



Genotype imputation from low-coverage data for medical and population genetic analyses

Simone Andrea Biagini, Sara Becelaere, Mio Aerden, et al.

Genome Res. published online July 22, 2025

Access the most recent version at doi:[10.1101/gr.280175.124](https://doi.org/10.1101/gr.280175.124)

| | |
|---------------------------------|--|
| P<P | Published online July 22, 2025 in advance of the print journal. |
| Accepted Manuscript | Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version. |
| Open Access | Freely available online through the <i>Genome Research</i> Open Access option. |
| Creative Commons License | This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution 4.0 International license), as described at http://creativecommons.org/licenses/by/4.0/ . |
| Email Alerting Service | Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here . |

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 Genotype imputation from low-coverage data for medical and 2 population genetic analyses

3 Simone Andrea Biagini,^{1,2,3} Sara Becelaere,¹ Mio Aerden^{1,4}, Tatjana Jatsenko¹, Laurens Hannes,^{1,4} Philip
4 Van Damme,^{5,6} Jeroen Breckpot,^{1,4} Koenraad Devriendt,^{1,4} Bernard Thienpont¹, Joris Robert Vermeesch,¹
5 Isabelle Cleynen,¹ Toomas Kivisild,^{1,7}

6 1 - Department of Human Genetics, KU Leuven, Leuven 3000, Belgium, 2 - Department of Archaeology and
7 Museology, Masaryk University, Brno, Czech Republic, 3 - Center of Molecular Medicine, Central European
8 Institute of Technology, Masaryk University, Brno, Czech Republic, 4 - Center for Human Genetics,
9 University Hospitals Leuven, University of Leuven, Leuven, Belgium, 5 - Laboratory of Neurobiology,
10 Neuroscience Department, KU Leuven, Leuven 3000, Belgium, 6 - Neurology Department, University
11 Hospitals Leuven, Leuven 3000, Belgium, 7 - Estonian Biocentre, Institute of Genomics, University of Tartu,
12 Tartu 51010, Estonia

13

14 **Corresponding Authors:** Prof. Toomas Kivisild, toomas.kivisild@kuleuven.be

15 **Running title:** Genotype Imputation from Low-Coverage Data

16

17 **ABSTRACT**

18 Genotype imputation from low-pass sequencing data presents unique opportunities for genomic
19 analyses but comes with specific challenges. In this study, we explore the impact of quality filters
20 on genetic ancestry and Polygenic Score (PGS) estimation after imputing 32,769 low-pass genome
21 wide sequences (LPS) from non-invasive prenatal screening (NIPS) with an average autosomal
22 sequence depth of $\sim 0.15\times$. In studies involving ultra-low coverage sequences, conventional
23 approaches to secure genotype accuracy may fail, especially when multiple samples are pooled.
24 To enhance the proportion of high-quality genotypes in large datasets we introduce a filtering
25 approach called GDI that combines genotype probability (GP), alternate allele dosage (DS), and
26 INFO score filters. We demonstrate that the imputation tools QUILT and GLIMPSE2 achieve
27 similar accuracy, which is high enough for broad-scale ancestry mapping but insufficient for high
28 resolution Principal Component Analysis (PCA), when applied without filters. With the GDI
29 approach we can achieve quality that is adequate for such purposes. Furthermore, we explored the
30 impact of imputation errors, choice of variants and filtering methods on PGS prediction for height in

31 1,911 subjects with height data. We show that polygenic scores predict 23.7% of variance in height
32 in our imputed data and that, contrary to the effect on PCA, the GDI filter does not improve the
33 performance of PGS in height prediction. These results highlight that imputed LPS data can be
34 leveraged for further biomedical and population genetic use but there is a need to consider each
35 downstream analysis tool individually for its imputation quality thresholds and filtering
36 requirements.

37

38 **INTRODUCTION**

39

40 Until recently, genotyping arrays were the only cost-effective approach for generating data for
41 medical genomics at the scale of tens of thousands to millions of individuals. With decreasing costs
42 of genome sequencing, low-pass (at coverage $<1\times$) whole-genome sequencing followed by
43 imputation has become a cost-effective alternative for trait mapping and estimation of polygenic
44 scores (PGS) (Martin et al. 2021; Wasik et al. 2021). For evolutionary and population genetic
45 studies, sequencing of large numbers of individuals at low depth provides a more comprehensive
46 representation of variation in the population than sequencing of a small number of individuals at a
47 higher depth (Fumagalli 2013). Importantly, imputation performed at sufficient accuracy can offer
48 opportunities to reuse large volumes of already generated data. Here, we explore the prospects to
49 leverage hundreds of thousands of LPS genomes generated from cell-free DNA (cfDNA) during
50 Non-Invasive Prenatal Screening (NIPS), (Wang et al. 2021; Van Riel, Stanley, and Vermeesch
51 2023), for the study of maternal genotypes via imputation.

52

53 Non-Invasive Prenatal Screening (NIPS) has become an integral part of prenatal care in many
54 countries, though adoption rates vary significantly across Europe. In most European countries,
55 fewer than 25% of individuals adopt NIPS, and in many cases, the rate is below 5% (Gadsbøll et
56 al. 2020). However, countries such as Italy, Spain, Austria, the Netherlands, and Belgium
57 demonstrate higher uptake. Belgium, in particular, stands out due to its progressive reimbursement
58 policy for NIPS, which has facilitated its widespread accessibility and seamless integration into
59 prenatal care (Bayindir et al., 2015; Van Den Bogaert et al., 2022). Approximately 80% of NIPS

60 performed in Belgium rely on low-pass genome-wide sequencing, making this technology a
61 cornerstone of prenatal care in the country. More than 200,000 LPS genomes have been
62 generated so far for NIPS in Belgium alone.

63

64 Imputation of maternal genotypes from NIPS data can offer unique genomic medicine opportunities
65 beyond fetal aneuploidy screening at the population level, such as genome-wide association
66 studies (GWAS) and PGS profiling of large cohorts. In this study, we focus on imputation strategies
67 scalable to large cohorts that ensure sufficient accuracy for downstream purposes by exploring
68 suitable post-imputation filtering strategies. In our study, we demonstrate how LPS data from NIPS
69 can be repurposed to support diverse biological analyses, showcasing the broader potential of
70 cfDNA-derived LPS data.

71

72 Several tools enable efficient and accurate genotype imputation from LPS data. Among these,
73 Beagle (Browning, Zhou, and Browning 2018) and Gencove's loimpute software (Wasik et al.,
74 2019) have been some of the most popular examples. However, they present limitations in
75 scalability for large cohorts and increased computational costs with large reference panels.
76 Recently, imputation methods like GLIMPSE (Rubinacci et al. 2021), GLIMPSE2 (Rubinacci et al.
77 2023), and QUILT (Davies et al. 2021), scalable for imputation tasks of tens of thousands of low-
78 coverage genomes, have been developed and benchmarked against other tools, showing high
79 accuracy and low computation costs with reference panels of tens of thousands individuals or
80 more.

81

82 The minimum coverage from which accurate genotype imputation can be achieved may depend on
83 various factors, including the sample quality (read length, contamination, damage), variant
84 frequency and representation in the reference panel, composition and size of reference panels, as
85 well as the error sensitivity of downstream analytical methods used. Gamba et al. (2014) showed

86 that imputation of common variants from European ancient genomes of 1× coverage can be
87 performed at 99% accuracy when filtering out a small proportion of variants that have low genotype
88 probability ($GP < 0.99$). Accuracy higher than 90% in down-sampled ancient genomes was
89 achievable also from 0.1× coverage with the cost of filtering out more variants (Hui et al. 2020).
90 Using the 1000 Genomes Project reference panel and GLIMPSE, Sousa Da Mota et al. (2023)
91 showed that the $GP \geq 0.99$ filter retained only 20-43% of the correctly imputed variants per
92 individual sample in low coverage samples. Individual-level filtering will thus produce, even in small
93 sample sizes, substantial cumulative missingness leading to a drastic reduction in the number of
94 sites available for downstream analyses. In search of the optimal balance between data loss and
95 quality with large batches of data, we explore different filtering options and quality metrics of
96 individuals and variants at the batch level. To do so, we implemented a filtering strategy called GDI
97 that incorporates filters on different post-imputation metrics: the posterior genotype probability, the
98 alternate allele dosage, and the INFO score, a value to measure the imputation accuracy of each
99 imputed variant that some tools for imputation usually emit (Marchini & Howie, 2010).

100

101 Genotypes imputed without reference panels from NIPS sequence data from large cohorts of
102 Chinese individuals have been used for the study of population structure in China, for replicating
103 and finding new genome-wide associations with height, BMI, twinning, and gestational diabetes,
104 and for calculation of polygenic scores (Homburger et al. 2019; Liu et al. 2018; Liu et al. 2023;
105 Zhen et al. 2024). Similarly, a GWAS approach applied on imputed LPS data of 140,000 NIPS
106 profiles from the Netherlands has been successful in identifying loci that affect the concentration
107 and fragmentation properties of cell-free DNA in plasma (Linthorst et al. 2024). However, the full
108 potential of this approach is yet to be revealed. Because of the availability of large haplotype

109 reference panels, such as HRC (the Haplotype Reference Consortium. 2016) with 27,165
110 individuals, imputation of European samples can be achieved with increased accuracy (Davies et
111 al. 2021; Rubinacci et al. 2021; Rubinacci et al. 2023). Yet, the impact of imputation errors,
112 particularly in the low minor allele frequency classes, on downstream population genetic
113 applications, such as PCA, or medical applications, such as PGS analysis, is not fully known. A
114 PGS estimates an individual's genetic risk to develop a certain complex trait or disease by
115 combining trait-associated variants into one predictive score. Adult height is a complex trait that is
116 highly polygenic in nature, with 12,111 independent Single Nucleotide Polymorphisms (SNPs)
117 collectively accounting for 40-45% of phenotypic variance in populations of European ancestry
118 (Yengo et al. 2022). Its high heritability, together with the fact that height is easily measured,
119 makes it an attractive model trait to evaluate the usability of LPS data with different imputation
120 options for genome-wide association scanning and PGS calculation. Here, we assess the accuracy
121 of PGS calculation from low-coverage imputed NIPS data for height and investigate the effect of
122 different filtering strategies, including the GDI filter on the correlation between predicted and
123 measured trait values.

124

125 In this study, we test the performance of QUILT (Davies et al. 2021) in genotype imputation on
126 large batches (>1,000 individuals) of LPS data (also referred to as NIPS data in the text) with an
127 average coverage of ~0.15x. As small insertions and deletions can be imputed with lower accuracy
128 (Nguyen et al. 2024), require separate handling (Rubinacci et al. 2023), and are omitted by QUILT
129 (Davies et al. 2021), we focus our analyses only on single nucleotide substitutions. We impute the
130 genotypes of 32,769 Belgian individuals and test the effect of different reference panels and
131 filtering steps on imputation accuracy of a single individual or a batch of many individuals. Using
132 the GDI approach we evaluate the effects of the imputation quality filters on downstream analyses,
133 including Principal Component Analysis (PCA) and PGS.

134

135 **RESULTS**

136

137 **Imputation accuracy from LPS data with QUILT and GLIMPSE2**

138

139 We first evaluated the performance of QUILT (Davies et al. 2021) on LPS data through sensitivity
140 values across multiple MAF bins for sites imputed as homozygous for the reference (REF) or the
141 alternate allele (ALT), and for heterozygous sites (HET) that tend to be more challenging to impute
142 (Ausmees et al. 2022). Additionally, we also assessed data filtered for maximum genotype
143 probability ($\max(\text{GP}) \geq 0.99$), a commonly employed approach to enhance the quality of imputed
144 data (Gamba et al. 2014; Hui et al. 2020), especially effective when working with few or individual
145 samples. We observed (Supplementary Figure 1) that the average sensitivity of raw imputed data
146 across three test samples varies among different MAF bins, with clear improvements after applying
147 a GP filter.

148

149 As detailed in Supplementary Table 1, the HET sensitivity shows an average improvement of 9.3%
150 across all observed MAF bins after GP filtering that allows sensitivity values to surpass ~97% for
151 common variants ($\text{MAF} > 0.05$). Similarly, the ALT sensitivity, which already exhibits values
152 exceeding ~93% across all MAF bins in raw imputed results, demonstrates an average
153 improvement of 3.33%, reaching values higher than ~97% across all MAF bins post GP filter
154 application. Moreover, the REF sensitivity presents an average improvement of 1.8%, with values
155 surpassing ~98% across all MAF bins following GP filter application.

