

## 1 **KaryoTap Enables Aneuploidy Detection in Thousands of Single Human Cells**

2 Joseph C. Mays; josephcmays@gmail.com; Institute for Systems Genetics and Department of  
3 Biochemistry and Molecular Pharmacology, New York University Grossman School of Medicine,  
4 New York, NY 10016, USA

5  
6 Sally Mei; Sally.Mei2@nyulangone.org; Institute for Systems Genetics and Department of  
7 Biochemistry and Molecular Pharmacology, New York University Grossman School of Medicine,  
8 New York, NY 10016, USA

9  
10 Manjunatha Kogenaru; Manjunatha.Kogenaru@nyulangone.org; Institute for Systems Genetics  
11 and Department of Biochemistry and Molecular Pharmacology, New York University Grossman  
12 School of Medicine, New York, NY 10016, USA

13  
14 Helberth M. Quysbertf; Helberth.Quisberth@nyulangone.org; Institute for Systems Genetics and  
15 Department of Biochemistry and Molecular Pharmacology, New York University Grossman  
16 School of Medicine, New York, NY 10016, USA

17  
18 Nazario Bosco; nazario.bosco.phd@gmail.com; Institute for Systems Genetics and Department  
19 of Biochemistry and Molecular Pharmacology, New York University Grossman School of  
20 Medicine, New York, NY 10016, USA. Current Address: Volastra Therapeutics, New York, NY  
21 10027, USA.

22  
23 Xin Zhao; Institute for Systems Genetics and Department of Biochemistry and Molecular  
24 Pharmacology, New York University Grossman School of Medicine, New York, NY 10016, USA

25  
26 Joy J. Bianchi; Joy.Bianchi@nyulangone.org; Institute for Systems Genetics and Department of  
27 Biochemistry and Molecular Pharmacology, New York University Grossman School of Medicine,  
28 New York, NY 10016, USA

29  
30 Aleah Goldberg; aleah.goldberg@nyulangone.org; Institute for Systems Genetics and  
31 Department of Biochemistry and Molecular Pharmacology, New York University Grossman  
32 School of Medicine, New York, NY 10016, USA

33  
34 Gururaj Rao Kidiyoor; Gururajrao.Kidiyoor@nyulangone.org; Institute for Systems Genetics and  
35 Department of Biochemistry and Molecular Pharmacology, New York University Grossman  
36 School of Medicine, New York, NY 10016, USA

37  
38 Liam J. Holt; liam.holt@nyulangone.org; Institute for Systems Genetics and Department of  
39 Biochemistry and Molecular Pharmacology, New York University Grossman School of Medicine,  
40 New York, NY 10016, USA

41  
42 David Fenyő; david.fenyo@nyulangone.org; Institute for Systems Genetics and Department of  
43 Biochemistry and Molecular Pharmacology, New York University Grossman School of Medicine,  
44 New York, NY 10016, USA

45  
46 Teresa Davoli (corresponding author); teresa.davoli@nyulangone.org; Institute for Systems  
47 Genetics and Department of Biochemistry and Molecular Pharmacology, New York University  
48 Grossman School of Medicine, New York, NY 10016, USA

49

50 **Abstract**

51 Investigating chromosomal instability and aneuploidy within tumors is essential for understanding  
52 tumorigenesis and developing diagnostic and therapeutic strategies. Single-cell DNA sequencing  
53 technologies have enabled such analyses, revealing aneuploidies specific to individual cells within  
54 the same tumor. However, it has been difficult to scale the throughput of these methods to detect  
55 rare aneuploidies while maintaining high sensitivity. To overcome this deficit, we developed  
56 KaryoTap, a method combining custom targeted DNA sequencing panels for the Tapestry platform  
57 with a computational framework to enable detection of chromosome- and chromosome arm-scale  
58 aneuploidy (gains or losses) and copy number neutral loss of heterozygosity in all human  
59 chromosomes across thousands of single cells simultaneously. KaryoTap allows detecting gains  
60 and losses with an average accuracy of 83% for arm events and 91% for chromosome events.  
61 Importantly, together with chromosomal copy number, our system allows us to detect barcodes  
62 and gRNAs integrated into the cells' genome, thus enabling pooled CRISPR- or ORF-based  
63 functional screens in single cells. As a proof of principle, we performed a small screen to expand  
64 the chromosomes that can be targeted by our recently described CRISPR-based KaryoCreate  
65 system for engineering aneuploidy in human cells. KaryoTap will prove a powerful and flexible  
66 approach for the study of aneuploidy and chromosomal instability in both tumors and normal  
67 tissues.

68

69 **Keywords**

70 Aneuploidy, Chromosomal Instability, Single-Cell DNA Sequencing, Targeted DNA Sequencing,  
71 Cancer genomics, Tumor heterogeneity, Copy Number Variants, CNVs

72

73

## 74 **Background**

75 A critical hallmark of cancer initiation and progression is the presence of aneuploidy, or gains and  
76 losses of whole chromosomes or chromosome arms, which arise due to mitotic missegregation  
77 events(1–3). In addition to aneuploidy, chromosomal instability, characterized by a continuously  
78 high rate of these missegregation events among tumor cells, has gained particular interest as a  
79 potential driver of tumor progression and metastasis(4,5). Chromosomal instability produces cell  
80 populations with heterogeneous aneuploid karyotypes that continuously evolve over time,  
81 granting tumor cells an opportunity to adapt to their environment and develop resistance to cancer  
82 therapies(6–8). Traditional methods for detecting aneuploidy, such as whole genome sequencing  
83 (WGS), rely on bulk averaging of cell populations, effectively masking heterogeneity among  
84 individual cells and preventing proper assessment of the extent or rate of chromosomal instability.  
85 Single-cell approaches, particularly single-cell DNA sequencing (scDNA-seq), overcome this  
86 limitation and can instead detect the full complement of distinct karyotypes present in a tumor(9–  
87 13). These data can be used to reassemble the evolution of the tumor's cells, which can provide  
88 insights into how aneuploidy and chromosomal instability may drive tumorigenesis or inform  
89 treatments for therapeutic resistance(8,11). Methods for modeling aneuploidies of specific  
90 chromosomes in cell culture such as KaryoCreate, have also emerged as powerful tools for  
91 studying the effects of aneuploidy in cancer(14–16). These methods have benefitted from scDNA-  
92 seq, as sequencing individual cells enables the evaluation of the specificity and accuracy of the  
93 engineered karyotypes(15,16).

94  
95 The fundamental challenge of scDNA-seq methods is faithfully and completely sequencing the 6  
96 picograms of DNA within a cell. The uniformity and depth of sequencing coverage determine the  
97 sensitivity with which aneuploidy can be detected. Conventional whole-genome amplification  
98 (WGA) methods, DOP-PCR(17,18), MDA(19), and MALBAC(20) amplify the genome prior to  
99 sequencing, introducing amplification biases and PCR errors that confound results. Finally, these

100 methods rely on the partitioning of cells into individual wells or tubes, which limits throughput,  
101 preventing the cost-efficient sequencing of enough cells to identify rare aneuploidy events in a  
102 large population(21). More recent methods vastly improve throughput by using microfluidic  
103 partitioning(22), combinatorial indexing(23), and liquid handling robots(24), allowing hundreds of  
104 thousands of cells to be sequenced at once. However, these methods suffer from uneven  
105 amplification across the genome and require sufficiently deep sequencing per cell, jeopardizing  
106 the confident detection of aneuploidy in individual cells(25). Furthermore, the need for custom  
107 hardware, expensive liquid handlers, and complicated protocols makes these methods difficult to  
108 adopt for most laboratories(21). We note that in scDNA-seq methods based on (untargeted) WGS  
109 there is a natural bias in the sensitivity of detecting whole chromosome or arm-level gains or  
110 losses across chromosomes depending on their size. In fact, the sensitivity of aneuploidy  
111 detection depends on the total number of reads per chromosome thus is lower for smaller  
112 chromosomes compared to larger ones as they contain a smaller proportion of the total reads.  
113 However, for the purpose of evaluating aneuploidy (gains and losses of whole chromosomes or  
114 chromosome arms) and chromosomal instability (rate of chromosome missegregation), each  
115 chromosome counts as a single entity. Thus, scDNA-seq methods based on WGS necessitate a  
116 high number of reads (and thus increase in cost) to achieve sufficient sensitivity of detection  
117 across all chromosomes.

118  
119 To address the need for a high-throughput method for detecting chromosome-scale aneuploidy  
120 across the human genome that maintains high sensitivity at a cost-efficient sequencing depth, we  
121 turned to the Tapestry platform from Mission Bio, a droplet-based targeted scDNA-seq solution  
122 that allows for the sequencing of hundreds of genomic loci across thousands of cells in one  
123 experiment(26). The platform is commonly used to detect tumor hotspot mutations in cancer driver  
124 genes and has not yet been utilized to identify aneuploidy (gains and losses of whole  
125 chromosomes or chromosome arms) across the genome. While targeted sequencing is typically

126 used to identify mutations in single nucleotides(27), we reasoned that we could use the relative  
127 sequencing depth of targeted loci to detect chromosome- and chromosome arm-scale aneuploidy  
128 in individual cells. Here, we describe KaryoTap, a method combining custom targeted Tapestri  
129 panels of PCR probes covering all human chromosomes with a Gaussian mixture model  
130 framework for DNA copy number detection, thus enabling the accurate detection of aneuploidy in  
131 all chromosomes in several thousand cells across different cell lines. Additionally, we included  
132 probes that detect lentiviral-integrated CRISPR guide RNAs to enable functional studies, and  
133 DNA barcodes to enable sample multiplexing. To enhance usability, we also developed a  
134 companion software package for R, *karyotapR*, which enables the straightforward copy number  
135 analysis, visualization, and exploration of the data produced by our custom panels. KaryoTap  
136 allows detecting gains and losses with an average accuracy of 83% for arm events and 91% for  
137 chromosome events. By overcoming the limitations of current methods, this system will be a  
138 valuable tool for investigating the evolution and consequences of aneuploidy and chromosomal  
139 instability in human tumors, in addition to other healthy and diseased tissues, such as normal  
140 tissues during physiological aging or clonal hematopoiesis of indeterminate potential.

141

## 142 **Results**

143

### 144 ***Design of a Custom Targeted Panel for Detecting Chromosome-Scale Aneuploidy***

145 To detect DNA copy number across the human genome in single cells, we designed a custom  
146 panel (Version 1; CO261) for the Tapestri system comprising 330 PCR probes that target and  
147 amplify specific loci across all 22 autosomes and the X chromosome (**Fig 1; Table S1; Additional**  
148 **File 1**; for the Version 3 panel, the Y chromosome was also included; see below). The number of  
149 probes targeting each chromosome was proportional to the size of the chromosome (e.g., 24  
150 probes for chr1, 5 for chr22) to achieve a roughly uniform density of ~1-2 probes per 100  
151 megabases across all chromosomes. Loci were selected to cover regions carrying single

152 nucleotide polymorphisms (SNPs) with major allele frequencies of 0.5-0.6 such that cells from  
153 different lines or individuals sequenced in the same experiment could be identified by their distinct  
154 genotypes (see Methods).

155

156

### 157 ***Population-Level Aneuploidy Detection Across Human Cell Lines***

158 The Tapestri system partitions individual cells into aqueous droplets and generates barcoded  
159 amplicons from the loci targeted by the supplied probe panel. Both cell-specific sequencing read  
160 counts and variant allele frequencies (VAFs) for each probe are generated by processing  
161 sequencing reads from these amplicons. To test whether our Version 1 panel could detect whole-  
162 chromosome aneuploidy in individual cells, we performed a Tapestri scDNA-seq experiment on  
163 a pool of five cell lines with varying karyotypes, mixed in equal proportion. The pool consisted of  
164 retinal pigment epithelial cells hTERT RPE-1 (hereafter RPE1; +10q, XX) as a near-diploid  
165 reference population(6,28), and four aneuploid colon cancer cell lines: LS513 (+1q, +5, +7, +9,  
166 +13, +13, XY), SW48 (+7, +14, XX), LoVo (+5p, +7, +12, XY), and CL11 (-6, +7, +7, -17, -18, -  
167 22, XY). The bulk (i.e., population-averaged) karyotypes for each cell line were determined by  
168 whole genome sequencing (WGS) (**Fig S1A**) and confirmed for RPE1 by G-banded karyotyping  
169 (**Fig S1B**).

170

171 Sequencing read counts and variant allele frequencies (VAFs) for 2,986 cells across the five cell  
172 lines were recovered from the Tapestri experiment. Dimensional reduction of the cells' VAFs by  
173 principal component analysis (PCA) and UMAP revealed 5 major clusters corresponding to the 5  
174 cell lines (**Fig S1C**). The remaining smaller clusters, representing composites of VAF profiles from  
175 multiple cells captured in the same droplet, were discarded from further analyses. We estimate  
176 the copy number for each probe in each cell as the ratio of read counts relative to a reference  
177 population. Here, we used the RPE1 cells as the reference population, which have 2 copies

178 (diploid) of each chromosome, except for a third copy (triploid) of the chr10q arm translocated to  
179 the X chromosome (**Fig S1B**). To identify the cell cluster comprising RPE1 in our data, we  
180 calculated the mean VAF across variants for each cluster and compared them to VAFs from  
181 published deep WGS of RPE1(29) by PCA (**Fig S1D**). The cells whose mean VAFs clustered  
182 closest to those of published RPE1 represent RPE1 cells. The copy number estimates for each  
183 probe in each cell for all cells (hereafter, cell-probe scores) were then calculated by taking the  
184 ratio of normalized read counts to the median normalized read counts of the RPE1 reference  
185 population (see Methods).

186  
187 The 330 target regions have varying base compositions (**Additional File 1**) and the probes have  
188 different optimal melting temperatures but are amplified under the same conditions and  
189 thermocycling parameters, introducing technical artifacts from amplification bias(30). The probe-  
190 level heatmap (**Fig S1E**) highlights such technical variation between copy number values from  
191 intra-chromosomal probes, suggesting that measurements from any individual probe are unlikely  
192 to reliably reflect a cell's copy number. To address this, we calculated a single copy number score  
193 for each chromosome in each cell (hereafter, cell-chromosome unit) by smoothing the cell-probe  
194 scores of all probes targeting the same chromosome. Smoothing was accomplished by  
195 calculating the weighted median of the cell-probe copy number values for all probes on a given  
196 chromosome; larger weights were assigned to probes whose copy number scores had smaller  
197 spreads (see Methods). Heatmap visualization of cell-chromosome copy number scores for the  
198 four colon cancer cell lines corresponds with the expected population-level copy number values  
199 from bulk WGS (**Fig 2A**). For example, the per-cell and average heatmap intensities indicate  
200 correctly that LS513, SW48, and LoVo carry 3 copies of chromosome 7, while CL11 carries 4  
201 copies. Similarly, 1 copy of chromosome 6 could be detected in CL11, indicating a chromosomal  
202 loss, and single copies of chromosome X could be detected in LS513, LoVo, and CL11, indicating  
203 XY sex chromosomes.

204

205 scDNAseq is often used on tumor cells to characterize the copy number heterogeneity that arises  
206 from chromosomal instability during tumorigenesis (intratumoral heterogeneity; ITH)(31). Here,  
207 we can consider the entire dataset a model of a heterogenous tumor carrying several major  
208 subclonal lineages, with each cell line representing a distinct major subclone. Unbiased clustering  
209 correctly groups the subclones/cell lines by copy number score into their respective cell lines (**FIG**  
210 **S2A**). In a real tumor, this could be used to distinguish subclones. Furthermore, clustering cells  
211 from each line can reveal subclones occurring within a given line. Here, we show that clustering  
212 of LoVo cells reveals two subpopulations hallmarked by exclusive gains in chr5p or chr15q (**FIG**  
213 **S2B**). Altogether, these data demonstrate that our Version 1 panel can resolve the average copy  
214 number of populations of cells at the whole-chromosome level and distinguish major subclones  
215 within a heterogenous cell population using copy number.

