

# 1 Imputation of low-coverage sequencing data from 2 150,119 UK Biobank genomes

3  
4 Simone Rubinacci <sup>1,2</sup>, Robin Hofmeister <sup>1,2</sup>, Bárbara Sousa da Mota <sup>1,2</sup>, Olivier Delaneau <sup>1,2,\*</sup>

5 <sup>1</sup>Department of computational Biology, University of Lausanne, Lausanne, Switzerland

6 <sup>2</sup>Swiss Institute of Bioinformatics, Lausanne, Switzerland

7 \* Corresponding author ([olivier.delaneau@unil.ch](mailto:olivier.delaneau@unil.ch))

8

## 9 Abstract

10 Recent work highlights the advantages of low-coverage whole genome sequencing (lcWGS), followed  
11 by genotype imputation, as a cost-effective genotyping technology for statistical and population  
12 genetics. The release of whole genome sequencing data for 150,119 UK Biobank (UKB) samples  
13 represents an unprecedented opportunity to impute lcWGS with high accuracy. However, despite  
14 recent progress<sup>1,2</sup>, current methods struggle to cope with the growing numbers of samples and  
15 markers in modern reference panels, resulting in unsustainable computational costs. For instance, the  
16 imputation cost for a single genome is 1.11£ using GLIMPSE v1.1.1 (GLIMPSE1) on the UKB research  
17 analysis platform (RAP) and rises to 242.8£ using QUILT v1.0.4. To overcome this computational  
18 burden, we introduce GLIMPSE v2.0.0 (GLIMPSE2), a major improvement of GLIMPSE, that scales  
19 sublinearly in both the number of samples and markers. GLIMPSE2 imputes a low-coverage genome  
20 from the UKB reference panel for only 0.08£ in compute cost while retaining high accuracy for both  
21 ancient and modern genomes, particularly at rare variants (MAF < 0.1%) and for very low-coverage  
22 samples (0.1x-0.5x).

## 23 Main

24 To demonstrate the benefits of using sequenced biobanks for lcWGS imputation, we phased the  
25 recent release of the UK Biobank (UKB) WGS data<sup>3,4</sup> using SHAPEIT5<sup>5</sup> and created a UKB reference  
26 panel of 280,238 haplotypes and 582,534,516 markers (**Supplementary Note S1**). We used the UKB  
27 panel to impute lcWGS samples with GLIMPSE2 and other recently released imputation methods:  
28 GLIMPSE1<sup>1</sup> and QUILT v1.0.4<sup>2</sup>. Compared to other reference panels, the UKB leads to considerable  
29 accuracy improvements for British samples across all tested depths of coverage. Furthermore,

1 GLIMPSE2 outperforms GLIMPSE1, particularly at rare variants (MAF<0.1%) and for very low-coverage  
2 (0.1-0.5x) and matches QUILT v1.0.4 accuracy, designed to condition on the full set of reference  
3 haplotypes (**Figure 1a, Supplementary Note S2**). To consider non-British populations, we imputed 276  
4 lcWGS samples from the Simons Genome Diversity Project and we show that the UKB panel drastically  
5 improves imputation accuracy of European samples, in particular of Northern Europe origin  
6 (**Supplementary Note S3**). Additionally, we imputed three ancient Europeans and a Yamnaya sample  
7 for which high-coverage data (>18x) is available, and find similar improvements (**Supplementary Note**  
8 **S4**), showing that some ancient populations, such as Viking, Western Hunter-Gatherer and Yamnaya  
9 could be well imputed from the UKB reference panel.