156 To provide a comparison with other available imputation tools, we also imputed the three test
157 samples with GLIMPSE (Rubinacci et al. 2021) and GLIMPSE2 (Rubinacci et al. 2023). The
158 comparison between QUILT and GLIMPSE (Figure S2, panel A) highlighted that QUILT provides a

159 higher number of correctly imputed genotypes (especially for heterozygous genotypes) after the
160 application of a GP filter $\max(\text{GP}) \geq 0.99$ (Supplementary Table 2). We observe that sensitivity
161 values show superior performance in QUILT for MAF bins below 5%, while for MAF bins exceeding
162 5%, QUILT maintains an average sensitivity value of 0.9895 ($\text{SD} \pm 0.28$) across the three test
163 samples. This is lower than with GLIMPSE, which presents an average sensitivity value of 0.994
164 ($\text{SD} \pm 0.14$) in the same frequency category ($\text{MAF} > 5\%$). Therefore, while GLIMPSE is, on average,
165 providing an additional 0.0045 sensitivity for the MAF bin including common variants and a lower
166 number of correctly imputed genotypes in any of the observed MAF bins, QUILT presents a better
167 compromise between the number of correctly imputed genotypes and their quality. In the
168 comparison with GLIMPSE2 (Figure S2, panel B), we observe that in each genotype category both
169 tools tend to display similar trends, both for the raw imputed and the GP-filtered data. Furthermore,
170 we also tested a dual-step imputation with Beagle 5 (Browning, Zhou, and Browning 2018) being
171 the second imputation step after applying GLIMPSE or GLIMPSE2. With this test we wanted to see
172 whether using Beagle 5 after running GLIMPSE or GLIMPSE2 could increase the quality of the
173 already imputed data. Overall, we observed that a dual-step imputation (Figure S2 C and D), with a
174 GP filter after the first step, resulted in an increased number of correctly imputed genotypes but at
175 the expense of reduced sensitivity compared to QUILT.

176 Based on these results, we found that among all the options tested, QUILT and GLIMPSE2
177 presented the highest performance levels, showing similar quality outcomes. Additionally, we
178 evaluated the effectiveness of a GP filter in enhancing sensitivity values for various MAF bins and

179 genotype categories at an individual level. However, GLIMPSE2 was not yet available at the time
180 of data production, and QUILT was chosen for the analysis of the bulk of the data. For details on
181 the methodology, please refer to the Materials and Methods section. The pipeline design we
182 employed is presented in Figure 1, panel A.

183 In terms of computational efficiency, a comparative analysis of QUILT and GLIMPSE2 showed that
184 the imputation of one NIPS genome on a single core with QUILT, using a reference panel that had
185 been pre-processed, took 21 hours and 38 minutes, compared to 11 hours and 23 minutes with
186 GLIMPSE2 (a ratio of approximately 1:2). This difference is smaller than the results reported by
187 Rubinacci et al. who observed an approximate 6-fold difference in computation cost in favor of
188 GLIMPSE2 when using the HRC panel (Rubinacci et al. 2023).

189 An additional observation drawn from our tests is that regardless of the imputation tool used, a
190 negative correlation is consistently observed between fetal fraction (SeqFF) and sensitivity (Figure
191 S3A), while a positive correlation exists between coverage and sensitivity (Figure S3B).

192

193 **Selection of the reference panel**

194

195 We compared the effect of different reference panels on imputation accuracy of NIPS data with
196 QUILT, using the three test samples. Specifically, we compared the 1000 Genomes reference
197 panel of 2,504 individuals and 32,140,179 variants, and the HRC panel of 27,165 individuals and
198 36,258,911 variants (Table 1, Table S3). While the sensitivities of heterozygote calls of common
199 variants (MAF >0.05) are similarly high (>0.965) with both panels, we find that the use of the bigger
200 HRC panel enables us to retain more correctly imputed common variants. In total, >65,000
201 heterozygous variants passed the $GP \geq 0.99$ filter with HRC among variants with MAF >0.05, which
202 is ~5% more than with the 1000 Genomes Project panel. In line with this, the total number of

203 missing genotype calls was greater after the GP ≥ 0.99 filters with the 1000 Genomes Project
204 panel.

205

206 **Cumulative** **missingness**

207

208 While the GP filter applied at the individual level can clearly help to enhance imputation accuracy, it
209 introduces cumulative missingness when combining individual data. With large batches, the
210 cumulative missingness can result in retention of only a small number, if any, of the sites that have
211 GP value ≥ 0.99 . Supplementary Figure 4A shows the distribution of variants removed by the GP
212 filter per sample, highlighting the high degree of variance in the batch of imputed NIPS samples,
213 with the majority of the individuals having 20-50% of the variants removed, with an average of
214 approximately 30%. In the full set of 32,769 NIPS individuals we found there are no variants that
215 pass the GP filter in all samples. The minimum proportion of missing calls per variant we observed
216 was 0.007 (210 individuals with GP < 0.99), with the distribution of the missingness by variants per
217 bins of individuals not passing the filter shown in Supplementary Figure 4B.

218 To address the issue of cumulative missingness, we attempted first to filter out variants with an
219 INFO score below 0.4 across all imputed batches, but this proved insufficient in removing the batch
220 effect from high-level PCA, where we observed the clustering together of all our imputed NIPS
221 batches separately from the reference data (Figure S5 A). Increasing the threshold to an INFO
222 score below 0.6 did not improve the result (Figure S5 B). Hence, we adopted a more refined
223 approach, we call GDI filtering, that removes only the most problematic sites by integrating multiple
224 post-imputation quality metrics, including posterior genotype probability, alternate allele dosage,

225 and the INFO score from QUILT (Figure 1, panel B). We applied GDI filtering dynamically on our
226 data to optimize the maximum data quality with minimal loss of variants observing a range of
227 metrics of imputation accuracy in light of the performance of downstream analyses tools, PCA and
228 PGS.

229 **Data filtering**

230

231 *Exclusion of Related and Duplicate Samples*

232 We used QUILT to impute 32,769 samples in batches of a maximum of 1,250 individuals (Figure 1,
233 panel A), as detailed in the Materials and Methods section. The first step after imputation was to
234 identify and remove duplicates and related individuals with the aim of avoiding possible allele
235 frequency biases in the downstream analyses. Duplicates may result from separate tests for
236 different or the same pregnancies. For the identification and removal of duplicates and closer than
237 third degree related individuals, we utilized IBIS (Seidman et al. 2020), which employs the same
238 bounds as KING (Manichaikul et al. 2010) for the degree of relatedness. A kinship coefficient value
239 of 0.088 defines the level of relatedness expected when members of the same family are within the
240 same healthcare system (Johnson et al. 2022). However, we applied a stricter threshold,
241 considering any pair of individuals with a kinship coefficient above 0.044 (up to third-degree
242 relatives) as related. Notably, when cross-checking our results with additional metadata obtained at
243 a later stage of the study, we found that our results with IBIS matched more than 90% of the
244 recorded information related to duplicates and relatives.

245 *GDI validation*

246 We compared the performance of QUILT across different post-imputation statistics on the three
247 test samples with high coverage data (mean coverage $\sim 44\times$) (Figure 2). Without any filters (raw
248 results in Figure 2), QUILT shows the lowest accuracy estimates for each of the statistics we
249 examined.

250

251 A notable improvement over unfiltered results was observed when applying INFO score filtering
252 alone (green line in Figure 2). Since the use of INFO score filtering with QUILT had not been
253 previously explored, we applied the method from Liu et al., originally designed for STITCH (Davies
254 et al. 2016), to evaluate INFO scores with QUILT. The most notable comparison is between the
255 GDI filtering approach and the GP-based filtering ($\max(\text{GP}) \geq 0.99$) (yellow versus purple line in
256 Figure 2). It is evident that both filtering approaches resulted in an improvement compared to the
257 raw data in all the statistics. Our tests confirmed that applying a GP filter at the individual level
258 enhances data quality. Conversely, the GDI filter shows a smaller improvement but is crucial for
259 retaining the majority of variant sites in large batches of data, thereby avoiding the introduction of
260 any missing data, as shown below in case of PCA analyses. Indeed, the GP filtering approach can
261 be effectively used only on a single genome or a few genomes, since in larger batches, most, if not
262 all, variants would have missing genotype calls introduced by the GP filter in at least some
263 individuals. To illustrate the difference on a small example, Figure 3 shows results from a batch of
264 three individual NIPS genomes: the GDI filter retains 83% of variants, while the GP filter, applied
265 cumulatively, retains only 49.2% of total sites. These differences increase with larger sample sizes.
266 In other words, Figure 3 illustrates on a small sample that GP filtering at individual level leads to
267 massive cumulative missingness when the data is pooled. In contrast, the GDI approach maintains
268 a higher proportion of non-missing variants across the batch, avoiding these issues.

269

270 To validate the effect of the GDI filtering approach on imputation accuracy in an independent
271 dataset, we used LPS and high-quality genotype data from 138 non-pregnant healthy controls that

272 were not part of the GDI filter development. Figure 4 summarizes the results of 138 individuals,
273 showing that for all MAF bins > 0.05, the GDI filtered LPS data had higher heterozygous sensitivity
274 than the raw data. Furthermore, the GDI filtered data scored better in other accuracy statistics as
275 well, such as heterozygous specificity and the non-reference concordance (Supp. Table 5). Metrics
276 from these additional samples validated our observations from the three NIPS samples. Moreover,
277 these results demonstrate that the GDI filter is robust when applied to other datasets.

278

279 *GDI application*

280 In light of the success of the GDI filtering approach in minimizing the presence of low-quality
281 variants across the largest possible subset of samples, the 28,512 NIPS samples that passed
282 initial quality checks and removal of related samples were further processed using this method. We
283 removed an additional 1,158 outliers based on the observed distribution of low-quality variants
284 (LQV scores, defined by the proportion of sites not passing the thresholds set for GP values, DS
285 values, and the INFO score provided by QUILT) (Figure S6). Notably, the average coverage for
286 these 1,158 samples is approximately 0.09x, confirming that the LQV metric effectively identified
287 low-coverage samples. These samples were classified as outliers based on their high LQV levels
288 prior to the application of GDI.

289 Furthermore, the distribution of SeqFF values revealed that samples removed due to the LQV filter
290 exhibited a higher median SeqFF compared to those retained for further analyses (Figure S7).
291 Specifically, the median SeqFF for LQV samples was 10.1, whereas for retained samples, it was
292 8.99 (Wilcoxon rank-sum test p-value = 4.147×10^{-13}). This confirmed the negative effect of SeqFF
293 on imputation accuracy.

294 *Evaluation of GDI effectiveness across coverage levels*

295 To evaluate the effectiveness of the GDI thresholds (INFO, GP, and DS) across different
296 sequencing coverage levels, we grouped the samples into coverage bins and assessed the
297 difference in data quality (LQV) before and after GDI application (Figure S8). The analysis revealed

298 a consistent positive impact on accuracy of the GDI method across all coverage bins, with the most
299 pronounced effect observed at lower coverage levels. For bins in the range ≥ 0.1 to $< 0.2x$, the
300 median indicates that the differences between pre- and post-GDI are more uniformly distributed.
301 However, the internal variability is quite high, with some samples showing much larger
302 improvements (outliers), which likely contributes to the observed overall uniformity in this bin. A
303 similar trend is observed in the ≥ 0.2 to $< 0.3x$ bin, where the differences remain generally uniform
304 but with no extreme outliers. In the ≥ 0.3 bins, the median indicates that the differences between
305 pre- and post-GDI are generally modest or small. It is important to note that within each coverage
306 bin, especially for those with coverage ranges of ≥ 0.1 to $< 0.2x$ and ≥ 0.2 to $< 0.3x$, the large number
307 of samples results in significant variability in the individual sample coverage. This variability may
308 influence the observed results and should be considered when interpreting the data.

309 *PCA analyses*

310 The final set of 27,354 samples was used alongside 498 samples from the Genome of the
311 Netherlands (GoNL) (The Genome of the Netherlands Consortium 2014), 3,643 samples from the
312 Project MinE (Project MinE ALS Sequencing Consortium 2018), and 2,495 samples from the 1000
313 Genomes Project (1KGP) (The 1000 Genomes Project Consortium et al. 2015) to visually explore
314 the distribution of genetic ancestry in the imputed data with PCA and UMAP approach, aiming to
315 assess the performance of downstream population genetic analysis method and to test whether
316 imputation is causing batch effects in the data.

