



# Data Availability of Open T-Cell Receptor Repertoire Data, a Systematic Assessment

Yu-Ning Huang<sup>1†</sup>, Naresh Amrat Patel<sup>2</sup>, Jay Himanshu Mehta<sup>2</sup>, Srishti Ginjala<sup>3†</sup>, Petter Brodin<sup>4,5†</sup>, Clive M. Gray<sup>6†</sup>, Yesha M. Patel<sup>1†</sup>, Lindsay G. Cowell<sup>7,8,9†</sup>, Amanda M. Burkhardt<sup>1†</sup> and Sergei Mangul<sup>1\*†</sup>

## OPEN ACCESS

### Edited by:

Harinder Singh,  
University of Pittsburgh, United States

### Reviewed by:

Anna Bernasconi,  
Politecnico di Milano, Italy

### \*Correspondence:

Sergei Mangul  
serghei.mangul@gmail.com

### †ORCID:

Yu-Ning Huang  
[orcid.org/0000-0003-1697-4267](https://orcid.org/0000-0003-1697-4267)  
Srishti Ginjala  
[orcid.org/0000-0001-7269-2382](https://orcid.org/0000-0001-7269-2382)  
Petter Brodin  
[orcid.org/0000-0002-8103-0046](https://orcid.org/0000-0002-8103-0046)  
Clive M. Gray  
[orcid.org/0000-0002-9293-901X](https://orcid.org/0000-0002-9293-901X)  
Yesha M. Patel  
[orcid.org/0000-0002-4983-1222](https://orcid.org/0000-0002-4983-1222)  
Lindsay G. Cowell  
[orcid.org/0000-0003-1617-8244](https://orcid.org/0000-0003-1617-8244)  
Amanda M. Burkhardt  
[orcid.org/0000-0002-7326-474X](https://orcid.org/0000-0002-7326-474X)  
Sergei Mangul  
[orcid.org/0000-0003-4770-3443](https://orcid.org/0000-0003-4770-3443)

### Specialty section:

This article was submitted to  
Integrative Systems Immunology,  
a section of the journal  
Frontiers in Systems Biology

Received: 12 April 2022

Accepted: 09 May 2022

Published: 06 June 2022

### Citation:

Huang Y-N, Patel NA, Mehta JH,  
Ginjala S, Brodin P, Gray CM,  
Patel YM, Cowell LG, Burkhardt AM  
and Mangul S (2022) Data Availability  
of Open T-Cell Receptor Repertoire  
Data, a Systematic Assessment.  
*Front. Syst. Biol.* 2:918792.  
doi: 10.3389/fsysb.2022.918792

<sup>1</sup>Department of Clinical Pharmacy, School of Pharmacy, University of Southern CA, Los Angeles, CA, United States, <sup>2</sup>Department of Pharmaceutical Sciences, School of Pharmacy, University of Southern CA, Los Angeles, CA, United States, <sup>3</sup>School of Computing and Electrical Engineering, Indian Institute of Technology, Mandi, India, <sup>4</sup>Department of Immunology and Inflammation, Faculty of Medicine, Imperial College London, London, United Kingdom, <sup>5</sup>Department of Women's and Children's Health, Karolinska Institute, Stockholm, Sweden, <sup>6</sup>Division of Molecular Biology and Human Genetics, Faculty of Medicine and Health Sciences, Stellenbosch University, Cape Town, South Africa, <sup>7</sup>Division of Biomedical Informatics, University of Texas Southwestern Medical Center at Dallas, Dallas, NC, United States, <sup>8</sup>Department of Immunology, Medical School, University of Texas Southwestern Medical Center, Dallas, NC, United States, <sup>9</sup>Department of Population and Data Sciences, University of Texas Southwestern Medical Center, Dallas, NC, United States

Modern data-driven research has the power to promote novel biomedical discoveries through secondary analyses of raw data. Therefore, it is important to ensure data-driven research with great reproducibility and robustness for promoting a precise and accurate secondary analysis of the immunogenomics data. In scientific research, rigorous conduct in designing and conducting experiments is needed, specifically in scientific writing and reporting results. It is also crucial to make raw data available, discoverable, and well described or annotated in order to promote future re-analysis of the data. In order to assess the data availability of published T cell receptor (TCR) repertoire data, we examined 11,918 TCR-Seq samples corresponding to 134 TCR-Seq studies ranging from 2006 to 2022. Among the 134 studies, only 38.1% had publicly available raw TCR-Seq data shared in public repositories. We also found a statistically significant association between the presence of data availability statements and the increase in raw data availability ( $p = 0.014$ ). Yet, 46.8% of studies with data availability statements failed to share the raw TCR-Seq data. There is a pressing need for the biomedical community to increase awareness of the importance of promoting raw data availability in scientific research and take immediate action to improve its raw data availability enabling cost-effective secondary analysis of existing immunogenomics data by the larger scientific community.

