

1 **Estimating the rate of index hopping on the Illumina HiSeq X platform**

2 *Tom van der Valk¹, Francesco Vezzi², Mattias Ormestad², Love Dalén^{3*}, Katerina Guschanski^{1*}*

3 ¹Animal Ecology, Department of Ecology and Genetics, Evolutionary Biology Centre, Uppsala
4 University, Norbyvägen 18D, 752 36, Uppsala, Sweden

5 ²Science for Life Laboratory, Tomtebodavägen 23A, 17165 Solna, Sweden

6 ³Department of Bioinformatics and Genetics, Swedish Museum of Natural History, SE-10405
7 Stockholm, Sweden

8 * These authors contributed equally

9 Corresponding authors: Katerina Guschanski (katerina.guschanski@ebc.uu.se), Tom van der Valk
10 (tom.vandervalk@ebc.uu.se)

11

12 **Abstract**

13 The high-throughput capacities of the Illumina sequencing platforms and possibility to label samples
14 individually have encouraged a wide use of sample multiplexing. However, this practice results in
15 read misassignment (usually <1%) across samples sequenced on the same lane. Alarming high rates
16 of read misassignment of up to 10% were reported for the latest generation of Illumina sequencing
17 machines. This potentially calls into question previously generated results and may make future use
18 of the newest generation of platforms prohibitive. In this study we rely on barcodes, short sequences
19 that are directly ligated to both ends of the DNA insert, which allows us to quantify the amount of
20 index hopping. Correcting for multiple sources of noise, we identify on average only 0.470% of reads
21 containing a hopped index. Multiplexing of samples on this platform is therefore unlikely to cause
22 markedly different results to those obtained from older platforms.

23

24 **Keywords**

25 Read misassignment, next generation sequencing, ExAmp chemistry, multiplexing

26

27

28 Introduction

29 Multiplexing samples for next-generation sequencing is a common practice in many biological and
30 medical applications [1–3]. During multiplexing, samples are individually labelled with unique
31 identifiers (indices) that are embedded within one (single indexing) or both (dual indexing)
32 sequencing platform-specific adapters [2,4]. The samples are subsequently pooled into a single DNA
33 library and sequenced on the same lane, greatly reducing per sample sequencing cost. Following
34 sequencing, computational demultiplexing, based on the sample-specific indices, allows for
35 assignment of the sequenced reads to the respective sample of origin. However, ever since
36 multiplexing approaches were introduced, low rates of read misassignment across samples
37 sequenced on the same lane have been reported on all Illumina platforms [4–8]. Read
38 misassignment is the results of reads carrying an unintended index and consequently being
39 erroneously attributed to the wrong original sample. The reported rate of read misassignment on
40 Illumina platforms that rely on the traditional bridge amplification for cluster generation is low (<1%)
41 [9,10] and therefore this source of error has been mostly ignored.

42 However, on the latest generation of Illumina sequencing platforms (HiSeq X, HiSeq 4000 and
43 NovaSeq) that rely on the exclusion amplification chemistry (ExAmp) in combination with patterned
44 flow cells, a wide range of read misassignment rates has been reported [11–13]. Whereas one study
45 reported no observable rate of read misassignment on neither HiSeq X and HiSeq 2500 platforms
46 [14], other studies have documented rates of up to 10% [13]. According to Illumina’s own estimates,
47 the rate of read misassignment on platforms with ExAmp chemistry is up to 2% [9].

48 As a consequence of conflicting results, the prevalence and severity of read misassignment on the
49 Illumina HiSeq X platforms remain unclear. This is partly due to the difficulties to reliably identify
50 misassigned reads in sequencing experiments, particularly if pooling similar samples types (e.g.
51 multiple individuals from the same population that have high sequence similarity). Recently, Illumina
52 introduced dual indices on all ExAmp sequencing platforms, allowing for the filtering of the majority

53 of reads that show signs of read misassignment. However, since indices can potentially be switched
54 at both ends of the molecule and the number of available indices is limited, it remains difficult to
55 obtain direct estimates of read misassignment on these platforms.

56 Some research questions require high confidence in read identity, as presence of rare sequence
57 variants can influence biological and medical conclusions. For instance, detection of low abundance
58 transcripts or rare mutations can influence diagnostic inferences [11,15–17]. Studies with low input
59 DNA quantities (e.g. single cell sequencing, ancient and historical DNA) are particularly susceptible to
60 such errors [4] . Similarly, population genomics studies frequently rely on low-coverage genomic
61 data, and presence of shared rare alleles across several populations or species can be interpreted as
62 evidence for gene flow or admixture [18–21].

63 Processes resulting in read misassignment, i.e. presence of reads with a switched index, are
64 numerous. The effect of sequencing errors that can convert one index sequence into another is well
65 known and has led to series of recommendations for designing highly distinct indices [2]. Jumping
66 PCR during bulk amplification of library molecules that carry different indices can generate chimeric
67 sequences and should be avoided [22–25]. Similarly, cross-contamination of indexing adapters during
68 oligonucleotide synthesis or laboratory work can lead to reads obtaining an unintended index.
69 Additionally, cluster misidentification due to “bleeding” of indices into neighbouring clusters have
70 been reported on all high throughput sequencing platforms [4–8,11,26]. However, for the latest
71 Illumina platforms with patterned flow cells and ExAmp chemistry, read misassignment has been
72 suggested to be caused by the presence of free-floating indexing primers in the final sequencing
73 library [9,13]. Such free-floating molecules can appear if sequencing libraries are not stored properly
74 and become fragmented or if final sequencing libraries are not properly size selected [27]. These
75 primers and molecules can then anneal to the pooled library molecules and get extended by DNA
76 polymerase before the rapid exclusion amplification on the flow cell, creating a new library molecule

77 with an erroneous index. We refer to this particular process of generating misassigned reads as index
78 hopping.

79 In this study we make use of inline barcodes, short unique seven base pair sequences ligated to both
80 ends of the DNA fragments [28], in combination with indexed primers that are traditionally used for
81 sample identification. The barcodes become part of the sequencing read and thus allow for accurate
82 identification of the read origin, even in the presence of index hopping. Therefore, the amount of
83 index hopping can be directly quantified by identifying reads with wrong barcode-index
84 combinations. We specifically use historical museum-preserved samples that are characterized by
85 low DNA quantity and quality (the DNA is degraded, chemically modified and shows single-strand
86 overhangs [29,30]). Libraries constructed from such low quality samples are prone to index hopping,
87 since generally these samples yield low coverage sequencing data and inferences are based on subtle
88 differences between limited sets of polymorphic sites [31,32]. Therefore, small quantities of
89 misassigned DNA fragments in aDNA studies can already cause erroneous inferences [33] and it is
90 thus crucial to distinguish genuine sample-derived ancient DNA fragments from false signals [34].
91 Additionally, estimates of ancient DNA preservation are often based on a small number of
92 sequencing reads [35,36] and index hopping could therefore result in an erroneous inference about
93 DNA preservation in a given sample. In the context of the study presented here, it is also worth
94 noting that the small insert size in such aDNA samples allows us to obtain sequence information from
95 both ends of the DNA fragment and thereby we can identify both barcodes with high accuracy.

96 Following sequencing on the HiSeq X platform, we identified a small fraction of reads (<1%) with a
97 wrong barcode-index combination. After excluding several possible explanations, we conclude that
98 index hopping happens in this system, but results in a similar rate of read misassignment as
99 previously reported for older versions of Illumina sequencing platforms that rely on traditional bridge
100 amplification for cluster generation. We therefore recommend using inline barcode-containing
101 sequencing adapters, independent of the sequencing platform in studies that rely on low-coverage

102 data, require absolute certainty, aim to characterize rare variants or combine a large number of
103 samples that exceed available index combinations.

