



**IEEE
COMPUTER
SOCIETY**

Guardians 2023 Conference

Proceedings of the 2nd Annual Guardians of Truth and Integrity Conference

IEEE CONFERENCE: RECORD # 59421, XPLORE # [10473739](#)

CONFERENCE DATE: 9 OCTOBER 2023

FINANCIAL SPONSOR: BRAIN HEALTH ALLIANCE

TECHNICAL SPONSORS: IEEE COMPUTER SOCIETY, IEEE TECHNICAL COMMUNITY ON MULTIMEDIA COMPUTING, IEEE TECHNICAL COMMUNITY ON SEMANTIC COMPUTING

© 2023 BRAIN HEALTH ALLIANCE (a 501c3 not-for-profit)

ISBN: 979-8-9900823-0-4

URL: WWW.BRAINIACSJOURNAL.ORG/ARC/PUB/GUARDIANS2023

PDP: NPDS.BRAINHEALTHALLIANCE.NET/NEXUS/BRAINIACSGUARDIANS2023

DOI: [10.48085/IBEBA9475](https://doi.org/10.48085/IBEBA9475)

Contents

| | |
|---------------------------------------------------------------------------------------------|-----|
| Guardians 2023 Program | 2 |
| Guardians 2023 Contributors | 3 |
| BHAVI 2023 Guardian: Anthony S. Fauci, MD (slides) | 4 |
| Fallacies and Pitfalls in Genome-Wide Association Studies (review) | 15 |
| Photoshop Fantasies (slides) | 23 |
| Academic Ghosting: Towards an Academy of Truth Telling (slides) | 64 |
| An Extended Active Learning Approach to Multiverse Analysis (report) | 71 |
| BabbleNewt: A Reference Citation Format for Bibliographic Metadata (report) | 85 |
| Managing Lexical-Semantic Hybrid Records of FAIR Metrics Analyses (report) | 90 |
| Reproducibility, Validity, and Integrity in Scholarly Research (commentary) | 99 |
| Who are the Guardians of Truth and Integrity? (summary) | 102 |

Guardians 2023 Program

Guardians 2023 was held on October 9th as a half-day online event with 3 invited speakers:

- Dr. Nan Laird, Harvard University, Boston MA
- Dr. Walter Scheirer, University of Notre Dame, Notre Dame IN
- Dr. Alicia Andrzejewski, William & Mary, Williamsburg VA

and a tribute to Dr. Anthony Fauci honoring him as our 2023 Guardian of Truth and Integrity.

Opening Remarks

- 09:00 Julie Neidich, BHA VI 2023 Guardian: Anthony S. Fauci (2023 Guardian [slides](#) and [video](#))

Invited Talks

- 09:15 Julian Hecker and Nan Laird, Fallacies and Pitfalls in Genome-Wide Association Studies ([JH slides](#), [NL slides](#), [JH+NL video](#))
- 10:15 Walter Scheirer, Photoshop Fantasies: Why is there so much fake stuff on the Internet? ([WS slides](#), [WS video](#))
- 11:15 Alicia Andrzejewski, Academic Ghosting: Towards an Academy of Truth-Telling ([AA slides](#), [AA video](#))

Technical Talks

- 12:30 Daniel Kristanto, Multiverse in Functional Magnetic Resonance Imaging Analysis ([DK slides](#), [DK video](#))
- 13:00 Koby Taswell, Consistent Bibliographic Data Formats with the BabbleNewt Project ([KT slides](#), [KT video](#))
- 13:30 Adam Craig, Managing Lexical-Semantic Hybrid Records of FAIR Metrics Analyses with the NPDS Cyberinfrastructure ([AC slides](#), [AC video](#))

Closing Remarks

- 14:00 Carl Taswell, Reproducibility, Reliability, and Integrity in Scholarly Research: What Accountability for Willful Disregard? ([CT slides](#), [CT video](#))

All slides and recordings of the talks are available from

- Guardians.BHAVI.us/Conf2023/Program
- www.BHAVI.us/Symposia/202310

Guardians 2023 Contributors

Andrzejewski, Alicia
Craig, Adam
Debener, Stefan
Hecker, Julian
Hildebrandt, Andrea
Hughes, Andrew
Gießing, Carsten
Kristanto, Daniel
Laird, Nan
Marek, Merle
Neidich, Julie
Scheirer, Walter
Taswell, Carl
Taswell, S. Koby
Thiel, Christiane
Zhou, Changsong

BHAVI 2023 Guardian: Anthony S. Fauci, MD

BHAVI Awards Committee

BHAVI 2023 Guardian: Anthony S. Fauci, MD

BHAVI Awards Committee

Brain Health Alliance, Ladera Ranch, CA, USA

BHAVI Symposium online 9 October 2023

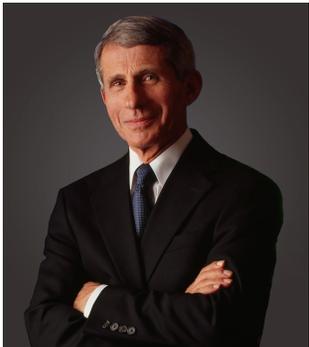


Anthony S. Fauci, MD



Honoring Dr. Anthony. S. Fauci

BHAVI 2023 Guardian of Truth and Integrity: We honor and thank Dr. Fauci as our 2023 Guardian in recognition of his fearless and tireless service as a physician, scientist, statistician and public health advocate during his lifetime of work in support of public health, societal health and the public good. Devoted to discovering the causes of infectious diseases, promoting the development of vaccines and use of vaccinations to prevent the transmission of contagious diseases from one person to the next, publicizing truthful information communicated with medical and scientific facts about public health measures and interventions to control and contain epidemics and pandemics, Dr. Fauci has contributed to stopping the spread of contagions from causing unnecessary deaths in communities around the world, and thus, has saved countless lives on planet earth.



Public Acclaim

- 2005 US Presidential Medal of Science [presented 27 Jul 2007](#)
- 2008 US Presidential Medal of Freedom [presented 19 Jun 2008](#)
- 2021 NAS Public Welfare Medal [presented 25 Apr 2021](#)
- Honorary Doctor of Science awards from numerous universities
- “Scientists Reflect on Anthony Fauci’s Impact: From the AIDS epidemic to the COVID-19 pandemic, the iconic medical chief has advised seven presidents on numerous outbreaks” at Scientific American [Nature Public Health 24 Aug 2022](#)
- PBS.org American Masters documentary “behind-the-scenes look at his 50-year career in public health” for [Dr. Tony Fauci 21 Mar 2023](#)
- Credited with saving the lives of as many as 25 million people around the world with his support of the [PEPFAR HIV Program](#)

A Very Brief History of Vaccination

- “Louis Pasteur, the Father of Immunology?” and the germ theory of disease 2012 K. A. Smith; doi:[10.3389/fimmu.2012.00068](https://doi.org/10.3389/fimmu.2012.00068)
- “Edward Jenner and the history of smallpox and vaccination” 2005 S. Riedel; doi:[10.1080/08998280.2005.11928028](https://doi.org/10.1080/08998280.2005.11928028)
- “Jonas Salk (1914–1995): A vaccine against polio” 2019 S. Y. Tan and N. Ponstein; doi:[10.11622/smedj.2019002](https://doi.org/10.11622/smedj.2019002)
- “Nobel Prize Awarded to Covid Vaccine Pioneers” Katalin Karikó and Drew Weissman; [New York Times 3 Oct 2023](#) by B. Mueller and G. Kolata
- The contributions of Dr. A. S. Fauci in the fight against viral disease must also be recognized in any [Brief History of Vaccination](#)

Highly Cited Research Authored/coauthored by Dr. Fauci

- Fauci (1983) “Wegener’s Granulomatosis: Prospective Clinical and Therapeutic Experience With 85 Patients for 21 Years”
- Pantaleo, Graziosi, and Fauci (1993) “The Immunopathogenesis of Human Immunodeficiency Virus Infection”
- Fauci (1996) “Host Factors and the Pathogenesis of HIV-induced Disease”
- Morens, Folkers, and Fauci (2004) “The Challenge of Emerging and Re-emerging Infectious Diseases”
- Fauci and Morens (2012) “The Perpetual Challenge of Infectious Diseases”
- Fauci (2022) “It ain’t over till it’s over. . . but it’s never over — Emerging and Re-emerging Infectious Diseases”

Biography

- Anthony S. Fauci was born in Brooklyn NY, worked in his father's pharmacy, graduated from Regis High School in Manhattan NY, then studied classics at College of the Holy Cross, and graduated with a Doctor of Medicine from Cornell University Medical College.
- Dr. Fauci served most of his career as Director of the NIH National Institute of Allergy and Infectious Diseases and provided leadership in the fight against viral diseases including HIV/AIDS, SARS, H1N1, MERS, Ebola, and COVID19.
- Detailed biographies available online:
 - anthonyfaucimd.com/bio-1
 - www.niaid.nih.gov/about/director
 - www.niaid.nih.gov/about/anthony-s-fauci-md-bio
 - en.wikipedia.org/wiki/Anthony_Fauci

Gallery

[Joan Baez](#) portrait of [Tony Fauci 2020](#) “Dear Dr. Fauci, I’ve painted your portrait to honor you and all you are doing for us and for the world.”
Other online galleries also tell the story of Dr. Fauci’s work:



- “[10 photos](#) exploring the many facets of Dr. Anthony Fauci” American Master at PBS.org
- “Dr. Anthony Fauci’s Esteemed [Career in Photos](#)” Person of the Year at People.com
- images via [Anthony Fauci](#) search at Google
- images via [NIAID](#) gallery at Flickr

Cited References



Fauci, Anthony S. (Jan. 1983). "Wegener's Granulomatosis: Prospective Clinical and Therapeutic Experience With 85 Patients for 21 Years". In: *Annals of Internal Medicine* 98.1, p. 76. DOI: [10.7326/0003-4819-98-1-76](https://doi.org/10.7326/0003-4819-98-1-76).



— (Dec. 1996). "Host Factors and the Pathogenesis of HIV-induced Disease". In: *Nature* 384.6609, pp. 529–534. DOI: [10.1038/384529a0](https://doi.org/10.1038/384529a0).



— (Dec. 2022). "It Ain't Over Till It's Over... but It's Never Over — Emerging and Reemerging Infectious Diseases". In: *New England Journal of Medicine* 387.22, pp. 2009–2011. DOI: [10.1056/nejmp2213814](https://doi.org/10.1056/nejmp2213814).



Fauci, Anthony S. and David M. Morens (Feb. 2012). "The Perpetual Challenge of Infectious Diseases". In: *New England Journal of Medicine* 366.5, pp. 454–461. DOI: [10.1056/nejmra1108296](https://doi.org/10.1056/nejmra1108296).



Morens, David M., Gregory K. Folkers, and Anthony S. Fauci (July 2004). "The challenge of emerging and re-emerging infectious diseases". In: *Nature* 430.6996, pp. 242–249. DOI: [10.1038/nature02759](https://doi.org/10.1038/nature02759).



Pantaleo, Giuseppe, Cecilia Graziosi, and Anthony S. Fauci (Feb. 1993). "The Immunopathogenesis of Human Immunodeficiency Virus Infection". In: *New England Journal of Medicine* 328.5. Ed. by Franklin H. Epstein, pp. 327–335. DOI: [10.1056/nejm199302043280508](https://doi.org/10.1056/nejm199302043280508).

Contact Info for Guardians

- guardians@bhavi.us
- www.BHAVI.us
- www.BrainiacsJournal.org
- www.BrainHealthAlliance.org

Fallacies and Pitfalls in Genome-Wide Association Studies

Julian Hecker, Adam Craig, Andrew Hughes, Julie Neidich, Carl Taswell, Nan Laird



Fallacies and Pitfalls in Genome-Wide Association Studies*

Julian Hecker, Adam Craig, Andrew Hughes, Julie Neidich, Carl Taswell, Nan Laird†

Abstract

Since the first genome-wide association study (GWAS) identifying variants associated with myocardial infarction was published over 20 years ago, GWASs have emerged as a powerful tool for exploring the genetic basis of complex traits. To date, hundreds of thousands of statistically significant associations have been reported across thousands of human phenotypes. Nevertheless, the design, implementation, and analysis of GWASs remain complex, and the results are easily misinterpreted. Common mistakes include 1) assuming that variants with the strongest statistical associations are causal instead of correlative, 2) believing that associated loci act through nearby genes, and 3) overemphasizing the contribution of individual loci to the total variability of particular traits. Clinical assays have been designed using the results of GWAS that rely on the contribution of such erroneous data interpretations to predict clinical phenotypes, reactions to medications or foods, and/or propensity to develop diseases. The failure to recognize these errors due to fallacies in logical reasoning and statistical inference presents problems for both the scientific community when the wrong targets may be prioritized in future research studies, as well as for communication with the general public when our understanding of the genetic basis of important traits may be misrepresented and overstated. Here, we review statistical data quality, analysis, and meta-analysis, of GWAS results with an emphasis on accurate and reliable interpretation. Placed in the appropriate context, GWASs enable genome-wide discovery of loci associated with diverse traits, but they constitute only a first step towards understanding the biological mechanism(s) underlying the observed associations. Scientific elucidation of these biological mechanisms must be required to establish causality with biochemical and pathophysiological explanations for any putative statistical correlations.

Keywords

Genome-wide association studies (GWAS), correlation-causation fallacy, meta-analysis, random effects model, fixed effects model, population stratification, family-based association studies (FBAS).