**Keywords:** data availability, T cell receptor, raw data, immunogenetics and immunogenomics, data analytics

## INTRODUCTION

Advanced high throughput sequencing technologies have reshaped the landscape of contemporary immunology research providing researchers with a rich set of tools and methods to study fundamental aspects of immune responses across a variety of disciplines. Raw TCR-Seq data is produced accompanying the development and progress of immunogenomics research. Availability of raw TCR-Seq data allows effective re-analysis of the data (Brito et al., 2020) (also known as

secondary analysis) which in turn may accelerate novel biomedical discoveries (Johnston, 2017). Using public repositories to share raw data substantially simplifies discovering and retrieving datasets of interest, as one has scalable, programmatic access to a large number of studies. Over the last decade, there has been tremendous progress to improve sharing of raw immunogenomics data, which allows researchers to easily obtain the various types of data stored in public repositories (Breden et al., 2017) (e.g. SRA (Kodama et al., 2012)). But to unlock the full potential of publicly available raw TCR-Seq data for the secondary analyses, it is crucial for a study to be conducted accurately with reproducible and reliable laboratory practices and accurate results (Miyakawa, 2020). Moreover, the research should report publicly available raw data accompanying metadata for secondary analysis (Rajesh et al., 2021; Rubelt et al., 2017). One of the merits of science is to be able to leverage published data for subsequent novel scientific discovery (Wilkinson et al., 2016). To ensure the results are reproducible and have sufficient quality for accurate and precise secondary analysis, the data should be shared under the FAIR principle—Findability, Accessibility, Interoperability, and Reusability—to ensure the quality of the raw data to be reused accurately and efficiently in research (Wilkinson et al., 2016; Wass et al., 2019). Scientific communities, peer-reviewed journals, funding agencies, and government agencies have emphasized the importance of making raw omics data publicly available with free access to the public (Rubelt et al., 2017). As such, establishing protocols, guidelines, and authors' checklists for reporting the raw immunogenomics data in research, and peer-review journals mandating authors to upload the data to public repositories are potential means to increase the raw data availability (Breden et al., 2017; Rubelt et al., 2017; Bishop, 2019).