216

### 217 ***Copy Number Estimation in Individual Cells***

218 To account for variation in the distributions of copy number scores, we classified the scores as  
219 integer copy number calls using a 5-component Gaussian mixture model (GMM)(32), where each  
220 component represents a possible copy number value of 1, 2, 3, 4 or 5 (**Fig 2B**). Using the known  
221 copy number for each chromosome in RPE1 and the corresponding distributions of copy number  
222 scores, we simulated the expected distributions of copy number scores that would be measured  
223 from chromosomes with actual copy numbers 1-5 for each chromosome. We then used Bayes  
224 theorem to calculate the posterior probability of each cell-chromosome score belonging to each  
225 of the 5 GMM component distributions and assign the cell-chromosome an integer copy number  
226 corresponding to the component with the highest posterior probability (**Fig S2C**) (see Methods for  
227 details).

228



229 To evaluate copy number calling performance, we determined the accuracy of the classifier  
230 model, calculated as the proportion of correct calls (i.e., true positives), using the known copy  
231 numbers for RPE1 as ground truth. We focused on RPE1 as its karyotype is stable and  
232 homogeneous across the population(6); karyotyping confirmed that the line is triploid for chr10q  
233 and diploid for all other chromosomes in all metaphases analyzed, with the exception of 3 copies  
234 of chr12 in 3% of metaphases (**Fig S1B**). Accuracy, or correctly identifying 2 copies of a  
235 chromosome in RPE1, varied between chromosomes and ranged from 95% for chr2 to 49% for  
236 chr22 with a mean of 82% (**Fig 2C**). Because we used the RPE1 cells to both fit and test the  
237 GMMs, we performed 5-fold cross validation by partitioning them randomly into 5 equally sized  
238 subsets and calculated accuracy five times, each time reserving one of the sets from the model  
239 generation and using it only to calculate accuracy. The mean absolute deviation of the 5 accuracy  
240 measurements for each chromosome ranged from 0.61 to 4.59 percentage points, suggesting  
241 that copy number calling performance would be maintained when classifying new data. 18 out of  
242 22 chromosomes had sensitivities of at least 75%; chr10 was excluded from the whole-  
243 chromosome analysis because the p and q arms have different copy numbers. Accuracy for each  
244 chromosome correlated strongly with chromosome length (Pearson  $r = 0.73$ ), and chromosome  
245 length itself correlated strongly with the number of probes targeting the chromosome (Pearson  $r$   
246  $= 0.93$ ) (**Fig S2D**), suggesting that classifier accuracy is related to the number of probes targeting  
247 a chromosome. As expected, linear regression of accuracy on the number of probes per  
248 chromosome indicated that the number of probes is predictive of copy number call accuracy ( $R^2$   
249  $= 0.81$ ;  $p = 1.12e-08$ ) (**Fig 2D**). This suggests that the classifier accuracy for poorly performing  
250 chromosomes, particularly the smaller chromosomes including 19, 21, and 22, could be improved  
251 by adding additional probes for those chromosomes to the panel.

252

253 The empirical accuracy for copy number calls in RPE1 only demonstrates the ability of our method  
254 to detect 2 copies of a chromosome. We can determine the theoretical or expected sensitivity for

255 detecting copy numbers of 1, 2, 3, 4, and 5 for each chromosome by calculating the proportion of  
256 each GMM component distribution that would be called correctly as belonging to that component.  
257 Overall, theoretical sensitivity was highest for 1 copy with an average of 97%, decreasing with  
258 each additional copy; 2 copies had an average sensitivity of 83% and 3 copies had 64% (**Table**  
259 **S2**). The theoretical sensitivity for 2 copies strongly correlated with the empirical accuracy for  
260 RPE1 calls (Pearson  $r = 0.97$ ). As expected, theoretical sensitivity at all 5 copy number levels  
261 decreased for chromosomes with fewer probes, as was the case with the empirical accuracy (**Fig**  
262 **2E**).

263  
264 WGS indicated aneuploidy restricted to one arm of a chromosome in RPE1 (chr10q), LS513  
265 (chr1p), and LoVo (chr5p) (**Fig S1A**). To determine if our Version 1 panel could also detect  
266 chromosome arm-level aneuploidy as well, we performed a similar analysis by smoothing cell-  
267 probe copy number scores across probes targeting each chromosome arm instead of across  
268 whole chromosomes (**Fig 3A**). Per-cell and average heatmap intensities indicate correctly that  
269 LS513 carries 3 copies of chr1q and 2 copies of chr1p, and LoVo carries 3 copies of chr5p. We  
270 called integer copy numbers using a GMM generated for each chromosome arm and evaluated  
271 the classifier accuracy for correctly calling copy numbers in RPE1 (**Fig S2E**). Accuracy ranged  
272 from 91% (chr8q) to 50% (chr19p, chr22q) with a mean of 73%. The mean absolute deviation of  
273 accuracy from 5-fold cross validation ranged from 0.63 to 6.2 percentage points. Only 22 out of  
274 41 arms had accuracy values of at least 75% (**Fig 3B**), demonstrating generally lower accuracy  
275 compared to whole chromosomes, likely because fewer probes typically target an arm than an  
276 entire chromosome. A positive relationship between the number of probes and accuracy was  
277 again revealed by linear regression ( $R^2 = 0.65$ ;  $p = 5.1e-09$ ) (**Fig 3C**). We calculated the  
278 theoretical sensitivity for detecting arm-level copy number across the 5 copy number levels.  
279 Sensitivity was again highest for 1 copy with an average of 95%; 2 copies had an average of 74%  
280 and 3 copies had 53% (**Fig 3D**; **Table S2**). Overall, we found that our system can accurately call

281 copy numbers for the majority of chromosomes and several chromosome arms, with less  
282 sensitivity for smaller chromosomes and arms.

283

### 284 ***Downsampling the Number of Probes Decreases Accuracy***

285 To confirm that copy number classification accuracy/sensitivity is dependent on the number of  
286 probes targeted to a chromosome and not the size of the chromosome itself, we downsampled  
287 the probes targeting chromosomes 2 (23 total probes) and 6 (18 total probes) and recalculated  
288 the classification accuracy for the RPE1 cells. 50 samples of  $n$  probes were evaluated for each  
289 value of  $n$ . Consistent with our findings above, median accuracy decreased from 95.8% (24  
290 probes) to 68.8% (4 probes) for chr2 and from 91.4% (18 probes) to 66.8% (4 probes) for chr6,  
291 indicating that probe number, not chromosome size, affects classification accuracy (**Fig 4A**).  
292 Furthermore, the interquartile range (IQR) of the accuracy distributions increased from 0.8  
293 percentage points (pp; 22 probes) to 9.14 pp (4 probes) for chr2 and from 1.5 pp (16 probes) to  
294 11.6 pp (4 probes) for chr6 indicating that having fewer probes per chromosome increases the  
295 variability of classification accuracy. Both the decrease in accuracy and increase in accuracy  
296 variance could be observed for the theoretical sensitivity at all 5 copy number levels as well (**Fig**  
297 **S3A**).

298

### 299 ***Additional Probes Increase Sensitivity***

300 Since a greater number of probes correlates with higher copy number call sensitivity, we reasoned  
301 that we could further increase sensitivity for all chromosomes by increasing the number of probes  
302 targeting each chromosome. To determine the number of probes required to approach 100%  
303 sensitivity for all 5 copy number levels, we simulated a panel using all probes targeting  
304 chromosomes 1 through 6 (120 probes total) and smoothed their RPE1 cell-probe copy number  
305 scores as if they were measurements from a single hypothetical chromosome. This is possible  
306 because all 6 chromosomes have 2 copies in RPE1. For each trial, we constructed a new panel

307 probe-by-probe by sampling the 120 probes without replacement, recalculating RPE1 copy  
308 number and sensitivity at every step until all 120 probes were added, repeated for 50 trials. As  
309 expected, mean theoretical sensitivity increased with probe number for all copy number levels  
310 (**Fig 4B**). The simulation achieved at least 90% sensitivity on average at 4 probes per  
311 chromosome for 1 copy, 16 probes for 2 copies, 42 probes for 3 copies, and 78 probes for 4  
312 copies. A maximum mean sensitivity of 91.5% was achieved for 5 copies at 120 probes. Again,  
313 the variability of the copy number call sensitivity decreased as the number of probes increased.  
314 In some cases, it may be sufficient for the user to detect either a gain or loss of an otherwise  
315 diploid chromosome rather than detect the specific copy number of the gained chromosome. In  
316 this circumstance, a GMM can be generated with only 3 components, representing states of loss  
317 (1 copy), neutral (2 copies), and gain (3 or more copies). We evaluated our simulation under this  
318 model, achieving at least 90% sensitivity at 4 probes for 1 copy, 16 probes for 2 copies, and 20  
319 probes for  $\geq 3$  copies (**Fig S3B**). Furthermore, 99% sensitivity could be achieved at 17 probes for  
320 1 copy, 55 probes for 2 copies, and 67 probes for  $\geq 3$  copies. These findings indicate that copy  
321 number call sensitivity can be increased for all copy number levels by adding probes to our panel.

322

323

#### 324 ***KaryoTap Version 2 Panel Increases Accuracy of Aneuploidy Detection***

325 As the accuracy of our custom panel increases with the number of probes targeting a  
326 chromosome, we attempted to improve the accuracy by adding probes to chromosomes with  
327 lower coverage. To balance the cost of producing a larger custom panel with meaningful  
328 sensitivity gains, we removed the 61 least efficient probes (by total read counts) and added 82  
329 probes such that each chromosome was targeted by at least 12 probes (**Table S1**). We also  
330 included 4 probes targeting chrY, which was not covered by Version 1, to enable the detection of  
331 all 24 chromosomes (**Fig S4A**). The new panel, Version 2 (v2; CO610) comprises 352 total  
332 probes, 270 of which are shared with Version 1 (**Fig 1; Additional File 2**).

333  
334 To evaluate performance for the Version 2 panel, we performed a Tapestry experiment using  
335 RPE1 cells (**Fig S4B**) and determined the empirical accuracy as the proportion of cells with correct  
336 copy number calls based on known copy number from karyotyping. We again performed 5-fold  
337 cross validation by partitioning the RPE1 cells randomly into 5 equally sized subsets and  
338 calculating accuracy 5 times, each time reserving one of the sets from the model generation and  
339 using it only to calculate accuracy. The mean absolute deviation of the 5 accuracy measurements  
340 for each chromosome ranged from 0.92 to 3.89, suggesting that copy number calling performance  
341 would be maintained when classifying new data. The poorest performing chromosome in Version  
342 1, chr22, had an accuracy of 70% with Version 2 compared to 49% for Version 1 (**Fig 5A-B, Fig**  
343 **S4C**). The mean accuracy across all chromosomes was 89%, increased from 82% for Version 1.  
344 Accuracy increased by 2.2 pp on average for each additional probe (**Fig 5C**). Similarly, for  
345 chromosome arms, the average accuracy across arms increased from 73% with Version 1 to 80%  
346 with Version 2 (**Fig 5D-E, Fig S4D**). We also calculated the theoretical sensitivity for copy number  
347 values of 1, 2, 3, 4, and 5 and saw increased average sensitivity compared to Version 1 (**Fig S4E,**  
348 **Table S3**). Furthermore, we calculated theoretical sensitivity for a simpler 3-component model  
349 representing chromosome loss, neutral, and gain states, which may be a more practical choice  
350 for certain users. The 3-component model had a mean sensitivity of 99% for losses, 90% for  
351 neutral states, and 87% for gains for whole chromosomes; in addition it showed a sensitivity of  
352 and 97% for losses, 81% for neutral states, and 78% for gains for chromosome arms (**Table S3**).  
353 These data provide strong evidence that increasing the number of probes targeting each  
354 chromosome improves the sensitivity of the panel in calling copy numbers in individual cells.

355  
356 ***Detection of Lentiviral Barcodes and gRNAs***  
357 To extend the capabilities of our system, we added two probes to the Version 2 panel that target  
358 and amplify either a DNA barcode sequence or CRISPR guide RNA (gRNA) sequence integrated

359 into a cell's genome by lentiviral transduction. DNA barcoding of cells can be used in situations  
360 where several samples from the same cell line or individual are sequenced in one experiment and  
361 are therefore unable to be distinguished by genotype. Similarly, CRISPR gRNAs can be used  
362 both for functional studies and as barcodes themselves, indicating the gRNA treatment received  
363 by an individual cell.

364  
365 As a proof-of-concept, we transduced RPE1 cells and human colorectal epithelial cells (hCECs)  
366 each with distinct gRNA constructs (gRNA1 and gRNA2, respectively; **Table S4**), and human  
367 Pancreatic Nestin-Expressing cells (hPNEs) with a mix of two DNA barcode constructs that drive  
368 expression of BFP. We used distinct cell lines for each construct so that the three populations  
369 could be distinguished by genotype without assuming successful barcoding. To enable panel  
370 Version 2 to detect gRNAs, we designed a probe, Probe AMP350, to target the region surrounding  
371 and including the gRNA sequence in the lentiviral vector. To enable the detection of DNA  
372 barcodes from the BFP-expressing vector, we similarly designed a probe, AMP351, to target the  
373 region surrounding and including the barcode sequence (**Fig 5F**).

374  
375 gRNAs were each transduced into target cells with a multiplicity of infection (MOI) of 1-1.5  
376 followed by puromycin selection to ensure that each cell had an average of ~1 integration and at  
377 least 1 integration. The BFP barcodes were transduced at a higher MOI and cells were enriched  
378 for BFP expression by FACS; a high MOI was used for the barcoding sequences to increase the  
379 chances of detection. The three cell populations were pooled and analyzed in a single TapeStri  
380 experiment using panel Version 2. The populations were distinguished by PCA, UMAP, and  
381 clustering of VAFs (**Fig S4B**) as done previously. To determine if gRNA1 could be detected in  
382 RPE1 cells, we took the aligned reads from Probe AMP350 associated with RPE1 cells and  
383 searched for the sequence of gRNA1. RPE1 cells had an average of 34 gRNA1 reads per cell,  
384 while hCECs and hPNEs had 0. Similarly, hCECs had an average of 30 gRNA2 reads per cell,

385 while RPE1 and hPNEs had 0 (**Fig 5G**). To determine if the BFP barcodes could be detected in  
386 hPNEs, we similarly took aligned reads from Probe AMP351 and searched for the sequence of  
387 the barcodes. hPNEs had an average of 81 barcode reads per cell, while RPE1 and hCECs had  
388 0. Altogether, these data indicate that gRNA sequences and specific DNA sequences can be  
389 recovered from transduced cells using panel Version 2.

390  
391 To determine the limit of detection for a gRNA in transduced cells, we analyzed the proportion of  
392 cells with 0 reads matching the appropriate gRNA sequence. 21% of RPE1 and 31% of hCECs  
393 had 0 counts per cell for gRNA1 and gRNA2, respectively (**Fig 5G**). We compared the read counts  
394 per cell for Probe AMP350 with the number of reads matching the appropriate gRNA sequence  
395 in both RPE1 and hCECs and found that virtually all of the counted reads from Probe AMP350  
396 matched the number of gRNA1 sequence reads for RPE1 and the gRNA2 sequence reads for  
397 hCECs (**Fig S4F**), indicating no contamination from other sequences. Altogether, these data  
398 indicate that Probe AMP350 in Panel Version 2 can detect at least one gRNA sequence in ~70-  
399 80% of cells, though the rate of detection may be improved by transduction at a higher MOI.