Contents

Introduction

1

*Presented 2023-10-09 with [JH slides](#), [NL slides](#), and [JH+NL video](#) at [Guardians 2023](#)

†JH and NL affiliated with Harvard University; AH and JN with Washington University St Louis; AC and CT with Brain Health Alliance; CT and NL contributed as senior co-authors; correspondence to [J Hecker at Harvard](#).

| | |
|--------------------------------------------|---|
| Multiple Testing | 2 |
| Linkage Disequilibrium | 2 |
| Study Design | 2 |
| Meta-Analyses | 3 |
| From GWAS to Biology | 3 |
| Direct-to-Consumer Testing | 4 |
| Conclusion | 5 |
| Citation | 5 |
| References | 5 |

Introduction

Genome-wide association studies (GWASs) aim to identify associations of genetic variants with phenotypes ([Visscher et al. 2017](#)). Most commonly, so-called single-nucleotide polymorphisms (SNPs) are considered in GWASs, and each available SNP is tested for association separately. After more than 15 years of GWASs, thousands of genetic associations were reported and partially replicated ([Abdellaoui et al. 2023](#)). Examples include identification of a female-specific association between SNPs at the *PAX1* enhancer locus and idiopathic scoliosis ([Sharma et al. 2015](#)), linking of *TAF3* to control of corpuscular hemoglobin concentration ([Pistis et al. 2013](#)), and discovery of the role of introns of the *FTO* gene in obesity ([Smemo et al. 2014](#); [Claussnitzer et al. 2016](#)). Many GWASs focus on so-called complex traits and diseases that are described by a polygenic architecture ([Visscher et al. 2017](#)). A trait with a polygenic architecture is influenced by thousands of causal genetic variants with rather small effect sizes ([Tam et al. 2019](#)). Examples include asthma, schizophrenia, body mass index, and human height ([Vicente et al. 2017](#); [Tam et al. 2019](#); [Yengo et al. 2022](#)).

Consequently, massive efforts by the research community to collect genetic and phenotypic data in large databases, such as the UK Biobank ([Bycroft et al. 2018](#)), led sample sizes in GWASs to grow rapidly over the years, enabling the identification of an increasing number of genetic risk loci ([Visscher et al. 2017](#)). According to one projection, use of GWASs to inform selection of drug targets and indications could double the number of drug candidates that successfully pass from phase I clinical trials to approval ([Nelson et al. 2015](#)). Nowadays, GWASs are considered to be a success story that identified several important genetic factors

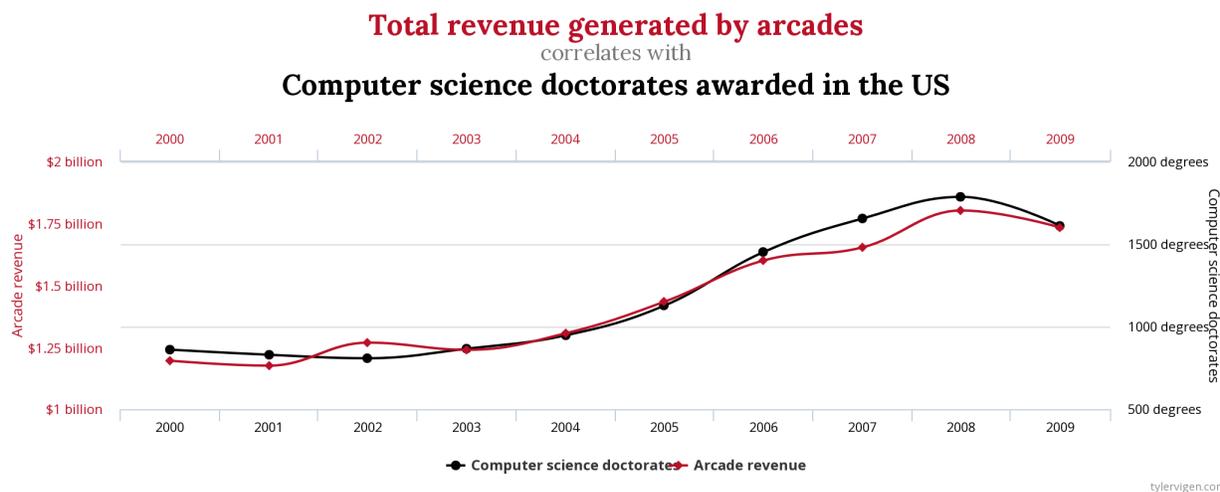


Figure 1: Correlation does not prove causation; see [Vigen \(2023\)](#) for this and other examples of the correlation-causation fallacy.

of complex diseases, but GWASs also face challenges, pitfalls, and limitations that we would like to discuss and review here.

Multiple Testing

As described above, each available SNP is tested for association with the phenotype of consideration. The density of available SNPs in a GWAS depends on the underlying platform used. Over time, the number of genes tested in a single GWAS has grown, starting from microarrays of a few thousand followed by genetic imputations (predicting genetic information based on reference panels), and expanding until, at present, whole-exome and whole-genome sequencing datasets are common ([Tam et al. 2019](#)). This implies that a typical GWAS incorporates more than a million common SNPs ([DerSimonian and Laird 2015](#)), leading to a substantial multiple testing burden. The established significance level for so-called genome-wide significance in a GWAS is $p=5e-08$ ([Tam et al. 2019](#)). This significance level corresponds to a nominal level of 0.05 corrected for 1 million independent statistical tests by the Bonferroni correction ([Tam et al. 2019](#)). More specific gene-based tests, such as the versatile gene-based association study (VEGAS) methodology, can arrive at appropriately corrected p -values by using statistics on common genetic variation and the corresponding linkage disequilibrium (LD) structures in reference panels such as the HapMap project ([Hecker et al. 2017](#)). LD describes nonrandom association of alleles at different loci resulting from complex interactions between recombination, mutation, selection, and genetic drift ([Slatkin 2008](#)). It provides the key patterns of information on which statistical methods of fine-scale gene mapping rely ([Slatkin 2008](#)). A pitfall in replication attempts is the winner's curse ([Zhong and Prentice 2010](#)). The winner's curse describes the phenomenon that the effect sizes of genetic variants that just passed the genome-wide significance level tend to be overestimated ([Zhong and Prentice 2010](#)). This in turn, leads to overestimated power calculations in replication GWASs and was one of the driving factors for the lack of replication in early GWASs ([Zhong and Prentice 2010](#)). A wide variety of statistical correction methods can partially account for the winner's curse effect, with empirical Bayesian, FDR inverse quantile transformation, and bootstrap resampling methods outperforming commonly used conditional likelihood methods ([Forde et al. 2023](#)).

Linkage Disequilibrium

Moreover, since nearby genetic variants are often in LD, the genetic information for local SNPs is usually correlated ([Tam et al. 2019](#); [Lappalainen and MacArthur 2021](#)). Consequently, the presence of a causal genetic variant resulting in a significant association with the phenotype also leads to a substantial number of significant associations for nearby SNPs ([Lappalainen and MacArthur 2021](#)). Therefore, GWASs typically report associated genetic risk loci that contain these multiple associations. A common pitfall is that a genetic variant with a genome-wide significant association p -value is interpreted to be causal, although it potentially only tags a causal variant through LD ([Visscher et al. 2017](#)). This is an example of the causation-correlation fallacy.

Fine-mapping is an approach to tackle this problem ([Schaid et al. 2018](#)). Fine-mapping prioritizes a set of genetic variants that most likely contains the causal variant at this genetic risk locus, often assuming the presence of at most one causal SNP ([Schaid et al. 2018](#)). This procedure incorporates LD information and the individual SNP association statistics ([Schaid et al. 2018](#)). Interestingly, there are potential scenarios with multiple genetic effects within the loci in which the most significant SNP is not causal, especially when the statistical power of the study is low ([Schaid et al. 2018](#)). Recent fine-mapping approaches are based on Bayesian computations ([Schaid et al. 2018](#); [Tam et al. 2019](#)). This purely statistical fine-mapping can be improved by incorporating external biological information such as functional annotations. This external information can be integrated into the prior distributions of the Bayesian models and therefore guides these analyses.

Study Design

Another challenge in GWAS concerns the selected study design. The design of a GWAS with regard to how the subjects are selected for participation can impact the results ([Heid, Huth, et al. 2009](#)). Two time-honored epidemiological designs for studying causality in the absence of randomization are the case-control study and the cohort study. The case-control design has been very successful in GWAS ([Clarke et al. 2011](#)). Likewise, many GWAS have taken advantage of existing cohorts, such as the Nurse's Health Study, provided that the phenotypes of interest can be obtained. An alternative approach chooses subjects as part of a random sample from a population. Random samples from

selected populations are generally difficult and expensive to obtain, but may be available in some countries as a part of ongoing research programs, for example, the Framingham Heart Study in the United States (Dawber 1980) and the KORA Study in Germany (Heid, Vollmert, et al. 2005).

Another paradigm developed for GWAS involves the use of genomic repositories or biobanks. The idea is to provide access to very large data sets. Genotypes are recorded in a central repository without regard to the phenotypic status of subjects. Selection bias may make it difficult to interpret the results of such studies. Even more importantly, allele frequencies vary with genetic ancestries (Derks et al. 2022). If phenotypic differences correlate with genetic ancestry in the study population, often because of specific participant sampling procedures, genetic association testing based on this phenotype can lead to false positive findings if not appropriately controlled for genetic ancestry (Derks et al. 2022). Furthermore, predictive models of quantitative traits based on data from a single ancestry group can generalize poorly to other populations, as with a predictive model of height found to account for 45% of variation in European populations but only 14–24% in others (Yengo et al. 2022).

The most established approach to adjust for genetic ancestry and therefore reducing the likelihood for false positive associations is to include principal components of genetic ancestry derived from genome-wide data as covariates in the statistical association tests (Price et al. 2006). However, this approach is not guaranteed to fully adjust for ancestry-induced signals, as even large data sets, such as that of the 1000 Genomes Project, may not have sufficient coverage of genetic diversity in some populations (D. Lu and Xu 2013), and the urge to increase sample sizes in recent GWAS can amplify this issue.

An approach to address the issue of population stratification is a family-based study design (Rabinowitz and Laird 2000; Tam et al. 2019). Here cases, or affected individuals, and their family members (ideally their parents) are chosen to be study participants. The family members serve as the controls. Family-based study designs allow genetic association tests for SNPs that are robust to population stratification (Derks et al. 2022). They are particularly useful for observing segregation of rare variants with very large effect sizes when those variants segregate within a family (Visscher et al. 2017).

The classical example is the transmission disequilibrium test (TDT) (Schaid 1998). The TDT considers affected offspring trios and tests the observed allele transmissions against Mendelian expectations (Schaid 1998). Since the test statistic conditions on parental genotypes, the test does not require any assumptions about the underlying allele frequencies and distributions (Ewens and Spielman 1995). This concept was extended to general pedigrees, general phenotypes, and groups of genetic variants in the Family-Based Association Test framework by Laird and Lange (2006).

Meta-Analyses

To achieve the desired large sample sizes, researchers combine their association results in meta-analyses across cohorts and studies (Abdellaoui et al. 2023; Mikolajewicz and Komarova 2019; Steel et al. 2021). Since meta-analyses can achieve the same results by combining summary statistics as with individual-level data (DerSimonian and Laird 2015), they also have the advantage that data resources can be combined without sharing protected individual genetic information across institutions and scientific groups. Several consortia of researchers and institutions have formed to pool data and set standards for studies to

be included in meta-analyses relevant to particular areas of health and wellness, including the Psychiatric Genomics Consortium, the Genetic Investigation of Anthropometric Traits (GIANT) Consortium, and the Global Lipids Genetics Consortium (O'Donovan 2015; H. Park et al. 2016; Klarin et al. 2018). The approaches to meta-analyses include fixed effects and random effects models (Steel et al. 2021). The latter explicitly allows for heterogeneity in the data (DerSimonian and Laird 2015; Steel et al. 2021). The combination of association results across studies with varying genetic ancestry has the advantage that the differences in the LD structure can lead to an improved resolution in the fine-mapping step of genetic risk loci since the LD effects dilute (DerSimonian and Laird 2015). However, careful interpretation is required. Since most GWAS so far were based on participants of European genetic ancestry, the analysis of other genetic ancestries and ethnicities has great potential to reveal a refined picture of genetic associations and generalizability across populations (Derks et al. 2022). Such meta-analyses rely on accurate and detailed metadata to ensure that results across different studies are comparable (Mikolajewicz and Komarova 2019; Steel et al. 2021). The NHGRI-EBI GWAS Catalog represents one attempt to compile such data in an online repository on a large scale (Sollis et al. 2022).

From GWAS to Biology

Even though GWAS publications have reported thousands of genetic associations with a plethora of complex diseases and traits, and recent advances in fine-mapping in combination with large sample sizes have pinpointed genetic variants with potential causal associations, the underlying biological mechanisms of these associations remain largely unknown. This is because the exact regulatory function of most GWAS hits, which are often located in non-coding regions of the genome, is poorly understood (Abdellaoui et al. 2023; Aguet et al. 2023). Therefore, the role of a SNP and its downstream effects on other genes and pathways is often unknown. Projects such as the Encyclopedia of DNA Elements (ENCODE) are working to fill this knowledge gap by compiling an extensive repository of millions of human and mouse functional elements, including protein-coding genes, regulatory RNA-coding genes, and non-coding regions with known mechanistic functions, such as promoters and enhancers (Moore et al. 2020). Similarly, the GENCODE project publishes extensive annotations of the human and mouse genomes, including protein-coding genes, pseudogenes, and long non-coding RNA genes (Frankish et al. 2020). Using such gene annotations enables new approaches to weighting the significance of association scores based on this prior knowledge in addition to LD and Bonferroni correction (Kichaev et al. 2019). The candidate causal gene is commonly inferred based on the smallest physical distance, but recent investigations showed that this might be misleading. One possibility to gain further insights into the identified genetic associations is to study molecular quantitative trait loci (QTLs) (Lappalainen and MacArthur 2021; Aguet et al. 2023). These SNPs are associated with molecular phenotypes such as RNA expression, DNA methylation, or metabolite levels (Aguet et al. 2023).

Colocalization analyses attempt to test if GWAS findings colocalize with both molecular and expression QTLs (i.e., the same genetic variant is implicated), and such successful colocalizations provide the basis for mediation hypotheses (Rheenen et al. 2021). A systematic version of this concept, a post-GWAS transcriptome-wide association study (TWAS), tests for associations between traits and gene expression levels imputed from eQTLs across the entire genome, while a proteome-wide



Figure 2: Mixed measures? When are meta-analyses reproducible and valid? (New Cuyama sign image by Gogulski 2007.)

association study (PWAS) tests for associations with protein abundance as predicted from population-level protein QTL (pQTL) data (Gedik et al. 2023). By studying the downstream effects of genetic variants, these approaches can identify genes that affect health via differences in quantitatively measured expression traits better correlated with phenotypes even when the direct association between the genotypes and phenotypes otherwise would be weak (Gedik et al. 2023).

