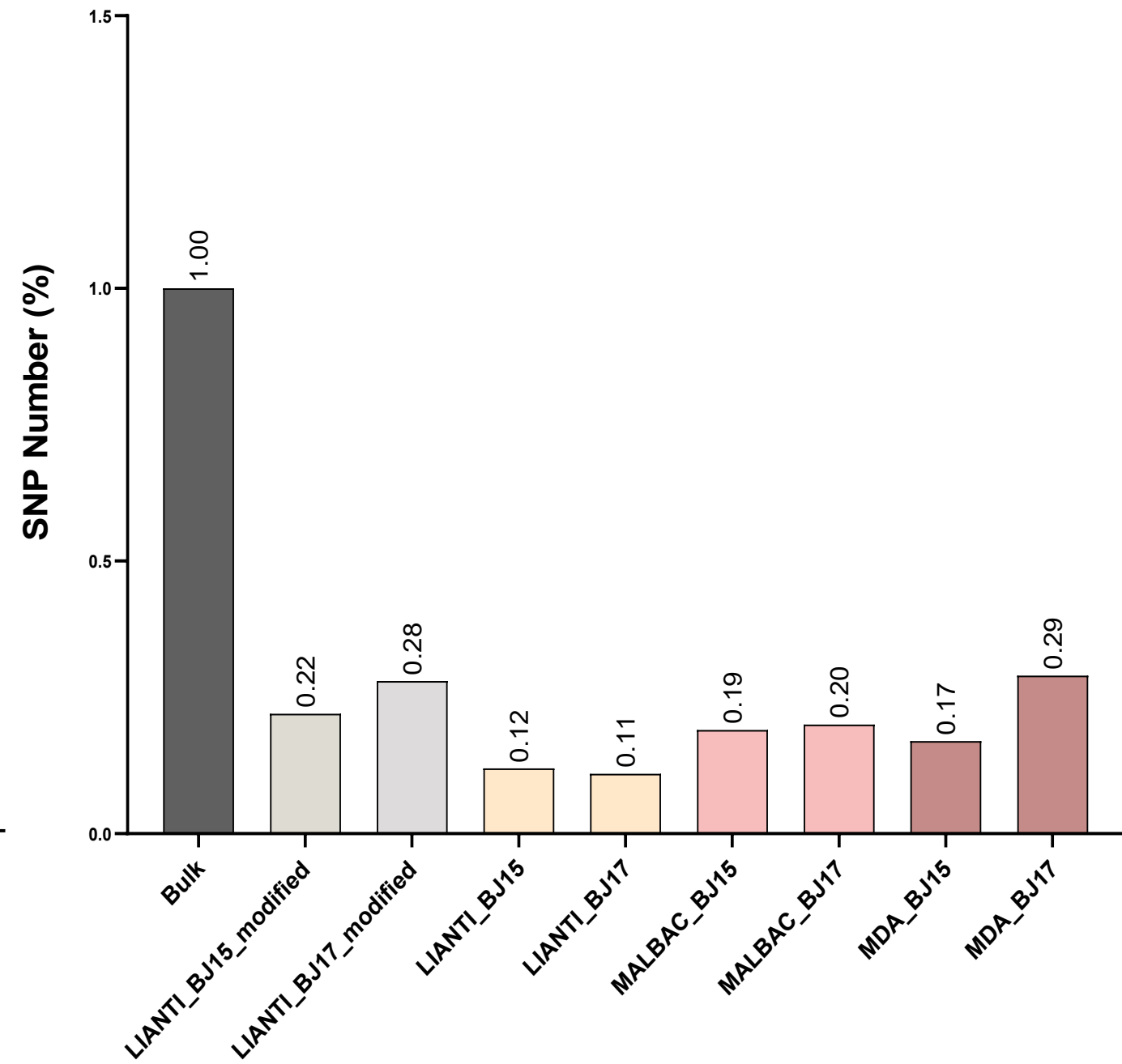


Reverse transcription



b



**a****b****c**

# BASEni-PGT-M

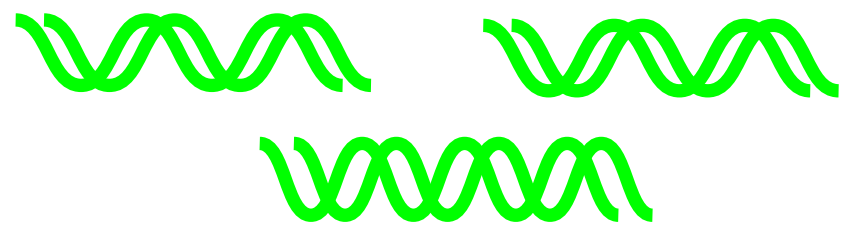
## DEFINITIONS

10MB (Million Base pairs) upstream and downstream of the disease-causing mutation site in the biological sample is defined as the adjacent region of the disease-causing mutation site.

Maternal Cell Contamination rate (MCC rate) is defined as the proportion of maternal DNA in the culture medium of an embryo compared to the total amount of DNA.

Large Fragment Loss of Haplotypes (LFLoH) is defined as the phenomenon that long and contiguous DNA fragments from one or both parents do not enter the culture medium, and the haplotype status of each SNP is defined to represent whether only DNA from one parent is present in the culture medium.

Amplify the DNA molecules from spent embryo culture medium, sequencing and performing mapping and SNP calling, filtering the SNPs of extremely low quality.



**INPUT**

Initially estimate the SER, the relative density of SNPs, the MCC rate and the haplotype status of each SNP.



Recursively estimate and update the MCC rate and the haplotype status of each SNP, using a conservative strategy with the single-SNP likelihood.



medRxiv preprint doi: <https://doi.org/10.1101/2024.11.10.24317049>; this version posted November 13, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](#).

