

How Imputation Can Mitigate Ascertainment Bias

Johannes Geibel (✉ johannes.geibel@uni-goettingen.de)

University of Goettingen, Department of Animal Sciences, Animal Breeding and Genetics 6 Group,
Albrecht-Thaer-Weg 3, 37075 Göttingen, Germany

Christian Reimer

University of Göttingen

Torsten Pook

University of Göttingen

Steffen Weigend

Friedrich-Loeffler-Institut

Annett Weigend

Friedrich-Loeffler-Institut

Henner Simianer

University of Göttingen

Research Article

Keywords: SNP ascertainment bias, imputation, chickens, population genetics

Posted Date: January 27th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-150008/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **How Imputation Can Mitigate Ascertainment Bias**

2

3 Johannes Geibel^{12*}, Christian Reimer¹², Torsten Pook¹², Steffen Weigend²³, Annett Weigend³
4 & Henner Simianer¹²

5

6 ¹University of Goettingen, Department of Animal Sciences, Animal Breeding and Genetics
7 Group, Albrecht-Thaer-Weg 3, 37075 Göttingen, Germany; johannes.geibel@uni-goettin-
8 gen.de; creimer@gwdg.de; torsten.pook@uni-goettingen.de; hsimian@gwdg.de

9 ²University of Goettingen, Center for Integrated Breeding Research, Albrecht-Thaer-Weg 3,
10 37075 Göttingen, Germany

11 ³Friedrich-Loeffler-Institut, Institute of Farm Animal Genetics, Höltystrasse 10, 31535 Neu-
12 stadt-Mariensee, Germany; steffen.weigend@fli.de; annett.weigend@fli.de

13

14 * johannes.geibel@uni-goettingen.de (Corresponding Author)

15

16 **Abstract**

17 **Background**

18 Population genetic studies based on genotyped single nucleotide polymorphisms (SNPs) are
19 influenced by a non-random selection of the SNPs included in the used genotyping arrays. The
20 resulting bias relative to whole genome sequencing (WGS) data is known as SNP ascertain-
21 ment bias. Correction for this bias requires detailed knowledge of the array design process
22 which is often not available in practice. This study intends to investigate an alternative ap-
23 proach to mitigate ascertainment bias of a large set of genotyped individuals by using infor-
24 mation of a small set of sequenced individuals via imputation without the need for prior
25 knowledge on the array design.

26 **Results**

27 The strategy was first tested by simulating additional ascertainment bias with a set of 1,566
28 chickens from 74 populations that were genotyped for the positions of the Affymetrix Axiom™
29 580k Genome-Wide Chicken Array. Imputation accuracy was shown to be consistently higher
30 for populations used for SNP discovery during the simulated array design process. Reference
31 sets of at least one individual per population in the study set led to a strong correction of
32 ascertainment bias for estimates of expected and observed heterozygosity, Wrights Fixation
33 Index and Nei's Standard Genetic Distance. In contrast, unbalanced reference sets introduced
34 a new bias towards the reference populations. Finally, the array genotypes were imputed to
35 WGS by utilization of reference sets of 74 individuals (one per population) to 98 individuals
36 (additional commercial chickens) and compared with a mixture of individually and pooled se-
37 quenced populations. The imputation reduced the slope between heterozygosity estimates of
38 array data and WGS data from 1.94 to 1.26 when using the smaller balanced reference panel

39 and to 1.44 when using the larger but unbalanced reference panel. This generally supported
40 the results from simulation but was less favorable, advocating for a larger reference panel
41 when imputing to WGS.

42 **Conclusions**

43 The results highlight the potential of using imputation for mitigation of SNP ascertainment
44 bias but also underline the need for unbiased reference sets.

45 **Keywords**

46 SNP ascertainment bias; imputation; chickens; population genetics

47 **Background**

48 To realize cost- and computational efficiency, many of the population genetic studies of the
49 last 10 years for humans [1, 2], as well as for model- [3, 4] and agricultural species [5–8] were
50 based on single nucleotide polymorphisms (SNP), which were genotyped by commercially
51 available SNP arrays. Those arrays are based on a non-random selection (ascertainment) of
52 SNPs, mostly performed in more intensively researched populations (e.g. commercially used
53 livestock breeds). This results in a shift of the allele frequency spectrum towards mean allele
54 frequencies and thereby an overestimation of heterozygosity compared to whole genome re-
55 sequencing (WGS) data. As this overestimation is stronger for populations involved in the ar-
56 ray design (discovery populations) than for those not involved [9], follow-up analyses can be
57 biased. This effect is widely known as SNP Ascertainment Bias [10–12].

58 Different population genetic estimators are affected by SNP ascertainment bias to a different
59 extent [11] and implementing bias-reduced estimators requires strong assumptions on the

60 design process of the used SNP array [13] which is often not public knowledge or too compli-
61 cated to be remodeled [14, 15]. Malomane *et al.* [16] therefore screened different raw data
62 filtering strategies on mitigation of ascertainment bias in SNP data and identified linkage prun-
63 ing to result in slightly decreasing ascertainment bias. Due to strongly decreasing sequencing
64 costs and the complexity of the ascertainment bias correction strategies, more and more stud-
65 ies started using WGS data for population genetic analysis during the last years [17–21]. How-
66 ever, costs for broad WGS based studies are still rather high, resulting in large-scale collabo-
67 rations such as the 1000 Genomes Project [22], the 1000 Bull Genomes Project [23], or the
68 1001 Arabidopsis Genomes Project [24].

69 A commonly used method to *in silico* increase the resolution of SNP data sets is imputation
70 [25]. Over the years a variety of imputation approaches [26–32] have been proposed that uti-
71 lize linkage, pedigree, and haplotype information. To increase the marker density, an addi-
72 tional reference panel of individuals that were genotyped/sequenced by the intended resolu-
73 tion is required to additionally infer information from SNPs missing on the respective lower
74 density study set.

75 Imputation-based studies mostly either used a reference panel of the same population as the
76 study set itself [33–35] or utilized large global reference panels as of the 1000 Genomes [22,
77 36, 37] or 1000 Bull genomes [23, 38] projects. Especially for admixed or small endangered
78 populations, the use of additional distantly related populations in the reference panel was
79 investigated. On one hand, Brøndum *et al.* [39], Ye *et al.* [40] and Rowan *et al.* [41] identified
80 multi-breed reference panels to increase imputation accuracy especially in admixed breeds
81 and for low frequent alleles when imputing from high-density genotypes to sequence data.
82 On the other hand, Berry *et al.* [42] observed that smaller within breed reference panels

83 (140 - 688 reference cattle individuals per breed) performed always superior compared to the
84 combined across breed reference panel when imputing from low density to high-density array
85 genotypes. Korkuć *et al.* [43] showed that adding 100 to 500 Holstein cattle sequences to a
86 reference panel of 30 German Black Pied cattle significantly decreased the imputation accu-
87 racy in comparison to the pure panel when imputing from array to sequence data. Adding the
88 same numbers of a multi-breed reference panel only outperformed the pure panel when at
89 least 300 reference animals were added. Pook *et al.* [44] investigated the inclusion of chicken
90 populations to the reference set which were differently distantly related to the study set.
91 While error rates generally decreased for rare alleles, the inclusion of distantly related popu-
92 lations slightly increased error rates for previously good imputed SNPs. Overall, the ideal setup
93 of a reference panel seems to be highly dependent on the application with positive effects for
94 some, but also potential harm in other cases.

95 In this context, the current study aims at assessing the influence of a study design on SNP
96 ascertainment bias, which uses a small number of sequenced chickens (the reference set) to
97 *in silico* correct SNP ascertainment bias in a broad multi-population set of genotyped chickens
98 (the study set) by imputation to sequence level. The general idea behind this design is to allow
99 for a large sample size, which reduces sampling bias while keeping sequencing costs affordable
100 as most individuals will only be genotyped. We, therefore, assessed the potential effects of
101 this design by imputing *in silico* created low-density array data to high-density array data, and
102 by imputing real high-density data to WGS data.

103 **Material and Methods**

104 **Data**

105 Three different sets of genomic data were used for this study:

- 106 1) Individual sequence data of 68 chickens from 68 different populations, sequenced
107 within the scope of the EU project Innovative Management of Animal Genetic Re-
108 sources (IMAGE; www.imageh2020.eu) [45]. They were complemented by 25 se-
109 quences (17 + 8) from two commercial white layer lines, 25 sequences (19 + 6) from
110 two commercial brown layer lines, and 40 sequences (20 each) from two commercial
111 broiler lines [20]. In total 158 sequences from 74 populations.
- 112 2) Pooled sequence data from 37 populations (9-11 chickens per population) [16]. All ex-
113 cept 4 chickens from two populations were part of set 3.
- 114 3) Genotypes of 1,566 chickens from 74 populations, either genotyped (sub-set of the
115 Synbreed Chicken Diversity Panel; SCDP) [46] with the Affymetrix Axiom™ 580k Ge-
116 nome-Wide Chicken Array [47], or complemented from set 1.

117 The intersection of the used data sets is shown in Figure 1 and accession information of the
118 raw data per sample can be found in Supplementary File 1. All three data sets came with their
119 own characteristics. While individual sequences are considered to be the gold standard
120 throughout this study, genotypes of the Affymetrix Axiom™ 580k Genome-Wide Chicken Array
121 [47] are biased towards variation which is common in the commercial chicken lines [9] and
122 pooled sequences only allow for an estimate of population allele frequencies and show a slight
123 bias due to sample size and coverage (Supplementary File 2) [48, 49].

124

125 Figure 1: UpSet plot showing the distinct intersections of chickens between the used sequenc-
126 ing/ genotyping technologies. The left bar plot contains the total number of individuals that
127 were genotyped (array), individually sequenced (indSeq), or pooled sequenced (poolSeq). The
128 upper bar plot contains the number of individuals within each distinct intersection, indicated
129 by the connected points below.

130

131 **Calling of WGS SNPs and generation of genotype set**

132 Alignment of the raw sequencing reads against the latest chicken reference genome GRCg6a
133 [50] and SNP calling was conducted for individual and pooled sequenced data following GATK
134 best practices [51, 52]. As the Affymetrix Axiom™ 580k Genome-Wide Chicken Array [47] does
135 not contain enough SNPs on chromosomes 30 – 33 for imputation (and chromosome 29 is not
136 annotated in the reference genome), only up to chromosome 28 was used. This resulted in
137 20,829,081 biallelic SNPs on chromosomes 1 - 28 which were used in further analyses. Addi-
138 tionally, all individual sequences were genotyped for the positions of the Affymetrix Axiom™
139 580k Genome-Wide Chicken Array [47].