104 **Methods**

105 **Library preparation and sequencing**

106 DNA extracts from seven historical gorilla samples were turned into sequencing libraries following
107 the strategy outlined in [28,37] (see supplementary material). All library preparation steps except
108 indexing PCR were performed in a dedicated ancient DNA facility to minimize contamination. Briefly,
109 20 μ l DNA extract was used in a 50 μ l blunting reaction together with USER enzyme treatment to
110 remove uracil bases resulting from aDNA damage [38]. DNA fragments within each sample were then
111 ligated to a unique combination of incomplete, partially double-stranded P5- and P7-adapters, each
112 containing a unique seven base pair sequence [37] (Table S1). We refer to these as the P5 and P7
113 barcodes from here on. All barcode sequences were at least three nucleotides apart from each other
114 to ensure high certainty during demultiplexing and avoid converting one barcode into another
115 through sequencing errors [37] (Table S1). To increase the complexity of the pooled sequencing
116 library, one sample (sample 7) was split in two fractions, each of which received a different barcode
117 combination (Table 1).

118 Indexing PCR was performed for 10 cycles using a unique P7 indexing primer for each sample, as in
119 [2] (Table S1). We refer to the unique sequence added during the indexing PCR as the P7 index. As
120 with the barcodes, all index sequences differed by at least three base pairs from each other (Table
121 S1). Indexing PCR for sample 7 was performed in a single reaction combining both fractions of this
122 sample. Following the indexing PCR, each DNA fragment contained three unique identifiers: the P5
123 and P7 barcodes directly ligated to the ends of the DNA fragments, and the P7 index, which becomes
124 part of the Illumina sequencing adapter (Figure 1). Sample libraries were cleaned using MinElute spin
125 columns, fragment length distribution and concentrations were measured on the Bioanalyzer. We
126 then pooled all seven sample libraries in a ratio of 2:1:2:1:1:1:2 for samples 1 to 7, and performed
127 two rounds of AMPure XP bead clean-up, using 0.5X and 1.8X bead:DNA ratio, respectively. We
128 confirmed that indexing primers were successfully removed during clean-up by running the final

129 library on a Bioanalyzer (Figure S2). The pooled library was sequenced on three HiSeq X lanes that
130 were part of independent runs with a 5% phiX spike-in, at the SciLife sequencing facility in Stockholm.

131

132 **Data processing**

133 All reads were demultiplexed based on their unique indices using Illumina's bcl2fastq (v2.17.1)
134 software with default settings, allowing for one mismatch per index and only retaining "pass filter"
135 reads (Illumina Inc.). All unidentified reads, i.e. reads containing indices not used in our experiment,
136 were subjected to the same filtering steps, as described below. We removed adapter sequences
137 using AdapterRemoval V2.1.7 with standard parameters [39]. Due to the fragmented nature of DNA
138 in historical samples, we could subsequently merge the reads, requiring a minimal overlap of 11bp
139 and allowing for a 10% error rate. The merging of reads allows us to obtain sequencing information
140 for the complete DNA molecule and thus accurate identification of the barcodes on both ends of the
141 DNA fragment (P5 and P7 barcodes, respectively, Figure 1). Unmerged reads and reads below 29 bp
142 were removed. To increase certainty, we only retained reads with error-free P5 and P7 barcodes and
143 an average quality score of at least 30 using prinseq V0.20.4 [40].

144

145 **Disentangling cross-contamination from index hopping**

146 Low rates of cross-contamination between barcodes and indexes can be expected, even if strict
147 measure, such as clean-room facilities, are taken during library preparation [4]. This can result in
148 reads containing a wrong index-barcode pair and could be falsely taken as evidence for index
149 hopping. Since the inline barcodes used in this experiment are unaffected by index-hopping (Figure
150 1), we can accurately estimate the rate of barcode cross-contamination as the fraction of reads
151 containing a P5-P7 barcode pair that was not used during library preparation. In rare cases, barcode
152 cross-contamination results in a read with a valid barcode pair (e.g. a barcode combination that was

153 intendedly used during library preparation) and thereby remain undetected in our estimate.
154 However, since we used every barcode only once, the probability of such an event is several orders
155 of magnitude lower than the fraction of reads containing an invalid barcode pair and does therefore
156 not significantly affect any of our estimates (see supplementary material).

157 As the Illumina HiSeq X platform did not support a dual-indexing design at the time of this
158 experiment, estimating the rate of index cross-contamination could not be based on invalid index
159 pairs. Therefore, we relied on the fact that only seven out of the 40 indices that are routinely used in
160 our laboratory, were implemented in this experiment (Table S3). Assuming a relative equal rate of
161 cross-contamination between all 40 indexes, we estimated index cross-contamination as the fraction
162 of reads containing any of the 33 indices that were not included during library preparation.

163 We then determined the raw rate of index hopping as the fraction of reads showing an index-
164 barcode combination not used during the library preparation. We accounted for the possibility of
165 barcode and index cross-contamination resulting in the same barcode-index combination by
166 subtracting the contamination estimates obtained above from the raw value of index hopping. All
167 statistical analyses were performed in R 2.15.3 [41] (see supplementary material).

168

169 **Results**

170 Our sequencing libraries were made from degraded historical samples and thus contained a large
171 proportion of short DNA fragments (Figure S3A). Therefore, the majority of reads could be
172 confidently merged for all three sequencing runs (95.3% SE \pm 1.0%). After all filtering steps (see
173 Methods), the final dataset contained 89.3% SE \pm 1.9% of the original sequence reads.

174

175 **Barcode and index cross-contamination**

176 We estimate the levels of barcode cross-contamination at 0.0276% SE \pm 0.0026 across all three runs
177 (Table 1, Table S2, Figure S1), with different rates observed between samples (global chi-square test,
178 $P < 10^{-15}$). The high observed level of incorrect barcode-pairings in sample seven, can be explained by
179 formation of chimeric reads during pooled amplification of the two different fractions of this sample
180 that were barcoded separately (see supplementary material). Assuming that adapter ligation of
181 barcodes is unbiased with respect to the barcode sequence [37], this low percentage of cross-
182 contamination will lead to a neglectable fraction of reads ($1.09 \cdot 10^{-8}\%$, see supplementary material)
183 with barcode pair that wrongly appear as having undergone index hopping.

184 The rate of index cross-contamination was estimated at 0.124% SE \pm 0.0023 (Table S3), by
185 quantifying the fraction of reads containing indices that were not intentionally used in our
186 experiment (see Methods, Table S3).