400  
401 We repeated a similar analysis for Probe AMP351 to determine the limit of detection for a DNA  
402 barcode in transduced cells. 6% of hPNEs had 0 barcode sequence counts per cell (**Fig 5G**) and  
403 >99% of the reads from Probe AMP351 matched the known barcode sequences in hPNEs (**Fig**  
404 **S4F**), indicating no contamination. Altogether, these data indicate that Probe AMP351 in Panel  
405 Version 2 can detect a DNA barcode in 94% of cells, which may be improved by increasing depth  
406 or MOI.

407  
408 ***Evaluation of Aneuploidy Induction by KaryoCreate***

409 To demonstrate the combined copy number detection and multiplexing capabilities of our system,  
410 we tested it on samples treated with KaryoCreate (Karyotype CRISPR Engineered Aneuploidy

411 Technology), a method we recently developed to induce chromosome-specific aneuploidy in  
412 cultured cells(14). KaryoCreate uses CRISPR gRNAs to target a mutant KNL1-dCas9 fusion  
413 protein to the centromere of a specific chromosome, causing missegregation in ~20% of cells.  
414 KaryoTap represents an ideal method to evaluate the efficiency of aneuploidy induction and  
415 chromosome-specificity of KaryoCreate. To do this, we performed a Tapestri experiment on  
416 hCECs that had been treated with 1 of 3 gRNAs previously tested using KaryoCreate: sgNC does  
417 not have a target and is used as a negative control, sgChr6-2 targets chr6 and sgChr7-1 targets  
418 chr7 (**Table S4**). The gRNA sequences amplified by AMP350 were used to identify the gRNA that  
419 each cell received. sgChr6-2 and sgChr7-1 induced gains and losses specifically in the intended  
420 chromosomes, but not others, compared to sgNC (**Fig 6A, Table S5**;  $p < 0.01$ , Fisher's exact  
421 test). sgChr6-2 induced 26.6% losses of chr6 compared to 0.5% with sgNC, and 8.3% gains  
422 compared to 2.0% with sgNC (**Table S6**). sgChr7-1 induced 6.1% losses of chr7 compared to  
423 0.5% with sgNC, and 5.1% gains compared to 4.3% with sgNC.

424  
425 In the same experiment, we also performed a small screen to address a current limitation of  
426 KaryoCreate in which we were unable to engineer aneuploidy of certain chromosomes, such as  
427 chromosome 20, one of the most frequently gained chromosomes in human cancer (2). In fact,  
428 while we could design gRNAs that are specific to the centromere of chromosome 20, it was not  
429 possible to visualize centromeric foci through the co-transduction of cells with gRNAs and  
430 fluorescently-tagged dCas9 by imaging, possibly due to the low (~700) number of gRNA binding  
431 sites (14). Furthermore, given the small size of chromosome 20, the single-cell RNA sequencing-  
432 based approach used in Bosco et al. does not have sufficient sensitivity to confidently assess  
433 gains and losses of this chromosome. Thus, using KaryoTap, we screened 5 sgRNAs targeting  
434 chromosome 20 (sgChr20-2, 20-3, 20-4, 20-6 and 20-7) that were previously described but not  
435 validated by imaging (14). sgChr20-2, 20-4, 20-6 and 20-7 did not induce changes in chr20 ( $p =$   
436 0.67-0.94). sgChr20-3 was able to induce 7.4% gains in chr20 compared to 3.0% with the sgNC



437 control, and 4.7% losses compared to 1.8% with sgNC ( $p = 0.006$ ). We also note that sgChr20-3  
438 induced 9.4% losses and 17.8% gains in chr2 ( $p < 0.001$ ), which we might not have observed if  
439 we had instead evaluated the effect of the gRNA using a chromosome-targeted method such as  
440 fluorescence in situ hybridization rather a method which covers all chromosomes. The sgChr20-  
441 3 sequence (**GGCAGCTTTGAGGATTTCTG**) matches 18 out of 20 base pairs for loci on the chr2  
442 centromere (**GATAGCTTTGAGGATTTCTG**) (14), suggesting an explanation for the off-target  
443 effect. These data indicate that KaryoTap successfully enables simultaneous detection of  
444 aneuploidy and gRNA/barcodes in the same cells and thus can be used to perform CRISPR-  
445 based (i.e., gRNA-based) or ORF-based (barcode-based) functional screens.

446

#### 447 ***Detection of Copy Number Neutral Loss of Heterozygosity (CNN-LOH)***

448 Gains and losses of diploid chromosomes result in a shift in VAF for their heterozygous SNPs  
449 from 50% in the direction of 100% or 0% depending on which parental chromosome copy (i.e.,  
450 haplotype) experienced a copy number change. In addition, a shift in VAF can also be observed  
451 in the absence of copy number changes in copy number neutral loss of heterozygosity (CNN-  
452 LOH), which has been observed in cancer as well as normal tissues(33,34) . Because each probe  
453 is sequenced at a high depth, KaryoTap should be able to detect this shift, allowing us to  
454 determine which of the two parental chromosomes/haplotypes was gained or lost. This is  
455 especially important for detecting loss of heterozygosity (LOH), a common event in cancer  
456 whereby a heterozygous-to-homozygous shift by chromosomal loss can inactivate tumor  
457 suppressor genes(35,36). To determine if KaryoTap could detect allele frequency shifts following  
458 chromosomal gains and losses, we examined the cells from the KaryoCreate experiment (**Fig 6A**)  
459 that had lost a copy of chromosome 6 after treatment with the sgChr6-2 gRNA. We identified 9  
460 heterozygous variants on chr6 called by the Tapestry Pipeline by identifying variants with a mean  
461 allele frequency between 20-80% in the sgNC control population. We then calculated a relative  
462 (i.e., haplotype-agnostic) allele frequency for sgChr6-2 treated cells with 1 or 2 copies of chr6

463 (called by GMM) by calculating the absolute difference between raw allele frequency and 50%  
464 such that 0 corresponded to heterozygous alleles and 50 corresponded to fully homozygous  
465 alleles (**Fig 6B**). 4 distinct clusters of cells emerged using K-means clustering on relative AFs.  
466 The cluster comprising cells with a copy number call of 2 for chr6 shows that the 9 variants are  
467 heterozygous in diploid cells as expected. The cluster with 1 copy of chr6 shows a shift across  
468 the variants from heterozygous to homozygous (i.e. a loss of heterozygosity), supporting the loss  
469 of one copy of each allele in these cells. There are also two smaller clusters representing the loss  
470 of either chromosome arm but not the other, supported by both the loss of heterozygosity in the  
471 variants on the affected arm and the copy number call of 1 for that arm. This indicates that  
472 KaryoTap can be used to detect loss of heterozygosity in single cells at the population level.

473  
474 It is possible that the loss of a chromosome can be followed by a duplication of the remaining  
475 chromosome, such that the copy number of the chromosome remains the same, but one allele is  
476 lost, i.e., a CNN-LOH. To determine if KaryoCreate can cause CNN-LOH in the targeted  
477 chromosome, we took the relative allele frequencies for sgChr6-2 treated cells with 1 or 2 copies  
478 of chr6 calculated above and averaged them such that each cell had one mean relative AF value  
479 for chr6. We also repeated this calculation for the cells treated with the sgNC control gRNA. If we  
480 consider relative AF between 40% and 50% to indicate homozygosity of chr6 alleles and 0% to  
481 40% to indicate heterozygosity, all sgNC-treated cells with a chr6 copy number of 2 were  
482 heterozygous (**Fig 6C**). Cells treated with sgChr6-2 that lost a copy were detected to be  
483 homozygous. 85% of cells treated with sgChr6-2 that had 2 copies of chr6 detected were  
484 heterozygous, while 15% were detected as homozygous, indicating a loss of heterozygosity with  
485 no change in the net copy number (2 copies of chr6). This indicates that KaryoCreate can induce  
486 CNN-LOH and KaryoTap can detect CNN-LOH events.

487

488 ***KaryoTap Version 3 Panel Further Improves Accuracy of Aneuploidy Detection***

489 To further improve the accuracy of our system for detection of aneuploidy, especially for  
490 chromosome arms, we modified panel Version 2 to create Version 3 (v3, CO810; **Fig 1;**  
491 **Additional File 3**). Version 3 comprises 399 total probes, 309 of which are shared with Version  
492 2. We removed 43 less-efficient probes from Version 2 and added 90 new probes, prioritizing  
493 chromosome arms with less coverage. The barcode and gRNA detecting probes described above  
494 were also included in the new design.

495  
496 To evaluate the performance of panel Version 3, we performed a Tapestry experiment using RPE1  
497 cells, made copy number calls using the GMM strategy as described above, and calculated  
498 accuracy as the proportion of cells with correct copy number calls based on known copy number.  
499 We again performed 5-fold cross validation by partitioning the RPE1 cells randomly into 5 equally  
500 sized subsets and calculating accuracy 5 times, each time reserving one of the sets from the  
501 model generation and using it only to calculate accuracy. The mean absolute deviation of the 5  
502 accuracy measurements for each chromosome ranged from 0.60 to 3.31. The poorest performing  
503 chromosome in Version 2, chr22, had an accuracy of 78% with Version 3 compared to 70% for  
504 Version 2 (**Fig 7A-B**). The mean accuracy across all chromosomes was 91%, increased from  
505 89% for Version 2. Similarly, for chromosome arms, the average accuracy across arms increased  
506 from 80% with Version 2 to 83% with Version 3 (**Fig 7C-D**). We also calculated the theoretical  
507 sensitivity for copy number values of 1, 2, 3, 4, and 5 and saw increased average sensitivity  
508 compared to Version 2 (**Fig S5A-B, Table S7**) for both whole chromosomes and chromosome  
509 arms. Furthermore, we calculated theoretical sensitivity for a simpler 3-component model  
510 representing chromosome loss, neutral, and gain states. The 3-component model had a mean  
511 sensitivity of 99% for losses, 91% for neutral states, and 88% for gains for whole chromosomes;  
512 in addition it showed a sensitivity of 97% for losses, 83% for neutral states, and 80% for gains for  
513 chromosome arms (**Table S7**). When compared qualitatively, the general degree of noisiness in

514 heatmaps of the GMM copy number calls for RPE1 decreases across chromosome arms between  
515 KaryoTap panels Version 1, 2 and 3, supporting improvement in the accuracy of copy number  
516 calling afforded by panel Version 3 (**Fig 7F**). Noisiness also decreases in heatmaps of copy  
517 number calls for the LoVo and LS513 cell lines between panels Version 1 and Version 3 (**Fig 7G-**  
518 **H**), supporting the improvement of copy number variant detection sensitivity for individual cells  
519 using our system. Altogether, these data indicate that panel Version 3 can deliver accurate copy  
520 number calls in thousands of single cells.

521

522

523

## 524 **Discussion**

525 We designed and improved two custom panels for the Tapestri platform that enable targeted  
526 scDNA-seq of 330-352 specific loci across all 24 human autosomes and sex chromosomes. This,  
527 coupled with a GMM-based copy number calling analysis pipeline, allowed us to identify  
528 chromosome- and chromosome arm-scale aneuploidy in thousands of individual cells in a single  
529 experiment with high accuracy and at greatly reduced sequencing depth compared to single-cell  
530 WGS methods. To increase the ease-of-use, we compiled the computational scripts used to  
531 analyze these data into an R package, *karyotapR*, which automates all steps for calling copy  
532 number and for basic visualization of the results (**Fig S6**).

533

534 While single-cell aneuploidy detection is not unique to KaryoTap, our design has several  
535 advantages, the most critical being the leveraging of Tapestri's throughput thus significantly  
536 reducing hands-on time, reagent cost, and sequencing cost per cell compared to low-throughput  
537 WGA methods. Additionally, targeting of specific loci allows us to forgo the typical technical  
538 difficulties of conventional WGS analysis, including correcting for mappability bias and GC bias,  
539 and the use of segmentation algorithms to make copy number calls(27). Furthermore, the total

540 number of sequencing reads needed to obtain high-accuracy copy number calls for KaryoTap is  
541 greatly reduced compared to WGS-based scDNA-seq methods and the targeted nature of the  
542 assay spreads the reads more evenly across the genome, preventing the biasing of detection  
543 sensitivity toward larger chromosomes that is seen with WGS. We expect this will be particularly  
544 important for assessing chromosomal instability as the smallest of chromosomes will be more  
545 equally represented in the data relative to the larger chromosomes. Finally, the commercial  
546 availability of the Tapestry system allows for easier adoption compared to non-commercialized  
547 "homebrew" methods that need to be established and optimized in each lab from scratch(23).

548  
549 As evidenced by the analysis of RPE1, our method will generate a range of smooth copy number  
550 measurements for chromosomes with the same discrete copy number, indicating some level of  
551 technical error. To account for this error and convert the continuous smooth copy number scores  
552 to discrete copy number values, we used a Gaussian mixture model (GMM) classification  
553 strategy, which has been previously used for copy number analysis of single-cell whole genome  
554 sequencing data(37). This allows each smooth copy number score to be associated with a set of  
555 (posterior) probabilities of being measured from a chromosome of a given range of copy numbers  
556 (e.g., 1, 2, 3, 4, or 5). While we assign each smooth score to the discrete copy number value for  
557 which its posterior probability of belonging is highest, the probabilities of belonging to the other  
558 copy number components of the model indicate the confidence the investigator can have that the  
559 call is accurate. As the number of probes increases, the variance of the model components  
560 decreases, resulting in an increase in classifier accuracy that we observe between our Version 1  
561 and Version 2 panel designs. While we use a near-diploid cell line as our ground truth, the GMM  
562 strategy also allows us to calculate the expected (theoretical) sensitivity that a chromosome with  
563 copy number 1, 2, 3, 4, or 5 would be correctly called using our system by calculating the  
564 proportion of overlap between the copy number components of the model. We used this in our  
565 panel simulation to extrapolate an optimal number of probes for detecting 1 copy (loss), 2 copies

566 (neutral) or 3 copies (gain) of a chromosome and determined that at least 90% sensitivity could  
567 be achieved at 4 probes for 1 copy, 20 probes for 2 copies, and 26 probes for  $\geq 3$  copies, and 99%  
568 sensitivity could be achieved at 22 probes for 1 copy, 66 probes for 2 copies, and 76 probes for  
569  $\geq 3$  copies. Further improving the panel by increasing the number of probes and thus reducing the  
570 technical variation will allow us to more confidently observe the karyotype heterogeneity in these  
571 samples as well as in tumors and other tissues. It is important to note here that our method  
572 requires a baseline sample with which to compare the other cells in the experiment. However,  
573 while we set the baseline copy number using near-diploid RPE1 to scale the read counts of each  
574 probe relative to 2 copies, it is not strictly necessary to spike a diploid control cell line into the  
575 sample preparation. Any distinguishable, largely homogenous subset of cells in an experiment  
576 can be used to set the baseline as long as the average copy number for each chromosome in  
577 that subset is known.

