

# 1 **PaSD-qc: Quality control for single cell whole-genome** 2 **sequencing data using power spectral density estimation**

3 Maxwell A. Sherman<sup>1</sup>, Alison R. Barton<sup>1</sup>, Michael Lodato<sup>2,3</sup>, Carl Vitzthum<sup>1</sup>, Michael E.  
4 Coulter<sup>2,3</sup>, Christopher A. Walsh<sup>2,3</sup>, and Peter J. Park<sup>1,4,\*</sup>

5 <sup>1</sup>Department of Biomedical Informatics, Harvard Medical School, 10 Shattuck St, Boston, MA 02115.

6 <sup>2</sup>Division of Genetics and Genomics, Boston Children's Hospital, 300 Longwood Avenue, Boston, MA 02115.

7 <sup>3</sup>Howard Hughes Medical Institute, Boston Children's Hospital, 300 Longwood Avenue, Boston MA, 02115.

8 <sup>4</sup>Ludwig Center at Harvard, 200 Longwood Ave, Boston, MA 02115.

9 \*To whom correspondence should be addressed (peter\_park@hms.harvard.edu)

10

11

12

13

14

15

16

17

18

19

20 **Abstract:**

21 Single cell whole-genome sequencing (scWGS) is providing novel insights into the nature of genetic  
22 heterogeneity in normal and diseased cells. However, scWGS introduces DNA amplification-related  
23 biases that can confound downstream analysis. Here we present a statistical method, with an  
24 accompanying package PaSD-qc (Power Spectral Density-qc), that evaluates the quality of single cell  
25 libraries. It uses a modified power spectral density to assess amplification uniformity, amplicon size  
26 distribution, autocovariance, and inter-sample consistency as well as identifies aberrantly amplified  
27 chromosomes. We demonstrate the usefulness of this tool in evaluating scWGS protocols and in selecting  
28 high-quality libraries from low-coverage data for deep sequencing.

29 **Keywords:**

30 Single cell whole-genome sequencing, data quality control, statistical signal processing

31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46

## 47 **Background:**

48 Whole-genome DNA sequencing of single cells (scWGS) has recently been made possible by the  
49 introduction of single cell amplification methods. Multiple displacement amplification (MDA) employs a  
50 highly processive polymerase which can synthesize new molecules (amplicons) of ~10-100 kb. High-  
51 quality MDA-derived data show that >90% of the human genome is amplified and 40-60% can be  
52 covered at >30X when the average depth is ~40-50X [1]. Copy number variations identified from low-  
53 coverage (<5X) MDA data have been used to elucidate tumor evolution [2] and to profile mosaic copy  
54 number variation [3]. With the decrease on cost of deep whole-genome sequencing, more recently, high-  
55 coverage (>30X) MDA data have allowed detection of somatic single nucleotide variants in the human  
56 brain [4]. Another protocol is Multiple annealing and looping based amplification cycles (MALBAC),  
57 which amplifies the genome in ~0.3-5 kb fragments and can cover ~50-90% of the human genome [5]. It  
58 has recently been proposed as a method for screening in-vitro fertilized embryos for genetic abnormalities  
59 prior to implantation [6, 7]. A third method based on DOP-PCR can amplify ~10% of the genome and is  
60 suitable for copy number variation detection but not single nucleotide variant detection [8].

61 All scWGS amplification methods induce biases and artifacts. These include non-uniform read  
62 depth that can appear as copy number aberrations, under and over amplification of entire chromosomes,  
63 uneven amplification of the two alleles, and correlation of features at the amplicon scale (e.g. ~10-100 kb  
64 for MDA) [9, 10], as well as single nucleotide and indel mutations and random ligation of fragments that  
65 are hard to distinguish from inversions. These biases fluctuate depending on the exact amplification  
66 protocol used and the state of the isolated cell (Figure 1A). For example, heat lysis during DNA  
67 extraction can increase the rate of artefactual C>T mutations compared to alkaline lysis [11], and cells in  
68 the G2/M phase amplify more uniformly than cells in the G1/G0 phase [12]. These biases in the data can  
69 affect the accuracy of variants detected in downstream analysis, and new protocols are frequently  
70 proposed claiming to mitigate these biases and provide superior variant detection [13, 14, 15]. It is thus  
71 important to characterize the biases computationally and assess the quality of single cell data.

72           Despite the growing popularity of scWGS, few methods exist to perform this evaluation, and the  
73 few that do are almost exclusively concerned with estimating the uniformity of amplification. This itself  
74 is a non-trivial task because the true amplification process is masked by non-unique mappability, locus  
75 dropout due to amplification failure, or sampling bias during sequencing; additionally, read depth is  
76 highly correlated at positions spanned by the same amplicon. Current methods fail to account for these  
77 challenges. For example, several methods estimate read depth variance by binning reads [15, 16]. Such  
78 methods evaluate dispersion at a fixed genomic scale (the bin size), which fails to capture the correlation  
79 patterns of scWGS; resolving this requires re-binning at many scales, which is time-intensive and  
80 computationally expensive. More recently, an autocovariance (ACF) method has been proposed [10]. In  
81 theory, ACF is an appealing choice to capture the patterns in scWGS data because it measures  
82 correlations between observations within a dataset; however, in practice algorithms to estimate the ACF  
83 cannot easily be modified to account for regions of low mappability or locus dropout. Additionally, no  
84 standard implementations of these tools are available for incorporation into an scWGS pipeline.

85           Here, we introduce a suite of tools to comprehensively measure scWGS data quality, in a package  
86 called PaSD-qc (Power Spectral Density-qc, pronounced “passed-qc”). Using techniques from digital  
87 signal processing to estimate the power spectral density (PSD) of a sample and correct for observation  
88 gaps due to non-unique mappability, assembly gaps, and locus dropout without the need for binning,  
89 PaSD-qc provides a robust assessment of amplification uniformity at all genomic scales simultaneously.  
90 Because our estimation method accounts for the uneven spacing of the data while concurrently reducing  
91 background noise, the PSD can be leveraged to obtain more accurate estimates of variance and  
92 autocovariance than other methods; to identify chromosomes which may be copy-aberrant due to  
93 amplification failure in a principled way; and to compare quality across jointly analyzed samples even at  
94 very low coverage ( $<0.1X$ ). Furthermore, our statistical method can estimate the full distribution of  
95 amplicon sizes in a sample, which has not previously been possible. PaSD-qc can easily be incorporated  
96 into existing pipelines and by default summarizes the quality and properties of each sample in an

97 interactive HTML report. We use the tool to profile several different scWGS protocols, compare different  
98 samples from the same protocol, and select high-quality libraries from an initial set of low coverage  
99 (<1X) data for full-depth sequencing.

100

## 101 **Results:**

### 102 **Characterizing the spatial correlation structure induced by whole-genome amplification**

103 Figure 1B provides an overview of PaSD-qc, and precise details of the algorithm are described in  
104 Methods. In brief, to mitigate issues of mappability, locus dropout, and sequencing bias, we extract read  
105 depth only at uniquely mappable positions covered by at least one read. The resulting signal is a time  
106 series (indexed by genomic position) with highly unevenly spaced observations. To infer the correlation  
107 patterns within this series, we apply the Lomb-Scargle algorithm [17, 18] to estimate the power spectral  
108 density (PSD) of the series. This method is one of the few which are capable of accurately analyzing  
109 correlation patterns of unevenly spaced time series data. We additionally apply a Welch correction [19] to  
110 minimize the noise of power spectral density estimation.

111         It is reasonable to ask (and was asked in [15]) whether the PSD is an appropriate approach given  
112 that it traditionally identifies periodic features when read depth is naturally an aperiodic signal. In fact, the  
113 PSD is mathematically equivalent to autocovariance, and for aperiodic signals, the PSD exactly estimates  
114 the variance of the generating process. See Supplemental Information (SI) for details. Thus, the smooth  
115 curves which result from our estimation method provide direct insight into the variance of scWGS  
116 amplification protocols at all length scales simultaneously.

117         Illustrative examples of a bulk sample, an MDA sample, and MALBAC sample are shown in  
118 Figure 2A. Below a genomic scale of ~1 kb, the samples show a characteristic pattern arising from  
119 paired-end sequencing. For a read pair with insert size  $k$  starting at position  $t$ , there will be an increase in

120 signal at  $x_t$  and  $x_{t+k}$  and a decrease in signal between the two reads. This results in periodicity at small  
121 genomic scales with the strongest periodicity at the mode of the insert size distribution (350 bp for the  
122 bulk sample shown). In fact, at small genomic scales, the PSD closely resembles the distribution of insert  
123 sizes in a sample (Figure S1). Above a genomic scale of  $\sim 1$  kb, the bulk sample is virtually flat with low  
124 amplitude, indicating that, as expected, the coverage profile from bulk sequencing has low-variance and  
125 has no large-scale correlation structure. The slight increase in the PSD at scales  $>100$  kb is an edge effect  
126 of the Welch correction. This edge effect is removed from scWGS PSDs by using an idealized bulk  
127 sample as a baseline (see Methods).

128 The MDA and MALBAC curves have a more complex shape above the pair-end scale. To  
129 interpret these curves, consider an amplicon of length  $h$  starting at position  $t$ . The read depth signal  $x_t$   
130 will be correlated with  $x_{t+i}$  for  $i < h$ . How often a correlation at length  $i$  is observed depends on the  
131 number of amplicons with length  $h \geq i$ . If  $i$  is less than the smallest amplicon, then read depth  $x_t$  and  
132  $x_{t+i}$  will almost always be correlated, resulting in small local variance and thus a lower amplitude PSD at  
133 sub-amplicon scales. For length  $i$  greater than the largest amplicon,  $x_t$  and  $x_{t+i}$  are necessarily  
134 independent, resulting in a higher amplitude PSD at supra-amplicon scales, reflecting the unevenness of  
135 the amplification. The PSD will smoothly transition from the sub- to supra-amplicon variances precisely  
136 following the cumulative distribution of amplicon sizes. These patterns are apparent in Figure 2A. The  
137 MDA curve rises from  $\sim 5$ -100 kb and the MALBAC curve rises from  $\sim 1$ -5 kb, consistent with expected  
138 amplicon sizes for these protocols. Additionally, the supra-amplicon variance of the MALBAC library is  
139 lower than the supra-amplicon variance of the MDA library while the opposite is true of the sub-amplicon  
140 variances, reflecting that MALBAC provides more consistent amplification at positions far apart but that  
141 MDA is locally more uniform since two positions close together are likely to be spanned by a single  
142 amplicon.

143 We are not the first to propose power spectral density estimation as a uniformity measure.  
144 However, prior estimation procedures [5, 13] require binning the data into 1 kb bins and do not take steps

145 to reduce background noise. This results in an inferior PSD estimate where resolution is limited to a  
146 minimum genomic scale of 2 kb (since the Nyquist frequency is  $5 \times 10^{-4}$ ), and fine scale differences  
147 between samples are obscured by the high level of background noise (Figure 2B). Additionally, the PSD  
148 was criticized as lacking reproducibility since a Fourier transform may not be stable in regions of zero  
149 read depth and of low mappability [15]. As stated before, PaSD-qc corrects for these regions, resulting in  
150 highly reproducible estimates (Figure S2).