317 We applied UMAP on 20 PCs estimated from the data as this approach has been shown to enable
318 the detection of regional clusters at higher resolution than the examination of only the first two PCs
319 (Diaz-Papkovich, Anderson-Trocmé, and Gravel 2021). Without the GDI filter (raw imputed data), a

320 clear batch effect can be observed in the UMAP (Figure S9, panel A); as expected, as UMAP
321 tends to allocate more space to the most represented group (Diaz-Papkovich, Anderson-Trocmé,
322 and Gravel 2021), LPS samples occupy most of the space, but no clear structure in the data can
323 be observed. A similar result is observed when plotting a UMAP after applying a filter on the INFO
324 score (<0.4) (Figure S5); the plot still shows a strong batch effect, similar to that seen with the raw
325 data in Figure S9. This suggests that, although filtering based on the INFO score improves various
326 post-imputation statistics, it does not effectively remove the batch effect. However, applying a GDI
327 filter and removing variants with low quality in more than 30% of the samples (Figure S9, panel B)
328 effectively helps to remove this bias and reorganize the data more clearly. In particular, we can see
329 that the larger group predominantly composed of LPS samples (bigger gray cloud in Figure S9,
330 panel A), is merged after filtering with the smaller separate cluster including individuals from the
331 1000 Genomes Project (GBR and CEU), most of the MinE dataset, and the GoNL group.
332 Furthermore, part of the LPS samples is clearly clustering with non-European groups. This batch
333 effect is also detectable in a milder form in the plot of PC1 and PC2 (Figure S9, panel C), where
334 the distribution of the Belgian NIPS samples appears not completely aligned with the distribution of
335 the other datasets of individuals with similar ancestry (MiNE and GoNL).

336 As a general observation, the removal of variants with low quality scores led to the elimination of
337 approximately 1,9 million variants, reducing the overall dataset of common variants to 3,355,663
338 variants. Following this removal, an average reduction of 54.1% was observed in the LQV scores
339 compared to pre-filter values (Figure S6). Upon removal of the low quality variants the imputed
340 LPS data appears to include less individuals outside the triangular structure of the PC plot (Figure
341 S9, panels C versus D).

342 *Genetic ancestry and population structure in Belgian NIPS data*

343

344 As genetic ancestry is an important confounding factor for many medical and population genetic
345 analysis tools, it is important to know how sensitive methods that cluster individuals by their
346 ancestry are to imputation quality. When focusing on the results of the first two PCs (Figure 5A),
347 we first confirmed that imputed NIPS samples exhibit a wide distribution in context of global

348 reference data: most of the samples are primarily clustered along the European edge of the plot,
349 representing 89.17% of the imputed data, while the remaining 10.83% showed dispersion towards
350 the African and the East Asian edges of the plot.

351

352 These results from PCA were also confirmed with supervised admixture, where 87.3% of the
353 imputed NIPS samples had a European ancestry proportion of $\geq 90\%$. Within the subset of NIPS
354 data where country of origin was known, 165 individuals from Iran formed a cluster right below the
355 European edge in the PCA, which becomes more distinct in the UMAP conversion (Figure 5B).
356 Similarly, after applying the GDI filter, within the European subset a cluster of 131 individuals from
357 Poland becomes more clearly pronounced (Figure 5B and Figure S9). The majority of NIPS
358 individuals whose birthplace was identified as being in Belgium mapped together with MinE and
359 GoNL individuals on PC plots (Figure S10), and between IBS, CEU and GBR populations of the
360 1000 GP data, consistent with previous observations (Van Den Eynden et al. 2018).

361

362 PGS for height

363

364 For 1,911 imputed NIPS samples, polygenic scores (PGS) for height were calculated at different p-
365 value thresholds (Pt) using four approaches: (i) data filtered by GDI and $\text{MAF} \geq 5\%$, (ii) data filtered
366 by $\text{MAF} \geq 5\%$ only, (iii) data filtered by $\text{MAF} \geq 1\%$, and (iv) data limited to HapMap SNPs. These are
367 referred to as PGS_{GDI} , $\text{PGS}_{\text{MAF}5\%}$, $\text{PGS}_{\text{MAF}1\%}$, and $\text{PGS}_{\text{HapMap}}$, respectively. When performing linear
368 regression, where the height of an individual is predicted by PGS, PGS_{GDI} consistently had a lower
369 variance explained for all Pt values than the other PGSs (Figure 6).

370

371 Whereas the three other PGSs mostly had similar values (Table S4). Overall, the best PGS model
372 was $\text{PGS}_{\text{HapMap}}$, with a variance explained of 0.237 at $\text{Pt} = 0.1$, including 62,059 SNPs. However, it
373 is worth noting that the variance explained for $\text{PGS}_{\text{HapMap}}$, $\text{PGS}_{\text{MAF}5\%}$ and $\text{PGS}_{\text{MAF}1\%}$ doesn't change

374 much from $P_t = 0.05$ onward, indicating that the addition of SNPs with P -value > 0.05 does not have
375 much added value towards prediction. We also calculated the Spearman's correlation coefficient (r)
376 of the correlation between genotypic (i.e. the PGS) and phenotypic height. For PGS_{HapMap} ,
377 $PGS_{MAF5\%}$ and $PGS_{MAF1\%}$ the correlation was the same, $r = 0.46$. This is in contrast with PGS_{GDI}
378 where $r = 0.36$. Together, these findings indicate that filtering the imputed NIPS data with the GDI
379 filter does not improve, at least in case of height, PGS performance in terms of either variance
380 explained or correlation with phenotypic height.

381

382 **DISCUSSION**

383

384 In this study, we explored the potential of genotype imputation of large cohorts of low-pass (0.1-
385 0.3 \times) sequence data from Belgium, for the population genetic study of genetic ancestry and PGS
386 estimation for medical risk stratification. In tests performed on an individual level, we found that
387 imputation can be achieved with comparable quality using QUILT (Davies et al. 2021) and
388 GLIMPSE2 (Rubinacci et al. 2023). When using ultra-low coverage ($< 0.2\times$) sequences without
389 quality filters, imputation accuracy appeared, however, to be too low for applications, such as PCA.
390 Our tests confirmed that imputed genotypes from 0.1-0.2 \times sequence data generated with QUILT
391 and GLIMPSE2 without further quality filters exhibit heterozygote sensitivity less than 0.95. For
392 downstream applications that require higher imputation accuracy, an alternative approach of
393 applying quality filters, such as GP, separately on each individual leads to high (> 0.98) individual
394 sensitivity scores but also introduces cumulative missingness when working on batches of data. As
395 a solution we introduced batch level filtering approach GDI that optimizes the balance imputation
396 accuracy and the retention of the largest possible subset of samples and variants. A similar result
397 could not be achieved by filtering solely based on INFO score (< 0.4), as this approach did not
398 resolve the batch effect observed in the UMAP. However, our results confirmed the validity of using
399 INFO score filtering as a foundation for improving data quality, which is why we adopted this metric
400 as the first step in our GDI filtering strategy. While we showed that the GDI approach was effective
401 in removing batch effects caused by low imputation accuracy in low coverage cohorts seen in PC
402 analyses, the GDI approach did not improve the PGS prediction of height. We observed negative

403 correlation between fetal fraction and sensitivity across the range of different coverage values
404 (Figure S3) suggesting that although higher coverage ensures higher imputation accuracy, this
405 relationship may be compromised in cases of high fetal fraction.

406 The best PGS for height was $\text{PGS}_{\text{HapMap}}$ with a variance explained of 0.237 and a correlation with
407 phenotypic height of 0.46 (Figure 8). $\text{PGS}_{\text{MAF5\%}}$ and $\text{PGS}_{\text{MAF1\%}}$ had similar results. In contrast, with
408 PGS_{GDI} only 14% of phenotypic variance of height could be explained (Table S4). A possible
409 reason for this finding could be that there were less SNPs available to build the PGS compared to
410 data that was only MAF 1% or MAF 5% filtered. However, the SNP count alone cannot fully explain
411 the observed difference in variance explained and correlation. This was demonstrated by
412 $\text{PGS}_{\text{HapMap}}$, which had the best performance even though there were less SNPs in the initial
413 dataset than the GDI filtered data. On the other hand, it is important to note that even though in the
414 HapMap filtered dataset there are less SNPs than in the GDI dataset, these SNPs have been
415 chosen to capture most of the common variability in the human genome (The International
416 HapMap Consortium 2005). Common SNPs, which are predominantly the signals detected by
417 GWAS, are consequently the ones used for PGS. Together, this could explain why the PGS_{GDI} is
418 not performing as well as the other filtered datasets. Of note, we only calculated a PGS for the
419 complex trait height, which is highly polygenic – or even omnigenic – in nature (Yengo et al. 2022;
420 Boyle, Li, and Pritchard 2017). Hence, these results could differ for other traits or diseases that
421 have a different genetic architecture. One commonly observed cause of the drop of variance
422 explained in PGS models is the mismatch between genetic ancestry of the test cohort and the
423 GWAS cohort on which the GPS model is based on. However, we cannot think of a reason why the
424 GDI filter would remove specifically variants that are informative for height PGS in Europeans.
425 Furthermore, the imputation reference panels we used were chosen to maximize the success of
426 genotype imputation in a cohort of predominantly European ancestry.

427 When selecting the most appropriate reference panel for imputation, we tried to ensure that it
428 would be representative of all major genetic ancestry components of the sampled population. As
429 the samples were predominantly from Flemish hospitals, we expected the majority of our cohort to
430 have European ancestry with the addition of ~10% of Asian or African ancestry. We compared the

431 performance of a global reference panel of the 1000 Genomes Project (The 1000 Genomes
432 Project Consortium et al. 2015) and the Haplotype Reference Consortium (HRC) panel (the
433 Haplotype Reference Consortium 2016) which includes the former and is enriched for individuals
434 with predominantly European ancestry. Given the larger number of individual sequences and the
435 number of variants the HRC panel provided, we achieved a higher imputation accuracy and found
436 a smaller number of variants that had to be filtered out (Table 1, Tables S2-3). We observed the
437 greatest improvement of accuracy for variants with MAF 0.01 to 0.05 which may have importance
438 for downstream applications that consider variants in this low frequency spectrum. The imputation
439 accuracy of rare variants (MAF <0.01) remained too low with both panels and filter settings for any
440 downstream purposes examined in this study.

441

442 One possible strategy to improve imputation quality at the batch level is to apply filtering by
443 summary quality measures calculated for the entire batch. INFO score is the only such measure
444 available in the QUILT output. However, with large cohorts that are constantly being supplemented
445 with new data, such as the 32,769 individuals examined here, it can be preferable for practical
446 reasons to carry out imputations in smaller batches and to merge them in the end for a common
447 set of variants. Only 131,908 variants (2.4% of common variants) had INFO score below the 0.4
448 threshold in all the batches of NIPS data we studied, which offered only a minor improvement of
449 quality over the option of using no filters. Furthermore, we observed that incorrectly imputed
450 variants included many with high INFO score, a subset of which had a wrongly imputed genotype
451 while their GP and DS values pointed consistently to the genotype called from the high coverage
452 data. This motivated us to carry out systematic revisions of conflicts between the GT and GP
453 values, which could be characteristic only to ultra-low coverage data, and to focus in our search for
454 optimal solutions in the GDI design on the combination of info score, GP and DS filters. We
455 acknowledge that in ideal circumstances sequencing to the minimum depth of 0.5x should be
456 encouraged to reach sufficient accuracy for downstream applications that would not require further
457 filtering (Sousa de Mota et al. 2023). However, in case of data sets that have already been
458 generated in the past, and/or where it is challenging, if not impossible, to generate more data, the
459 GDI approach may offer leverages for data that otherwise cannot be used for genomic scale

460 genotype analyses. Application of the GDI filter after imputation enabled us to achieve imputation
461 accuracy in the large batch of 32,769 individuals on par, though less superior than that achieved
462 with the GP filter applied at individual level (Figure 2). However, the advantage represented by the
463 GDI approach over the GP filter at a batch level was significant as It allowed us to prevent the
464 cumulative introduction of missing data by each individual, keeping the majority (83%) of the
465 imputed common variants for downstream analyses. Furthermore, when we assessed the
466 effectiveness of GDI across varying coverage levels, we observed that the method generally
467 improved data quality in all coverage categories. The greatest improvements were noted in
468 samples with lower coverage, while higher coverage bins displayed more variability in the degree
469 of improvement. This variability may reflect factors such as differences in sample characteristics
470 and coverage distribution within each bin. These results demonstrate the broad applicability of
471 GDI, while suggesting that fine-tuning of the thresholds may be needed for datasets with higher
472 coverage.

473

474 The study by Sousa de Mota (2023) showed that individuals imputed from 0.1-0.25x coverage
475 were correctly placed in the expected continental clusters in PCA while showing at the same time
476 significant deviations from the precise placement within those clusters. Our results confirm this:
477 when using imputed genotype data without filters as the input we find that the majority of
478 individuals map in the PC1 and PC2 plot to the context of the reference sources from Belgium and
479 the Netherlands, as expected, while in the higher resolution UMAP analyses of P1-20 data we find
480 that the unfiltered data cluster separately from the others (Figure S9). With our final settings of the
481 GDI filter, which retained 3,355,663 out of 5,421,789 common variants (61.9%), we were able to
482 remove the separate clustering of NIPS data as a batch. Regarding the determination of the cutoff
483 for variant removal in the GDI filter, we leveraged visual exploration of the sample distribution on a
484 UMAP plot after a series of filtering steps, with increments of 10% of LQV sites being removed at
485 each step, to guide our decision-making. However, in scenarios where such visual exploration is
486 not feasible, the fundamental principle guiding this approach is to strike a balance between being
487 overly conservative and excessively aggressive in variant removal, avoiding the loss of a
488 substantial number of variants. While it might seem arbitrary, custom-oriented handling of this step

489 is essential, recognizing that, in principle, each combination of a dataset and a downstream data
490 analysis method can be unique and necessitate tailored consideration.