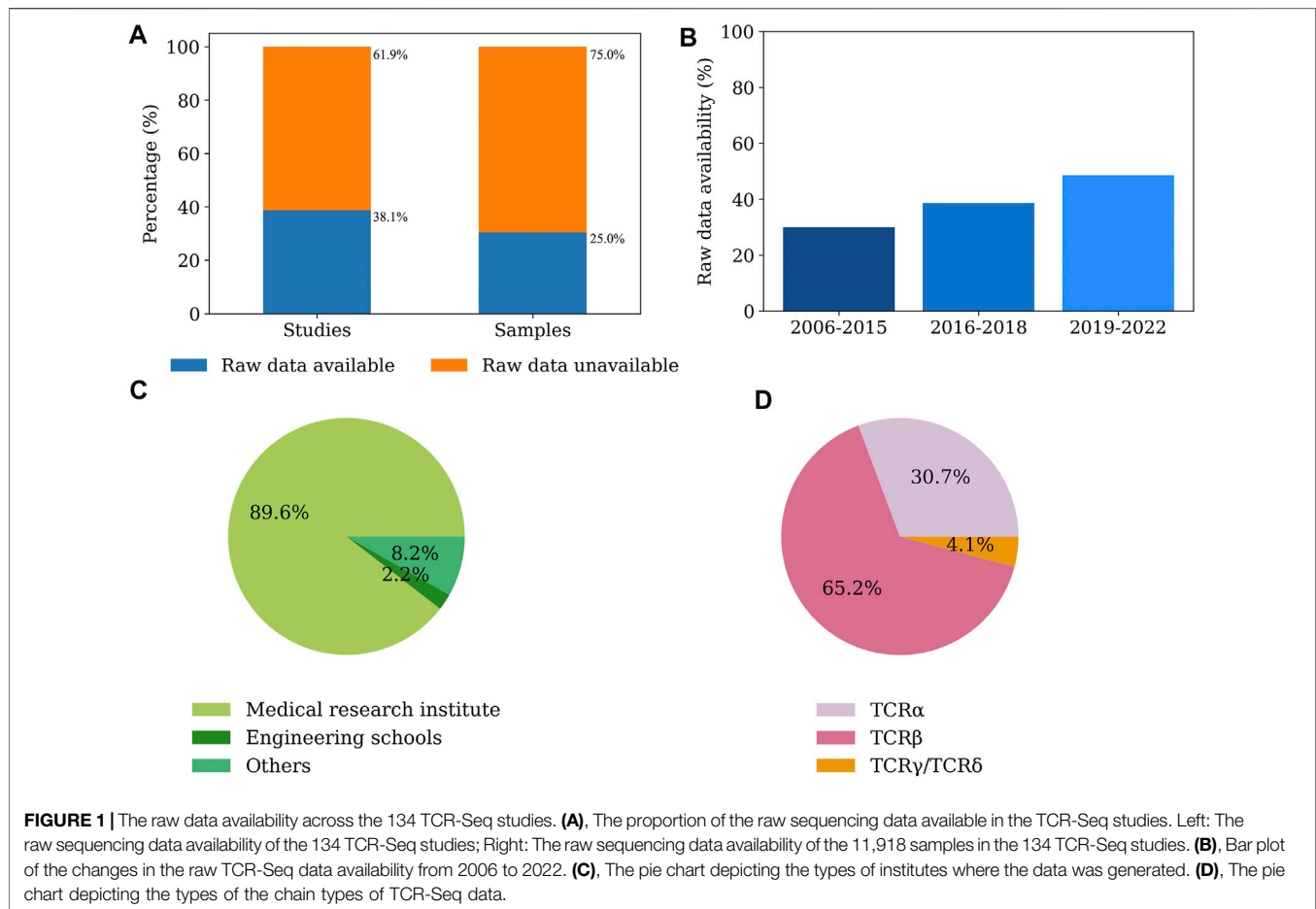
It is also crucial to ensure studies share the raw data used for publication and analyses. In the fields of immunogenetics studies since the development of next-generation sequencing tools promotes research on genetic sequencing and production of raw immunogenomics data due to the reduction in price for genetic sequencing (van Dijk et al., 2014). T cell receptor sequencing (TCR-Seq) studies allow researchers to profile human immune systems providing updated insights into T cell receptors. TCR-Seq data allows researchers to examine a unique individual's immune status, immune responses, and compare the T cell populations between individuals in healthy or disease states, such as autoimmune diseases, infectious diseases, and cancer (Benichou et al., 2012; Dziubianau et al., 2013; Hou et al., 2016). Additionally, TCR-Seq studies allow the development of novel therapeutics and biomarkers, including diagnostics (Arnaout et al., 2021) for autoimmune disease (Ostmeyer et al., 2017) and cancer (Linette et al., 2019; Cowell, 2020; Ostmeyer et al., 2020), CAR-T cell therapy (Sheih et al., 2020), vaccines (Lee et al., 2016), and monoclonal and therapeutic antibodies (Richardson et al., 2021). Due to the production of vast TCR-Seq data in the field of TCR repertoire studies, there is a pressing need to ensure that raw sequencing data from TCR-Seq studies is provided with full access to the public to facilitate further computational and bioinformatics research in the field.

In this study, we examined the raw data availability across 134 TCR-Seq studies among the 1,195 existing studies. We also examined the public genomic repositories where the researchers stored the generated TCR-Seq data. We found, only 38.1% of TCR-Seq studies have made the overall raw data available, and that the majority of studies (35 studies) with available raw data stored the raw data in Sequence Read Archive (Kodama et al., 2012). Further, 61.9% of the TCR-Seq studies did not share the raw data in the original publications or that the raw data will only be available upon request. We discussed the potential barriers researchers are facing that deter them from sharing the raw data and the potential ways to improve the availability of the raw data in the field of TCR-Seq studies.

## The Availability of Raw Sequencing Data in TCR-Seq Studies Is Limited

We investigated 134 published TCR-Seq studies across 11,918 samples in PubMed for the raw sequencing data availability ranging from 2006 to 2022. The studies were considered as having available raw sequencing data if we could acquire the samples of the studies' raw FASTQ or FASTA files. According to our results, only 38.1% (51 out of 134 studies) of the TCR-Seq studies shared raw TCR-Seq data in the original publications at public genomic repositories (Figure 1A). Conversely, 61.9% (83 out of 134 studies) of the TCR-Seq studies have unavailable raw RNA-Seq data or have raw data available upon request (Figure 1A). We observed a similar trend of raw data availability among the 11,918 samples of the 134 TCR-Seq studies, in which 25% of the samples among the TCR-Seq studies had available raw data (Figure 1A). We also observed that the raw TCR-Seq data availability has increased over the past decade (Figure 1B). Among the 134 studies, 89.6% (120 studies) of the raw data were generated in medical research institutes, including medical schools, hospitals, medical centers, private health/disease research institutes, and government health research institutes, such as the National Institutes of Health (Figure 1C). The remaining raw data were generated in engineering schools and other research institutes (Figure 1C). Additionally, we examined the types of TCR-Seq chains in the available raw sequencing data. Among the samples with available raw sequencing data, 65.2% of the TCR-Seq data were TCR beta (TCR $\beta$ ) chain, 30.7% of the TCR-Seq data were TCR alpha (TCR $\alpha$ ) chain, and 4.1% of the TCR-Seq data were TCR gamma (TCR $\gamma$ ) or TCR delta (TCR $\delta$ ) chain (Figure 1D).