### 151 **Estimating the distribution of amplicon sizes in scWGS data**

152 Since the dynamic region of the scWGS PSD curve reflects the cumulative distribution of amplicon sizes,  
153 this distribution can be estimated by fitting a properly scaled probability function to the PSD. The error  
154 function (erf) provides a particularly good fit (Figure 3A) and defines a density over the log amplicon  
155 sizes of the form  $\mathcal{N}(\mu, \sigma^2)$  (Figure 3B) where  $\mu$  and  $\sigma$  are parameters estimated from the erf curve. In  
156 standard coordinates, the distributions are skewed with heavy tails extending into larger genomic ranges  
157 (Figure 3C). To confirm the accuracy of this estimation, we simulated an idealized amplification process  
158 on the p-arm of chromosome 3 using the amplicon size distribution estimated from the MDA curve as the  
159 generative model (see Methods). The resulting read depth signal was then analyzed using the PaSD-qc  
160 algorithm. Figure 3D shows the results of the simulation (red curve) along with the true estimate (green  
161 curve); Figure S3 shows the simulated curve for the MALBAC sample. In theory, any properly scaled and  
162 shifted sigmoidal cumulative distribution function can be used for this density estimation. We  
163 additionally tested the logistic distribution and gamma distribution as possible candidates, but found that  
164 the erf function produced the most consistent estimates (Figure S3).

165 The median and percentiles of the amplicon sizes per sample are inferred using Monte Carlo  
166 simulation (see Methods). For the MDA sample, the median amplicon size is 16.7 kb and 95% of all  
167 amplicons fall between 3.3 and 89.0 kb in size; for the MALBAC sample, the median amplicon size is 2.5  
168 kb and 95% of all amplicons fall between 1.1 and 5.6 kb in size (Figure 3C). We additionally profiled the  
169 amplicon distribution of 35 samples from [4] and 14 samples from [20] and found that while distributions

170 are consistent between samples amplified with the same protocol, they are divergent between different  
171 protocols (Figure 3E); samples using the Qiagen REPLI-g Single Cell Kit with heat lysis [20] have a  
172 smaller median amplicon size than samples using Epicenter RepliPhi Phi-29 with alkaline lysis [4] ( $17.5$   
173  $\pm 1.3$  kb vs.  $5.9 \pm 0.5$  kb, p-value:  $1.2e-9$  by Kolmogorov-Smirnov test). We profiled 4 samples also  
174 profiled in that study and found that our estimated median amplicon size was consistent with their  
175 characteristic length scale estimate (Table S1). Although the characteristic length scale of correlation has  
176 been calculated before [10], no other method estimates the full distribution of amplicon sizes in scWGS  
177 data.

### 178 **Comparison to existing scWGS quality control metrics**

179 The autocovariance function (ACF) of scWGS data has previously been proposed as a quality metric.  
180 While the ACF can be calculated directly from unevenly spaced time series data in theory, no  
181 computationally efficient algorithm exists to perform the estimation, and implementations are either time  
182 intensive, memory intensive, or both. Additionally, the statistical power at each lag varies and no  
183 theoretical results exist on the consistency of the unevenly spaced ACF estimator. However, it is possible  
184 and theoretically justified to calculate the ACF from the PSD (see SI). PaSD-qc implements an efficient  
185 algorithm based on this principle (see Methods).

186 To compare the performance of the PaSD-qc ACF against the directly calculated estimate, we  
187 analyzed all 16 single cell samples from the “1465” individual in [4] using both methods (Figure 4A).  
188 These samples were pair-end sequenced with an average insert size of 350 bp. The PaSD-qc ACF  
189 estimate consistently identifies the peak in correlation expected at this scale; the direct estimation fails to  
190 capture this feature. Additionally, the autocorrelation should oscillate around zero beyond the largest  
191 amplicon size. While this behavior is present in the PaSD-qc ACF, the direct estimation remains positive  
192 beyond 1 mb, a genomic scale far larger than the upper amplicon size limit of the Phi-29 polymerase used  
193 in MDA. This empirically demonstrates the potential inaccuracy of directly calculating the ACF from  
194 highly unevenly spaced observations and illustrates how PaSD-qc surmounts this limitation.



195           Additionally, the ACF at lag zero (equivalently the integral of the PSD) provides an estimate of  
196 the overall variance. This dispersion estimate outperforms the other commonly used dispersion estimate,  
197 median absolute pairwise difference (MAPD) [16, 21]. MAPD is calculated by binning the read depth  
198 signal into fixed-width bins, calculating the normalized copy number in each bin, and taking the median  
199 of the pair-wise differences between all neighboring bins. We calculated MAPD scores at a range of bin  
200 sizes (Figure 4B) and the PaSD-qc PSD estimates (Figure 4C) for all “1465” and “4643” samples from  
201 [4]. Both reveal “1465” samples have higher supra-amplicon variance than “4643” samples. However,  
202 calculating MAPD even at a single bin size is computationally intensive; as such, it is usually calculated  
203 only for a single bin size, often 50 kb. At this scale, MAPD fails to distinguish a difference between the  
204 sets of samples (Figure 4D, p-value: 0.11 by Kolmogorov-Smirnov test). However, the PaSD-qc variance  
205 readily discriminates the two sets (Figure 4E, p-value: 1.7e-6 by Kolmogorov-Smirnov test).

#### 206 **Identification of chromosomes with copy number altered due to aberrant amplification**

207 The close relationship between a power spectral density estimate and a normal distribution [22] permits  
208 the calculation of a statistical distance measure, the symmetric Kullback-Leibler (KL) divergence,  
209 between two spectra (see Methods). For a given sample, PaSD-qc identifies chromosomes with aberrant  
210 amplification patterns by calculating the distance of each chromosome’s PSD from the sample-average  
211 PSD. A chromosome is considered aberrant if it’s KL-divergence is two standard deviations beyond the  
212 sample median across all chromosomes.

213           We demonstrate how this method can identify false-positive chromosomal copy alterations by  
214 analyzing the “1465” neurons from [4]. This set of samples is informative as it includes a high coverage  
215 bulk sample from the same tissue to establish a true copy profile. PaSD-qc identifies chromosomes 15-17  
216 and 19-22 as inconsistently amplified in at least half of the samples (Figure 5A, Table S2); sex  
217 chromosomes are ignored in this analysis. Except for chromosome 15, each of these chromosomes is  
218 called as significantly copy-altered in at least 25% of samples (Figure 5B, Table S2) by the BICseq2  
219 algorithm [23]. If a whole chromosome deletion is present in at least 25% of cells, that deletion should be

220 apparent in bulk sequencing; however, bulk analysis of the tissue reveals all chromosomes to be copy  
221 neutral, indicating that the single cell copy alterations are artifacts. Different scWGS protocols show  
222 different patterns of aberrantly amplified chromosomes (Figure S4).

### 223 **Discriminating high- and low-quality samples**

224 We additionally profiled three newly amplified samples from the “1465” individual. Prior analysis  
225 showed these samples to be of low quality (Figure S5). Comparing them to high-quality samples from  
226 “1465” and “4638” provide an illustrative example of how PaSD-qc distinguishes high- and low-quality  
227 samples. Not only are the PSDs distinguishable by eye (Figure 6A), but the poor-quality samples also  
228 have a wider distribution of amplicon sizes and smaller median amplicon size (Figure 6B). Additionally,  
229 using the symmetric KL-divergence, PaSD-qc clusters the libraries based on amplification behavior  
230 (Figure 6C). The clustering correctly groups samples by high- and low-quality and further by biological  
231 origin. Finally, PaSD-qc can use the symmetric KL-divergence to probabilistically assign samples to  
232 different categories (e.g., high- and low-quality) using pre-computed gold-standard spectra. The toolbox  
233 includes methods which allow users to generate these gold-standard spectra from their own data. PaSD-qc  
234 can fully and accurately profile samples with coverage as low as 0.5X, and it provides accurate sample  
235 clustering and category assignment with coverage as low as 0.1X (Figure S6).

236

### 237 **Discussion:**

238 Here we have demonstrated the effectiveness of PaSD-qc to comprehensively evaluate the quality and  
239 amplification properties of scWGS data. Although several studies have recently compared the uniformity  
240 of different scWGS protocols [21, 24, 25], each study uses its own collection of statistics, making the task  
241 of determining the superior protocol difficult. We believe PaSD-qc represents an important step forward  
242 for the field as it provides a standardized suite of analyses that researchers can easily insert into any  
243 pipeline. In particular, PaSD-qc introduces novel methods to estimate the full distribution of amplicon

244 sizes in a sample, identify individual chromosomes which were poorly amplified, and compare samples  
245 based on amplification behavior.

246         These analyses not only allow comparisons across amplification protocols but also provide an  
247 important starting point for variant analysis. For example, it was recently demonstrated that the  
248 correlation in allelic balance induced by the large amplicons of MDA can be exploited to increase the  
249 accuracy of single cell single nucleotide variant (SNV) calling [11]. Dong et al. proposed a method  
250 employing a kernel smoothing algorithm that requires a user-defined bandwidth to compute the expected  
251 balance at a given genomic locus. The length of the bandwidth reflects the user's belief about the  
252 maximum distance at which informative correlation exists, and the authors suggest using a fixed  
253 bandwidth of 10 kb for all samples. However, PaSD-qc provides a principled, data-driven strategy to  
254 assign a tailored bandwidth to each individual sample as the 95% upper bound on amplicon sizes  
255 naturally defines a maximum correlation distance.

256         Additionally, our results address the question of whether *in vitro* amplification of the human  
257 genome by the Phi-29 MDA polymerase [26] produces amplicons of 10-100 kb as documented in  
258 bacterial genomes [27]. We found that some protocols approach the upper bound while others produce far  
259 smaller amplicons, with the lower bound in the 1-5 kb range. This has important consequences for PacBio  
260 or 10X Genomics sequencing on single cells in which fragments many kilobases in length are required. In  
261 particular, only some protocols may consistently produce large enough amplicons to make long-insert or  
262 haplotype-based sequencing possible.

263         Our results demonstrate that single cell amplification methods can artifactually induce whole  
264 chromosome copy alterations due to systematic under-amplification. Patterns of under amplification  
265 appear to be consistent across the same protocol but to vary between different protocols. As scWGS is  
266 becoming an increasingly popular choice to characterize copy number alterations in both research and  
267 clinical settings [6, 16, 20], the ability to identify false-positive copy changes is important. In addition,  
268 our results suggest that high-quality MDA data are likely yield accurate calls for small CNVs (smaller

269 than its amplicon sizes), as its sub-amplicon variance approaches that of bulk sequencing. Prior studies  
270 have focused on the detection of large copy alterations; none have specifically examined suitability for  
271 very small CNV calling.

272 Lastly, full mutational analysis at the single cell level requires high-coverage (>30X) sequencing,  
273 but the uneven quality of scWGS data, primarily due to the variable quality of cells, has often resulted in  
274 only a portion of the data generated being usable. The ability to accurately characterize data quality from  
275 low-coverage data suggests that a cost-effective approach in scWGS data generation is to screen a large  
276 number of cells at very low coverage (e.g., <0.1X) and select only a small number of high-quality  
277 candidates for additional sequencing. PaSD-qc provides an efficient computational framework to perform  
278 this evaluation.