187

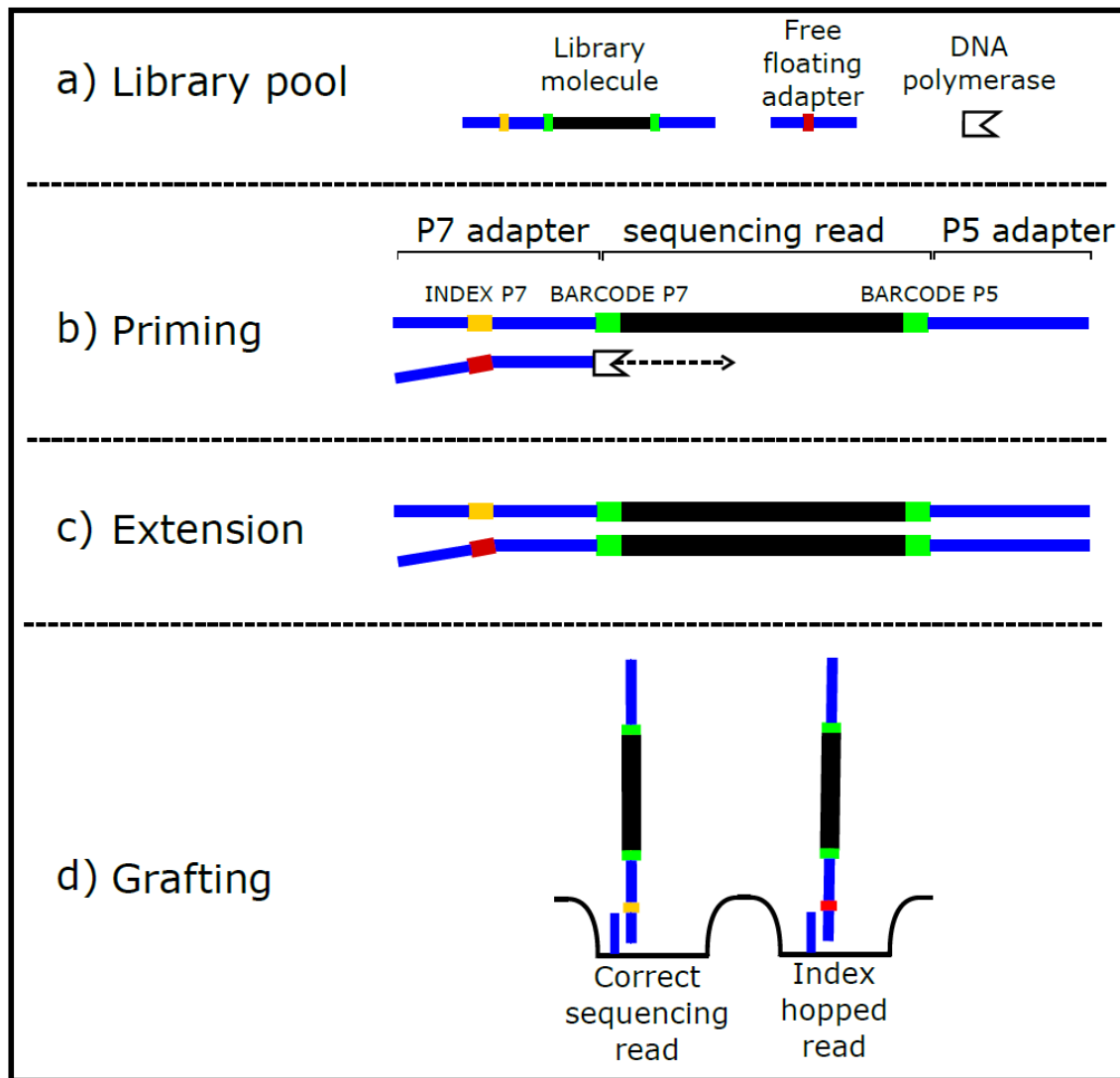
188 **Index hopping**

189 Index hopping will not affect the barcodes that are directly ligated to the DNA fragments. Therefore,
190 it can be readily distinguished from barcode cross-contamination by identifying reads containing a
191 combination between an index and a barcode pair that was not used during the library preparation.
192 Across all three sequencing runs, we detected a low proportion of such reads (mean=0.594%, SE \pm

193 0.0434%, Table 1). As previously estimated $\sim 0.124\%$ of these reads are a result of index and barcode
194 cross-contamination. Therefore, the estimated rate of index hopping in our experiment across all
195 three sequencing runs is $\sim 0.470\% \text{ SE} \pm 0.044$ (0.594% minus 0.124%). The proportion of hopped
196 reads differed significantly by sample (chi-square test, $P < 10^{-15}$) and was positively correlated with the
197 number of sequenced reads per sample (Pearson's $r = 0.96$, $P = 0.0005$). This suggests that in
198 multiplexed sequencing runs the samples with higher number of sequenced reads will serve as the
199 dominant source of hopped reads (Figure 2). Even though the overall rate of index hopping is low,
200 samples with proportionally few sequenced reads are thus considerably more affected by index
201 hopping. In our experiments, this resulted in $2.49\% \text{ SE} \pm 0.29\%$ of index hopped reads in the sample
202 with the lowest number of sequenced reads (Table 1, Table S4, Figure 2).

203 In addition, the rate of index hopping differed significantly by read length and GC content. Reads
204 shorter than 90 bp and reads with GC content above 40% showed significantly higher proportion of
205 hopped reads than expected under a random distribution (chi-square test, both $P < 10^{-15}$, Figure S3).

206

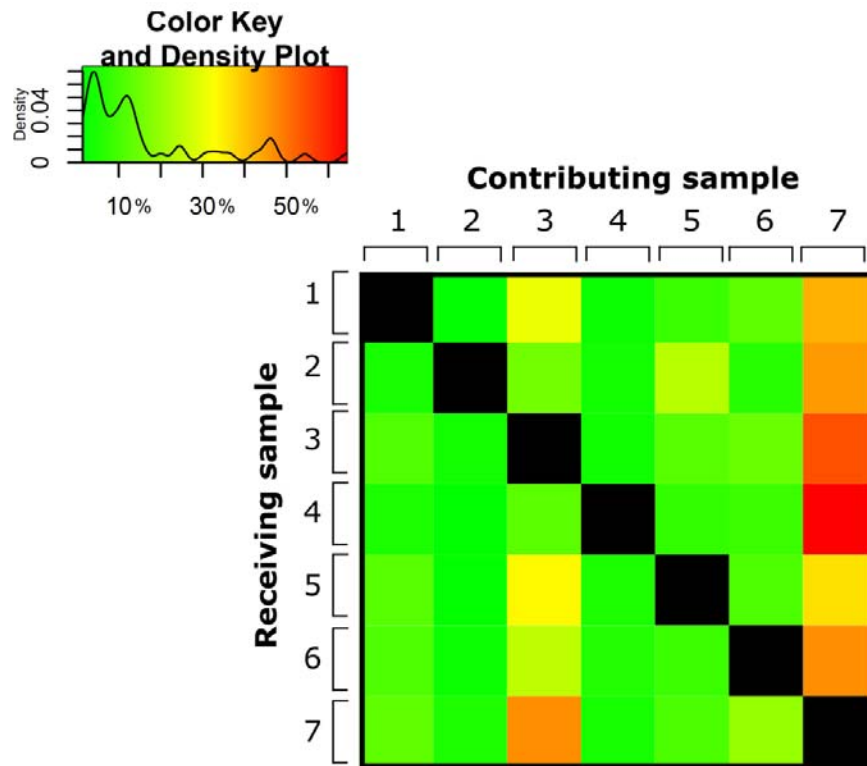


207

208 *Figure 1: Outcome of index hopping. A) The library pool, containing barcoded and indexed library*
209 *molecules and free-floating indexing primers, is mixed with ExAmp reagents before loading on the*
210 *patterned flow cell. B) Free-floating adapters anneal to the adapter sequence of a library molecule*
211 *and C) the library molecule subsequently gets extended by DNA polymerase forming a new library*
212 *molecule containing a wrong index. D) The library molecules are denatured, separating the strands,*
213 *and each library molecule is allowed to graft onto a nanowell on the patterned flow cell.*

214

215



216

217

218 *Figure 2, Proportion of hopped reads per sample out of all hopped reads. Samples in the top row*

219 *contribute hopped reads, whereas samples on the left receive hopped reads. Samples with high*

220 *number of reads (e.g. 3 and 7) are also the main contributors of hopped reads.*

221

222 *Table 1: Sequencing statistics and estimates of contamination and index hopping.*

