

1 Removing array-specific batch effects in GWAS mega-analyses  
2 by applying a two-step imputation workflow reveals  
3 new associations for thyroid volume and goiter  
4

5 M. Kamal Nasr<sup>1,2</sup>, Eva König<sup>3</sup>, Christian Fuchsberger<sup>3</sup>, Sahar Ghasemi<sup>4</sup>, Uwe  
6 Völker<sup>2,5</sup>, Henry Völzke<sup>2,6</sup>, Hans J. Grabe<sup>1,7</sup>, Alexander Teumer<sup>1,2</sup>

7

8 <sup>1</sup> Department of Psychiatry and Psychotherapy, University Medicine Greifswald,  
9 Greifswald, Germany

10 <sup>2</sup> DZHK (German Centre for Cardiovascular Research), Partner Site Greifswald,  
11 Greifswald, Germany

12 <sup>3</sup> Institute for Biomedicine, Eurac Research, Bolzano, Italy

13 <sup>4</sup> Institute of Genetic Epidemiology, Medical Center – University of Freiburg, Faculty  
14 of Medicine, University of Freiburg, Germany

15 <sup>5</sup> Interfaculty Institute for Genetics and Functional Genomics, University Medicine  
16 Greifswald, Greifswald, Germany

17 <sup>6</sup> Institute for Community Medicine, University Medicine Greifswald, Greifswald,  
18 Germany

19 <sup>7</sup> German Center for Neurodegenerative Diseases (DZNE), Site Rostock/Greifswald,  
20 Greifswald, Germany

21

22 Keywords: genetic imputation, association studies, mega-analysis

23 Correspondence to:

24 Dr. Alexander Teumer

25 University Medicine Greifswald

26 Department of Psychiatry and Psychotherapy

27 Ellernholzstr. 1-2

28 17475 Greifswald, Germany

29 [ateumer@uni-greifswald.de](mailto:ateumer@uni-greifswald.de)

30 Phone: +49 (0) 3834-86-6918

31 ORCID: <https://orcid.org/0000-0002-8309-094X>

32

33

34 **Abstract**

35 **Background:** Combining individual-level data in genetic association studies (mega-  
36 analyses) enhances statistical power for identifying gene-trait associations. However,  
37 batch effects from combining variants of different arrays pose a major limitation.  
38 Here, we developed a two-step imputation workflow to overcome the array type bias.

39 **Methods:** Genotype data of 10,647 individuals generated using five different arrays  
40 were included. Intermediate array-specific panels were generated and subsequently  
41 imputed against the 1000 Genomes Project Phase3 reference panel. Genetic  
42 principal component (PC) analysis assessed batch effects in the cohort-combined  
43 imputed data. The workflow's performance was evaluated by comparing imputation  
44 quality  $r^2$  and allele frequency difference of the proposed two-step imputation to the  
45 conventional array-specific imputation as well as its matching with a whole-genome  
46 sequenced subgroup for further validation. We performed a genome-wide association  
47 study (GWAS) to test for genetic associations with goiter risk and thyroid gland  
48 volume, comparing summary statistics of both approaches.

49 **Results:** The proposed workflow eliminated the batch effect from the first twenty  
50 genetic PCs. The outcome of the workflow also showed high correlation with the  
51 conventional approach for allele frequencies ( $r^2 > 0.99$ ). GWAS results from the two-  
52 step imputation confirmed known associations on thyroid traits and revealed novel  
53 loci for thyroid volume (*TG*, *PAX8*, *IGFBP5*, *NRG1*), and one novel locus for goiter  
54 (*XKR6*), which was not statistically significant following the GWAS meta-analysis of  
55 conventional imputation.

56 **Conclusion:** Our imputation workflow provides high-quality imputation results without  
57 technical batch effects, fostering mega-analysis involving multiple genotyping arrays  
58 for different genetic association analysis.

59

60

## 61 **Introduction**

62 Genome-wide association studies (GWAS) represent an agnostic approach for  
63 identifying genetic associations with common traits and diseases by testing millions  
64 of variants with continuous outcomes or between groups. The testing checks allele  
65 frequency differences between individuals in a selected population that show  
66 different representations of the trait value <sup>1</sup>. GWAS analyses have been utilized in  
67 more than 5,700 studies, exploring 3,300 different traits <sup>1,2</sup>. The results of these  
68 analysis enriched our knowledge about disease risk variants, as well as identifying  
69 individuals with high disease-risk profiles through risk scores for complex heritable  
70 traits <sup>1,3</sup>.

71 The power to detect associations increases with sample size and number of variants  
72 tested <sup>3,4</sup>, rare and low frequency single nucleotide variants (SNVs) are of high  
73 interest. This requires the inclusion of large sample sizes in the experimental design,  
74 which is complicated or even unrealistic for rare diseases and low number of cases in  
75 the investigated populations. One alternative approach is a meta-analysis of  
76 summary statistics from different GWAS analyses of the same trait, performed either  
77 on same or different populations <sup>4,5</sup>. Another approach is to combine the individual-  
78 level data of samples from different cohorts to perform a mega-analysis <sup>6</sup>.

79 Meanwhile, genetic imputation is a reliable method to estimate alleles of variants not  
80 directly genotyped on an array. Based on linkage disequilibrium (LD), genetic  
81 imputation can significantly increase coverage of the human genome when using  
82 different commercial genotyping arrays <sup>7,8</sup>. This method is utilized as an alternative  
83 for expensive whole genome sequencing <sup>9</sup>.

84 GWAS on imputed variants and subsequent meta-analysis are frequently utilized  
85 cost-effective approaches for conducting large GWAS<sup>10</sup>, however, they are also  
86 associated with technical limitations <sup>11</sup>. If the sample size or number of cases in a  
87 specific cohort especially in a (nested) case-control design is low, the analysis results  
88 are less reliable because the effective number of samples is too small for the  
89 association models like linear regression, logistic regression or mixed-effects models,  
90 particularly when it comes to low-frequency variants. In addition, the analysis  
91 workload substantially increases with the number of cohorts included in a project <sup>12</sup>.

92 Former studies have shown that both meta-analyses and mega-analyses using  
93 individual participant data are mathematically equivalent <sup>13</sup>, and also comparable

94 when using imputed genotypes<sup>14</sup>. However, previous studies have mainly focused  
95 on data from cohorts genotyped on the same array type. Highlighting challenges in  
96 mega-analysis, where individual-level data from multiple cohorts genotyped on  
97 different arrays are combined.

98 We identified a concerning technical bias when combining the imputed data of  
99 individuals from genetically homogeneous northeastern German cohorts that were  
100 genotyped using diverse array types<sup>15,16</sup>. Principal component analysis (PCA)  
101 conducted on the quality-controlled genotype data revealed variation influenced by  
102 the array type. Notably, this variation was detected upon both imputation against the  
103 haplotype reference consortium (HRC)<sup>17</sup> and the 1000 Genomes v5 (1000G)  
104 reference panels<sup>18</sup>.

105 Here, we propose a newly developed workflow, to minimize the technical bias due to  
106 batch effect when combining genotype data from different arrays. The workflow is  
107 composed of two imputation steps. First we impute the included genotype datasets  
108 pairwise against each other and then create a panel of overlapping variants for each  
109 imputation outcome. Finally, we impute the generated intermediate panel against one  
110 of the commonly available large panels. We evaluate the outcome of the new  
111 workflow in comparison to conventional imputation approaches. In an application  
112 example, we conducted a GWAS analysis on thyroid gland volume and goiter,  
113 identifying new associations with goiter while demonstrating the robustness of our  
114 imputation workflow by validating known associations.

## 115 **Materials and Methods**

### 116 *Workflow design overview*

117 In this project, three different approaches for genetic imputation have been  
118 conducted. The first proposed approach is a new two-step imputation process. The  
119 second and third approaches are conventional single-step imputation for comparison  
120 with the newly proposed method. We modified the third approach to be a single-step  
121 imputation using only the intersecting variants genotyped on all included array types.  
122 The genotype data that were used for imputation were the same in both imputation  
123 approaches. We used data from five different arrays (Affymetrix SNP 6.0, Affymetrix  
124 Axiom [Thermo Fisher Scientific, Santa Clara, CA, USA], and Illumina Omni 2.5,  
125 Illumina GSA, and Illumina PsychArray [Illumina, Inc., San Diego, CA, USA])  
126 obtained from samples of the German GANI\_MED (n= 2410), SHIP-START (n=

127 4070) and SHIP-TREND (n= 4119) (Figure 1) <sup>15,16</sup>. The individuals genotyped on the  
128 Affymetrix Axiom array (n= 48) were a subset of the Affymetrix SNP 6.0 samples,  
129 whereas all other individuals were genotyped only once. Detailed information  
130 regarding the included datasets, sample sizes, and genotype arrays are provided in  
131 Figure 1 and the Supplementary Methods. For evaluation of the imputation  
132 performance, a whole-genome sequenced subset of 192 individuals from SHIP-  
133 TREND has been used for genotype matching concordance checking.

#### 134 *Data pre-processing and phasing*

135 Genotype quality control was performed using PLINK <sup>19</sup> as reported previously <sup>20,21</sup>.  
136 Briefly, arrays with genotype call rate <94%, as well as variants with missing call rate  
137 >5%, Hardy-Weinberg equilibrium p-value <10<sup>-4</sup>, and monomorphic SNVs and  
138 singletons were excluded. All the included genotype datasets were aligned to  
139 reference genome build (GRCh37) using BCFtools software followed by retaining  
140 only sites with at least one alternative allele <sup>22</sup>. Haplotype phasing preceding the first  
141 imputation round was performed for each genotype dataset using Eagle2 (5 phasing  
142 iterations and number of conditioning haplotypes = 10<sup>4</sup>) and estimated to hg19  
143 mapping <sup>23</sup>. No external reference panel was used for phasing.

#### 144 *Two-step imputation*

145 In the first step, the intermediate imputation, each phased panel was used as a  
146 reference panel after compression with the minimac4 imputation software <sup>24</sup>. Each of  
147 the five pre-processed genotype arrays was imputed pairwise against the other four,  
148 yielding twenty imputed datasets. Using an R scripted algorithm with VCFtools and  
149 BCFtools for variants filtration <sup>22</sup>, each of these imputed variant sets were subset to  
150 the variants overlapping between all generated panels with high imputation quality  
151 ( $R^2 \geq 0.9$ ), selecting the source panel with the highest imputation quality score for  
152 each variant. The outcome of the first-step imputation was one VCF file for each  
153 cohort dataset, which was then used as input for the second imputation step against  
154 1000G reference panel using the Michigan imputation server <sup>24</sup>. Eagle2 was selected  
155 for phasing without (additional) imputation quality filters applied at this stage of the  
156 imputation.

157 *Conventional imputation*

158 To validate the outcome of the proposed workflow, we imputed the genotype data of  
159 each array directly against 1000G reference panel using the same parameters and  
160 quality control procedures as described in the two-step imputation.

161 *Conventional imputation with overlapping genotypes*

162 To evaluate a simple approach for removing the array type specific batch effect, we  
163 ran a single-step imputation by restricting the input variants to those assessed on all  
164 array types. This approach was tested for three (Affymetrix SNP 6.0, Illumina Omni  
165 2.5 and Illumina GSA) and all five array types. We compared the imputation quality of  
166 these imputations to the conventional and the two-step imputation exemplarily for the  
167 SHIP-TREND.

168 *Batch effect assessment*

169 Technical bias was assessed by estimating genetic principal components of the  
170 imputed and quality-controlled genotype data. Quality control filters included missing  
171 variant call rates above 5%, Hardy-Weinberg equilibrium p-values less than  $10^{-4}$ ,  
172 minor allele frequency (MAF) less than 1%, correlated variants were removed by LD  
173 pruning with a window size of 50, step size of 5 SNVs and  $R^2$  threshold of 0.2.  
174 Genetic principal components, along with their explained variances were compared  
175 between two-step imputation and conventional imputation. To check the robustness  
176 of the two-step imputation approach for rare variants, we further ran PCA on rare  
177 variants (MAF < 1%). All PCAs are performed using PLINK 2.0 software <sup>19</sup>. Principal  
178 components (PC) of the first twenty components were plotted with ggplot2 package  
179 of R software.