10

11 The imputation of a single lcWGS genome using the UKB reference panel is expensive using existing  
12 methods. On the RAP, the cost is 1.11£ and 242.80£ for GLIMPSE1 and QUILT v1.0.4, respectively. In  
13 contrast, the same task performed with GLIMPSE2 only costs 0.08£, due to major algorithmic  
14 improvements that drastically reduce the imputation time for rare variants (**Fig 1b, Supplementary**  
15 **Note S2**). We confirm this trend for up to 2 million reference haplotypes, using simulated data. This  
16 keeps lcWGS close to SNP arrays in terms of computational costs<sup>6</sup> (**Supplementary Note S3**) while  
17 having the major advantage of providing better genotype calls. Indeed, we find that imputation of 0.5x  
18 data yields to more accurate results than the UKB Axiom array, specifically designed for the British  
19 population, with a notable difference at rare variants (for 0.5x coverage, accuracy improvement of  
20  $r^2 > 0.1$  for variants with a MAF < 0.01%, **Figure 1c**).

21

22 To assess the impact of these improvements on Genome-Wide Association Studies (GWAS), we  
23 imputed 10,000 UKB samples that we used to test 22 quantitative traits for association, comparing  
24 the respective abilities of lcWGS and SNP array data to recover the signals found with high-coverage  
25 sequencing data. We find that 0.25x leads to p-values and effect size estimates as accurate as those  
26 obtained from Axiom array data (**Figure 1d**) while delimiting regions of association with matching  
27 sensitivity and specificity (**Supplementary Note S6**). We also look at rare loss-of-function, missense  
28 and synonymous variants<sup>7</sup>, and show that 1.0x significantly outperforms the Axiom array in burden-  
29 test analysis (**Supplementary Note S7**). Altogether, this shows that lcWGS constitutes a powerful  
30 alternative to SNP array for downstream GWAS and rare variant analysis.

31

32 In this work, we introduce several major improvements for GLIMPSE that solve the computational  
33 problem of imputing lcWGS data from the 150,119 WGS samples in the UK Biobank. We demonstrate  
34 that this new reference panel leads to striking accuracy improvements across all allele frequencies,

1 sample ancestries, and depths of coverages. Our study further confirms the advantage of lcWGS over  
2 SNP arrays for GWAS, by showing that using imputed data with coverage as low as 0.5x is enough to  
3 outperform a SNP array specifically designed for the target population, particularly at rare variants.  
4 Our work can be applied to other sequenced and diverse biobanks, such as Trans-Omics for Precision  
5 Medicine<sup>8</sup>, gnomAD<sup>9</sup> or AllofUs<sup>10</sup>, thereby facilitating lcWGS imputation of non-European individuals.  
6 We believe that the difference between low-coverage and high-coverage WGS will become  
7 increasingly smaller as large reference panels will keep collecting more human haplotype diversity.

8

### 9 **Code availability**

10 GLIMPSE2 source code is available with MIT licence from <https://github.com/odelaneau/GLIMPSE> and  
11 <https://odelaneau.github.io/GLIMPSE/>. Pre-compiled binaries and docker images are available at  
12 <https://odelaneau.github.io/GLIMPSE/release>. Scripts to produce all figures of the paper are available  
13 on Github.

14

### 15 **Data availability**

16 The 1000 Genomes Project phase 3 dataset sequenced at high coverage by the New York Genome  
17 Center is available on the European Nucleotide Archive under accession no. PRJEB31736, the  
18 International Genome Sample Resource (IGSR) data portal and the University of Michigan school of  
19 public health ftp site (URL: <ftp://share.sph.umich.edu/1000g-high-coverage/freeze9/phased/>). The  
20 publicly available subset of the HRC dataset is available from the European Genome-phenome Archive  
21 at the European Bioinformatics Institute (EBI) under accession no. EGAS00001001710. The publicly  
22 available Simons Genome Diversity project is available on the IGSR data portal and Cancer Genomics  
23 Cloud, powered by Seven Bridges. The UK Biobank genetic and phenotypic data are available under  
24 restricted access. Access can be obtained by application via the UK Biobank Access Management  
25 System (URL: <https://www.ukbiobank.ac.uk/>).