279

## 280 **Conclusion:**

281 High-coverage scWGS enables identification of single nucleotide variants and other mutations at the  
282 single cell level, but mitigating the biases arising from whole-genome amplification remains a challenge.  
283 The proposed statistical method allows a detailed characterization of the data quality for scWGS datasets,  
284 aiding in selection of appropriate protocols and ensuring the fidelity of downstream analysis.

285

## 286 **Methods:**

### 287 **Data:**

288 MDA data for “4638” (Brain A), “1465” (Brain B), and “4643” (Brain C) were previously obtained by  
289 our group [4]. Additionally, three muscles cells from the “1465” individual were isolated, amplified, and  
290 sequenced as in that study. Fourteen additional MDA samples (C1a/b, C2a/b, C3a, N1a/b, N2a/b, N3a,  
291 N4a/b) were obtained from [20] (Short Read Archive accession number SRP052954). MALBAC samples  
292 were obtained from [5]. In Figure 2, the bulk sample is bulk cortex from “1465”, the MDA sample is cell  
293 30 from “1465”, and the MALBAC sample has the SRA accession number SRX204745. All data were  
294 downsampled to 1X using SAMtools prior to analysis.

### 295 **Power spectral density estimation**

296 Starting with a BAM file, read depth for each arm of each chromosome is extracted as the time series  
297  $\mathbf{x}_{c_a}(t) = (x_{t_1}, x_{t_2}, \dots, x_{t_n})$  where  $c$  is the chromosome,  $a$  is the chromosome arm and  $t_i$  is the start  
298 position of a uniquely mappable read. Uniquely mappable positions for the hg19 genome were download  
299 from the UCSC genome browser. By default, PaSD-qc uses mappability tracks calculated for 100 bp  
300 reads. Any series with fewer than 10 million observations is removed from further analysis. Each series  
301  $\mathbf{x}_{c_a}(t)$  is then divided into  $M$  windows of length  $L$  overlapping by  $D$  positions. By default,  $L = 1 \times 10^6$   
302 and  $D = 5 \times 10^5$ . The Lomb-Scargle algorithm [17, 18] is used to calculate the power spectral density,  
303  $f_{c_a,m}$ , for each  $\mathbf{x}_{c_a,m}(t)$  at eight thousand frequencies,  $\omega$ , evenly spaced from 1e-6 to 5e-3. The PSD for  
304 each chromosome is then estimated as

$$f_c(\omega) = \frac{\sum_{m=1}^M f_{c_p,m}(\omega) + \sum_{n=1}^N f_n(\omega)}{M + N}$$

305 where  $M$  and  $N$  are the number of windows on the  $p$  and  $q$  arms of chromosome  $c$ , respectively. The  
306 average PSD for an individual sample is then calculated as  $f(\omega) = \text{median}\{f_c(\omega)\}$ .

307           The mathematical details of Lomb-Scargle PSD estimation are described in SI. The theoretical  
308 justification for the power spectral density as a measure of variance in an aperiodic signal is also given in  
309 SI.

### 310 **Normalizing and plotting power spectral densities**

311 To remove edge effects and effects arising purely from sequencing, we take an idealized bulk sample as  
312 the baseline for the read depth power spectral density. The idealized bulk PSD,  $f_b$ , was derived by fitting  
313 a lowess curve to the bulk PSD shown in Figure 2A. The spectral density for each single cell sample is  
314 then normalized using the decibel transform as

$$dB(\omega) = 10 \times \log_{10} \frac{f(\omega)}{f_b(\omega)}.$$

315 This transform is standard in digital signal processing to remove a background signal.

316           Traditionally, power spectral densities are plotted as a function of frequency. However, for the  
317 genomic read depth signal, frequency takes on the unintuitive units of inverse genomic scale (1/bp). We  
318 instead choose to plot the PSD as a function of period,  $1/\omega$ . This results in the familiar units of genomic  
319 scale (bp) on the x-axis. We believe this eases interpretation, especially for those unfamiliar with power  
320 spectral densities.

### 321 **Estimating the distribution of amplicon sizes from the power spectral density**

322 As motivated in “Results”, the dynamic portion of the scWGS PSD curve reflects the cumulative  
323 distribution of the amplicon sizes in that sample. This distribution can thus be estimated by fitting a  
324 linearly scaled cumulative distribution function to this dynamic region. In practice, which distribution  
325 function should be fit is governed by two principles: 1) how tractable is fitting the curve using modern  
326 gradient descent algorithms, and 2) how well does the estimated distribution reproduce the original data.  
327 The first problem is one purely of computation and amounts to whether the distribution function has a  
328 closed-form solution or easily approximated integral solution. We tested three distributions which fit this

329 criterion: the normal (erf), logistic, and gamma distributions. To solve the second problem, we used the  
330 estimated density to simulate an idealized amplification process and compared the PSD of the idealized  
331 process to that of the original sample. The simulation procedure is described in the section below. We  
332 found the normal (erf) distribution best reproduced the data.

333 Let  $y = 10 \times \log_{10} \frac{f(\omega)}{f_b(\omega)}$  and  $x = -\log_{10} \omega$ . The dynamic region of the curve is fit as

$$y \approx A + \frac{B}{\sqrt{\pi}} \int_{-x}^x e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt.$$

334 The log-transformed density of the amplicon sizes is then estimated as  $\mathcal{N}(\mu, \sigma^2)$ . To estimate the median  
335 and 95% bounds, we draw 100,000 observations from the above distribution and calculate the median and  
336 percentiles of  $\{10^\theta\}_{\theta=1}^{1e5}$ , where  $\theta$  is a simulated observation.

### 337 **Simulating an idealized amplification process**

338 Let  $p(\cdot | \Theta)$  be the log-distribution of amplicon sizes estimated using the above method. For a given  
339 chromosome arm, an idealized amplification process is simulated using the following algorithm:

- 340 1. Initialize a vector,  $v$ , of length equal to the length of chromosome arm with all entries zero.
- 341 2. Randomly simulate an amplicon size as  $l = 10^\theta$  where  $\theta \sim p(\cdot | \Theta)$ .
- 342 3. Randomly choose a starting position  $s$ , where  $s \sim \text{Unif}(a, b)$  where  $a$  and  $b$  are the start and end  
343 coordinates of the chromosome arm
- 344 4. Increase the values of the entries of  $v$  overlapped by the amplicon by one
  - 345 a. Note: if  $s + l > b$ , the simulated amplicon is discarded
- 346 5. Repeat 2-5 until the desired average depth of coverage is reached.
  - 347 a. Depth of coverage is calculated as  $\sum_{i=0}^{b-a} v_i / (b - a)$ .
- 348 6. Randomly choose  $K$  non-zero observations from  $v$  where  $K$  is the number of non-zero  
349 observations from the chromosome arm in the original data.

350 The PSD of the resulting simulated read depth signal is then estimated and normalized as described  
351 above. To account for total power differences and mean shifts between the simulated data and the true  
352 data due to the idealized nature of the above algorithm, we normalize each curve by the maximum  
353 observed power and mean shift each curve such that  $f(10^{-3}) = 0$ . We chose to use the p arm of  
354 chromosome 3 for simulation purposes as in our experience it is a large arm with highly consistent  
355 amplification across samples.

### 356 **Estimating the autocovariance function**

357 The autocovariance function,  $\gamma$ , estimates the covariance of a time series against itself at lags  $k$ . As  
358 derived in SI, the real-valued sample autocovariance can be estimated from the PSD as

$$\gamma(k) = \int_{-\frac{1}{2}}^{\frac{1}{2}} \cos(2\pi\omega k) f(\omega) d\omega.$$

359 This integral can be quickly and accurately estimated numerically using any modern quadrature  
360 technique. We use Simpson's rule.

361 To directly calculate the ACF from unevenly space time series data, we define the "observation"  
362 function as

$$\mathbb{1}(t) = \begin{cases} 1, & \text{if } x_t \text{ observed} \\ 0, & \text{otherwise.} \end{cases}$$

363 For lag  $h$  we construct the set  $S_h = \{x_t \mid \mathbb{1}(t+h) = 1\}$ , which is the set of all observations such that an  
364 observation at a distance of  $h$  is also present. The sample autocovariance is then calculated as

$$\hat{\gamma}(h) = \frac{1}{|S_h|} \sum_{x_t \in S_h} (x_t - \bar{x})(x_{t+h} - \bar{x})$$

365 where  $\bar{x}$  is the sample mean of the time series and  $|S_h|$  denotes the size of  $S_h$ .

### 366 **Comparing the behavior of different spectra**

367 Given two probability densities,  $p_1$  and  $p_2$  and a vector of observations,  $\mathbf{X}$ , the Kullback-Leibler  
368 divergence is an informatic dissimilarity measure between the two densities and is defined as



$$KL(p_1, p_2) = \mathbb{E}_{p_1} \left[ \ln \frac{p_1(\mathbf{X})}{p_2(\mathbf{X})} \right].$$

369 It can be shown (see SI) that the Kullback-Leibler (KL) divergence between two PSDs is

$$KL(f_1, f_2) = \sum_{0 < \omega_i < \frac{1}{2}} -\ln \frac{|f_1(\omega_i)|}{|f_2(\omega_i)|} + f_2(\omega_i)^{-1} f_1(\omega_i) - 1.$$

370 The KL-divergence is not a true distance metric as  $KL(f_1, f_2) \neq KL(f_2, f_1)$ . Following [22], we define the

371 symmetric divergence between to spectra as

$$\begin{aligned} d(f_1, f_2) &\equiv \frac{1}{N} [KL(f_1, f_2) + KL(f_2, f_1)] \\ &= \frac{1}{N} \sum_{0 < \omega_i < \frac{1}{2}} \frac{f_1(\omega_i)}{f_2(\omega_i)} + \frac{f_2(\omega_i)}{f_1(\omega_i)} - 2 \end{aligned}$$

372 where  $N$  is the total number of frequencies in the sum. This value is reflexive and always non-negative

373 (see SI); thus  $d$  is a principled statistical distance metric between two spectra.

374 To identify aberrantly amplified chromosomes, we calculate  $d(f, f_c)$  for each chromosome of a

375 sample. We then calculate the median divergence and the median absolute difference of the divergences.

376 A chromosome is considered aberrant if its divergence is greater than the sum of the median and two

377 times the median absolute difference. To cluster samples by behavior, the pairwise divergence is

378 calculated between each pair of sample PSDs. The resulting symmetric distance matrix is then used to

379 perform hierarchical clustering.

### 380 **Estimating median absolute pairwise difference**

381 The BICseq2 algorithm [23] was used to calculate the copy number in bins of 1 kb, 5kb, 10 kb, 50 kb,

382 100 kb, 500 kb, and 1 mb for all “1465” samples. Estimates were corrected for mappability and GC

383 content. For each bin size, MAPD is calculated as  $\text{median}\{|CN_i - CN_{i+1}|\}_{i=2}^n$ , where  $CN_i$  is the copy

384 number in the  $i^{\text{th}}$  bin and  $n$  is the total number of bins.