140 To ensure compatibility between Array- and WGS data, the genotypes of the Synbreed Chicken
141 Diversity panel were lifted over from galGal5 to galGal6 and corrected for switches of refer-
142 ence and alternate alleles. Only SNPs with known autosomal position, call rates > 0.95 and
143 genotype recall rates > 0.95 were further considered. MAF filters were later used when sub-
144 sampling the different sets and thus not considered in this step. Further, missing genotypes
145 were imputed using Beagle 5.0 [32] with ne=1000 [44] and the genetic map taken from
146 Groenen *et al.* [53]. This resulted in a final set of 1,566 animals from 74 populations (18 - 37
147 animals per population) and 462,549 autosomal SNPs, further referred to as the **genotype set**.

148 As Malomane *et al.* [16] described LD-based pruning as an effective filtering strategy to mini-
149 mize the impact of ascertainment bias in SNP array data, the genotype set was additionally LD
150 pruned using plink 1.9 [54] with --indep 50 5 2 flag. This reduced the genotype set to 136,755
151 SNPs (30 %) and will be referred to as **pruned genotype set**.

152 The description of the detailed pipeline can be found in Supplementary File 2.

153 **Analyses based on simulation of ascertainment bias within the genotype set**

154 A first comparison was based solely on the 15,868 SNPs of chromosome 10 of the genotype
155 set which allowed for a high number of repetitions while still being based on a sufficiently
156 sized chromosome. To simulate an ascertainment bias of known strength, a stronger biased
157 array was designed *in silico* from the genotype set for each of the 74 populations (further
158 called discovery population) by using only SNPs with MAF > 0.05 within this population. Then,
159 reference samples for imputation were chosen in five different ways with ten different num-
160 bers of reference samples and three repetitions per sampling:

- 161 1) allPop_74_740: Equally distributed across all populations by sampling one to ten chick-
162 ens per population (74 - 740 reference samples).
- 163 2) randSamp_5_50: 5, 10, ..., 50 randomly sampled chickens (5-50 reference samples).
- 164 3) randPop_5_50: Five chickens from each of one to ten randomly sampled populations
165 (5 - 50 reference samples).
- 166 4) minPop_5_50: Five chickens from each of one to ten populations which were closest
167 related to the discovery population, based on Nei's Distance ([55]; 5 - 50 reference
168 samples).
- 169 5) maxPop_5_50: Five chickens from each of one to ten populations which were most
170 distantly related to the discovery population, based on Nei's Distance ([55]; 5 - 50 ref-
171 erence samples).

172 This resulted in 2,200 repetitions of *in silico* array development and re-imputation per sam-
173 pling strategy. The reference set was formed by sub-setting the total genotype matrix to SNPs
174 with MAF > 0.01 within the reference samples and the reference samples chosen via the

175 above-mentioned strategies. Imputation of the *in silico* arrays to the reference set was per-
176 formed by running Beagle 5.0 [32] with ne=1000 [44], the genetic distances taken from
177 Groenen *et al.* [53] and the according reference set. The schematic workflow can be found in
178 Figure 2.

179

180 Figure 2: Schematic representation of the workflow of creating and re-imputing the *in silico*
181 arrays. The starting point was a 0/1/2 coded marker matrix with SNPs in rows and individuals
182 in columns (different populations separated by vertical lines). In a first step, an array (light
183 blue rows) was constructed *in silico* from known data by setting all SNPs to missing which were
184 invariable (MAF < 0.05, red rows) in the discovery population (first three columns). In a second
185 step, a reference set (dark blue columns) was set up from animals for which complete
186 knowledge of all SNPs was assumed. This Reference set was then used in a third step to impute
187 the missing SNPs in the study set using Beagle 5.0 and resulting in a certain amount of impu-
188 tation errors (red numbers).

189

190 Analyses were then based on comparisons between the *in silico* ascertained and later imputed
191 sets and the genotype set, which was considered as the 'true' set for those comparisons.

192 **Imputation of genotype set to sequence level**

193 After the initial tests of the imputation strategies by the *in silico* designed arrays, we imputed
194 the complete genotype set to sequence level, using the available individual sequences as the
195 reference panel. In the first run, one reference sample per sequenced population was chosen
196 (74 reference samples; 74_1perLine) which is equivalent to the first scenario allPop_74 of the
197 *in silico* array imputation. As we had more than one sequenced individual for the commercial
198 lines, the number of reference samples for the commercial lines was subsequently increased
199 to five reference samples per line (up to 98 reference samples; 98_5perLine). Finally, we used
200 all available individually sequenced animals as reference samples (158 reference samples;

201 158_all), which resulted in a strong imbalance towards the two broiler lines (20 reference
202 samples per broiler line).

203 Parameter settings in Beagle were further tweaked by increasing the window parameter to
204 200 cM to ensure enough overlap between reference and study SNPs. This was needed as we
205 observed low assembly quality and insufficient coverage of the array on the small chromo-
206 somes. Analyses were then based on comparisons between the genotype set, the pruned set
207 or the imputed sets and the gold standard, the WGS data.

208 **Assessment of imputation accuracy**

209 Assessment of imputation accuracy was done by using Pearson correlation (r) between true
210 and imputed genotypes [42, 56] for the *in silico* designed arrays. Pearson correlation puts a
211 higher relative weight on imputation errors in rare alleles than plain comparison of allele- or
212 genotype concordance rates [56]. In case of the imputation to sequence level, no ‘true’ set
213 was available for the evaluation of imputation accuracy. Instead, we used the internal Beagle
214 quality measure, the dosage r-squared (DR2) [57]. This, however, has the drawback that it only
215 shows the theoretical imputation accuracy, cannot capture biases due to biased reference sets
216 and also does not allow for a per population evaluation of imputation accuracy.

217 **Comparison of population genetic estimators**

218 Ascertainment bias is most crucial when comparing populations with varying relatedness to
219 the populations used for the SNP ascertainment [9]. To investigate these effects, we concen-
220 trated on two heterozygosity estimates: expected (H_E) and observed (H_O) heterozygosity; and
221 two distance measurements: Wright’s fixation index (F_{ST}) [58] and Nei’s distance (D) [55].

222 H_O , as the proportion of heterozygous genotypes in a population, could only be calculated
223 when the genotypic status of a population was known (individual sequences or genotypes). In

224 contrast, H_E could also be calculated from pooled sequences which allow the estimation of
 225 allele frequencies (p). Thereby, H_O and H_E (equation (1)) are calculated as average over all loci
 226 ($l = 1, \dots, L$).

$$227 \quad H_E = \frac{\sum_l 2p_l(1-p_l)}{L} \quad (1)$$

228 As pooled sequence data comes with a slight but systematic underestimation of H_E ([48]; Sup-
 229 plementary File 2), H_E for pooled sequences was multiplied with the correction factor $\frac{n}{n-1}$,
 230 introduced by Futschik and Schlötterer [48], where n is the number of haplotypes in the pool.
 231 This partially corrected the H_E estimates for the bias introduced by pooled sequencing (Sup-
 232 plementary File 2).

233 D was calculated as given by equation (2), where D_{xy} accounts for the genetic distance be-
 234 tween populations X and Y , while x_{il} and y_{il} represent the frequency of the i^{th} allele at the l^{th}
 235 locus in population X and Y , respectively.

$$236 \quad D_{xy} = -\ln \left(\frac{\sum_l \sum_i x_{il} y_{il}}{\sqrt{\sum_l \sum_i x_{il}^2 \sum_l \sum_i y_{il}^2}} \right) \quad (2)$$

237 Pairwise F_{ST} values between populations X and Y were estimated using equation (3), where
 238 HT_l accounts for the H_E within the total population at locus l and \overline{HS}_l for the mean H_E within
 239 the two subpopulations at locus l [58].

$$240 \quad F_{ST} = \frac{\sum_l (HT_l - \overline{HS}_l)}{\sum_l HT_l} \quad (3)$$

241 D and F_{ST} both show a downward bias that is comparable to HE when estimated from pooled
 242 data (Supplementary File 2). The effect of ascertainment bias is much larger than the effect of
 243 pooling for D. In contrast, F_{ST} is generally robust against the effects of ascertainment bias when
 244 a sufficiently large discovery panel was used for array development [11]. Therefore, it shows
 245 underestimation when calculated from pooled sequence data which is larger than the effect
 246 of ascertainment bias (Supplementary File 2). We therefore could not dissect the effects of
 247 the two biases in the comparisons on sequence level and did not include F_{ST} there.

248 Having no ascertainment bias would mean that estimates of a respective set would lie on the
 249 line of identity (diagonal) when regressing the set against the true values. The magnitude of
 250 the bias can therefore be defined as the distance of the estimates to that line. We therefore
 251 regressed the estimates from biased data (y_{ij}) on the unbiased ones (x_{ij}) while fitting group
 252 specific intercepts ($group_i$) as well as group-specific slopes ($group_i \times \beta_i$) and a random error
 253 ($\hat{\delta}_{ij}, \delta \sim N(0, \mathbf{I}\sigma_e^2)$) as in equation (4).

$$254 \quad y_{ij} = group_i + group_i \times \beta_i x_{ij} + \hat{\delta}_{ij} \quad (4)$$

255 The definition of a group describes for within-population estimators (e.g. H_E) whether a pop-
 256 ulation was used for SNP discovery (discovery population), samples from that population were
 257 used as reference set (reference population) or none of both (application population). Note
 258 that in scenarios where reference individuals were present for every population, we only di-
 259 vided them into discovery and application populations. For between population estimators
 260 (F_{ST} , D), a group describes the according combination of the two involved population groups.
 261 Differences of the estimated slopes from one and the correlation between heterozygosity and

262 distance estimates from biased and true set within groups were used as indicators for the
263 magnitude of bias and random estimation error.

264 To get a measure for a fixed estimation error, we also calculated the mean overestimation
265 across populations ($j = 1 \dots J$) as in equation (5).

$$266 \quad \text{mean overestimation} = \frac{\sum_j \frac{\text{biased estimate}_j - \text{true estimate}_j}{\text{true estimate}_j}}{J} \quad (5)$$

267 Note, that we did not have a population-wide estimate of allele frequency on sequence level
268 for each population we did have single individual sequences for. Comparisons of population
269 estimates on sequence level are therefore only based on 45 populations instead of all 74 pop-
270 ulations which were used as study and reference set for the imputation process.

271 **Results**

272 **In silico array to genotype**

273 For the *in silico* array to genotype set imputation, median per-animal imputation accuracies
274 (r) were already high with 0.94 when using one reference individual per population (all-
275 Pop_74). Increasing the number of reference individuals subsequently increased the accuracy
276 up to 0.99 for 10 reference individuals per population (allPop_740). The accuracy was consist-
277 ently higher for individuals which were part of the discovery population (Figure 3).