26

### 27 **Acknowledgements**

28 This work was funded by a Swiss National Science Foundation project grant 373 (PP00P3\_176977) and  
29 conducted under UK Biobank project 66995. The New York Genome Center 1000 Genomes data were  
30 generated at the New York Genome Center with funds provided by a National Human Genome  
31 Research Institute grant no. 3UM1HG008901-03S1.

32

### 33 **Authors and Affiliations**

1 Department of Computational Biology, University of Lausanne, Lausanne, Switzerland and  
2 Swiss Institute of Bioinformatics, University of Lausanne, Lausanne, Switzerland  
3 Simone Rubinacci, Robin J. Hofmeister, Bárbara Sousa da Mota & Olivier Delaneau

4

5 **Contributions**

6 S.R. and O.D. designed the study. S.R. and O.D. developed the algorithms and wrote the software.

7 R.J.H. performed the GWAS experiments. S.R. and B.S.M performed imputation of ancient samples.

8 B.S.M. provided interpretation regarding imputed ancient samples. S.R. performed the remaining  
9 experiments. O.D. supervised the project. All authors reviewed the final manuscript.

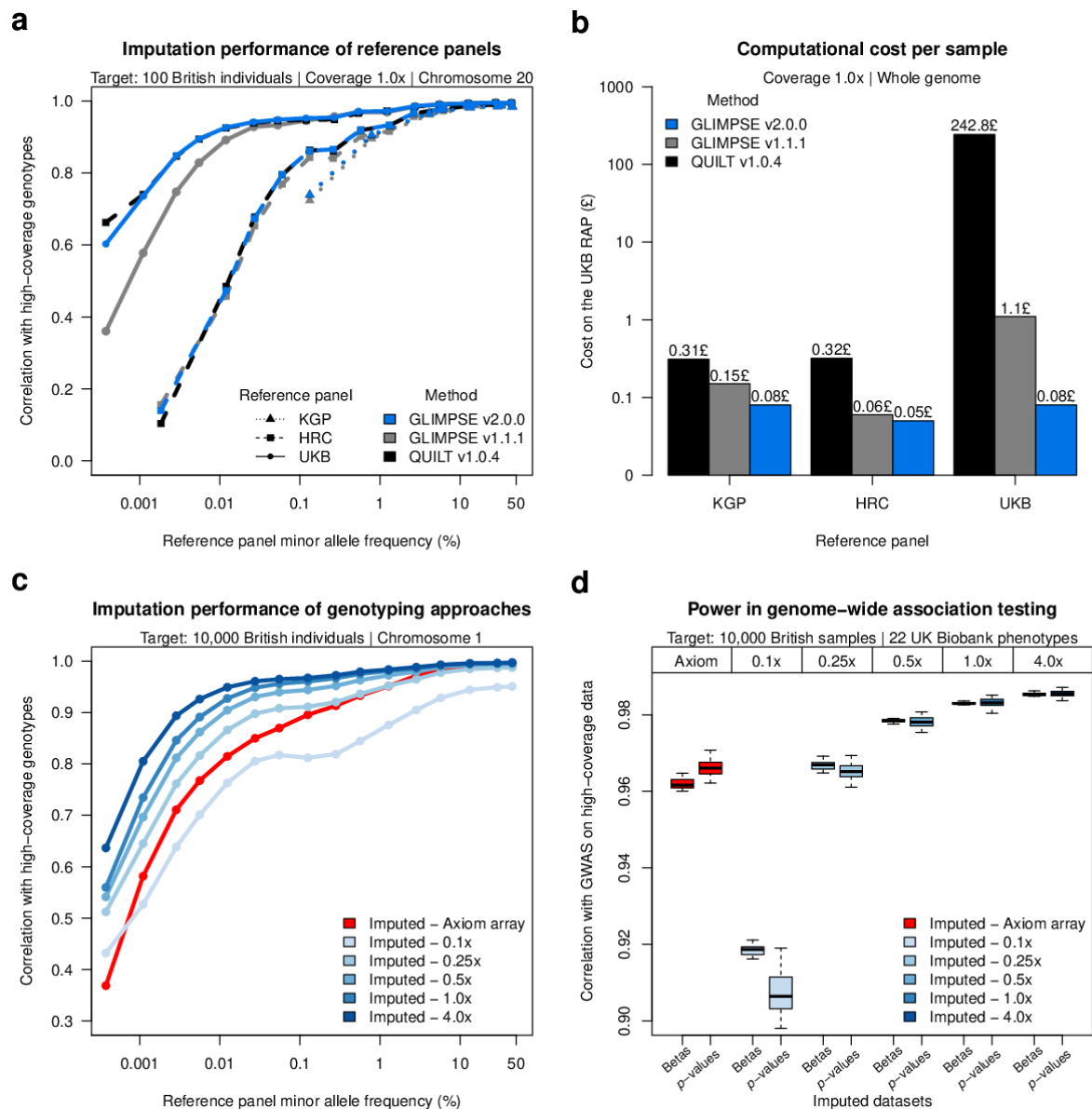
10

11 **Corresponding author**

12 Correspondence to Olivier Delaneau.

## References

1. Rubinacci, S., Ribeiro, D. M., Hofmeister, R. J. & Delaneau, O. Efficient phasing and imputation of low-coverage sequencing data using large reference panels. *Nat. Genet.* **53**, 120–126 (2021).
2. Davies, R. W. *et al.* Rapid genotype imputation from sequence with reference panels. *Nat. Genet.* **53**, 1104–1111 (2021).
3. Halldorsson, B. V. *et al.* The sequences of 150,119 genomes in the UK Biobank. *Nature* **607**, 732–740 (2022).
4. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).
5. Hofmeister, R. J., Ribeiro, D. M., Rubinacci, S. & Delaneau, O. Accurate rare variant phasing of whole-genome and whole-exome sequencing data in the UK Biobank. *bioRxiv* (2022) doi:10.1101/2022.10.19.512867.