Sample	Used P5 barcode	Used P7 barcode	Original reads (Millions)			Reads after quality filtering (Millions)			Reads with wrong barcode pairs (cross-contamination)			Reads with wrong index-barcode combination			Cross contaminated reads (%)			Index hopped reads (%)		
			Run 1	Run 2	Run 3	Run 1	Run 2	Run 3	Run 1	Run 2	Run 3	Run 1	Run 2	Run 3	Run 1	Run 2	Run 3	Run 1	Run 2	Run 3
1	3	3	40.63	14.74	53.95	34.49	13.01	50.43	2042	971	3580	158905	87697	280808	0.0059	0.0075	0.0071	0.4587	0.6693	0.5538
2	4	4	11.28	4.83	14.16	9.80	4.37	13.51	1447	644	1398	130125	77502	186512	0.0148	0.0147	0.0103	1.3100	1.7435	1.3622
3	5	5	127.21	44.98	157.84	104.79	40.40	147.48	12939	4184	6937	265347	163502	572787	0.0123	0.0104	0.0047	0.2526	0.4031	0.3869
4	6	6	13.80	4.53	18.76	11.31	3.88	17.05	1831	865	2314	262764	127417	349123	0.0162	0.0223	0.0136	2.2700	3.1824	2.0061
5	7	7	22.69	10.27	34.67	20.00	9.13	32.23	7555	3676	9563	166238	86408	308164	0.0378	0.0402	0.0297	0.8245	0.9374	0.9472
6	8	8	30.46	13.57	39.58	27.56	12.65	38.09	2034	1245	2449	78427	39006	161356	0.0074	0.0098	0.0064	0.2838	0.3073	0.4219
7	9	9	125.50	49.63	142.72	108.26	44.33	130.73	63867	26727	62677	481435	256394	881862	0.0590	0.0603	0.0479	0.4427	0.5751	0.6700
Unidentified	14	14	21.55	9.43	33.00	14.44	7.19	18.14	7860	4145	5609	-	-	-	0.0544	0.0576	0.0309	-	-	-
Total	-	-	393	152	495	331	135	448	99575	42457	94527	1543241	837926	2740612	-	-	-	-	-	-
Average barcode cross-contamination (%)	-	-	-	-	-	-	-	-	0.0301	0.0315	0.0211	-	-	-	-	-	-	-	-	-
Average index hopping (%)	-	-	-	-	-	-	-	-	-	-	-	0.488	0.656	0.638	-	-	-	-	-	-

223
224

225 Discussion

226 We show that index hopping occurs on the Illumina HiSeq X platform, but its rate is low in our study
227 (0.470% SE \pm 0.044). Although multiple sources of error such as jumping PCR, barcode and index
228 cross-contamination, sequencing errors, and index hopping can result in read misassignment we
229 could systematically address each of them through a careful experimental design. Jumping PCR can
230 be eliminated as explanation for wrong index-barcode combinations, as we avoided amplification of
231 pooled libraries from different samples. We further show that the rate of barcode and index cross-
232 contamination is very low (0.027% SE \pm 0.0026 and 0.124% \pm 0.0023, respectively) and therefore not
233 the primary cause of the observed proportion of reads with wrong index-barcode pairs.

234 Read misassignment is not a novel phenomenon for the Illumina sequencing platforms. Reported
235 error rates range from 0.1% to 0.582% for HiSeq 2500 [4,6,42] and from 0.06% to 0.51% for the
236 MiSeq platforms [5,7,8]. It is therefore noteworthy that the fraction of hopped reads as estimated in
237 our study (0.470%) is similar to that reported for other platforms. However, it markedly differs from
238 previous estimates for the Illumina HiSeq platforms with ExAmp chemistry [12–14]. Since the
239 sequencing chemistry of the Illumina NovaSeq platform is identical to that used for the HiSeq X, this
240 platform is likely to be affected at a similar rate as reported here.

241 We propose that the low rate of index hopping achieved in our experiment is the result of rigorous
242 removal of free-floating adapters through size selection and library clean-up (Figure S2). Therefore,
243 as previously suggested, strict library clean-up, such as those performed in this study, should be
244 performed in multiplexed sequencing studies [27].

245 A novel observation in our study is that the number of hopped reads within a sample is proportional
246 to the total number of reads contributed by the sample to the pooled sequencing library.
247 Importantly, pooling samples in different quantities leads to a greater proportion of hopped reads
248 into samples with fewer sequenced reads. In this study, libraries with the lowest number of
249 sequenced reads (e.g. sample 2 and 4) displayed up to 3.2% of misassigned reads (Table 1), an order

250 of magnitude higher than the average rate within a lane. The effect of this skewed rate of index
251 hopping becomes even more severe if the endogenous content is markedly different between
252 samples, as is often observed in aDNA studies [43–45]. Since the endogenous content is usually not
253 known beforehand, pooling samples in equal quantities can lead to large differences in the number of
254 endogenous reads between samples. In these cases, the proportion of false assigned endogenous
255 reads can reach rates above 10% (Figure S4). Therefore, variation in sample DNA quantities should be
256 minimized when sequencing a sample pool on the same lane. Ideally, pre-pooling qPCR quantification
257 of sample DNA should be performed to balance the sequencing libraries. Additionally, when samples
258 are sequenced at high depth (i.e. across multiple lanes/flowcells), re-pooling should be considered
259 after the first sequencing run if high variation in (endogenous) read numbers is observed. This is
260 especially relevant for the NovaSeq platform, the most powerful sequencing platform currently
261 available, since it has been specifically designed for the multiplexing of up to hundreds of samples.
262 Additionally, we detected a higher rate of index hopping among the shorter reads and small
263 differences in the fraction of index-hopped reads related to read GC-content. This suggests that the
264 annealing of free floating adapters present in the sequencing libraries does not occur randomly. The
265 underlying mechanisms are not yet well understood but could be related to differences in the DNA
266 denaturation temperatures between DNA fragments of different size. Due to the lower denaturation
267 temperature, short fragments might be occurring at a higher rate in single stranded conformation
268 and are thereby more accessible to free floating index primers.

269 The recent addition of dual-indexing on the HiSeq X platform allows for the mitigation of the majority
270 of errors resulting from index hopping. Nonetheless, in cases where low coverage data is generated
271 or absolutely certainty is required, even a low remaining rate of misassigned reads might represent a
272 major problem. Our study outlines a strategy for identification, quantification and removal of
273 misassigned reads that result from index hopping or other sources of error. We therefore
274 recommend the use of short barcoded in-line adapters when preparing pooled libraries for next

275 generation sequencing, independently of sequencing platform, if a very high degree of certainty is
276 required.

277 **Author contributions:** TvdV and MO performed the wetlab experiments. TvdV and FV performed
278 computational data-analysis. TvdV and KG established the experimental design. TvdV, LD and KG
279 conceived the study, interpreted the results and wrote the manuscript (with contribution of all co-
280 authors).

281

282 **Acknowledgments:**

283 We acknowledge the support from the Science for Life Laboratory, the Knut and Alice Wallenberg
284 Foundation, the National Genomics Infrastructure funded by the Swedish Research Council, and
285 Uppsala Multidisciplinary Center for Advanced Computational Science for assistance with massively
286 parallel sequencing and access to the UPPMAX computational infrastructure. We thank Illumina for
287 providing sequencing reagents. Illumina had no role in study design, data collection and analysis,
288 decision to publish or preparation of the manuscript.

289 **Data access:** Raw sequencing data is available at the European nucleotide archive under accession
290 number XXXX

291 **Funding sources:** This project was supported by FORMAS grant 2015-676 to LD, FORMAS grant 2016-
292 00835 to KG and the Jan Löfqvist Endowments of the Royal Physiographic Society of Lund to KG.