One approach that has built on this idea further is the use of colocalization in conjunction with similarity of annotations from single-cell gene expression, protein-protein interaction, and pathway participation features to compute a polygenic priority score to identify associations between non-coding loci and protein-coding genes that are likely to be causal (Weeks et al. 2023). As noted in (Weeks et al. 2023), combining such similarity-based methods with complementary locus-based methods can achieve better results than either one can alone. Taking that reasoning even further, (Gazal et al. 2022) propose a framework for arriving at combinations SNP-to-gene strategies and apply it to select seven such strategies that together achieve higher recall than attainable with any one strategy alone. Ultimately, GWAS alone cannot determine the causal mechanisms behind human health and diversity, which requires taking the next step of analyzing the GWAS-identified

candidate genes through both statistical and bench-based functional testing (Gallagher and Chen-Plotkin 2018).

Direct-to-Consumer Testing

Some privately held laboratories, especially those offering direct-to-consumer testing (US Food & Drug Administration 2019; Malgorzata et al. 2021), have used GWAS data for the interpretation of genomic tests for a variety of diverse indications including fear of heights to cat allergy to anxiety to dietary advice. These kinds of indications often have vague symptoms, little evidence of heritability, or are common disorders that may have a multifactorial pattern of inheritance without a specific genotype-phenotype association. The commercial labs often describe the tests as available for personal amusement and not for diagnosis of any specific condition. For the FDA-approved assays that also run at these labs, the reported results are not based on GWAS data. For example, some offer FDA-approved tests for pathogenic variants in the genes associated with increased risk of the development of breast or other cancers (National Human Genome Research Institute 2023). Studies have identified ethical and legal concerns with the DTC testing modality (Martins et al. 2022; Panacer 2023), and specifically with the use of polygenic risk scores that are based solely on GWAS data (J. K.

Park and C. Y. Lu 2023).

Conclusion

Genome-wide association studies identified thousands of genetic associations with a wide range of phenotypes. As a consequence of the polygenic architecture of complex traits and diseases, recent GWASs reached sample sizes of 1 million samples to identify novel genetic risk loci. However, the interpretation of GWAS results requires careful consideration. Technical artefacts such as population stratification can introduce false positive findings in GWAS and identified genetic associations should always be replicated in independent studies. A significant GWAS signal does not imply causality and the identification of causal genetic variants within a genetic risk locus remains a challenge. Furthermore, most GWAS hits are in non-coding regions of the genome and mapping genetic associations to candidate genes for functional follow-up analyses is non-trivial and of limited success so far. Overall, the underlying mechanisms of genetic associations remain poorly understood and there is a risk of overinterpreting their individual relevance in clinical risk prediction and other complex traits such as educational attainment (Okbay et al. 2022; Cesarini and Visscher 2017). While there is great potential in utilizing the findings from GWAS to support the development of new drugs and approach the reality of personalized medicine based on individual risk evaluation, the application of GWAS as a research tool comes with ethical and social responsibility.

Citation

Brainiacs 2023 Volume 4 Issue 2 Edoc GFA4E8812

Title: "Fallacies and Pitfalls in Genome-Wide Association Studies"

Authors: Julian Hecker, Adam Craig, Andrew Hughes, Julie Neidich, Carl Taswell, Nan Laird

Dates: created 2023-10-01, received 2023-10-03, presented 2023-10-09, updated 2023-12-21, published 2023-12-21, endorsed 2023-12-30

Copyright: © 2023 Brain Health Alliance

Contact: J Hecker at Harvard

URL: [Brainiacsjournal.org/arc/pub/Hecker2023FPGWAS](https://brainiacsjournal.org/arc/pub/Hecker2023FPGWAS)

PDP: [/Nexus/Brainiacs/Hecker2023FPGWAS](https://nexus.brainiacs/Hecker2023FPGWAS)

DOI: [/10.48085/GFA4E8812](https://doi.org/10.48085/GFA4E8812)

References

- [1] A. Abdellaoui, L. Yengo, K. J. Verweij, and P. M. Visscher. "15 years of GWAS discovery: Realizing the promise." *The American Journal of Human Genetics* 110.2 (Feb. 2023), pp. 179–194. ISSN: 0002-9297. DOI: [10.1016/j.ajhg.2022.12.011](https://doi.org/10.1016/j.ajhg.2022.12.011) (cited pp. 1, 3).
- [2] F. Aguet, K. Alasoo, Y. I. Li, A. Battle, H. K. Im, S. B. Montgomery, and T. Lappalainen. "Molecular quantitative trait loci." *Nature Reviews Methods Primers* 3.1 (Jan. 2023). ISSN: 2662-8449. DOI: [10.1038/s43586-022-00188-6](https://doi.org/10.1038/s43586-022-00188-6) (cited p. 3).
- [3] C. Bycroft, C. Freeman, D. Petkova, G. Band, et al. "The UK Biobank resource with deep phenotyping and genomic data." *Nature* 562.7726 (Oct. 2018), pp. 203–209. ISSN: 1476-4687. DOI: [10.1038/s41586-018-0579-z](https://doi.org/10.1038/s41586-018-0579-z) (cited p. 1).
- [4] D. Cesarini and P. M. Visscher. "Genetics and educational attainment." *npj Science of Learning* 2.1 (Feb. 2017). ISSN: 2056-7936. DOI: [10.1038/s41539-017-0005-6](https://doi.org/10.1038/s41539-017-0005-6) (cited p. 5).
- [5] G. M. Clarke, C. A. Anderson, F. H. Pettersson, L. R. Cardon, A. P. Morris, and K. T. Zondervan. "Basic statistical analysis in genetic case-control studies." *Nature Protocols* 6.2 (Feb. 2011), pp. 121–133. ISSN: 1750-2799. DOI: [10.1038/nprot.2010.182](https://doi.org/10.1038/nprot.2010.182) (cited p. 2).
- [6] M. Claussnitzer, S. N. Dankel, K.-H. Kim, G. Quon, W. Meuleman, et al. "FTO Obesity Variant Circuitry and Adipocyte Browning in Humans." *New England Journal of Medicine* 374.2 (Jan. 2016), pp. 190–193. ISSN: 1533-4406. DOI: [10.1056/nejmc1513316](https://doi.org/10.1056/nejmc1513316) (cited p. 1).
- [7] T. R. Dawber. *The Framingham Study. The Epidemiology of Atherosclerotic Disease*. Harvard University Press, 1980. ISBN: 9780674492080. DOI: [10.4159/harvard.9780674492097](https://doi.org/10.4159/harvard.9780674492097) (cited p. 3).
- [8] E. M. Derks, J. G. Thorp, and Z. F. Gerring. "Ten challenges for clinical translation in psychiatric genetics." *Nature Genetics* 54.10 (Sept. 2022), pp. 1457–1465. ISSN: 1546-1718. DOI: [10.1038/s41588-022-01174-0](https://doi.org/10.1038/s41588-022-01174-0) (cited p. 3).
- [9] R. DerSimonian and N. Laird. "Meta-analysis in clinical trials revisited." *Contemporary Clinical Trials* 45 (Nov. 2015), pp. 139–145. ISSN: 1551-7144. DOI: [10.1016/j.cct.2015.09.002](https://doi.org/10.1016/j.cct.2015.09.002) (cited pp. 2, 3).
- [10] W. J. Ewens and R. S. Spielman. "The transmission/disequilibrium test: history, subdivision, and admixture." *American journal of human genetics* 57.2 (1995), p. 455 (cited p. 3).
- [11] A. Forde, G. Hemani, and J. Ferguson. "Review and further developments in statistical corrections for Winner's Curse in genetic association studies." *PLoS Genetics* 19.9 (2023), e1010546 (cited p. 2).
- [12] A. Frankish, M. Diekhans, I. Jungreis, J. Lagarde, et al. "GENCODE 2021." *Nucleic Acids Research* 49.D1 (Dec. 2020), pp. D916–D923. ISSN: 1362-4962. DOI: [10.1093/nar/gkaa1087](https://doi.org/10.1093/nar/gkaa1087) (cited p. 3).
- [13] M. D. Gallagher and A. S. Chen-Plotkin. "The Post-GWAS Era: From Association to Function." *The American Journal of Human Genetics* 102.5 (May 2018), pp. 717–730. ISSN: 0002-9297. DOI: [10.1016/j.ajhg.2018.04.002](https://doi.org/10.1016/j.ajhg.2018.04.002) (cited p. 4).
- [14] S. Gazal, O. Weissbrod, F. Hormozdiari, K. K. Dey, et al. "Combining SNP-to-gene linking strategies to identify disease genes and assess disease omnigenicity." *Nature Genetics* 54.6 (June 2022), pp. 827–836. ISSN: 1546-1718. DOI: [10.1038/s41588-022-01087-y](https://doi.org/10.1038/s41588-022-01087-y) (cited p. 4).
- [15] H. Gedik, R. E. Peterson, B. P. Riley, V. I. Vladimirov, and S.-A. Bacanu. "Integrative Post-Genome-Wide Association Study Analyses Relevant to Psychiatric Disorders: Imputing Transcriptome and Proteome Signals." *Complex Psychiatry* 9.1-4 (2023), pp. 130–144 (cited p. 4).
- [16] M. Gogulski. *Image of sign at 4923 Primero St, New Cuyama, CA 93254*. Ed. by Wikipedia. Aug. 4, 2007. URL: https://en.wikipedia.org/wiki/New_Cuyama,_California#/media/File:US-CA,_New_cuyama.jpg (cited p. 4).
- [17] J. Hecker, A. Maaser, D. Prokopenko, H. L. Fier, and C. Lange. "Reporting Correct p Values in VEGAS Analyses." *Twin Research and Human Genetics* 20.3 (Mar. 2017), pp. 257–259. ISSN: 1839-2628. DOI: [10.1017/tbg.2017.16](https://doi.org/10.1017/tbg.2017.16) (cited p. 2).
- [18] I. M. Heid, C. Vollmert, A. Hinney, A. Döring, et al. "Association of the 1031 MC4R allele with decreased body mass in 7937 participants of two population based surveys." *Journal of Medical Genetics* 42.4 (Apr. 2005), e21. ISSN: 1468-6244. DOI: [10.1136/jmg.2004.027011](https://doi.org/10.1136/jmg.2004.027011) (cited p. 3).
- [19] I. M. Heid, C. Huth, R. J. F. Loos, F. Kronenberg, et al. "Meta-Analysis of the INSIG2 Association with Obesity Including 74,345 Individuals: Does Heterogeneity of Estimates Relate to Study Design?" *PLoS Genetics* 5.10 (Oct. 2009). Ed. by D. B. Allison, e1000694. ISSN: 1553-7404. DOI: [10.1371/journal.pgen.1000694](https://doi.org/10.1371/journal.pgen.1000694) (cited p. 2).