We further investigated the specific reasons that the raw data was unavailable among the 83 studies without available raw TCR-Seq data. Among the 83 studies with unavailable raw data, rather than sharing raw sequencing data, 44 studies only shared summary data, such as summary data on ImmuneACCESS (immunoSEQ<sup>®</sup>)<sup>®</sup> (Adaptive Biotechnologies), VDJDdb (Shugay et al., 2018), and supplementary files. Other 34 studies did not include statements about the raw data or provide raw sequencing data in the original publications and five studies indicated that the raw sequencing data will be available upon making direct requests to the authors (Figure 2A).



## The Platforms Used to Store the Raw TCR-Seq Data

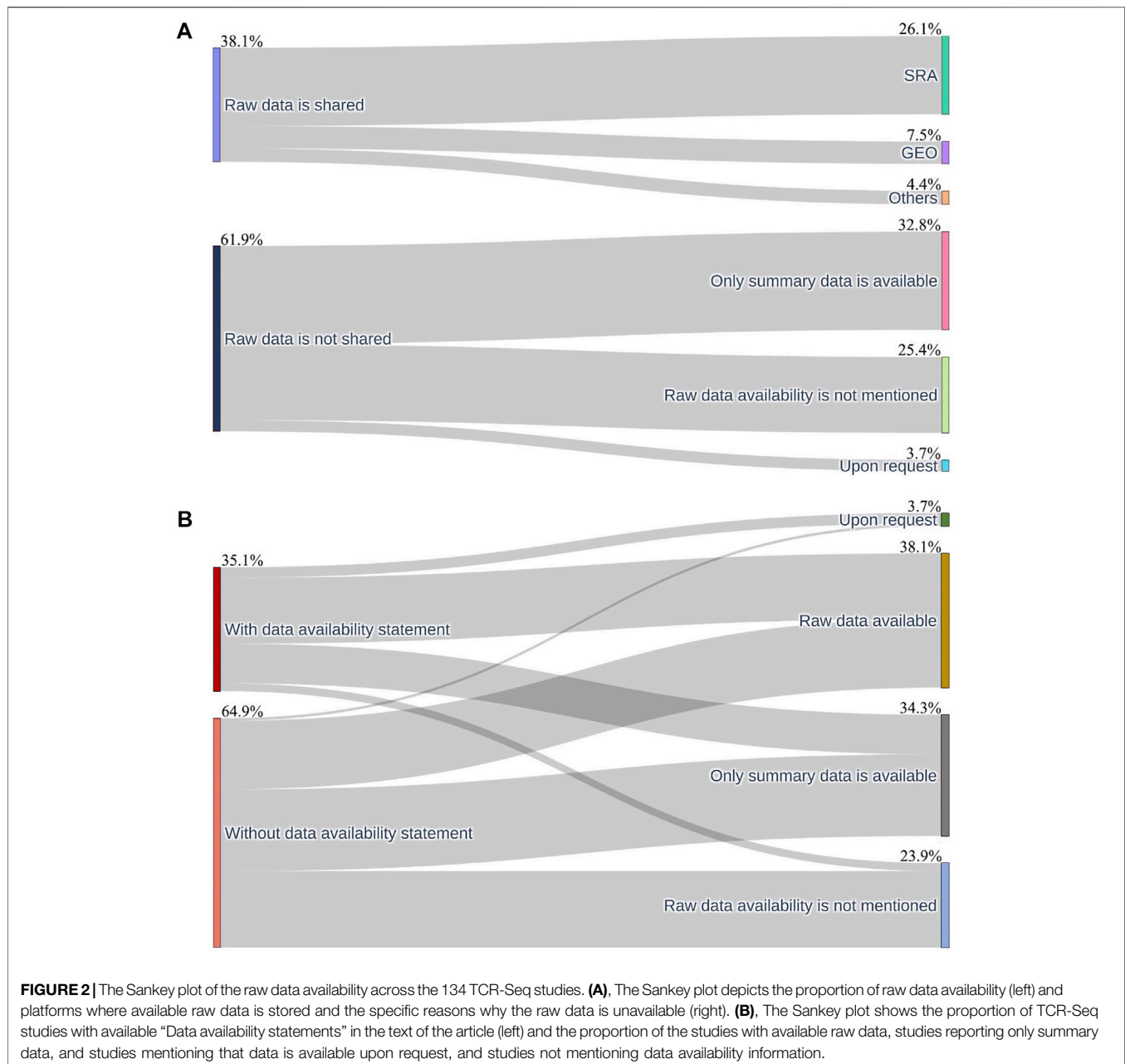
We further examined the platforms researchers used to share raw TCR-Seq data. Among the TCR-Seq studies that shared the raw sequencing data, 35 of the 51 studies (68.6%) of the raw data are shared in Sequence Read Archive (SRA) (Kodama et al., 2012) and ten of the 51 studies (19.6%) shared the raw data in the Gene Expression Omnibus (GEO) (Barrett et al., 2013). The remaining studies (11.8%) shared the raw sequencing data in various online repositories, including European Genome-phenome Archive (Freeberg et al., 2022) (EGA) (3 studies), VDJSerVer (Christley et al., 2018) (1 study), and National Genomics Data Center, China (CNCB-NGDC Members and Partners, 2021) (NGDC) (2 studies) (Figure 2A).

