

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 Li<sup>1</sup>, Fei Xue<sup>1</sup>, Bingxuan Li<sup>2</sup>, Yilin Yang<sup>3</sup>, Zirui Fan<sup>4</sup>, Juan Shu<sup>1</sup>, Xiaochen Yang<sup>1</sup>, Xiyao  
7 Wang<sup>2</sup>, Jinjie Lin<sup>5</sup>, Carlos Copana<sup>1</sup>, and Bingxin Zhao<sup>4,6-10\*</sup>

8

9 <sup>1</sup>Department of Statistics, Purdue University, West Lafayette, IN 47907, USA.

10 <sup>2</sup>Department of Computer Science, Purdue University, West Lafayette, IN 47907, USA.

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

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

15 <sup>5</sup>Yale School of Management, Yale University, New Haven, CT 06511, USA.

16 <sup>6</sup>Applied Mathematics and Computational Science Graduate Group, University of Pennsylvania,  
17 Philadelphia, PA 19104, USA.

18 <sup>7</sup>Center for AI and Data Science for Integrated Diagnostics, Perelman School of Medicine,  
19 University of Pennsylvania, Philadelphia, PA 19104, USA.

20 <sup>8</sup>Penn Institute for Biomedical Informatics, Perelman School of Medicine, University of  
21 Pennsylvania, Philadelphia, PA 19104, USA.

22 <sup>9</sup>Population Aging Research Center, University of Pennsylvania, Philadelphia, PA 19104, USA.

23 <sup>10</sup>Institute for Translational Medicine and Therapeutics, University of Pennsylvania, Philadelphia,  
24 PA 19104, USA.

25

26 *\*Corresponding to:*

27 Bingxin Zhao

28 413 Academic Research Building

29 265 South 37th Street, Philadelphia, PA 19104.

30 E-mail: [bxzhao@upenn.edu](mailto: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 (<https://bigagwas.org/>), a website that aims  
10 to offer unified data analysis pipelines and processed data resources for cross-trait  
11 genetic architecture analyses using GWAS summary statistics. We have developed a  
12 framework to implement statistical genetics tools on a cloud computing platform,  
13 combined with extensive curated GWAS data resources. Through BIGA, users can upload  
14 data, submit jobs, and share results, providing the research community with a convenient  
15 tool for 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 globally<sup>1</sup>. 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 treatment<sup>2</sup>. Publicly available GWAS summary-level  
6 data (or GWAS summary statistics) encompass thousands of phenotypes<sup>3-8</sup>. 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, such as LDSC<sup>9</sup>, LAVA<sup>10</sup>, SumHer<sup>11</sup>, and  
12 Popcorn<sup>12</sup>. These methods offer insights into genetic links from various perspectives and  
13 have been widely applied to clinical biomarkers and outcomes<sup>13,14</sup>.

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 biobanks<sup>15</sup>, addressing  
20 these limitations becomes increasingly urgent. For example, the UK Biobank (UKB)  
21 imaging study<sup>16</sup> collected multimodal brain imaging data, generating over 5,000 imaging-  
22 derived phenotypes using different imaging modalities and processing pipelines<sup>17-21</sup>.  
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. Such challenges are also present in centralized GWAS databases, such as  
28 GWAS Catalog<sup>3</sup> and IEU OpenGWAS<sup>7</sup>, where users are expected to download and manage  
29 large datasets locally to conduct most analyses. Several online research platforms based  
30 on cloud computing have been developed, most of which focus on one database (such as  
31 the UKB study, <https://ukbiobank.dnanexus.com/>), univariate trait GWAS analysis (such  
32 as FUMA<sup>22</sup>), or single data analysis method/function (such as LD Hub<sup>23</sup> and Locus

1 Compare<sup>24</sup>). Developing an integrated platform for cross-trait analyses of GWAS summary  
2 data resources will make existing large-scale GWAS summary data more accessible to  
3 researchers.

4  
5 To address these limitations, we developed BIGA (<https://bigagwas.org/>), an online cloud-  
6 based platform that offers unified data harmonization and analysis pipelines and  
7 processed data resources for cross-trait analyses using GWAS summary statistics. BIGA  
8 aims to provide various tools for quantifying cross-trait genetic architectures, such as  
9 genome-wide genetic correlation methods (e.g., LDSC<sup>9</sup>, Popcorn<sup>12</sup>, and SumHer<sup>11</sup>) and  
10 local genetic correlation analysis (e.g., LAVA<sup>10</sup>). We have also aggregated and harmonized  
11 GWAS summary statistics from various resources, including the GWAS Catalog<sup>3</sup>, UKB  
12 study<sup>15</sup>, Psychiatric Genomics Consortium<sup>25</sup>, FinnGen<sup>6</sup>, Biobank Japan<sup>8</sup>, CHIMGEN<sup>26</sup>, UKB-  
13 PPP<sup>27</sup>, BIG-KP<sup>18,19,21</sup>, and Oxford BIG40<sup>17,20</sup>. These curated datasets, currently including  
14 over 15,000 traits, have been integrated with multiple methods, facilitating easy online  
15 analysis for users. With our established infrastructure in place, we are committed to the  
16 continuous development and growth of BIGA, aiming to broaden its capabilities by  
17 consistently including new tools and data resources.