491

492 As we sought a solution for using 0.1-0.3x sequence data in ancestry analyses we developed the
493 specific combination GDI filters to be optimized for the performance of high resolution PCA on
494 Belgian NIPS data and the same parameter settings we used may not offer the most optimal
495 solution in other settings. We expect, however, that when appropriate reference panels are
496 available, e.g. for other European cohorts with similar coverage data, the GDI filtering following
497 QUILT (or GLIMPSE) imputation can be a useful approach for genetic ancestry analyses with PCA.
498 Other downstream analyses, including those focused on PGS, can, most likely, benefit from
499 individual filter optimizations that are different, either more conservative or stringent, from those
500 described and explored here. In cases where increasing the sequence coverage is not an option, a
501 generic strategy that may prove effective should involve exploration of the distribution of the
502 number of low-quality variants against the number of samples in which they fail, aiding in defining
503 the most appropriate cutoffs and optimal reduction in LQV scores. A further point worth stressing
504 here is that since the imputation of heterozygous sites is always the most challenging, it is highly
505 likely that such sites are the ones with lower GP and DS values and therefore do not pass the GDI
506 filter. As a consequence, more stringent filters can lead to a reduction of heterozygosity in the
507 sample.

508 Applying a genome-wide SNP array of around 300,000 sites on 189 individual samples from
509 Belgium, Van den Eynden et al. (2018) showed typical European genetic constitution of the
510 Belgian population. Our exploratory analyses of the downstream effects of imputation quality filters
511 on principal component analysis have facilitated the study of the genetic ancestry of 32,769
512 pregnant individuals recently collected from a hospital in present-day Flanders at the genomic
513 scale for more than 3 Million common variants. This has provided a higher resolution view on the
514 genetic ancestry of present-day Flanders and has enabled better quantification of the extent of
515 recent migration. The majority of the samples group within regional clusters of Belgian and Dutch
516 individuals from the MinE and GoNL projects, while a minor subset clusters with other European
517 sub-groups and individuals with non-European genetic background. These results are not

518 surprising given Belgium's long history of immigration and integration. Historically, Belgium is
519 known to have been a nation characterized by persistent immigration, a phenomenon that has
520 undoubtedly contributed significantly to shaping its demographic landscape. Indeed, this
521 continuous demographic evolution has increasingly enriched what is now Belgium's population,
522 which to date welcomes individuals from different corners of the world (Martiniello 2013). Notably,
523 10.83% of the imputed samples mapped outside the main European cluster in PCA (Figure 4A),
524 which is similar to a recent estimate according to which 11% of the population living in the Flemish
525 Region was born outside the EU. Given that our samples come primarily from Flemish hospitals,
526 they represent the demographic diversity of individuals residing in Belgium, including not only
527 those of Belgian origin but also a broader spectrum. As a result, our findings offer an accurate
528 portrayal of the diverse population found in contemporary Belgium, truly reflecting the cosmopolitan
529 nature of the country and providing a realistic snapshot of its current demographics.

530

531 **MATERIALS AND METHODS**

532

533 A concise description of the materials and methods used in this study is provided below. For a
534 comprehensive and detailed version of all protocols and procedures, please refer to the Extended
535 Supplementary Methods.

536

537 *Data collection*

538

539 Peripheral blood samples were collected from 32,769 pregnant individuals undergoing NIPS and
540 from 140 non-pregnant healthy controls, with ethics approval from KU Leuven and University
541 Hospitals Leuven (S63253, S66621, S66450). cfDNA was extracted from plasma and sequenced
542 at low-pass coverage (median $\sim 0.15\times$) using Illumina platforms, following protocols previously
543 described (Bayindir et al. 2015). Sequence data were aligned to hg38, and standard QC,
544 deduplication, and sorting procedures were applied. For the 140 controls, high-quality genotyping
545 was also performed for benchmarking imputation accuracy.

546

547 *Imputation tool testing*

548

549 To select the optimal imputation tool, we evaluated GLIMPSE, GLIMPSE2, QUILT, and a dual-step
550 strategy combining either with Beagle. Imputation accuracy was benchmarked using three low-
551 pass NIPS samples (~0.16x) with known fetal fractions and corresponding high-coverage whole-
552 blood genomes. All tools were applied using the HRC reference panel, and accuracy was
553 assessed per genotype class and MAF bin using metrics such as sensitivity and dosage r^2 .

554

555 *Reference panel setup*

556

557 All imputation strategies used either the HRC or 1000 Genomes Project reference panels, both of
558 which were preprocessed using standard liftover and variant filtering procedures.

559

560 *Imputation with QUILT*

561

562 Genome-wide genotype imputation of the 32,769 samples was performed using QUILT with the
563 HRC panel. Samples were processed in parallel batches using a custom pipeline (Figure 1A), with
564 chromosomes divided into 5 Mb windows. Post-imputation filtering retained 5.42 million variants
565 with MAF >5% in the HRC panel for downstream analysis.

566 **Post imputation filters**

567

568 *Removal of duplicates and related individuals*

569

570 Related individuals were identified and removed based on kinship coefficients estimated with IBIS,
571 resulting in the exclusion of 4,257 samples.

572

573 *Filter on variants at batch level*

574

575 To enhance the quality of the imputed data, we implemented a filtering strategy we called GDI
576 (Figure 1, panel B) that combines filters on genotype-related metrics, including the posterior
577 genotype probability (GP), alternate allele dosage (DS) and the INFO score provided by QUILT.
578 The INFO score provides a quantitative measure of the certainty associated with genotype
579 imputation based on the distribution and uniformity of genotype posteriors; a low score indicates a
580 flat, non-informative distribution, while a score near 1 suggests concentrated, confident genotypes.
581 The DS field represents the expected number of alternate alleles for a given genotype. For diploid
582 genotypes, the DS values should range from 0 to 2. Ideally, a DS equal 0 means that both alleles
583 are the reference allele (0/0), a DS equal 1 means that one allele is the reference allele and the
584 other is the alternate allele (0/1), and a DS equal 2 means that both alleles are the alternate allele
585 (1/1). However, it is possible for the DS values to assume intermediate values between 0 and 2 for
586 diploid genotypes. This can happen when there is uncertainty in the genotype call, such as when
587 the genotype posterior probabilities in the GP field are not clearly in favor of one genotype over
588 another. In this light, the dosage value provides an indication of how well the genotype is
589 supported by imputation. As for the GP field, it represents a measure of how likely each possible
590 genotype at a site is after imputation, with values closer to 1 indicating a higher likelihood and
591 greater confidence in the accuracy of the predicted genotypes following imputation.

592

593 When imputing multiple samples with QUILT, the INFO score associated with each variant
594 represents a consensus score. Due to imputing our samples in separate batches, different INFO
595 scores were obtained for the same variants. To address this variability, we aggregated all INFO
596 scores for the same variants across multiple imputed batches. Variants consistently tagged with an
597 INFO score < 0.4 across all batches totaled 131,908 and were consequently identified for removal.
598 The initial step of the GDI strategy involved generating a file containing the GT, GP, and DS fields
599 for each variant and individual. Subsequently, variants tagged for removal from the INFO score
600 screening were excluded from further analysis in subsequent steps. For each sample, a list of
601 variants to be removed is generated, determined by a GP < 0.99 and DS ranges defined by the GT
602 field (DS > 0.1 for GT 0/0, DS < 1.8 for GT 1/1, and $0.8 > DS > 1.01$ for GT 0/1). The DS
603 thresholds are determined based on theoretical expectations and practical considerations, rather

604 than empirical testing. These values account for minor variations and provide a margin of safety in
605 variant classification. By setting these thresholds, we aim to ensure accurate classification despite
606 potential fluctuations in DS values. The proportion of variants earmarked for removal defines the
607 fraction of low-quality variants (LQV score) for each sample, and an observation of the distribution
608 based on LQV scores was conducted to establish a cutoff. Samples with an LQV score exceeding
609 40% were identified as outliers, leading to the removal of a total of 1.158 samples from subsequent
610 analyses.

611

612 The initial step concludes by generating, for each sample, a list reporting sites that do not pass the
613 DS and GP filters. In the second step of the GDI strategy, these lists are used to gather information
614 on the quality of sites for all individuals (except those identified as outliers in the first GDI step and
615 those flagged as duplicates or relatives by IBIS). Ultimately, variants exhibiting low quality in more
616 than 30% of the samples were excluded, totaling 1,933,956 variants. The choice of this cutoff was
617 defined by the visual observation of the distribution of the imputed data in the context of a UMAP
618 plot based on 20 PCs including also samples from the 1KGP (lifted over to build GRCh38 from the
619 HRC panel), the MinE (lifted over to build GRCh38 with Crossmap), and the GoNL datasets (lifted
620 over to build GRCh38 with Crossmap). The filtered dataset comprised 27,354 samples and
621 3,355,663 variants, reflecting a reduction of approximately 11% in the sample count and about
622 38% of the total variants.

623

624 The average reduction of the proportion of LQV sites was calculated by taking the percentage
625 difference for each sample, followed by calculating the mean and standard deviation of these
626 differences. The standard deviation of 6.2% suggests that variations between individual samples
627 and the mean are relatively consistent.

628

629 After applying the GDI filter (Figure S12), we observed that 86.4% of the variants had successfully
630 passed all the filters, 4.3% failed the DS filter, 5.1% failed the GP filter, and 4.2% failed both the
631 DS and the GP filters. However, when focusing only on the GP values, we can see that the 90.7%

632 of the variants remaining in the dataset have a GP ≥ 0.99 , the 3.06% have a GP between 0.9 and
633 0.99, and only a 6.24% have a GP < 0.9 . Overall, the majority of retained variants in our dataset
634 exhibit a GP ≥ 0.99 , highlighting a significant presence of reliable genetic variations. This
635 underscores the effectiveness of the GDI filter in ensuring data integrity and enhancing overall
636 variant quality.

637

638 Additionally, when observing the correlation between the LQV scores and the coverage of each
639 sample, we observed that as the coverage increases, the LQV score tends to decrease. This
640 suggests that samples with higher coverage tend to have fewer low-quality variants, indicating a
641 potential association between higher coverage and improved quality of imputed variants. Before
642 applying the filter, we calculated the correlation between the coverage and the LQV scores for
643 each sample and observed a strongly negative value (Spearman's correlation is -0.75). This
644 indicates that higher coverage values correspond to lower LQV scores. In practice, this is
645 suggesting that samples with higher coverage will be associated with a lower proportion of low-
646 quality variants. After the filter, we observed a reduction in the correlation value (Spearman's
647 correlation is -0.66) indicating that despite the reduction in LQV scores, the negative correlation
648 with the coverage persists. In addition, the observed mean LQV score decreased from 0.28 to 0.13
649 after the application of the filter, indicating a consistent reduction in the proportion of low-quality
650 variants among the samples. Furthermore, the standard deviation changed from ± 0.0524 to
651 ± 0.0396 , indicating an increase of the consistency and homogeneity of the post-filter data. Overall,
652 these results show that the application of the filter affected the distribution of low-quality variants
653 among the samples, contributing to increased uniformity of LQV values and a persistent correlation
654 between coverage and LQV scores, albeit slightly attenuated.

655

656 ***Comparison with other filtering strategies***

657

658 To assess the effectiveness of our GDI method compared to other filtering strategies, we employed
659 the three test samples and compared the results using different post imputation statistics for the
660 raw imputed data, the GDI strategy, and a direct filter on $\max(\text{GP}) \geq 0.99$. Additionally, we
661 conducted a cross-method comparison by applying imputation through QUILT, GLIMPSE, and
662 GLIMPSE2. Imputation performances were observed over different allele frequency ranges (MAF
663 bins).

664

665 **Principal components analysis and UMAP**

666

667 Population structure was assessed through PCA and UMAP embedding using LD-pruned variants
668 across imputed NIPS samples and external references (1KGP, MinE, GoNL).

669

670 ***PGS calculation***

671

672 To calculate the polygenic scores (PGS), we used the effect sizes provided by the genome-wide
673 association study (GWAS) for height by Yengo et al. (Nature 2022). Duplicate, ambiguous and
674 multi-allelic SNPs were removed and a MAF filter of 1% was applied.

675

676 Out of the 32,769 NIPS samples, only 2,698 had a reported height available. After excluding
677 samples because of relatedness or being an outlier (see above), non-European samples were
678 excluded as well. Samples were included as European based on a supervised admixture analysis,
679 where samples with a proportion of ≥ 0.95 European ancestry were retained. Based on these
680 criteria, 1,911 samples were included for the PGS calculation.