180 *Evaluation of imputation performance*

181 Genotype data attributes, including differences in MAF and imputation quality  $R^2$   
182 measure between the two-step and conventional imputation were compared for each  
183 included cohort after stratification by minor allele frequency. We additionally checked  
184 the concordance (number of matching genotypes/total number of genotypes <sup>25</sup>) of the  
185 imputed best-guess (hard call) genotypes with their whole-genome sequencing data  
186 of a subset of 192 individuals from SHIP-TREND (Omni 2.5) for which whole genome  
187 sequencing (WGS) data was also available.

## 188 *GWAS and results comparison*

189 We evaluated the results of GWAS analyses using genotype dosage information from  
190 both two-step and conventional imputation approaches using thyroid gland volume  
191 and goiter risk as outcomes as described before <sup>26</sup>, and summarized in the  
192 Supplementary Methods.

193 All GWAS were conducted with the EPACTS 3.3.2 software using dosages <sup>27</sup>.  
194 Inverse-variance weighted meta-analysis of the single cohort GWAS summary  
195 statistics based on the conventional imputation approach was done using the METAL  
196 software <sup>4</sup>. We compared the summary statistics p-values, effect estimates and its  
197 standard errors. A detailed description of the GWAS analysis plan is available in the  
198 Supplementary Methods. Independent (lead) variants associated with thyroid traits  
199 were identified using the clumping function of PLINK with a threshold of  $p < 5 \times 10^{-8}$ ,  
200  $r^2 > 0.01$  and 1 Mb distance, and compared across the imputation approaches. For  
201 biological validation of the findings, lead variants were investigated by checking their  
202 association with thyrotropin to support their biological plausibility <sup>20</sup>.

## 203 **Results**

204 Genotype datasets of 10,647 individuals genotyped on five different array types were  
205 included in the workflow (Figure 1). After the intermediate imputation, 1,942,499 high-  
206 quality overlapping autosomal variants were generated for each of the included five  
207 datasets, and subsequently used for the second imputation step against the 1000G  
208 reference panel. In contrast, only 58,091 autosomal genotyped variants were  
209 available as imputation input in the three cohorts using the conventional imputation  
210 approach with overlapping variants. This number further dropped to 12,265 variants  
211 when including the Affymetrix Axiom and Illumina PsychArrays. Supplementary Table  
212 1 provides the distribution of these variants for each chromosome in the two-step  
213 imputation, the conventional imputation and conventional imputation using  
214 overlapped variants.

## 215 *Batch effect investigation (PCA analysis)*

216 For the conventional imputation, 939,802 variants with a total genotyping rate of  
217 0.988 were included in the PCA analysis, with 2,314,205 variants removed due to  
218 missing genotype threshold, 9,118 variants due to Hardy-Weinberg equilibrium, and  
219 38,008,458 variants due to the MAF threshold. Projecting the first two principle



220 components of the imputed genotypes clustered the samples by array types with a  
221 combined explained variance of 0.242 (Figure 2a, b, Supplementary Figure 1a).

222 The two-step imputation included 1,360,834 variants with a total genotyping rate of  
223 0.989 in the PCA analysis after removing 725,008 variants due to missing genotype  
224 threshold, 5,363 variants due to Hardy-Weinberg equilibrium, and 37,753,448  
225 variants due to the MAF threshold. Using this approach, the first 20 principal  
226 components of the imputed genotypes did not show any array type specific clusters  
227 (Figure 2c, d, Supplementary Figure 1b). PCA restricted to the rare genotypes  
228 imputed with the two-step approach did not show clustering by the array type  
229 (Supplementary Figure 2).

230 The conventional imputation using only overlapping variants also removed the batch  
231 effect (Supplementary Figure 3), but led to a drastic decrease in the quality of the  
232 imputed variants for all allele frequency groups in comparison to the other  
233 approaches (Supplementary Figure 4). Median  $R^2$  of the overlapping variants was  
234 0.011, while for two-step imputation and conventional imputation using all variants  
235 was 0.993 for variants with  $MAF > 0.05$ , which was the main reason for not  
236 considering this approach any further in the performance evaluation and GWAS  
237 analysis, focusing only on the two-step imputation against conventional imputation  
238 using all variants.

### 239 *Two-step imputation outcome parameters comparison with conventional imputation*

240 Our newly developed two-step imputation approach showed higher overall quality for  
241 the imputation outcome of rarer variants ( $MAF < 0.01$ ) than the conventional  
242 imputation (Figure 3, Supplementary Figure 5). For less-common variants with MAF  
243 between 0.04 and 0.05, the median  $R^2$  increased by 0.19 when using the two-step  
244 imputation approach for the GANI\_MED Illumina PsychArray. It also shows  
245 comparable allele frequencies to the conventional approach, judging by the absolute  
246 difference in allele frequencies for each imputed variant in both approaches (Figure  
247 4). Detailed statistics about median  $R^2$  and MAF differences are presented in  
248 Supplementary Table 2.

249 Hard call genotype concordance results in the subgroup of SHIP-TREND with  
250 matching WGS data showed strong correlation between two-step imputation and  
251 sequenced genotypes. The Pearson correlation coefficients with the homozygous



252 reference, heterozygous, and homozygous alternative genotypes were 0.996, 0.981,  
253 and 0.979, respectively (Table 1, Supplementary Figure 6).

#### 254 *GWAS on thyroid traits*

255 A total number of 6,894 individuals from SHIP-START and SHIP-TREND (both array  
256 types) with thyroid measurement information were included in the GWAS analysis, all  
257 participants were of European ancestry. Detailed information about cohort  
258 characteristics is presented in Table 2. The estimated genomic control for all  
259 conducted GWAS analysis showed no signs of inflation, with a minimum and  
260 maximum  $\lambda_{GC} = 1.001$  and 1.038, respectively (Supplementary Figure 7).

261 Meta-analysis of the GWAS analysis on goiter risk using conventionally imputed  
262 genotypes revealed four significantly associated loci ( $p < 5 \times 10^{-8}$ ). Confirming the  
263 results from previously conducted GWAS analysis<sup>26</sup>. Two of the loci are located at  
264 the *CAPZB* region in chromosome 1, the other two at the *FAM227B* and *MAFTRR*  
265 regions on chromosome 15 and 16, respectively (Figure 5b). However, the GWAS  
266 analysis of the combined two-step imputed genotype data revealed another  
267 associated locus at the *XKR6* region on chromosome 8 (Figure 5c), which did not  
268 attain statistical significance using the conventional meta-analysis approach (Figure  
269 5a).

270 The conventional approach of the GWAS meta-analysis of thyroid volume revealed  
271 four novel associations at the *PAX8* region on chromosome 2, at *IGFBP5*, *NRG1* and  
272 *TG* (Figure 6b,c), and confirmed all known associations with this trait<sup>20</sup>. Genome-  
273 wide significance for these regions were also obtained with the GWAS using the  
274 combined two-step imputed genotype data, with the exception of *NRG1* ( $p$ -value =  
275  $5.22 \times 10^{-8}$ ). (Figure 6a). Except the *PAX8* region, all associated loci were also  
276 associated with thyroxin in recently published multi-trait GWAS meta-analysis  
277 analysis for thyroid function<sup>20</sup>.

278 Table 3 summarizes the results of the SNVs with the strongest association of both  
279 GWAS approaches and traits. The significant GWAS results of both approaches were  
280 comparable, judging by the magnitude of the estimates. However, the standard errors  
281 of the GWAS results obtained from the two-step imputed genotypes were generally  
282 lower, where the natural log p-values were slightly higher in the linear regression and  
283 lower in the logistic regression based analyses compared to the conventional  
284 imputation and subsequent meta-analysis. (Supplementary Figures 8 and 9).

## 285 **Discussion**

286 Genome-wide association studies have helped identify genetic factors for many traits  
287 and diseases. Since the 1960s, collecting genotype samples from diverse  
288 populations has grown in importance<sup>28,29</sup>. Imputation techniques now address gaps  
289 in whole-genome sequencing<sup>30</sup>. However, variations in these methods across  
290 cohorts reduce the accuracy of genetic analyses due to biases from differing array  
291 technologies<sup>31,32</sup>. For instance, in our genotyped samples where we compared the  
292 allele frequencies of the SHIP-TREND subgroup genotyped using Illumina Omnia 2.5  
293 with the corresponding whole genome sequencing variants (Supplementary Figure  
294 10), this array-specific variation in the allele frequency, both in common and rare  
295 variants, can affect LD estimation in haplotype phasing and genotype imputation. Our  
296 newly developed imputation method addresses the additional variation induced by  
297 including multiple array types. By forming an intermediate panel of high-quality  
298 variants, the method enables high imputation accuracy while removing the array type  
299 induced batch effects.

300 The developed workflow was inspired by the use of only overlapping variants in the  
301 included cohorts (Supplementary Table 1), which eliminated the observed bias, yet  
302 led to a significant decrease in the imputation quality due to removal of informative  
303 tag SNVs (Supplementary Figure 3), the influence of informative SNVs on imputation  
304 quality has been shown in previous research work<sup>33</sup>. Based on that, we introduced  
305 an intermediate imputation step to generate a panel that can retain the same  
306 genotype information of the included arrays and thus preserving its LD structure.

307 Selecting an appropriate threshold for the imputation  $R^2$  ( $\geq 0.9$ ) of the overlapping  
308 variants in the intermediate panel was essential for having reliable genotype  
309 information upon imputing against 1000G panel. The value of the threshold was  
310 decided following the output of several imputations using different thresholds. We  
311 aimed to use the highest possible value for imputation quality  $R^2$  without affecting  
312 haplotype phasing results due to removing too many variants. Supplementary Table  
313 1 shows the number of included variants per chromosome when the  $R^2$  threshold  
314 was adjusted to 0.8 as well as 0.9. Both approaches led to similar imputation outputs  
315 as seen after plotting genetic PCs of the imputation outcomes (Figure 2 and  
316 Supplementary Figure 11).

317 To evaluate the existence of the array type bias, we used genetic PCs to capture the  
318 main variance in the allele frequencies, which can be seen as the influence of the  
319 array type after projecting the PCs of the eigenvectors and evaluate the homogeneity  
320 of the projected points on PC axis<sup>34</sup>. The imputation following the developed two-  
321 step approach showed its capability to overcome the array type differentiation  
322 whenever existing in the first twenty components, compared to conventional  
323 imputation outcomes where a clear clustering effect by the array type was observed  
324 (Figure 2).

325 To test the performance of the imputation outcomes, we focused on comparing allele  
326 frequency and imputation quality parameters of the cohort genotypes. These  
327 parameters are particularly relevant for conducting trait-association studies. The  
328 developed imputation flow showed strong matching of the allele frequency,  
329 represented by minimized difference in allele frequency for each variants between  
330 the imputation approaches (Figure 4 and Supplementary Table 2). The developed  
331 imputation workflow was successful in providing more reliable genetic predictions for  
332 imputed rarer alleles (Figure 3). Moreover, a strong correlation in the median  $R^2$  of  
333 the two imputation approaches using the SHIP-TREND Omni 2.5 arrays was  
334 observed in comparison to other cohorts. This could be due to the high variant  
335 coverage of the array used for genotyping the cohort's samples. Further analysis of  
336 the allele frequency variance for each variant from both imputation outcomes of this  
337 specific cohort shows the strong correlation in the autosomal allele frequency  
338 (correlation  $r^2 > 0.999$ ) (Supplementary Figure 12). Nevertheless, SHIP-TREND  
339 imputed genotypes showed strong concordance with the corresponding genotyped  
340 variants in the whole-genome sequenced subgroup, indicating the representation of  
341 well estimated genotypes in both rarer and common variants (Table 1 and  
342 Supplementary Table 3). Using the allele frequencies of the WGS subgroup as a gold  
343 standard for comparison of the imputation results is somewhat misleading because  
344 differences in the genotype frequencies exist already between genotyped variants on  
345 the arrays and the WGS (Supplementary Figure 10).