578  
579 Deep sequencing ( $\sim 80$ - $100\times$  on average per cell) of each target region allows for robust single  
580 nucleotide variant (SNV) calling that is not possible at the lower genomic coverage afforded by  
581 other high-throughput methods(23). This enabled us to resolve and identify 5 multiplexed cell lines  
582 in a single experiment by clustering cells by variant allele frequencies. Since our panels  
583 specifically target loci known to harbor SNPs across the human population, we can extend sample  
584 multiplexing to clinical samples (e.g., tumor tissue) from different individuals without the need for  
585 barcodes. While we demonstrated sample identification using a clustering approach here, sample  
586 identities for each cell can be determined directly from known SNPs that occur at sequenced loci.  
587 Additionally, while our panels were designed for copy number analysis, additional probes could  
588 be added that cover tumor suppressor genes and oncogenes of interest, thus revealing  
589 consequential point mutations alongside chromosomal copy number. Mission Bio offers several  
590 ready-made panels covering mutational hotspots and genes relevant to a broad range of tumor  
591 types, allowing for a great degree of customizability. Furthermore, we demonstrated that SNV-

592 associated allele frequency shifts detected using KaryoTap could be used to infer loss of  
593 heterozygosity (LOH), a common event in cancer where the germline heterozygous state of a  
594 chromosome changes to a homozygous state in tumor cells(38). LOH has been demonstrated to  
595 promote tumorigenesis by inactivating tumor suppressor genes through chromosomal  
596 loss(35,36). In the context of chromosomal instability, the chromosome remaining after the loss  
597 of its homologue can be duplicated, resulting in LOH with a net-neutral copy number change (copy  
598 number neutral (CNN)-LOH). The deep sequencing depth and copy number detection enabled  
599 by KaryoTap allow for discovery of CNN-LOH events, which would otherwise be difficult to detect  
600 with the shallow coverage typical of other scDNA-seq methods(38).

601  
602 To enable experimental design flexibility when using our custom panels, we added a set of probes  
603 that can detect DNA barcode and CRISPR gRNA sequences integrated into the genome. DNA  
604 barcodes can be used to multiplex and resolve cells belonging to different samples in a single  
605 Tapestry experiment that otherwise could not be distinguished by genotype. Through barcoding,  
606 users can compare samples from the same individual or compare experimental and control  
607 conditions in the same cell lines while minimizing batch effects. Regardless of design, combining  
608 several samples into one experimental run greatly reduces the per-sample reagent and  
609 sequencing costs in addition to the hands-on time required to process the samples. Detecting  
610 barcodes in thousands of cells is made possible by exploiting the targeted nature of the  
611 sequencing assay. Single-cell DNA sequencing methods with comparable throughput rely on  
612 inefficient and random transposon insertion, which would only detect a randomly inserted barcode  
613 in about 20% of cells(23). By specifically targeting the barcoded insert, we can reliably recover  
614 the barcode sequence in over 90% of cells. Including a probe that targets inserted CRISPR gRNA  
615 sequences allows for an additional layer of experimental design flexibility where the gRNA-  
616 mediated treatment each cell receives can be identified by the gRNA sequence itself. Here, we  
617 demonstrated the gRNA detection and multiplexing capabilities of our system by evaluating the

618 efficiency and specificity of KaryoCreate, our method for inducing chromosome specific  
619 aneuploidy. Since the gRNAs can be detected in 70-80% of cells when transduced at low MOI,  
620 this system could also be used for CRISPR screen applications where cells are randomly treated  
621 with one gRNA from a library of hundreds of possible gRNAs and thus require high detection  
622 sensitivity(39).

623

624 Our system in its current form is limited to calling whole chromosome aneuploidy, and, with less  
625 confidence, chromosome arm-level aneuploidy. Further optimization of the panels will be required  
626 to achieve sufficient confidence in copy number detection for some chromosome arms. The  
627 sensitivities of either measurement vary according to the number of probes used, and thus  
628 confidently calling aneuploidy in increasingly smaller regions becomes challenging. Sub-arm (i.e.,  
629 focal) aneuploidy could potentially be detected with a greater density of probes, though  
630 manufacturing increasingly larger panels also increases the cost of the panels. Fortunately, this  
631 cost can be offset by lowering sequencing depth, as we demonstrated that sequencing at as low  
632 as ~35 average reads per cell per probe does not significantly affect aneuploidy call accuracy.

633

634 Coupled with the Tapestry platform, KaryoTap shows considerable promise as an easily  
635 adoptable, flexible, and highly scalable method for detecting chromosome- and chromosome arm-  
636 scale aneuploidy in thousands of single cells. Here we demonstrated population-level copy  
637 number detection in several cell lines and, most significantly, highly accurate copy number  
638 classification in individual cells using a Gaussian mixture model framework, which is otherwise  
639 unattainable using currently available methods. We identified the number of PCR probes per  
640 chromosome as a dominant factor affecting copy number classification performance and  
641 calculated the number of probes necessary to sufficiently improve detection sensitivity for various  
642 applications. Finally, we applied our method to the aneuploidy-engineering tool KaryoCreate to  
643 demonstrate sample multiplexing capabilities and the ability to detect gRNAs in transduced cells.



644 We believe this system lays the groundwork for a new class of tools for studying aneuploidy and  
645 chromosomal instability in healthy and diseased tissues and tumors.

646

## 647 **Methods**

648

### 649 *Cell Culture*

650

651 All cells were grown at 37°C with 5% CO<sub>2</sub> levels. All cell media was supplemented with 1X pen-  
652 strep, and 1X L-glutamine. hTERT human retinal pigment epithelial cells (RPE-1; ATCC CRL-  
653 4000) and SW48 cells (ATCC CCL-231) were incubated in DMEM, supplemented with 10% FBS.  
654 LoVo cells (ATCC CCL-229) were incubated in Ham's F12-K media with 10% FBS. LS513 cells  
655 (ATCC CRL-2134) and hTERT human pancreatic nestin-expressing cells (hPNEs; ATCC CRL-  
656 4023) were incubated in RPMI media with 10% FBS. CL11 cells (Cellosaurus CVCL\_1978) were  
657 incubated in DMEM:F12 and 20% FBS. hTERT p53<sup>-/-</sup> human colonic epithelial cells (hCECs; Ly  
658 et al.(40)) were cultured in a 4:1 mix of DMEM:Medium 199, supplemented with 2% FBS, 5 ng/mL  
659 EGF, 1 µg/mL hydrocortisone, 10 µg/mL insulin, 2 µg/mL transferrin, 5 nM sodium selenite, pen-  
660 strep, and L-glutamine. For long-term storage, cells were cryopreserved at -80°C in 70% cell  
661 medium, 20% FBS, and 10% DMSO. All cell lines were tested for mycoplasma.

662

### 663 *Custom Tapestry Panel Design*

664

665 Panel Version 1 (CO261) comprises 330 probes across the 22 human autosomes and the X  
666 chromosome. To identify candidate target regions for the panel, we used the Common SNP files  
667 downloaded from UCSC(41,42) (snp151Common, hg19), and considered only synonymous  
668 variant SNPs with a major allele frequency at >0.5 and <0.6. For cytobands with more than 4  
669 synonymous variants, we split the cytoband into 4 subregions based on the percentile of the

670 cytoband coordinates (0-25th percentile, 25-50th percentile, 50-75th percentile and 75-100th  
671 percentile). From each subregion, we randomly selected 1 SNP as a representative candidate. In  
672 cases where there were less than 5 synonymous variant SNPs, all SNPs were used. We  
673 submitted all candidate SNPs to the Tapestry Panel Designer to generate a panel design and  
674 ensured that the designed probes targeted the candidate SNPs and had similar GC contents.  
675 Next, randomly selected probes such that each chromosome had a probe density of ~1 per 10MB.  
676 Panel Version 2 (CO610) comprises 352 probes across all 24 human chromosomes. This panel  
677 was generated using Panel v1 as a base: first, we removed 61 probes that had low PCR  
678 amplification efficiency based on total read counts per probe. Then we added 82 probes such that  
679 each chromosome was targeted by at least 12 probes and included 4 probes targeting chrY. To  
680 enable the detection of lentiviral-delivered gRNAs, we added one probe targeting the region of  
681 the construct containing the gRNA sequence and one probe targeting a region upstream as a  
682 vector control. Similarly for the detection of lentiviral-delivered DNA barcodes, we added one  
683 probe targeting the region of the construct surrounding the barcode sequence, and one probe  
684 targeting a region downstream as a vector control. Support for the custom panel design and  
685 synthesis of the panel was provided by Mission Bio (San Francisco, CA, USA). Panel maps were  
686 created using the karyoploteR R package(43).

687

### 688 *Tapestry Single Cell DNA Sequencing*

689

690 Cell lines were trypsinized for 2-3 minutes, washed in room temperature  $Mg^{2+}/Ca^{2+}$ -free DPBS,  
691 centrifuged at 300g for 5 minutes, and resuspended in DPBS at a concentration of 3K cells/uL.  
692 For the experiment using the RPE1, SW48, LS513, LoVo, and CL11 cell lines, 600K cells from  
693 each cell line were combined, centrifuged at 300g for 5 minutes, and resuspended in Tapestry  
694 Cell Buffer at a concentration of 3.5K cells/uL. For the experiment using the RPE1, hPNE, and  
695 hCEC cell lines, 45K cells from each cell line were combined, centrifuged at 300 x g for 5 minutes,

696 and resuspended in Tapestri Cell Buffer at a concentration of 4K cells/uL. For the KaryoCreate  
697 experiment, ~100K cells from each condition were combined, centrifuged at 300 x g for 5 minutes,  
698 and resuspended in Tapestri Cell Buffer at a concentration of 3.4K cells/uL. Cell droplet  
699 encapsulation, barcoding, and sequencing library preparation were performed using the Tapestri  
700 instrument according to the manufacturer's instructions (Mission Bio, San Francisco, CA, USA).  
701 Sequencing was performed using an Illumina NovaSeq 6000 or NextSeq 500 in 2x150bp paired-  
702 end format. After sequencing, deconvolution of barcodes, read counting, and variant calling were  
703 handled by the online Tapestri Pipeline (v2.0.2)(26). The pipeline outputs both read counts per  
704 probe for each cell and variant allele frequencies for called variants for each cell.

705

#### 706 *Low Pass Whole Genome Sequencing & Karyotyping*

707

708 Genomic DNA was extracted from cell pellets using 0.3 µg/µL Proteinase K (QIAGEN #19131) in  
709 10mM Tris pH 8.0 for 1 hour at 55°C, following heat inactivation at 70°C for 10 minutes. DNA was  
710 digested using NEBNext dsDNA Fragmentase (NEB #M0348S) for 25 minutes at 37°C followed  
711 by magnetic DNA bead cleanup with 2X Sera-Mag Select Beads (Cytiva #29343045). Library  
712 prep was performed using NEBNext Ultra II DNA Library Prep Kit for Illumina (NEB #E7103)  
713 according to the manufacturer's instructions, generating DNA libraries with an average library size  
714 of 320 bp. Quantification was performed using a Qubit 2.0 fluorometer (Invitrogen #Q32866) and  
715 the Qubit dsDNA HS kit (Invitrogen #Q32854). Libraries were sequenced on an Illumina NextSeq  
716 500 at a target depth of 4-8 million reads. Reads were trimmed using trimmomatic(44), aligned to  
717 the hg38 genome using bwa-mem(45), and analyzed for copy number variants using the  
718 CopywriteR(46) R package. G-banded karyotyping of 100 RPE-1 cells was performed by WiCell  
719 Research Institute, Inc. (Madison, WI).

720

721 *Cloning of sgRNAs*

722

723 We modified the scaffold sequence of pLentiGuide-Puro (Addgene #52963) by Gibson assembly  
724 to contain the A-U flip (F) and hairpin extension (E) described by Chen et al(47). for improved  
725 sgRNA-dCas9 assembly, obtaining pLentiGuide-Puro-FE. sgRNAs were designed and cloned  
726 into this pLentiGuide-Puro-FE vector according to the Zhang Lab General Cloning Protocol(48).  
727 To be suitable for cloning into *BbsI*-digested vectors, sense oligos were designed with a CACC 5'  
728 overhang and antisense oligos were designed with an AAAC 5' overhang. The sense and  
729 antisense oligos were annealed, phosphorylated, and ligated into *BbsI*-digested pLentiGuide-  
730 Puro-FE for KaryoCreate purposes. Sequences were confirmed by Sanger sequencing.

731

732 *Gateway Recombination Cloning for Generation of Barcode Library*

733

734 pHAGE-CMV-DEST-PGKpuro-C-BC was a library of lentiviral vectors containing 24-bp random  
735 barcodes that was built as described in Sack & Davoli et al., 2018(49). Destination vector pHAGE-  
736 CMV-DEST-PGKpuro-C-BC and entry vector pDONR223\_BFP (Addgene: 25891) are  
737 recombined following the manufacturer's protocol. Briefly, 50 ng of entry vector and 100 ng of  
738 destination vector are mixed with LR Clonase™ enzyme and incubated overnight at room  
739 temperature. The next day, the reaction mixture is incubated with Proteinase K at 37°C for 10  
740 minutes, followed by inactivation at 75°C for 15 minutes. The reaction is then transformed into  
741 stb13 bacterial competent cells, plated onto LB agar plates, and incubated overnight at 37°C.  
742 Individual clones are collected into 96 well plates and expanded. Plasmid is extracted from the  
743 bacterial culture using a 96-well mini-prep kit (Zymo kit, Zippy 96 plasmid kit). All clones are  
744 sequenced by Sanger sequencing at the site of the barcode using primer  
745 ACTTGTGTAGCGCCAAGTGC. Duplicates are eliminated, and unique barcodes are retained in

746 the final library. BFP expression and Puromycin selection are validated by transfecting randomly  
747 selected clones into HEK293T cells.

748

#### 749 *Lentivirus Production and Nucleofection*

750

751 For transduction of cells, lentivirus was generated as follows: 1 million 293T cells were seeded in  
752 a 6-well plate 24 hours before transfection. The cells were transfected with a mixture of gene  
753 transfer plasmid (2 µg) and packaging plasmids including 0.6 µg ENV (VSV-G; addgene #8454),  
754 1 µg Packaging (pMDLg/pRRE; addgene #12251), and 0.5 µg pRSV-REV (addgene #12253)  
755 along with CaCl<sub>2</sub> and 2x HBS or using Lipofectamine 3000 (Thermo #L3000075). The medium  
756 was changed 6 hours later and virus was collected 48 hours after transfection by filtering the  
757 medium through a 0.45-µm filter. Polybrene (1:1000) was added to the filtered medium before  
758 infection.

759

#### 760 *KaryoCreate Experiments*

761

762 KaryoCreate experiments were performed as described in Bosco et al., 2023(14). Briefly, p53-/-  
763 hCEC were first lentivirally transduced with pHAGE-KNL1Mut-dCas9 and selected with  
764 blasticidin. The cells were then lentivirally transduced with the indicated sgRNAs and selected  
765 with puromycin. scDNA seq was performed ~10 days after transduction with the gRNAs. The  
766 sequences of the gRNAs targeting the centromeres of specific chromosomes are listed in Table  
767 S4 and were designed as described in Bosco et al., 2023(14). To compare conditions, Fisher's  
768 exact test was performed in R using the `fisher.test()` function, comparing the proportion of cells  
769 for each chromosome and sample that are diploid (copy number = 2) and aneuploid (copy number  
770 = {1, 3, or 4}) between the sgNC control and sample and the given experiment sample. The

771 Benjamini-Hochberg correction for multiple comparisons was applied to p-values using  
772 ``p.adjust()``.

773

#### 774 *Parsing and Counting of Barcoded Reads*

775

776 To detect specific gRNA or DNA barcode sequences, we searched for the known sequences  
777 against the cell-associated aligned reads (cells.bam file) generated from the Tapestri Pipeline.  
778 Search queries were conducted `vcountPattern()` using the Biostrings R package(50), with  
779 tolerance for up to 2 base mismatches. BAM files were manipulated using the Rsamtools R  
780 package(51).

781

#### 782 *Cell Line Demultiplexing and Identification*

783

784 To demultiplex cells from different cell lines in the initial scDNA-seq experiment, we use the allele  
785 frequency (AF) of variants that are called by GATC as part of the Tapestri Pipeline. Variants were  
786 filtered by selecting those with standard deviations of AF >20 to select variants whose allele  
787 frequencies vary the most across all cell lines. PCA was used to reduce the dimensions of the  
788 remaining variants. The top 4 principal components were embedded in two dimensions by UMAP  
789 and then clustered using the dbscan method. The 5 clusters with the greatest number of cells  
790 were kept, corresponding to the 5 expected cell lines. The remaining clusters, likely representing  
791 cell doublets, were discarded from further analyses. This method was repeated for subsequent  
792 Tapestri experiments, adjusting for the expected number of cell populations.