278

279 Figure 3: Development of the per-animal imputation accuracy for the *in silico* array to geno-
280 type set imputation with an increasing number of reference animals per population. Individu-
281 als are grouped on whether they belong to the population used for SNP discovery or not and
282 reference individuals were chosen as in scenario allPop_74_740. The lines show the trend of

283 the median and outliers are not shown in the plot as they do not add valuable information
284 due to the high number of repetitions.

285

286 Accuracies were lower for the other strategies, mainly due to fewer reference individuals
287 (maximum 50). However, Figure S 2 indicates that the best imputation accuracies were
288 achieved for populations which contained at least one reference individual. The discovery
289 populations showed slightly to strongly increased accuracies compared to the application pop-
290 ulations for all scenarios besides the scenario where the reference populations were chosen
291 to be maximum distant to the discovery population (scenario maxPop_5_50; Figure S 2). Ran-
292 domly choosing reference individuals across the complete population set (scenario rand-
293 Samp_5_50; Figure S 2) resulted in the overall best imputation accuracies. When assigning the
294 reference samples to randomly chosen populations in blocks of five samples per population
295 (scenario randPop_5_50, Figure S 2), the accuracies for the reference populations were only
296 slightly increased, while accuracies for discovery- and application populations dropped mas-
297 sively. Interestingly, the accuracies for application populations were higher in scenario
298 minPop_5_50 (Figure S 2), in which the reference populations were chosen according to their
299 minimum distance to the discovery population, than with the random reference population
300 selection from scenario randPop_5_50. A selection of reference populations to be maximum
301 distant to the discovery population (scenario maxPop_5_50, Figure S 2) performed overall
302 worst.

303 As expected, the *in silico* ascertained sets showed a strong overestimation of the H_E for nearly
304 all populations in all cases. The overestimation was much stronger for populations used for

305 SNP discovery (Figure 4 A). Imputation using an equal number of reference samples per pop-
306 ulation (scenario allPop_74_740) massively decreased this bias (Figure 4 B). The strength of
307 the correction was increased by an increase in the number of reference populations.

308

309 Figure 4: True H_E vs. ascertained H_E (A) and imputed H_E (B) by population group. For the im-
310 puted case, the strategy of using the same number of reference samples per population (all-
311 Pop_74_740) is shown, an increase in the number of reference samples per population (1-10)
312 is marked by an increasing color gradient and the line of identity is marked by a solid black
313 line.

314

315 To get an impression on the strength of the correction and the needed size of the reference
316 panel, Figure 5 compares the correlation by population group, the slope for the within-group
317 regression of the true H_E and H_O vs. the ascertained/ imputed cases and mean overestimation
318 for strategy allPop_74_740. It shows that the effects of ascertainment bias were stronger for
319 H_E than for H_O . Imputation when using the reference set with just one individual per popula-
320 tion corrects the initially much lower correlation within population group to > 0.99 . While
321 slope and mean overestimation are also pushed nearly ultimately towards the intended values
322 of one respectively zero for the non-discovery populations, there remains a small bias for the
323 discovery populations, which decreases with an increasing number of reference samples.

324

325 Figure 5: Development of correlation within population group (A), slope (B) and mean overes-
326 timation (C) of the regression lines for the two heterozygosity estimates when distributing the
327 reference samples equally across all populations (allPop_74_740). The intended value for un-
328 biasedness and minimum variance is marked as dense black horizontal line. Note that the case
329 without imputation is consistent with zero reference samples.

330

331 The effects were observed in a comparable manner for the other imputation strategies (Figure
332 S 3). Due to smaller reference panels, the correction effect of the imputation was generally
333 worse than for strategy allPop_74_740. Interestingly, when limiting the reference samples to
334 a small number of populations (strategies randPop_5_50, minPop_5_50, maxPop_5_50), we
335 observed a newly introduced bias towards the reference populations (Figure S 3). This effect
336 was strongest for strategy maxPop_5_50, where we chose the populations due to a maximum
337 distance from the discovery population. However, increasing the number of reference sam-
338 ples minimized the bias of reference and discovery populations with all strategies.

339 **Population distances**

340 The effects of ascertainment bias were also present but less pronounced for the distance
341 measurements (D and F_{ST} ; Figure S 4). The bias was thereby mainly present, when estimating
342 the distances between populations which belong to differently strong biased population
343 groups (Figure S 5). Note that F_{ST} was, all in all, less affected than D .

344 **Genotype to sequence**

345 For the imputation from array to sequence, the smallest reference panel (74_1perLine)
346 showed a median DR2 of 0.93 with 5 % and 25 % quantiles being 0.73 and 0.88. Increasing the
347 number of reference samples from the commercial populations to 158_all increased the 5%
348 quantile of DR2 by 0.05 while increasing all higher quantiles only by 0.01 (Table S 1). Note that
349 chromosomes 16, 22 and 25 clearly showed a higher proportion of badly imputed SNPs than
350 the other chromosomes (Figure S 6; Figure S 7).

351 The effect of imputation to WGS on ascertainment bias of HE is shown in Figure 6. Given the
352 situation that we cannot completely exclude pooling bias for the pooled sequenced samples
353 (Supplementary File 2), only the effect on the individually sequenced samples can be discussed

354 with adequate reliability. While the array-based estimates showed a slope of 1.94, the linkage
355 pruning slightly reduced this slope to 1.71. The clearly best result was achieved with imputa-
356 tion to WGS (slope = 1.26; 74_1perLine; Figure 6 A). However, the effect was also observed
357 for all samples and they additionally showed a slight increase of correlation (0.97 to 0.98),
358 which however is also influenced by pooling bias. Slightly increasing the reference panel (Fig-
359 ure 6 B) up to five samples per commercial line (98_5perLine) does not show any effect, while
360 using all commercial samples in the reference panel (158_all) and thereby clearly biasing the
361 reference panel towards the broiler samples increases HE again for all samples (slope = 1.44).

362

363 Figure 6: Effect of different correction strategies on ascertainment bias for expected hetero-
364 zygosity (HE). A – uncorrected array, linkage pruned array and imputed array (reference set
365 74_1perLine) based vs. sequence-based HE. B – array imputed with different reference sets
366 vs. sequence-based HE. The solid black line represents the line of identity, the solid colored
367 lines are regression lines within the individually sequenced populations (larger points) and the
368 dashed lines regression lines within all populations which include individually and pooled
369 (small points) sequenced populations.

370

371 The results for Nei's standard genetic distance (D; Figure 7) showed the same pattern as the
372 results for HE. The slope for distances between individually sequenced populations decreased
373 from 2.86 (array) and 1.77 (array_pruned) to 1.38 (imputed, 74_1perLine). The unbalanced
374 reference panel 158_all then again increased the slope to 1.56. The correlation for all dis-
375 tances, besides being also influenced by pooling bias and therefore being a rough estimate,
376 was increased from 0.93 (array) respectively 0.95 (array_pruned) to 0.98 (all reference sets).

377

378 Figure 7: Effect of different correction strategies on ascertainment bias for Nei's standard ge-
379 netic distance (D). A – uncorrected array, linkage pruned array and imputed array (reference
380 set 74_1perLine) based vs. sequence-based D. B – array imputed with different reference sets

381 vs. sequence-based HE. The solid black line represents the line of identity, the solid colored
382 lines are regression lines within distances between individually sequenced populations (larger
383 points) and the dashed lines regression lines within distances between all populations which
384 include individually and pooled (small points) sequenced populations.

385

386 **Discussion**

387 **Overall performance of the correction method**

388 Imputation of SNP data sets from lower to higher density is a commonly used technique to
389 either increase the resolution of data sets [33, 34, 38] or make them comparable across dif-
390 ferent platforms [59, 60]. The according studies mostly use a relatively homogeneous study
391 set and a closely related and large reference set [33, 34]. However, studies exist which inves-
392 tigate the effect of increasing the reference set to a multi-population reference set to make
393 use of the increased number of reference haplotypes [39–43]. To our knowledge, we here
394 present the first study which investigates the use of a relatively small and diverse reference
395 set on a large and diverse study set to correct for a genotyping platform-specific bias, the SNP
396 ascertainment bias.

397 This approach intends that single imputation errors do not harm, if the mean across the ge-
398 nome, presented by different population genetic estimators, shows unbiased results with min-
399 imum variance. Therefore, imputation to WGS level using a comparably small reference panel
400 can be used to correct for the ascertainment bias of commercial arrays.

401 Especially the *in silico* ascertained SNP arrays showed that even a very small reference panel
402 consisting of one individual of each population showed very good results for all investigated
403 estimators (Figure 5, Figure S 4) and became better with an increasing number of reference
404 populations. The results were less beneficial for the real WGS data, but also showed a strong

405 decrease of the slope towards one. From the imputed *in silico* arrays, we could additionally
406 realize a fast closing of the gap of the stronger overestimation of heterozygosity within dis-
407 covery populations and the less severe overestimation in non-discovery populations. This also
408 seemed to be the case when imputing to WGS level where we observed that the slope within
409 the commercial populations (closely related to discovery populations of the real array) de-
410 creased more than the slope within all populations due to imputation. However, this observa-
411 tion in the WGS data has to be regarded with caution, as we additionally identified a non-
412 negligible bias due to pooled sequencing which interfered with the assessment of ascertain-
413 ment bias and which was, in our study, confounded with the difference between commercial
414 populations (sequenced individually) and non-commercial populations (sequenced as pools).

415 The use of WGS information via imputation also consistently showed better results in regard
416 of reduction of ascertainment bias than using linkage pruned array SNPs which was reported
417 to be an effective filtering strategy for ascertainment bias mitigation by Malomane *et al.* [16].

418 Generally, the effect of imputation on the investigated estimators was shown to be compara-
419 ble across estimators, regardless of their initial reaction to ascertainment bias. An interesting
420 side observation was that F_{ST} did not show any ascertainment bias on the real array data (Fig-
421 ure S 10) when calculated in the form of summing the numerator across SNPs and dividing by
422 the sum of the denominator as calculated in this study. F_{ST} was only affected when used to
423 estimate differentiation between the, in regard of heterozygosity differently strong affected,
424 discovery- and non-discovery populations in the simulated array data. This strongly supports
425 the findings of Albrechtsen *et al.* [11], who showed F_{ST} to be relatively robust against the ef-
426 fects of ascertainment bias.