## The Presence of the “Data Availability Statement” Increases the Availability of Raw TCR-Seq Data

As part of promoting more transparent and reproducible research, many journals do require data availability statements for studies to be published in the journals (Rajesh et al., 2021), (Stodden et al., 2018; Schriml et al., 2020; Gozashti and Corbett-

Detig, 2021; Tedersoo et al., 2021). We examined the impact of the presence of the “Data availability statement” in research articles on the availability of raw TCR-Seq data. According to the results, the presence of data availability statements improves the raw data availability from 29.9% (26 out of 87 studies without data availability statements) to 53.2% (25 out of 47 studies with data availability statements), and increase the raw data availability by 23.3%. (Supplementary Table S1).

There are 47 studies containing “Data availability statements” in the corresponding publications. Among the 47 studies with “Data availability statements”, 53.2% (25 out of the 47 studies) studies shared the raw TCR-Seq data in the publications while the rest of the studies did not share raw TCR-Seq data in the original publications or would be available upon request. Conversely, there are 87 studies that did not have “Data availability statements” in the original publications. For the 87 publications that do not have the “Data availability statement”, only 29.9% (26 out of the 87 studies) of the studies shared the raw TCR-Seq data in the publications. There are three studies provided data availability statement in the article but with incorrect accession number or mentioning that will share the raw data while we cannot access the raw data of the studies. Therefore, the three studies are categorized into the category “With data availability statement” and “Raw data availability is



not mentioned”. We conducted a Pearson’s chi-squared test ( $\chi^2$ ) to examine the associations between the presence of data availability statements and the availability of raw data and found a statistically significant relationship between the two ( $p = 0.014$ ) - the data availability statements are improving access to raw data availability.

We also examined parts of the publications where the raw data availability is mentioned. Among the studies that shared the raw TCR-Seq data, authors included the data availability statements in various locations of the research articles. 49.0% (25 out of 51 studies) had the data availability statements in the “Data availability statement” portion of the research articles. Eight out of 51 studies (15.7%) had the data availability statements

in the footnote of the studies, whereas the remaining 18 out of 51 studies (35.3%) mentioned their raw data availability directly in the main text of the studies.

## DISCUSSION

We produced the first study to assess the raw data availability of TCR-Seq data. According to the results, the raw data availability is alarmingly low, with only 38.1% of the TCR-Seq studies sharing the raw data (**Figure 2B**). We investigated the association between the presence of data availability statements in the articles and the raw data availability. We discovered that the

data availability statements presented have a statistically significant association with the raw data availability ( $p = 0.014$ ). The presence of data availability statements improves the raw data availability by 23.3%. While data availability statements help to improve data availability, we noticed that it is still possible to have such a statement and not share raw data or only share summary data or only will be available upon request (46.8%) (**Supplementary Table S1**). Some studies only shared the summary data in the data availability statement section of the publications. Additionally, three of the studies shared an erroneous SRA accession number in the data availability statement sections of the publications so we were not able to access the raw data of the articles.