793

794 The cluster containing RPE1 cells in each experiment was identified by clustering with published  
795 deep WGS of RPE1. Published RPE1 WGS data was obtained from SRA Accession  
796 ERR7477340(29,52). Reads were aligned to the hg19 genome using `bwa`. `MarkDuplicatesSpark`

797 and HaplotypeCaller from GATK were used to mark duplicate reads and get AFs from called  
798 variants(53). The vcfR R package(54) was used to extract the AFs for called variants common to  
799 the published data and our dataset. The mean AF for each variant was calculated for each of the  
800 5 cell lines. PCA was used to cluster our mean AF dataset with the RPE1 AFs. The cell line that  
801 clustered most closely with the published RPE1 data was labeled as RPE1 cells.

802

### 803 *Copy Number Calling*

804

805 Copy number scores for each probe in each cell (cell-probe scores) were calculated relative to  
806 RPE1, for which we know the copy number of each chromosome: The raw count matrix was  
807 normalized by scaling each cell's mean to 1 (Equation 1) and then each probe's median to 1  
808 (Equation 2). The normalized counts were then scaled such that the value of the median  
809 normalized RPE1 counts for each probe was set to 2 for all probes except those targeting chr10q,  
810 which were set to 3 (Equation 3). The identities of the remaining 4 populations of cells were  
811 identified by comparing their overall copy number profile with matched bulk WGS data.

812

$$813 \quad RC_{intermediate} = \frac{RC_{cell,probe}}{\text{mean}(RC_{cell}) + 1} \quad (1)$$

814

$$815 \quad RC_{norm} = \frac{2 \cdot RC_{intermediate}}{\text{median}(RC_{intermediate,probe}) + 0.05} \quad (2)$$

816

$$817 \quad CN_{cell,probe} = \frac{CN_{WGS,probe} \cdot RC_{norm}}{\text{median}(RC_{norm,RPE1})} \quad (3)$$

818

819 Smooth copy number scores for each chromosome in each cell (cell-chromosome scores) were  
820 generated by taking the weighted median of the probe-specific copy number values for probes

821 targeting a common chromosome (Equation 4). Weights for each probe were calculated as the  
822 proportion of RPE1 cell-probe scores that fell within  $\pm 0.5$  of the known copy number (3 for chr10q,  
823 2 for all others). This was also modified to calculate cell-chromosome-arm scores for probes  
824 common to a chromosome arm.

825

826  $CN_{smooth,cell,chr} = weightedMedian(CN_{cell,probe})$ , for all probes on chromosome *chr* (4)

827

828 Integer copy number values for each cell-chromosome were classified using Gaussian mixture  
829 models (GMMs) with either five components representing possible copy number values of 1, 2,  
830 3, 4, and 5, or three components representing copy number values 1, 2, and 3. To generate the  
831 GMMs, the normalized counts for each probe for the RPE1 cells were fitted to Weibull distributions  
832 using the `fitdistrplus` R package(55). These Weibull parameters represented parameters for copy  
833 number = 2 for all probes except those targeting chr10q, which has 3 copies in RPE1. The scale  
834 parameters were then scaled for possible copy number values 1 through 6, relative to the RPE1  
835 copy number: probes with RPE1 copy number = 2 were scaled by 50%, 100%, 150%, 200%,  
836 250%, and 300% for copy number = 1, 2, 3, 4, 5 and 6; probes with RPE1 copy number = 3 were  
837 scaled by 33%, 67%, 100%, 133%, 167%, and 200% for copy number = 1, 2, 3, 4, 5 and 6. 500  
838 Weibull-distributed values are drawn using each of the six parameter sets for each probe to  
839 simulate six matrices of 500 simulated cells. For each cell, the values were smoothed across the  
840 probes belonging to each chromosome to simulate cell-chromosome copy number values. The  
841 distribution of the scores for each chromosome was then fit to Gaussian (normal) distributions,  
842 separately for each copy number level. The result is a set of normal parameters (mean  $\mu$  and  
843 standard deviation  $\sigma$ ) for each chromosome for each value of copy number  $k = \{1, 2, 3, 4, 5, 6\}$ .  
844 The six copy number Gaussian components for each chromosome were combined into a GMM,  
845 representing the probability densities for each copy number value for that chromosome (Equation



846 5). Using Bayes rule and assuming equal priors, the posterior probability of a cell-chromosome  
847 copy number score being generated under each component  $k$  is given by Equation 6.

848

849 
$$pdf(x = CN_{smooth,cell,chr}) = \frac{1}{\sigma\sqrt{2\pi}} \cdot e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \quad (5)$$

850

851 
$$P(k|x) = pdf_k(x) / \sum_{k=1}^6 pdf_k(x) \quad (6)$$

852

853 Decision boundaries for the GMMs are calculated by finding the transitions between components,  
854 i.e., the point  $x$  where the PDFs of the components are equal. We evaluated copy numbers using  
855 GMMs including copy number components 1-5, throughout the study and 1-3 where indicated.  
856 Upper boundaries for component 5 were calculated using components 1-6. Theoretical sensitivity  
857 for each copy number component was calculated as the proportion of the component PDF that  
858 falls within its decision boundaries (i.e., true positive rate). 5-fold cross validation was performed  
859 by partitioning RPE1 cells into 5 equally sized groups and using each group once to evaluate a  
860 model generated using the remaining 4 groups. R scripts for copy number calling were compiled  
861 into an R package, karyotapR. karyotapR version 0.1 was used for analyses in this study.

862

### 863 *Panel Simulations*

864

865 For the probe downsampling simulation of chromosome 2, 50 samples each of  $n$  probes from the  
866 set of 24 probes targeting chr2 were generated where  $n = \{4, 6, 8 \dots 20, 22\}$ . Copy numbers were  
867 called for each set of probes for each cell. The sensitivity of the copy number calls was  
868 recalculated for the RPE1 cells as well as the theoretical sensitivity for all GMM components. This  
869 analysis was repeated for the set of probes targeting chr6 where  $n = \{4, 6, 8 \dots 14, 16\}$ .

870

871 For the simulation of an expanded custom panel, the set of 120 probes targeting chromosomes  
872 1, 2, 3, 4, 5, and 6 were sampled 50 times to produce 50 sets of the 120 probes in unique orders.  
873 Starting with the first 4 probes of each set, copy numbers were called for RPE1 cells using those  
874 4 probes, and again as each additional probe was added to the set until all 120 probes were used  
875 for the calculation. This procedure was repeated for each of the 50 sets. The sensitivity of copy  
876 number cells was recalculated for the RPE1 cells at each step as well as the theoretical sensitivity  
877 for all GMM components using a model with 5 components (with the upper boundary of the 5th  
878 component being calculated using a 6-component model), and with 3 components where noted.

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

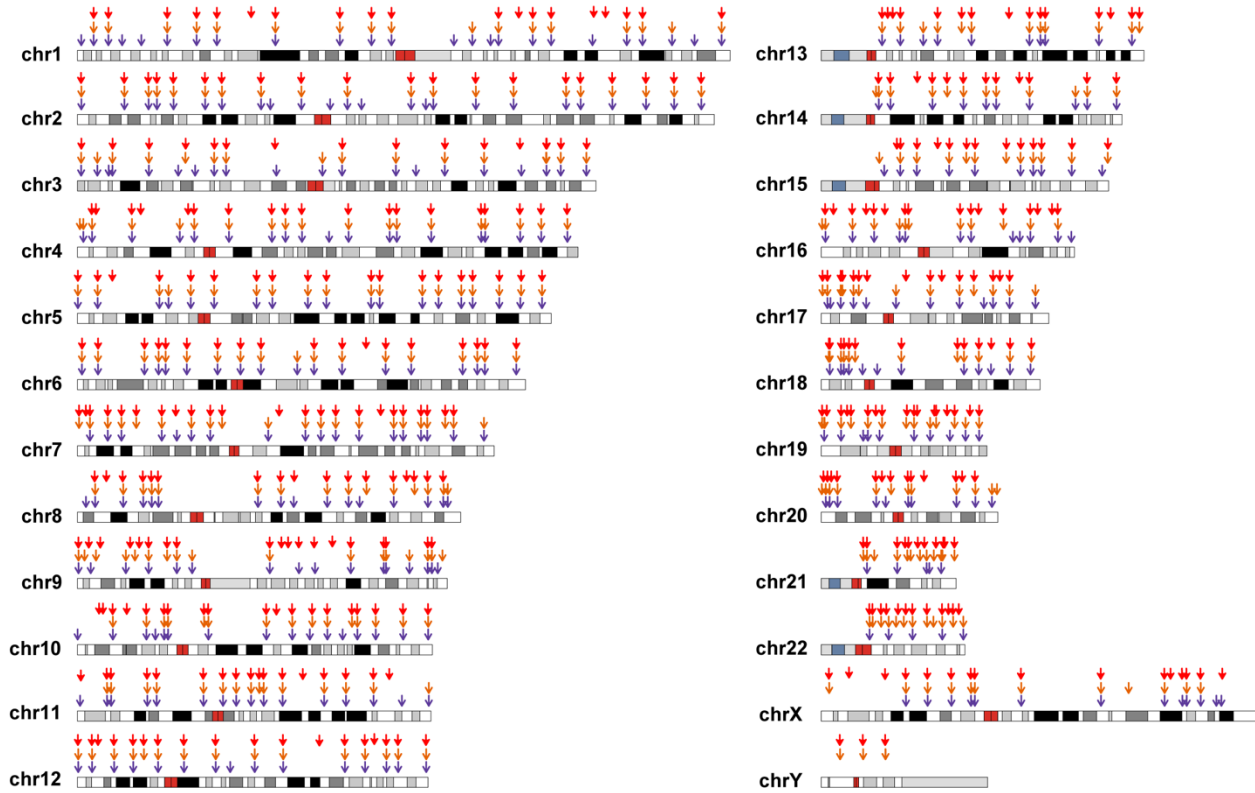
897 **Figure and Table Legends**

**Custom Tapestri Panel Probes**

**Version 3 (CO810)**

**Version 2 (CO610)**

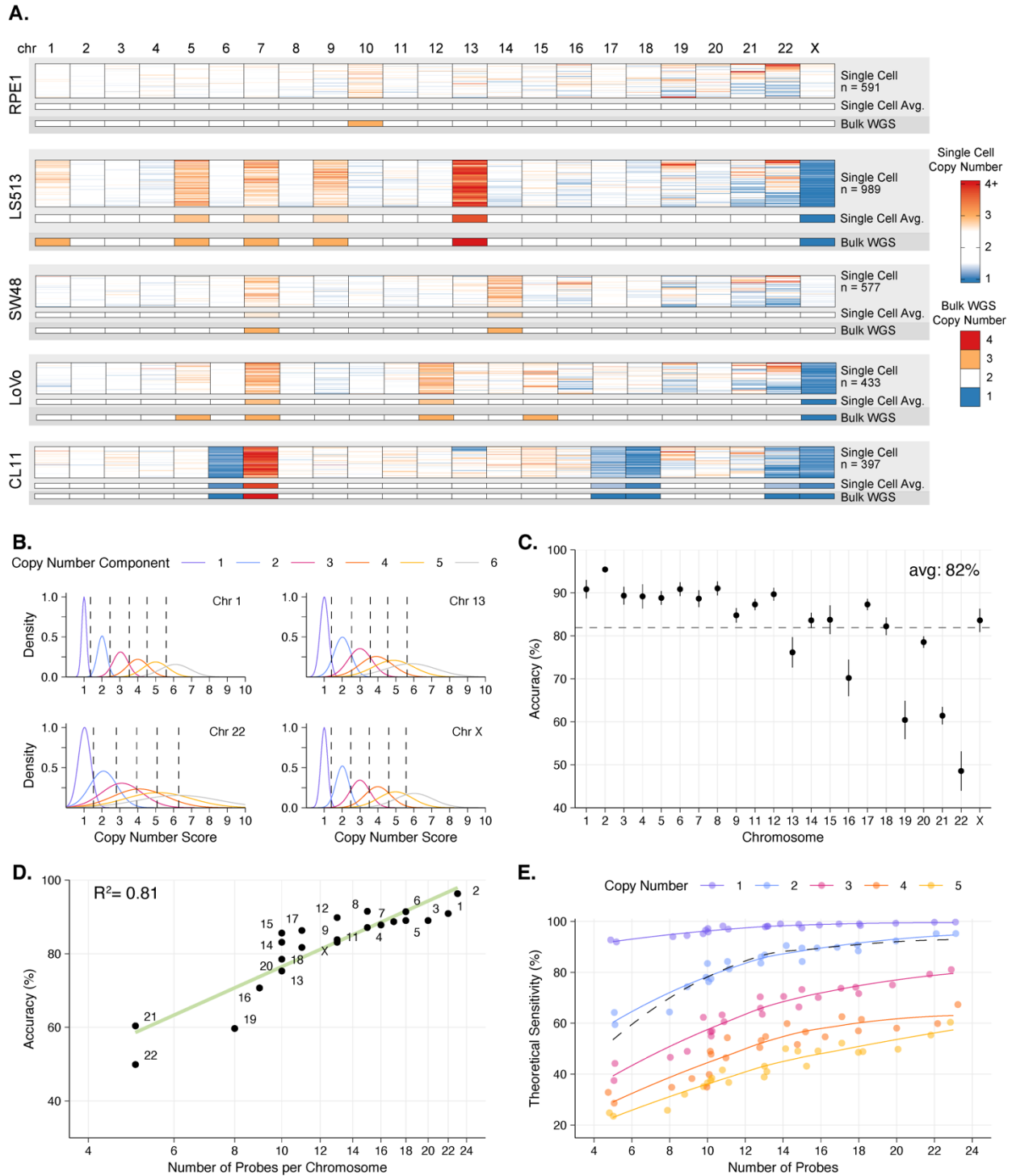
**Version 1 (CO261)**



898

899 **Figure 1**

900 Map of PCR probe locations (arrows) for custom Tapestri panels for KaryoTap Version 1 (CO261)  
901 and Version 2 (CO610) on human genome hg19. Red blocks indicate centromeres, grayscale  
902 blocks indicate the G-band intensity of cytobands, and blue blocks indicate acrocentric  
903 chromosome arms.



