

1 **Analyzing bivariate cross-trait genetic architecture in GWAS summary statistics with**
2 **the BIGA cloud computing platform**

3

4 **Running title: BIGA GWAS cloud platform**

5

6 Yujue Li¹, Fei Xue¹, Bingxuan Li², Yilin Yang³, Zirui Fan⁴, Juan Shu¹, Xiyao Wang², Jinjie Lin⁵,
7 Carlos Copana¹, and Bingxin Zhao^{4*}

8

9 ¹Department of Statistics, Purdue University, West Lafayette, IN 47907, USA.

10 ²Department of Computer Science, Purdue University, West Lafayette, IN 47907, USA.

11 ³Department of Computer and Information Science and Electrical and Systems Engineering,
12 School of Engineering & Applied Science, University of Pennsylvania, Philadelphia, PA 19104, USA.

13 ⁴Department of Statistics and Data Science, University of Pennsylvania, Philadelphia, PA 19104,
14 USA.

15 ⁵Yale School of Management, Yale University, New Haven, CT 06511, USA.

16

17 **Corresponding to:*

18 Bingxin Zhao

19 413 Academic Research Building

20 265 South 37th Street, Philadelphia, PA 19104.

21 E-mail: bxzhao@upenn.edu Phone: (215) 898-8222

1 **Abstract**

2 As large-scale biobanks provide increasing access to deep phenotyping and genomic data,
3 genome-wide association studies (GWAS) are rapidly uncovering the genetic architecture
4 behind various complex traits and diseases. GWAS publications typically make their
5 summary-level data (GWAS summary statistics) publicly available, enabling further
6 exploration of genetic overlaps between phenotypes gathered from different studies and
7 cohorts. However, systematically analyzing high-dimensional GWAS summary statistics
8 for thousands of phenotypes can be both logistically challenging and computationally
9 demanding. In this paper, we introduce BIGA (<http://bigagwas.org/>), a website that offers
10 unified data analysis pipelines and centralized data resources for cross-trait genetic
11 architecture analyses using GWAS summary statistics. We have developed a framework
12 to implement statistical genetics tools on a cloud computing platform, combined with
13 extensive curated GWAS data resources. Through BIGA, users can upload data, submit
14 jobs, and share results, providing the research community with a convenient tool for
15 consolidating GWAS data and generating new insights.

16

17 **Keywords:** GWAS; Cross-trait analysis; Cloud computing; Online platform.

1 The rapid development of biobank-scale biomedical databases, encompassing
2 phenotyping and genomic data, has occurred globally¹. Numerous genome-wide
3 association studies (GWAS) have been conducted to determine the genetic architecture
4 underlying a wide range of complex traits and clinical outcomes, with the aim of
5 improving disease prevention and treatment². Publicly available GWAS summary-level
6 data (or GWAS summary statistics) encompass thousands of phenotypes³⁻⁸. These
7 summary statistics, derived from large-scale studies, provide valuable opportunities for
8 in-depth investigations into genetic overlaps and shared architectures between
9 phenotypes across studies and cohorts. Various statistical genetic tools have been
10 developed to analyze GWAS summary statistics and examine the shared genetic
11 components between pairs of phenotypes, including LDSC⁹, LAVA¹⁰, Coloc¹¹, and
12 Mendelian randomization¹². These methods offer insights into genetic links from various
13 perspectives and have been widely applied to clinical biomarkers and outcomes^{13,14}.

14

15 However, implementing and batch-running these tools often requires robust computing
16 and data infrastructure, which may not always be available to all researchers.
17 Consequently, systematic bivariate cross-trait analyses using massive GWAS summary
18 statistics for thousands of phenotypes can be logistically and computationally challenging.
19 As more complex and deep phenotyping data are obtained from biobanks¹⁵, addressing
20 these limitations becomes increasingly urgent. For example, the UK Biobank (UKB)
21 imaging study¹⁶ collected multimodal brain imaging data, generating over 5,000 imaging-
22 derived phenotypes using different imaging modalities and processing pipelines¹⁷⁻²¹.
23 Researchers interested in a specific disease and its genetic connections with imaging
24 biomarkers have traditionally downloaded all the GWAS summary statistics for over 5,000
25 imaging biomarkers from the Oxford BIG40 Project (<http://big.stats.ox.ac.uk>) and the BIG-
26 KP project (<https://bigkp.org/>), and run their statistical tools in local clusters, which can
27 be inefficient. Several online research platforms based on cloud computing have been
28 developed, most of which focus on one database (such as the UKB study,
29 <https://ukbiobank.dnanexus.com/>), univariate trait GWAS analysis (such as FUMA²²), or
30 single data analysis method/function (such as LD Hub²³ and Locus Compare²⁴). Developing
31 an integrated platform for cross-trait analyses of GWAS summary data resources will
32 make existing large-scale GWAS summary data more accessible to researchers.