681

682 PRSice-2 (Choi and O'Reilly 2019) was used for the calculation of PGS for height for a range of
683 predefined p-value thresholds (Pt) (5×10^{-8} , 1×10^{-5} , 1×10^{-4} , 1×10^{-3} , 0.05, 0.1, 0.5, 1)

684 with the following clumping parameters: distance to both ends from the index SNP = 250 kb, r^2 =
685 0.1 and p-value threshold = 1. Scores were calculated for individuals from the non-Finnish
686 European 1000 Genomes Project (1KG-NFE), as well as for 1,911 NIPS samples, using the same
687 SNPs per Pt as selected for the 1KG-NFE scores. Four different filtering methods were used to
688 calculate the scores: (1) 1KG-NFE were MAF 5% filtered, then the selected SNPs were used to
689 calculate scores per Pt in the GDI filtered NIPS (PGSGDI); (2) 1KG-NFE were MAF 5% filtered,
690 then the selected SNPs were used to calculate scores per Pt in MAF 5% filtered NIPS
691 (PGSMAF5%); (3) 1KG-NFE were MAF 1% filtered, then the selected SNPs were used to calculate
692 scores per Pt in MAF 1% filtered NIPS (PGSMAF1%); and (4) an overlap was taken between 1KG-
693 NFE and HapMap SNPs, then the selected SNPs were used to calculate scores per Pt in NIPS
694 (PGSHapMap). The 1KG-NFE dataset contained 6,064,728 SNPs, 8,743,364 SNPs, or 1,116,280
695 SNPs after filtering for MAF 5%, MAF 1% or HapMap SNPs respectively.

696

697 After calculation of the raw scores, the first ten principal components (PCs) were regressed out of
698 the scores. Subsequently, the PC-corrected scores were then standardized against the PC-
699 corrected scores from the 1KG-NFE individuals. Specifically, the scores of each individual were
700 standardized by subtracting the mean score of individuals from the 1KG-NFE group from their own
701 scores and then dividing the resulting value by the standard deviation of the scores within 1KG-
702 NFE. When we use the term PGS, we refer to the standardized PC-corrected scores.

703

704 **Declaration of Generative AI in the Writing Process**

705

706 During the preparation of this manuscript, the authors used ChatGPT to improve the language and
707 readability of the text. After using this tool, the authors reviewed and edited the content as needed
708 and take full responsibility for the final version of the manuscript.

709

710 **Software availability**

711 All source code and custom scripts used in this study are available in the GitHub repositories
712 <https://github.com/SABiagini/GDI> and <https://github.com/SABiagini/PostImputationStats>.

713

714 DATA access

715 The NIPS sequence data used in this study are available under restricted and controlled access, in
716 compliance with the GDPR. The genomic data are locally stored at UZ Leuven and access can be
717 obtained with permission from the local UZ Leuven data access committee (DAC:
718 <https://www.uzleuven.be/en/dac>, dac@uzleuven.be; cc: joris.vermeesch@uzleuven.be). The
719 access process is further described in [https://www.uzleuven.be/en/data-access-committee-](https://www.uzleuven.be/en/data-access-committee-dac/external-applicant-requests-access-ku-leuven-researchers-dataset)
720 [dac/external-applicant-requests-access-ku-leuven-researchers-dataset](https://www.uzleuven.be/en/data-access-committee-dac/external-applicant-requests-access-ku-leuven-researchers-dataset). Sequence data of the 140
721 healthy non-pregnant individuals, used as controls for imputation accuracy, is available through the
722 European Nucleotide Archive (EGA) at EMBL-EBI under Accession No. EGAS50000001114.

723

724 Competing interest statement

725 The authors declare no competing interests.

726

727 ACKNOWLEDGMENTS

728 Our thanks go primarily to all the individuals who participated in this study as donors. We also want
729 to express our gratitude to Robert Davies for the invaluable conversations that helped us
730 determine the most appropriate use of QUILT. Additionally, we are grateful to Ilse Parijs for
731 generously providing the height data of the individuals. This study has been supported by FWO
732 SBO grant S003422N 'MICADO' (SB, IC, TK, JV, BT), KU Leuven start-up grant STG/18/021 (TK),
733 KU Leuven BOF-C24 grant ZKD6488 C24M/19/075 (TK, SAB, LH, JB), and FWO grants
734 G050822N (TK) and G0B4822N (BT). MA is supported by a PhD fellowship of FWO Flanders. JB
735 is funded by a senior clinical investigator fellowship of FWO Flanders. SAB was also supported by
736 the Czech Ministry of Education, Youth and Sports (CZ.02.01.01/00/22_008/0004593, RES-HUM:
737 Ready for the Future: Understanding the Long-Term Resilience of Human Culture grant). The
738 computational resources and services used in this work were provided by the VSC (Flemish
739 Supercomputer Center), funded by the Research Foundation Flanders (FWO) and the Flemish
740 Government – department EWI, and made available as a Rainbow Project of the Biobanking and
741 BiomolecularResearch Infrastructure Netherlands (BBMRI-NL).

742

743 This study makes use of data generated by the Genome of the Netherlands Project. A full list of the
744 investigators is available from www.nlgenome.nl. Funding for the project was provided by the
745 Netherlands Organization for Scientific Research under award number 184021007, dated July 9,
746 2009 and made available as a Rainbow Project of the Biobanking and Biomolecular Research
747 Infrastructure Netherlands (BBMRI-NL). a. The LifeLines Cohort Study (<http://www.lifelines.nl>), and
748 generation and management of GWAS genotype data for it, is supported by the Netherlands
749 Organization of Scientific Research (NWO, grant 175.010.2007.006), the Dutch government's
750 Economic Structure Enhancing Fund (FES), the Ministry of Economic Affairs, the Ministry of
751 Education, Culture and Science, the Ministry for Health, Welfare and Sports, the Northern
752 Netherlands Collaboration of Provinces (SNN), the Province of Groningen, the University Medical
753 Center Groningen, the University of Groningen, the Dutch Kidney Foundation and Dutch
754 Diabetes Research Foundation; and b. For sponsorship of the EMC Ergo Study please refer to
755 (<http://www.ergo-onderzoek.nl/wp/>); and c. The LUMC Longevity Study was supported by a grant
756 from the Innovation-Oriented Research Program on Genomics (SenterNovem IGE01014 and
757 IGE05007), the Centre for Medical Systems Biology and the National Institute for Healthy Ageing
758 (Grant 05040202 and 05060810), all in the framework of the Netherlands Genomics
759 Initiative/Netherlands Organization for Scientific Research.; d. For sponsorship of the VU
760 Netherlands Twin Register please refer to www.tweelingenregister.org. Project MinE Belgium was
761 supported by a grant from IWT (n° 140935), the ALS Liga België, the National Lottery of Belgium
762 and the KU Leuven Opening the Future Fund.

763

764 Author contributions: J.R.V., I.C, and T.K. designed and supervised the study. J.R.V., M.A., L.H.,
765 T.J., J. B., K. D., and B. T. organized and performed the recruitment of the case samples. P. V. D.
766 provided genomic data of the MinE cohort. S.A.B. and S. B. performed data analyses. S.A.B., T.K.,
767 I.C. and J.R.V. drafted and revised the manuscript. All authors reviewed and contributed to the
768 writing of the final manuscript.

769

770 **REFERENCES**

- 771 Abraham, Gad, Yixuan Qiu, and Michael Inouye. 2017. "FlashPCA2: Principal Component Analysis
772 of Biobank-Scale Genotype Datasets" ed. Oliver Stegle. *Bioinformatics* 33(17): 2776–78.
773 doi:10.1093/bioinformatics/btx299.
- 774 Ausmees, Kristiina, Federico Sanchez-Quinto, Mattias Jakobsson, and Carl Nettelblad. 2022. "An
775 Empirical Evaluation of Genotype Imputation of Ancient DNA" ed. A Sethuraman. *G3*
776 *Genes/Genomes/Genetics* 12(6): jkac089. doi:10.1093/g3journal/jkac089.
- 777 Bayindir, Baran, Luc Dehaspe, Nathalie Brison, Paul Brady, Simon Ardui, Molka Kammoun, Lars
778 Van Der Veken, et al. 2015. "Noninvasive Prenatal Testing Using a Novel Analysis Pipeline to
779 Screen for All Autosomal Fetal Aneuploidies Improves Pregnancy Management." *European*
780 *Journal of Human Genetics* 23(10): 1286–93. doi:10.1038/ejhg.2014.282.
- 781 Boyle, Evan A., Yang I. Li, and Jonathan K. Pritchard. 2017. "An Expanded View of Complex
782 Traits: From Polygenic to Omnigenic." *Cell* 169(7): 1177–86. doi:10.1016/j.cell.2017.05.038.
- 783 Browning, Brian L., Ying Zhou, and Sharon R. Browning. 2018. "A One-Penny Imputed Genome
784 from Next-Generation Reference Panels." *The American Journal of Human Genetics* 103(3): 338–
785 48. doi:10.1016/j.ajhg.2018.07.015.
- 786 Choi, Shing Wan, and Paul F O'Reilly. 2019. "PRSice-2: Polygenic Risk Score Software for
787 Biobank-Scale Data." *GigaScience* 8(7): giz082. doi:10.1093/gigascience/giz082.
- 788 Cingolani, P., Platts, A., Wang, L.L., Coon, M., Nguyen, T., Wang, L., Land, S.J., Lu, X., Ruden,
789 D.M., 2012. A program for annotating and predicting the effects of single nucleotide
790 polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2;
791 iso-3. *Fly* 6, 80–92. <https://doi.org/10.4161/fly.19695>
- 792 Davies, Robert W, Jonathan Flint, Simon Myers, and Richard Mott. 2016. "Rapid Genotype
793 Imputation from Sequence without Reference Panels." *Nature Genetics* 48(8): 965–69.
794 doi:10.1038/ng.3594.