### 385 **Estimating chromosome-level copy number**

386 The BICseq2 algorithm was used to calculate the copy number in bins of 500 kb normalized for  
387 mappability and GC content. The copy number for each chromosome was taken to be the median copy  
388 number over all bins overlapping that chromosome. Additionally, BICseq2 automatically assigns a p-  
389 value to the significance of the copy change, and a chromosome was considered significantly copy altered  
390 if the assigned p-value was less than 0.05.

### 391 **Implementation**

392 PaSD-qc is implemented in python. It uses SAMtools to extract coverage from bam files and the astropy  
393 package [28] to implement an  $O(n \cdot \log n)$  Lomb-Scargle algorithm. The function `curve_fit` in the  
394 `scipy` module is used to fit the modified erf function to the scWGS PSD. Clustering of samples is  
395 performed by the `linkage` function also in the `scipy` module. PaSD-qc parallelizes across samples for  
396 efficient multi-sample analysis. Source code, documentation, and examples – including all data and code  
397 to reproduce the figures in this manuscript – are available at <https://github.com/parklab/PaSD-qc>.

398

### 399 **Abbreviations:**

400 **scWGS:** single cell whole-genome sequencing

401 **PSD:** power spectral density

402 **ACF:** autocovariance function

403 **MDA:** multiple displacement amplification

404 **MALBAC:** multiple looping and annealing based amplification cycles

405 **KL-divergence:** Kullback-Leibler divergence

406 **CNV:** copy number variation

407 **SNV:** single nucleotide variation

408 **Erf:** error function

409

410 **Declarations:**

411 **Acknowledgements:**

412 We thank Joe Luquette, Doga Gulhan, and Alon Galor for their valuable input in preparing this  
413 manuscript. We thank the Brain Somatic Mosaicism Network for supporting this research.

414 **Funding:**

415 This work has been supported by a grant from NIH (1U01MH106883) and the Harvard Ludwig Center.  
416 CAW is an investigator of the Howard Hughes Medical Institute.

417 **Availability of data and materials:**

418 PaSD-qc is implemented as a python package which is freely available at  
419 <https://github.com/parklab/PaSD-qc> along with all data and scripts necessary to reproduce the figures in  
420 this paper. BAM files were downloaded from the Short Read Archive with accession numbers  
421 SRP042470 for “1465” from [4], SRP061939 for “4638” and “4643” from [4], SRA060929 for samples  
422 from [5], and SRP052954 for samples from [20]. The BAMs for the three additional “1465” samples  
423 shown in Figure 6 are available from the corresponding author upon reasonable request.

424 **Author’s contributions:**

425 MAS and PJP conceived the idea for the paper. MAS designed the statistical methods, implemented the  
426 algorithms, performed the analysis, and wrote the manuscript. ARB and CV aligned the samples and  
427 helped prepare the manuscript. ML and MEC collected the data. CAW provided resource support and  
428 read the manuscript.

429 **Competing interest:**

430 The authors declare that they have no competing interest.

431 **Consent for publication:**

432 Not applicable.

433 **Ethics approval and consent to participate:**

434 Not applicable

435 **Additional files:**

436 Additional file 1: supplementary information, tables S1-S2, and figures S1-S6

437

438

439

440

441

442

443

444

445

446

447

## 448 **References:**

- 449 [1] G. D. Evrony, E. Lee, B. K. Mehta, Y. Benjamini, R. M. Johnson, X. Cai, L. Yang, P. Haseley,  
450 H. S. Lehmann, P. J. Park, and C. A. Walsh, “Cell lineage analysis in human brain using endogenous  
451 retroelements.” *Neuron*, vol. 85, pp. 49–59, Jan. 2015.
- 452 [2] Y. Wang, J. Waters, M. L. Leung, A. Unruh, W. Roh, X. Shi, K. Chen, P. Scheet, S. Vattathil,  
453 H. Liang, A. Multani, H. Zhang, R. Zhao, F. Michor, F. Meric-Bernstam, and N. E. Navin, “Clonal  
454 evolution in breast cancer revealed by single nucleus genome sequencing.” *Nature*, vol. 512, pp. 155–160,  
455 Aug. 2014.
- 456 [3] M. J. McConnell, M. R. Lindberg, K. J. Brennand, J. C. Piper, T. Voet, C. Cowing-Zitron,  
457 S. Shumilina, R. S. Lasken, J. R. Vermeesch, I. M. Hall, and F. H. Gage, “Mosaic copy number variation  
458 in human neurons.” *Science*, vol. 342, pp. 632–637, Nov. 2013.
- 459 [4] M. A. Lodato, M. B. Woodworth, S. Lee, G. D. Evrony, B. K. Mehta, A. Karger, S. Lee, T. W.  
460 Chittenden, A. M. D’Gama, X. Cai, L. J. Luquette, E. Lee, P. J. Park, and C. A. Walsh, “Somatic  
461 mutation in single human neurons tracks developmental and transcriptional history.” *Science*, vol. 350,  
462 pp. 94–98, Oct. 2015.
- 463 [5] C. Zong, S. Lu, A. R. Chapman, and X. S. Xie, “Genome-wide detection of single-nucleotide and  
464 copy-number variations of a single human cell.” *Science*, vol. 388, pp. 1622–1626, December 2012.
- 465 [6] W. Liu, H. Zhang, D. Hu, S. Lu, and X. Sun, “The performance of MALBAC and MDA methods  
466 in the identification of concurrent mutations and aneuploidy screening to diagnose beta-thalassaemia  
467 disorders at the single- and multiple-cell levels.” *Journal of clinical laboratory analysis*, May 2017.
- 468 [7] J. Xu, R. Fang, L. Chen, D. Chen, J.-P. Xiao, W. Yang, H. Wang, X. Song, T. Ma, S. Bo, C. Shi,  
469 J. Ren, L. Huang, L.-Y. Cai, B. Yao, X. S. Xie, and S. Lu, “Noninvasive chromosome screening of  
470 human embryos by genome sequencing of embryo culture medium for in vitro fertilization.” *Proceedings*

- 471 *of the National Academy of Sciences of the United States of America*, vol. 113, pp. 11907–11912, Oct.  
472 2016.
- 473 [8] T. Baslan, J. Kendall, L. Rodgers, H. Cox, M. Riggs, A. Stepansky, J. Troge, K. Ravi,  
474 D. Esposito, B. Lakshmi, M. Wigler, N. Navin, and J. Hicks, “Genome-wide copy number analysis of  
475 single cells.” *Nature protocols*, vol. 7, pp. 1024–1041, May 2012.
- 476 [9] Y. Wang and N. E. Navin, “Advances and applications of single-cell sequencing technologies.”  
477 *Molecular cell*, vol. 58, pp. 598–609, May 2015.
- 478 [10] C.-Z. Zhang, V. A. Adalsteinsson, J. Francis, H. Cornils, J. Jung, C. Maire, K. L. Ligon,  
479 M. Meyerson, and J. C. Love, “Calibrating genomic and allelic coverage bias in single-cell sequencing.”  
480 *Nature communications*, vol. 6, p. 6822, Apr. 2015.
- 481 [11] X. Dong, L. Zhang, B. Milholland, M. Lee, A. Y. Maslov, T. Wang, and J. Vijg, “Accurate  
482 identification of single-nucleotide variants in whole-genome-amplified single cells.” *Nature methods*,  
483 vol. 14, pp. 491–493, May 2017.
- 484 [12] M. L. Leung, Y. Wang, J. Waters, and N. E. Navin, “SNES: single nucleus exome sequencing.”  
485 *Genome biology*, vol. 16, p. 55, Mar. 2015.
- 486 [13] K. Leung, A. Klaus, B. K. Lin, E. Laks, J. Biele, D. Lai, A. Bashashati, Y.-F. Huang, R. Aniba,  
487 M. Moksa, A. Steif, A.-M. Mes-Masson, M. Hirst, S. P. Shah, S. Aparicio, and C. L. Hansen, “Robust  
488 high-performance nanoliter-volume single-cell multiple displacement amplification on planar substrates.”  
489 *Proceedings of the National Academy of Sciences of the United States of America*, vol. 113, pp. 8484–  
490 8489, July 2016.
- 491 [14] M. Rhee, Y. K. Light, R. J. Meagher, and A. K. Singh, “Digital droplet multiple displacement  
492 amplification (ddMDA) for whole genome sequencing of limited DNA samples.” *PloS one*, vol. 11, 2016.

- 493 [15] C. Chen, D. Xing, L. Tan, H. Li, G. Zhou, L. Huang, and X. S. Xie, “Single-cell whole-genome  
494 analyses by linear amplification via transposon insertion (LIANTI).” *Science*, vol. 356, pp. 189–194, Apr.  
495 2017.
- 496 [16] X. Cai, G. D. Evrony, H. S. Lehmann, P. C. Elhosary, B. K. Mehta, A. Poduri, and C. A. Walsh,  
497 “Single-cell, genome-wide sequencing identifies clonal somatic copy-number variation in the human  
498 brain.” *Cell reports*, vol. 8, pp. 1280–1289, Sept. 2014.
- 499 [17] N. R. Lomb, “Least-squares frequency analysis of unequally spaced data.” *Astrophysics and*  
500 *Space Science*, vol. 39, pp. 447–462, Feb. 1976.
- 501 [18] J. D. Scargle, “Studies in astronomical time series analysis. ii - statistical aspects of spectral  
502 analysis of unevenly spaced data.” *Astrophysical Journal*, vol. 263, pp. 835–853, Dec. 1982.
- 503 [19] P. Welch, “The use of fast Fourier transform for the estimation of power spectra: A method based  
504 on time averaging over short, modified periodograms.” *IEEE Transactions on Audio and*  
505 *Electroacoustics*, vol. 15, pp. 70–73, June 1967.
- 506 [20] C.-Z. Zhang, A. Spektor, H. Cornils, J. M. Francis, E. K. Jackson, S. Liu, M. Meyerson, and  
507 D. Pellman, “Chromothripsis from DNA damage in micronuclei.” *Nature*, vol. 522, pp. 179–184, June  
508 2015.
- 509 [21] L. Ning, Z. Li, G. Wang, W. Hu, Q. Hou, Y. Tong, M. Zhang, Y. Chen, L. Qin, X. Chen, H.-Y.  
510 Man, P. Liu, and J. He, “Quantitative assessment of single-cell whole genome amplification methods for  
511 detecting copy number variation using hippocampal neurons.” *Scientific reports*, vol. 5, p. 11415, June  
512 2015.
- 513 [22] R. H. Shumway and D. S. Stoffer, “Statistical methods in the frequency domain,” in *Time Series*  
514 *Analysis and Its Applications*, Springer, Jan. 2011.