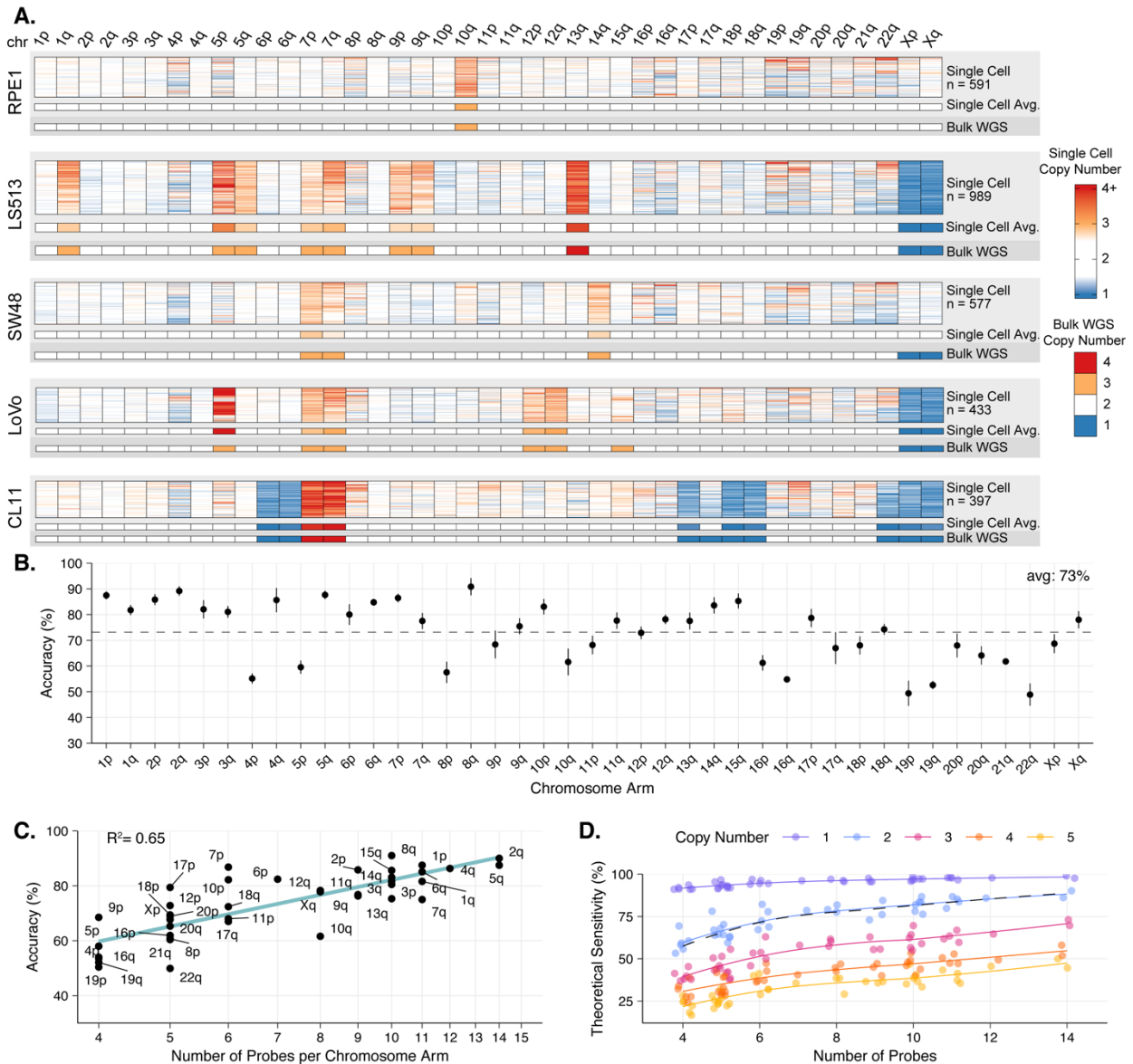
904

905 **Figure 2**

906 Copy number calling for whole chromosomes using KaryoTap panel Version 1. A) Heatmap of

907 copy number scores for each cell-chromosome unit for 5 cell lines using custom Tapestri panel

908 Version 1. Upper blocks indicate copy number scores for each cell, middle blocks indicate  
909 average intensity of upper blocks, and lower blocks indicate copy number from bulk WGS (see  
910 Fig S1A). Half-filled lower blocks indicate chromosome arm-level aneuploidy. Number of cells  
911 included in single-cell blocks is indicated. B) Probability density functions of Gaussian mixture  
912 models (GMM) fit for chromosomes 1, 13, 22, and X using RPE1 cells. Dotted lines indicate  
913 decision boundaries between GMM components. C) 5-fold cross validation of copy number call  
914 accuracy for RPE1 cells by chromosome. Chr10 is omitted. Dot indicates mean accuracy and  
915 lines indicate  $\pm$  mean absolute deviation. Horizontal dotted line indicates average (avg) accuracy  
916 across chromosomes. D) Linear regression of RPE1 copy number call accuracy for each  
917 chromosome on number of probes per chromosome. X-axis is log-scaled to reflect log  
918 transformation of number of probes in regression. E) Theoretical copy number call sensitivity for  
919 each chromosome and copy number level calculated from GMMs. Points are slightly jittered  
920 horizontally to decrease overlapping.



921

922

923

924 Figure 3

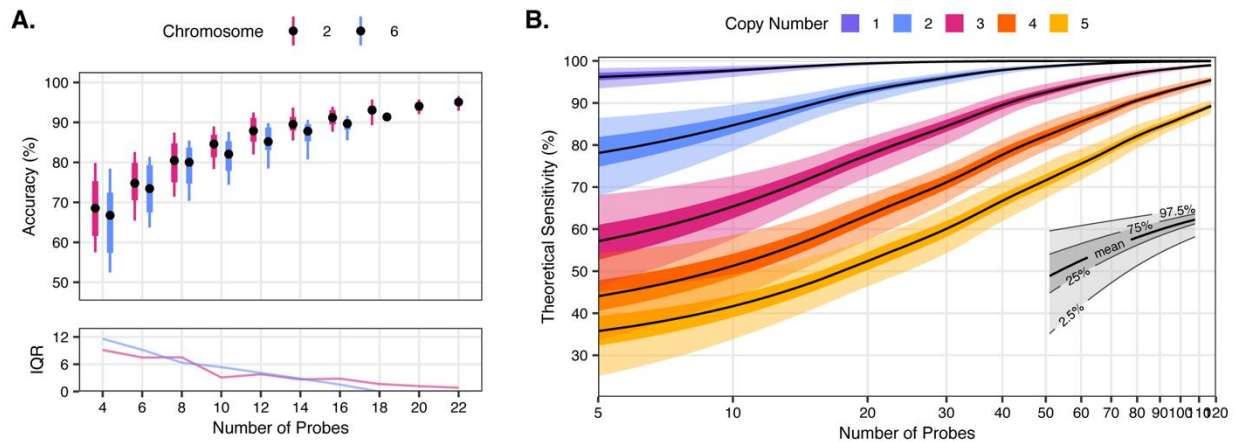
925 Copy number calling for chromosome arms using custom KaryoTap panel Version 1. A) Heatmap

926 of copy number scores for each cell, smoothed across chromosome arms, for five cell lines using

927 custom Tapestri panel Version 1. Upper blocks indicate copy number scores for each cell, middle

928 blocks indicate average intensity of upper blocks, and lower blocks indicate copy number from

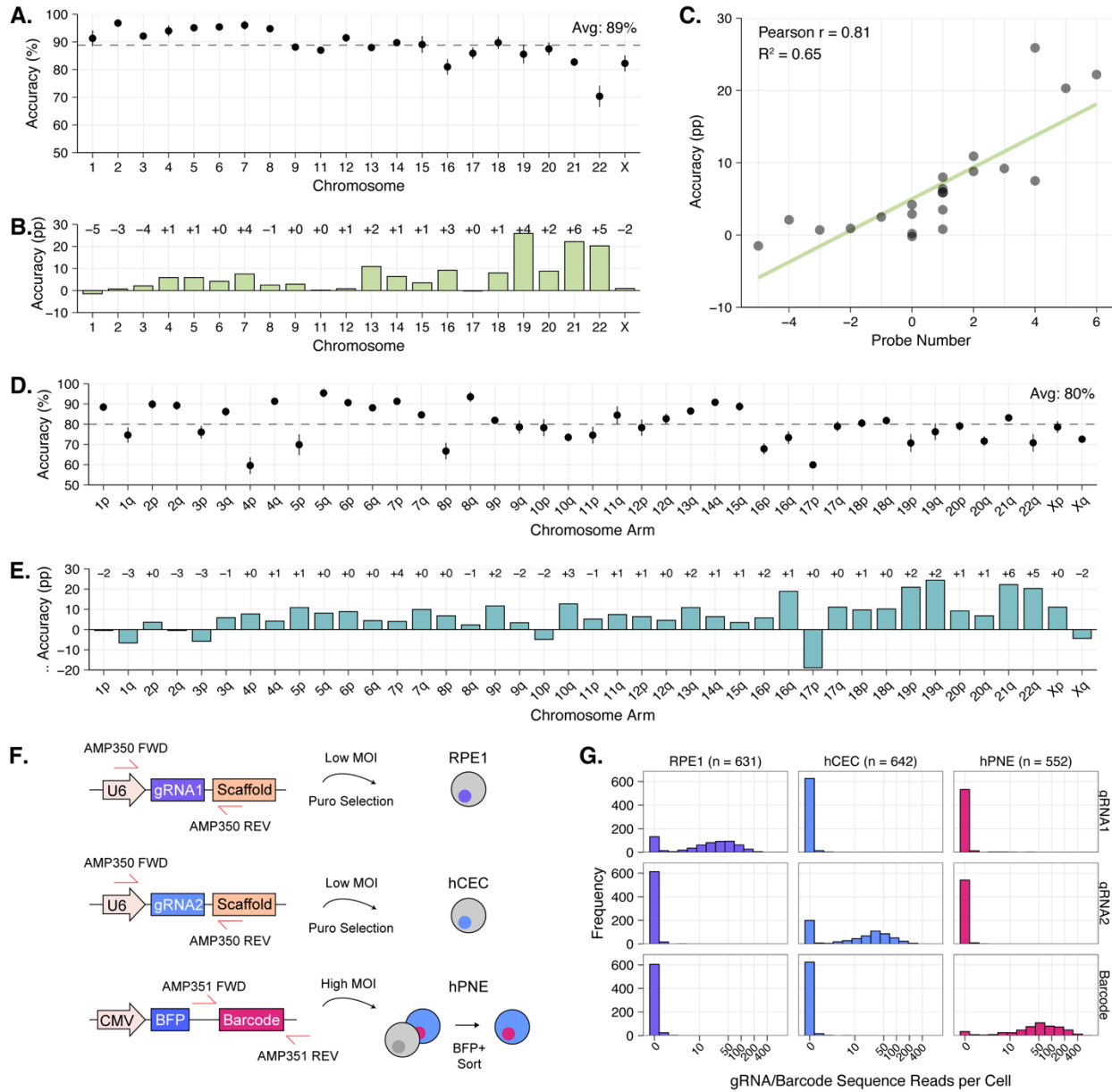
929 bulk WGS (see Fig S1A). B) 5-fold cross validation of copy number call accuracy for RPE1 cells  
930 by chromosome arm. Dot indicates mean accuracy and lines indicate  $\pm$  mean absolute deviation.  
931 Horizontal dotted line indicates average accuracy across chromosome arms. C) Linear regression  
932 of RPE1 copy number call accuracy for each chromosome arm on number of probes per  
933 chromosome arm. X-axis is log-scaled to reflect log transformation of number of probes in  
934 regression. D) Theoretical copy number call sensitivity for each chromosome arm and copy  
935 number level calculated from GMMs. Points are slightly jittered horizontally to decrease  
936 overlapping.



937

#### 938 Figure 4

939 Effects on copy number call accuracy on probe sampling simulations. A) Box plots and inter  
940 quartile range (IQR) of accuracy from 50 probe downsampling simulations for chr2 and chr6.  
941 Boxes encompass middle 50%, whiskers encompass middle 95%, dot indicates median. B)  
942 Theoretical sensitivity of 50 panel simulations. Values for each copy number level were smoothed  
943 by Loess regression. The black line represents the mean, the darker inner shading indicates the  
944 middle 50% of the data, and the lighter outer shading represents the middle 95% of the data.



945

946

947 **Figure 5**

948 Copy number calling and lentiviral barcoding for custom KaryoTap panel Version 2. A) 5-fold cross

949 validation of copy number call accuracy for n=631 RPE1 cells by chromosome using panel

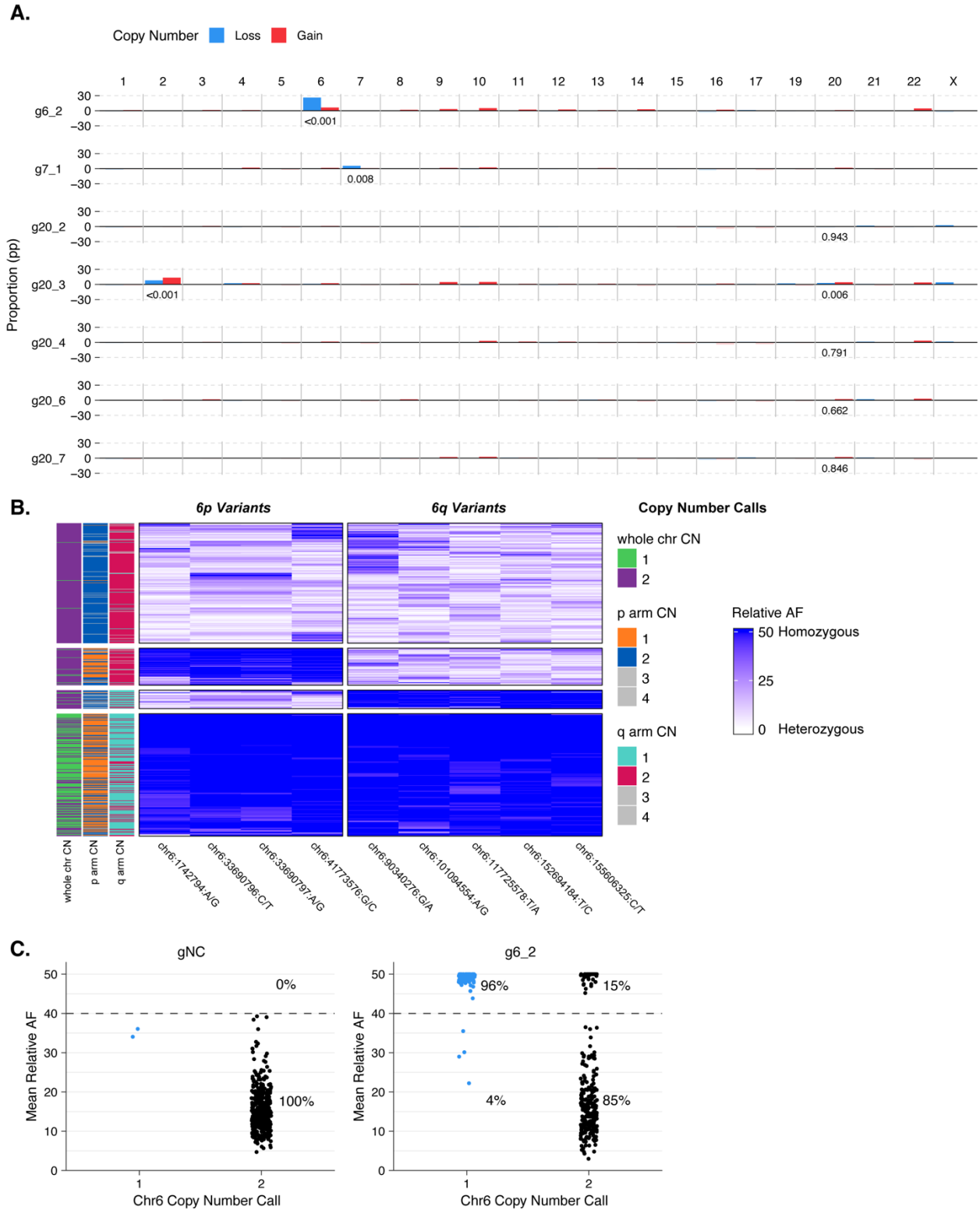
950 Version 2. Chr10 is omitted. Dot indicates mean accuracy and lines indicate  $\pm$  mean absolute

951 deviation. Horizontal dotted line indicates average (avg) accuracy across chromosomes. B)