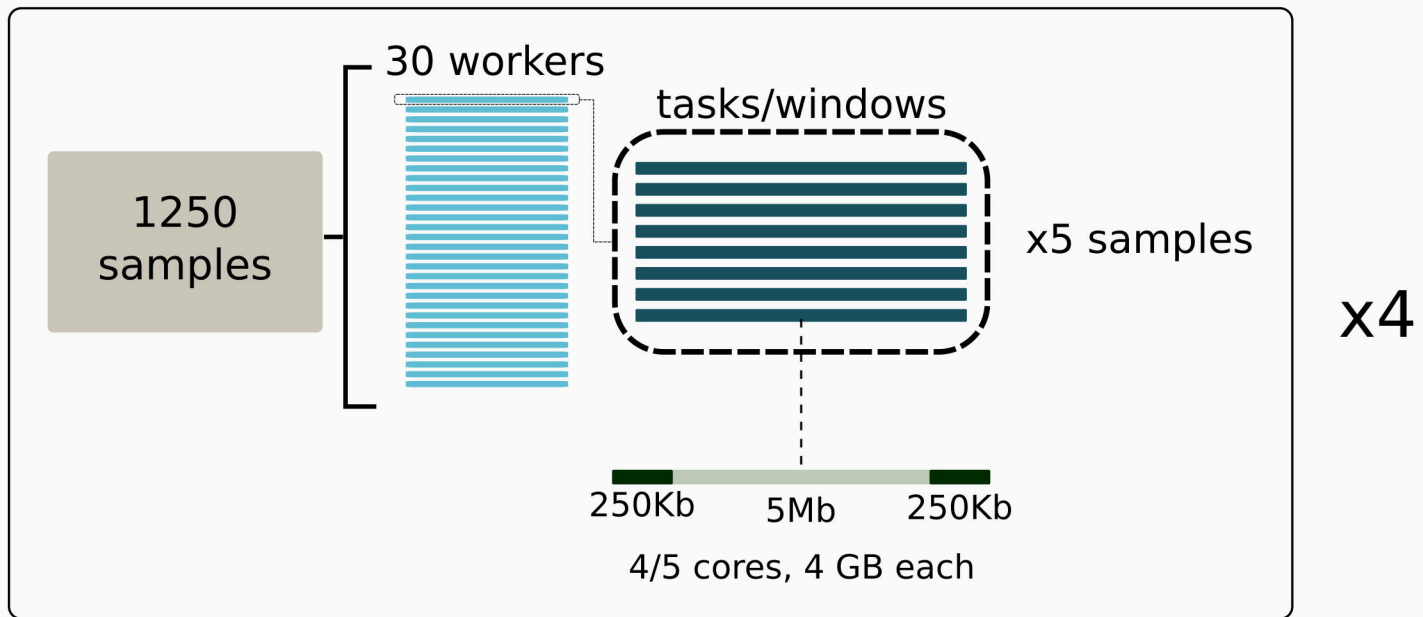
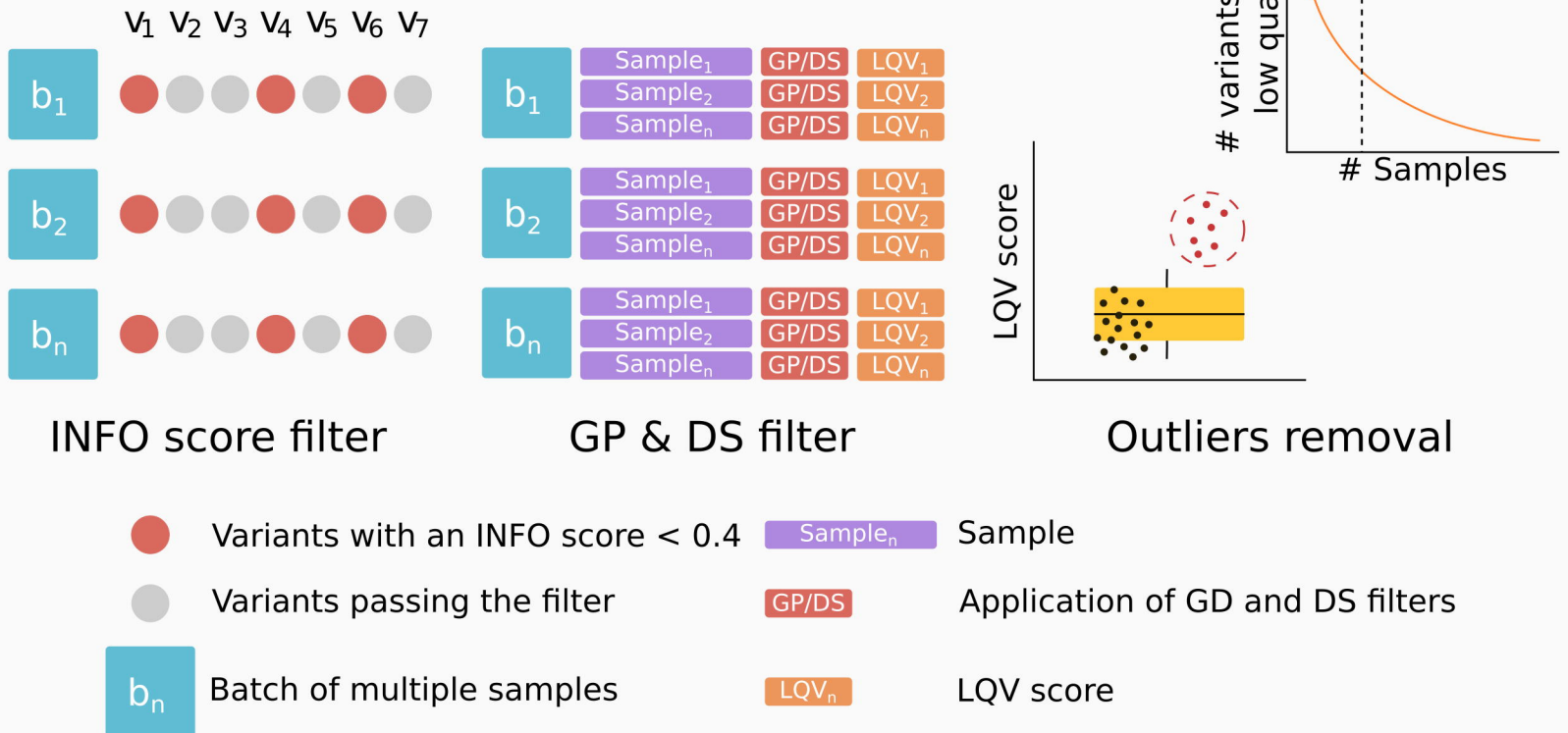
- 795 Davies, Robert W. et al. 2021. "Rapid Genotype Imputation from Sequence with Reference
796 Panels." *Nature Genetics* 53(7): 1104–11. doi:10.1038/s41588-021-00877-0.
- 797 Diaz-Papkovich, Alex, Luke Anderson-Trocmé, and Simon Gravel. 2021. "A Review of UMAP in
798 Population Genetics." *Journal of Human Genetics* 66(1): 85–91. doi:10.1038/s10038-020-00851-4.
- 799 Fumagalli, Matteo. 2013. "Assessing the Effect of Sequencing Depth and Sample Size in
800 Population Genetics Inferences" ed. Ludovic Orlando. *PLoS ONE* 8(11): e79667.
801 doi:10.1371/journal.pone.0079667.
- 802 Gadsbøll, Kasper et al. 2020. "Current Use of Noninvasive Prenatal Testing in Europe, Australia
803 and the USA: A Graphical Presentation." *Acta Obstetrica et Gynecologica Scandinavica* 99(6):
804 722–30. doi:10.1111/aogs.13841.
- 805 Gamba, Cristina, Eppie R. Jones, Matthew D. Teasdale, Russell L. McLaughlin, Gloria Gonzalez-
806 Fortes, Valeria Mattiangeli, László Domboróczki, et al. 2014. "Genome Flux and Stasis in a Five
807 Millennium Transect of European Prehistory." *Nature Communications* 5(1): 5257.
808 doi:10.1038/ncomms6257.
- 809 Homburger, Julian R., Cynthia L. Neben, Gilad Mishne, Alicia Y. Zhou, Sekar Kathiresan, and Amit
810 V. Khera. 2019. "Low Coverage Whole Genome Sequencing Enables Accurate Assessment of
811 Common Variants and Calculation of Genome-Wide Polygenic Scores." *Genome Medicine* 11(1):
812 74. doi:10.1186/s13073-019-0682-2.
- 813 Hui, Ruoyun et al. 2020. "Evaluating Genotype Imputation Pipeline for Ultra-Low Coverage Ancient
814 Genomes." *Scientific Reports* 10(1): 18542. doi:10.1038/s41598-020-75387-w.
- 815 Johnson, Ruth et al. 2022. "Leveraging Genomic Diversity for Discovery in an Electronic Health
816 Record Linked Biobank: The UCLA ATLAS Community Health Initiative." *Genome Medicine* 14(1):
817 104. doi:10.1186/s13073-022-01106-x.
- 818 Kim, Sung K., Gregory Hannum, Jennifer Geis, John Tynan, Grant Hogg, Chen Zhao, Taylor J.
819 Jensen, et al. 2015. "Determination of Fetal DNA Fraction from the Plasma of Pregnant Women

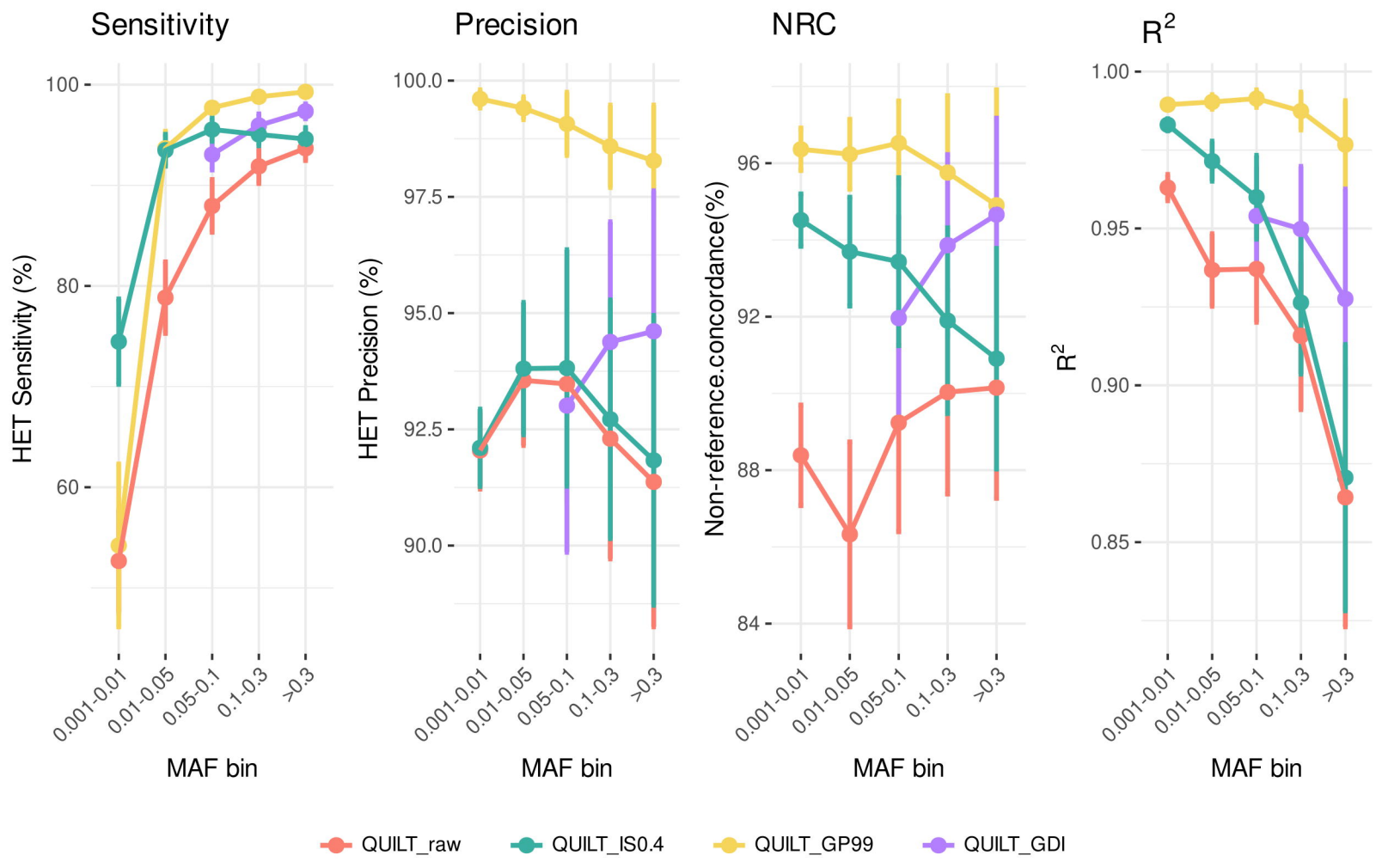
- 820 Using Sequence Read Counts: Determination of Fetal DNA Fraction from the Plasma of Pregnant
821 Women Using Sequence Read Counts.” *Prenatal Diagnosis* 35(8): 810–15. doi:10.1002/pd.4615.
- 822 Li, Heng. 2011. “A Statistical Framework for SNP Calling, Mutation Discovery, Association
823 Mapping and Population Genetical Parameter Estimation from Sequencing Data.” *Bioinformatics*
824 27(21): 2987–93. doi:10.1093/bioinformatics/btr509.
- 825 Linthorst, J., Nivard, M., Sistermans, E.A., 2024. GWAS shows the genetics behind cell-free DNA
826 and highlights the importance of p.Arg206Cys in DNASE1L3 for non-invasive testing. *Cell Reports*
827 43, 114799. <https://doi.org/10.1016/j.celrep.2024.114799>
- 828 Liu, Siyang et al. 2018. “Genomic Analyses from Non-Invasive Prenatal Testing Reveal Genetic
829 Associations, Patterns of Viral Infections, and Chinese Population History.” *Cell* 175(2): 347-
830 359.e14. doi:10.1016/j.cell.2018.08.016.
- 831 Liu, Siyang, Shujia Huang, Yanhong Liu, Yuqin Gu, Xingchen Lin, Huanhuan Zhu, Hankui Liu, et
832 al. 2023. “Utilizing Non-Invasive Prenatal Test Sequencing Data Resource for Human Genetic
833 Investigation.” *bioRxiv*: 2023.12.11.570976. doi:10.1101/2023.12.11.570976.
- 834 Manichaikul, Ani et al. 2010. “Robust Relationship Inference in Genome-Wide Association
835 Studies.” *Bioinformatics* 26(22): 2867–73. doi:10.1093/bioinformatics/btq559.
- 836 Marchini, J., Howie, B. “Genotype imputation for genome-wide association studies”. *Nat Rev Genet*
837 11, 499–511 (2010). <https://doi.org/10.1038/nrg2796>
- 838 Martin, Alicia R. et al. 2021. “Low-Coverage Sequencing Cost-Effectively Detects Known and
839 Novel Variation in Underrepresented Populations.” *The American Journal of Human Genetics*
840 108(4): 656–68. doi:10.1016/j.ajhg.2021.03.012.
- 841 Martiniello, Marco. 2013. “Belgium, Migration, 1946 to Present.” In *The Encyclopedia of Global*
842 *Human Migration*, ed. Immanuel Ness. Wiley. doi:10.1002/9781444351071.wbeghm063.
- 843 Melville J. 2023. “Uwot: The Uniform Manifold Approximation and Projection (UMAP) Method for
844 Dimensionality Reduction.”