- 515 [23] R. Xi, S. Lee, Y. Xia, T.-M. Kim, and P. J. Park, “Copy number analysis of whole-genome data  
516 using bic-seq2 and its application to detection of cancer susceptibility variants.” *Nucleic acids research*,  
517 vol. 44, pp. 6274–6286, July 2016.
- 518 [24] C. F. A. de Bourcy, I. De Vlaminck, J. N. Kanbar, J. Wang, C. Gawad, and S. R. Quake, “A  
519 quantitative comparison of single-cell whole genome amplification methods.” *PloS one*, vol. 9, 2014.
- 520 [25] E. Borgström, M. Paterlini, J. E. Mold, J. Frisen, and J. Lundeberg, “Comparison of whole  
521 genome amplification techniques for human single cell exome sequencing.” *PloS one*, vol. 12, 2017.
- 522 [26] F. B. Dean, S. Hosono, L. Fang, X. Wu, A. F. Faruqi, P. Bray-Ward, Z. Sun, Q. Zong, Y. Du,  
523 J. Du, M. Driscoll, W. Song, S. F. Kingsmore, M. Egholm, and R. S. Lasken, “Comprehensive human  
524 genome amplification using multiple displacement amplification.” *Proceedings of the National Academy  
525 of Sciences of the United States of America*, vol. 99, pp. 5261–5266, Apr. 2002.
- 526 [27] L. Blanco, A. Bernad, J. M. Lázaro, G. Martín, C. Garmendia, and M. Salas, “Highly efficient  
527 DNA synthesis by the phage phi 29 DNA polymerase. symmetrical mode of DNA replication.” *The  
528 Journal of biological chemistry*, vol. 264, pp. 8935–8940, May 1989.
- 529 [28] A. Collaboration, T. P. Robitaille, E. J. Tollerud, P. Greenfield, M. Droettboom, E. Bray,  
530 T. Aldcroft, M. Davis, A. Ginsburg, A. M. Price-Whelan, W. E. Kerzendorf, A. Conley, N. Crighton,  
531 K. Barbary, D. Muna, H. Ferguson, F. Grollier, M. M. Parikh, P. H. Nair, H. M. Unther, C. Deil,  
532 J. Woillez, S. Conseil, R. Kramer, J. E. H. Turner, L. Singer, R. Fox, B. A. Weaver, V. Zabalza, Z. I.  
533 Edwards, K. Azalee Bostroem, D. J. Burke, A. R. Casey, S. M. Crawford, N. Dencheva, J. Ely,  
534 T. Jenness, K. Labrie, P. L. Lim, F. Pierfederici, A. Pontzen, A. Ptak, B. Refsdal, M. Servillat, and  
535 O. Streicher, “Astropy: A community python package for astronomy,” , vol. 558, p. A33, Oct. 2013.
- 536
- 537



538 **Figure legends:**

539 **Figure 1: Overview of single cell whole-genome sequencing and sources of artifacts, and the PaSD-**

540 **qc pipeline. A.** Schematic overview of single-cell whole genome sequencing and the artifacts created by

541 whole-genome amplification. The extent and patterns of the biases depend on the cell condition (high- or

542 low-integrity) and on the scWGS protocol used (protocol A or protocol B). The pink triangles in the

543 “Large-Scale Feature Correlation” represent genomic events (e.g., single nucleotide variants) which are

544 spanned by a single amplicon and are thus correlated. The only correlation pattern present in bulk

545 sequencing is due to paired-end sequencing, represented by positions marked “A” and “T” spanned by the

546 mate pair. **B.** Schematic overview of the PaSD-qc pipeline. Read depth is extracted from bam files at

547 uniquely mappable positions. Red rectangles represent regions where the true read depth is unknown due

548 to low mappability, locus dropout, or sequencing bias. PaSD-qc uses a custom power spectral density

549 estimation procedure to accurately estimate the correlation patterns in the data, and these patterns are then

550 used to assess amplification properties and quality control measures. By default, the results are

551 summarized in an interactive HTML report.

552 **Figure 2: Using power spectral density to infer sample-specific amplification properties of scWGS**

553 **data. A.** PaSD-qc power spectral densities for a bulk sample (blue), MDA sample (green), and MALBAC

554 sample (purple). The very low noise of the estimates allows amplification properties of the three samples

555 to be inferred, including the paired-end insert size distributions (Figure S1), the range of amplicon sizes

556 for MDA (~5-100 kb) and MALBAC (~1-5 kb), and the sub- and supra-amplicon variances of the two

557 amplification protocols. Interestingly, whereas MDA has a higher supra-amplicon variance than

558 MALBAC, its sub-amplicon variance is considerably lower. **B.** Power spectral density estimates using

559 the algorithm from Leung et al, 2016 [13]. A similar algorithm is used in Zong et al, 2012 [5].

560 Background noise dominates the estimates making feature extraction infeasible. Resolution was limited to

561 2 kb because the data were binned into 1 kb bins as suggested per those algorithms.

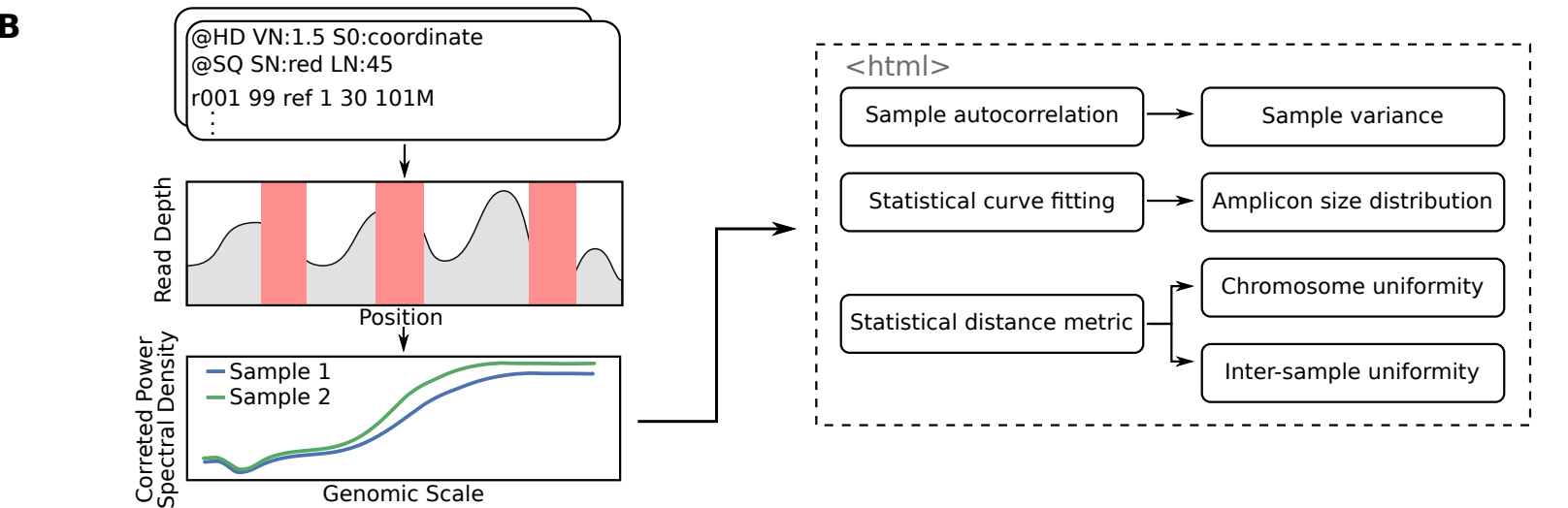
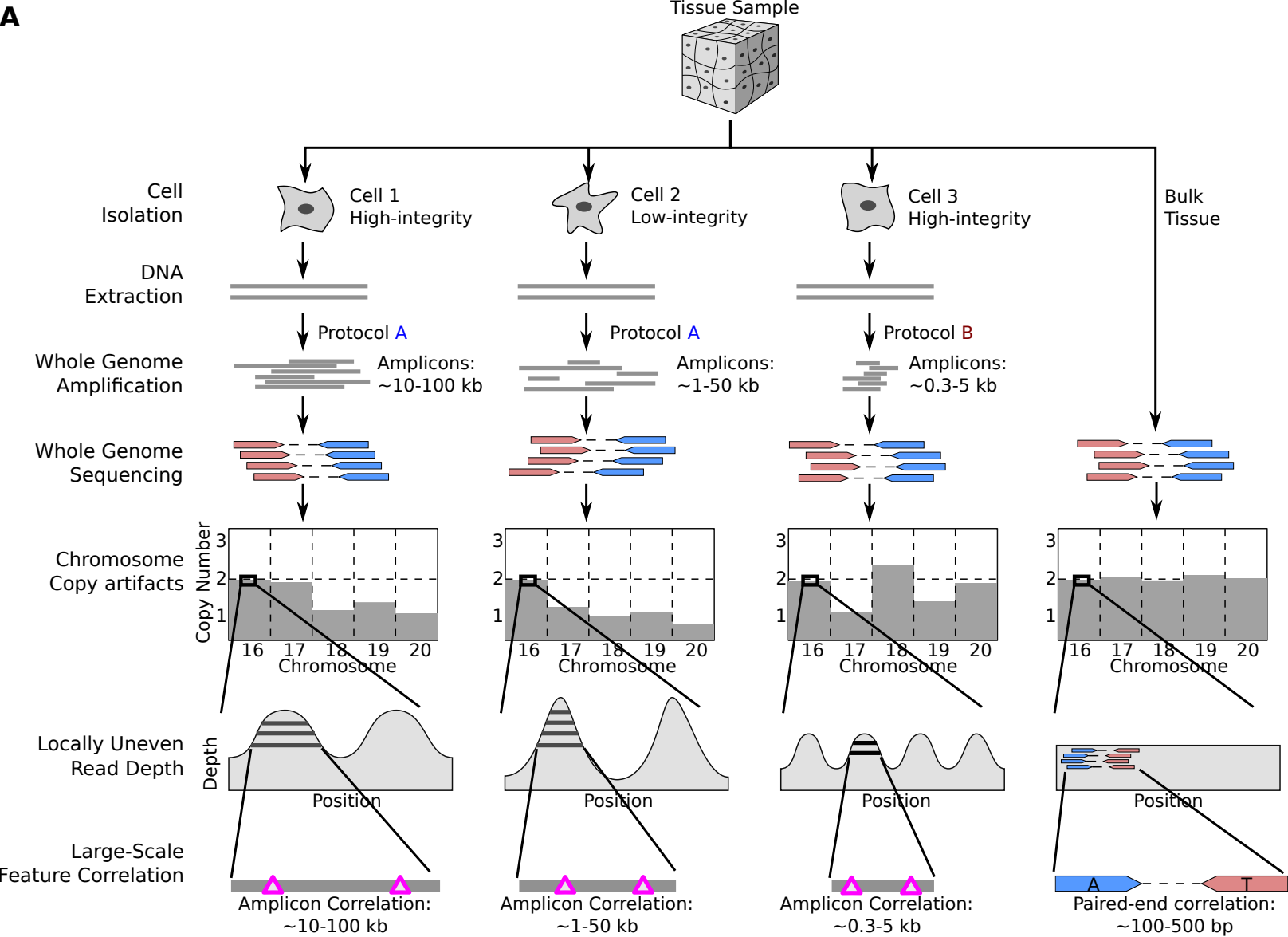
562 **Figure 3: The distribution of amplicon sizes can be directly estimated from the power spectral**  
563 **density. A.** MDA (green) and MALBAC (purple) curves as in Figure 2 along with the inferred error  
564 function (erf) fit of the dynamic region (red), the median amplicon size (pink stars), and 95% bounds on  
565 amplicon sizes (yellow stars). **B,C.** Distributions of inferred amplicon sizes in the MDA and MALBAC  
566 sample. Densities are normally distributed in a log scale (B), but highly skewed in standard coordinates  
567 (C). **D.** The average power spectral density (red) resulting from ten simulated amplification processes  
568 using the MDA density shown in C as the generative distribution. The shaded region represents the 95%  
569 confidence interval and the green curve corresponds to the original data. The MALBAC fit and fits using  
570 other distributions are shown in Figure S3. **E.** The average amplicon size distributions for 35 samples  
571 from Lodato et al, 2015 [4] (green) and 14 samples from Zhang et al, 2015 [20] (blue) reveal that  
572 different MDA protocols produce different amplicon size distributions, but a single protocol produces  
573 consistent amplicon size distributions across samples (shaded regions represent 95% confidence intervals  
574 around the average).