346 The aim of conducting GWAS as part of the imputation workflow evaluation is to  
347 compare regression analyses results of the combined two-step imputed datasets to  
348 the conventional inverse-variance meta-analysis of the same samples. Goiter risk  
349 and thyroid volume are both suitable traits for our evaluation as they represent  
350 different trait datatypes with a true positive genetic association in SHIP<sup>26</sup>. The results

351 of GWAS analysis for the two-step imputed data showed two key advantages. First,  
352 we were able to identify a new association in the *XKR6* gene region with goiter risk  
353 as a dichotomous trait which did not reach statistical significance in the conventional  
354 meta-analysis. Variants in this locus were also associated with thyroxin in a former  
355 GWAS, making this association biologically reasonable<sup>20</sup>. A slightly better  
356 performance of the mega-analysis vs. a meta-analysis for a logistic regression  
357 GWAS is in line with the results of a former study<sup>14</sup>. The second point was its  
358 consistency with conventional imputation GWAS meta-analysis in revealing a novel  
359 locus associated with thyroid gland volume (*PAX8*). *PAX8* is a paired box family gene  
360 member that was found to be associated with the development of thyroid gland in  
361 embryonic development, and its transcription is a diagnostic marker for anaplastic  
362 thyroid carcinoma<sup>35-37</sup>, and thus represents also a plausible association with thyroid  
363 volume in adults. All other significant GWAS findings represent also plausible true  
364 positive associations as they confirmed former findings of these traits (including  
365 replication in an independent cohort)<sup>26</sup>, or were associated with thyroid function<sup>20</sup>.

366 The comparison of the estimates and standard errors of both GWAS approaches did  
367 not show signs of p-value inflation. However, it showed that the GWAS of the two-  
368 step imputation data had a slight decrease in standard errors, in comparison to the  
369 meta-analysis outcomes.

370 Besides its role in discovering more hits in case-control GWAS analysis, our  
371 developed two-step imputation workflow is not restricted to a specific array, the  
372 inclusion of different arrays with different sequencing technologies in the workflow  
373 has proven its robustness in overcoming the array bias regardless of the number or  
374 the type of the array. It is also applicable for other generalized imputation panels like  
375 TOPMed or HRC reference panels<sup>24</sup>. These two strengths shall enable the utilization  
376 of combined genotype information for better understanding of rare diseases or  
377 genetic associations in populations that are represented in small sample sizes.

378 Although the array type specific batch effect in GWAS mega-analyses might be  
379 reduced in specific scenarios by adjusting for genetic PCs, such a correction will not  
380 be possible in all analyses. Such analyses include the ones using polygenic scores.  
381 Our two-step imputation provides a powerful solution for combining datasets while  
382 reducing technical bias also for analyses of polygenic scores.

383 Comparing the imputation outcomes to the whole genome-sequenced data had  
384 limitations for evaluating the proposed workflow, either due to the relatively small  
385 sample size of individuals who underwent WGS ( $n = 192$  after QC) and its exclusivity  
386 for one of the included five cohorts. Although the total sample size provided us the  
387 possibility of analysing low frequency alleles, it was likely too small for evaluating the  
388 performance of the different approaches for very rare alleles. However, our two-step  
389 imputation approach seems to be on average superior to the classical imputation with  
390 regard to the imputation quality measure (Figure 3), while the difference in allele  
391 frequency is small particularly for rarer variants (Figure 4 and Supplementary Table  
392 2). As indicated also in these results, the average difference in allele frequency  
393 seems to depend also on the density and design of the underlying genotyping array.

394 While the imputed genotypes generated from the proposed workflow will be utilized  
395 for identifying novel genetic associations in SHIP and GANI-MED, it will be interesting  
396 to evaluate the impact of the developed workflow on other cohorts, especially those  
397 comprising diverse ancestries. Thus far, all the included cohorts are of European  
398 ancestry from the North-East of Germany, highlighting the need to test the workflow's  
399 efficacy on genotyped cohorts from other or multiple genetic ancestries. Such  
400 evaluations will help to determine the generalizability and robustness of the  
401 imputation method across different genetic backgrounds, ultimately enhancing its  
402 utility in global genomic research.

403 In conclusion, our developed two-step imputation workflow aims to overcome the  
404 array type bias, by creating an intermediate panel of high-quality overlapping imputed  
405 variants. This approach enables the conduction of mega-analysis by combining  
406 genotype information from different arrays without inducing a technical array type  
407 effect. Our workflow will increase statistical power for conducting large-scale GWAS  
408 mega-analyses and other genetic analyses like polygenic risk score calculations,  
409 playing an important role in genetic research and its application in individualized  
410 medicine.

411

#### 412 **Acknowledgements:**

413 SHIP is part of the Community Medicine Research net of the University of Greifswald,  
414 Germany, which is funded by the Federal Ministry of Education and Research (grants  
415 no. 01ZZ9603, 01ZZ0103, and 01ZZ0403), the Ministry of Cultural Affairs as well as

416 the Social Ministry of the Federal State of Mecklenburg-West Pomerania, and the  
417 network 'Greifswald Approach to Individualized Medicine (GANI\_MED)' funded by the  
418 Federal Ministry of Education and Research (grant 03IS2061A). Genome-wide data  
419 have been partly supported by the Federal Ministry of Education and Research (grant  
420 no. 03ZIK012) and a joint grant from Siemens Healthineers, Erlangen, Germany and  
421 the Federal State of Mecklenburg- West Pomerania. The project is funded by the  
422 Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) –  
423 455978266 (A.T.). The authors are grateful to Linda Garvert for testing and helpful  
424 feedback on the two-step imputation workflow.

#### 425 **Author Contributions:**

426 Project design and supervision: A.T.

427 Analyses: M.K.N., E.K., S.G.

428 Interpretation of the results: M.K.N., A.T., E.K., C.F.

429 Drafting of the manuscript: M.K.N., A.T.

430 Providing genotype and phenotype data: U.V., H.V., H.J.G

431 Critical review of the manuscript: all authors

#### 432 **Data and code availability**

433 Developed scripts for the workflow are available on github  
434 ([https://github.com/GenEpi-psych-UMG/Two\\_Step\\_Imputation](https://github.com/GenEpi-psych-UMG/Two_Step_Imputation)). The data of the SHIP  
435 study cannot be made publically available due to the informed consent of the study  
436 participants, but it can be accessed through a data application form available at  
437 <https://transfer.ship-med.uni-greifswald.de/> for researchers who meet the criteria for  
438 access to confidential data. The full results of the GWAS summary statistics are  
439 available on the ThyroidOmics Consortium website (<http://www.thyroidomics.com>).

#### 440 **Consents and approvals**

441 The study followed the recommendations of the Declaration of Helsinki. The medical  
442 ethics committee of the University of Greifswald approved the study protocol, and  
443 oral and written informed consents were obtained from each of the study participants.

#### 444 **Conflicts of interest**

445 The authors have no affiliation with any organization with a direct or indirect financial  
446 interest in the subject matter discussed in the manuscript. HJG received travel grants

447 and speakers honoraria from Neuraxpharm, Servier, Indorsia and Janssen Cilag not  
448 related to the current project. HV received travel grants and speakers honoraria from  
449 Sanofi-Aventis not related to the current project.

450

451



## 452 References

- 453 1. Uffelmann, E., Huang, Q.Q., Munung, N.S., de Vries, J., Okada, Y., Martin,  
454 A.R., Martin, H.C., Lappalainen, T., and Posthuma, D. (2021). Genome-wide  
455 association studies. *Nat. Rev. Methods Prim.* *1*, 59.  
456 <https://doi.org/10.1038/s43586-021-00056-9>.
- 457 2. Watanabe, K., Stringer, S., Frei, O., Umičević Mirkov, M., de Leeuw, C.,  
458 Polderman, T.J.C., van der Sluis, S., Andreassen, O.A., Neale, B.M., and  
459 Posthuma, D. (2019). A global overview of pleiotropy and genetic architecture  
460 in complex traits. *Nat. Genet.* *51*, 1339–1348. [https://doi.org/10.1038/s41588-](https://doi.org/10.1038/s41588-019-0481-0)  
461 [019-0481-0](https://doi.org/10.1038/s41588-019-0481-0).
- 462 3. Khera, A. V., Chaffin, M., Aragam, K.G., Haas, M.E., Roselli, C., Choi, S.H.,  
463 Natarajan, P., Lander, E.S., Lubitz, S.A., Ellinor, P.T., et al. (2018). Genome-  
464 wide polygenic scores for common diseases identify individuals with risk  
465 equivalent to monogenic mutations. *Nat. Genet.* *50*, 1219–1224.  
466 <https://doi.org/10.1038/s41588-018-0183-z>.
- 467 4. Willer, C.J., Li, Y., and Abecasis, G.R. (2010). METAL: fast and efficient meta-  
468 analysis of genomewide association scans. *Bioinformatics* *26*, 2190–2191.  
469 <https://doi.org/10.1093/bioinformatics/btq340>.
- 470 5. Mägi, R., Horikoshi, M., Sofer, T., Mahajan, A., Kitajima, H., Franceschini, N.,  
471 McCarthy, M.I., COGENT-Kidney Consortium, T2D-GENES Consortium, A.P.,  
472 and Morris, A.P. (2017). Trans-ethnic meta-regression of genome-wide  
473 association studies accounting for ancestry increases power for discovery and  
474 improves fine-mapping resolution. *Hum. Mol. Genet.* *26*, 3639–3650.  
475 <https://doi.org/10.1093/hmg/ddx280>.
- 476 6. Abou-Khalil, B., Auce, P., Avbersek, A., Bahlo, M., Balding, D.J., Bast, T.,  
477 Baum, L., Becker, A.J., Becker, F., Berghuis, B., et al. (2018). Genome-wide  
478 mega-analysis identifies 16 loci and highlights diverse biological mechanisms  
479 in the common epilepsies. *Nat. Commun.* *9*, 5269.  
480 <https://doi.org/10.1038/s41467-018-07524-z>.
- 481 7. Burton, P.R., Clayton, D.G., Cardon, L.R., Craddock, N., Deloukas, P.,  
482 Duncanson, A., Kwiatkowski, D.P., McCarthy, M.I., Ouwehand, W.H., Samani,  
483 N.J., et al. (2007). Genome-wide association study of 14,000 cases of seven  
484 common diseases and 3,000 shared controls. *Nature* *447*, 661–678.  
485 <https://doi.org/10.1038/nature05911>.
- 486 8. De Marino, A., Mahmoud, A.A., Bose, M., Bircan, K.O., Terpolovsky, A.,  
487 Bamunusinghe, V., Bohn, S., Khan, U., Novković, B., and Yazdi, P.G. (2022). A  
488 comparative analysis of current phasing and imputation software. *PLoS One*  
489 *17*, e0260177. <https://doi.org/10.1371/journal.pone.0260177>.
- 490 9. Quick, C., Anugu, P., Musani, S., Weiss, S.T., Burchard, E.G., White, M.J.,  
491 Keys, K.L., Cucca, F., Sidore, C., Boehnke, M., et al. (2020). Sequencing and  
492 imputation in GWAS: Cost-effective strategies to increase power and genomic  
493 coverage across diverse populations. *Genet. Epidemiol.* *44*, 537–549.  
494 <https://doi.org/10.1002/gepi.22326>.
- 495 10. Pasaniuc, B., Rohland, N., McLaren, P.J., Garimella, K., Zaitlen, N., Li, H.,  
496 Gupta, N., Neale, B.M., Daly, M.J., Sklar, P., et al. (2012). Extremely low-  
497 coverage sequencing and imputation increases power for genome-wide  
498 association studies. *Nat. Genet.* *44*, 631–635. <https://doi.org/10.1038/ng.2283>.
- 499 11. Palmer, C., and Pe'er, I. (2016). Bias Characterization in Probabilistic  
500 Genotype Data and Improved Signal Detection with Multiple Imputation. *PLOS*  
501 *Genet.* *12*, e1006091. <https://doi.org/10.1371/journal.pgen.1006091>.
- 502 12. Winkler, T.W., Day, F.R., Croteau-Chonka, D.C., Wood, A.R., Locke, A.E.,