293 **Conflicts of Interest**

294 The authors declare no conflicts of interests.

295

296 **References**

- 297 1. Craig, D.W., Pearson, J. V, Szeling, S., Sekar, A., Redman, M., Corneveaux, J.J., Pawlowski,
298 T.L., Laub, T., Nunn, G., Stephan, D.A., *et al.* (2008). Identification of genetic variants using
299 bar-coded multiplexed sequencing. *Nat. Methods* 5, 887–893.
- 300 2. Meyer, M., and Kircher, M. (2010). Illumina sequencing library preparation for highly
301 multiplexed target capture and sequencing. *Cold Spring Harb. Protoc.* 5, pdb.prot5448.
- 302 3. Smith, A.M., Heisler, L.E., St.Onge, R.P., Farias-Hesson, E., Wallace, I.M., Bodeau, J., Harris,
303 A.N., Perry, K.M., Giaever, G., Pourmand, N., *et al.* (2010). Highly-multiplexed barcode
304 sequencing: an efficient method for parallel analysis of pooled samples. *Nucleic Acids Res.* 38,
305 e142–e142.
- 306 4. Kircher, M., Sawyer, S., and Meyer, M. (2012). Double indexing overcomes inaccuracies in
307 multiplex sequencing on the Illumina platform. *Nucleic Acids Res.* 40, e3.
- 308 5. Renaud, G., Stenzel, U., Maricic, T., Wiebe, V., and Kelso, J. (2015). deML: robust
309 demultiplexing of Illumina sequences using a likelihood-based approach. *Bioinformatics* 31,
310 770–2.
- 311 6. Wright, E.S., and Vetsigian, K.H. (2016). Quality filtering of Illumina index reads mitigates
312 sample cross-talk. *BMC Genomics* 17, 876.
- 313 7. Nelson, M.C., Morrison, H.G., Benjamino, J., Grim, S.L., and Graf, J. (2014). Analysis,
314 Optimization and Verification of Illumina-Generated 16S rRNA Gene Amplicon Surveys. *PLoS*
315 *One* 9, e94249.
- 316 8. D’Amore, R., Ijaz, U.Z., Schirmer, M., Kenny, J.G., Gregory, R., Darby, A.C., Shakya, M., Podar,
317 M., Quince, C., and Hall, N. (2016). A comprehensive benchmarking study of protocols and
318 sequencing platforms for 16S rRNA community profiling. *BMC Genomics* 17, 55.

- 319 9. Illumina Inc. (2017). Effects of Index Misassignment on Multiplexing and Downstream
320 Analysis.
- 321 10. Bentley, D.R., Balasubramanian, S., Swerdlow, H.P., Smith, G.P., Milton, J., Brown, C.G., Hall,
322 K.P., Evers, D.J., Barnes, C.L., Bignell, H.R., *et al.* (2008). Accurate whole human genome
323 sequencing using reversible terminator chemistry. - Supplement. *Nature* 456, 53–9.
- 324 11. Vodak, D., Lorenz, S., Nakken, S., Aasheim, L.B., Holte, H., Bai, B., Myklebost, O., Meza-
325 Zepeda, L.A., and Hovig, E. (2017). Sample-Index Misassignment Impacts Tumor Exome
326 Sequencing. *bioRxiv*.
- 327 12. Griffiths, J.A., Lun, A.T.L., Richard, A.C., Bach, K., and Marioni, J.C. (2017). Detection and
328 removal of barcode swapping in single-cell RNA-seq data. *bioRxiv*.
- 329 13. Sinha, R., Stanley, G., Gulati, G.S., Ezran, C., Travaglini, K.J., Wei, E., Chan, C.K.F., Nabhan, A.N.,
330 Su, T., Morganti, R.M., *et al.* (2017). Index Switching Causes “Spreading-Of-Signal” Among
331 Multiplexed Samples In Illumina HiSeq 4000 DNA Sequencing. *bioRxiv*.
- 332 14. Owens, G.L., Todesco, M., Drummond, E.B.M., Yeaman, S., and Rieseberg, L.H. (2017). A Novel
333 Post Hoc Method For Detecting Index Switching Finds No Evidence For Increased Switching On
334 The Illumina HiSeq X. *bioRxiv*.
- 335 15. Schmitt, M.W., Kennedy, S.R., Salk, J.J., Fox, E.J., Hiatt, J.B., and Loeb, L.A. (2012). Detection of
336 ultra-rare mutations by next-generation sequencing. *Proc. Natl. Acad. Sci. U. S. A.* 109, 14508–
337 13.
- 338 16. Flaherty, P., Natsoulis, G., Muralidharan, O., Winters, M., Buenrostro, J., Bell, J., Brown, S.,
339 Holodniy, M., Zhang, N., and Ji, H.P. (2012). Ultrasensitive detection of rare mutations using
340 next-generation targeted resequencing. *Nucleic Acids Res.* 40, e2–e2.
- 341 17. Trapnell, C., Hendrickson, D.G., Sauvageau, M., Goff, L., Rinn, J.L., and Pachter, L. (2013).
342 Differential analysis of gene regulation at transcript resolution with RNA-seq. *Nat. Biotechnol.*

- 343 31, 46–53.
- 344 18. Green, R.E., Krause, J., Briggs, A.W., Maricic, T., Stenzel, U., Kircher, M., Patterson, N., Li, H.,
345 Zhai, W.W., Fritz, M.H.Y., *et al.* (2010). A Draft Sequence of the Neandertal Genome. *Science*
346 328, 710–722.
- 347 19. Nielsen, R., Korneliussen, T., Albrechtsen, A., Li, Y., and Wang, J. (2012). SNP calling, genotype
348 calling, and sample allele frequency estimation from new-generation sequencing data. *PLoS*
349 *One* 7, e37558.
- 350 20. Fumagalli, M., Vieira, F.G., Korneliussen, T.S., Linderoth, T., Huerta-Sánchez, E., Albrechtsen,
351 A., and Nielsen, R. (2013). Quantifying population genetic differentiation from next-
352 generation sequencing data. *Genetics* 195, 979–992.
- 353 21. Therkildsen, N.O., and Palumbi, S.R. (2017). Practical low-coverage genomewide sequencing
354 of hundreds of individually barcoded samples for population and evolutionary genomics in
355 nonmodel species. *Mol. Ecol. Resour.* 17, 194–208.
- 356 22. Meyerhans, A., Vartanian, J.P., and Wain-Hobson, S. (1990). DNA recombination during PCR.
357 *Nucleic Acids Res.* 18, 1687–1691.
- 358 23. Odelberg, S.J., Weiss, R.B., Hata, A., and White, R. (1995). Template-switching during DNA
359 synthesis by *Thermus aquaticus* DNA polymerase I. *Nucleic Acids Res.* 23, 2049–57.
- 360 24. Carlsen, T., Aas, A.B., Lindner, D., Vrålstad, T., Schumacher, T., and Kausrud, H. (2012). Don't
361 make a mista(g)ke: Is tag switching an overlooked source of error in amplicon pyrosequencing
362 studies? *Fungal Ecol.* 5, 747–749.
- 363 25. Esling, P., Lejzerowicz, F., and Pawlowski, J. (2015). Accurate multiplexing and filtering for
364 high-throughput amplicon-sequencing. *Nucleic Acids Res.* 43, 2513–2524.
- 365 26. Mitra, A., Skrzypczak, M., Ginalski, K., and Rowicka, M. (2015). Strategies for achieving high

- 366 sequencing accuracy for low diversity samples and avoiding sample bleeding using Illumina
367 platform. *PLoS One* 10, e0120520.
- 368 27. Illumina Inc. (2017). Effects of Index Misassignment on Multiplexing and Downstream
369 Analysis.
- 370 28. Rohland, N., and Reich, D. (2012). Cost-effective, high-throughput DNA sequencing libraries
371 for multiplexed target capture. *Genome Res.* 22, 939–946.
- 372 29. Mulligan, C.J. (2005). Isolation and Analysis of DNA from Archaeological, Clinical, and Natural
373 History Specimens. *Methods Enzymol.* 395, 87–103.
- 374 30. Sawyer, S., Krause, J., Guschanski, K., Savolainen, V., and Paabo, S. (2012). Temporal patterns
375 of nucleotide misincorporations and DNA fragmentation in ancient DNA. *PLoS One* 7, e34131.
- 376 31. Allentoft, M.E., Sikora, M., Sjögren, K.-G., Rasmussen, S., Rasmussen, M., Stenderup, J.,
377 Damgaard, P.B., Schroeder, H., Ahlström, T., Vinner, L., *et al.* (2015). Population genomics of
378 Bronze Age Eurasia. *Nature* 522, 167–172.
- 379 32. Haak, W., Lazaridis, I., Patterson, N., Rohland, N., Mallick, S., Llamas, B., Brandt, G.,
380 Nordenfelt, S., Harney, E., Stewardson, K., *et al.* (2015). Massive migration from the steppe
381 was a source for Indo-European languages in Europe. *Nature* 522, 207–211.
- 382 33. Wall, J.D., and Kim, S.K. (2007). Inconsistencies in Neanderthal Genomic DNA Sequences. *PLoS*
383 *Genet.* 3, e175.
- 384 34. Skoglund, P., Northoff, B.H., Shunkov, M. V, Derevianko, A.P., Pääbo, S., Krause, J., Jakobsson,
385 M., and Klein, R.G. Separating endogenous ancient DNA from modern day contamination in a
386 Siberian Neandertal.
- 387 35. Slon, V., Hopfe, C., Weiß, C.L., Mafessoni, F., de la Rasilla, M., Lalueza-Fox, C., Rosas, A.,
388 Soressi, M., Knul, M. V, Miller, R., *et al.* (2017). Neandertal and Denisovan DNA from