427 We also investigated the effect of differently sized and constructed reference sets for imputation. Generally, larger reference sets increased the accuracy of imputation and therefore
428 decreased the effects of ascertainment bias. The best results were achieved when the reference set was as evenly distributed across the study set as possible. When reference populations were closely related to the discovery population, reduction in imputation quality and
431 increase in ascertainment bias were less severe in case of unbalanced reference sets than if
432 distantly related reference populations were used. This suggests that variation within study-
433 and reference set needs to show enough overlap to achieve sufficient imputation accuracy
434 and therefore reduction of ascertainment bias.

436 Results from literature suggest that multi-breed reference panels generally increase imputation accuracy especially for rare variants and within admixed populations [39–41]. Additionally, Rowan *et al.* [41] argue that they do not seem to introduce variation at a relevant scale
438 for markers for which the breeds are actually fixed. However, some studies also showed that
439 strongly unbalanced reference sets can reduce imputation accuracy [42, 43]. In this study, including additional reference samples in a biased way when going from reference set 74_1per-
441 Line to 158_all increased the effects of ascertainment bias. However, theoretical imputation
442 accuracies rather increased than decreased (Figure S 6; Table S 1) for previously poorly im-
443 puted SNPs. On one hand, this effect supports the findings of Brøndum [39], Rowan *et al.* [41]
444 and Ye *et al.* [40] and on the other hand, nicely highlights the main reason for ascertainment
445 bias. One can only identify variation which is present in the investigated sample. When developing an array, it is the variation in the discovery set, while in our case it is the reference set
446 used for imputation. Therefore, it is crucial to use a reference set for imputation which covers
447 the intended range of variation.

450 Besides the previously described effects of imputation on ascertainment bias, we also identi-
451 fied an effect of array design on imputation accuracy. Discovery populations show higher im-
452 putation accuracies than non-discovery populations (Figure 3). As markers on arrays are more
453 representative for discovery populations than non-discovery populations, relatively more of
454 the genetic variability in discovery populations is explained by the array and imputation is
455 more accurate on average.

456 **Conclusion**

457 Imputation is generally able to mitigate ascertainment bias. The effect is already present when
458 using a very small reference set of only one sequenced individual per population. It also per-
459 forms better than simple filtering strategies based on the array data alone. However, care has
460 to be taken in designing an evenly spaced reference panel to not introduce a new bias towards
461 variation present in the reference panel. We also suggest using a larger reference panel than
462 the one which was available for this study to achieve better results. Additionally, we observed
463 an effect of array design on imputation accuracy as discovery populations showed a higher
464 imputation accuracy than non-discovery populations. This should be taken into account when
465 designing studies based on imputed SNPs by choosing an appropriate genotyping array for the
466 intended study populations.

467 **Abbreviations**

468 D – Nei's Standard Genetic Distance

469 F_{ST} – Wright's Fixation Index

470 H_E – expected heterozygosity

471 H_0 – observed heterozygosity
472 MAF – minor allele frequency
473 r – Pearson Correlation
474 SNP – single nucleotide polymorphism
475 WGS – whole genome sequencing

476

477 **Declarations**

478 **Ethics approval and consent to participate**

479 The study did not involve new treatment of animals as only published data was used. DNA
480 samples for all already published raw data were taken from a data base established during the
481 project AVIANDIV (EC Contract No. BIO4-CT98_0342; 1998 – 2000; www.aviandiv.fli.de) and
482 later extended by samples of the project SYNBREED (FKZ 0315528E; 2009 – 2014; www.synbreed.tum.de). Blood sampling was done in strict accordance to the German animal welfare
483 regulations, with written consent of the animal owners and was approved by the at the ac-
484 cording times ethics responsible persons of the Friedrich-Loeffler-Institut. According to Ger-
485 man animal welfare regulations, notice was given to the responsible institution, the Lower
486 Saxony State Office for Consumer Protection and Food Safety (33.9-42502-05-10A064).

488 **Consent for publication**

489 Not applicable

490 **Availability of data and materials**

491 Raw sequencing and genotyping data were previously published by different studies. The re-
492 pository information for each sample can be found in Supplementary File 1. All datasets gen-
493 erated by analyses during this study from the raw data are additionally available from the
494 corresponding author on reasonable request.

495 **Competing interests**

496 The authors declare that they have no competing interests

497 **Funding**

498 The positions of JG and TP were partly funded by the project IMAGE - Innovative Management
499 of Animal Genetic Resources (www.imageh2020.eu, funded by the EU Horizon 2020 research
500 and innovation program No 677353). The sampling, genotyping and sequencing of the chicken
501 populations involved funding by the EC project AVIANDIV (EC Contract No. BIO4-CT98_0342;
502 1998 – 2000; www.aviandiv.fli.de), the German Federal Ministry of Education and Research
503 (BMBF) via the SYNBREED project (FKZ 0315528E; 2009 – 2014; www.synbreed.tum.de) and
504 the EU project IMAGE - Innovative Management of Animal Genetic Resources (www.im-
505 ageh2020.eu, funded by the EU Horizon 2020 research and innovation program No 677353).
506 Further, the publication fees were covered by the Open Access Publication Funds of the Uni-
507 versity of Goettingen.

508 **Authors' contributions**

509 JG conceptualized and designed the study, analyzed and interpreted the data, wrote the initial
510 draft and revised the manuscript. CR and TP substantially contributed to design of the study,
511 interpretation of the data and revision of the manuscript. SW substantially contributed to ac-

512 quisation and interpretation of the data and revised the manuscript. AW substantially contrib-
513 uted to acquisition and curation of the data. HS substantially contributed to conception and
514 design of the study, interpretation of the data and revision of the manuscript. All authors have
515 red and agreed on the final draft.

516 **Acknowledgement**

517 Not applicable

518 **Supplementary**

519 Supplementary File 1: Accession Information of raw data per sample

520 Supplementary File 2: Supplementary Methods

521 Figure S 1: Recall rates for samples which were genotyped as well as sequenced per SNP (A;
522 B) and per animal (C; D); before (A; C; red) and after (B; D; blue) correction of potential refer-
523 ence allele switches in the genotype data.

524 Figure S 2: Development of the per-animal imputation accuracy with an increasing number of
525 reference animals per population. A – scenario randSamp_5_50; B – scenario randPop_5_50;
526 C – scenario minPop_5_50; D – scenario maxPop_5_50. Individuals are grouped on whether
527 they belong to the population which contains reference individuals, was used as for SNP dis-
528 covery or none of them (application). the lines show the trend of the median.

529 Figure S 3: Development of correlations within population group (r), slope and mean overes-
530 timation of the regression lines for H_E and H_O estimates and different reference panel strate-
531 gies. The intended value for unbiasedness and minimum variance is marked as dense black
532 horizontal line. Note that the case without imputation is consistent with zero reference sam-
533 ples.

534 Figure S 4: Development of correlation within population group (A), slope (B) and intercept (C)
535 of the regression lines for D and F_{ST} when distributing the reference samples equally over all
536 populations (allPop_74_740). The intended value for unbiasedness and minimum variance is
537 marked as dense black horizontal line. Note that the case without imputation is consistent
538 with zero reference samples.

539 Figure S 5: Development of correlation within population group (r), slope and mean overesti-
540 mation of the regression lines for Nei's Distance (D) and F_{ST} estimates and different reference
541 panel strategies. The intended value for unbiasedness and minimum variance is marked as
542 dense black horizontal line. Note that the case without imputation is consistent with zero ref-
543 erence samples.

544 Figure S 6: Distribution of DR2 values by chromosome and reference set. Note that outliers
545 are not shown due to a large number of underlying values.

546 Figure S 7: Two-dimensional distributions of DR2 values vs. MAF by chromosome when im-
547 puted with the reference set 74_1perLine. The red line represents the median within 0.05
548 MAF bins.

549 Figure S 8: Effect of pooled sequencing and the correction factor of Futschik and Schlötterer
550 [48] on expected heterozygosity (HE) and ascertainment bias. A – HE estimated from array
551 positions of the sequencing data vs. HE directly estimated from array data. The color indicates
552 the state before and after correcting the pooled sequence estimates and the accordingly col-
553 ored solid lines the group specific regression lines while the black solid line indicates the line
554 of identity in all three plots. The plot therefore shows the magnitude of the bias introduced
555 by pooled sequencing and the according effect of the correction factor. B – HE estimated from
556 the array data vs. HE estimated from the complete sequence data. The color again shows the
557 values before and after implementing the correction of the pooled sequence estimates. While
558 the solid regression line and dense circles indicate the individually sequenced samples, the
559 dashed regression lines and triangles indicate pooled sequenced samples. The plot therefore
560 shows the combined effect of ascertainment bias and pooled sequencing bias. C – HE esti-
561 mated from array positions of the sequencing data vs. HE estimated from all positions of the
562 sequencing data. The plot therefore shows the pure ascertainment bias.

563 Figure S 9: Effect of pooled sequencing on the expression of the ascertainment bias in Nei's
564 standard genetic distance (D). The biased D was either estimated directly from the array gen-
565 otypes (D_{arr} , pooled bias + ascertainment bias) or from the array positions of the sequenc-
566 ing data ($D_{arr.seq}$, pure ascertainment bias), while the estimates from the complete sequence
567 were assumed to be the true estimates. The black solid line represents the line of identity,
568 solid colored regression lines and dense points represent estimates between individually se-
569 quenced populations and dashed lines and triangles represent estimates between two popu-
570 lations of which at least one was pooled sequenced.

571 Figure S 10: Effect of pooled sequencing on the expression of the ascertainment bias in
572 Wright's fixation index (F_{ST}). The biased F_{ST} was either estimated directly from the array gen-
573 otypes ($F_{ST.arr}$, pooled bias + ascertainment bias) or from the array positions of the sequenc-
574 ing data ($F_{ST.arr.seq}$, pure ascertainment bias), while the estimates from the complete se-
575 quence were assumed to be the true estimates. The black solid line represents the line of
576 identity, solid colored regression lines and dense points represent estimates between individ-
577 ually sequenced populations and dashed lines and triangles represent estimates between two
578 populations of which at least one was pooled sequenced.

579 Table S 1: Quantiles of theoretical imputation accuracies (DR2) by reference set

580

581 **References**

582 1. Novembre J, Johnson T, Bryc K, Kutalik Z, Boyko AR, Auton A, et al. Genes mirror geogra-
583 phy within Europe. *Nature*. 2008;456:98. doi:10.1038/nature07331.

584 2. Patterson N, Moorjani P, Luo Y, Mallick S, Rohland N, Zhan Y, et al. Ancient admixture in
585 human history. *Genetics*. 2012;192:1065–93. doi:10.1534/genetics.112.145037.

586 3. Laurie CC, Nickerson DA, Anderson AD, Weir BS, Livingston RJ, Dean MD, et al. Linkage
587 Disequilibrium in Wild Mice. *PLoS Genet*. 2007;3:e144. doi:10.1371/jour-

588 nal.pgen.0030144.