952 Change in copy number call accuracy by chromosome for RPE1 cells between Version 2 and



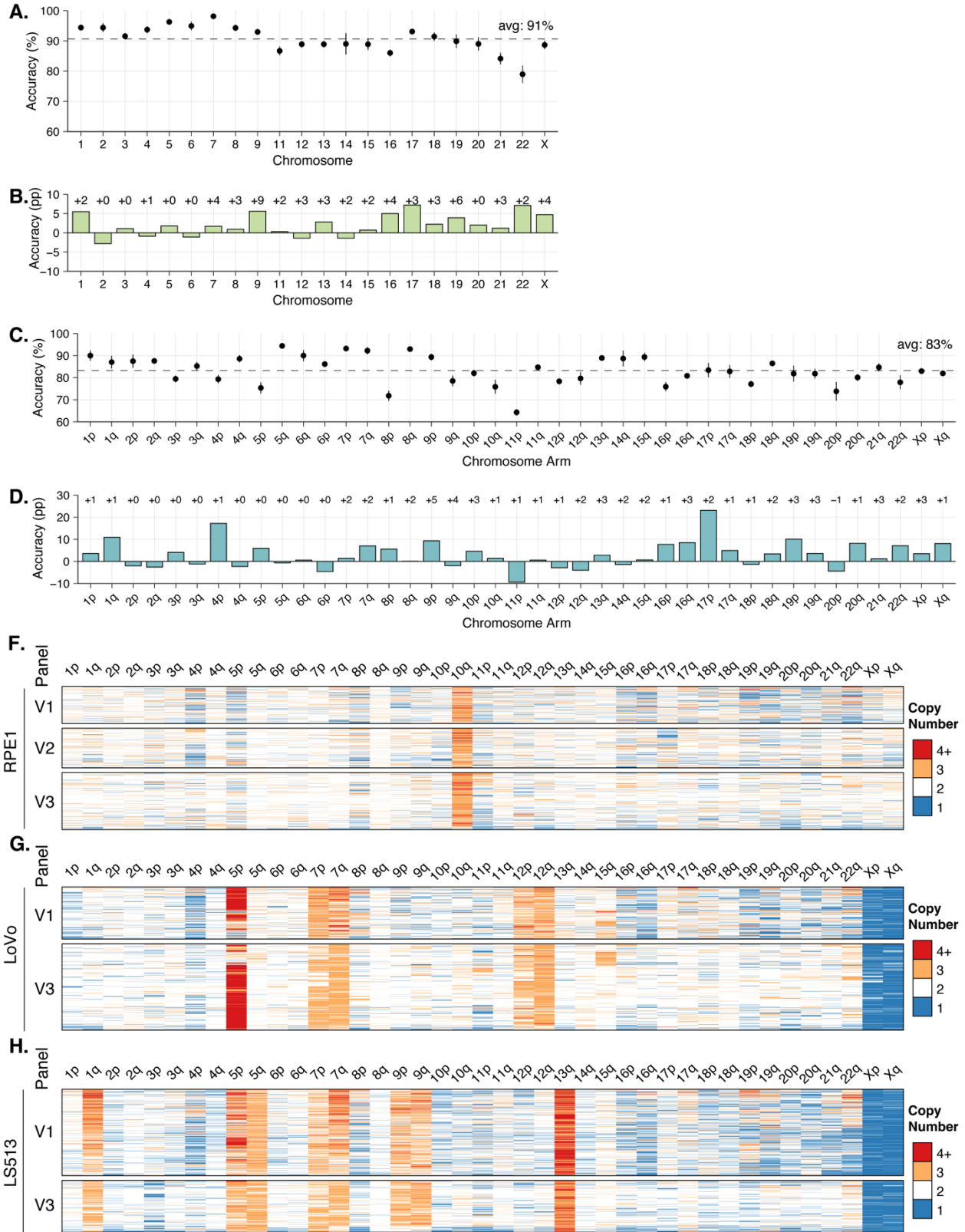
953 Version 1 panels.  $\Delta$  Probe Number is the difference in number of probes targeting a given  
954 chromosome between Version 2 and Version 1. pp: percentage points. C) Linear regression of  
955 the change in copy number call accuracy for RPE1 cells between Version 2 and Version 1 panels  
956 on the change in probe number targeting each chromosome. D) 5-fold cross validation of copy  
957 number call accuracy for RPE1 cells by chromosome arms using custom Tapestri panel Version  
958 2. Dot indicates mean accuracy and lines indicate  $\pm$  mean absolute deviation. Horizontal dotted  
959 line indicates average across chromosome arms. E) Change in copy number call accuracy by  
960 chromosome arm for RPE1 cells between Version 2 and Version 1 panels. F) Plasmid constructs  
961 for lentiviral transduction of RPE1, hCEC D29, and hPNE cell lines. Probe pair AMP350 amplifies  
962 a 253 bp region including a CRISPR gRNA sequence and part of the U6 promoter and F+E  
963 scaffold. Probe pair AMP351 amplifies a 237 bp region including a barcode sequence. G) Number  
964 of reads in each cell that match the expected sequence of gRNA1, gRNA2, or barcodes in 3  
965 transduced cell lines. X-axis is log-scaled. Number of cells of each cell line is indicated.



966

967 **Figure 6**

968 Evaluation of KaryoCreate and Loss of Heterozygosity. A) Evaluation of KaryoCreate technology  
969 by KaryoTap panel Version 2. Bars represent the percentage point (pp) change ( $\Delta$ ) in the  
970 proportion of chromosomal losses (1 copy) and gains (3+ copies), compared to sgNC negative  
971 control (n = 398 cells). Chr18 was omitted due to additional copies of the chromosome in the cell  
972 line. p-values from Fisher's Exact test comparing the proportion of cells with copy number = 2 to  
973 copy number = {1, 3, or 4} (i.e., diploid vs. aneuploid) in each chromosome in each sample to the  
974 corresponding chromosome in the sgNC negative control sample are shown where  $p < 0.1$ .  
975 Additional p-values  $> 0.1$  are reported where relevant. All p-values are reported in Table S2. Bars  
976 with negative values and  $p > 0.1$  have reduced opacity for clarity. B) Heatmap of relative VAFs  
977 for sgChr6-2 treated cells with 1 or 2 copies of chr6 as called by GMM, for 9 originally  
978 heterozygous variants. Relative allele frequencies calculated as the absolute difference between  
979 raw allele frequency and 50%. 0 corresponds to balanced heterozygous alleles and 50  
980 corresponds to fully homozygous alleles. Heatmap rows split by k-means clustering where  $k=4$   
981 and sorted by hierarchical clustering. CN: copy number. C) Mean relative VAFs for sgNC and  
982 sgChr6-2 treated cells with 1 or 2 copies of chr6 as called by GMM. 0%-40% AF indicates  
983 heterozygous haplotype, 40%-50% AF indicates homozygous haplotype.



984

985 Figure 7

986 Copy number calling and lentiviral barcoding for custom KaryoTap panel Version 3. A) 5-fold cross  
987 validation of copy number call accuracy for n=908 RPE1 cells by chromosome using panel  
988 Version 3. Chr10 is omitted. Dot indicates mean accuracy and lines indicate  $\pm$  mean absolute  
989 deviation. Horizontal dotted line indicates average (avg) accuracy across chromosomes. B)  
990 Change in copy number call accuracy by chromosome for RPE1 cells between Version 3 and  
991 Version 2 panels.  $\Delta$  Probe Number is the difference in number of probes targeting a given  
992 chromosome between Version 3 and Version 2. pp: percentage points. C) 5-fold cross validation  
993 of copy number call accuracy for RPE1 cells by chromosome arms using custom Tapestri panel  
994 Version 3. Dot indicates mean accuracy and lines indicate  $\pm$  mean absolute deviation. Horizontal  
995 dotted line indicates average across chromosome arms. D) Change in copy number call accuracy  
996 by chromosome arm for RPE1 cells between Version 3 and Version 2 panels. E-G: Heatmaps of  
997 GMM copy number calls for chromosome arms using panels Version 1, 2, and 3 for cell lines  
998 RPE1, LoVo, and LS513.

999

#### 1000 Supplemental Figure 1

1001 A) Bulk low-pass whole genome sequencing of RPE1, LS513, SW48, LoVo, and CL11 cell lines.  
1002 Red highlight indicates amplification of at least one copy of highlighted segment; blue similarly  
1003 indicates deletion. B) Representative g-banded karyogram of RPE1, indicating additional copy of  
1004 chr10q translocated to the X chromosome (red arrow). C) UMAP projection of top 4 principal  
1005 components of allele frequencies for N=2,986 cells representing 5 cell lines. Clustering was  
1006 performed using the dbscan method. Cells were considered doublets if they were not members  
1007 of the 5 largest clusters. D) PCA plot of first two principal components of mean allele frequencies  
1008 for previously published deep sequencing of RPE1 and the 5 cell lines analyzed by scDNA-seq.  
1009 E) Heatmap of cell-probe copy number values for five cell lines using custom Tapestri panel  
1010 Version 1. Probes are organized by chromosome arm in genomic order.

1011

1012 Supplemental Figure 2

1013 A) k-means clustering of arm-level copy number scores for n=2,986 using Panel Version 1. Color  
1014 annotation indicates which cell line each row belongs to as determined by cell line SNPs. B) k-  
1015 means clustering of arm-level copy number scores for LoVo cells (n=433) using Panel Version 1,  
1016 at k = 2. C) Heatmap of GMM calls for whole chromosomes. D) Chromosome length (in 100  
1017 megabases) vs. accuracy of RPE1 panel Version 1 copy number calls chromosome. Chr10 is  
1018 omitted. Trendline fit by linear regression. E) Heatmap of GMM calls for chromosome arms. using  
1019 Panel Version 1 across 5 cells lines.

1020

1021

1022 Supplemental Figure 3

1023 A) Theoretical copy number call sensitivity from 50 probe downsampling simulations for chr2 and  
1024 chr6 across five copy number levels. Boxes encompass middle 50%, whiskers encompass middle  
1025 95%, dot indicates median. B) Theoretical copy number call sensitivity of 50 panel simulations,  
1026 using a 3 component GMM. Values for each copy number level were smoothed by Loess  
1027 regression. The black line represents the mean, the darker inner shading indicates the middle  
1028 50% of the data, and the lighter outer shading represents the middle 95% of the data.

1029

1030 Supplemental Figure 4

1031 A) Read counts per cell of four probes targeting chrY across 3 cell lines. B) UMAP projection of  
1032 top 2 principal components of allele frequencies for N=2,347 cells representing 3 cell lines.  
1033 Clustering was performed using the dbscan method. Cells were considered doublets if they were  
1034 not members of the 3 largest clusters. C) Heatmap of GMM calls for whole chromosomes across  
1035 3 cell lines using panel Version 2. D) Heatmap of GMM calls for chromosome arms across 3 cell  
1036 lines using panel Version 2. E) Mean theoretical accuracy of panel Version 2 copy number calls  
1037 for each chromosome (left panel) or arm (right panel) at copy number values of 1, 2, 3, 4 and 5,

1038 compared to those from panel Version 1. Dotted red line indicates  $x = y$ . F) Relationship  
1039 between read counts per cell for either the gRNA or DNA Barcode probes (x-axis) and the number  
1040 of reads per cell that match the specific gRNA or DNA Barcode sequences (y-axis). Dotted red  
1041 line indicates  $x = y$ .

1042

#### 1043 Supplemental Figure 5

1044 A) Mean theoretical accuracy of panel Version 3 copy number calls for each chromosome at copy  
1045 number values of 1, 2, 3, 4 and 5, compared to those from panel Version 2. Dotted red line  
1046 indicates  $x = y$ . B) Mean theoretical accuracy of panel Version 3 copy number calls for each  
1047 chromosome arm at copy number values of 1, 2, 3, 4 and 5, compared to those from panel Version  
1048 2.

1049

#### 1050 Supplemental Figure 6

1051 Overview of workflow for copy number analysis using karyotapR R package.

1052

#### 1053 **List of Abbreviations**

1054 BFP Blue florescent protein

1055 CIN Chromosomal instability

1056 CNV Copy number variant

1057 DLP Direct library preparation

1058 DOP-PCR Degenerate-oligonucleotide-primed PCR

1059 FACS Fluorescence associated cell sorting

1060 GMM Gaussian mixture model

1061 gRNA guide RNA

1062 hCECs Human colorectal epithelial cells

1063 hPNEs Human pancreatic nestin-expressing cells

1064 MALBAC Multiple annealing and looping-based amplification cycles

1065 MDA Multiple displacement amplification

1066 MOI Multiplicity of infection

1067 ORF Open reading frame

1068 PCA Principal components analysis

1069 PCR Polymerase chain reaction

1070 RPE1 Retinal pigment epithelial (cells) 1

1071 scDNA-seq Single-cell DNA sequencing

1072 scRNA-seq Single-cell RNA sequencing

1073 sgRNA Single guide RNA

1074 SNP Single-nucleotide polymorphism

1075 SNV Single nucleotide variant

1076 VAF Variant allele frequency

1077 WGA Whole genome amplification

1078 WGS Whole genome sequencing

1079 **Declarations**

1080

1081 **Ethics Approval and Consent to Participate**

1082 Not applicable

1083

1084 **Consent for Publication**

1085 Not applicable

1086

1087 **Availability of Data and Materials**

1088 Sequencing data for Tapestri experiments are available in the SRA repository under NCBI

1089 BioProject accession PRJNA950110, <https://www.ncbi.nlm.nih.gov/bioproject/950110> (56). The



1090 karyotapR package is available on GitHub at <http://github.com/joeymays/karyotapR> (57).The  
1091 source code for karyotapR version 0.1 used for this study is archived on Zenodo under DOI  
1092 <https://doi.org/10.5281/zenodo.8305561>(58). All data analysis scripts used in this study are  
1093 available at <https://github.com/joeymays/karyotap-publication> (59) and are archived on Zenodo  
1094 under DOI <https://doi.org/10.5281/zenodo.8329277>. Tapestri Pipeline output files used in this  
1095 study are available on Zenodo under DOI <https://doi.org/10.5281/zenodo.8305841>(60).

1096

### 1097 **Competing Interests**

1098 TD is on the Scientific Advisory Board of io9 and founder of KaryoVerse Therapeutics.

1099

### 1100 **Funding**

1101 LJH and GRK were supported by NIH R37 CA240765. TD and members of her lab were  
1102 supported by the Cancer Research UK Grand Challenge and the Mark Foundation for Cancer  
1103 Research (C5470/A27144), NIH R37 R37CA248631, R01 R01HG012590, R01 R01DK135089,  
1104 the MRA Young Investigator Award, the Breast Cancer Alliance Young Investigator Award and a  
1105 grant from The National Foundation for Cancer Research.

1106

### 1107 **Authors' contributions**

1108 JCM designed the study, performed wet lab experiments with help from others, performed  
1109 bioinformatics analyses, developed the karyotapR software package, and wrote the manuscript.  
1110 SM performed lentiviral preparation and transductions. NB and AG performed wet lab  
1111 experiments. XZ designed Tapestri panels and performed bioinformatics analyses. JJB  
1112 performed and supported cell culture. GRK prepared barcode vectors. LJH and DF supervised  
1113 research. TD designed the study, edited the manuscript, and supervised research. MK and HMQ  
1114 edited the manuscript. All authors read and approved the final manuscript.

1115

## 1116 **Acknowledgments**

1117 NYU Langone's Genome Technology Center (RRID: SCR\_017929) is supported by the Cancer  
1118 Center Support Grant P30CA016087 at the Laura and Isaac Perlmutter Cancer Center. The  
1119 computational requirements for this work were supported in part by the NYU Langone High  
1120 Performance Computing (HPC) Core's resources and personnel. We thank Paolo Mita, Pan  
1121 Cheng, Lizabeth Katznelson, and Maria Trifas for their helpful suggestions and insights in the  
1122 preparation of this work.