1

2 To address these limitations, we developed BIGA (<http://bigagwas.org/>), an online cloud-
3 based platform that offers unified data analysis pipelines and centralized data resources
4 for cross-trait analyses using GWAS summary statistics. BIGA aims to provide various tools
5 for quantifying cross-trait genetic architectures, such as genome-wide genetic correlation
6 methods (e.g., LDSC⁹, Popcorn²⁵, and SumHer²⁶), polygenic overlap estimation between
7 two traits (e.g., MiXeR²⁷), genomic structural equation modeling (e.g., GenomicSEM²⁸),
8 local genetic correlation analysis (e.g., LAVA¹⁰), Mendelian randomization (e.g., MR-
9 Base²⁹, MR-RAPS³⁰, DIWV³¹, and GRAPPLE³²), and colocalization (e.g., Coloc¹¹). We have
10 also aggregated and preprocessed GWAS summary statistics from various resources (e.g.,
11 the UKB¹⁵, GWAS Atlas⁴, FinnGen⁶, Biobank Japan⁸, BIG-KP^{18,19,21}, and Oxford BIG40^{17,20})
12 and provided curated datasets. Our framework can easily be extended to incorporate
13 additional methods and GWAS summary statistics.

14

15 **Figure 1** provides an overview of the BIGA architecture. We offer users several options
16 for inputting GWAS summary statistics data, including uploading their own data, querying
17 harmonized data from public databases (such as the IEU OpenGWAS⁷ and Pan-UK
18 Biobank, <https://pan.ukbb.broadinstitute.org/>), and querying data from the BIGA in-
19 house database. Users can specify the tools and job types they are interested in and
20 submit their requests. After submission, the job request will be passed to the back-end
21 and executed on our cloud computing platform using the specified tools and datasets.
22 Once completed, the results will be presented to the user through the front-end interface.
23 Users have the option to share their results and/or data with our public database.

24

25 Here we present an LDSC data analysis example for schizophrenia³³. Specifically, we
26 investigated the genetic correlation between schizophrenia and complex traits and
27 diseases available in our current curated datasets. These datasets include 4,575 UKB
28 phenotypes provided by the Neale Lab (<http://www.nealelab.is/uk-biobank>), 3,095
29 clinical outcomes from FinnGen⁶, 199 brain-related phenotypes in GWAS Atlas⁴, 3,905
30 brain imaging phenotypes from the Oxford BIG40^{17,20}, and 3,016 multi-organ imaging
31 phenotypes from the BIG-KP^{18,19,21}. **Figure S1** displays the job information for these LDSC
32 data analyses, including Job ID, Job Name, Job Type, Status, Download Link, Result link,

1 and a link to share your results with the public. Once the job is complete, users can either
2 download the full set of results or view them in our online table. For example, among the
3 199 brain-related phenotypes in GWAS Atlas⁴, we discovered that the top-ranking
4 phenotypes exhibiting strong genetic correlations with schizophrenia included
5 neuroticism and its subclusters, daytime napping, and daytime sleepiness (**Fig. S2**). These
6 results can be effortlessly shared on our Public Results Portal (**Fig. S3**). This data analysis
7 demonstrates that our platform facilitates rapid analysis of extensive GWAS summary
8 statistics and enables seamless result sharing with the community.

9

10 In summary, our platform enables researchers to easily perform multiple cross-trait
11 analyses without needing access to a local research computing cluster, implementing
12 methods locally, or downloading large datasets. BIGA will help reduce the imbalance in
13 the research community caused by unequal computing resources. The source code to
14 build the BIGA platform will be made publicly available on GitHub, which can serve as an
15 example for creating more cloud computing-based genomic data analysis websites. The
16 BIGA website welcomes user feedback and requests, which aids in improving the project
17 and implementing new tools and functions to better meet the needs of the research
18 community.

19

20 **ADDITIONAL INFORMATION**

21 *One supplementary information pdf file is available, where the figures of an example data*
22 *analysis using BIGA are displayed.*

23

24 **ACKNOWLEDGEMENTS**

25 The study has been partially supported by funding from the Wharton Dean's Research
26 Fund and Analytics at Wharton, as well as start-up funds from Purdue Statistics
27 Department. We would like to thank the research computing team at the Wharton School
28 of the University of Pennsylvania and Purdue University for providing computational
29 resources and support that have contributed to this project. We would like to thank all
30 the developers of the tools and methods implemented in our project. We gratefully
31 acknowledge all the studies and databases that made GWAS summary data available and
32 thank the individuals who represented these studies for their participation and the

1 research teams for their work in collecting, processing, and disseminating these datasets
2 for analysis.