- 389 Pleistocene sediments. *Science* 356, 605–608.
- 390 36. Pečnerová, P., Díez-del-Molino, D., Dussex, N., Feuerborn, T., von Seth, J., van der Plicht, J.,
391 Nikolskiy, P., Tikhonov, A., Vartanyan, S., and Dalén, L. (2017). Genome-Based Sexing Provides
392 Clues about Behavior and Social Structure in the Woolly Mammoth. *Curr. Biol.* 27, 3505–
393 3510.e3.
- 394 37. Rohland, N., Harney, E., Mallick, S., Nordenfelt, S., and Reich, D. (2015). Partial uracil-DNA-
395 glycosylase treatment for screening of ancient DNA. *Philos. Trans. R. Soc. Lond. B. Biol. Sci.*
396 370, 20130624.
- 397 38. Briggs, A.W., Stenzel, U., Meyer, M., Krause, J., Kircher, M., and Pääbo, S. (2010). Removal of
398 deaminated cytosines and detection of in vivo methylation in ancient DNA. *Nucleic Acids Res.*
399 38, e87.
- 400 39. Schubert, M., Lindgreen, S., and Orlando, L. (2016). AdapterRemoval v2: rapid adapter
401 trimming, identification, and read merging. *BMC Res. Notes* 9, 88.
- 402 40. Schmieder, R., and Edwards, R. (2011). Quality control and preprocessing of metagenomic
403 datasets. *Bioinformatics* 27, 863–864.
- 404 41. Team R Core (2016). R: A language and environment for statistical computing. R Foundation
405 for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, URL [http://www.R-](http://www.R-project.org/)
406 [project.org/](http://www.R-project.org/).
- 407 42. Wright, E.S., and Vetsigian, K.H. (2016). Inhibitory interactions promote frequent bistability
408 among competing bacteria. *Nat. Commun.* 7, 11274.
- 409 43. Damgaard, P.B., Margaryan, A., Schroeder, H., Orlando, L., Willerslev, E., and Allentoft, M.E.
410 (2015). Improving access to endogenous DNA in ancient bones and teeth. *Sci. Rep.* 5, 11184.
- 411 44. Pinhasi, R., Fernandes, D., Sirak, K., Novak, M., Connell, S., Alpaslan-Roodenberg, S., Gerritsen,

- 412 F., Moiseyev, V., Gromov, A., Raczky, P., *et al.* (2015). Optimal ancient DNA yields from the
413 inner ear part of the human petrous bone. *PLoS One* *10*, e0129102.
- 414 45. van der Valk, T., Lona Durazo, F., Dalén, L., and Guschanski, K. (2017). Whole mitochondrial
415 genome capture from faecal samples and museum-preserved specimens. *Mol. Ecol. Resour.*
416 *17*, e111–e121.
- 417

418 **Supplementary materials for:**

419 **Estimating the rate of index hopping on the Illumina HiSeq X platform**

420 *Tom van der Valk, Francesco Vezzi, Mattias Ormestad, Love Dalén, Katerina Guschanski*

421

422 **This PDF includes:**

423	Material and methods: Library preparation	2 – 4
424	Material and methods: Data processing	5
425	Material and methods: Barcode cross-contamination	6
426	Material and methods: References	7
427	Figure S1	8
428	Figure S2	9
429	Figure S3	10
430	Figure S4	11

431

432 **Methods**

433 **Library preparation and sequencing**

434 DNA extracts from seven historical eastern gorilla samples that previously yielded good sequencing
435 results on the Illumina HiSeq 2500 platform and showed high endogenous content were turned into
436 sequencing libraries following the strategy outlined in (1,2) as detailed below. All library preparation
437 steps except indexing PCR were performed in a dedicated ancient DNA facility to minimize
438 contamination.

439 **1 USER-TREATMENT AND BLUNT-ENDING**

440 We used 20 μ l DNA extract in a 50 μ l blunting reaction together with USER enzyme treatment to
441 remove uracil bases resulting from aDNA damage (3) as follows:

	<i>Volume per sample</i>	<i>Final concentration</i>
H ₂ O	9.89 μ l	
Tango Buffer (10x)	4 μ l	1x
dNTPs (25 mM each)	0.16 μ l	100 μ M each
ATP (100 mM)	0.4 μ l	1 mM
T4 PNK (10 U/ μ L)	2 μ l	0.5 U/ μ l
USER enzyme (1U/ μ L)	2.75 μ l	0.06U/ μ l
T4 DNA Pol (5 U/ μ L)	0.8 μ l	0.1 U/ μ l
DNA	20 μ l	
TOTAL VOLUME	40 μl	

442

443 Samples were incubated for 3 h at 37°C, followed by the addition of 1 μ l T4 DNA polymerase (Thermo
444 Scientific) and incubation at 25°C for 15 min and 12°C for 5 min. We then MinElute purified the
445 reaction according to the manufacturer's protocol and eluted in 22 μ l of EB buffer. DNA fragments
446 within each sample were then ligated to a unique combination of incomplete, partially double-
447 stranded P5- and P7-adapters (10 μ M each), each containing a unique seven base pair sequence as
448 follows:

449

450 **2) ADAPTER LIGATION AND BARCODING**

451

	<i>Volume per sample</i>	<i>Final concentration</i>
H ₂ O	10 µl	
T4 DNA ligase Buffer (10x)	4 µl	1x
PEG-4000 (50%)	4 µl	5%
Barcoded adapters (10 µM)	1 µl each	1.25 µM
T4 DNA ligase (5 U/ul)	1 µl	0.125 U/µl
Blunt-ended DNA	20 µl	
TOTAL VOLUME	40 µl	

452 *1 µl of the barcode P5.F-mix (10 µM) and 1 µL of the barcode P7.F-mix (10 µM).

453 Samples were incubated for 30 minutes at room temperature. We refer to the 7 basepair adapters as
454 the P5 and P7 barcodes from here on. All barcode sequences were at least three nucleotides apart
455 from each other to ensure high certainty during demultiplexing and avoid converting one barcode
456 into another through sequencing errors (2, Table S1). To increase the complexity of the pooled
457 sequencing library, one sample received two different barcode combinations (Table 1). Following
458 adapter ligation reactions were MinElute purified according to the manufacturer's protocol and
459 eluted in 22 µl of EB buffer. We then performed adapter fill-in as follows:

460

461 **3) ADAPTER FILL-IN**

	<i>Volume per sample</i>	<i>Final concentration</i>
H ₂ O	14.1 µl	
Isothermal Amp. buffer 10x	4 µl	1x
dNTPs (25 mM each)	0.4 µl	250 uM each
Bst polymerase 2.0, Large Fragment (8 U/µL)	1.5 µl	0.4 U/µL
Adapter-ligated DNA	20 µl	
TOTAL VOLUME	40 µl	