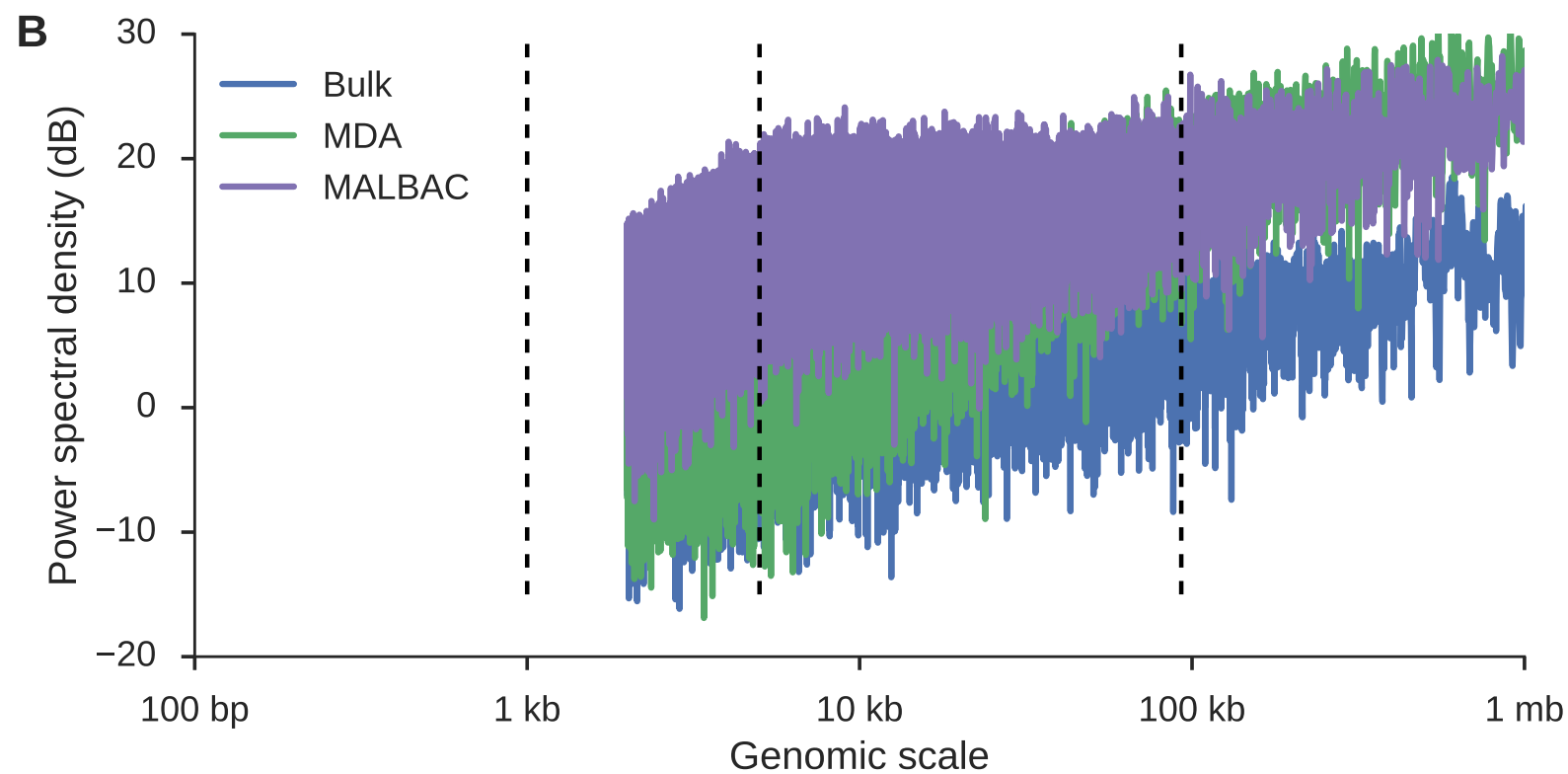
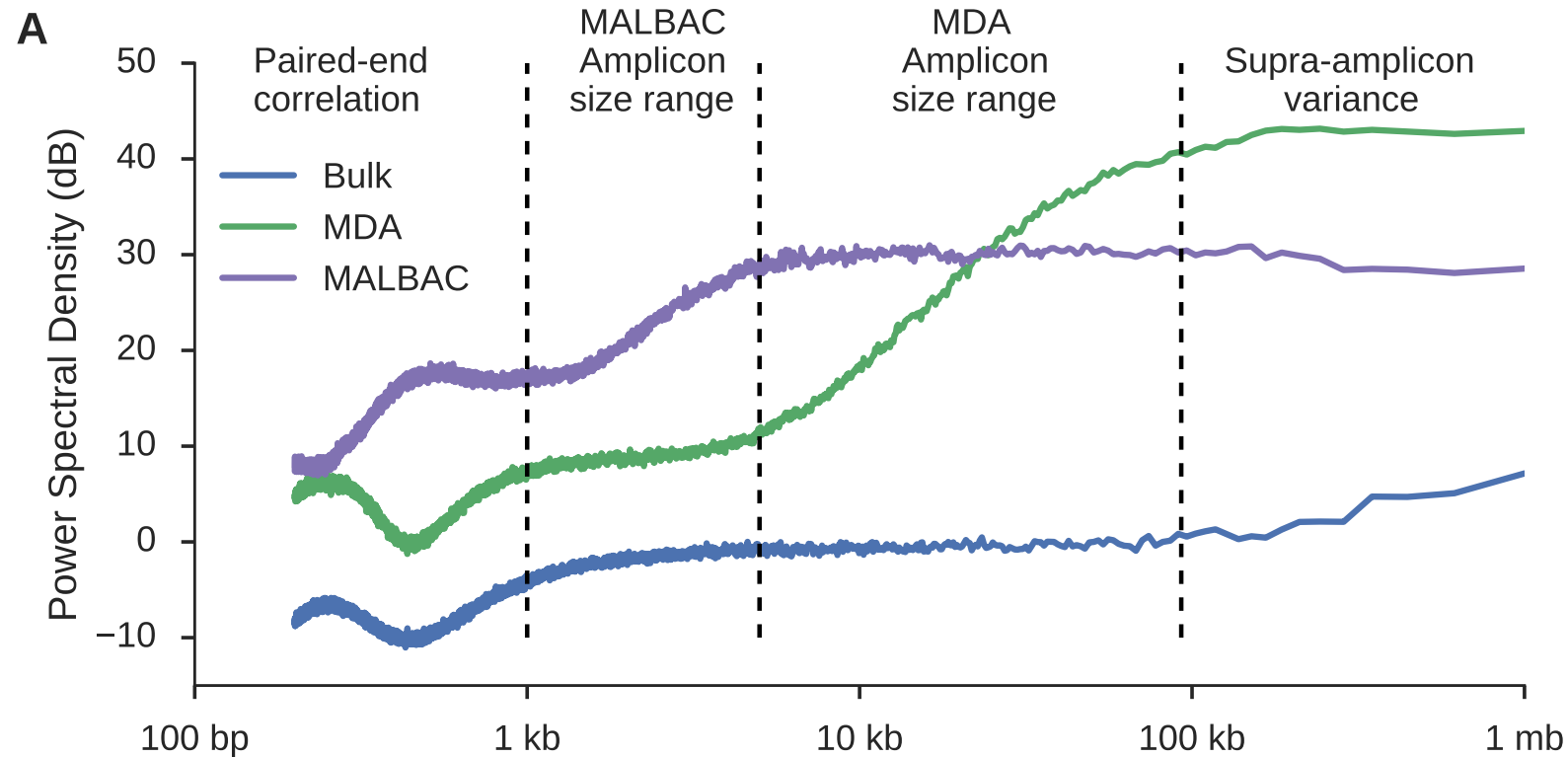
575 **Figure 4: The PaSD-qc variance measure outperforms prior dispersion estimates. A.** Average  
576 sample autocovariance with 95% confidence intervals for the 16 “1465” samples from Lodato et al, 2015  
577 [4] as calculated by PaSD-qc (blue) and by direct estimation (gold). See text for a comparison. **B.**  
578 Average MAPD scores with 95% confidence intervals calculated for seven bin sizes ranging from 1 kb to  
579 1 mb for 16 “1465” samples and 11 “4638” samples from Lodato et al, 2015. **C.** The average power  
580 spectral density with 95% confidence intervals for the same samples. **D.** Densities for the MAPD scores  
581 of the two sets of samples at 50 kb, the standard bin size at which the score is calculated. At this bin size,  
582 MAPD cannot distinguish behavior of the two sets of samples. **E.** Densities of PaSD-qc variance for the  
583 two sets of samples are significantly different.