According to our analysis, the primary reason for unavailable raw data is that the studies did not indicate the statements of raw data in the research articles nor shared the raw data in the article (34 studies), which did not share the raw sequencing data of the analyses. Secondly, most authors only shared summary data in the original publications (44 studies) reckoning that they shared the data in their research. For example, 23 studies shared the summary data in the corporate-owned repository, ImmuneACCESS<sup>®</sup> by Adaptive Biotechnologies. Unfortunately, the ImmuneACCESS<sup>®</sup> repository only offers access to the summary data of the studies, thus the raw data files are still inaccessible to the public, making the researchers unable to re-analyze the raw data generated by Adaptive Biotechnologies for novel biomedical discoveries. Thirdly, several studies mentioned that the data will only be available by making direct requests of the data from the authors (5 studies). However, it has been previously shown that the statements mentioning that the raw data will be available upon request did not guarantee data availability and instead were not a sustainable and practical way to improve research reproducibility and raw data availability (Stodden et al., 2018; Tedersoo et al., 2021). The specific reasons for not sharing raw TCR-Seq data are beyond the scope of this manuscript and need to be investigated in future studies. The perceptual and technical barriers researchers are facing when sharing the data are yet to be determined. Previous work has suggested that the reasons for authors to not share the raw data might be because they are unaware of the importance of sharing the raw data, there may be culturally related or technical barriers that deter or prevent authors from sharing the raw data (Tedersoo et al., 2021).

Individual researchers, research institutes, and journals should all take part in ensuring raw data availability (Tedersoo et al., 2021). Individual researchers and research institutes more specifically should hold the burden of making the raw data publicly available. Journals also have a critical role in promoting raw data availability, and many journals are already taking actions to promote data availability via data availability statement requests within articles upon

publication. However, more stringent measures need to be imposed by the journals to ensure raw data availability (Grant and Hrynaszkiewicz, 2018). It is known that many journals have already mandated the authors to share the raw data of the studies (Kim et al., 2020). Journal's policies in mandating authors to share the raw data might be a feasible way to improve the raw data availability (Deshpande et al., 2021). In conclusion, the pressing need to increase awareness of enhancing raw data availability in scientific research can enable cost-effective secondary analysis of existing immunogenomics data for novel biomedical discoveries. In addition to immunogenomics data, metadata (Rajesh et al., 2021), (Field et al., 2008), raw sequencing data (Caspar. et al., 2018), and open human health data (Peters and Zeeb, 2022), such as medical history, should also be made publicly available for secondary analysis. Therefore, we recommend that all members of the biomedical community, including individual researchers, research institutions, and journals, should contribute to increasing the awareness of raw data sharing and improve the raw data availability in future studies.

## AUTHOR CONTRIBUTIONS

Y-NH analyzed the data and wrote the manuscript with input from all the authors. NP, JM, and SG collected the data. All authors discussed the text and commented on the manuscript. All authors read and approved the final manuscript. SM conceived and supervised the study.

## FUNDING

SM is supported by the National Science Foundation grants 2041984 and 2135954.

## ACKNOWLEDGMENTS

Our paper is dedicated to all freedom-loving people around the world, and to the people of Ukraine who fight for our freedom.

## SUPPLEMENTARY MATERIAL

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fsysb.2022.918792/full#supplementary-material>

## REFERENCES

- Arnaout, R. A., Prak, E. T. L., Schwab, N., Rubelt, F., and Immune, A. (2021). The Future of Blood Testing Is the Immunome. *Front. Immunol.* 12, 626793. doi:10.3389/fimmu.2021.626793
- Barrett, T., Wilhite, S. E., Ledoux, P., Evangelista, C., Kim, I. F., Tomashevsky, M., et al. (2013). NCBI GEO: Archive for Functional Genomics Data Sets—Update. *Nucleic Acids Res.* 41, D991–D995. doi:10.1093/nar/gks1193
- Benichou, J., Ben-Hamo, R., Louzoun, Y., and Efroni, S. (2012). Rep-Seq: Uncovering the Immunological Repertoire through Next-Generation Sequencing. *Immunology* 135, 183–191. doi:10.1111/j.1365-2567.2011.03527.x