589 4. Platt A, Horton M, Huang YS, Li Y, Anastasio AE, Mulyati NW, et al. The Scale of Popula-
590 tion Structure in *Arabidopsis thaliana*. *PLoS Genet*. 2010;6:e1000843. doi:10.1371/jour-

591 nal.pgen.1000843.

592 5. Travis AJ, Norton GJ, Datta S, Sarma R, Dasgupta T, Savio FL, et al. Assessing the genetic
593 diversity of rice originating from Bangladesh, Assam and West Bengal. *Rice*. 2015;8:35.

594 doi:10.1186/s12284-015-0068-z.

595 6. Mayer M, Unterseer S, Bauer E, Leon N de, Ordas B, Schön C-C. Is there an optimum
596 level of diversity in utilization of genetic resources? *Theor Appl Genet*. 2017;130:2283–

597 95. doi:10.1007/s00122-017-2959-4.

598 7. Muir WM, Wong GK-S, Zhang Y, Wang J, Groenen MAM, Crooijmans RPMA, et al. Ge-

599 nome-wide assessment of worldwide chicken SNP genetic diversity indicates significant

600 absence of rare alleles in commercial breeds. *Proc. Natl. Acad. Sci*. 2008;105:17312–7.

601 doi:10.1073/pnas.0806569105.

- 602 8. Gibbs RA, Taylor JF, van Tassell CP, Barendse W, Eversole KA, Gill CA, et al. Genome-wide
603 survey of SNP variation uncovers the genetic structure of cattle breeds. *Science*.
604 2009;324:528–32. doi:10.1126/science.1167936.
- 605 9. Geibel J, Reimer C, Weigend S, Weigend A, Pook T, Simianer H. How Array Design Affects
606 SNP Ascertainment Bias. *bioRxiv*. 2019:833541. doi:10.1101/833541.
- 607 10. Clark AG, Hubisz MJ, Bustamante CD, Williamson SH, Nielsen R. Ascertainment bias in
608 studies of human genome-wide polymorphism. *Genome Res*. 2005;15:1496–502.
- 609 11. Albrechtsen A, Nielsen FC, Nielsen R. Ascertainment biases in SNP chips affect measures
610 of population divergence. *Mol. Biol. Evol.* 2010;27:2534–2547.
- 611 12. Lachance J, Tishkoff SA. SNP ascertainment bias in population genetic analyses: why it is
612 important, and how to correct it. *Bioessays*. 2013;35:780–6.
- 613 13. Nielsen R. Population genetic analysis of ascertained SNP data. *Hum. Genomics*.
614 2004;1:1.
- 615 14. Nielsen R, Hubisz MJ, Clark AG. Reconstituting the frequency spectrum of ascertained
616 single-nucleotide polymorphism data. *Genetics*. 2004;168:2373–82.
- 617 15. Quinto-Cortés CD, Woerner AE, Watkins JC, Hammer MF. Modeling SNP array ascertain-
618 ment with Approximate Bayesian Computation for demographic inference. *Sci. Rep*.
619 2018;8:10209. doi:10.1038/s41598-018-28539-y.
- 620 16. Malomane DK, Reimer C, Weigend S, Weigend A, Sharifi AR, Simianer H. Efficiency of dif-
621 ferent strategies to mitigate ascertainment bias when using SNP panels in diversity stud-
622 ies. *BMC Genomics*. 2018;19:22. doi:10.1186/s12864-017-4416-9.
- 623 17. Qanbari S, Pausch H, Jansen S, Somel M, Strom T-M, Fries R, et al. Classic selective
624 sweeps revealed by massive sequencing in cattle. *PLoS Genet*. 2014;10:e1004148.
625 doi:10.1371/journal.pgen.1004148.

- 626 18. Qanbari S, Seidel M, Strom T-M, Mayer KFX, Preisinger R, Simianer H. Parallel Selection
627 Revealed by Population Sequencing in Chicken. *Genome Biol. Evol.* 2015;7:3299–306.
628 doi:10.1093/gbe/evv222.
- 629 19. Lawal RA, Al-Atiyat RM, Aljumaah RS, Silva P, Mwacharo JM, Hanotte O. Whole-Genome
630 Resequencing of Red Junglefowl and Indigenous Village Chicken Reveal New Insights on
631 the Genome Dynamics of the Species. *Front Genet.* 2018;9:264.
632 doi:10.3389/fgene.2018.00264.
- 633 20. Qanbari S, Rubin C-J, Maqbool K, Weigend S, Weigend A, Geibel J, et al. Genetics of ad-
634 aptation in modern chicken. *PLoS Genet.* 2019;15:e1007989. doi:10.1371/jour-
635 nal.pgen.1007989.
- 636 21. Peripolli E, Reimer C, Ha N-T, Geibel J, Machado MA, Panetto, João Cláudio do Carmo, et
637 al. Genome-wide detection of signatures of selection in indicine and Brazilian locally
638 adapted taurine cattle breeds using whole-genome re-sequencing data. *BMC Genomics.*
639 2020;21:624. doi:10.1186/s12864-020-07035-6.
- 640 22. Auton A, Abecasis GR, Altshuler DM, Durbin RM, Bentley DR, Chakravarti A, et al. A
641 global reference for human genetic variation. *Nature.* 2015;526:68–74. doi:10.1038/na-
642 ture15393.
- 643 23. Hayes BJ, Daetwyler HD. 1000 Bull Genomes Project to Map Simple and Complex Genetic
644 Traits in Cattle: Applications and Outcomes. *Annual Review of Animal Biosciences.*
645 2019;7:89–102. doi:10.1146/annurev-animal-020518-115024.
- 646 24. Alonso-Blanco C, Andrade J, Becker C, Bemm F, Bergelson J, Borgwardt KM, et al. 1,135
647 Genomes Reveal the Global Pattern of Polymorphism in *Arabidopsis thaliana*. *Cell.*
648 2016;166:481–91. doi:10.1016/j.cell.2016.05.063.

- 649 25. Marchini J, Howie B. Genotype imputation for genome-wide association studies. *Nat.*
650 *Rev. Genet.* 2010;11:499 EP -. doi:10.1038/nrg2796.
- 651 26. Li N, Stephens M. Modeling linkage disequilibrium and identifying recombination
652 hotspots using single-nucleotide polymorphism data. *Genetics.* 2003;165:2213–33.
- 653 27. Marchini J, Howie B, Myers S, McVean G, Donnelly P. A new multipoint method for ge-
654 nome-wide association studies by imputation of genotypes. *Nat. Genet.* 2007;39:906–
655 13.
- 656 28. Howie BN, Donnelly P, Marchini J. A Flexible and Accurate Genotype Imputation Method
657 for the Next Generation of Genome-Wide Association Studies. *PLoS Genet.* 2009;5:1–15.
658 doi:10.1371/journal.pgen.1000529.
- 659 29. Delaneau O, Marchini J, Zagury J-F. A linear complexity phasing method for thousands of
660 genomes. *Nature Methods.* 2012;9:179–81.
- 661 30. Sargolzaei M, Chesnais JP, Schenkel FS. A new approach for efficient genotype imputa-
662 tion using information from relatives. *BMC Genomics.* 2014;15:478. doi:10.1186/1471-
663 2164-15-478.
- 664 31. Money D, Gardner K, Migicovsky Z, Schwaninger H, Zhong G-Y, Myles S. LinkImpute: Fast
665 and Accurate Genotype Imputation for Nonmodel Organisms. *G3.* 2015;5:2383.
666 doi:10.1534/g3.115.021667.
- 667 32. Browning BL, Zhou Y, Browning SR. A One-Penny Imputed Genome from Next-Genera-
668 tion Reference Panels. *Am. J. Hum. Genet.* 2018;103:338–48.
669 doi:10.1016/j.ajhg.2018.07.015.
- 670 33. Pausch H, Aigner B, Emmerling R, Edel C, Götz K-U, Fries R. Imputation of high-density
671 genotypes in the Fleckvieh cattle population. *Genet Sel Evol.* 2013;45:3.
672 doi:10.1186/1297-9686-45-3.

- 673 34. Heidaritabar M, Calus MPL, Megens H-J, Vereijken A, Groenen MAM, Bastiaansen JWM.
674 Accuracy of genomic prediction using imputed whole-genome sequence data in white
675 layers. *J Anim Breed Genet*. 2016;133:167–79. doi:10.1111/jbg.12199.
- 676 35. van den Berg S, Vandenplas J, van Eeuwijk FA, Bouwman AC, Lopes MS, Veerkamp RF.
677 Imputation to whole-genome sequence using multiple pig populations and its use in ge-
678 nome-wide association studies. *Genetics Selection Evolution*. 2019;51:2.
679 doi:10.1186/s12711-019-0445-y.
- 680 36. Huang J, Ellinghaus D, Franke A, Howie B, Li Y. 1000 Genomes-based imputation identi-
681 fies novel and refined associations for the Wellcome Trust Case Control Consortium
682 phase 1 Data. *European Journal Of Human Genetics*. 2012;20:801–5.
683 doi:10.1038/ejhg.2012.3.
- 684 37. Artigas MS, Wain LV, Miller S, Kheirallah AK, Huffman JE, Ntalla I, et al. Sixteen new lung
685 function signals identified through 1000 Genomes Project reference panel imputation.
686 *Nat. Commun*. 2015;6:8658. doi:10.1038/ncomms9658.
- 687 38. Raymond B, Bouwman AC, Schrooten C, Houwing-Duistermaat J, Veerkamp RF. Utility of
688 whole-genome sequence data for across-breed genomic prediction. *Genetics Selection
689 Evolution*. 2018;50:27. doi:10.1186/s12711-018-0396-8.
- 690 39. Brøndum RF, Guldbbrandtsen B, Sahana G, Lund MS, Su G. Strategies for imputation to
691 whole genome sequence using a single or multi-breed reference population in cattle.
692 *BMC Genomics*. 2014;15:728. doi:10.1186/1471-2164-15-728.
- 693 40. Ye S, Yuan X, Huang S, Zhang H, Chen Z, Li J, et al. Comparison of genotype imputation
694 strategies using a combined reference panel for chicken population. *Animal*.
695 2019;13:1119–26. doi:10.1017/S1751731118002860.