462

463

464 After setting up the reaction samples were incubated for 20 min. at 37°C, followed by 20 min. at 80°C
465 to inactivate the Bst polymerase and cleaned using MinElute spin columns as above. Indexing PCR
466 was performed for 10 cycles in 125 µl volume using a unique P7 indexing primer for each sample, as
467 follows:

468

<i>PCR cocktail</i>	<i>Volume per 15 µL DNA</i>	<i>Final concentrations</i>
H ₂ O	84 µL	
Pfu Turbo Cx Hotstart DNA Buffer 10x	12.5 µL	1x
dNTPs (25 mM each)	1.25 µL	250 µM each
Indexing primer P7 (10 µM)	5 µL	0.4 µM
Primer IS4 (10 µM)	5 µL	0.4 µM
Pfu Turbo Cx Hotstart DNA Polymerase (2.5 U/µL)	2.5 µL	5 U
Adapter-ligated DNA	15 µL	
TOTAL VOLUME	125 µL	

469 Cycling conditions:

470 *2 min. 95°C, (95°C- 30sec. 59°C- 30sec. 72°C- 1min), 7min. 72°C, hold 8°C*

471

472

473

474 **Data processing**

475 All reads were demultiplexed based on their unique indices using Illumina's bcl2fastq (v2.17.1)
476 software with defaults settings, allowing for one mismatch per index and only retaining "pass filter"
477 reads (Illumina Inc.). All unidentified reads, i.e. reads with indices that were not used in our
478 experiment, were subjected to the same filtering steps, as described below. We removed adapter
479 sequences using AdapterRemoval V2.1.7 using standard parameters and subsequently merged the
480 reads, requiring a minimal overlap of 11bp and allowing for a 10% sequencing error rate ⁽⁴⁾.
481 Unmerged reads and reads below 29 bp were removed leaving only merged reads with an original
482 insert size of at least 15 bp (7 bp barcodeP7 + 7 bp barcodeP5 + 15 bp DNA fragment = 29 bp). To
483 increase certainty, we only retained reads with intact and error-free P5 and P7 barcodes (assessed
484 using an in-house python script available upon request) and an average quality score of at least 30
485 using prinseq V0.20.4 ⁽⁵⁾.

486

487 **Statistical analyses**

488 Statistical analyses were performed in R 2.15.3 ⁽⁶⁾. Significant global chi-square tests were followed
489 by a post hoc procedure as implemented in the R package polytomous
490 (<https://artax.karlin.mff.cuni.cz/r-help/library/polytomous/html/00Index.html>). The minimum value
491 of the chi-squared test statistic for the given degrees of freedom was used to assess if individual
492 observed values differ significantly from an overall hypothetical homogeneous distribution. The test
493 also identified the direction of these differences.

494

495 **Estimating barcode cross-contamination**

496 We estimated the rate of barcode cross-contamination as the fraction of reads containing a P5-P7
497 barcode pair that was not used during library preparation. There is a small probability that cross-
498 contaminating changes both the P5 and P7 barcode into another valid barcode pair (e.g. a pair that
499 was used during library preparation). The chance of this occurring is proportional to the number of
500 used barcodes and can be estimated as $((n - 1) \cdot \left(\frac{x}{n-1}\right)^2$, where n represents the number of used
501 barcodes (eight in this study) and x the estimated average barcode cross contamination). The fraction
502 of reads with invalid barcode-pairs was 0.0276%, which results in an estimate of $(8 - 1) \cdot$
503 $\frac{0.000276}{8-1}^2 \cdot 100\% = 1.09 \cdot 10^{-8} \%$ of undetected cross-contaminated reads.

504

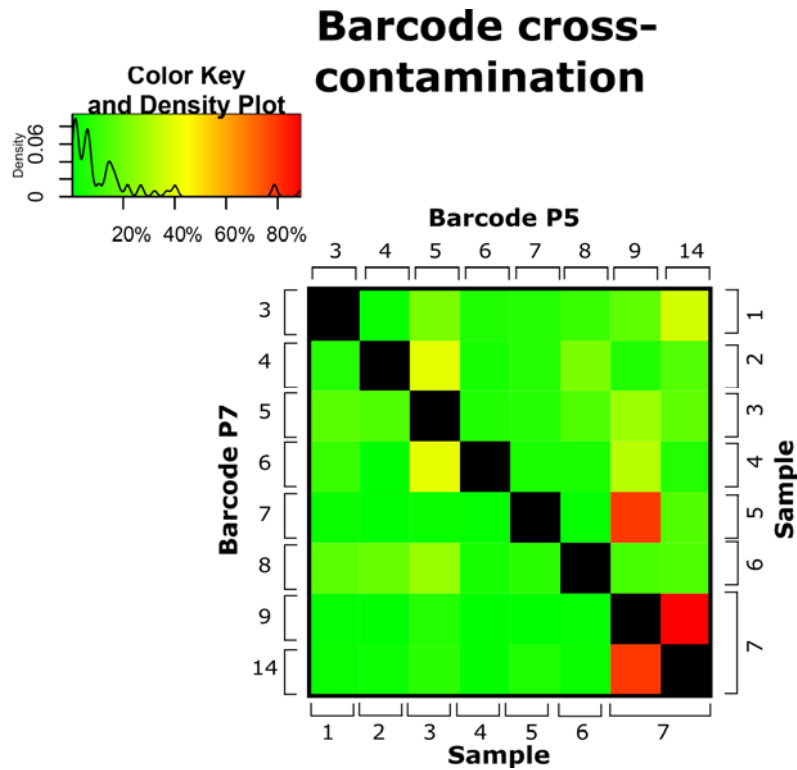
505 **Difference in barcode cross-contamination between samples**

506 The rate of barcode cross-contamination differed significantly by sample (global chi-square test,
507 $P < 10^{-15}$). The implemented posthoc procedure suggested that samples 5 and 7 had significantly more
508 reads with wrong barcode combinations than expected, whereas all the other samples had
509 significantly fewer such reads. Among reads with barcode cross contamination we found an
510 overrepresentation of incorrectly paired barcodes #9 and #14 (Table S2), both of which were used for
511 sample 7 in the following combinations: P5-#9 with P7-#9 and P5-#14 with P7-#14 (Table 1). Elevated
512 cross-contamination between these two barcodes during laboratory procedures could explain the
513 results. However, the observed high rate of wrong barcode pairs (P5-#9 with P7-#14, P5-#14 with P7-
514 #9, Figure S3) is more likely the result of jumping PCR during the 10 rounds of indexing PCR, as both
515 fractions of sample 7 were indexed in a pooled reaction. Equal frequency of wrong barcode pairs is
516 further supporting this notion (Table S2) and can be explained by jumping PCR happening randomly
517 among the reads. In contrast, it is rather unlikely that all four barcodes would have received equal
518 amounts of cross-contamination during laboratory procedures.