- 503 Mägi, R., Ferreira, T., Fall, T., Graff, M., Justice, A.E., et al. (2014). Quality  
504 control and conduct of genome-wide association meta-analyses. *Nat. Protoc.* 9,  
505 1192–1212. <https://doi.org/10.1038/nprot.2014.071>.
- 506 13. Lin, D.Y., and Zeng, D. (2010). On the relative efficiency of using summary  
507 statistics versus individual-level data in meta-analysis. *Biometrika* 97, 321–332.  
508 <https://doi.org/10.1093/biomet/asq006>.
- 509 14. Gorski, M., Günther, F., Winkler, T.W., Weber, B.H.F., and Heid, I.M. (2019).  
510 On the differences between mega- and meta-imputation and analysis  
511 exemplified on the genetics of age-related macular degeneration. *Genet.*  
512 *Epidemiol.* 43, 559–576. <https://doi.org/10.1002/gepi.22204>.
- 513 15. Völzke, H., Schössow, J., Schmidt, C.O., Jürgens, C., Richter, A., Werner, A.,  
514 Werner, N., Radke, D., Teumer, A., Ittermann, T., et al. (2022). Cohort Profile  
515 Update: The Study of Health in Pomerania (SHIP). *Int. J. Epidemiol.*  
516 <https://doi.org/10.1093/ije/dyac034>.
- 517 16. Grabe, H.J., Assel, H., Bahls, T., Dörr, M., Endlich, K., Endlich, N., Erdmann,  
518 P., Ewert, R., Felix, S.B., Fiene, B., et al. (2014). Cohort profile: Greifswald  
519 approach to individualized medicine (GANI\_MED). *J. Transl. Med.* 12, 144.  
520 <https://doi.org/10.1186/1479-5876-12-144>.
- 521 17. McCarthy, S., Das, S., Kretzschmar, W., Delaneau, O., Wood, A.R., Teumer,  
522 A., Kang, H.M., Fuchsberger, C., Danecek, P., Sharp, K., et al. (2016). A  
523 reference panel of 64,976 haplotypes for genotype imputation. *Nat. Genet.* 48,  
524 1279–1283. <https://doi.org/10.1038/ng.3643>.
- 525 18. 1000 Genomes Project Consortium, A., Auton, A., Brooks, L.D., Durbin, R.M.,  
526 Garrison, E.P., Kang, H.M., Korbel, J.O., Marchini, J.L., McCarthy, S., McVean,  
527 G.A., et al. (2015). A global reference for human genetic variation. *Nature* 526,  
528 68–74. <https://doi.org/10.1038/nature15393>.
- 529 19. Chang, C.C., Chow, C.C., Tellier, L.C., Vattikuti, S., Purcell, S.M., and Lee, J.J.  
530 (2015). Second-generation PLINK: rising to the challenge of larger and richer  
531 datasets. *Gigascience* 4, 7. <https://doi.org/10.1186/s13742-015-0047-8>.
- 532 20. Sterenborg, R.B.T.M., Steinbrenner, I., Li, Y., Bujnis, M.N., Naito, T., Marouli,  
533 E., Galesloot, T.E., Babajide, O., Andreassen, L., Astrup, A., et al. (2024). Multi-  
534 trait analysis characterizes the genetics of thyroid function and identifies causal  
535 associations with clinical implications. *Nat. Commun.* 15, 888.  
536 <https://doi.org/10.1038/s41467-024-44701-9>.
- 537 21. Teumer, A., Tin, A., Sorice, R., Gorski, M., Yeo, N.C., Chu, A.Y., Li, M., Li, Y.,  
538 Mijatovic, V., Ko, Y.-A., et al. (2016). Genome-wide Association Studies Identify  
539 Genetic Loci Associated With Albuminuria in Diabetes. *Diabetes* 65, 803–817.  
540 <https://doi.org/10.2337/db15-1313>.
- 541 22. Danecek, P., Bonfield, J.K., Liddle, J., Marshall, J., Ohan, V., Pollard, M.O.,  
542 Whitwham, A., Keane, T., McCarthy, S.A., Davies, R.M., et al. (2021). Twelve  
543 years of SAMtools and BCFtools. *Gigascience* 10.  
544 <https://doi.org/10.1093/gigascience/giab008>.
- 545 23. Loh, P.-R., Danecek, P., Palamara, P.F., Fuchsberger, C., A Reshef, Y., K  
546 Finucane, H., Schoenherr, S., Forer, L., McCarthy, S., Abecasis, G.R., et al.  
547 (2016). Reference-based phasing using the Haplotype Reference Consortium  
548 panel. *Nat. Genet.* 48, 1443–1448. <https://doi.org/10.1038/ng.3679>.
- 549 24. Das, S., Forer, L., Schönherr, S., Sidore, C., Locke, A.E., Kwong, A., Vrieze,  
550 S.I., Chew, E.Y., Levy, S., McGue, M., et al. (2016). Next-generation genotype  
551 imputation service and methods. *Nat. Genet.* 48, 1284–1287.  
552 <https://doi.org/10.1038/ng.3656>.
- 553 25. Wuttke, M., König, E., Katsara, M.-A., Kirsten, H., Farahani, S.K., Teumer, A.,

- 554 Li, Y., Lang, M., Göcmen, B., Pattaro, C., et al. (2023). Imputation-powered  
555 whole-exome analysis identifies genes associated with kidney function and  
556 disease in the UK Biobank. *Nat. Commun.* *14*, 1287.  
557 <https://doi.org/10.1038/s41467-023-36864-8>.
- 558 26. Teumer, A., Rawal, R., Homuth, G., Ernst, F., Heier, M., Evert, M.,  
559 Dombrowski, F., Völker, U., Nauck, M., Radke, D., et al. (2011). Genome-wide  
560 association study identifies four genetic loci associated with thyroid volume and  
561 goiter risk. *Am J Hum Genet* *88*, 664–673.  
562 <https://doi.org/10.1016/j.ajhg.2011.04.015>.
- 563 27. Kang HM (2016). EFACTS: Efficient and Parallelizable Association Container  
564 Toolbox.
- 565 28. Nelson, M.R., Bryc, K., King, K.S., Indap, A., Boyko, A.R., Novembre, J., Briley,  
566 L.P., Maruyama, Y., Waterworth, D.M., Waeber, G., et al. (2008). The  
567 Population Reference Sample, POPRES: A Resource for Population, Disease,  
568 and Pharmacological Genetics Research. *Am. J. Hum. Genet.* *83*, 347–358.  
569 <https://doi.org/10.1016/j.ajhg.2008.08.005>.
- 570 29. Evaluating Human Genetic Diversity (1997). (National Academies Press)  
571 <https://doi.org/10.17226/5955>.
- 572 30. TRECCANI, M., LOCATELLI, E., PATUZZO, C., and MALERBA, G. (2023). A  
573 broad overview of genotype imputation: Standard guidelines, approaches, and  
574 future investigations in genomic association studies. *BIOCELL* *47*, 1225–1241.  
575 <https://doi.org/10.32604/biocell.2023.027884>.
- 576 31. Lachance, J., and Tishkoff, S.A. (2013). SNP ascertainment bias in population  
577 genetic analyses: Why it is important, and how to correct it. *BioEssays* *35*,  
578 780–786. <https://doi.org/10.1002/bies.201300014>.
- 579 32. Geibel, J., Reimer, C., Weigend, S., Weigend, A., Pook, T., and Simianer, H.  
580 (2021). How array design creates SNP ascertainment bias. *PLoS One* *16*,  
581 e0245178. <https://doi.org/10.1371/journal.pone.0245178>.
- 582 33. Wojcik, G.L., Fuchsberger, C., Taliun, D., Welch, R., Martin, A.R.,  
583 Shringarpure, S., Carlson, C.S., Abecasis, G., Kang, H.M., Boehnke, M., et al.  
584 (2018). Imputation-Aware Tag SNP Selection To Improve Power for Large-  
585 Scale, Multi-ethnic Association Studies. *G3 Genes|Genomes|Genetics* *8*,  
586 3255–3267. <https://doi.org/10.1534/g3.118.200502>.
- 587 34. Elhaik, E. (2022). Principal Component Analyses (PCA)-based findings in  
588 population genetic studies are highly biased and must be reevaluated. *Sci.*  
589 *Rep.* *12*, 14683. <https://doi.org/10.1038/s41598-022-14395-4>.
- 590 35. Bishop, J.A., Sharma, R., and Westra, W.H. (2011). PAX8 immunostaining of  
591 anaplastic thyroid carcinoma: a reliable means of discerning thyroid origin for  
592 undifferentiated tumors of the head and neck. *Hum. Pathol.* *42*, 1873–1877.  
593 <https://doi.org/10.1016/j.humpath.2011.02.004>.
- 594 36. Suzuki, A., Hirokawa, M., Takada, N., Higuchi, M., Yamao, N., Kuma, S., Daa,  
595 T., and Miyauchi, A. (2015). Diagnostic significance of PAX8 in thyroid  
596 squamous cell carcinoma. *Endocr. J.* *62*, 991–995.  
597 <https://doi.org/10.1507/endocrj.EJ15-0226>.
- 598 37. Ozcan, A., Shen, S.S., Hamilton, C., Anjana, K., Coffey, D., Krishnan, B., and  
599 Truong, L.D. (2011). PAX 8 expression in non-neoplastic tissues, primary  
600 tumors, and metastatic tumors: a comprehensive immunohistochemical study.  
601 *Mod. Pathol.* *24*, 751–764. <https://doi.org/10.1038/modpathol.2011.3>.
- 602  
603

## 604 **Figure Legends**

605 **Figure 1.** Overview of the designed workflow. Variants colored in dark green  
606 represent the variants imputed in the intermediate imputation phase (first imputation),  
607 while variants colored in white are imputed from general 1000G reference panel in  
608 both two-step and conventional imputation.

609 **Figure 2.** Genetic PCA of the included cohorts in the workflow. The samples are  
610 colored by the cohorts with their unique array type. Panel A and B show the first four  
611 genetic components of the conventional imputation approach (variance explained for  
612 the components combined is 0.34). Panels C and D show the first four components  
613 of the proposed approach (variance explained for the components combined is 0.17).

614 **Figure 3.** Median R<sup>2</sup> of the imputation outcomes of genotype data of the included  
615 cohorts, using conventional (dotted) and two-step imputation (full line).

616 **Figure 4.** Boxplots grid of the absolute difference in allele frequency (AF) between  
617 conventional and two-step imputation outcomes for the included cohorts. Each  
618 column represents an allele frequency group (rare, low, common) and each row  
619 represents one of the included cohorts.

620 **Figure 5.** Manhattan plot of the GWAS analysis of goiter risk using combined two-  
621 step imputation genotypes (A) and conventional imputation and meta-analysis  
622 approach (B). Variants are plotted on the x axis and  $-\log_{10}$  p-values of the  
623 association testing on the y axis. Associations significant after correction for multiple  
624 testing ( $p < 5 \times 10^{-8}$ ) are colored in red. Regional association results and  
625 recombination rates for the *XKR6* gene from two-step imputation GWAS are  
626 presented in part C,  $-\log_{10}$  p-values (y-axis) of the single nucleotide variants  
627 according to their chromosomal positions (x-axis) with lead variant (rs7005680) is  
628 shown as a purple diamond.

629 **Figure 6.** Manhattan plot of the GWAS analysis of log thyroid volume using  
630 combined two-step imputation genotypes (A) and conventional imputation and meta-  
631 analysis approach (B). Variants are plotted on the x axis and  $-\log_{10}$  p-values of the  
632 association testing on the y axis. Associations significant after correction for multiple  
633 testing ( $p < 5 \times 10^{-8}$ ) are colored in red. Regional association results and  
634 recombination rates for the *NRG1* gene from conventional imputation GWAS are  
635 presented in part C,  $-\log_{10}$  p-values (y-axis) of the single nucleotide variants

636 according to their chromosomal positions (x-axis) with lead variant (rs7000397) is

637 shown as a purple diamond.