- 696 41. Rowan TN, Hoff JL, Crum TE, Taylor JF, Schnabel RD, Decker JE. A multi-breed reference
697 panel and additional rare variants maximize imputation accuracy in cattle. *Genetics Se-*
698 *lection Evolution*. 2019;51:77. doi:10.1186/s12711-019-0519-x.
- 699 42. Berry DP, McClure MC, Mullen MP. Within- and across-breed imputation of high-density
700 genotypes in dairy and beef cattle from medium- and low-density genotypes. *J Anim*
701 *Breed Genet*. 2014;131:165–72. doi:10.1111/jbg.12067.
- 702 43. Korkuć P, Arends D, Brockmann GA. Finding the Optimal Imputation Strategy for Small
703 Cattle Populations. *Front Genet*. 2019;10:52. doi:10.3389/fgene.2019.00052.
- 704 44. Pook T, Mayer M, Geibel J, Weigend S, Cavero D, Schoen CC, Simianer H. Improving Im-
705 putation Quality in BEAGLE for Crop and Livestock Data. *G3*. 2019:g3.400798.2019.
706 doi:10.1534/g3.119.400798.
- 707 45. Bortoluzzi C, Megens H-J, Bosse M, Derks MFL, Dibbitts B, Laport K, et al. Parallel Genetic
708 Origin of Foot Feathering in Birds. *Mol. Biol. Evol*. 2020;37:2465–76. doi:10.1093/mol-
709 bev/msaa092.
- 710 46. Malomane DK, Simianer H, Weigend A, Reimer C, Schmitt AO, Weigend S. The SYNBREED
711 chicken diversity panel: A global resource to assess chicken diversity at high genomic res-
712 olution. *BMC Genomics*. 2019;20:345. doi:10.1186/s12864-019-5727-9.
- 713 47. Kranis A, Gheyas AA, Boschiero C, Turner F, Le Yu, Smith S, et al. Development of a high
714 density 600K SNP genotyping array for chicken. *BMC Genomics*. 2013;14:59.
715 doi:10.1186/1471-2164-14-59.
- 716 48. Futschik A, Schlötterer C. The Next Generation of Molecular Markers From Massively
717 Parallel Sequencing of Pooled DNA Samples. *Genetics*. 2010;186:207–18.
718 doi:10.1534/genetics.110.114397.

- 719 49. Schlötterer C, Tobler R, Kofler R, Nolte V. Sequencing pools of individuals [mdash] mining
720 genome-wide polymorphism data without big funding. *Nat. Rev. Genet.* 2014;15:749–
721 63.
- 722 50. Genome Reference Consortium GRCg6a. GRCg6a chicken reference genome. 2018.
723 <http://hgdownload.soe.ucsc.edu/goldenPath/galGal6/bigZips/galGal6.fa.gz>. Accessed 9
724 Apr 2019.
- 725 51. DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, et al. A framework for
726 variation discovery and genotyping using next-generation DNA sequencing data. *Nat.*
727 *Genet.* 2011;43:491. doi:10.1038/ng.806.
- 728 52. van der Auwera GA, Carneiro MO, Hartl C, Poplin R, Del Angel G, Levy-Moonshine A, et
729 al. From FastQ data to high confidence variant calls: The Genome Analysis Toolkit best
730 practices pipeline. *Curr Protoc Bioinformatics.* 2013;43:11.10.1-11.10.33.
731 doi:10.1002/0471250953.bi1110s43.
- 732 53. Groenen MAM, Wahlberg P, Foglio M, Cheng HH, Megens H-J, Crooijmans RPMA, et al. A
733 high-density SNP-based linkage map of the chicken genome reveals sequence features
734 correlated with recombination rate. *Genome Res.* 2009;19:510–9.
- 735 54. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK:
736 Rising to the challenge of larger and richer datasets. *Gigascience.* 2015;4:7.
737 doi:10.1186/s13742-015-0047-8.
- 738 55. Nei M. Genetic Distance between Populations. *The American Naturalist.* 1972;106:283–
739 92.
- 740 56. Hickey JM, Crossa J, Babu R, los Campos G de. Factors Affecting the Accuracy of Geno-
741 type Imputation in Populations from Several Maize Breeding Programs. *Crop Science.*
742 2012;52:654. doi:10.2135/cropsci2011.07.0358.

- 743 57. Browning BL, Browning SR. A unified approach to genotype imputation and haplotype-
744 phase inference for large data sets of trios and unrelated individuals. *Am. J. Hum. Genet.*
745 2009;84:210–23. doi:10.1016/j.ajhg.2009.01.005.
- 746 58. Wright S. The genetical structure of populations. *Ann. Eugen.* 1949;15:323–54.
747 doi:10.1111/j.1469-1809.1949.tb02451.x.
- 748 59. Al-Tassan NA, Whiffin N, Hosking FJ, Palles C, Farrington SM, Dobbins SE, et al. A new
749 GWAS and meta-analysis with 1000Genomes imputation identifies novel risk variants for
750 colorectal cancer. *Sci. Rep.* 2015;5:10442. doi:10.1038/srep10442.
- 751 60. Bouwman AC, Daetwyler HD, Chamberlain AJ, Ponce CH, Sargolzaei M, Schenkel FS, et al.
752 Meta-analysis of genome-wide association studies for cattle stature identifies common
753 genes that regulate body size in mammals. *Nat. Genet.* 2018;50:362–7.
754 doi:10.1038/s41588-018-0056-5.
- 755
- 756

Figures

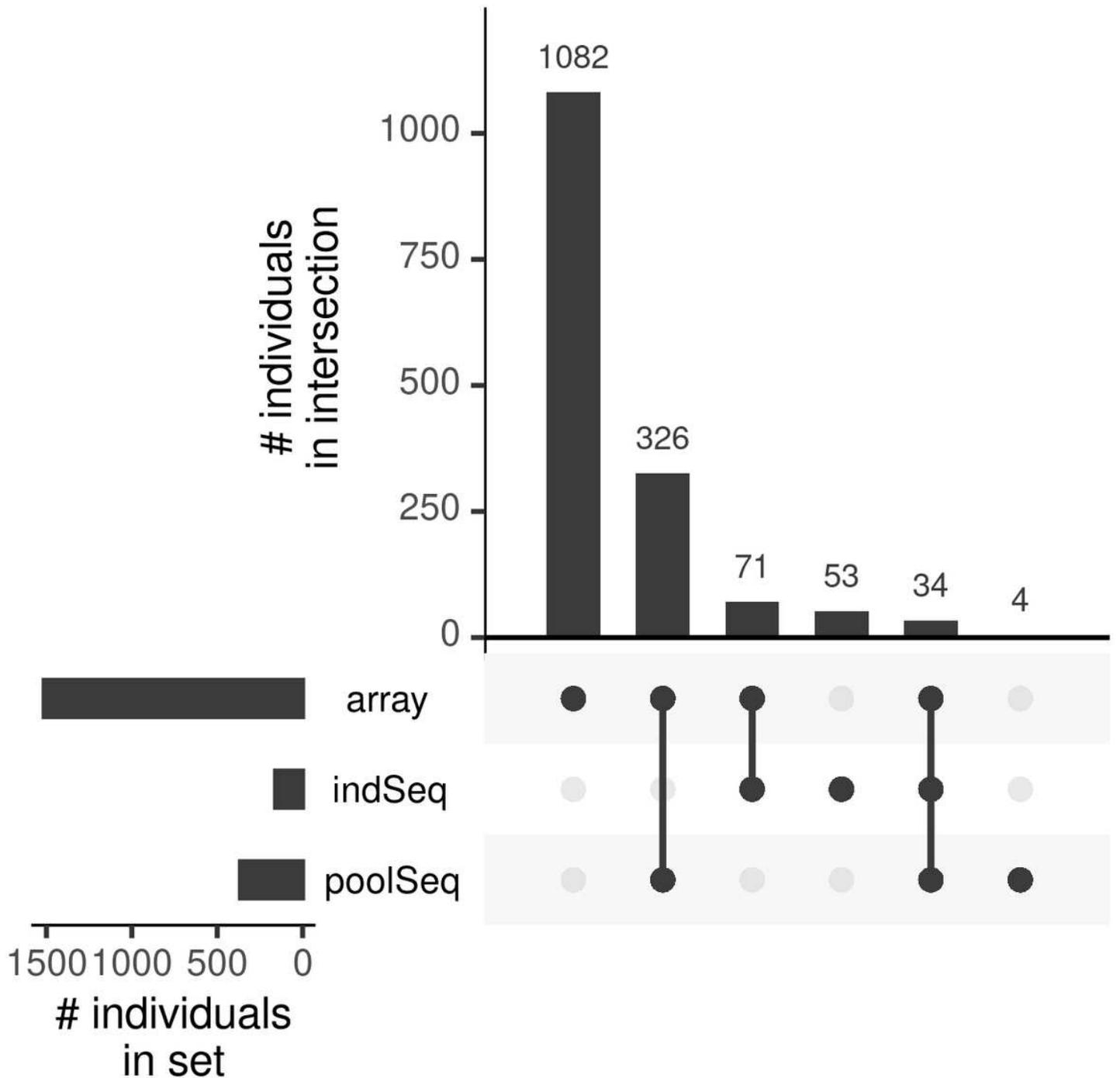


Figure 1

UpSet plot showing the distinct intersections of chickens between the used sequencing/ gen-otyping technologies. The left bar plot contains the total number of individuals that were genotyped (array), individually sequenced (indSeq), or pooled sequenced (poolSeq). The up-per bar plot contains the number of individuals within each distinct intersection, indicated by the connected points below.