519

520

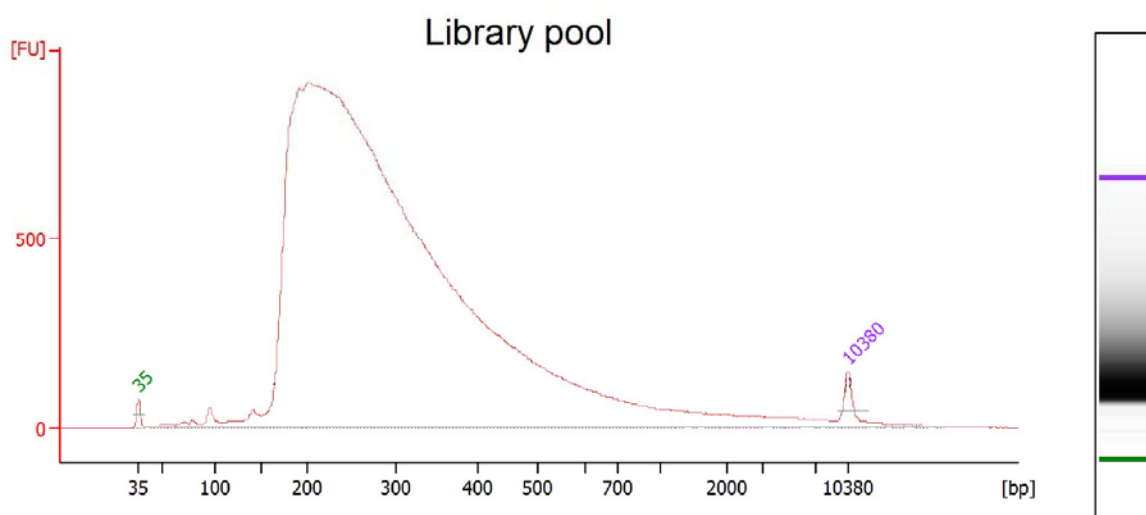
- 521 1. Rohland, N. & Reich, D. Cost-effective, high-throughput DNA sequencing libraries for
522 multiplexed target capture. *Genome Res.* **22**, 939–946 (2012).
- 523 2. Rohland, N., Harney, E., Mallick, S., Nordenfelt, S. & Reich, D. Partial uracil-DNA-glycosylase
524 treatment for screening of ancient DNA. *Philos. Trans. R. Soc. Lond. B. Biol. Sci.* **370**, 20130624
525 (2015).
- 526 3. Briggs, A. W. *et al.* Removal of deaminated cytosines and detection of in vivo methylation in
527 ancient DNA. *Nucleic Acids Res.* **38**, e87 (2010).
- 528 4. Schubert, M., Lindgreen, S. & Orlando, L. AdapterRemoval v2: rapid adapter trimming,
529 identification, and read merging. *BMC Res. Notes* **9**, 88 (2016).
- 530 5. Schmieder, R. & Edwards, R. Quality control and preprocessing of metagenomic datasets.
531 *Bioinformatics* **27**, 863–864 (2011).
- 532 6. Team R Core. R: A language and environment for statistical computing. R Foundation for
533 Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0, URL <http://www.R-project.org/>.
534 (2016).
- 535
- 536



537

538 *Figure S1: Proportion of a given wrong barcode pair in the data out of all cross-contaminated barcode*
539 *pairs. Note that the overall rate of barcode cross-contamination is only 0.0276% (Table 1) Barcodes 9*
540 *and 14 are paired significantly more often due to the formation of chimeric reads during pooled*
541 *amplification of two fractions of sample 7 (see methods).*

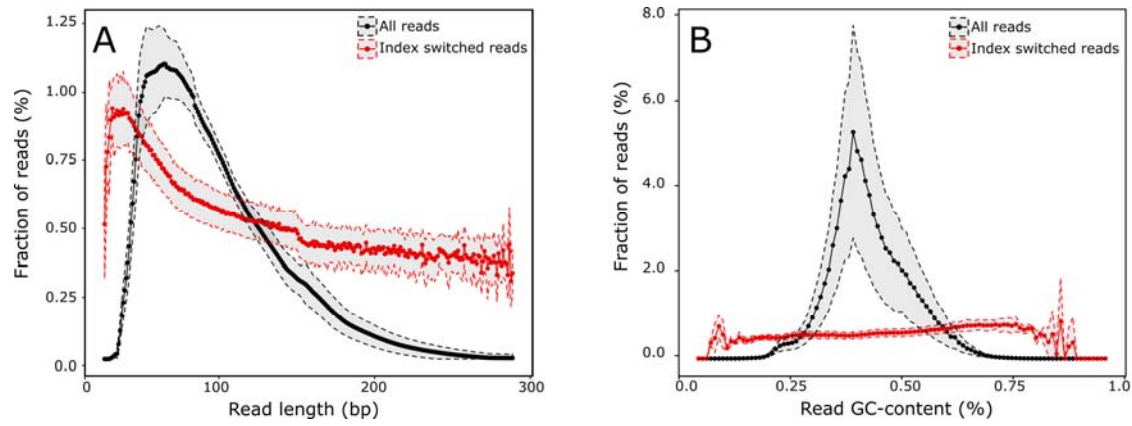
542



543

544 *Figure S2: Bioanalyzer profile of the final pooled library. Note that during library preparation,*
545 *sequencing adapters are attached to the DNA fragments, adding an additional 136 bp to the original*
546 *DNA fragments. The insert size of the DNA is therefore 136 bp lower than what the Bioanalyzer*
547 *shows.*

548



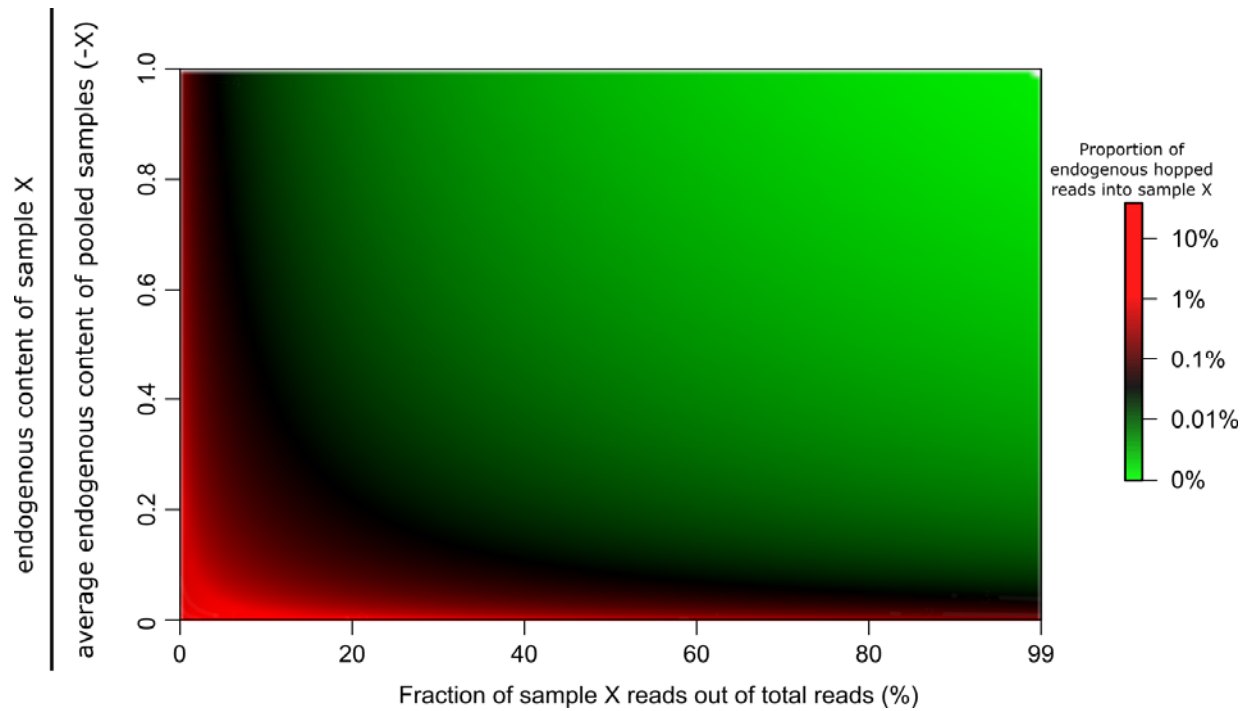
549

550 *Figure S3: A) Read length distribution and the proportion of index hopping by read length. B) Read*
551 *GC-content distribution and the proportion of index hopping by read GC content. Shaded area depicts*
552 *95% confidence interval.*

553

554

555



556

557 *Figure S4. Theoretical relationship between endogenous content, fraction of total reads contributed*
558 *by a given sample (referred to as sample X) to the pooled sequencing library, and index hopping. The*
559 *lower the proportion of reads coming from sample X and the lower its endogenous content compared*
560 *to other samples in the pooled sequencing library, the higher the proportion of endogenous hopped*
561 *reads that sample X will receive from other samples.*

562

563

564