Iterately calculate the likelihood of observing the SNP data within regions adjacent to the disease-causing mutation site, combining the likelihood at each single SNP and the recombination probabilities.

$$P(AD_{n:end} | F_n = i, M_n = j) = L_{ij}^{(n)} \sum_{k,l=1}^2 P(AD_{n+1:end} | F_{n+1} = k, M_{n+1} = l) \cdot P(F_{n+1} = k, M_{n+1} = l | F_n = i, M_n = j)$$

$AD$  is the measured Allele Depth data of SNPs.

Conduct haplotype pre-phasing, for the blood samples of both parents or solely the disease-carrying parent.

Incorporate the maternal cell contamination (MCC) rate and haplotype status into the likelihood of the four possible inheritance scenarios at each single SNP representing whether the genetic chains are inherited from the paternal or maternal alleles of the parents.

$$\hat{L}_{ij}^{(n)} = \frac{AD_n!}{AD_{0,n}! \cdot AD_{1,n}!} \left[ \hat{p}_{ij}^{(n)} \widehat{SER} + (1 - \hat{p}_{ij}^{(n)}) (1 - \widehat{SER}) \right]^{AD_{0,n}} \cdot \left[ \hat{p}_{ij}^{(n)} (1 - \widehat{SER}) + (1 - \hat{p}_{ij}^{(n)}) \widehat{SER} \right]^{AD_{1,n}}$$

$\widehat{SER}$  is the estimated sequencing error rate, and  $\hat{p}_{ij}^{(n)}$  is the theoretical probability that an ALT gene is obtained by randomly sampling one read at n-th SNP.