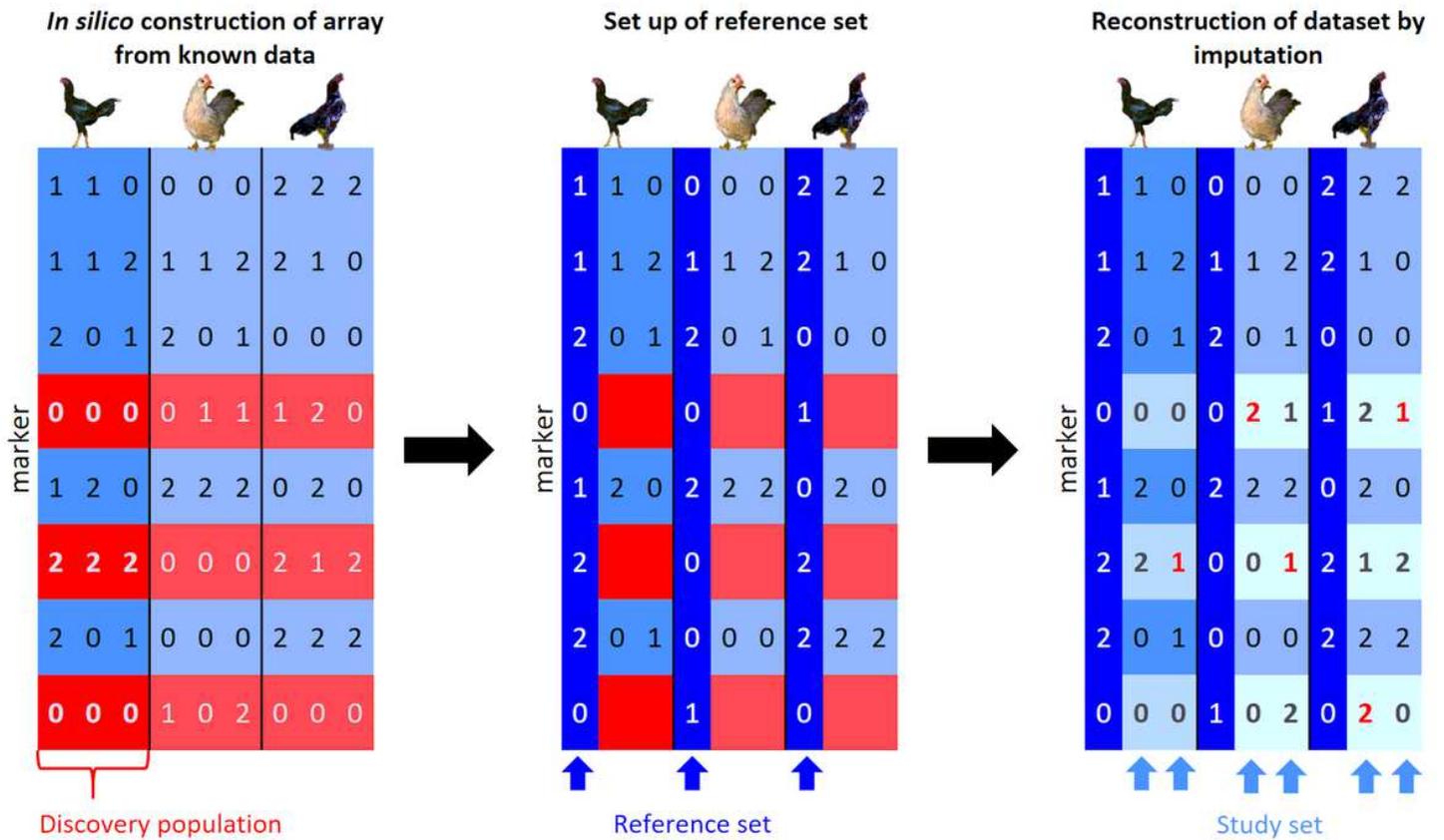


Figure 2

Schematic representation of the workflow of creating and re-imputing the in silico arrays. The starting point was a 0/1/2 coded marker matrix with SNPs in rows and individuals in columns (different populations separated by vertical lines). In a first step, an array (light blue rows) was constructed in silico from known data by setting all SNPs to missing which were invariable (MAF < 0.05, red rows) in the discovery population (first three columns). In a second step, a reference set (dark blue columns) was set up from animals for which complete knowledge of all SNPs was assumed. This Reference set was then used in a third step to impute the missing SNPs in the study set using Beagle 5.0 and resulting in a certain amount of imputation errors (red numbers).

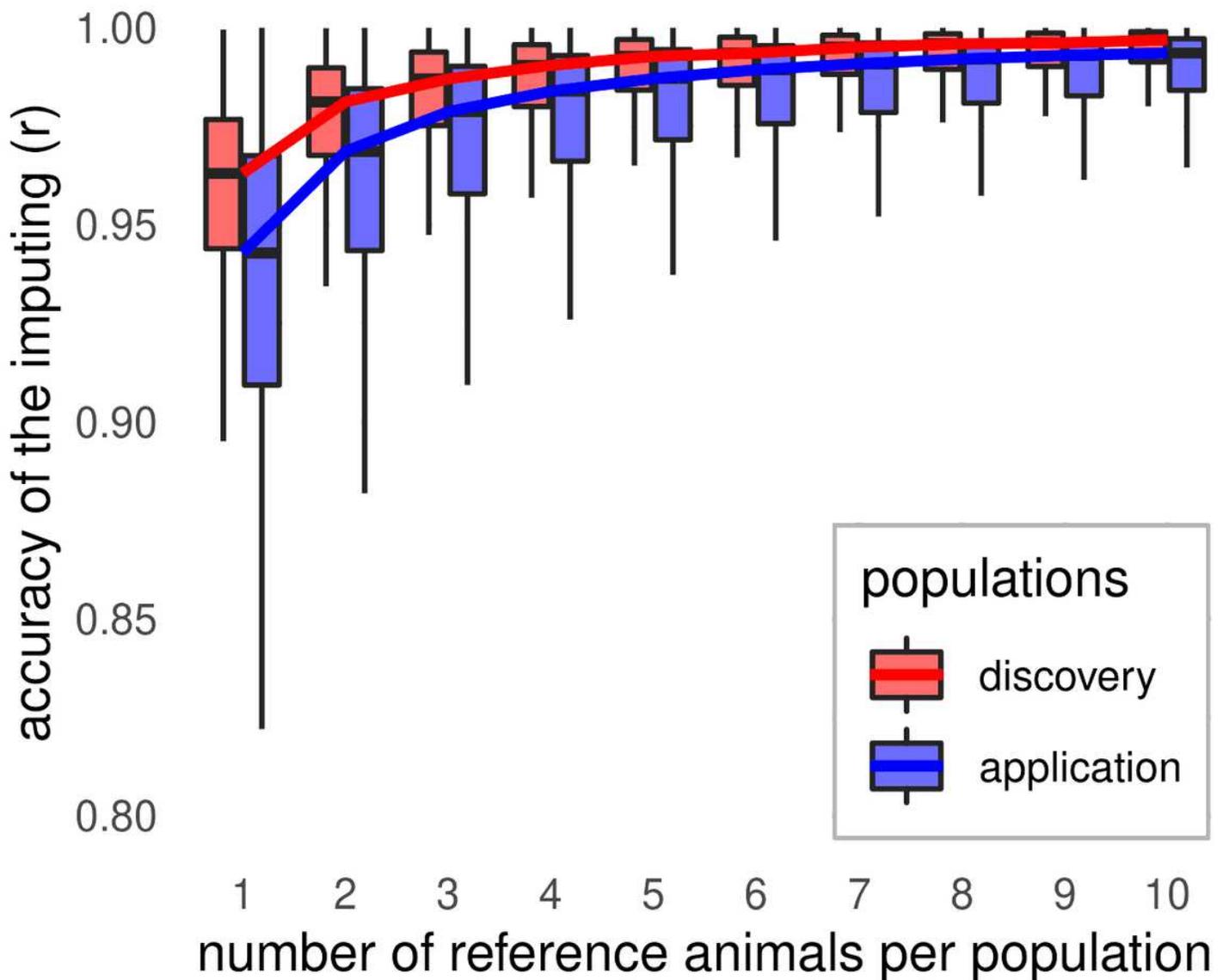


Figure 3

Development of the per-animal imputation accuracy for the in silico array to genotype set imputation with an increasing number of reference animals per population. Individuals are grouped on whether they belong to the population used for SNP discovery or not and reference individuals were chosen as in scenario allPop_74_740. The lines show the trend of the median and outliers are not shown in the plot as they do not add valuable information due to the high number of repetitions.

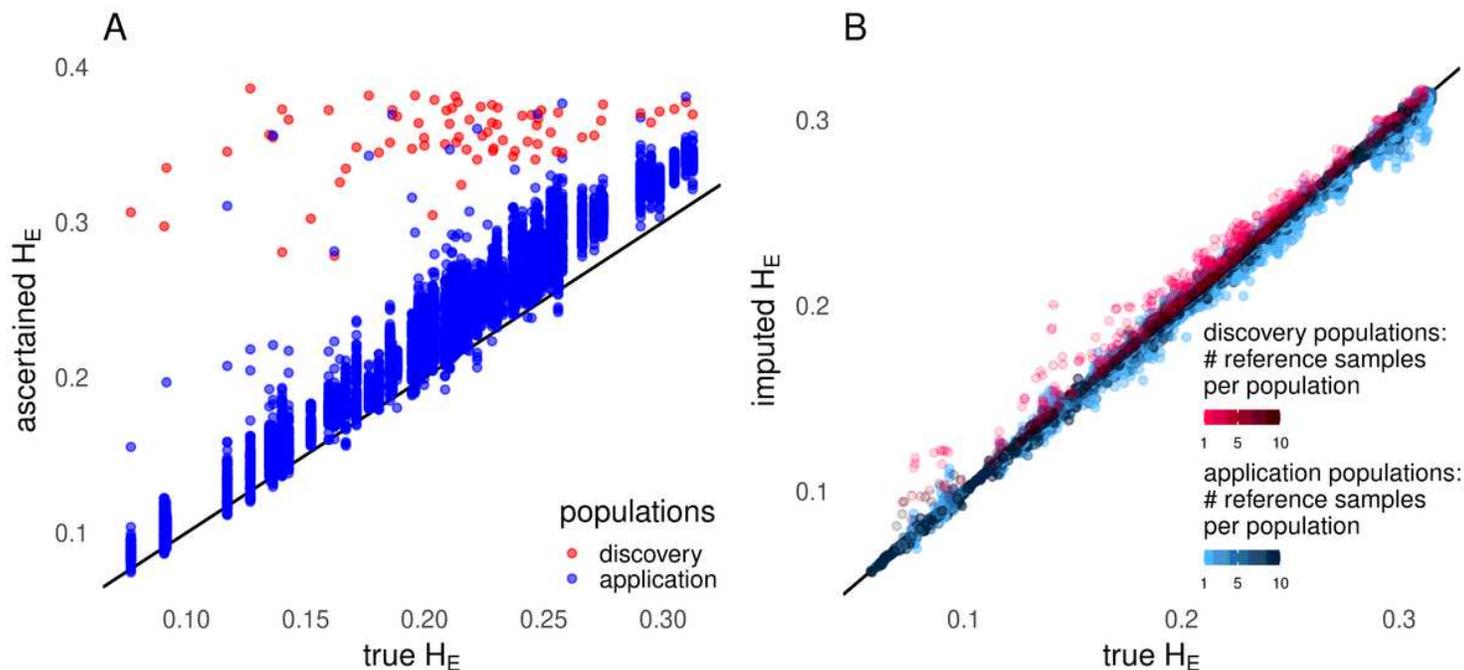


Figure 4

True HE vs. ascertained HE (A) and imputed HE (B) by population group. For the imputed case, the strategy of using the same number of reference samples per population (allPop_74_740) is shown, an increase in the number of reference samples per population (1-10) is marked by an increasing color gradient and the line of identity is marked by a solid black line.