18  
19 **Figure 1** provides an overview of the BIGA architecture. We offer users several options  
20 for inputting GWAS summary statistics data with user-friendly features, including  
21 uploading their own data, querying data from public databases (such as the IEU  
22 OpenGWAS<sup>7</sup>, GWAS Catalog<sup>3</sup>, and Neale Lab (<http://www.nealelab.is/uk-biobank>), and  
23 reusing data from recent previous jobs (Supplementary Text). Users can specify the tools  
24 and job types they are interested in and submit their requests. After submission, the job  
25 request will be passed to the back-end and executed on our cloud computing platform  
26 using the specified tools and datasets. Briefly, we have developed a thorough pipeline for  
27 harmonizing user-input data, similar to procedures used in the GWAS Catalog  
28 (<https://github.com/EBISPOT/gwas-sumstats-harmoniser>). After harmonization, datasets  
29 will have a standard format with column names outlined in **Table S1**. Considering the  
30 specific data format needed by the user-requested analysis, we will accordingly adapt the  
31 data to fulfill these requirements and execute the analysis (**Fig. S1**) Once completed, users  
32 will receive email notification and the results will be presented to the users through the

1 front-end interface. A quick-start tutorial and comprehensive documentation are  
2 available on our website for users.

3

4 BIGA uses a powerful and efficient computational framework for automated analysis.  
5 Every step, from the initial data input to the final results output, is organized by a  
6 standardized pipeline, offering the flexibility to incorporate new methods. For example,  
7 BIGA operates on the Django 3.2 web framework (<https://www.djangoproject.com/>) to  
8 accommodate various tasks and tools, and we use Redis (<https://redis.io/>) and Celery  
9 (<https://docs.celeryq.dev>) for task management and queuing system. BIGA's  
10 computational infrastructure is efficient, currently supporting 20 concurrent user jobs  
11 running with just 128GB of RAM and 16 Intel vCPUs. Notably, cloud computing services  
12 provide a flexible management system for CPU and RAM, enabling us to easily modify our  
13 resource allocation for scaling up or down as needed. Even with only 16GB of RAM, BIGA  
14 can execute 3 jobs concurrently using our efficient configuration. We have conducted  
15 large-scale tests to validate the stability and computational efficiency of BIGA (**Figs. S2-3**  
16 and Supplementary Text).

17

18 To showcase the extensive genetic analyses that BIGA can conduct, we present a blood  
19 pressure data analysis example, aiming to explore its genetic correlation with over 15,000  
20 complex traits and diseases curated on BIGA. We initiated the analysis by searching for  
21 blood pressure data on the IEU OpenGWAS database and used the BIGA query function  
22 to directly query systolic blood pressure<sup>28</sup> summary statistics. BIGA performed  
23 harmonization and then used the harmonized data to run LDSC massive analysis, spanning  
24 over all groups of traits from European population on BIGA. As expected, at a false  
25 discovery rate 5% level, systolic blood pressure was widely associated with complex traits  
26 and diseases, such as hypertension, atrial fibrillation, stroke, brain and body imaging  
27 traits, as well as plasma proteomics ( $P$  range =  $(5.44 \times 10^{-244}, 4.00 \times 10^{-2})$ , **Fig. S4**). We further  
28 examined the diastolic blood pressure<sup>28</sup> and found similar association patterns to systolic  
29 blood pressure (**Fig. S5**). We applied SumHer to repeat the analysis (**Fig. S6**) and observed  
30 that the results from LDSC and SumHer were generally consistent (**Fig. S7**, Pearson's  
31 correlation = 0.9273). In addition, we performed local genetic correlation analysis using  
32 LAVA and cross-population genetic correlation using Popcorn. More details can be found

1 in the Supplementary Text (**Tables S2-S8**). This data analysis example demonstrates that  
2 BIGA facilitates efficient analysis of extensive GWAS summary statistics with different  
3 methods.

4  
5 In summary, our platform enables researchers to easily perform multiple cross-trait  
6 analyses without needing access to a local research computing cluster, implementing  
7 methods locally, or downloading large datasets. BIGA will help reduce the imbalance in  
8 the research community caused by unequal computing resources and attract a wider user  
9 base to these developed methods. The source code to build the BIGA platform will be  
10 made publicly available on GitHub. The BIGA website welcomes user feedback and  
11 requests, which aids in improving the project and implementing new tools and functions  
12 to better meet the needs of the research community.