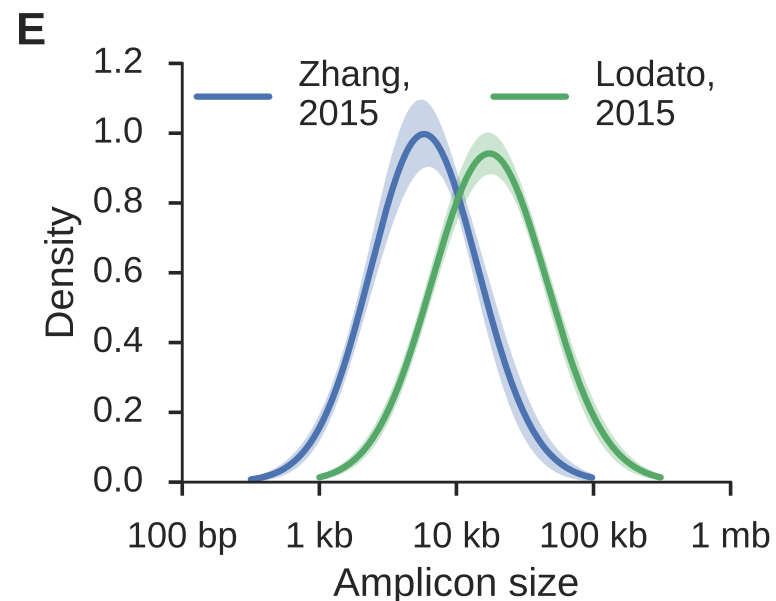
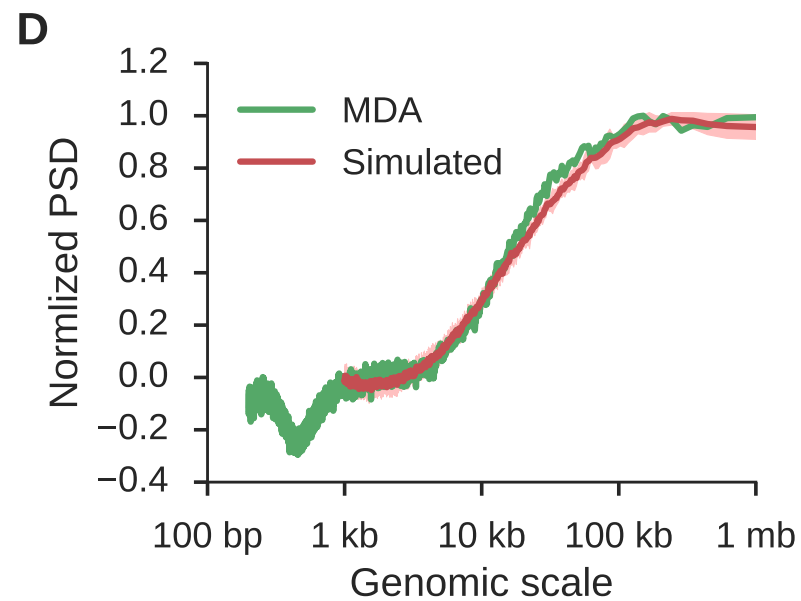
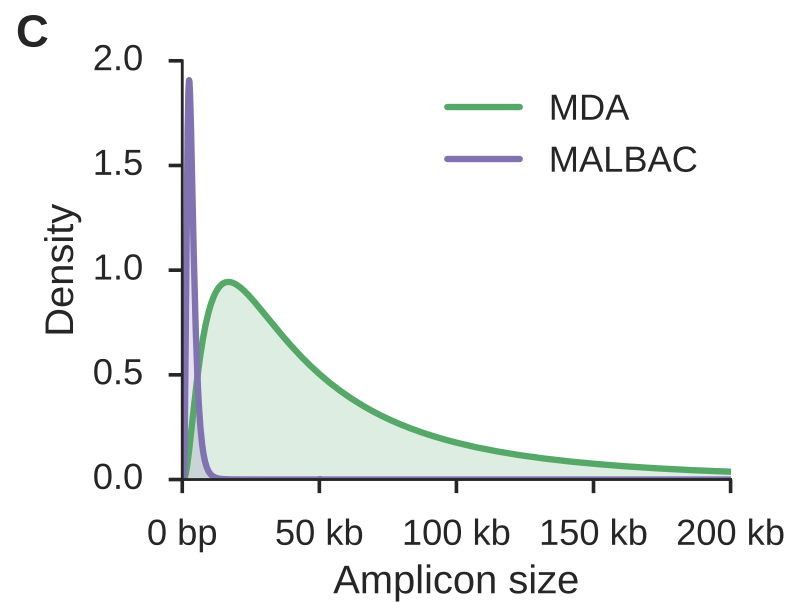
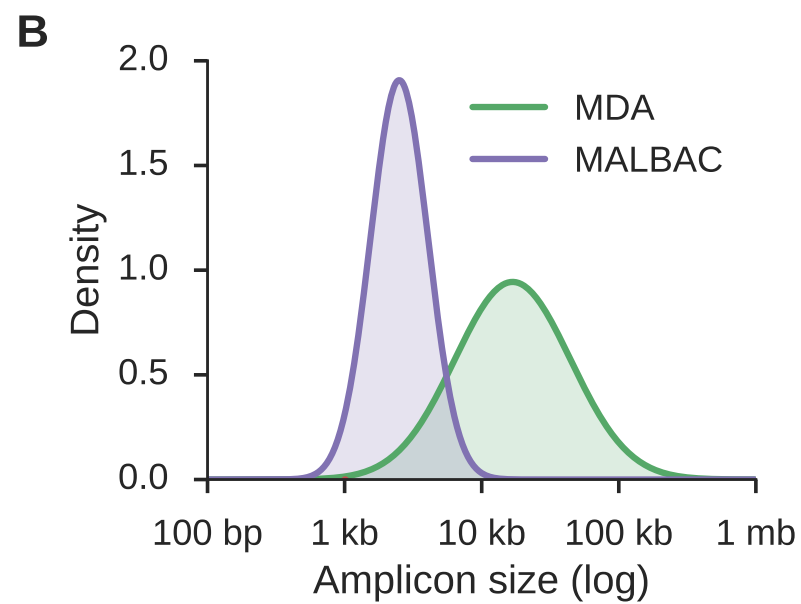
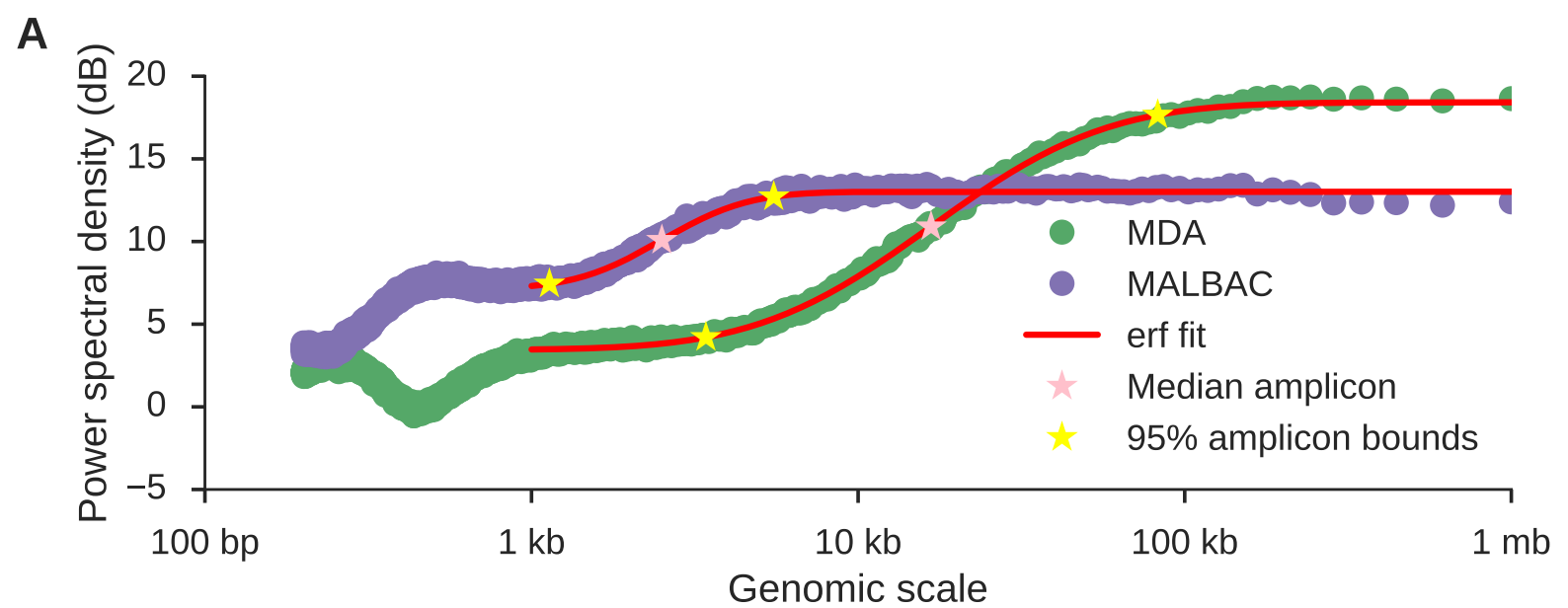
584 **Figure 5: Identification of false-positive chromosomal copy changes due to poor amplification. A.**  
585 boxplots of the KL divergence of each autosome from the sample-average PSD for the 16 “1465”  
586 samples. Chromosomes are labeled as failed (red) if PaSD-qc identifies the chromosome as aberrantly