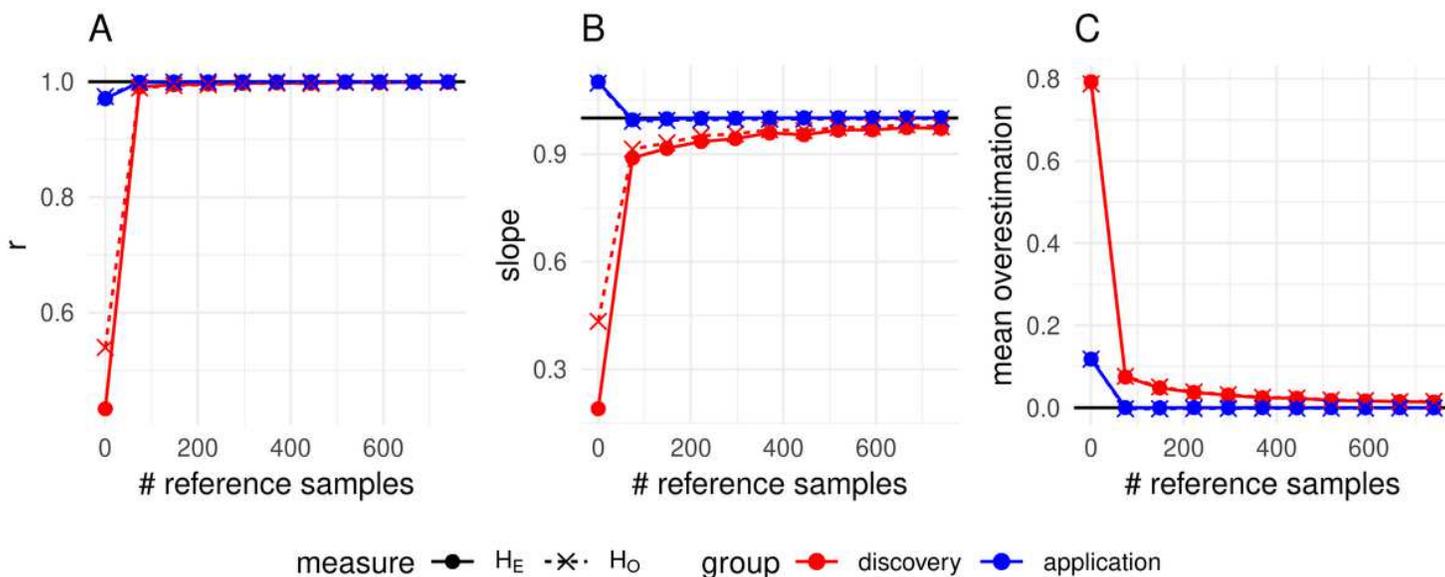


Figure 5

Development of correlation within population group (A), slope (B) and mean overestimation (C) of the regression lines for the two heterozygosity estimates when distributing the reference samples equally across all populations (allPop_74_740). The intended value for unbiasedness and minimum variance is

marked as dense black horizontal line. Note that the case without imputation is consistent with zero reference samples.

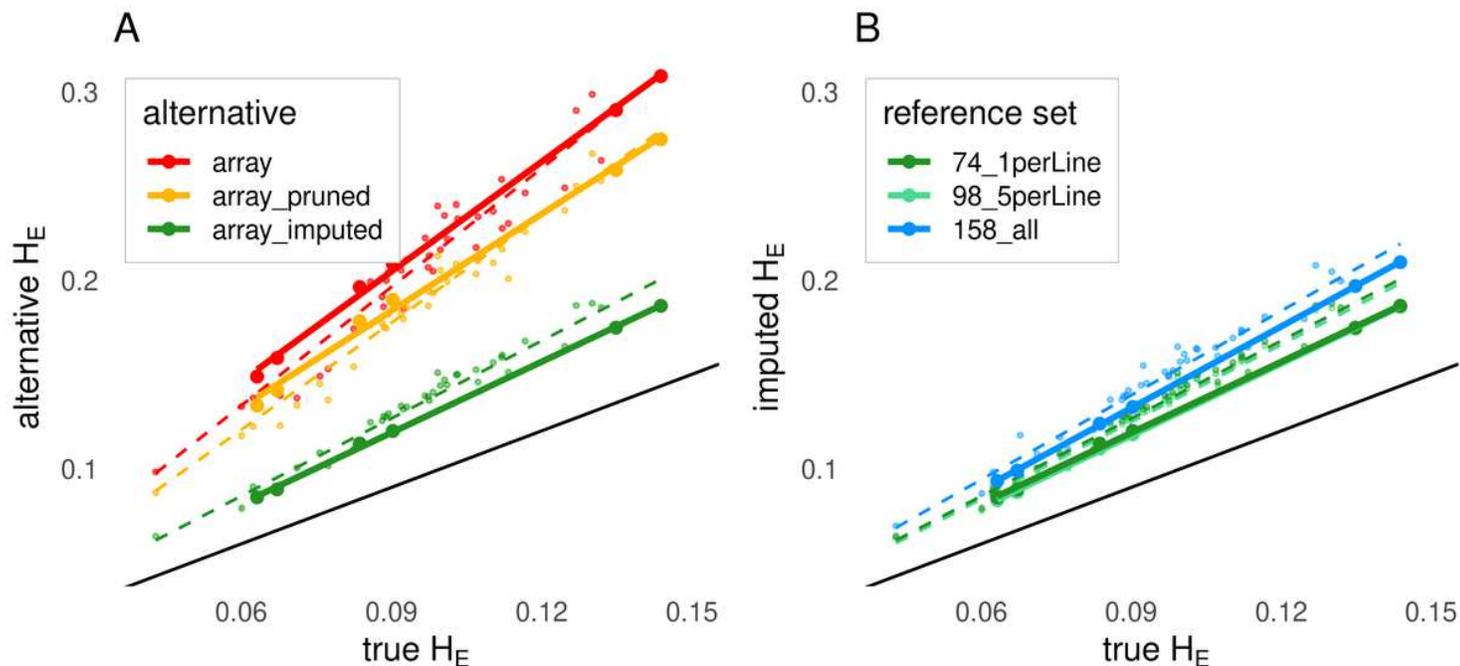


Figure 6

Effect of different correction strategies on ascertainment bias for expected heterozygosity (HE). A – uncorrected array, linkage pruned array and imputed array (reference set 74_1perLine) based vs. sequence-based HE. B – array imputed with different reference sets vs. sequence-based HE. The solid black line represents the line of identity, the solid colored lines are regression lines within the individually sequenced populations (larger points) and the dashed lines regression lines within all populations which include individually and pooled (small points) sequenced populations.

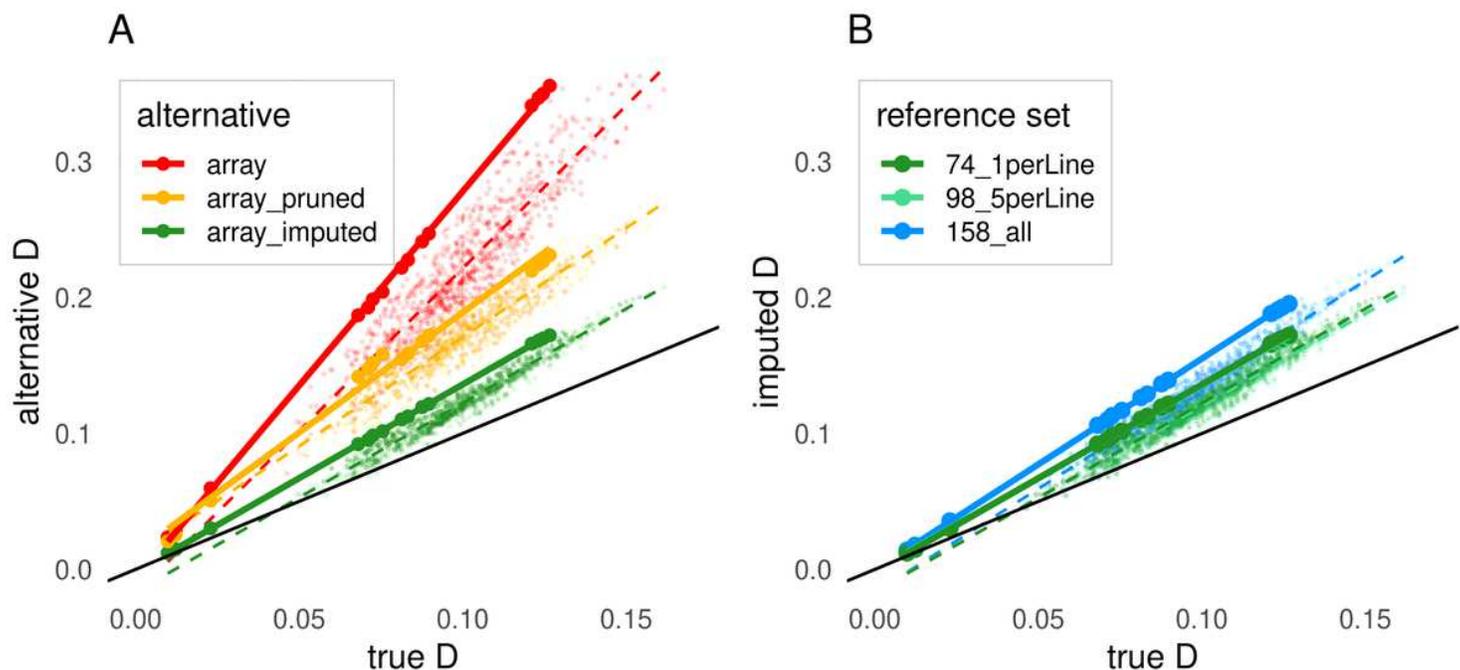


Figure 7

Effect of different correction strategies on ascertainment bias for Nei's standard genetic distance (D). A – uncorrected array, linkage pruned array and imputed array (reference set 74_1perLine) based vs. sequence-based D . B – array imputed with different reference sets vs. sequence-based HE. The solid black line represents the line of identity, the solid colored lines are regression lines within distances between individually sequenced populations (larger points) and the dashed lines regression lines within distances between all populations which include individually and pooled (small points) sequenced populations.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [FigureS1.tiff](#)
- [FigureS2.tiff](#)
- [FigureS3.tiff](#)
- [FigureS4.tiff](#)
- [FigureS5.tiff](#)
- [FigureS6.tiff](#)
- [FigureS7.tiff](#)
- [FigureS8.tiff](#)
- [FigureS9.tiff](#)
- [FigureS10.tiff](#)
- [SupplementaryFile1.csv](#)
- [SupplementaryFile2.pdf](#)
- [TableS1.docx](#)