638

639 **Tables**

640 **Table 1.** Genotype concordance (number of matching genotypes/total number of  
 641 genotypes) of the hard call imputed genotypes with sequenced data for homozygous  
 642 reference (HomRef), homozygous alternative (HomAlt), and heterozygous (Het) calls,  
 643 all variants followed by stratification by minor allele frequency (MAF) group of the  
 644 imputed data. (n) represents number of variants represented per group in both  
 645 imputed and sequenced genotypes.  
 646

Group	All variants		MAF 0.05-0.5		MAF 0.01-0.05		MAF 0.0-0.01	
	Conventional	Two-step	Conventional	Two-step	Conventional	Two-step	Conventional	Two-step
<b>HomRef</b>	0.9963 (618,453,288)	0.9962 (613,997,680)	0.9913 (224,034,629)	0.9910 (222,769,137)	0.9959 (171,621,821)	0.9988 (183,771,150)	0.9993 (222,756,838)	0.9994 (207,703,823)
<b>Het</b>	0.9814 (116,274,486)	0.981 (116,144,961)	0.9897 (105,647,496)	0.9893 (105,617,755)	0.9424 (8,530,966)	0.9418 (8,903,479)	0.7875 (2,139,474)	0.7525 (1,695,994)
<b>HomAlt</b>	0.9783 (21,587,792)	0.9787 (21,408,323)	0.9802 (21,463,336)	0.9807 (21,290,637)	0.7979 (128,442)	0.7862 (124,359)	0.2512 (5,149)	0.1389 (2,369)

647

648 **Table 2.** Cohort characteristics for GWAS analysis of thyroid volume and goiter risk.  
 649 BSA: body surface area.

Cohort	N	Age (mean ± SD)	Sex	BSA (mean ± SD)	Current smoker	Log thyroid volume (mean ± SD)	Goiter cases(%)
SHIP START	3,611	49.1 (16.3)	47.4% Male, 52.6% Female	1.9 (0.2)	31.6%	2.9 (0.5)	1322 (36.6%)
SHIP-TREND	784	49.2 (13.9)	51.5% Male, 48.5% Female	1.9 (0.2)	22.2%	2.9 (0.4)	253 (32.3%)
SHIP-TREND batch II	2,499	51.3 (16)	56.9% Male, 43.1% Female	1.9 (0.2)	29.8%	2.9 (0.4)	737 (29.5%)

650

651

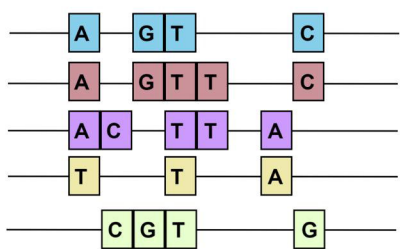
652 **Table 3.** Loci lead SNVs with the strongest association with thyroid volume (A), and  
 653 goiter risk (B) in GWAS analysis using both two-step imputed and conventionally  
 654 imputed genotypes.  
 655 Novel loci are marked in bold. A1: coded allele, AF: frequency of coded allele, SE:  
 656 standard errors, R2: linkage disequilibrium in the top SNVs of the same locus in both  
 657 imputation outcomes  
 658

A) Volume		Two-step imputation GWAS						Conventional imputation GWAS meta-analysis							
Locus	Lead SNP	CHR	A1	AF	Effect	SE	P-value	Lead SNP	CHR	A1	AF	Effect	SE	P-value	R2
CAPZB	rs4911994	1	A	0.64	0.057	0.01	1.63E-18	rs4911994	1	A	0.63	0.055	0.01	5.93E-19	Same SNP
	rs10799824	1	A	0.15	0.086	0.01	1.82E-22	rs12410532	1	T	0.14	0.092	0.01	6.14E-24	0.978
<b>PAX8</b>	rs7560701	2	C	0.48	-0.036	0.01	1.12E-08	rs1110839	2	T	0.48	0.038	0.01	6.35E-10	0.576
<b>IGFBP5</b>	rs2712172	2	A	0.27	0.041	0.01	1.46E-08	rs2712172	2	A	0.26	0.039	0.01	1.95E-08	Same SNP
<b>NRG1</b>	rs7000397	8	G	0.35	0.035	0.01	5.22E-08	rs7000397	8	G	0.34	0.04	0.01	1.42E-08	Same SNP
<b>TG</b>	rs114322847	8	T	0.03	0.122	0.02	1.90E-09	rs79676842	8	T	0.03	0.135	0.02	2.06E-11	1
<b>FAM227B</b>	rs17477923	15	C	0.26	0.054	0.01	9.55E-14	rs73398264	15	T	0.75	-0.052	0.01	7.73E-14	1
<b>MAFTRR</b>	rs3813579	16	A	0.53	0.056	0.01	2.57E-19	rs562609617	16	G	0.66	-0.064	0.01	1.85E-22	0.413
B) Goiter		Two-step imputation GWAS						Conventional imputation GWAS meta-analysis							
Locus	Lead SNP	CHR	A1	AF	Effect	SE	P-value	Lead SNP	CHR	A1	AF	Effect	SE	P-value	R2
CAPZB	rs4911994	1	A	0.64	0.614	0.68	8.89E-14	rs4911994	1	A	0.62	0.302	0.04	1.13E-13	Same SNP
	rs10799824	1	A	0.15	0.135	0.19	6.68E-16	rs12410532	1	T	0.14	0.450	0.06	1.25E-15	0.978
<b>XKR6</b>	rs7005680	8	T	0.35	-0.218	0.04	7.98E-09	rs11778398	8	T	0.46	-0.200	0.04	4.35E-07	0.761
<b>FAM227B</b>	rs75929244	15	T	0.24	0.321	0.04	2.17E-13	rs73398264	15	T	0.75	-0.313	0.04	1.61E-12	0.897
<b>MAFTRR</b>	rs3813579	16	GT	0.35	0.29	0.04	4.44E-13	rs562609617	16	G	0.34	0.064	0.01	1.85E-22	Same SNP

659

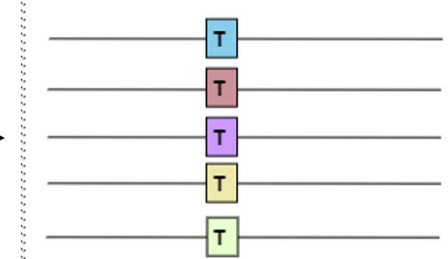


### Preprocessed genotyped data

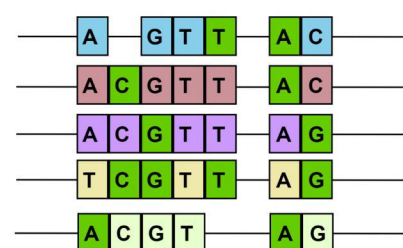


Cohort	N. Samples	Array Type	N. SNPs on the array
GANI_MED	2410	Illumina PsychArray	570.000
SHIP START	4070	Affymetrix SNP 6.0	900.000
SHIP START II	48	Affymetrix Axiom	560.000
SHIP-TREND	986	Illumina Omni 2.5	2.500.000
SHIP-TREND batch II	3133	Illumina GSA	655.000

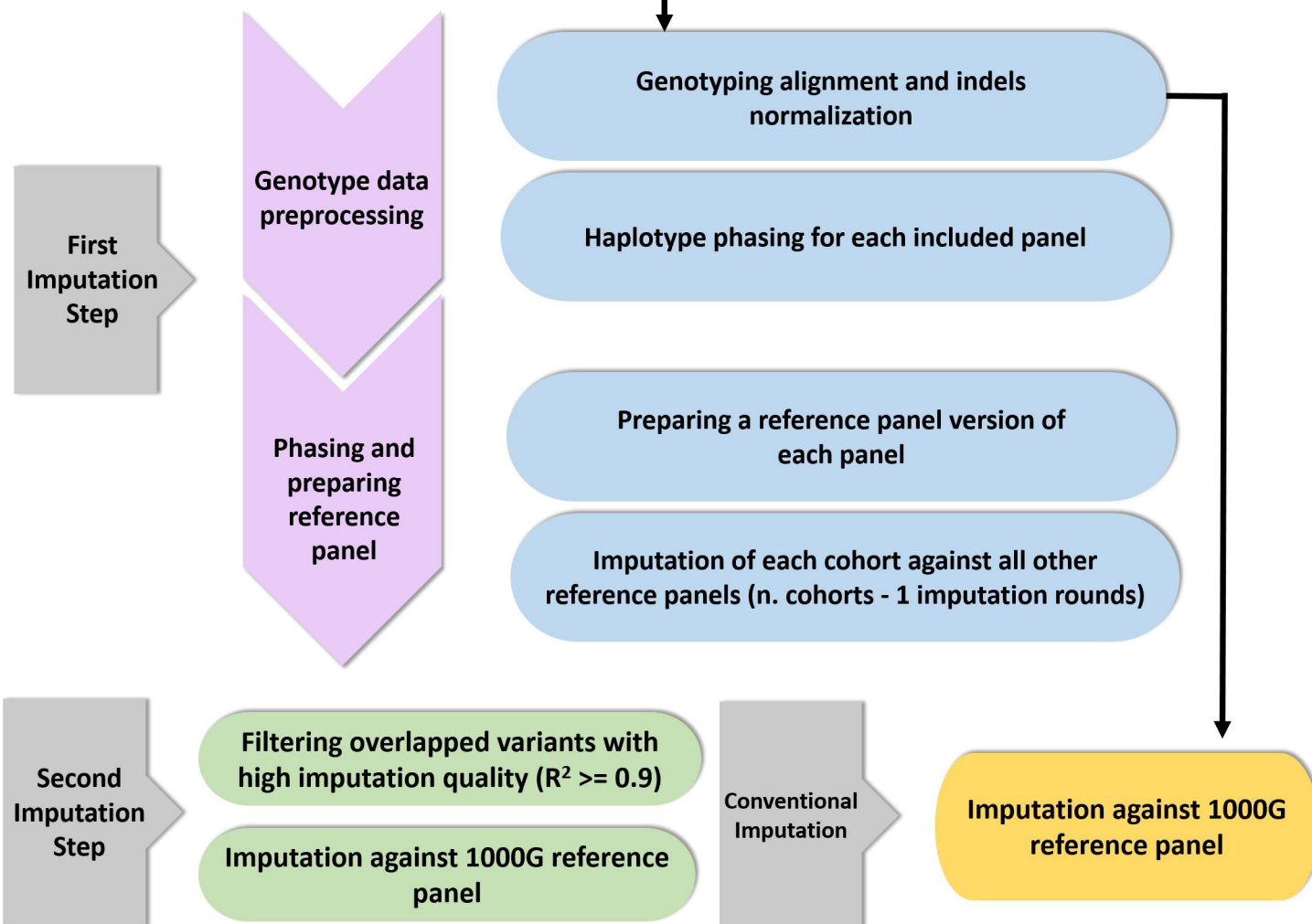
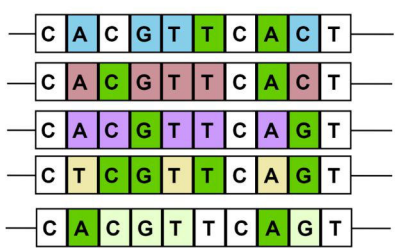
### Overlapped genotypes only