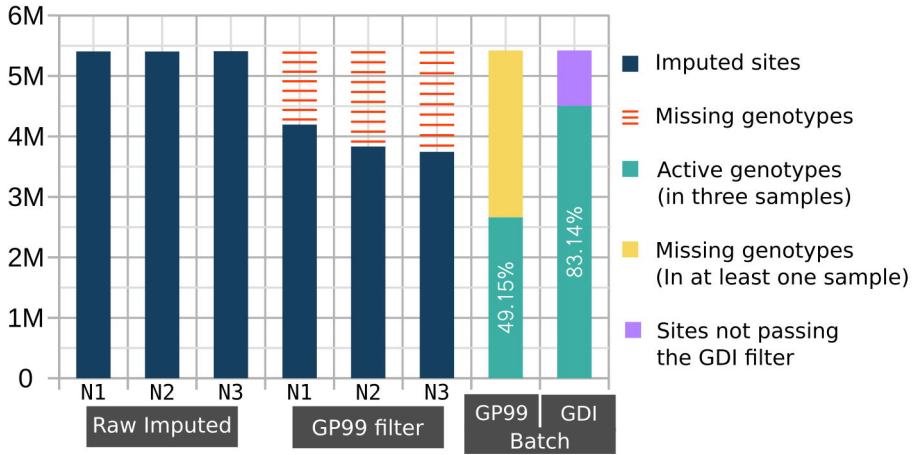
- 845 Mokveld, Tom, Zaid Al-Ars, Erik A. Sistermans, and Marcel Reinders. 2023. "A Comprehensive
846 Performance Analysis of Sequence-Based within-Sample Testing NIPT Methods" ed. Hao Sun.
847 *PLOS ONE* 18(4): e0284493. doi:10.1371/journal.pone.0284493.
- 848 Nguyen, T.V., Bolormaa, S., Reich, C.M., Chamberlain, A.J., Vander Jagt, C.J., Daetwyler, H.D.,
849 MacLeod, I.M., 2024. Empirical versus estimated accuracy of imputation: optimising filtering
850 thresholds for sequence imputation. *Genet Sel Evol* 56, 72. [https://doi.org/10.1186/s12711-024-](https://doi.org/10.1186/s12711-024-00942-2)
851 00942-2
- 852 Pedersen, Brent S, and Aaron R Quinlan. 2018. "Mosdepth: Quick Coverage Calculation for
853 Genomes and Exomes" ed. John Hancock. *Bioinformatics* 34(5): 867–68.
854 doi:10.1093/bioinformatics/btx699.
- 855 Poplin, Ryan et al. 2017. *Scaling Accurate Genetic Variant Discovery to Tens of Thousands of*
856 *Samples*. Genomics. preprint. doi:10.1101/201178.
- 857 Project MinE ALS Sequencing Consortium. 2018. "Project MinE: Study Design and Pilot Analyses
858 of a Large-Scale Whole-Genome Sequencing Study in Amyotrophic Lateral Sclerosis." *European*
859 *Journal of Human Genetics* 26(10): 1537–46. doi:10.1038/s41431-018-0177-4.
- 860 Qiao, Longwei et al. 2019. "Sequencing of Short cfDNA Fragments in NIPT Improves Fetal
861 Fraction with Higher Maternal BMI and Early Gestational Age." *American Journal of Translational*
862 *Research* 11(7): 4450–59.
- 863 Rava, Richard P, Anupama Srinivasan, Amy J Sehnert, and Diana W Bianchi. 2014. "Circulating
864 Fetal Cell-Free DNA Fractions Differ in Autosomal Aneuploidies and Monosomy X." *Clinical*
865 *Chemistry* 60(1): 243–50. doi:10.1373/clinchem.2013.207951.
- 866 Rubinacci, Simone, Robin J. Hofmeister, Bárbara Sousa Da Mota, and Olivier Delaneau. 2023.
867 "Imputation of Low-Coverage Sequencing Data from 150,119 UK Biobank Genomes." *Nature*
868 *Genetics* 55(7): 1088–90. doi:10.1038/s41588-023-01438-3.

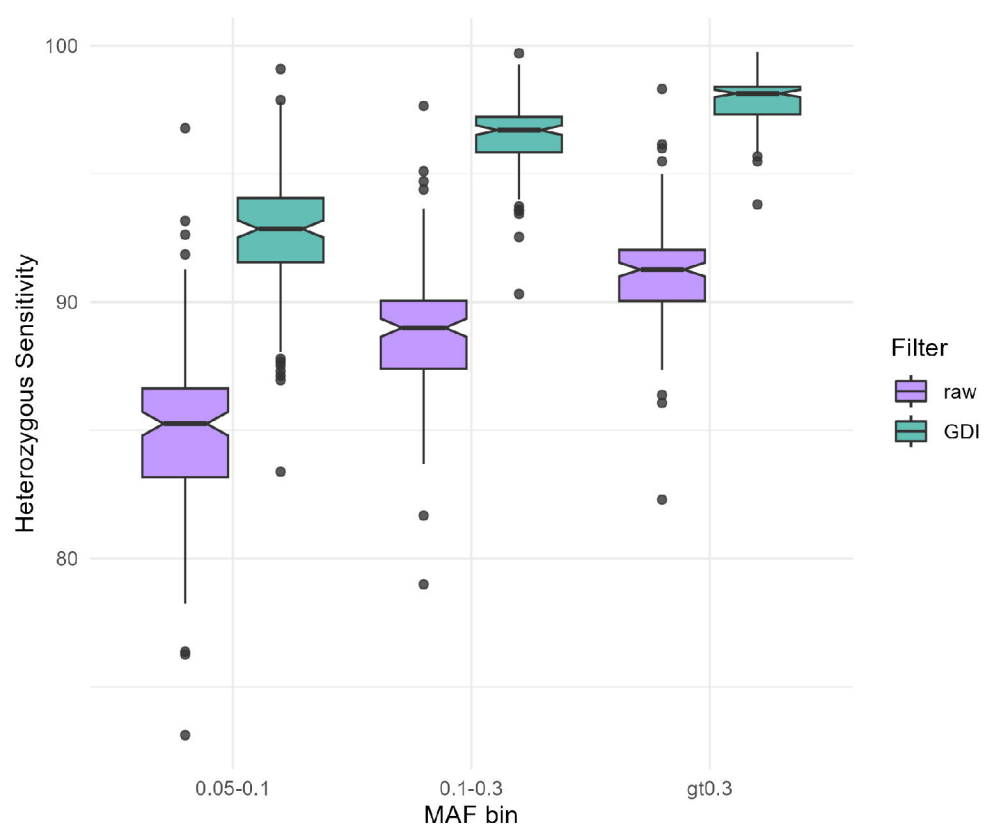
- 869 Rubinacci, Simone, Diogo M. Ribeiro, Robin J. Hofmeister, and Olivier Delaneau. 2021. "Efficient
870 Phasing and Imputation of Low-Coverage Sequencing Data Using Large Reference Panels."
871 *Nature Genetics* 53(1): 120–26. doi:10.1038/s41588-020-00756-0.
- 872 Seidman, Daniel N. et al. 2020. "Rapid, Phase-Free Detection of Long Identity-by-Descent
873 Segments Enables Effective Relationship Classification." *The American Journal of Human*
874 *Genetics* 106(4): 453–66. doi:10.1016/j.ajhg.2020.02.012.
- 875 Sousa Da Mota, Bárbara et al. 2023. "Imputation of Ancient Human Genomes." *Nature*
876 *Communications* 14(1): 3660. doi:10.1038/s41467-023-39202-0.
- 877 The 1000 Genomes Project Consortium et al. 2015. "A Global Reference for Human Genetic
878 Variation." *Nature* 526(7571): 68–74. doi:10.1038/nature15393.
- 879 The Genome of the Netherlands Consortium. 2014. "Whole-Genome Sequence Variation,
880 Population Structure and Demographic History of the Dutch Population." *Nature Genetics* 46(8):
881 818–25. doi:10.1038/ng.3021.
- 882 the Haplotype Reference Consortium. 2016. "A Reference Panel of 64,976 Haplotypes for
883 Genotype Imputation." *Nature Genetics* 48(10): 1279–83. doi:10.1038/ng.3643.
- 884 The International HapMap Consortium. 2005. "A Haplotype Map of the Human Genome." *Nature*
885 437(7063): 1299–1320. doi:10.1038/nature04226.
- 886 Van Den Eynden, Jimmy et al. 2018. "The Genetic Structure of the Belgian Population." *Human*
887 *Genomics* 12(1): 6. doi:10.1186/s40246-018-0136-8.
- 888 Van Riel, Margot, Kate Stanley, and Joris R. Vermeesch. 2023. "Noninvasive Prenatal
889 Testing/Screening by Circulating Cell-Free DNA." In *Human Reproductive and Prenatal Genetics*,
890 Elsevier, 823–51. doi:10.1016/B978-0-323-91380-5.00013-7.
- 891 Vermeesch, Joris Robert, Thierry Voet, and Koenraad Devriendt. 2016. "Prenatal and Pre-
892 Implantation Genetic Diagnosis." *Nature Reviews Genetics* 17(10): 643–56.
893 doi:10.1038/nrg.2016.97.

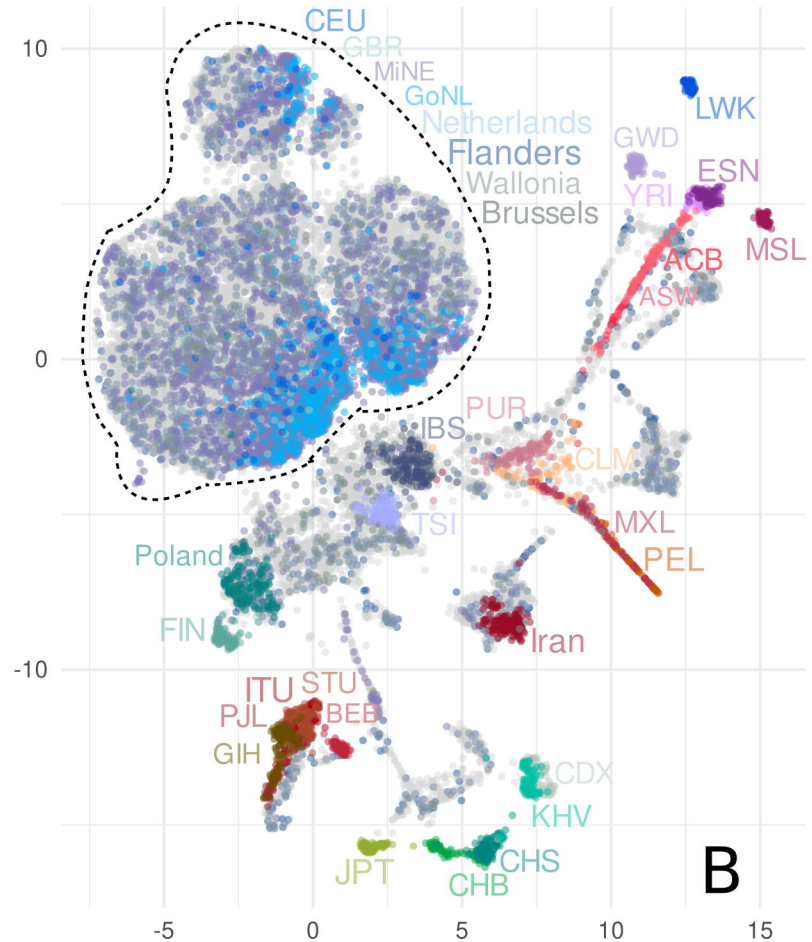
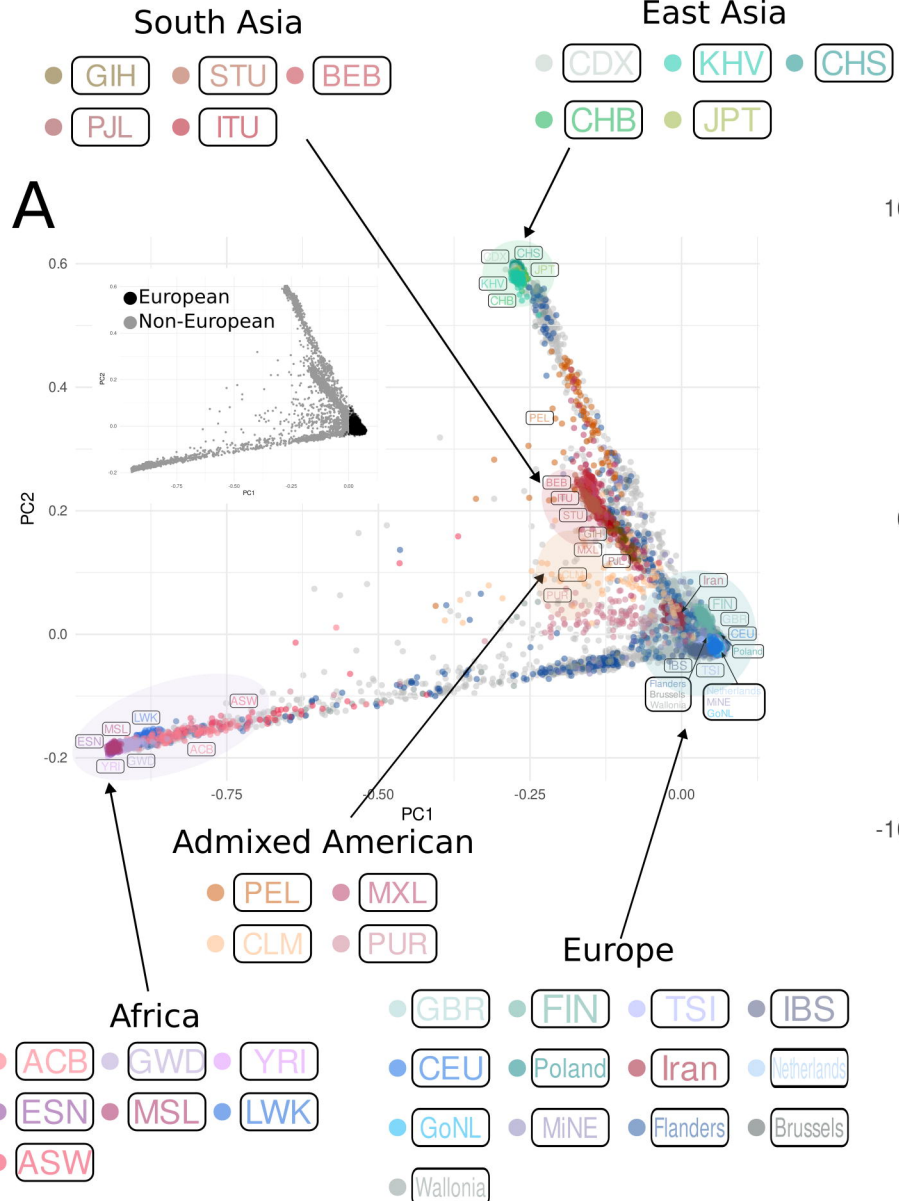
- 894 Wang, Jing-wei et al. 2021. "Cell-Free Fetal DNA Testing and Its Correlation with Prenatal
895 Indications." *BMC Pregnancy and Childbirth* 21(1): 585. doi:10.1186/s12884-021-04044-5.
- 896 Wasik, Kaja et al. 2021. "Comparing Low-Pass Sequencing and Genotyping for Trait Mapping in
897 Pharmacogenetics." *BMC Genomics* 22(1): 197. doi:10.1186/s12864-021-07508-2.
- 898 Yengo, Loïc, Sailaja Vedantam, Eirini Marouli, Julia Sidorenko, Eric Bartell, Saori Sakaue,
899 Marielisa Graff, et al. 2022. "A Saturated Map of Common Genetic Variants Associated with
900 Human Height." *Nature* 610(7933): 704–12. doi:10.1038/s41586-022-05275-y.
- 901 Zhao, Hao et al. 2014. "CrossMap: A Versatile Tool for Coordinate Conversion between Genome
902 Assemblies." *Bioinformatics* 30(7): 1006–7. doi:10.1093/bioinformatics/btt7
- 903 Zhen, Jianxin, Yuqin Gu, Piao Wang, Weihong Wang, Shengzhe Bian, Shujia Huang, Hui Liang, et
904 al. 2024. "Genome-Wide Association and Mendelian Randomisation Analysis among 30.699
905 Chinese Pregnant Women Identifies Novel Genetic and Molecular Risk Factors for Gestational
906 Diabetes and Glycaemic Traits." *Diabetologia* 67(4): 703–13. doi:10.1007/s00125-023-06065