- [20] G. Kichaev, G. Bhatia, P.-R. Loh, S. Gazal, et al. "Leveraging Polygenic Functional Enrichment to Improve GWAS Power." *The American Journal of Human Genetics* 104.1 (Jan. 2019), pp. 65–75. ISSN: 0002-9297. DOI: [10.1016/j.ajhg.2018.11.008](https://doi.org/10.1016/j.ajhg.2018.11.008) (cited p. 3).
- [21] D. Klarin, S. M. Damrauer, K. Cho, Y. V. Sun, et al. "Genetics of blood lipids among ~300,000 multi-ethnic participants of the Million Veteran Program." *Nature genetics* 50.11 (2018), pp. 1514–1523 (cited p. 3).
- [22] N. Laird and C. Lange. "Family-based designs in the age of large-scale gene-association studies." *Nature Reviews Genetics* 7.5 (May 2006), pp. 385–394. ISSN: 1471-0064. DOI: [10.1038/nrg1839](https://doi.org/10.1038/nrg1839) (cited p. 3).
- [23] T. Lappalainen and D. G. MacArthur. "From variant to function in human disease genetics." *Science* 373.6562 (Sept. 2021), pp. 1464–1468. ISSN: 1095-9203. DOI: [10.1126/science.abi8207](https://doi.org/10.1126/science.abi8207) (cited pp. 2, 3).
- [24] D. Lu and S. Xu. "Principal component analysis reveals the 1000 Genomes Project does not sufficiently cover the human genetic diversity in Asia." *Frontiers in Genetics* 4 (2013). ISSN: 1664-8021. DOI: [10.3389/fgene.2013.00127](https://doi.org/10.3389/fgene.2013.00127) (cited p. 3).
- [25] M. Malgorzata, S. Maria, and W. Michał. "Genetic testing—whether to allow complete freedom? Direct to consumer tests versus genetic tests for medical purposes." *Journal of Applied Genetics* 63.1 (Nov. 2021), pp. 119–126. ISSN: 2190-3883. DOI: [10.1007/s13353-021-00670-z](https://doi.org/10.1007/s13353-021-00670-z) (cited p. 4).
- [26] M. F. Martins, L. T. Murry, L. Telford, and F. Moriarty. "Direct-to-consumer genetic testing: an updated systematic review of healthcare professionals' knowledge and views, and ethical and legal concerns." *European Journal of Human Genetics* 30.12 (Oct. 2022), pp. 1331–1343. ISSN: 1476-5438. DOI: [10.1038/s41431-022-01205-8](https://doi.org/10.1038/s41431-022-01205-8) (cited p. 4).
- [27] N. Mikolajewicz and S. V. Komarova. "Meta-Analytic Methodology for Basic Research: A Practical Guide." *Frontiers in Physiology* 10 (Mar. 2019). ISSN: 1664-042X. DOI: [10.3389/fphys.2019.00203](https://doi.org/10.3389/fphys.2019.00203) (cited p. 3).
- [28] J. E. Moore, M. J. Purcaro, H. E. Pratt, C. B. Epstein, et al. "Expanded encyclopaedias of DNA elements in the human and mouse genomes." *Nature* 583.7818 (2020), pp. 699–710 (cited p. 3).
- [29] National Human Genome Research Institute. *Direct-to-Consumer Genetic Testing FAQ for Healthcare Professionals*. June 14, 2023. URL: <https://www.genome.gov/For-Health-Professionals/Provider-Genomics-Education-Resources/Healthcare-Provider-Direct-to-Consumer-Genetic-Testing-FAQ> (cited p. 4).
- [30] M. R. Nelson, H. Tipney, J. L. Painter, J. Shen, et al. "The support of human genetic evidence for approved drug indications." *Nature Genetics* 47.8 (June 2015), pp. 856–860. ISSN: 1546-1718. DOI: [10.1038/ng.3314](https://doi.org/10.1038/ng.3314) (cited p. 1).
- [31] M. C. O'Donovan. "What have we learned from the Psychiatric Genomics Consortium." *World Psychiatry* 14.3 (2015), p. 291 (cited p. 3).
- [32] A. Okbay, Y. Wu, N. Wang, H. Jayashankar, et al. "Polygenic prediction of educational attainment within and between families from genome-wide association analyses in 3 million individuals." *Nature Genetics* 54.4 (Mar. 2022), pp. 437–449. ISSN: 1546-1718. DOI: [10.1038/s41588-022-01016-z](https://doi.org/10.1038/s41588-022-01016-z) (cited p. 5).
- [33] K. S. Panacer. "Ethical Issues Associated With Direct-to-Consumer Genetic Testing." *Cureus* (June 2023). ISSN: 2168-8184. DOI: [10.7759/cureus.39918](https://doi.org/10.7759/cureus.39918) (cited p. 4).
- [34] H. Park, X. Li, Y. E. Song, K. Y. He, and X. Zhu. "Multivariate analysis of anthropometric traits using summary statistics of genome-wide association studies from GIANT Consortium." *PLoS one* 11.10 (2016), e0163912 (cited p. 3).
- [35] J. K. Park and C. Y. Lu. "Polygenic Scores in the Direct-to-Consumer Setting: Challenges and Opportunities for a New Era in Consumer Genetic Testing." *Journal of Personalized Medicine* 13.4 (Mar. 2023), p. 573. ISSN: 2075-4426. DOI: [10.3390/jpm13040573](https://doi.org/10.3390/jpm13040573) (cited p. 4).
- [36] G. Pistis, S. U. Okonkwo, M. Traglia, C. Sala, et al. "Genome Wide Association Analysis of a Founder Population Identified TAF3 as a Gene for MCHC in Humans." *PLoS ONE* 8.7 (July 2013). Ed. by F. Di Cunto, e69206. ISSN: 1932-6203. DOI: [10.1371/journal.pone.0069206](https://doi.org/10.1371/journal.pone.0069206) (cited p. 1).
- [37] A. L. Price, N. J. Patterson, R. M. Plenge, M. E. Weinblatt, N. A. Shadick, and D. Reich. "Principal components analysis corrects for stratification in genome-wide association studies." *Nature genetics* 38.8 (2006), pp. 904–909 (cited p. 3).
- [38] D. Rabinowitz and N. Laird. "A Unified Approach to Adjusting Association Tests for Population Admixture with Arbitrary Pedigree Structure and Arbitrary Missing Marker Information." *Human Heredity* 50.4 (2000), pp. 211–223. ISSN: 1423-0062. DOI: [10.1159/000022918](https://doi.org/10.1159/000022918) (cited p. 3).
- [39] W. van Rheenen, R. A. A. van der Spek, M. K. Bakker, J. J. F. A. van Vugt, et al. "Common and rare variant association analyses in amyotrophic lateral sclerosis identify 15 risk loci with distinct genetic architectures and neuron-specific biology." *Nature Genetics* 53.12 (Dec. 2021), pp. 1636–1648. ISSN: 1546-1718. DOI: [10.1038/s41588-021-00973-1](https://doi.org/10.1038/s41588-021-00973-1) (cited p. 3).
- [40] D. J. Schaid. "Transmission disequilibrium, family controls, and great expectations." *The American Journal of Human Genetics* 63.4 (1998), pp. 935–941 (cited p. 3).
- [41] D. J. Schaid, W. Chen, and N. B. Larson. "From genome-wide associations to candidate causal variants by statistical fine-mapping." *Nature Reviews Genetics* 19.8 (May 2018), pp. 491–504. ISSN: 1471-0064. DOI: [10.1038/s41576-018-0016-z](https://doi.org/10.1038/s41576-018-0016-z) (cited p. 2).
- [42] S. Sharma, D. Londono, W. L. Eckalbar, X. Gao, et al. "A PAX1 enhancer locus is associated with susceptibility to idiopathic scoliosis in females." *Nature Communications* 6.1 (Mar. 2015). ISSN: 2041-1723. DOI: [10.1038/ncomms7452](https://doi.org/10.1038/ncomms7452) (cited p. 1).
- [43] M. Slatkin. "Linkage disequilibrium—understanding the evolutionary past and mapping the medical future." *Nature Reviews Genetics* 9.6 (June 2008), pp. 477–485. ISSN: 1471-0064. DOI: [10.1038/nrg2361](https://doi.org/10.1038/nrg2361) (cited p. 2).
- [44] S. Smemo, J. J. Tena, K.-H. Kim, E. R. Gamazon, et al. "Obesity-associated variants within FTO form long-range functional connections with IRX3." *Nature* 507.7492 (Mar. 2014), pp. 371–375. ISSN: 1476-4687. DOI: [10.1038/nature13138](https://doi.org/10.1038/nature13138) (cited p. 1).
- [45] E. Sollis, A. Mosaku, A. Abid, A. Buniello, et al. "The NHGRI-EBI GWAS Catalog: knowledgebase and deposition resource." *Nucleic Acids Research* 51.D1 (Nov. 2022), pp. D977–D985. ISSN: 1362-4962. DOI: [10.1093/nar/gkac1010](https://doi.org/10.1093/nar/gkac1010) (cited p. 3).
- [46] P. Steel, S. Beugelsdijk, and H. Aguinis. "The anatomy of an award-winning meta-analysis: Recommendations for authors, reviewers, and readers of meta-analytic reviews." *Journal of International Business Studies* 52.1 (Jan. 2021), pp. 23–44. ISSN: 1478-6990. DOI: [10.1057/s41267-020-00385-z](https://doi.org/10.1057/s41267-020-00385-z) (cited p. 3).
- [47] V. Tam, N. Patel, M. Turcotte, Y. Bossé, G. Paré, and D. Meyre. "Benefits and limitations of genome-wide association studies." *Nature Reviews Genetics* 20.8 (May 2019), pp. 467–484. ISSN: 1471-0064. DOI: [10.1038/s41576-019-0127-1](https://doi.org/10.1038/s41576-019-0127-1) (cited pp. 1–3).
- [48] US Food & Drug Administration. *Direct-to-Consumer Tests*. Dec. 20, 2019. URL: <https://www.fda.gov/medical-devices/in-vitro-diagnostics/direct-consumer-tests> (cited p. 4).

- [49] C. T. Vicente, J. A. Revez, and M. A. Ferreira. "Lessons from ten years of genome-wide association studies of asthma." *Clinical & translational immunology* 6.12 (2017), e165 (cited p. 1).
- [50] T. Vigen. *Spurious Correlations*. 2023. URL: <https://tylervigen.com/spurious-correlations> (cited p. 2).
- [51] P. M. Visscher, N. R. Wray, Q. Zhang, P. Sklar, M. I. McCarthy, M. A. Brown, and J. Yang. "10 Years of GWAS Discovery: Biology, Function, and Translation." *The American Journal of Human Genetics* 101.1 (July 2017), pp. 5–22. ISSN: 0002-9297. DOI: [10.1016/j.ajhg.2017.06.005](https://doi.org/10.1016/j.ajhg.2017.06.005) (cited pp. 1–3).
- [52] E. M. Weeks, J. C. Ulirsch, N. Y. Cheng, B. L. Trippe, et al. "Leveraging polygenic enrichments of gene features to predict genes underlying complex traits and diseases." *Nature Genetics* 55.8 (July 2023), pp. 1267–1276. ISSN: 1546-1718. DOI: [10.1038/s41588-023-01443-6](https://doi.org/10.1038/s41588-023-01443-6) (cited p. 4).
- [53] L. Yengo, S. Vedantam, E. Marouli, J. Sidorenko, et al. "A saturated map of common genetic variants associated with human height." *Nature* 610.7933 (Oct. 2022), pp. 704–712. ISSN: 1476-4687. DOI: [10.1038/s41586-022-05275-y](https://doi.org/10.1038/s41586-022-05275-y) (cited pp. 1, 3).
- [54] H. Zhong and R. L. Prentice. "Correcting "winner's curse" in odds ratios from genomewide association findings for major complex human diseases." *Genetic Epidemiology: The Official Publication of the International Genetic Epidemiology Society* 34.1 (2010), pp. 78–91 (cited p. 2).

Photoshop Fantasies

Walter Scheirer



Photoshop Fantasies

Walter J. Scheirer
Dept. of Computer Science and Engineering
University of Notre Dame



Notre Dame Institute for
ADVANCED STUDY

TEC Technology
Ethics Center

Why is there so much fake stuff on the Internet?



Donald J. Trump
47.9K Tweets

Donald J. Trump
@realDonaldTrump

45th President of the United States of America 🇺🇸
Washington, DC [instagram.com/realDonaldTrump](#)
Joined March 2009

47 Following 70.9M Followers
Followed by [asad@revolution](#), [Ethan Azad](#), and 384 others you follow

Tweets Tweets & replies Media Likes

Donald J. Trump Retweeted
[@DDwn_Linder](#) · 4h
The corrupted Dems trying their best to come to the Ayatollah's rescue.
[#NancyPelsOffAsHeres](#)

1.3K 3.8K 9.5K

Team Trump (Text TRUMP to 88022) @TeamTrump

Deep in the heart of Delaware, Joe Biden sits in his basement.
Alone. Hiding. Diminished.



89.1K views
Delaware
Text TRUMP to 88022

Eric Trump @EricTrump

Two great, courageous, Americans! 🇺🇸
[@icecube](#) [@50cent](#)

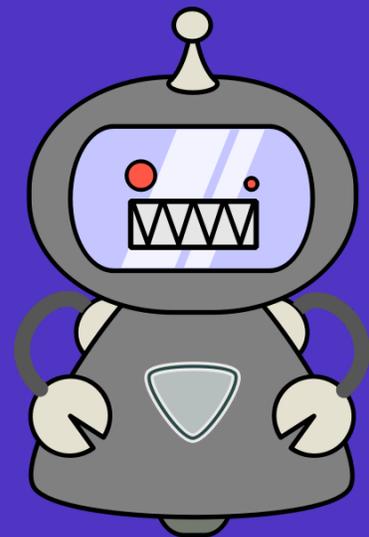
Manipulated media

1:59 PM · Oct 20, 2020 · Twitter for iPhone

How much fake stuff is out there?

- Less than 60% of all web traffic is human
- The bulk of users on social media platforms are bots
- The veracity of most of the content people are consuming is in question

Source: "How Much of the Internet Is Fake? Turns Out, a Lot of It, Actually." *New York Magazine*. Dec. 2018.



Deep Fakes



Restyling Reality



Surely we should just ban all of this stuff, right?



Maybe a reflection on photography itself is warranted

A very brief history of photographic propaganda

“If you want knowledge, you must take part in the practice of changing reality.”

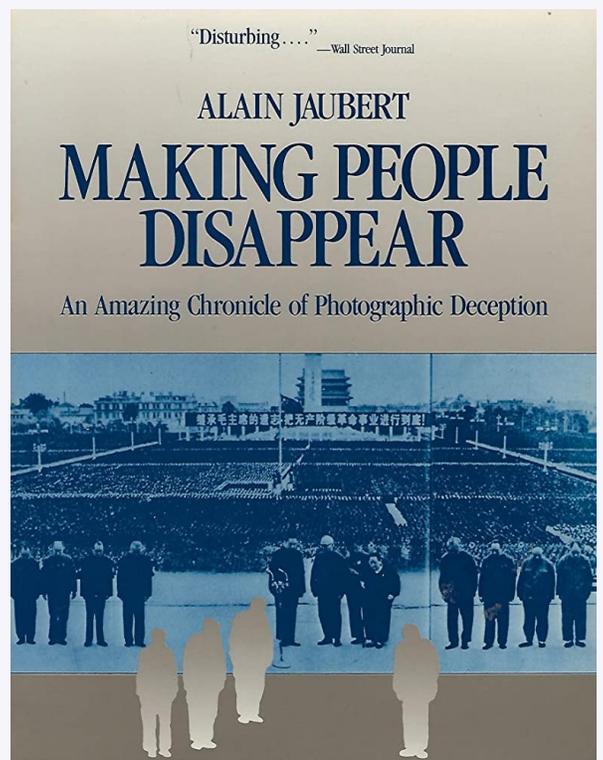
Mao Zedong, “On Practice”

Mao's Funeral

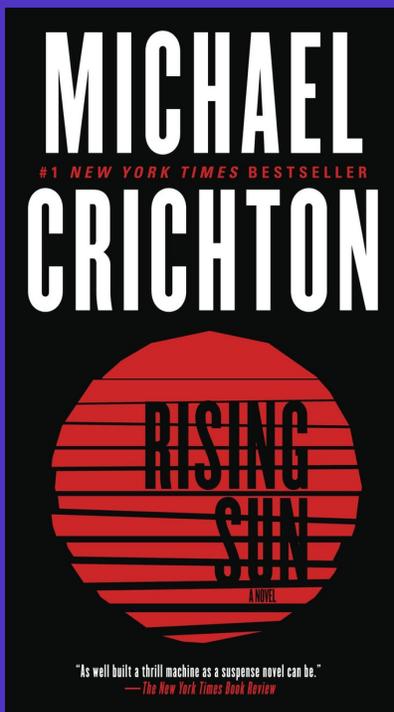


“Jaubert points out that the systematic faking or falsification of history are encountered ‘anywhere a serious effort is made to copy the methods of totalitarian propaganda characterized by its three major styles: Fascist, Nazi, and Communist...’ The falsification of photographs comes easily to those governments and elites that seek to be the sole interpreters of history and have a monopoly on the information media...”