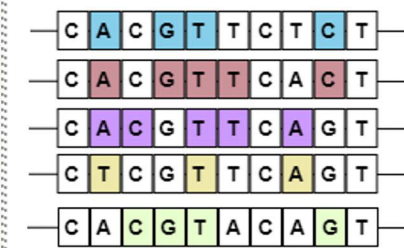
### First (intermediate) imputation



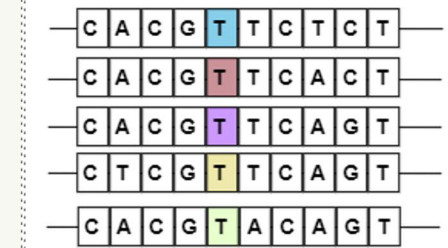
### Second imputation outcome



### Conventional imputation outcome



### Imputation outcome



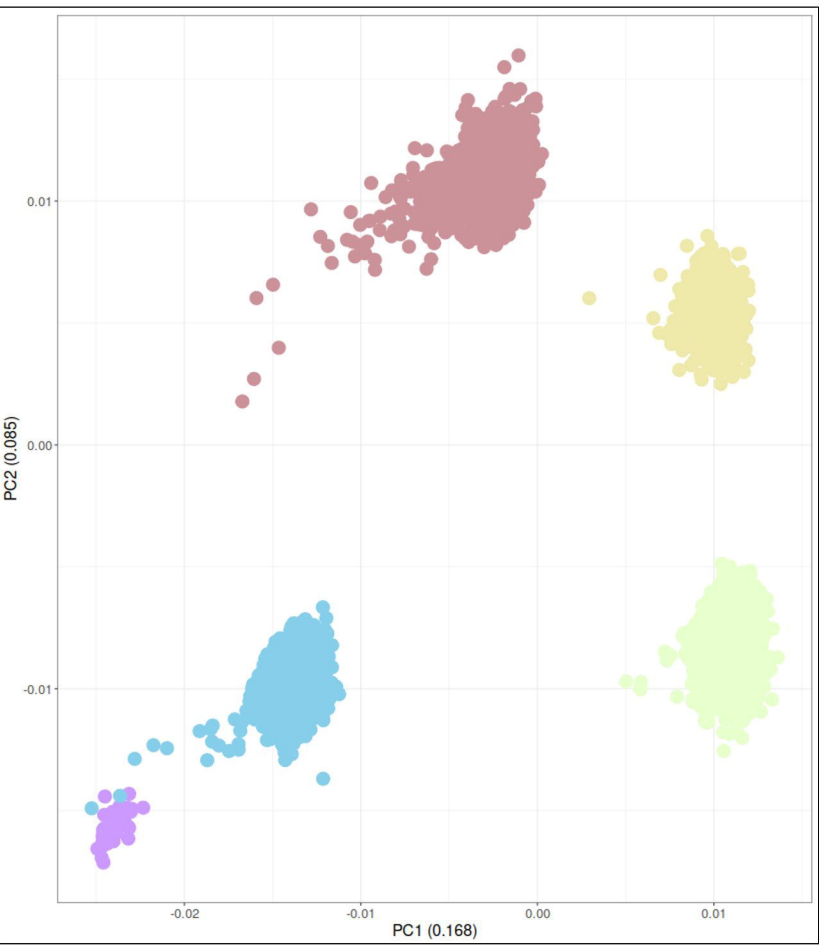
### Imputation Outcome Evaluation

**A** Batch effect (PCA)

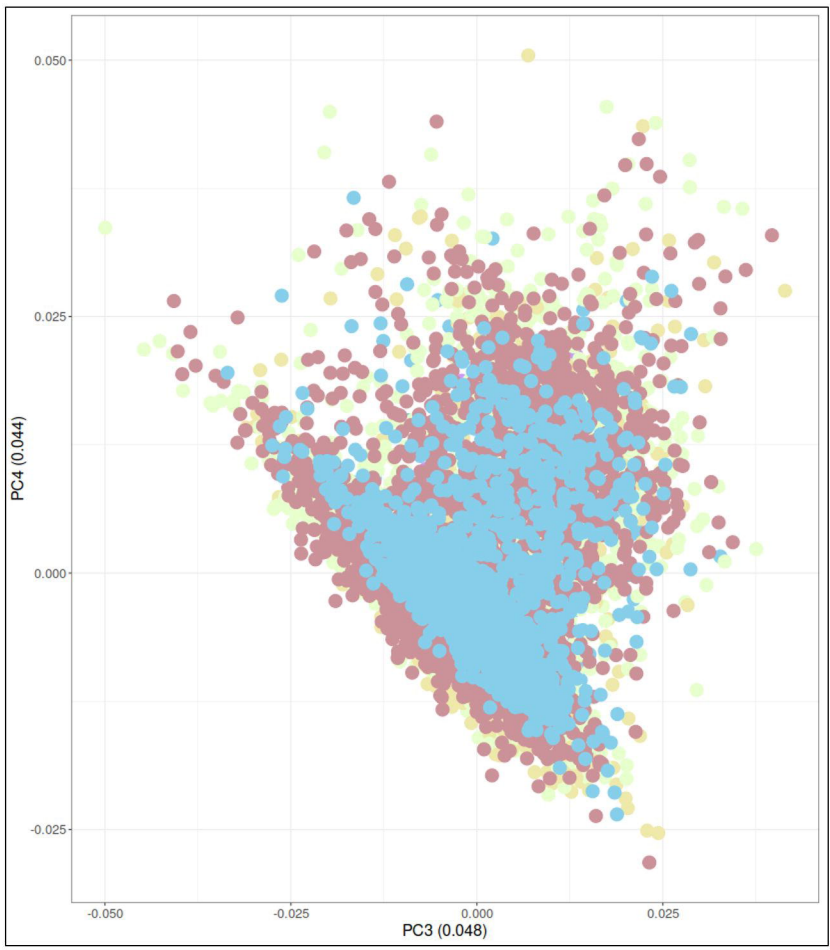
**B** Comparative analysis of variants information

**C** GWAS results comparison (Thyroid function)

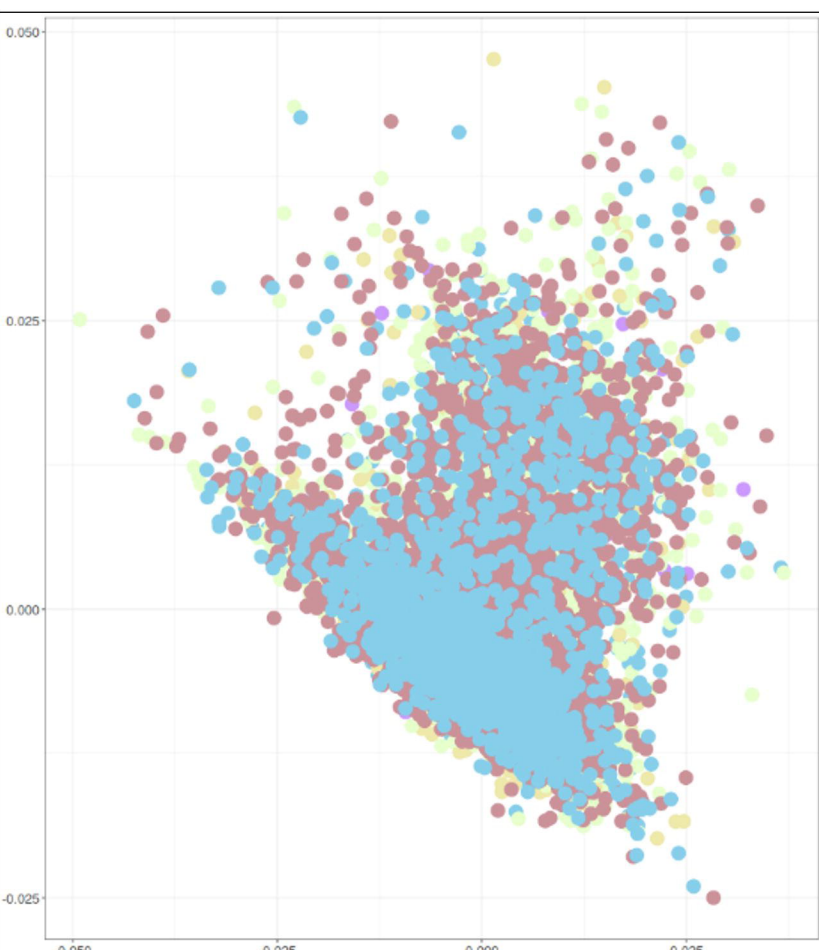
A)



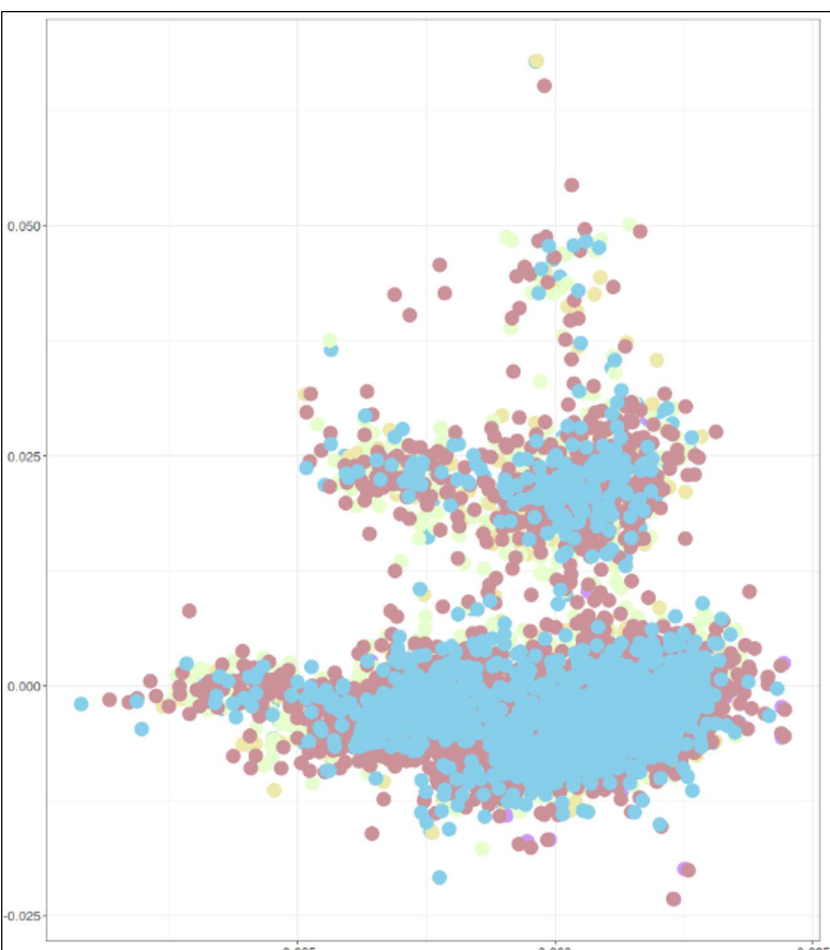
B)



C)