3

4 **AUTHOR CONTRIBUTIONS**

5 Y.L. and B.Z. designed the study, developed the BIGA website, and wrote the manuscript
6 with feedback from all authors. B.L. helped with the implementation of statistical genetic
7 methods and website functions. Y.Y., Z.F., J.S., X.W., B.L., and C.C. processed the GWAS
8 summary statistics, developed the curated datasets, and contributed to the development
9 of the website. F.X. and J.L. provided feedback on the study design and website.

10

11 **CORRESPONDENCE AND REQUESTS FOR MATERIALS** should be addressed to B.Z.

12

13 **COMPETING FINANCIAL INTERESTS**

14 The authors declare no competing financial interests.

15

16 **REFERENCES**

- 17 1. Zhou, W. *et al.* Global Biobank Meta-analysis Initiative: Powering genetic
18 discovery across human disease. *Cell Genomics* **2**, 100192 (2022).
- 19 2. Uffelmann, E. *et al.* Genome-wide association studies. *Nature Reviews Methods*
20 *Primers* **1**, 1-21 (2021).
- 21 3. Sollis, E. *et al.* The NHGRI-EBI GWAS Catalog: knowledgebase and deposition
22 resource. *Nucleic Acids Research* (2022).
- 23 4. Watanabe, K. *et al.* A global overview of pleiotropy and genetic architecture in
24 complex traits. *Nature genetics* **51**, 1339-1348 (2019).
- 25 5. Canela-Xandri, O., Rawlik, K. & Tenesa, A. An atlas of genetic associations in UK
26 Biobank. *Nature genetics* **50**, 1593-1599 (2018).
- 27 6. Kurki, M.I. *et al.* FinnGen provides genetic insights from a well-phenotyped
28 isolated population. *Nature* **613**, 508-518 (2023).
- 29 7. Elsworth, B. *et al.* The MRC IEU OpenGWAS data infrastructure. *BioRxiv* (2020).
- 30 8. Sakaue, S. *et al.* A cross-population atlas of genetic associations for 220 human
31 phenotypes. *Nature genetics* **53**, 1415-1424 (2021).

- 1 9. Bulik-Sullivan, B. *et al.* An atlas of genetic correlations across human diseases
2 and traits. *Nature Genetics* **47**, 1236-1241 (2015).
- 3 10. Werme, J., van der Sluis, S., Posthuma, D. & de Leeuw, C.A. An integrated
4 framework for local genetic correlation analysis. *Nature genetics* **54**, 274-282
5 (2022).
- 6 11. Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of
7 genetic association studies using summary statistics. *PLoS genetics* **10**, e1004383
8 (2014).
- 9 12. Burgess, S., Scott, R.A., Timpson, N.J., Davey Smith, G. & Thompson, S.G. Using
10 published data in Mendelian randomization: a blueprint for efficient
11 identification of causal risk factors. *European journal of epidemiology* **30**, 543-
12 552 (2015).
- 13 13. Romero, C. *et al.* Exploring the genetic overlap between twelve psychiatric
14 disorders. *Nature Genetics*, 1-8 (2022).
- 15 14. Lee, P.H. *et al.* Genomic relationships, novel loci, and pleiotropic mechanisms
16 across eight psychiatric disorders. *Cell* **179**, 1469-1482. e11 (2019).
- 17 15. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic
18 data. *Nature* **562**, 203-209 (2018).
- 19 16. Littlejohns, T.J. *et al.* The UK Biobank imaging enhancement of 100,000
20 participants: rationale, data collection, management and future directions.
21 *Nature communications* **11**, 1-12 (2020).
- 22 17. Elliott, L.T. *et al.* Genome-wide association studies of brain imaging phenotypes
23 in UK Biobank. *Nature* **562**, 210-216 (2018).
- 24 18. Zhao, B. *et al.* Common genetic variation influencing human white matter
25 microstructure. *Science* **372**, eabf3736 (2021).
- 26 19. Zhao, B. *et al.* Genome-wide association analysis of 19,629 individuals identifies
27 variants influencing regional brain volumes and refines their genetic co-
28 architecture with cognitive and mental health traits. *Nature genetics* **51**, 1637-
29 1644 (2019).
- 30 20. Smith, S.M. *et al.* An expanded set of genome-wide association studies of brain
31 imaging phenotypes in UK Biobank. *Nature neuroscience* **24**, 737-745 (2021).