Roy Godson



Rising Sun (1992)



“The case law isn't entirely clear. But it's coming. All photographs are suspect these days. Because now, with digital systems, they can be changed perfectly. Perfectly.”

What is the truth?

“The idea that photographs hand us an objective piece of reality, that they by themselves provide us with the truth, is an idea that has been with us since the beginnings of photography. But photographs are neither true nor false in and of themselves. They are only true or false with respect to statements that we make about them or the questions that we might ask of them.”

Errol Morris

<https://opinionator.blogs.nytimes.com/2007/07/10/pictures-are-supposed-to-be-worth-a-thousand-words>

What is the truth?

“Conceptually, we may call the truth what we cannot change; metaphorically, it is the ground on which we stand and the sky that stretches above us.”

Hannah Arendt

Arendt, Hannah. "Truth and politics." *Truth: Engagements across philosophical traditions* 295 (1967).

Manipulation in the era of film photography

Self Portrait as a Drowned Man (1840)



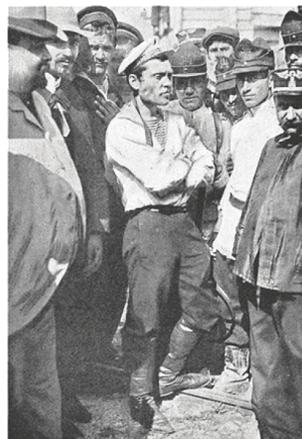
Self Portrait as a Drowned Man (Reverse)

“The corpse which you see here is that of M. Bayard, inventor of the process that you have just seen, or the marvellous results of which you are soon going to see. To my knowledge, this ingenious and indefatigable researcher has been working for about three years to perfect his invention. The Academy, the King, and all those who have seen his pictures, that he himself found imperfect, have admired them as you do at this moment. This has brought him much honour but has not yielded him a single farthing. The government, having given too much to M. Daguerre, said it could do nothing for M. Bayard and the unhappy man drowned himself. Oh! The fickleness of human affairs! Artists, scholars, journalists were occupied with him for a long time, but here he has been at the morgue for several days, and no-one has recognized or claimed him. Ladies and Gentlemen, you’d better pass along for fear of offending your sense of smell, for as you can observe, the face and hands of the gentleman are beginning to decay. H.B. 18 October 1840.”

Photo Retouching



Facial Retouching

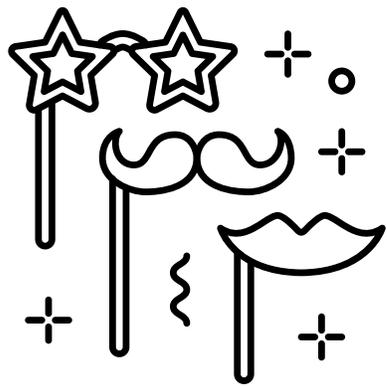


Cropping

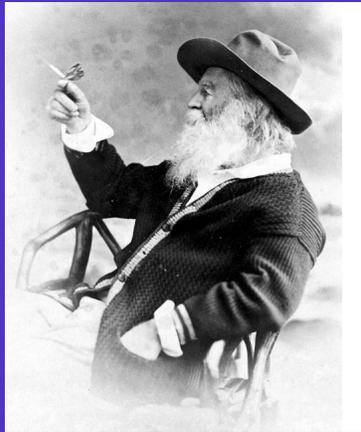


Face Swapping With Cutouts





Props



Photomontages

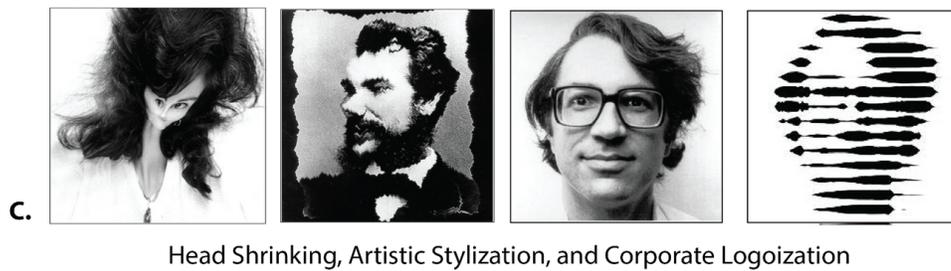
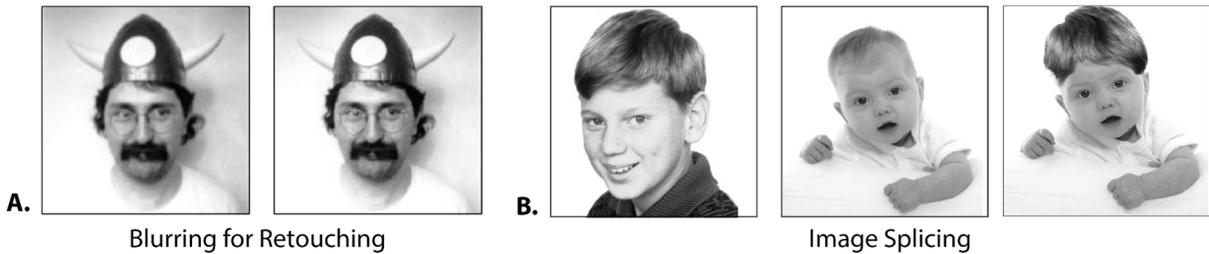


Double Exposure Photography



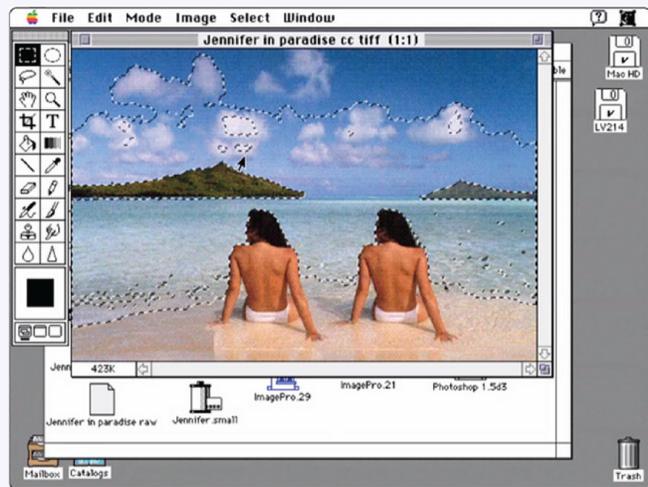
Manipulation in the early era of digital photography

Digital Darkroom (1988)

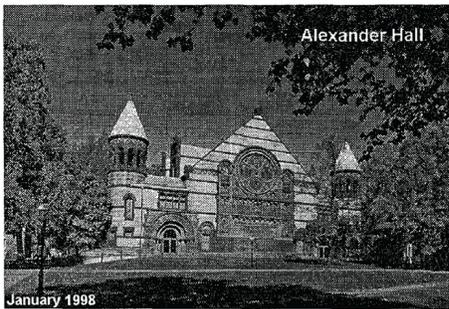
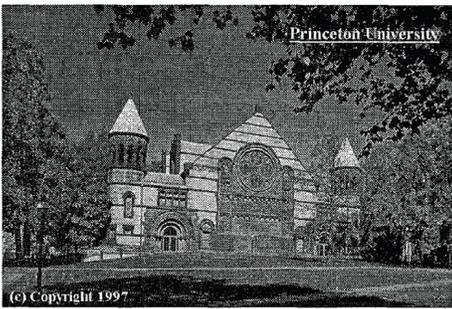


Holzmann, Gerard J. 1988. *Beyond Photography*. Hoboken: Prentice Hall.

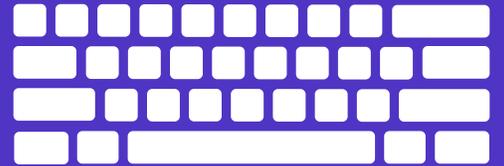
Photoshop (1988)



<https://www.youtube.com/watch?v=Tda7jCwvSzg>

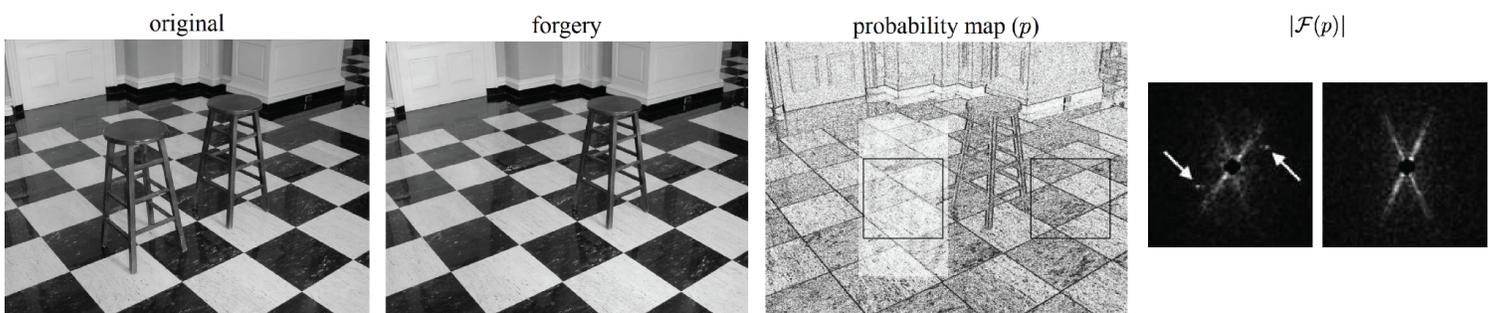


Wu, Min and Bede Liu. 1998. "Watermarking for Image Authentication." IEEE International Conference on Image Processing 2: 437-441.



Text Overlay Editing

Deletion of Objects (2004)



Popescu, Alin C. 2004. Statistical Tools for Digital Image Forensics. Dartmouth College Ph.D. Dissertation.

Face Swapping (2004)



Blanz, Volker, Kristina Scherbaum, Thomas Vetter, and Hans-Peter Seidel. 2004. "Exchanging Faces in Images." *Computer Graphics Forum* 23 (3): 669-676.

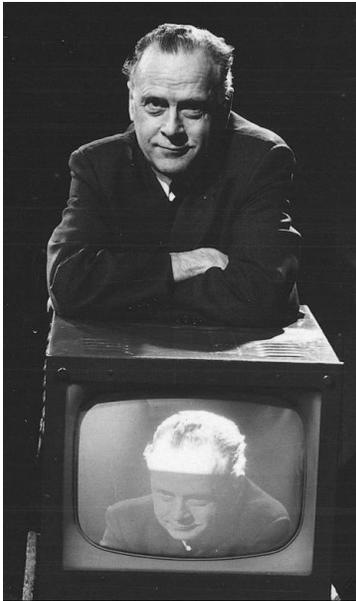
Memes and Manipulation (2001)



<https://www.wired.com/2001/11/hes-the-real-tourist-guy/>

The Frontier of the Imagination

Marshall McLuhan



"electric circuitry, *an extension of the central nervous system*"

Fiore, Quentin, and Marshall McLuhan. *The Medium is the Message*. Vol. 9. New York: Random House, 1967.

VS.

Information Superhighway



"Wander through a distant library. Turn your corner store into a multinational. Curious? *IBM can get you there.*"

<https://www.youtube.com/watch?v=GcoBY-GBBQM>



Something Awful's Photoshop Fridays (Early 2000s)



Artist: Lowtax (2001)

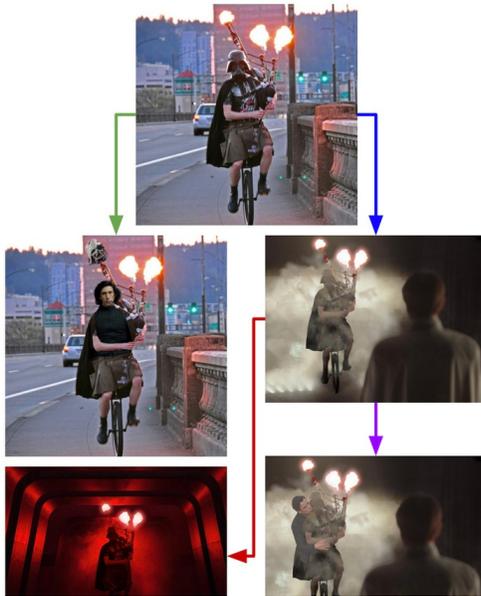


Artist: darthfunk (2002)



Artist: Groggie (2006)

Reddit's r/photoshopbattles (2012)



27.6k



Darth Vader Playing Bagpipe On Unicycle With Fire
Submitted 2 days ago by [blokmotion](#)
368 comments share save hide give gold report

[Shashakeitup](#) 1 8k 1588 points 1 day ago
Deleted scene from Rogue One

[Pwnagelad](#) 255 points 1 day ago
Should be the hallway scene

[Shashakeitup](#) 1 8k 1588 points 1 day ago
This one?

[LaerPoweredDeviltry](#) 125 points 1 day ago
Thats terrifying

[Owny_McOwnerton](#) 55 points 1 day ago
Red=Scary

[Artunitinc](#) P5 10 9k 330 points 1 day ago
<https://i.imgur.com/s8gbDxJ.jpg>

[spatulababy](#) 26 points 1 day ago
It's unfortunate how buried this one is

[snugglypatch](#) 889 points 1 day ago
I'm gonna finish what he started...
<https://i.imgur.com/OepNzqZ.png>

Generative Art Communities (Present)



<https://huggingface.co/spaces/akhaliq/ArcaneGAN>

Want to learn more?



Out Fall 2023 From Stanford University Press!