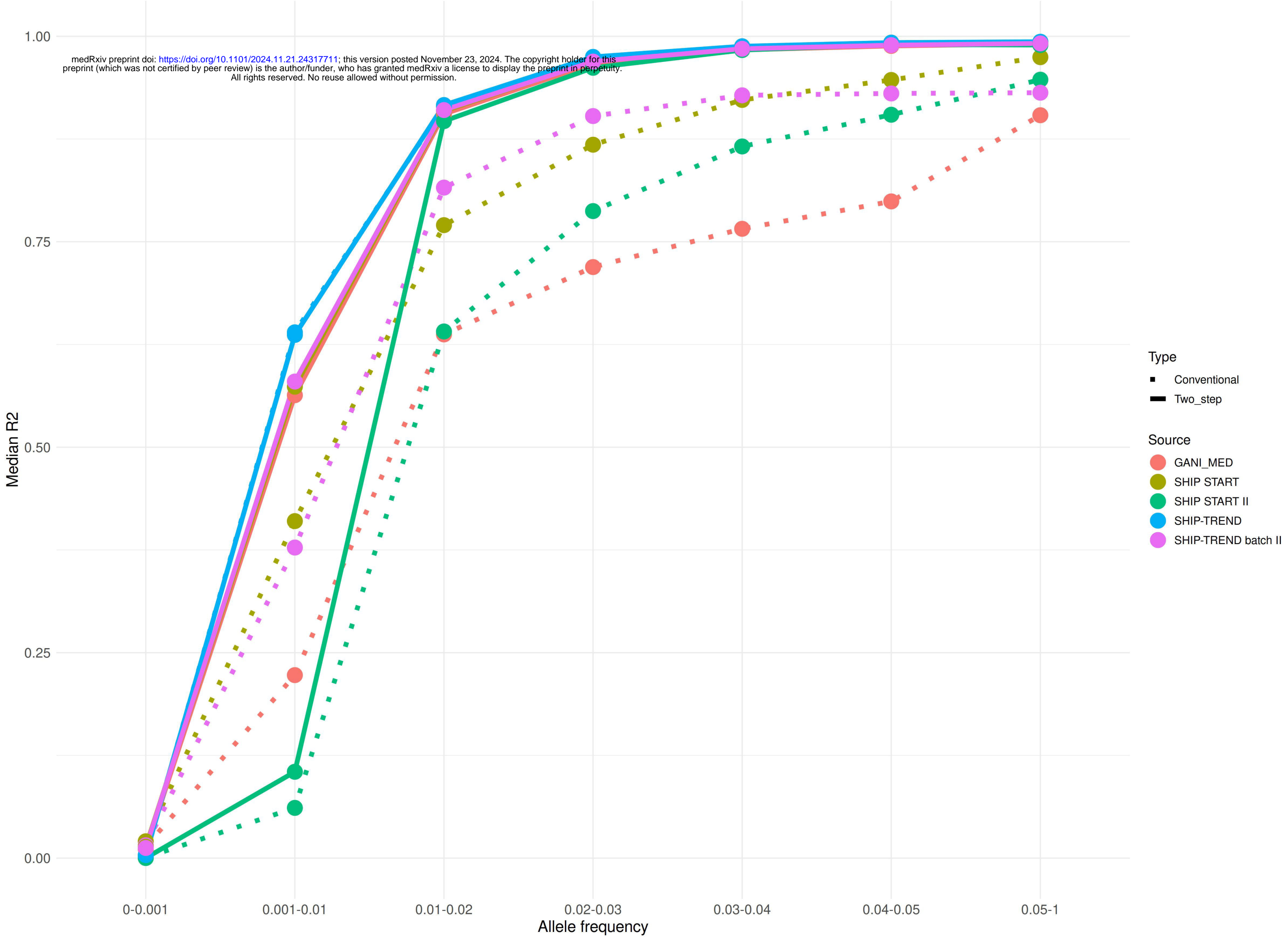


D)



# Median R2 with Quantiles by Type

medRxiv preprint doi: <https://doi.org/10.1101/2024.11.21.24317711>; this version posted November 23, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. All rights reserved. No reuse allowed without permission.

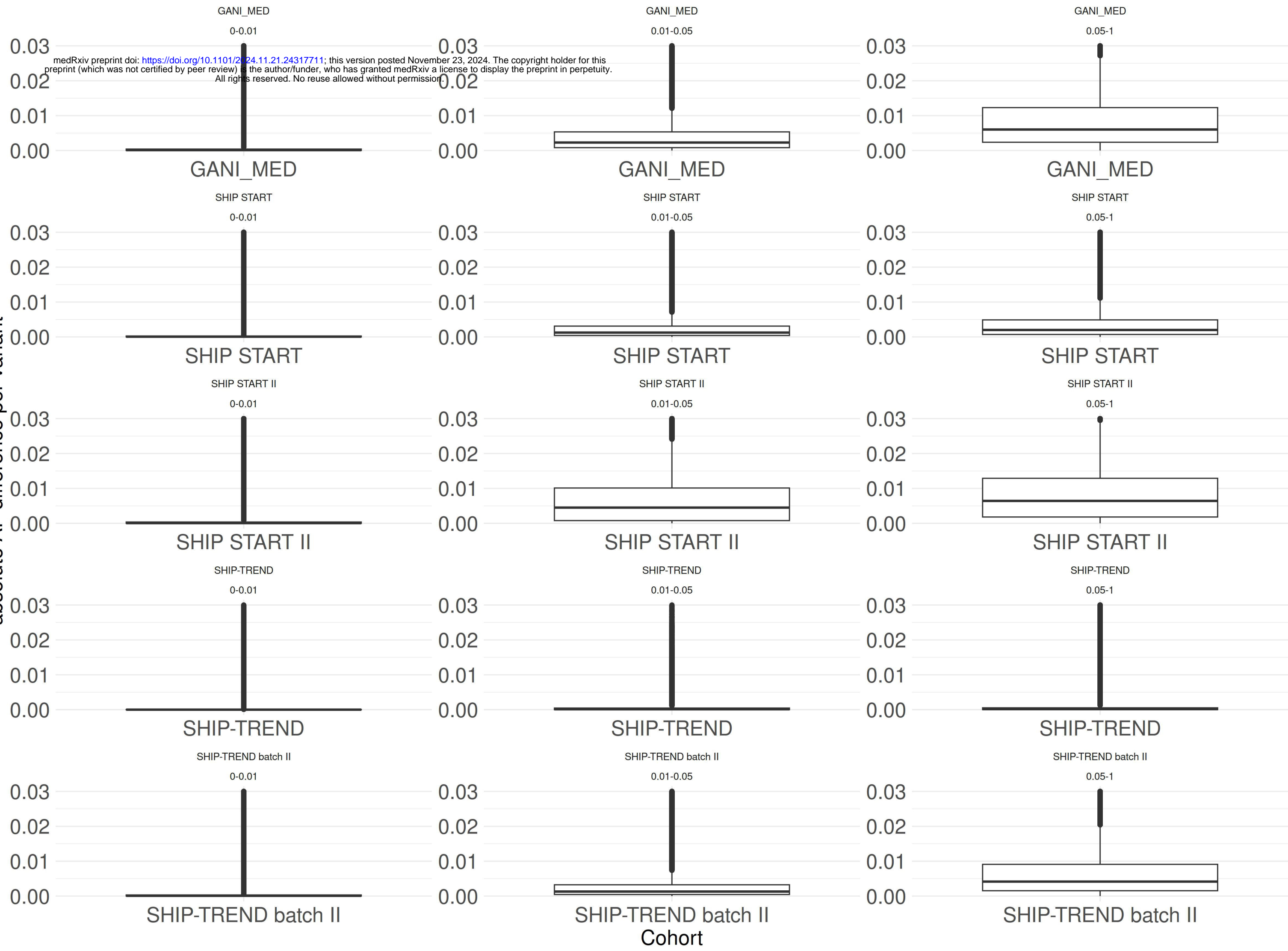




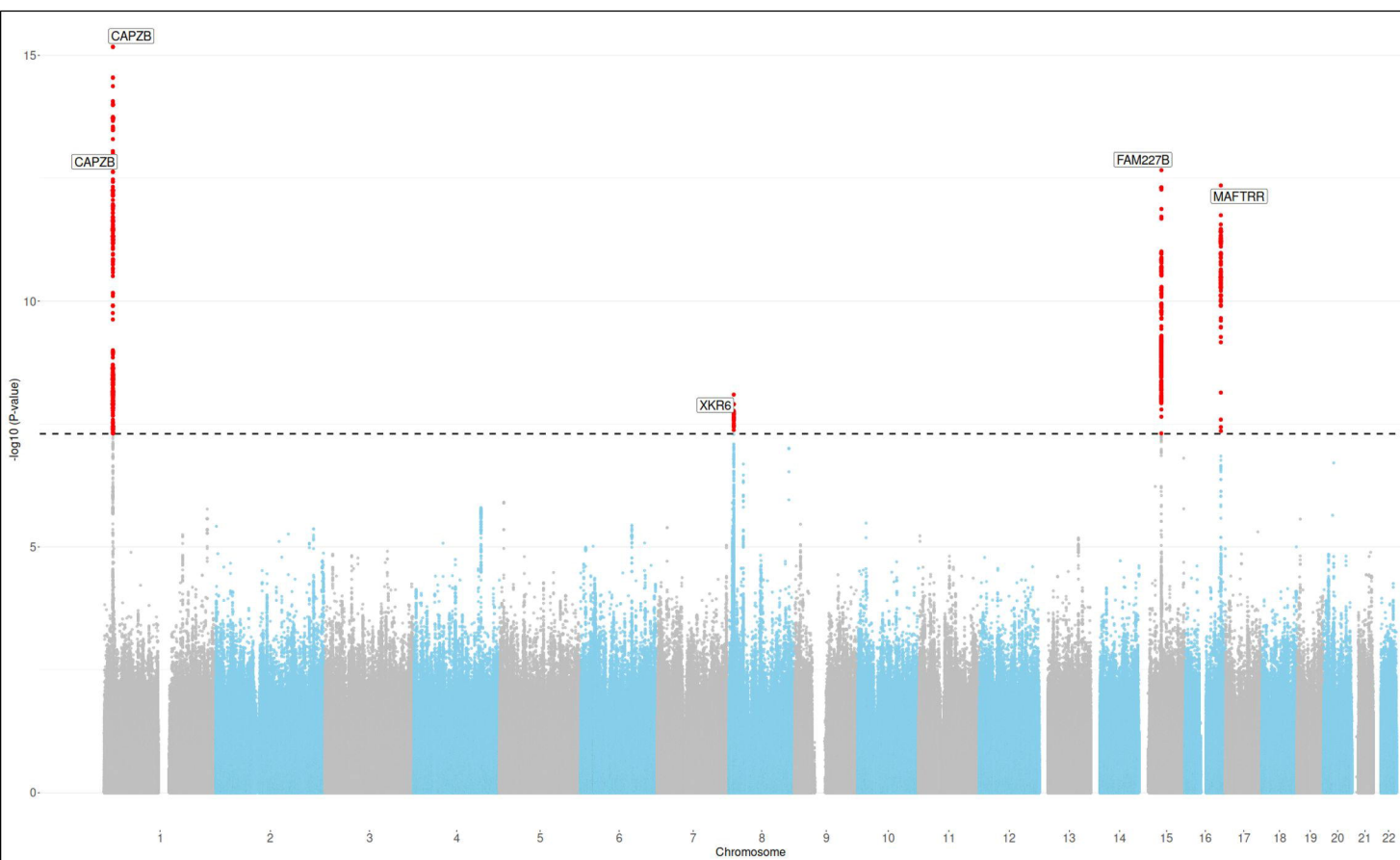
# boxplot Plot of the absolute difference in AF between conventional and two-step imputation outcomes for the included cohorts

medRxiv preprint doi: <https://doi.org/10.1101/2024.11.21.24317711>; this version posted November 23, 2024. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. All rights reserved. No reuse allowed without permission.