6. Browning, B. L., Zhou, Y. & Browning, S. R. A One-Penny Imputed Genome from Next-Generation Reference Panels. *Am. J. Hum. Genet.* **103**, 338–348 (2018).
7. Karczewski, K. J. *et al.* Systematic single-variant and gene-based association testing of thousands of phenotypes in 394,841 UK Biobank exomes. *Cell Genomics* **2**, 100168 (2022).
8. Taliun, D. *et al.* Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nature* **590**, 290–299 (2021).
9. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **581**, 434–443 (2020).
10. All of Us Research Program Investigators *et al.* The ‘All of Us’ Research Program. *N. Engl. J. Med.* **381**, 668–676 (2019).



**Figure 1: Accuracy, running time and power of low-coverage imputation using the UK Biobank WGS data**

(a-b) Imputation performance of different imputation methods: QUILT v1.0.4 (black), GLIMPSE1 (grey) and GLIMPSE2 (blue); across different reference panels (KGP, HRC and UKB) for 100 UKB British samples at 1.0x coverage. (a) Accuracy on chromosome 20 (Pearson  $r^2$ , y-axis), of imputation methods and reference panels: KGP (triangle dotted line), HRC (squared dashed line) and UKB (full line). Accuracy is plotted against minor allele frequency of the appropriate reference panel (x-axis, log-scale). (b) Cost per sample on the RAP for whole-genome imputation (y-axis, log scale) across different reference panels (x-axis).

(c-d) Performance of imputed data using the UKB reference panel across coverages (0.1-4.0x, different shades of blue, GLIMPSE2 imputation), and Axiom array data (red). (c) Accuracy on chromosome 1 of 10,000 UKB British samples (Pearson  $r^2$ , y-axis) against minor allele frequency of the appropriate reference panel (x-axis, log-scale). (d) Power in GWAS in association testing of 10,000 UKB British samples compared to high-coverage data. Correlation of betas and p-values (y-axis) of different imputed datasets (x-axis) across 22 UKB phenotypes.