Academic Ghosting: Towards an Academy of Truth Telling

Alicia Andrzejewski



Academic Ghosting: Towards an Academy of Truth- Telling

Alicia Andrzejewski
The College of William & Mary
apandrzejewski@wm.edu @aliciaandr

THE CHRONICLE OF HIGHER EDUCATION

[NEWS](#) | [ADVICE](#) | [THE REVIEW](#) | [TOPICS](#) ∨ | [CURRENT ISSUE](#) | [VIRTUAL EVENTS](#) | [STORE](#) ∨ | [JOBS](#) ∨ | [Q](#)

“When Students Harass Professors”(2022)



“I do not dream of being able to swiftly remove students from my classroom. I dream about an academy where I can teach authentically and without fear. An academy where complaints from disempowered members of our community, whether instructors or students, are freely spoken. And an academy where all of us – not just those being scraped away – are invested in hearing, and addressing, these complaints.”

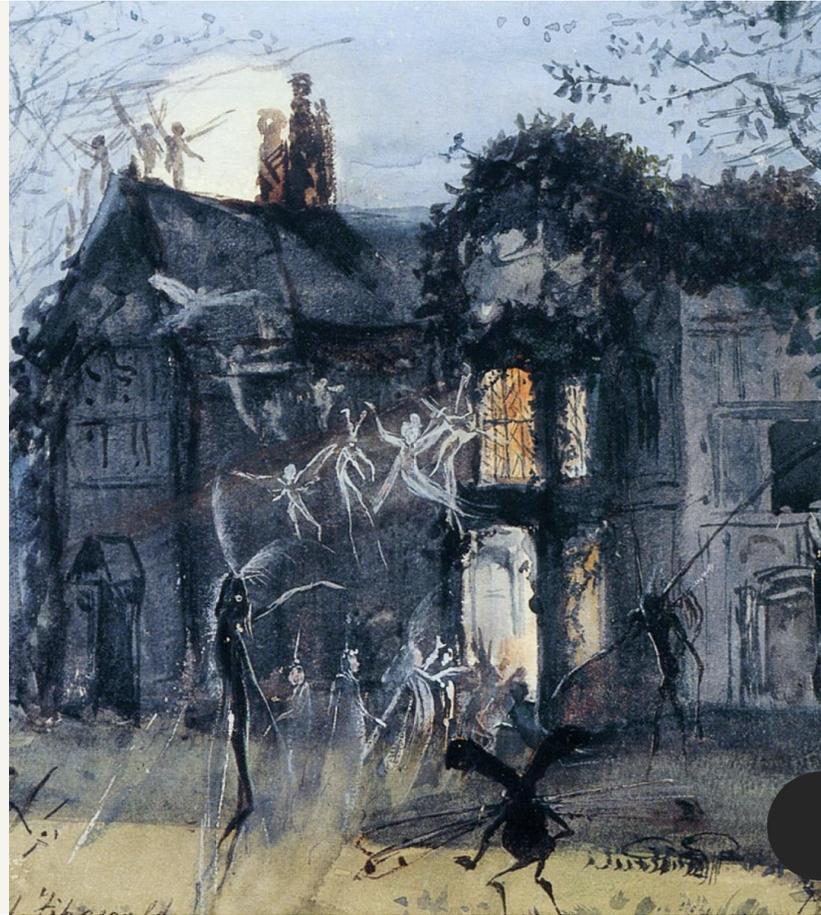
“Academics Don’t Talk About Our Mental Illnesses. We Should” (2023)



“Those of us who struggle with mental illness are desperate to fight the stigmas associated with our respective diagnoses, for others to know that our need for accommodation is not a personal failure, that we just operate, think, differently. And this is why it’s crucial to accommodate us: to diverge is to differ, to deviate, to continually depart from a standard, a norm. It follows, then, that neurodivergent academics are trailblazers; they offer the academy something more than already exists.”

“The Sad Humiliations of Academic Ghosting” (2023)

“How to hold ghosts accountable? How to relearn how to communicate with one another? How to face conflict in clear, direct, and truthful ways?”



Academic Life Podcast

Alicia Andrzejewski

Jun 8, 2023

Academic Ghosting

A Discussion with Alicia Andrzejewski

NEW BOOKS NETWORK 2023

Have you ever been ghosted in academia? The mentor who no longer replies when you reach out, the collaborators who mysteriously stopped collaborating with you, the search committee that said you were a top candidate and then stopped communicating with you—these are academic ghosts. They are people who are important to your career and suddenly stop responding to you without warning or explanation. What makes academic ghosting different than romantic ghosting? And why does it seem to hurt so much more? Dr. Andrea Andrzejewski joins us to explain:

- The systems in academia that make some forms of ghosting inevitable.
- What to do about it.
- The lingering pain and shame that being ghosted causes.
- Why your ghoster may reappear but they won't apologize.
- The ethical and financial reasons to address the issues that perpetuate ghosting.

An Extended Active Learning Approach to Multiverse Analysis

Daniel Kristanto, Carsten Gießing, Merle Marek, Changsong Zhou, Stefan Debener, Christiane Thiel, Andrea Hildebrandt



An Extended Active Learning Approach to Multiverse Analysis: Predictions of Latent Variables from Graph Theory Measures of the Human Connectome and Their Direct Replication^{*}

Daniel Kristanto, Carsten Gießing, Merle Marek, Changsong Zhou, Stefan Debener, Christiane Thiel, Andrea Hildebrandt[†]

Abstract

Multiverse analysis has been proposed as a powerful technique to disclose the large number of degrees of freedom in data preprocessing and analysis that strongly contribute to the current replication crisis in science. However, in the field of imaging neuroscience, where multidimensional, complex and noisy data are measured, multiverse analysis may be computationally infeasible. The number of possible forking paths given by different methodological decisions and analytical choices is immense. Recently, Dafflon et al. (2022) proposed an active learning approach as an alternative to exhaustively exploring all forking paths. Here, we aimed to extend their active learning pipeline by integrating latent underlying variables which are not directly observable. The extension to latent outcomes is particularly valuable for computational psychiatry and neurocognitive psychology, where latent traits are conceptualized as common cause of a variety of observable neural and behavioral symptoms. To illustrate our approach and to test its direct replicability, we analyzed the individual organization and topology of functional brain networks of two relatively large samples from the ABCD study dataset ($N = 1491$) and HCP dataset ($N = 833$). Graph-theoretical parameters that take into account both brain-wide and region-specific network properties were used as predictors of a latent variable reflecting general cognition. Our results demonstrate the ability of the extended method to selectively explore the multiverse when predicting a latent variable. First, the low-dimensional space created with the proposed approach was able to cluster the forking paths according to their similarity. Second, the active learning approach successfully estimated the prediction performance of all pipelines in both datasets. To interactively explore the multiverse of results, we developed a [Shiny app](#) to visualize the predictive accuracy resulting from each forking path and to illustrate the similarity between pipelines created by different combinations of data processing choice. The code for active learning and the app are available at the Github repository [ExtendedAL](#).

Keywords

Multiverse analysis, latent variable modeling, active learning, Shiny application.

Contents

| | |
|-------------------------------------------------------------------|----|
| Introduction | 1 |
| Methods | 2 |
| Brain Data | 2 |
| Behavioral Data | 2 |
| Multiverse Analysis Approach | 3 |
| An Active Learning Approach | 3 |
| Extension of the Method | 3 |
| The Multiverse of the Present Study | 5 |
| Implementation | 6 |
| Results | 6 |
| The search space from the proposed method | 6 |
| Active learning for guided multiverse analysis with SEM | 6 |
| Interactive visualization of multiverse outcome | 7 |
| Discussion | 7 |
| Conclusion | 10 |
| Acknowledgments | 10 |
| Code and Data Availability | 12 |
| Citation | 12 |
| References | 12 |

Introduction

The large number of options available to researchers for preprocessing and analyzing their data has been cited as one of the reasons for the replication crisis in science (Paul et al. 2022). A huge heterogeneity in data analysis has recently been reported in cognitive neuroscience based on functional magnetic resonance imaging (fMRI) data (Botvinik-Nezer et al. 2020). The complexity arising from the nature of such data,

^{*}Presented 2023-10-09 with [DK slides](#) and [DK video](#) at [Guardians 2023](#)

[†]DK, CG, MM, SD, CT, and AH affiliated with University of Oldenburg; CZ with Hong Kong Baptist University; correspondence to [D Kristanto at Univ Oldenburg](#).

characterized by inherent noise and multidimensionality, requires extensive pre-processing to remove machine and physiological artefacts. It further offers many different ways to parameterize the properties of such multidimensional data. This means, for example, that there are many ways to define the characteristics of brain networks from fMRI time-series data, increasing the variability of research results.

Multiverse analysis has been proposed as a promising approach to address this problem (Steegen et al. 2016), because it allows researchers to systematically explore different analytical choices, called forking paths, and report the multiplicity of their findings. The primary goal of a multiverse analysis is thus to assess the robustness of research findings, thereby reducing the likelihood of false-positive discoveries and mitigating the replication crisis. However, performing multiverse analysis presents its own challenges, particularly in network neuroscience, which deals with high-dimensional fMRI data and where the number of forking paths can be excessive (Dafflon et al. 2022; Botvinik-Nezer et al. 2020).

Recently, Dafflon and colleagues (2022) proposed an active learning-based approach to estimate the outcomes of multiple forking paths without the need for exhaustive computation of the multiverse (Dafflon et al. 2022). Their algorithm uses Bayesian optimization to sample a subset of forking paths and manually compute their outcome, and it uses Gaussian processes to estimate the outcome of the remainder. Dafflon et al.'s (2022) work applied this active learning approach to 1) predict brain age and 2) classify individuals with autism, using graph measures derived from fMRI-based whole brain networks. Both of these supervised learning problems are concerned with predicting an observed outcome variable.

However, in computational psychiatry and neurocognitive psychology, many outcome variables of interest cannot be measured directly and therefore reflect "latent" variables. To facilitate multiverse analyses in these fields, we extended Dafflon et al.'s (2022) active learning-based approach in two key ways. First, we augmented the approach with a predictive model that includes an endogenous variable that is latent and can be indicated by quantitative or ordinal measures. To accomplish this, we combined the proposed method by Dafflon et al. (2022) with Structural Equation Modeling (SEM, with latent variables) to infer predictive accuracy of brain measures with respect to a latent variable.

The original study by Dafflon used active learning to infer the prediction performance of each forking path without exhaustively sampling each of them. In short, active learning is an approach in machine learning where the model, during learning, can select the data that need to be labeled with the desired output (Settles 2009). SEM is a statistical analysis tool used to model the relationships between observed and latent variables (Kline 2015). In this study, the latent, non-directly measurable variable is general cognition g , which is estimated from various directly observable measures of performance on cognitive tasks (e.g., memory, reasoning and processing speed). In 1904, Spearman found that all indicators of cognition were positively correlated, referred to as the positive manifold, which is interpreted as general intelligence g (Spearman 1904). Recent studies have found that g is strongly associated with school achievement (for a review see Kriegebaum et al. 2018). Given the importance of g , a growing body of research in neuroscience has investigated the neural basis of g from the perspective of neurons (Bruton 2021), brain areas' activation patterns (Kovacs and Conway 2016; Jung and Haier 2007), and more recently, brain networks (Barbey 2018; Barabási et al. 2023).

In addition to combining Dafflon's method with SEM, we propose an

alternative approach to estimating forking path (pipeline) similarity that incorporates both brain-wide and region-specific graph measures. Note that the approach proposed by Dafflon et al. (2022) is only applicable to region-specific graph measures. Importantly, global graph measures, such as global efficiency and modularity, are relevant to predict behavioral outcomes in many applications (Alavash et al. 2015). To assess the effectiveness of our extended multiverse analysis approach and to test its direct replicability across datasets, we conducted a study on two large samples of NABCD = 1491 individuals from ABCD study and NHPC = 833 from HCP study (see Materials for the details of the datasets). Moreover, to better and more dynamically explore the multitude of results, we created an interactive visualization of the multiplicity of outcomes resulting from an exhaustive multiverse analysis using the Shiny app platform. The corresponding code is openly available at the Github repository [ExtendedAL](#).

Methods

We illustrate and test our multiverse analysis approach to predicting a latent outcome variable by examining the relationship between graph measures of the functional human connectome and g . We use data from two open-access studies. The first dataset was derived from the Adolescent Brain Cognitive Development (ABCD) study, the largest shared neuroimaging dataset to date. In order to obtain both brain and behavioral data, we used the ABCD data release 2.0.1. The second dataset was obtained from the Human Connectome Project (HCP) Young-Adult study.

Brain Data

Functional connectivity (FC) between brain regions was analyzed using functional magnetic resonance imaging (fMRI) data. Specifically, the blood oxygenation level dependent (BOLD) time series of different brain regions were measured as an indicator of underlying neuronal activations. The pairwise correlations between BOLD time series were calculated as an estimator of functional brain connectivity. For the ABCD dataset, we used previously preprocessed resting-state fMRI data (J. Chen et al. 2022) available at [NDA repository](#). Mean BOLD time series across voxels were extracted from a total of 419 brain regions, 400 cortical regions of interest (ROIs) from Schaefer's atlas (Schaefer et al. 2018) and 19 subcortical ROIs (Fischl et al. 2002). We computed the pairwise correlations between time series which resulted in functional connectivity matrices (419 x 419 brain areas) for 1491 individuals. The data have been preprocessed to remove motion-related, machine-related, and physiological noise (see J. Chen et al. 2022 for details). For the HCP dataset, we used data from our previously published study with NHPC = 833 individuals (Kristanto et al. 2023). In contrast to the ABCD dataset, time series of 360 brain regions were extracted according to the multimodal parcellation atlas and functional connectivity matrices with a dimensionality of 360 x 360 were calculated (Glasser, Coalson, et al. 2016). The MR data were also cleaned from artifacts using a minimal preprocessing pipeline from HCP (Glasser, Sotiropoulos, et al. 2013).

Behavioral Data

The aim of the present analysis was to explore the relationship between graph measures and a latent variable of general cognition, g . We used performance scores of five behavioral tasks available in both ABCD and HCP datasets: Picture Vocabulary (PicVocab), List Sorting Memory (ListSort), Pattern Comparison (PattComp), Picture Sequence (PicSeq),

and Reading Comprehension (Reading) (Fig. 1B). Picture Vocabulary and Reading Comprehension are known as indicators of crystallized intelligence, List Sorting memory is an indicator of reasoning ability, Pattern Comparison is an indicator of processing speed ability, and Picture Sequence is an indicator of memory.