13

#### 14 **ADDITIONAL INFORMATION**

15 *One supplementary pdf file and one supplementary table zip file are available.*

16

#### 17 **ACKNOWLEDGEMENTS**

18 We thank for the helpful conversations with Doug Speed regarding SumHer. Research  
19 reported in this publication was supported by the National Institute On Aging of the  
20 National Institutes of Health under Award Number RF1AG082938. The content is solely  
21 the responsibility of the authors and does not necessarily represent the official views of  
22 the National Institutes of Health. The study has also been partially supported by funding  
23 from the Wharton Dean's Research Fund, Analytics at Wharton, NSF Grant DMS 2210860,  
24 and Purdue Statistics Department. We would like to thank the research computing groups  
25 at Purdue University, the Wharton School of the University of Pennsylvania, and the  
26 University of North Carolina at Chapel Hill for providing computational resources and  
27 support that have contributed to these research results. We would like to thank all the  
28 developers of the tools and methods implemented in our project. We gratefully  
29 acknowledge all the studies and databases that made GWAS summary data available and  
30 thank the individuals who represented these studies for their participation and the  
31 research teams for their work in collecting, processing, and disseminating these datasets  
32 for analysis.

1

## 2 **AUTHOR CONTRIBUTIONS**

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

9

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

11

## 12 **COMPETING FINANCIAL INTERESTS**

13 The authors declare no competing financial interests.

14

## 15 **REFERENCES**

- 16 1. Zhou, W. *et al.* Global Biobank Meta-analysis Initiative: Powering genetic  
17 discovery across human disease. *Cell Genomics* **2**, 100192 (2022).
- 18 2. Uffelmann, E. *et al.* Genome-wide association studies. *Nature Reviews Methods*  
19 *Primers* **1**, 1-21 (2021).
- 20 3. Sollis, E. *et al.* The NHGRI-EBI GWAS Catalog: knowledgebase and deposition  
21 resource. *Nucleic Acids Research* (2022).
- 22 4. Watanabe, K. *et al.* A global overview of pleiotropy and genetic architecture in  
23 complex traits. *Nature genetics* **51**, 1339-1348 (2019).
- 24 5. Canela-Xandri, O., Rawlik, K. & Tenesa, A. An atlas of genetic associations in UK  
25 Biobank. *Nature genetics* **50**, 1593-1599 (2018).
- 26 6. Kurki, M.I. *et al.* FinnGen provides genetic insights from a well-phenotyped  
27 isolated population. *Nature* **613**, 508-518 (2023).
- 28 7. Elsworth, B. *et al.* The MRC IEU OpenGWAS data infrastructure. *BioRxiv* (2020).
- 29 8. Sakaue, S. *et al.* A cross-population atlas of genetic associations for 220 human  
30 phenotypes. *Nature genetics* **53**, 1415-1424 (2021).
- 31 9. Bulik-Sullivan, B. *et al.* An atlas of genetic correlations across human diseases  
32 and traits. *Nature Genetics* **47**, 1236-1241 (2015).

- 1 10. Werme, J., van der Sluis, S., Posthuma, D. & de Leeuw, C.A. An integrated  
2 framework for local genetic correlation analysis. *Nature genetics* **54**, 274-282  
3 (2022).
- 4 11. Speed, D. & Balding, D.J. SumHer better estimates the SNP heritability of  
5 complex traits from summary statistics. *Nature Genetics* **51**, 277-284 (2019).
- 6 12. Brown, B.C., Ye, C.J., Price, A.L., Zaitlen, N. & Consortium, A.G.E.N.T.D.  
7 Transethnic genetic-correlation estimates from summary statistics. *The American*  
8 *Journal of Human Genetics* **99**, 76-88 (2016).
- 9 13. Romero, C. *et al.* Exploring the genetic overlap between twelve psychiatric  
10 disorders. *Nature Genetics*, 1-8 (2022).
- 11 14. Lee, P.H. *et al.* Genomic relationships, novel loci, and pleiotropic mechanisms  
12 across eight psychiatric disorders. *Cell* **179**, 1469-1482. e11 (2019).
- 13 15. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic  
14 data. *Nature* **562**, 203-209 (2018).
- 15 16. Littlejohns, T.J. *et al.* The UK Biobank imaging enhancement of 100,000  
16 participants: rationale, data collection, management and future directions.  
17 *Nature communications* **11**, 1-12 (2020).
- 18 17. Elliott, L.T. *et al.* Genome-wide association studies of brain imaging phenotypes  
19 in UK Biobank. *Nature* **562**, 210-216 (2018).
- 20 18. Zhao, B. *et al.* Common genetic variation influencing human white matter  
21 microstructure. *Science* **372**, eabf3736 (2021).
- 22 19. Zhao, B. *et al.* Genome-wide association analysis of 19,629 individuals identifies  
23 variants influencing regional brain volumes and refines their genetic co-  
24 architecture with cognitive and mental health traits. *Nature genetics* **51**, 1637-  
25 1644 (2019).
- 26 20. Smith, S.M. *et al.* An expanded set of genome-wide association studies of brain  
27 imaging phenotypes in UK Biobank. *Nature neuroscience* **24**, 737-745 (2021).
- 28 21. Zhao, B. *et al.* Common variants contribute to intrinsic human brain functional  
29 networks. *Nature Genetics* **54**, 508-517 (2022).
- 30 22. Watanabe, K., Taskesen, E., Bochoven, A. & Posthuma, D. Functional mapping  
31 and annotation of genetic associations with FUMA. *Nature Communications* **8**,  
32 1826 (2017).