- 1 21. Zhao, B. *et al.* Common variants contribute to intrinsic human brain functional
2 networks. *Nature Genetics* **54**, 508-517 (2022).
- 3 22. Watanabe, K., Taskesen, E., Bochoven, A. & Posthuma, D. Functional mapping
4 and annotation of genetic associations with FUMA. *Nature Communications* **8**,
5 1826 (2017).
- 6 23. Zheng, J. *et al.* LD Hub: a centralized database and web interface to perform LD
7 score regression that maximizes the potential of summary level GWAS data for
8 SNP heritability and genetic correlation analysis. *Bioinformatics* **33**, 272-279
9 (2017).
- 10 24. Liu, B., Gludemans, M.J., Rao, A.S., Ingelsson, E. & Montgomery, S.B. Abundant
11 associations with gene expression complicate GWAS follow-up. *Nature genetics*
12 **51**, 768-769 (2019).
- 13 25. Brown, B.C., Ye, C.J., Price, A.L., Zaitlen, N. & Consortium, A.G.E.N.T.D.
14 Transethnic genetic-correlation estimates from summary statistics. *The American*
15 *Journal of Human Genetics* **99**, 76-88 (2016).
- 16 26. Speed, D. & Balding, D.J. SumHer better estimates the SNP heritability of
17 complex traits from summary statistics. *Nature genetics* **51**, 277-284 (2019).
- 18 27. Frei, O. *et al.* Bivariate causal mixture model quantifies polygenic overlap
19 between complex traits beyond genetic correlation. *Nature communications* **10**,
20 1-11 (2019).
- 21 28. Grotzinger, A.D. *et al.* Genomic structural equation modelling provides insights
22 into the multivariate genetic architecture of complex traits. *Nature human*
23 *behaviour* **3**, 513-525 (2019).
- 24 29. Hemani, G. *et al.* The MR-Base platform supports systematic causal inference
25 across the human phenome. *elife* **7**(2018).
- 26 30. Zhao, Q., Wang, J., Hemani, G., Bowden, J. & Small, D.S. Statistical inference in
27 two-sample summary-data Mendelian randomization using robust adjusted
28 profile score. *The Annals of Statistics* **48**, 1742-1769 (2020).
- 29 31. Ye, T., Shao, J. & Kang, H. Debiased inverse-variance weighted estimator in two-
30 sample summary-data Mendelian randomization. *The Annals of statistics* **49**,
31 2079-2100 (2021).

- 1 32. Wang, J. *et al.* Causal inference for heritable phenotypic risk factors using
2 heterogeneous genetic instruments. *PLoS genetics* **17**, e1009575 (2021).
3 33. Trubetskov, V. *et al.* Mapping genomic loci implicates genes and synaptic biology
4 in schizophrenia. *Nature* **604**, 502-508 (2022).

5

6 **Code availability**

7 All source code to develop the BIGA platform will be made publicly available on the BIGA
8 GitHub repository. The statistical tools and methods implemented in the BIGA platform
9 are also open source, and their source code has already been made available to the public
10 by their authors.

11

12 **Data availability**

13 GWAS summary statistics used in the BIGA platform are publicly available and can be
14 found in several public databases, such as the
15 Pan-UK Biobank (<https://pan.ukbb.broadinstitute.org/>),
16 IEU OpenGWAS (<https://gwas.mrcieu.ac.uk/>),
17 GWAS Atlas (<https://atlas.ctglab.nl/>),
18 FinnGen (https://www.finngen.fi/en/access_results),
19 Biobank Japan (<https://pheweb.jp/>),
20 BIG-KP (<https://bigkp.org/>), and
21 Oxford BIG40 (<https://open.win.ox.ac.uk/ukbiobank/big40/>).

22

23 **Figure Legend**

24 **Fig. 1 Overview of BIGA GWAS cloud computing platform.**

25 **(A)** The motivation of this project is to address the substantial logistical and
26 computational challenges associated with implementing and batch-running the
27 constantly evolving tools for cross-trait genetic architecture analysis. Our aim is to offer a
28 cloud computing-based solution that can effectively overcome these challenges. **(B)**
29 Overview of the BIGA GWAS platform. Users can easily upload their GWAS summary level
30 data and submit data analysis jobs through the front-end interface. These jobs are then
31 processed on the back-end, and the results are subsequently returned to the users. In

1 addition, we have established public data and result centers where users can share their
2 data and outcomes with the public. **(C)** The front-end interface of the BIGA GWAS
3 platform offers users a comprehensive set of options to manage their data resources,
4 choose the appropriate tools, and select the desired mode of data analysis. **(D)** Details of
5 the back-end of the BIGA GWAS platform. **(E)** Overview of the analysis workflow.