All the tasks are part of the National Institute of Health Toolbox Cognition Battery NIH Toolbox. The details of the behavioral tasks are available from Casey et al. 2018 for ABCD dataset and Barch et al. 2009 for HCP dataset.

Multiverse Analysis Approach

Since the goal of this study is to extend a previously proposed method, in this section we first briefly describe the published method and point out the aspects that we aim to extend. Second, we explain our proposals for extending the aspects we point out from the original study. Third, we explain the design of our multiverse analysis and finally, we briefly explain the implementation of the extended method.

An Active Learning Approach

An active learning-based method for exploring the results in a multiverse analysis of predicting observed age of individuals from brain-based graph measures has been recently proposed (Dafflon et al. 2022). This method aims to estimate the analysis results (e.g., prediction performance) from a series of possible forking paths (544 and 384 forking paths in total for age prediction and autism classification, respectively) by sampling only a small fraction of them and inferring the other forking paths. The approach has been shown to be comparable to an exhaustive analysis, where each fork is executed sequentially to obtain the prediction.

Dafflon's et al. active learning approach for multiverse analysis consists of several steps. The first step is to prepare the brain and behavioral data and partition them into three sets. The first dataset is used to create the low dimensional space for embedding the similarity of graph measures obtained by all forking paths. The second dataset is used for prediction/classification, and the third dataset is for evaluation, thus to assess the performance of the best pipelines identified by active learning.

In the second step of the multiverse analysis method an embedding into a two-dimensional space of the forking path is created based on the similarity of the graph measures obtained by each forking path. This step is performed in the first dataset. In detail, the output of each forking path is a vector containing a graph measure of all brain regions. To obtain the forking paths' similarity, the cosine similarity of these vectors between any two individuals is computed, resulting in a similarity matrix of $(N) \times ((N)-1)/2 \times (NF)$, where (N) is the number of individuals and (NF) is the number of forking paths (for a detailed illustration, please refer to Fig. 8 of the original article (Dafflon et al. 2022)). Note that this step requires a vector as the output of all forking paths and is, therefore only applicable to region-specific graph measures. Next, the matrix is submitted to a dimension reduction algorithm to obtain a two-dimensional space of forking paths' similarity, which we refer to as the "search space" with the dimension of $2 \times (NF)$.

The third step – which aims at method evaluation – is to perform an exhaustive analysis to obtain a "true prediction accuracy" of each forking path. This step is performed on the second dataset. In the case of age prediction, this is done by predicting the age from the graph measures as the output of each forking path. The end result of this step is a vector of the true prediction accuracy of each forking path (e.g.,

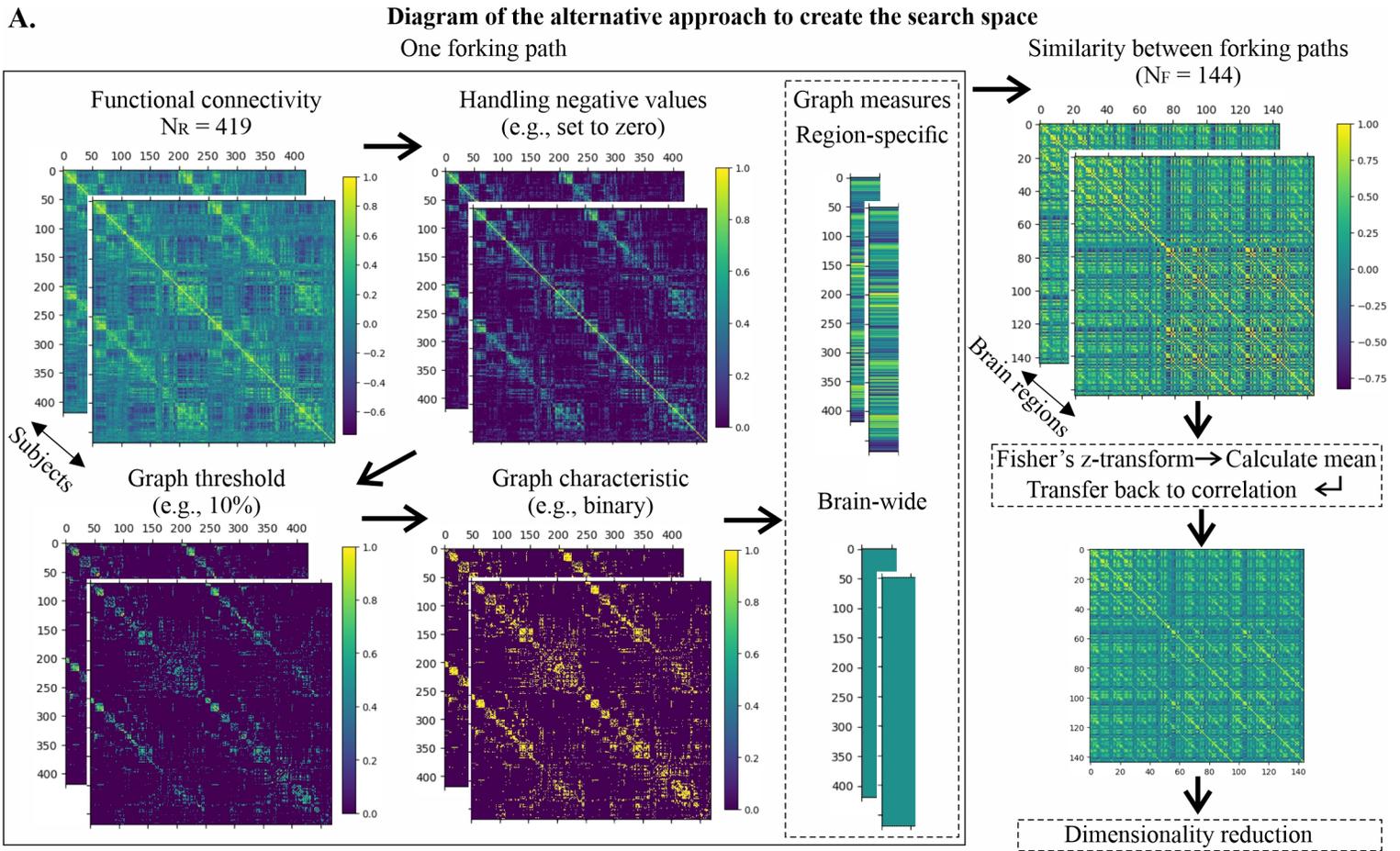
evaluated by the mean absolute error between predicted and actual age).

Finally, an active learning algorithm based on Bayesian optimization and Gaussian processes is implemented in the search space to infer the prediction accuracy of each forking path and to compare it with the true prediction accuracy obtained from the exhaustive search. The search space is a 2-dimensional space where the forking paths are represented as points in the space. The active learning first performs a burn-in phase where it randomly selects 10 points and evaluates their prediction accuracy. After this phase, more points (i.e., 40 points) are selected using Bayesian optimization and the prediction accuracy of those points is evaluated by predicting an observed outcome variable (age in the application) using graph measures derived from the corresponding forking paths. Finally, the Gaussian process is implemented to estimate the prediction accuracy of the other points/forking paths based on the selected points.

To test the robustness of the active learning, the whole analysis is repeated 20 times, each with different starting points. The third dataset is finally used to evaluate the prediction performance of the best pipeline identified by the active learning in different repetitions. However, we emphasize that here we are interested in using the active learning to estimate the prediction performance of all possible forking paths and in reporting the results from all forking paths.

Extension of the Method

Developing the search space that handles both brain-wide and region-specific graph measures: As a first extension, we proposed a different approach to generating the search space that allows the use of brain-wide and region-specific graph measures. Specifically, the use of brain-wide graph measures is in line with the current trend in behavioral neuroscience that aims to associate a measure from the whole brain to more general abilities such as *g*. It is important to note that the search space is the low-dimensional representation of the similarity (or dissimilarity) matrix between the forking paths. In the original study, the cosine similarity of the graph measures for the brain regions between all possible pairs of individuals was used to define the matrix. However, this approach is not applicable when the output of the forking path is a single value, which is the case for the brain-wide graph measures (e.g., global efficiency and modularity). The cosine similarity between two single values is always 1 ($\cos 180^\circ$) because they overlap and are on the same line. Therefore, no matter how much the forking paths differ when computing the brain-wide graph measures, the cosine similarity will always be 1 for any given pair of individuals. Notably, the absolute differences and mean absolute differences can also be used to replace cosine similarity, since they can handle graph measures with single and multiple values. However, since we include the options of graph measures in the forking paths, and thus different forking paths can have different graph measures, using (mean) absolute difference may not work when comparing forking paths with different graph measures. To overcome this challenge, we propose a different approach to construct the search space from the similarity matrix of the forking paths with brain-wide and region-specific graph measures, as shown in Fig. 1A. In detail, performing the analysis on the left side of Fig. 1A (all steps within the black box) for all forking paths results in a 3-dimensional matrix of size $(NF) \times (N) \times (NR)$, where (NF) is the number of forking paths, (N) is the number of subjects, and (NR) is the number of brain regions. Next, the similarity between a pair of forking paths is calculated for each brain region (right side of Fig. 1A). This is done by computing the



B. Diagram of the SEM for predicting general intelligence (g)

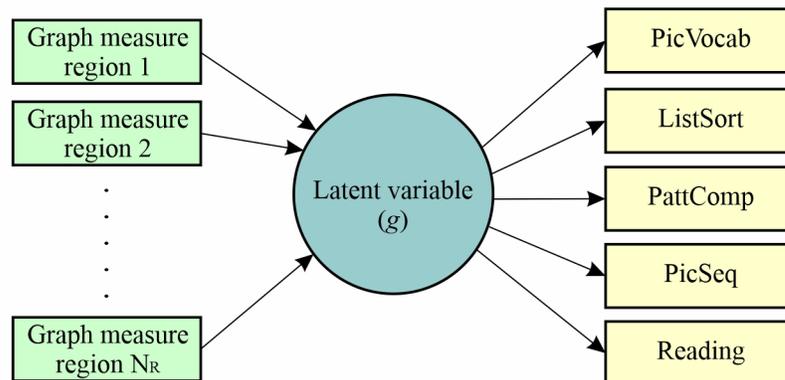


Figure 1: (A) Illustration of an alternative approach that handles both brain-wide and region-specific graph measures to create the search space. (B) A diagram of the SEM for predicting general cognition (g).

Pearson's correlation coefficient between two graph measure vectors of two forking paths in each brain region across individuals. The correlation coefficient fills the cells in the matrix at the top right of Fig. 1A. There are (NR) matrices, where each matrix has the size of $(NF) \times (NF)$. Importantly, the differences between forking paths with brain-wide and region-specific brain measures appear in these matrices of $(NF) \times (NF)$. Note that columns and rows represent forking paths. For the cells with both column and row representing forking paths with brain-wide graph measures, the values are identical across brain regions ((NR)) since the brain-wide graph measures do not vary across regions. In contrast, for the cells with column and/or row representing forking paths with region-specific graph measures, the values are different across brain regions since the graph measures are different across regions. Next, the average of these matrices across brain regions is computed by first transforming them using Fisher's z-transformation to account for the nonlinearity of the correlation coefficients (Silver and Dunlap 1987). Notably, the averaging only affects the cells containing forking paths with region-specific graph measures, but not the forking paths with brain-wide graph measures, since the values are identical across brain regions. The averaged matrix is transformed back into a correlation matrix before being subjected to the dimension reduction approach to obtain the search space of size $2 \times (NF)$.

In particular, a step to reduce the dimension (i.e., number of features) of each forking path is a key step to perform the active learning algorithm with small number of observations (i.e., 144 forking paths). The step allows the creation of a low-dimensional search space for the active learning. Using the similarity matrix (with a size of 144×144 , where 144 is the number of forking paths) as the active learning search space will require much more observations/forking paths which are distributed across space in order to implement the active learning, since each forking path will have 144 features or dimensions. Here, we showed that, although we reduced the number of dimensions, the results of the search space in the lower 144×2 space, still retains the information related to the similarity between forking paths, where similar forking paths are located close to each other. In addition, we also found that the search space with only two dimensions allows active learning to mimic the results of the exhaustive search where we computed the performance of each forking path (see Results).

Integrating active learning approach with SEM: The second extension of the multiverse analysis proposed here involves integrating the SEM model with active learning to infer the explained variance of a latent outcome variable from graph measures. Unlike the method in the original multiverse analysis (Botvinik-Nezer et al. 2020) which used only observed variables, the extension allows for inferring latent outcome variables commonly of interest in computational psychiatry and neurocognitive psychology. To do so, we replace the prediction model with SEM implemented using the *semopy* package (Igolkina and Meshcheryakov 2020; Meshcheryakov et al. 2021). The predictive performance of each forking path is evaluated by the explained variance of the latent variable from graph measures. Therefore, the exhaustive search will output a "true prediction accuracy" vector, which is indicated by the explained variance of each forking path. Active learning is then performed to infer this value. Fig. 1B shows a schematic of the SEM used in the present study to evaluate the proposed approach. The latent variable reflects general cognition, indicated by 5 items behavioral measures (see Materials). The latent variable is then predicted by graph measures of the connectome as described above. The explained variance of the latent variable becomes the predictive performance of each forking

path.

Interactive visualization of the multiverse analysis results: An interactive visualization of the results of all forking paths is the final extension of the multiverse analysis approach proposed here. For this we used the *Shiny* package (Chang et al. 2023) in R the software for statistical computing (R Core Team 2021). A graph visualization based on force network *uner networkD3* package (Allaire et al. 2022) was created to represent the multiverse where the nodes are the forking paths and the edges are the relationships between the forking paths. Notably, we set the node size to represent the prediction performance of the corresponding forking path. A larger node indicates that the corresponding forking path has higher prediction performance. Moreover, the relationship between forking paths is represented by the similarity between them, which was taken from the matrix of average similarity between forking paths (Fig. 1A). We also added some features to the shiny application. First, hovering the mouse over the nodes will trigger the name of the forking paths and all other connected forking paths. Clicking on a node brings up a dialog box with the corresponding forking path and its prediction performance. Hovering over the edges will show the degree of similarity between connected forking paths. Finally, we also incorporated some option buttons where the user can select a specific forking path and explore other paths connected to it. In addition, a slider option allows users to specify the threshold of similarity between forking paths (e.g., to find the forking paths that are connected by a correlation coefficient of at least 0.8).