- 1 23. Zheng, J. *et al.* LD Hub: a centralized database and web interface to perform LD  
2 score regression that maximizes the potential of summary level GWAS data for  
3 SNP heritability and genetic correlation analysis. *Bioinformatics* **33**, 272-279  
4 (2017).
- 5 24. Liu, B., Gludemans, M.J., Rao, A.S., Ingelsson, E. & Montgomery, S.B. Abundant  
6 associations with gene expression complicate GWAS follow-up. *Nature genetics*  
7 **51**, 768-769 (2019).
- 8 25. Sullivan, P. A framework for interpreting genome-wide association studies of  
9 psychiatric disorders. *Molecular psychiatry* (2009).
- 10 26. Xu, Q. *et al.* CHIMGEN: a Chinese imaging genetics cohort to enhance cross-  
11 ethnic and cross-geographic brain research. *Molecular Psychiatry* **25**, 517-529  
12 (2020).
- 13 27. Sun, B.B. *et al.* Plasma proteomic associations with genetics and health in the UK  
14 Biobank. *Nature* **622**, 329-338 (2023).
- 15 28. Evangelou, E. *et al.* Genetic analysis of over 1 million people identifies 535 new  
16 loci associated with blood pressure traits. *Nature Genetics* **50**, 1412-1425 (2018).

17

### 18 **Code availability**

19 All source code to develop the BIGA platform will be made publicly available on the BIGA  
20 GitHub repository. The statistical tools and methods implemented in the BIGA platform  
21 are also open source, and their source code has already been made available to the public  
22 by their authors. A summary of our implemented tools and data resources can be found  
23 at <https://bigagwas.org/documentation>.

24

### 25 **Data availability**

26 GWAS summary statistics used in the BIGA platform are publicly available and can be  
27 found in several public databases, such as the  
28 Neale Lab UK Biobank Results (<http://www.nealelab.is/uk-biobank>),  
29 Psychiatric Genomics Consortium (<https://pgc.unc.edu/>),  
30 IEU OpenGWAS (<https://gwas.mrcieu.ac.uk/>),  
31 FinnGen ([https://www.finngen.fi/en/access\\_results](https://www.finngen.fi/en/access_results)),  
32 Biobank Japan (<https://pheweb.jp/>),

- 1 BIG-KP (<https://bigkp.org/>),
- 2 Oxford BIG40 (<https://open.win.ox.ac.uk/ukbiobank/big40/>),
- 3 UKB-PPP (<https://metabolomics.org.uk/bbpgwas/>),
- 4 GWAS Catalog (<https://www.ebi.ac.uk/gwas/downloads/summary-statistics>), and
- 5 CHIMGEN (<http://chimgen.tmu.edu.cn/en/index.php?c=article&id=2036>).

6

## 7 **Figure Legend**

### 8 **Fig. 1 Overview of BIGA GWAS cloud computing platform.**

9 **(A)** The motivation of this project is to address the substantial logistical and  
10 computational challenges associated with implementing and batch-running the  
11 constantly evolving tools for cross-trait genetic architecture analysis. Our aim is to offer a  
12 cloud computing-based solution that can effectively overcome these challenges. **(B)**  
13 Overview of the BIGA GWAS platform. Users can easily upload or query GWAS summary  
14 statistics and submit data analysis jobs through the front-end interface. These jobs are  
15 then processed on the back-end, and the results are subsequently returned to the users.  
16 **(C)** The front-end interface of the BIGA GWAS platform offers users a comprehensive set  
17 of options to manage their data resources, choose the appropriate tools, and select the  
18 desired mode of data analysis. **(D)** Details of the back-end of the BIGA GWAS platform. **(E)**  
19 Overview of the analysis workflow.