- Bishop, M. (2019). *Building the Foundation for Future Research through Open Data, Code and Protocols*. The Official PLOS Blog Available at: <https://theplosblog.plos.org/2019/12/building-the-foundation-for-future-research-through-open-data-code-and-protocols/>.
- Breden, F., Luning Prak, E. T., Peters, B., Rubelt, F., Schramm, C. A., Busse, C. E., et al. (2017). Reproducibility and Reuse of Adaptive Immune Receptor Repertoire Data. *Front. Immunol.* 8, 1418. doi:10.3389/fimmu.2017.01418
- Brito, J. J., Li, J., Moore, J. H., Greene, C. S., Nogoy, N. A., Garmire, L. X., et al. (2020). Recommendations to Enhance Rigor and Reproducibility in Biomedical Research. *GigaScience*, 9, g1aa056. doi:10.1093/gigascience/g1aa056
- Caspar, S. M., Dubacher, N., Kopps, A. M., Meienberg, J., Henggeler, C., and Matyas, G. (2018). Clinical Sequencing: From Raw Data to Diagnosis with Lifetime Value. *Clin. Genet.* 93, 508–519. doi:10.1111/cge.13190
- Christley, S., Scarborough, W., Salinas, E., Rounds, W. H., Toby, I. T., Fonner, J. M., et al. (2018). VDJServer: A Cloud-Based Analysis Portal and Data Commons for Immune Repertoire Sequences and Rearrangements. *Front. Immunol.* 9, 976. doi:10.3389/fimmu.2018.00976
- CNCB-NGDC Members and Partners (2021). Database Resources of the National Genomics Data Center, China National Center for Bioinformatics in 2021. *Nucleic Acids Res.* 49, D18–D28. doi:10.1093/nar/gkaa1022
- Cowell, L. G. (2020). The Diagnostic, Prognostic, and Therapeutic Potential of Adaptive Immune Receptor Repertoire Profiling in Cancer. *Cancer Res.* 80, 643–654. doi:10.1158/0008-5472.can-19-1457
- Deshpande, D., Sarkar, A., Guo, R., Moore, A., Darci-Maher, N., and Mangul, S. (2021). A Comprehensive Analysis of Code and Data Availability in Biomedical Research. *Mapping Intimacies* 67. doi:10.31219/osf.io/uz7m5
- Dziubianau, M., Hecht, J., Kuchenbecker, L., Sattler, A., Stervbo, U., Rödelsperger, C., et al. (2013). TCR Repertoire Analysis by Next Generation Sequencing Allows Complex Differential Diagnosis of T Cell-Related Pathology. *Am. J. Transpl. Off. J. Am. Soc. Transpl. Am. Soc. Transpl. Surg.* 13, 2842–2854. doi:10.1111/ajt.12431
- Field, D., Garrity, G., Gray, T., Morrison, N., Selengut, J., Sterk, P., et al. (2008). The Minimum Information about a Genome Sequence (MIGS) Specification. *Nat. Biotechnol.* 26, 541–547. doi:10.1038/nbt1360
- Freeberg, M. A., Fromont, L. A., D'Altri, T., Romero, A. F., Ciges, J. I., Jene, A., et al. (2022). The European Genome-Phenome Archive in 2021. *Nucleic Acids Res.* 50, D980–D987. doi:10.1093/nar/gkab1059
- Gozashti, L., and Corbett-Detig, R. (2021). Shortcomings of SARS-CoV-2 Genomic Metadata. *BMC Res. Notes* 14, 189. doi:10.1186/s13104-021-05605-9
- Grant, R., and Hrynaskiewicz, I. (2018). The Impact on Authors and Editors of Introducing Data Availability Statements at Nature Journals. *Int. J. Digit. Curation* 13, 195–203. doi:10.2218/ijdc.v13i1.614
- Hou, D., Chen, C., Seely, E. J., Chen, S., and Song, Y. (2016). High-Throughput Sequencing-Based Immune Repertoire Study during Infectious Disease. *Front. Immunol.* 7, 336. doi:10.3389/fimmu.2016.00336
- immunoSEQ® | *The Gold Standard of Immunosequencing*. Available at: <https://www.immunoseq.com/>.
- Johnston, M. P. (2017). Secondary Data Analysis: A Method of Which the Time Has Come. *Qual. Quant. Methods Libr.* 3, 619–626.
- Kim, J., Kim, S., Cho, H.-M., Chang, J. H., and Kim, S. Y. (2020). Data Sharing Policies of Journals in Life, Health, and Physical Sciences Indexed in Journal Citation Reports. *PeerJ* 8, e9924. doi:10.7717/peerj.9924
- Kodama, Y., Shumway, M., and Leinonen, R. (2012). The Sequence Read Archive: Explosive Growth of Sequencing Data. *Nucleic Acids Res.* 40, D54–D56. doi:10.1093/nar/gkr854
- Lee, J., Boutz, D. R., Chromikova, V., Joyce, M. G., Vollmers, C., Leung, K., et al. (2016). Molecular-level Analysis of the Serum Antibody Repertoire in Young Adults before and after Seasonal Influenza Vaccination. *Nat. Med.* 22, 1456–1464. doi:10.1038/nm.4224