Recombination probability derived from database

$$P(F_{n+1} = k, M_{n+1} = l | F_n = i, M_n = j)$$

The linkage state of the n-th SNP:

$$F_n = 1(P1), F_n = 2(P2), M_n = 1(M1), M_n = 2(M2).$$

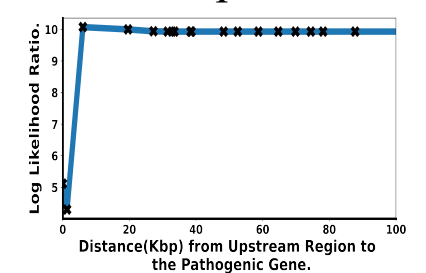
**INPUT**

Plot the log-likelihood-ratio(LLR) curves as a function of the physical distance to the disease-causing mutation site, using the SNPs within this region.

$$LLR = \log \left( \frac{\sum_{j=1}^2 \mathbb{P}(AD_{0:N} | F_0 = 1, M_0 = j)}{\sum_{j=1}^2 \mathbb{P}(AD_{0:N} | F_0 = 2, M_0 = j)} \right)$$

$$LLR = \log \left( \frac{\sum_{i=1}^2 \mathbb{P}(AD_{0:N} | F_0 = i, M_0 = 1)}{\sum_{i=1}^2 \mathbb{P}(AD_{0:N} | F_0 = i, M_0 = 2)} \right)$$

One example:



Determine whether the embryo carries the disease-causing chromosome, and providing a level of confidence.

**OUTPUT**

**DIAGNOSE:**  
P1 OR P2? M1 OR M2?

UNDETERMINED

HIGHLY confident

POSSIBLY confident

MODERATELY confident

Phasing informative SNPs of SCM						Building haplotype of SCM				Reconstructing genotype of SCM			
Paternal haplotype		Maternal haplotype		Initial genotype of SCM	SNP type	Paternal Haplotype of SCM	Maternal Haplotype of SCM	Genetic distance	Paternal Haplotype of SCM	Maternal Haplotype of SCM	Paternal allele of SCM	Maternal allele of SCM	Imputed genotype of SCM
P1	P2	M1	M2										
1	0	0	0	1/1	Haploid	P1		0.12	P1	M1	1	0	1 0
1	1	0	1	0/0	Haploid		M1	0.48	P1	M1	1	0	1 0
0	0	1	0	1/1	Haploid		M1	0.61	P1	M1	0	1	0 1
0	1	1	1	0/0	Haploid	P1	[REDACTED]	0.85	P1	M1	0	1	0 1
1	1	0	1	1/1	Uncertain			1.02	P1	M1	1	0	1 0
0	1	1	1	0/0	Haploid	P1		1.16	P1	M1	0	1	0 1
1	0	1	1	0/0	Haploid	P2		1.22	P1	M1	1	1	1 1
0	1	1	1	0/0	Haploid	P1		1.34	P1	M1	0	1	0 1
1	1	0	1	0/1	Diploid		M1	1.98	P1	M1	1	0	1 0
1	0	0	0	0/1	Diploid	P1		2.39	P1	M1	1	0	1 0
0	0	1	0	1/1	Haploid	[REDACTED]	M1	2.64	P1	M1	0	1	0 1
1	1	1	0	0/0	Haploid		M2	2.73	P1	M1	1	1	1 1
0	0	1	0	1/1	Haploid		M1	2.85	P1	M1	0	1	0 1
0	0	1	0	1/1	Haploid				2.96	P1	M1	0	1
1	0	0	0	0/1	Diploid	P1		3.15	P1	M1	1	0	1 0
0	0	1	0	0/1	Diploid		M1	3.21	P1	M1	0	1	0 1
0	0	0	1	0/1	Diploid		M2	4.37	P1	M2	0	1	0 1
1	0	0	0	0/1	Diploid	P1		4.93	P1	M2	1	0	1 0
0	1	1	0	0/1	Diploid			5.30	?	M2	?	0	./.
1	0	1	1	0/1	Diploid	P2		5.67	P2	M2	0	1	0 1
0	1	0	0	0/1	Diploid	P2		5.69	P2	M2	1	0	1 0
0	0	0	1	1/1	Haploid		M2	5.74	P2	M2	0	1	0 1
0	1	0	0	1/1	Haploid	P2		5.88	P2	M2	1	0	1 0



1033

## **Supplementary Information**

1034

1035 **Supplementary Notes**

1036

1037 **Supplementary Tables**

1038