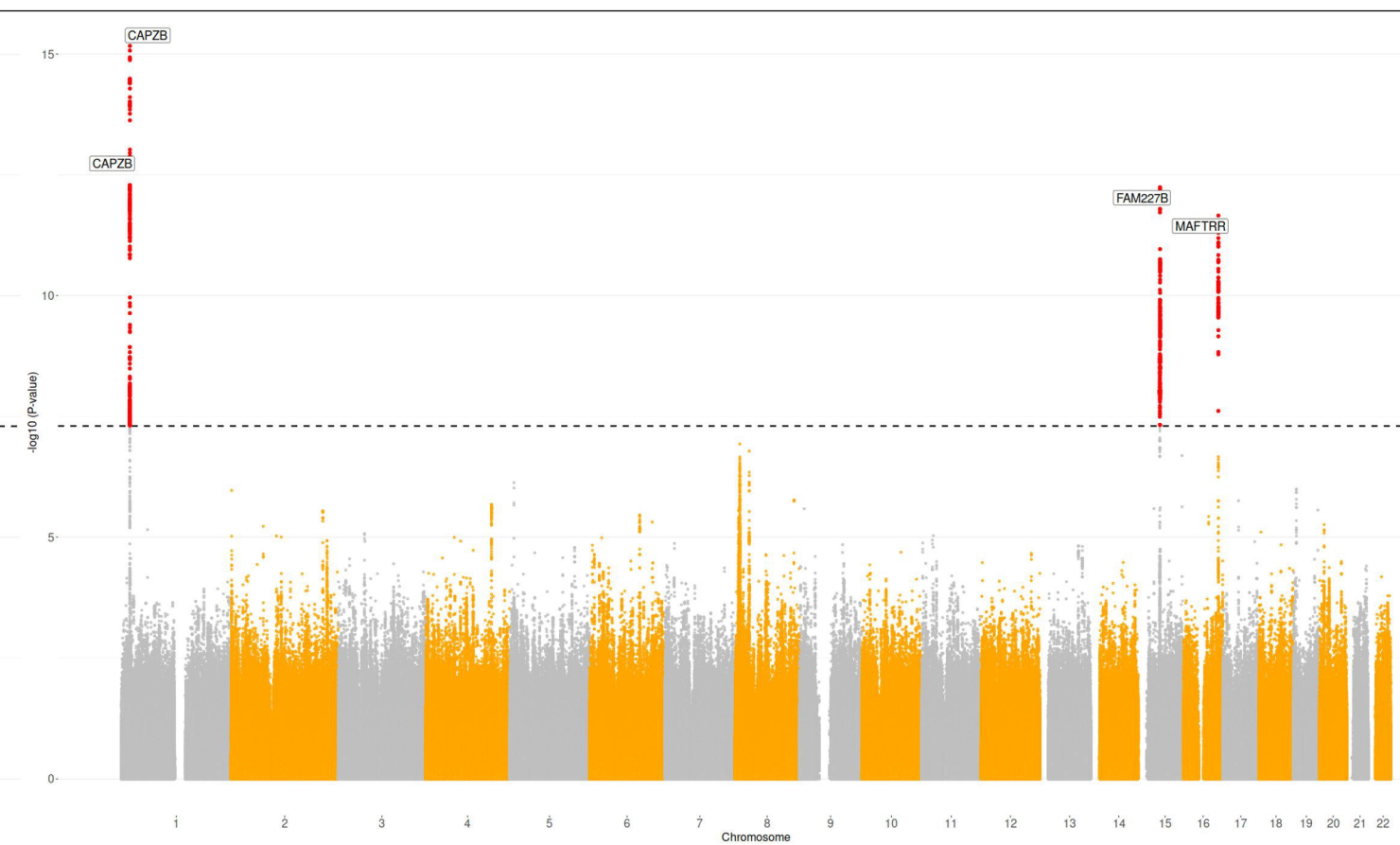
absolute AF difference per variant



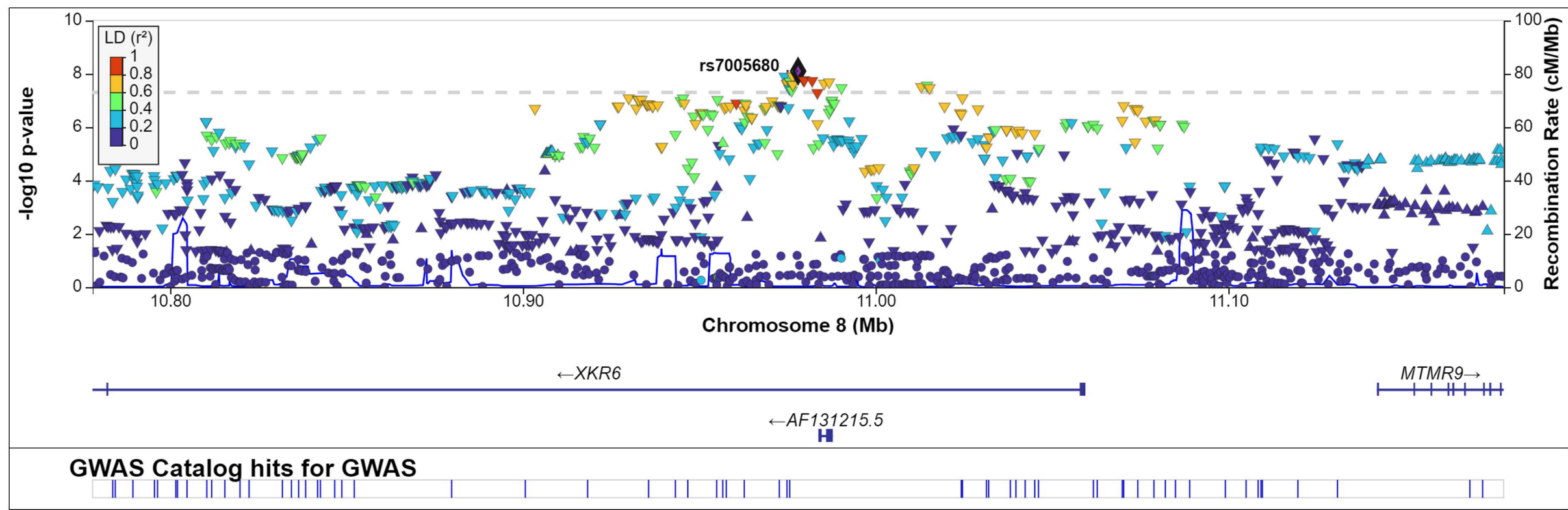
A)



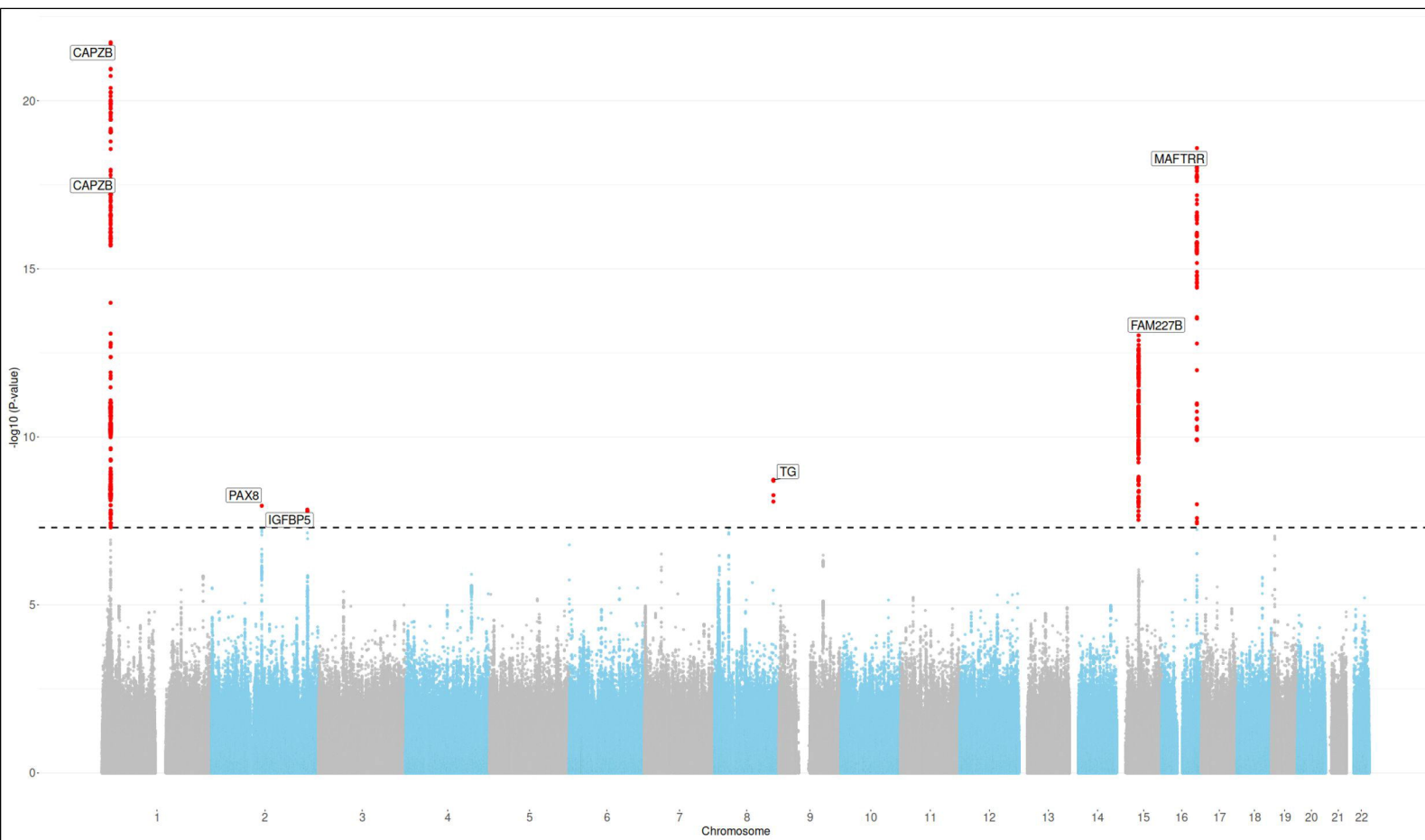
B)



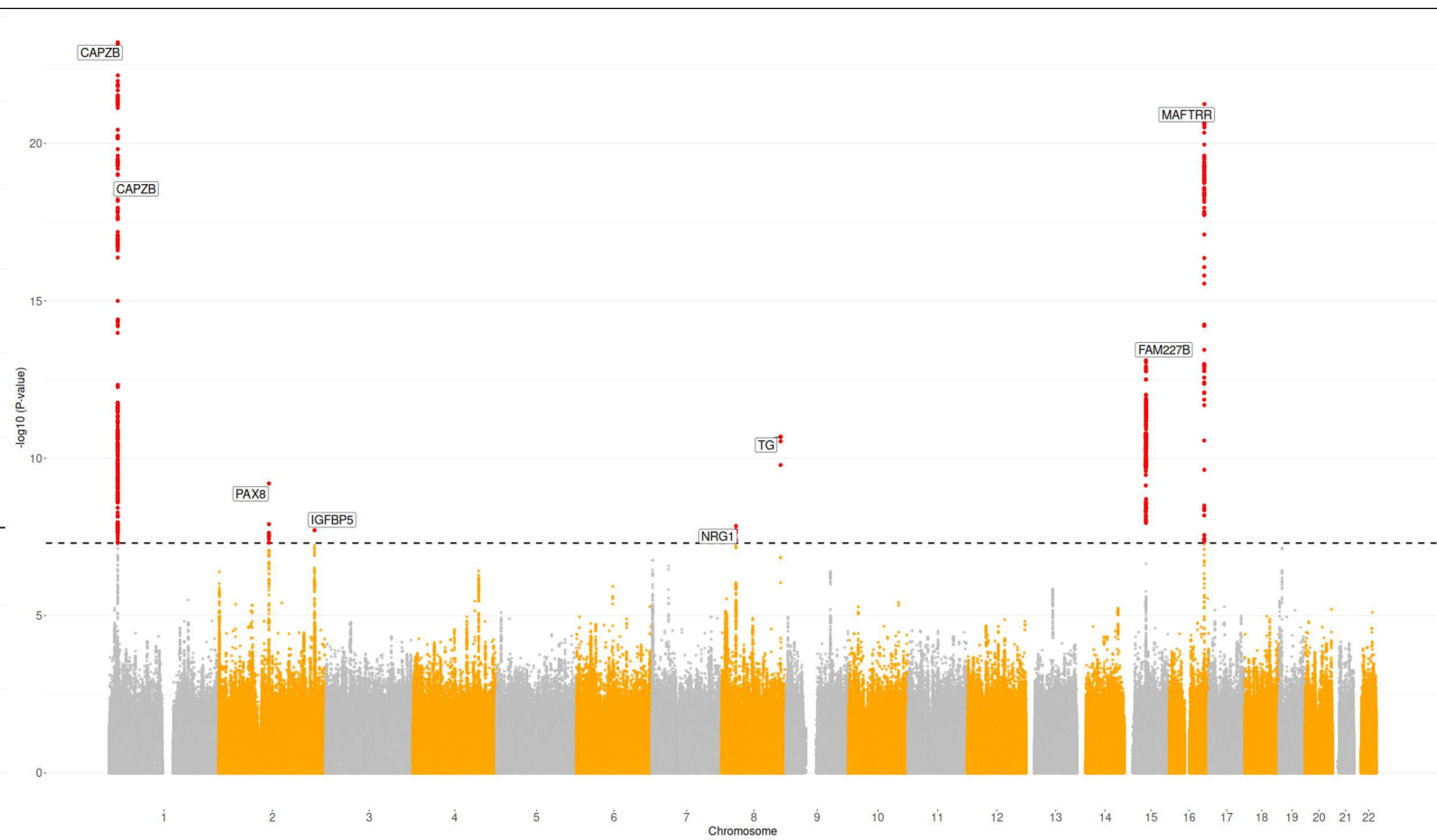
C)



A)



B)



C)