- Linette, G. P., Becker-Hapak, M., Skidmore, Z. L., Baroja, M. L., Xu, C., Hundal, J., et al. (2019). Immunological Ignorance Is an Enabling Feature of the Oligo-Clonal T Cell Response to Melanoma Neoantigens. *Proc. Natl. Acad. Sci. U.S.A.* 116, 23662–23670. doi:10.1073/pnas.1906026116
- Miyakawa, T. (2020). No Raw Data, No Science: Another Possible Source of the Reproducibility Crisis. *Mol. Brain* 13, 24. doi:10.1186/s13041-020-0552-2
- Ostmeyer, J., Lucas, E., Christley, S., Lea, J., Monson, N., Tiro, J., et al. (2017). Statistical Classifiers for Diagnosing Disease from Immune Repertoires: a Case Study Using Multiple Sclerosis. *BMC Bioinforma.* 18, 401. doi:10.1186/s12859-017-1814-6
- Ostmeyer, J., Lucas, E., Christley, S., Lea, J., Monson, N., Tiro, J., et al. (2020). Biophysicochemical Motifs in T Cell Receptor Sequences as a Potential Biomarker for High-Grade Serous Ovarian Carcinoma. *PLoS One*, 15, e0229569. doi:10.1371/journal.pone.0229569
- Peters, M., and Zeeb, H. (2022). Availability of Open Data for Spatial Public Health Research. *GMS Ger. Med. Sci.*, 20, Doc01.
- Rajesh, A., Chang, Y., Abedalthagafi, M. S., Wong-Beringer, A., Love, M. I., and Mangul, S. (2021). Improving the Completeness of Public Metadata Accompanying Omics Studies. *Genome Biol.* 22, 106. doi:10.1186/s13059-021-02332-z
- Richardson, E., Galson, J. D., Kellam, P., Kelly, D. F., Smith, S. E., Palser, A., et al. (2021). A Computational Method for Immune Repertoire Mining that Identifies Novel Binders from Different Clonotypes, Demonstrated by Identifying Anti-pertussis Toxoid Antibodies. *mAbs* 13, 1869406. doi:10.1080/19420862.2020.1869406
- Rubelt, F., Busse, C. E., Bukhari, S. A. C., Bürckert, J. P., Ferrandiz, E. M., and Cowell, L. G. (2017). Adaptive Immune Receptor Repertoire Community Recommendations for Sharing Immune-Repertoire Sequencing Data. *Nat. Immunol.* 18, 1274–1278. doi:10.1038/ni.3873
- Schriml, L. M., Chuvochina, M., Davies, N., Eloie-Fadrosch, E. A., Finn, R. D., Hugenholtz, P., et al. (2020). COVID-19 Pandemic Reveals the Peril of Ignoring Metadata Standards. *Sci. Data* 7, 188. doi:10.1038/s41597-020-0524-5
- Sheih, A., Voillet, V., Hanafi, L.-A., DeBerg, H. A., Yajima, M., Hawkins, R., et al. (2020). Clonal Kinetics and Single-Cell Transcriptional Profiling of CAR-T Cells in Patients Undergoing CD19 CAR-T Immunotherapy. *Nat. Commun.* 11, 219. doi:10.1038/s41467-019-13880-1
- Shugay, M., Bagaev, D. V., Zvyagin, I. V., Vroomans, R. M., Crawford, J. C., Dolton, G., et al. (2018). VDJdb: a Curated Database of T-Cell Receptor Sequences with Known Antigen Specificity. *Nucleic Acids Res.* 46, D419–D427. doi:10.1093/nar/gkx760
- Stodden, V., Seiler, J., and Ma, Z. (2018). An Empirical Analysis of Journal Policy Effectiveness for Computational Reproducibility. *Proc. Natl. Acad. Sci. U.S.A.* 115, 2584–2589. doi:10.1073/pnas.1708290115
- Tedersoo, L., Küngas, R., Oras, E., Köster, K., Eenmaa, H., Leijen, Ä., et al. (2021). Data Sharing Practices and Data Availability upon Request Differ across Scientific Disciplines. *Sci. Data* 8, 192. doi:10.1038/s41597-021-00981-0
- van Dijk, E. L., Jaszczyszyn, Y., Thermes, C., and Thermes, C. (2014). Ten Years of Next-Generation Sequencing Technology. *Trends Genet.* 30, 418–426. doi:10.1016/j.tig.2014.07.001
- Wass, M. N., Ray, L., and Michaelis, M. (2019). Understanding of Researcher Behavior Is Required to Improve Data Reliability. *GigaScience*, 8(5), g1z017. doi:10.1093/gigascience/g1z017
- Wilkinson, M. D., Dumontier, M., Aalbersberg, I. J., Appleton, G., Axton, M., Baak, A., et al. (2016). The FAIR Guiding Principles for Scientific Data Management and Stewardship. *Sci. Data* 3, 160018. doi:10.1038/sdata.2016.18

**Conflict of Interest:** The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

**Publisher's Note:** All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

Copyright © 2022 Huang, Patel, Mehta, Ginjala, Brodin, Gray, Patel, Cowell, Burkhardt and Mangul. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.