A**B**

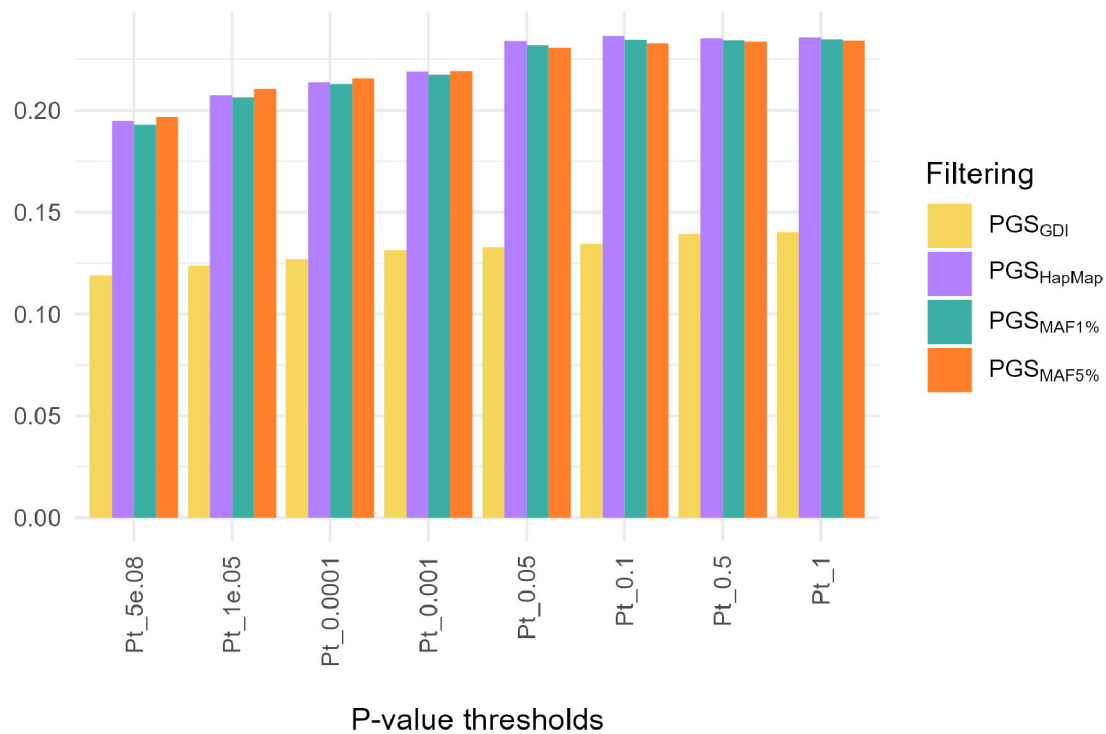








Adjusted R^2



Captions for Main Figures

Figure 1 A) Imputation pipeline with QUILT. The image illustrates the design of imputation of a batch of 1,250 individuals described in further detail in the Materials and Methods section. **B)** GDI filtering scheme: *INFO score filter*. Each column (e.g., v1,v2,v3, etc.) represents a distinct variant, and each row (e.g., b1, b2, bn) corresponds to a batch. Variants with an INFO score < 0.4 in all batches (red dots) are excluded from subsequent steps. This exclusion applies only when the same variant consistently displays this condition across all batches (e.g., three out of seven variants in the example cartoon are excluded). Variants that pass the filter (gray dots) proceed to the next step. *GP & DS filter*. Using the variants that passed the previous step, for each individual within a batch, an LQV score (Low-quality Variant score) is calculated, indicating the proportion of variants that do not meet both GP and DS thresholds. Each sample's LQV score is then used to assess the quality of the data. *Outliers removal*: The distribution of LQV scores is analyzed to identify and remove outlier samples (represented by the red dots in the boxplot). For the remaining samples (represented by the black dots in the boxplot), the distribution of variants with low quality is assessed relative to the total number of samples. The curve in the top-right subpanel illustrates how a cutoff is determined to exclude variants that consistently exhibit low quality in more than a specified percentage of samples. Visual inspection of UMAP plots was used to determine the optimal cutoff (additional details are provided in the main text).

Figure 2 Comparative analysis of post-imputation filters on QUILT performance by four measures of accuracy. Each panel displays results for different post-imputation statistics for those genotypes imputed as heterozygous: Sensitivity ($TP/(TP+FN)$), Precision ($TP/(TP+FP)$), Non-reference concordance ($NRC = 1 - (Err + Era + Eaa) / (Err + Era + Eaa + Mra + Maa)$). Where Err, Era, and Eaa are the counts of the mismatches for the homozygous reference, heterozygous and homozygous alternative genotypes, while Mra and Maa are the counts of the matches at the heterozygous and homozygous alternative genotypes), and the dosage r-squared (Pearson's r^2 on dosage values calculated using BCFtools stats). Lines are colored based on different filtering approaches. Results are average across the three test samples, with SD bars showing the variability around the mean values. (Please, refer to Supplementary Table 2 for detailed information on the values). The lack of QUILT_GDI data points in the 0.001-0.01 and 0.01-0.05 frequency bins reflects our focus on common variants with $MAF > 5\%$ in the HRC panel.

Figure 3 Cumulative missingness of filtered-out variants across multiple imputed genomes. The number of imputed common variants (MAFHRC > 0.05) is displayed in millions on the Y axis for each of the three imputed NIPS samples (N1, N2, and N3). The raw imputed data of each individual had 5.4M variants. Applying GP<0.99 filter (GP99 filter) on each individual separately results in a variable number of missing variants (red bars). The GP99 and GDI Batch bars to the right show the batch effect of cumulative missingness when applying the respective filters. The green bar shows the number of retained sites after removing variants with GP<0.99 (yellow) in at least one sample or variant sites not passing the GDI filter (lilac).

Figure 4. Boxplots per MAF showing the heterozygous genotype imputation sensitivity before (raw) and after GDI filters. Data are derived from 138 non-pregnant healthy individuals sequenced at 0.1-0.2x coverage and genotyped, in parallel, on Illumina Infinium Global Screening Array-24 v3.0 for ~700,000 SNPs. The array-based genotype calls were considered as truth.

Figure 5 Principal Component Analysis of imputed NIPS data in context of regional and global reference populations. **(A)** The plot of PC1 and PC2. The smaller PCA plot highlights European and non-European ancestries; **(B)** UMAP plot of the first 20 PC-s. Both plots display the distribution of 27.354 imputed NIPS samples (gray dots in the background) together with samples from the 1000 Genomes Project, the GoNL, and the MinE datasets, after applying a GDI filter removing variants with low quality in more than 30% of the imputed NIPS samples. Population labels are displayed next to the group they represent (samples from Iran and Poland are part of the Belgian data), with colors of the labels matching the colors of the corresponding samples. Unlabeled imputed samples are depicted as gray dots in the background. The dotted line includes the group represented by samples from the GoNL, the MinE, the 1000 Genomes Project (CEU and GBR only), unlabelled NIPS samples, and NIPS individuals labeled as Flanders, Wallonia, Brussels, and the Netherlands.

Figure 6 Phenotypic variance in height explained by the PGS per p-value threshold and for different filtering approaches. The overall best performing PGS in terms of variance explained is PGSHapMap at p-value threshold (Pt) 0.1 with an r^2 of 0.237 (beta = 0.03, p-value = 2.83×10^{-114}).

Table 1

Average imputation sensitivities of heterozygous genotypes obtained using HRC and 1000GP reference panels.

| Filter | MAF bin | Het sensitivity | | Number of Hets | | Missing Hets | |
|---------------------|------------|-----------------|--------|----------------|--------|--------------|--------|
| | | HRC | 1000GP | HRC | 1000GP | HRC | 1000GP |
| No filter | 0.001-0.01 | 52.7 | 50.9 | 17188 | 14039 | 0 | 0 |
| | 0.01-0.05 | 78.8 | 74.0 | 83397 | 77676 | 0 | 0 |
| | 0.05-0.1 | 88.0 | 84.1 | 127130 | 120522 | 0 | 0 |
| | 0.1-0.3 | 91.9 | 88.5 | 677999 | 647528 | 0 | 0 |
| | >0.3 | 93.7 | 90.4 | 789923 | 755478 | 0 | 0 |
| max(GP) \geq 0.99 | 0.001-0.01 | 54.2 | 50.4 | 5046 | 5669 | 71.5 | 59.3 |
| | 0.01-0.05 | 93.7 | 88.1 | 48631 | 46943 | 50.9 | 49.3 |
| | 0.05-0.1 | 97.7 | 96.1 | 88139 | 82656 | 37.6 | 40 |
| | 0.1-0.3 | 98.8 | 98.3 | 479079 | 452659 | 34.3 | 37 |
| | >0.3 | 99.3 | 99.1 | 559099 | 531081 | 33.2 | 35.9 |
| GDI | 0.001-0.01 | na | na | na | na | 0 | na |
| | 0.01-0.05 | na | na | na | na | 0 | na |
| | 0.05-0.1 | 93.1 | na | 101500 | na | 0 | na |
| | 0.1-0.3 | 96.0 | na | 432852 | na | 0 | na |
| | >0.3 | 97.4 | na | 417039 | na | 0 | na |

Note, Het sensitivity - heterozygote genotype sensitivities are expressed in percentage.

Number of hets - number of correctly imputed heterozygous genotypes in the given MAF bin#

Missing hets - proportion of imputed heterozygous genotypes that were set to missing by the GP99 filter in the given MAF bin