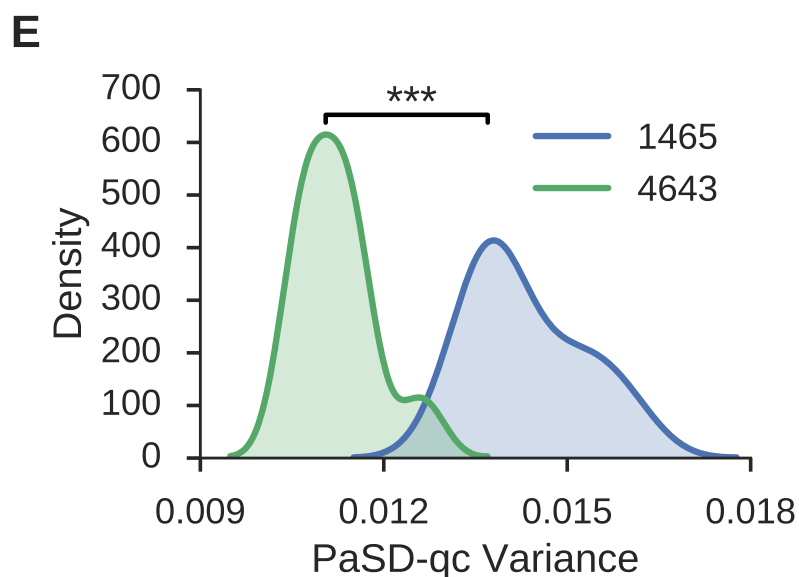
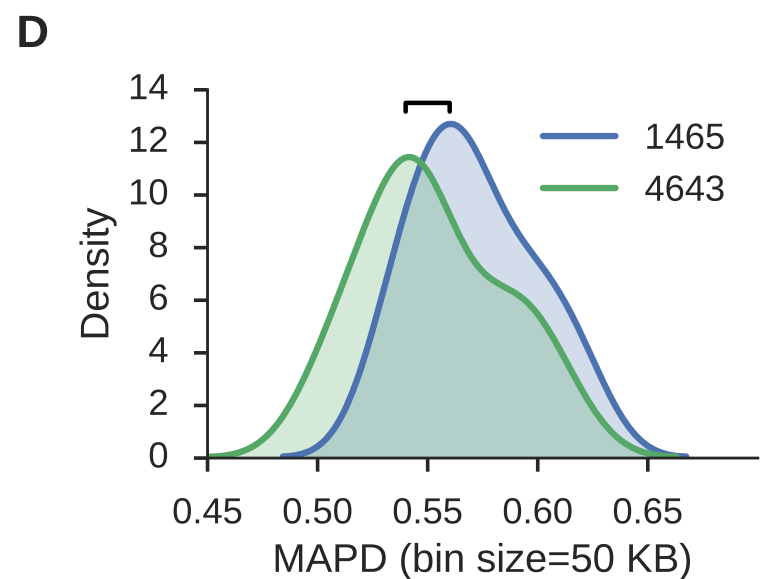
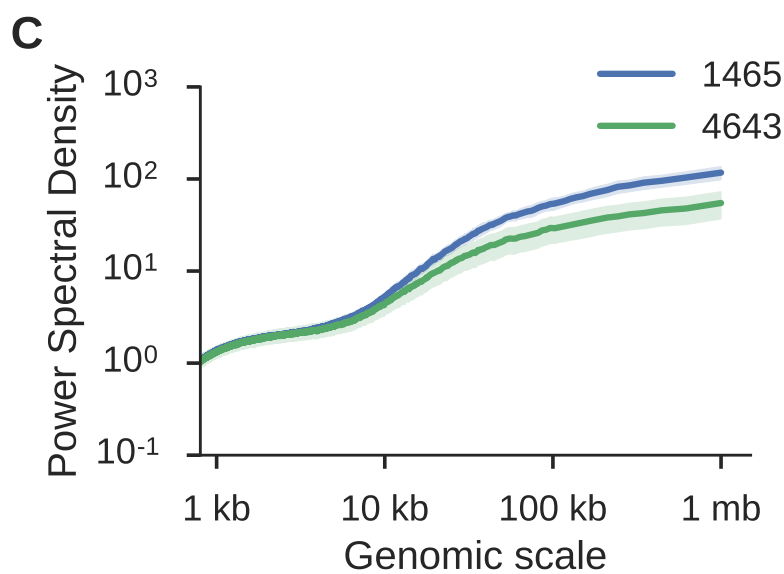
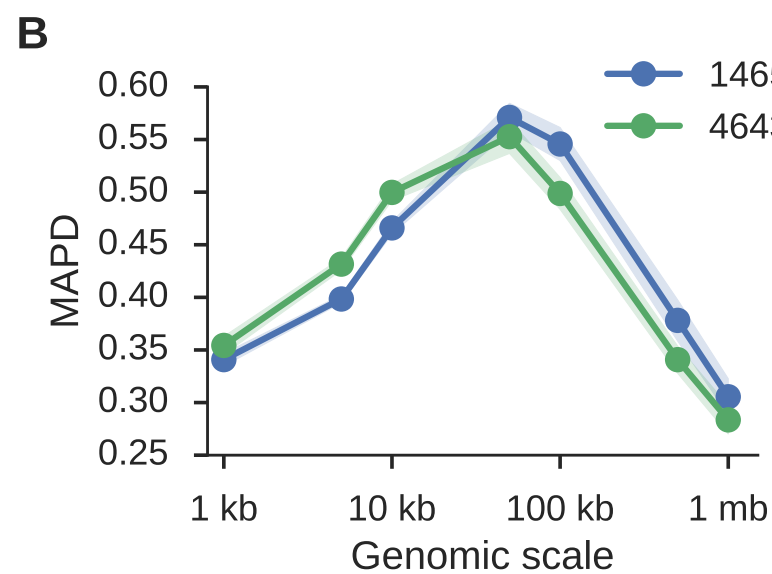
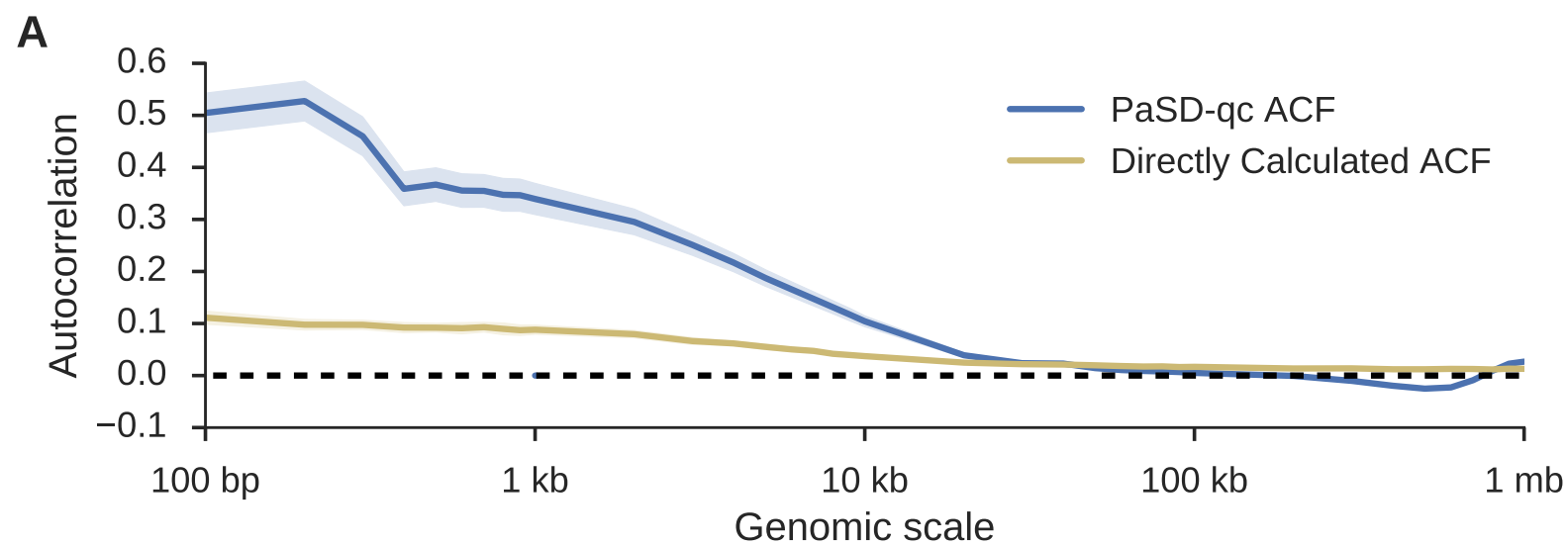
587 amplified in at least half of the samples (Table S2). **B.** the average copy number across all samples as  
588 inferred by the BICseq2 algorithm. Errorbars represent standard deviation across all samples.  
589 Chromosomes are considered copy aberrant if BICseq2 identifies a significant ( $p$ -value  $< 0.05$ ) alteration  
590 in at least 25% of samples (Table S2). A chromosomal deletion in at least 25% of cells should be  
591 identifiable in bulk sequencing. **C.** chromosome copy profile of a bulk sample from the same tissue as the  
592 single cell samples. All chromosomes are copy neutral.

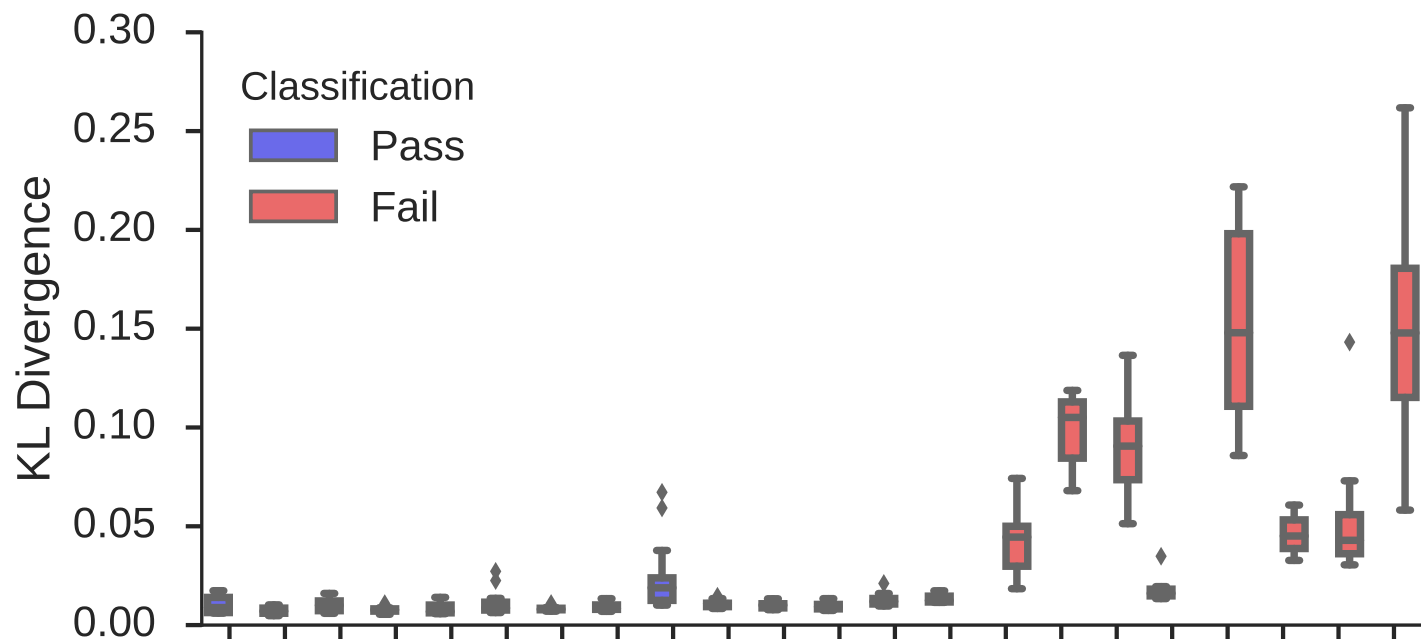
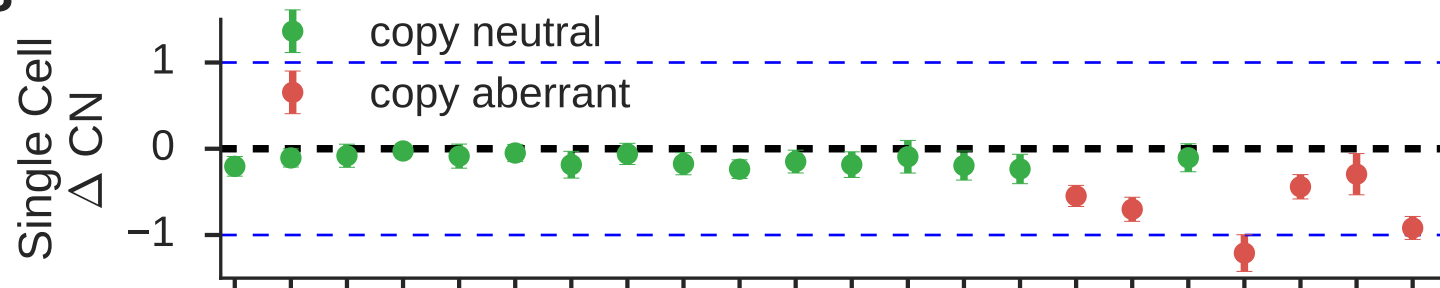
593 **Figure 6: PaSD-qc separates high-quality from low-quality samples and groups similarly behaving**  
594 **libraries.** **A.** power spectral densities for three low-quality (red) and six high-quality libraries (green). **B.**  
595 Amplicon size density plots for the nine samples. **C.** Hierarchical clustering using the symmetric KL-  
596 divergence correctly groups the samples based on both quality and biological origin.









**A****B****C**