1123

## 1124 **Additional Files**

- 1125 ● AdditionalFile01.csv: Probe Design for Panel Version 1 (CO216)
- 1126 ● AdditionalFile02.csv: Probe Design for Panel Version 2 (CO610)
- 1127 ● AdditionalFile03.csv: Probe Design for Panel Version 3 (CO810)

1128

1129

## 1130 **References**

- 1131 1. Navin N, Krasnitz A, Rodgers L, Cook K, Meth J, Kendall J, et al. Inferring tumor progression  
1132 from genomic heterogeneity. *Genome Res.* 2010 Jan;20(1):68–80.
- 1133 2. Ben-David U, Amon A. Context is everything: aneuploidy in cancer. *Nature Reviews*  
1134 *Genetics.* 2020 Jan;21(1):44–62.
- 1135 3. Weaver BAA, Cleveland DW. Aneuploidy: Instigator and Inhibitor of Tumorigenesis. *Cancer*  
1136 *Research.* 2007 Nov 1;67(21):10103–5.
- 1137 4. Bakhoun SF, Ngo B, Laughney AM, Cavallo JA, Murphy CJ, Ly P, et al. Chromosomal  
1138 instability drives metastasis through a cytosolic DNA response. *Nature.* 2018  
1139 Jan;553(7689):467–72.
- 1140 5. Bakhoun SF, Cantley LC. The Multifaceted Role of Chromosomal Instability in Cancer and  
1141 Its Microenvironment. *Cell.* 2018 Sep 6;174(6):1347–60.
- 1142 6. Ippolito MR, Martis V, Martin S, Tjihuis AE, Hong C, Wardenaar R, et al. Gene copy-number  
1143 changes and chromosomal instability induced by aneuploidy confer resistance to  
1144 chemotherapy. *Developmental Cell.* 2021 Sep 13;56(17):2440-2454.e6.

- 1145 7. Lee AJX, Endesfelder D, Rowan AJ, Walther A, Birkbak NJ, Futreal PA, et al. Chromosomal  
1146 Instability Confers Intrinsic Multidrug Resistance. *Cancer Research*. 2011 Mar  
1147 1;71(5):1858–70.
- 1148 8. Anagnostou V, Smith KN, Forde PM, Niknafs N, Bhattacharya R, White J, et al. Evolution of  
1149 Neoantigen Landscape during Immune Checkpoint Blockade in Non–Small Cell Lung  
1150 Cancer. *Cancer Discovery*. 2017 Mar 5;7(3):264–76.
- 1151 9. Lynch A, Bradford S, Burkard ME. The reckoning of chromosomal instability: past, present,  
1152 future. *Chromosome Res*. 2024 Feb 17;32(1):2.
- 1153 10. Bakker B, Taudt A, Belderbos ME, Porubsky D, Spierings DCJ, de Jong TV, et al. Single-  
1154 cell sequencing reveals karyotype heterogeneity in murine and human malignancies.  
1155 *Genome Biology*. 2016 May 31;17(1):115.
- 1156 11. Baslan T, Hicks J. Unravelling biology and shifting paradigms in cancer with single-cell  
1157 sequencing. *Nat Rev Cancer*. 2017 Sep;17(9):557–69.
- 1158 12. Navin N, Kendall J, Troge J, Andrews P, Rodgers L, McIndoo J, et al. Tumour evolution  
1159 inferred by single-cell sequencing. *Nature*. 2011 Apr;472(7341):90–4.
- 1160 13. Wang Y, Waters J, Leung ML, Unruh A, Roh W, Shi X, et al. Clonal evolution in breast  
1161 cancer revealed by single nucleus genome sequencing. *Nature*. 2014 Aug;512(7513):155–  
1162 60.
- 1163 14. Bosco N, Goldberg A, Zhao X, Mays JC, Cheng P, Johnson AF, et al. KaryoCreate: A  
1164 CRISPR-based technology to study chromosome-specific aneuploidy by targeting human  
1165 centromeres. *Cell [Internet]*. 2023 Apr 18 [cited 2023 Apr 20]; Available from:  
1166 <https://www.sciencedirect.com/science/article/pii/S0092867423003264>
- 1167 15. Truong MA, Cané-Gasull P, de Vries SG, Nijenhuis W, Wardenaar R, Kapitein LC, et al. A  
1168 kinesin-based approach for inducing chromosome-specific mis-segregation in human cells.  
1169 *The EMBO Journal*. 2023 May 15;42(10):e111559.
- 1170 16. Tovini L, Johnson SC, Guscott MA, Andersen AM, Spierings DCJ, Wardenaar R, et al.  
1171 Targeted assembly of ectopic kinetochores to induce chromosome-specific segmental  
1172 aneuploidies. *The EMBO Journal*. 2023 May 15;42(10):e111587.
- 1173 17. Telenius H, Carter NP, Bebb CE, Nordenskjöld M, Ponder BAJ, Tunnacliffe A. Degenerate  
1174 oligonucleotide-primed PCR: General amplification of target DNA by a single degenerate  
1175 primer. *Genomics*. 1992 Jul 1;13(3):718–25.
- 1176 18. Arneson N, Hughes S, Houlston R, Done S. Whole-Genome Amplification by Degenerate  
1177 Oligonucleotide Primed PCR (DOP-PCR). *Cold Spring Harb Protoc*. 2008 Jan  
1178 1;2008(1):pdb.prot4919.
- 1179 19. Dean FB, Hosono S, Fang L, Wu X, Faruqi AF, Bray-Ward P, et al. Comprehensive human  
1180 genome amplification using multiple displacement amplification. *Proceedings of the National  
1181 Academy of Sciences*. 2002 Apr 16;99(8):5261–6.

- 1182 20. Zong C, Lu S, Chapman AR, Xie XS. Genome-Wide Detection of Single-Nucleotide and  
1183 Copy-Number Variations of a Single Human Cell. *Science*. 2012 Dec 21;338(6114):1622–6.
- 1184 21. Mallory XF, Edrisi M, Navin N, Nakhleh L. Methods for copy number aberration detection  
1185 from single-cell DNA-sequencing data. *Genome Biology*. 2020 Aug 17;21(1):208.
- 1186 22. Zahn H, Steif A, Laks E, Eirew P, VanInsberghe M, Shah SP, et al. Scalable whole-genome  
1187 single-cell library preparation without preamplification. *Nat Methods*. 2017 Feb;14(2):167–  
1188 73.
- 1189 23. Yin Y, Jiang Y, Lam KWG, Berletch JB, Disteche CM, Noble WS, et al. High-Throughput  
1190 Single-Cell Sequencing with Linear Amplification. *Molecular Cell* [Internet]. 2019 Sep 5  
1191 [cited 2019 Sep 6]; Available from:  
1192 <http://www.sciencedirect.com/science/article/pii/S1097276519306185>
- 1193 24. Laks E, McPherson A, Zahn H, Lai D, Steif A, Brimhall J, et al. Clonal Decomposition and  
1194 DNA Replication States Defined by Scaled Single-Cell Genome Sequencing. *Cell*. 2019 Nov  
1195 14;179(5):1207-1221.e22.
- 1196 25. Evrony GD, Hinch AG, Luo C. Applications of Single-Cell DNA Sequencing. *Annu Rev*  
1197 *Genomics Hum Genet*. 2021 Aug 31;22:171–97.
- 1198 26. Mission Bio. Mission Bio Tapestri. 2020 [cited 2020 May 7]. Tapestri: The Precision  
1199 Genomics Platform. Available from: <https://missionbio.com/tapestri/>
- 1200 27. Gawad C, Koh W, Quake SR. Single-cell genome sequencing: current state of the science.  
1201 *Nat Rev Genet*. 2016 Mar;17(3):175–88.
- 1202 28. Crasta K, Ganem NJ, Dagher R, Lantermann AB, Ivanova EV, Pan Y, et al. DNA breaks  
1203 and chromosome pulverization from errors in mitosis. *Nature*. 2012 Feb;482(7383):53–8.
- 1204 29. Reijns MAM, Parry DA, Williams TC, Nadeu F, Hindshaw RL, Rios Szwed DO, et al.  
1205 Signatures of TOP1 transcription-associated mutagenesis in cancer and germline. *Nature*.  
1206 2022 Feb;602(7898):623–31.
- 1207 30. Aird D, Ross MG, Chen WS, Danielsson M, Fennell T, Russ C, et al. Analyzing and  
1208 minimizing PCR amplification bias in Illumina sequencing libraries. *Genome Biology*. 2011  
1209 Feb 21;12(2):R18.
- 1210 31. Wang K, Kumar T, Wang J, Minussi DC, Sei E, Li J, et al. Archival single-cell genomics  
1211 reveals persistent subclones during DCIS progression. *Cell*. 2023 Aug 31;186(18):3968-  
1212 3982.e15.
- 1213 32. Jothilakshmi S, Gudivada VN. Chapter 10 - Large Scale Data Enabled Evolution of Spoken  
1214 Language Research and Applications. In: Gudivada VN, Raghavan VV, Govindaraju V, Rao  
1215 CR, editors. *Handbook of Statistics* [Internet]. Elsevier; 2016 [cited 2023 Aug 22]. p. 301–  
1216 40. (Cognitive Computing: Theory and Applications; vol. 35). Available from:  
1217 <https://www.sciencedirect.com/science/article/pii/S0169716116300463>

- 1218 33. Loh PR, Genovese G, Handsaker RE, Finucane HK, Reshef YA, Palamara PF, et al.  
1219 Insights into clonal haematopoiesis from 8,342 mosaic chromosomal alterations. *Nature*.  
1220 2018 Jul;559(7714):350–5.
- 1221 34. Melcher R, Hartmann E, Zopf W, Herterich S, Wilke P, Müller L, et al. LOH and copy neutral  
1222 LOH (cnLOH) act as alternative mechanism in sporadic colorectal cancers with  
1223 chromosomal and microsatellite instability. *Carcinogenesis*. 2011 Apr 1;32(4):636–42.
- 1224 35. Knudson AG. Mutation and Cancer: Statistical Study of Retinoblastoma. *Proceedings of the*  
1225 *National Academy of Sciences*. 1971 Apr;68(4):820–3.
- 1226 36. Shah JB, Poeschl D, Wubbenhorst B, Fan M, Pluta J, D’Andrea K, et al. Analysis of  
1227 matched primary and recurrent BRCA1/2 mutation-associated tumors identifies recurrence-  
1228 specific drivers. *Nat Commun*. 2022 Nov 7;13(1):6728.
- 1229 37. Dong X, Zhang L, Hao X, Wang T, Vijg J. SCCNV: A Software Tool for Identifying Copy  
1230 Number Variation From Single-Cell Whole-Genome Sequencing. *Frontiers in Genetics*  
1231 [Internet]. 2020 [cited 2023 Jun 2];11. Available from:  
1232 <https://www.frontiersin.org/articles/10.3389/fgene.2020.505441>
- 1233 38. Ryland GL, Doyle MA, Goode D, Boyle SE, Choong DYH, Rowley SM, et al. Loss of  
1234 heterozygosity: what is it good for? *BMC Medical Genomics*. 2015 Aug 1;8(1):45.
- 1235 39. Shalem O, Sanjana NE, Zhang F. High-throughput functional genomics using CRISPR–  
1236 Cas9. *Nat Rev Genet*. 2015 May;16(5):299–311.
- 1237 40. Ly P, Eskiocak U, Kim SB, Roig AI, Hight SK, Lulla DR, et al. Characterization of Aneuploid  
1238 Populations with Trisomy 7 and 20 Derived from Diploid Human Colonic Epithelial Cells.  
1239 *Neoplasia*. 2011 Apr 1;13(4):348-IN17.
- 1240 41. Sherry ST, Ward MH, Kholodov M, Baker J, Phan L, Smigielski EM, et al. dbSNP: the NCBI  
1241 database of genetic variation. *Nucleic Acids Res*. 2001 Jan 1;29(1):308–11.
- 1242 42. Nassar LR, Barber GP, Benet-Pagès A, Casper J, Clawson H, Diekhans M, et al. The  
1243 UCSC Genome Browser database: 2023 update. *Nucleic Acids Res*. 2023 Jan  
1244 6;51(D1):D1188–95.
- 1245 43. Gel B, Serra E. karyoploteR: an R/Bioconductor package to plot customizable genomes  
1246 displaying arbitrary data. *Bioinformatics*. 2017 Oct 1;33(19):3088–90.
- 1247 44. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data.  
1248 *Bioinformatics*. 2014 Aug 1;30(15):2114–20.
- 1249 45. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.  
1250 2013.
- 1251 46. Kuilman T, Velds A, Kemper K, Ranzani M, Bombardelli L, Hoogstraat M, et al. CopywriteR:  
1252 DNA copy number detection from off-target sequence data. *Genome Biol*. 2015 Feb  
1253 27;16:49.

- 1254 47. Chen B, Gilbert LA, Cimini BA, Schnitzbauer J, Zhang W, Li GW, et al. Dynamic Imaging of  
1255 Genomic Loci in Living Human Cells by an Optimized CRISPR/Cas System. *Cell*. 2013 Dec  
1256 19;155(7):1479–91.
- 1257 48. Ran FA, Hsu PD, Wright J, Agarwala V, Scott DA, Zhang F. Genome engineering using the  
1258 CRISPR-Cas9 system. *Nat Protoc*. 2013 Nov;8(11):2281–308.
- 1259 49. Sack LM, Davoli T, Li MZ, Li Y, Xu Q, Naxerova K, et al. Profound Tissue Specificity in  
1260 Proliferation Control Underlies Cancer Drivers and Aneuploidy Patterns. *Cell*. 2018 Apr  
1261 5;173(2):499-514.e23.
- 1262 50. Pagès H, Aboyou P, Gentleman R, DebRoy S. Biostrings: Efficient manipulation of  
1263 biological strings [Internet]. 2022. Available from:  
1264 <https://bioconductor.org/packages/Biostrings>
- 1265 51. Morgan M, Pagès H, Obenchain V, Hayden N. Rsamtools: Binary alignment (BAM), FASTA,  
1266 variant call (BCF), and tabix file import [Internet]. 2022. Available from:  
1267 <https://bioconductor.org/packages/Rsamtools>
- 1268 52. Reijns MAM. WT1\_ancestral\_culture [Internet]. SRA; 2021 [cited 2023 Jan 18]. Available  
1269 from: <https://www.ncbi.nlm.nih.gov/sra/?term=ERR7477340>
- 1270 53. GATK [Internet]. [cited 2023 Mar 21]. Available from: <https://gatk.broadinstitute.org/hc/en-us>
- 1271 54. Knaus BJ, Grünwald NJ. vcfr: a package to manipulate and visualize variant call format data  
1272 in R. *Molecular Ecology Resources*. 2017;17(1):44–53.
- 1273 55. Delignette-Muller ML, Dutang C. fitdistrplus: An R Package for Fitting Distributions. *Journal*  
1274 *of Statistical Software*. 2015 Mar 20;64:1–34.
- 1275 56. Mays J, Davoli T. ID 950110 - BioProject - NCBI [Internet]. 2023 [cited 2023 Sep 12].  
1276 Available from: <https://www.ncbi.nlm.nih.gov/bioproject/950110>
- 1277 57. Mays JC. karyotapR: CNV Analysis for Tapestry [Internet]. 2023 [cited 2023 Sep 8].  
1278 Available from: <https://github.com/joeymays/karyotapR>
- 1279 58. Mays J. joeymays/karyotapR: Version 0.1 [Internet]. Zenodo; 2023 [cited 2023 Sep 8].  
1280 Available from: <https://zenodo.org/record/8305561>
- 1281 59. Mays J. KaryoTap Publication Codebook [Internet]. 2023 [cited 2023 Sep 8]. Available from:  
1282 <https://github.com/joeymays/karyotap-publication>
- 1283 60. Mays J, Davoli T. KaryoTap Enables Aneuploidy Detection in Thousands of Single Human  
1284 Cells [Internet]. Zenodo; 2023 [cited 2023 Sep 8]. Available from:  
1285 <https://zenodo.org/record/8305841>
- 1286