The Multiverse of the Present Study

In this study, a multiverse analysis was performed in which a latent variable of g was predicted from graph measures derived from fMRI data. First, the investigated forking paths were identified through a systematic literature review on the multiverse of fMRI data preprocessing and fMRI graph analysis steps (paper in preparation). This multiverse covers a wide range of pre-processing and analysis steps in fMRI-based graph analysis including structural image pre-processing, functional image preprocessing, noise/artifact removal, functional connectivity definition, and network definition. In this study, we focus only on the small fraction of pre-processing paths which we call data multiverse. Analysis paths were pre-dominantly selected if the corresponding options are variable across studies and are highly controversial:

- *Paths for handling negative correlations:* Use absolute values, keep negative values, assign 0 values to the negative values, discussed in G. Chen et al. 2011;
- *Paths for Controlling graph sparsity:* 50%, 30%, and 10%, discussed in Franco 2022. Notably, these options cover a variety of network sparseness where '50%' represents a relatively dense network, while 10% represents a sparse network with only 10 percent of all possible connections;
- *Paths for defining graph edges:* weighted and binarized, discussed in Xiang et al. 2020;
- *Paths for computing graph measures:* strength, betweenness centrality, clustering coefficient, eigenvector centrality, local efficiency, global efficiency, modularity, and participation coefficient. for detail of each measure please refer to Brain Connectivity Toolbox (Rubinov et al. 2009).

In total, there are $3 \times 3 \times 2 \times 8 = 144$ different forking paths. Note that modularity, global efficiency, and participation coefficient are brain-wide graph measures, while the rest are region-specific graph measures.

Implementation

Overall, the implementation of the proposed method is similar to the original study. In the first step we divided the data into 3 sets by keeping the similar ratio of the original study (seimilar number for set 1 and set 3, and higher number for set 2). Since we had a larger sample size in both datasets than the original study, for ABCD dataset, we used 350 individuals to define the search space, 791 individuals to perform the prediction based on the SEM, and 350 individuals to validate the best-performing forking paths identified by active learning. We kept the ratio for the HCP dataset: 200 individuals to define the search space, 433 individuals to perform the prediction based on the SEM, and 200 individuals to validate the best-performing forking paths identified by active learning. Notably, we are more interested in reporting the results from all possible forking paths and not in finding the best-performing forking paths. The validation here further confirms the robustness of the results obtained by active learning from different repetitions using different forking paths in burn-in/initialization phase.

To create the search space, we followed the pipeline shown in Fig. 1A. In line with the result of the original study, we applied multidimensional scaling (MDS) as a dimension reduction method, because this embedding approach was shown to perform best in Dafflon et al. 2022. Further embedding methods can be explored in the future. The exhaustive search was then performed to obtain the "true prediction performance" vector of the explained variance of the latent variable for all forking paths. Notably, for the SEM model to predict the latent variable by the graph measures, we only used the brain areas in dorsal attention and fronto-parietal networks (Schaefer et al. 2018; Yeo et al. 2011). The areas of these networks were found to be associated with g (Jung and Haier 2007; Hilger et al. 2020). In total, we have 77 and 89 brain areas for the ABCD and the HCP datasets, respectively, to predict g . Active learning was then performed to quantify the explained variance in the latent variable, which was further compared with the result of the exhaustive search. Since we had a smaller number of pipelines as compared with the original study, we used different numbers of forking paths to train the active learning (i.e., 10 points randomly selected for the burn-in phase and 20 points selected using Bayesian optimization). We set the active learning to be exploratory, with a kappa of 10, following the result of the original study. We also run the active learning 20 times with different starting points to evaluate its robustness. For each iteration, we identified the best performing forking path. The robustness is indicated by the replicated best-performing forking paths identified across repetitions. Finally, we used the third dataset to validate the prediction performance of those best-performing forking paths.

Results

The search space from the proposed method

The first result of the present study was the creation of a low-dimensional space (search space) by using the newly proposed approach to deal with both brain-wide and region-specific brain measures. As described in the original study (Botvinik-Nezer et al. 2020), the creation of the search space mainly aims to capture the similarity of the forking paths in the 2-dimensional space with the constraint that similar forking paths should stay close to each other.

The search space generated by our proposed approach shows that the forking paths are well distributed in the space, as shown in Fig. 2 for both datasets. More importantly, there is considerable structure in the location of the different forking paths, meaning that especially the similar types of graph measures, illustrated by different shapes, are generally proximal. Moreover, we see a clear distinction between the forking paths calculating graph measures related to integration (e.g., global efficiency and participation coefficient) versus segregation (e.g., modularity and local efficiency). This observation is true for both data sets. For the ABCD dataset (Fig. 2A), the forking paths of the integration graph measures are mostly located in the upper right part of the space, while the forking paths of the segregation graph measures are mostly located in the lower left part of the space. For the HCP dataset (Fig. 2B), the forking paths of integration graph measures are mostly located in the upper part of the space, while the forking paths of segregation graph measures are mostly located in the lower part of the space. This finding suggests that in both dataset, the proposed method to create a two-dimensional search space, was able to cluster the forking paths according to their similarity. In order to assess the similarity between the spaces from both datasets, we performed clustering analysis in the low-dimensional spaces from both datasets and computed the Rand Index ((RI)) of the clustering results, which is the ratio between the number of matching pairs and the number of pairs (Hubert and Arabie 1985). The (RI) value of '0' indicates that two clustering results are completely different, while a value closer to '1' represents high agreement between two clusterings. We defined the optimal number of clusters using the Elbow method based on the intra-cluster sum of squares, also known as inertia (Thorndike 1953). For both datasets, we found that the optimal number of clusters is 3. An (RI) value of 0.71 was determined for these three clusters, indicating that the clustering results in the low-dimensional spaces of the two data sets were very similar.

Given these results, first we conclude that the proposed similarity evaluation approach is suitable to create a low-dimensional space that serves as the search space for active learning where similar forking paths tend to be close to each other. Second, we found that the proposed approach can be replicated in a different dataset in terms of clustering results from the low-dimensional space. It is also important to note that the slight differences in terms of clustering results from the low-dimensional space between the data sets may be due to the differences in terms of the number of individuals available to create the low-dimensional space and the number of nodes in the FC (419 and 360 nodes for ABCD and HCP datasets, respectively).

Active learning for guided multiverse analysis with SEM

We implemented the proposed extensions of the guided multiverse analysis on the ABCD study and HCP datasets to predict the latent variable reflecting g using graph measures from fMRI data. The results are shown in Fig. 3. First, Fig. 3A captures the prediction performance of all forking paths when the exhaustive search (i.e., manual execution of all forking paths) was performed, left panel is for ABCD dataset and right panel is for HCP dataset. It can be seen that region-specific graph measures outperform brain-wide graph measures in explaining the variance of the latent variable. The prediction performance shown in Fig. 3A serves as the "true prediction performance" to be used to evaluate the performance of the active learning to guide the multiverse analysis.

Next, Fig. 3B shows how active learning selects the training points (=forking paths) in 5, 10, 15, 20, and 30 iterations and infers the predic-

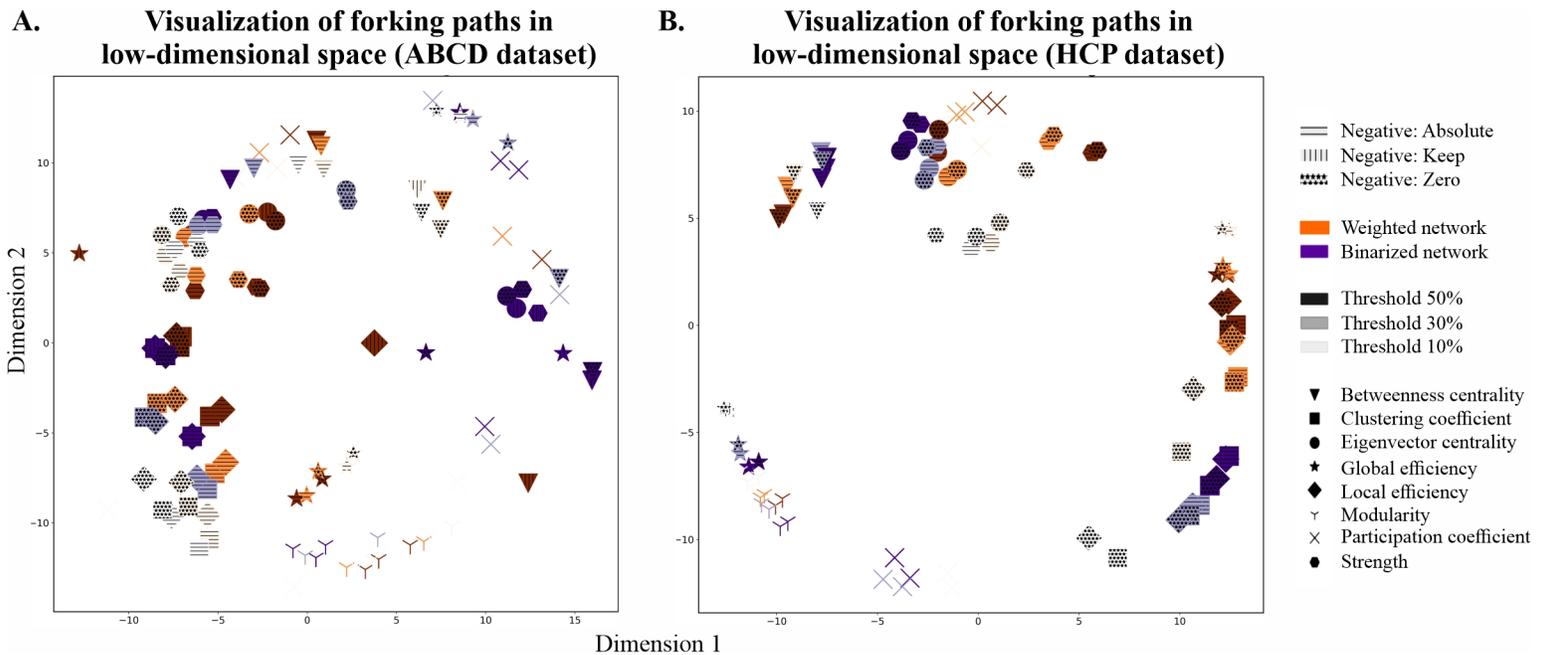


Figure 2: The visualization of forking paths in low-dimensional space for the ABCD dataset (A) vs. the HCP dataset (B). The items with different features (i.e. texture, color, shade, and symbols) represent the forking paths (NF = 144) in two dimensional space (see Methods).

tion performance of the search space in different iterations, left panel is for ABCD dataset and right panel is for HCP dataset. After 30 iterations (e.g., 30 training forking paths are selected), the active learning can satisfactorily mimic the prediction performance of all the forking paths from the exhaustive search. Notably, the Spearman correlations of the prediction performance of all forking paths between exhaustive search and active learning are 0.69 and 0.75 for ABCD and HCP datasets, respectively. This shows that the rank orders of the forking paths in terms of prediction performance obtained by exhaustive search and active learning are sufficiently similar. Next, we ran the active learning over 20 repetitions, where each repetition randomly selects different training points. Fig. 4 shows the comparison of prediction performance between exhaustive search and active learning across repetitions in the ABCD (Fig. 4A) and the HCP (Fig. 4B) datasets. The figures on the left panel are line plot of the Mean Absolute Error (MAE) of the prediction performance (explained variance of g) between the active learning and exhaustive search in 20 repetitions with different points for training. For both datasets, we found the MAE is around 0.05 for all repetitions. The figures on the right panel show the distribution of Spearman correlations of the prediction performance of all forking paths between exhaustive search and active learning across 20 repetitions. The correlations, which range from 0.37 to 0.75 for ABCD dataset, and from 0.58 to 0.77 for HCP dataset, indicate that active learning robustly mimics the prediction performance of all forking paths obtained from exhaustive search. The robustness of the active learning is also supported by the identification of similar best-performing forking paths across 20 repetitions. For the ABCD dataset the best-performing forking path is the following: keep the negative correlations, use a sparse network with a threshold of 10%, use a weighted network, and compute local efficiency as the graph measure. The analysis on HCP dataset identified a similar best-performing forking path across 20 repetitions: keep the negative correlations or set them to zero, use a sparse network with a threshold of 10%, use a weighted network, and compute either local

efficiency or betweenness centrality as the graph measures.

Interactive visualization of multiverse outcome

A screenshot of the interactive application (available online at <https://meteor-oldenburg.shinyapps.io/ExtendedAL/>) is shown in Fig. 5. Note that the results displayed in the interactive visualization originated from the ABCD dataset. When the threshold for the relationship (similarity, right bottom matrix in Fig. 1) between the forking paths is set to a higher threshold (e.g., 0.7), the clusters of the forking paths are shown based on the corresponding graph measures. Consistent with the search space discovery, the interactive visualization also shows that the similar forking paths (e.g. with similar graph measures) are highly correlated (connected) with each other. The user can also explore different thresholds to investigate the relationships between the forking paths. Moreover, the user can also select a particular forking path to explore how it is connected to other forking paths. For example, as also illustrated in the red box in Fig. 5, selecting the “betweenness centrality weighted 0.1abs” forking path will show the other forking paths that are connected to it.

Discussion

The present study has extended a previously proposed method for a guided multiverse analysis of fMRI based graph theory measures to predict behavioral outcomes (Dafflon et al. 2022). We show that our extensions perform well, allowing the use of both brain-wide and region-based graph measures and the prediction of latent variables from fMRI-based graph measures. In addition, since our goal is to report on the multiplicity of findings, we propose an interactive visualization of the results of the multiverse analysis, where the user can explore the outcomes obtained by all possible forking paths.

We first discuss the methodological contribution of the present study. The originally proposed method (Dafflon et al. 2022) is a valuable contribution to research when it comes to addressing the issue of the

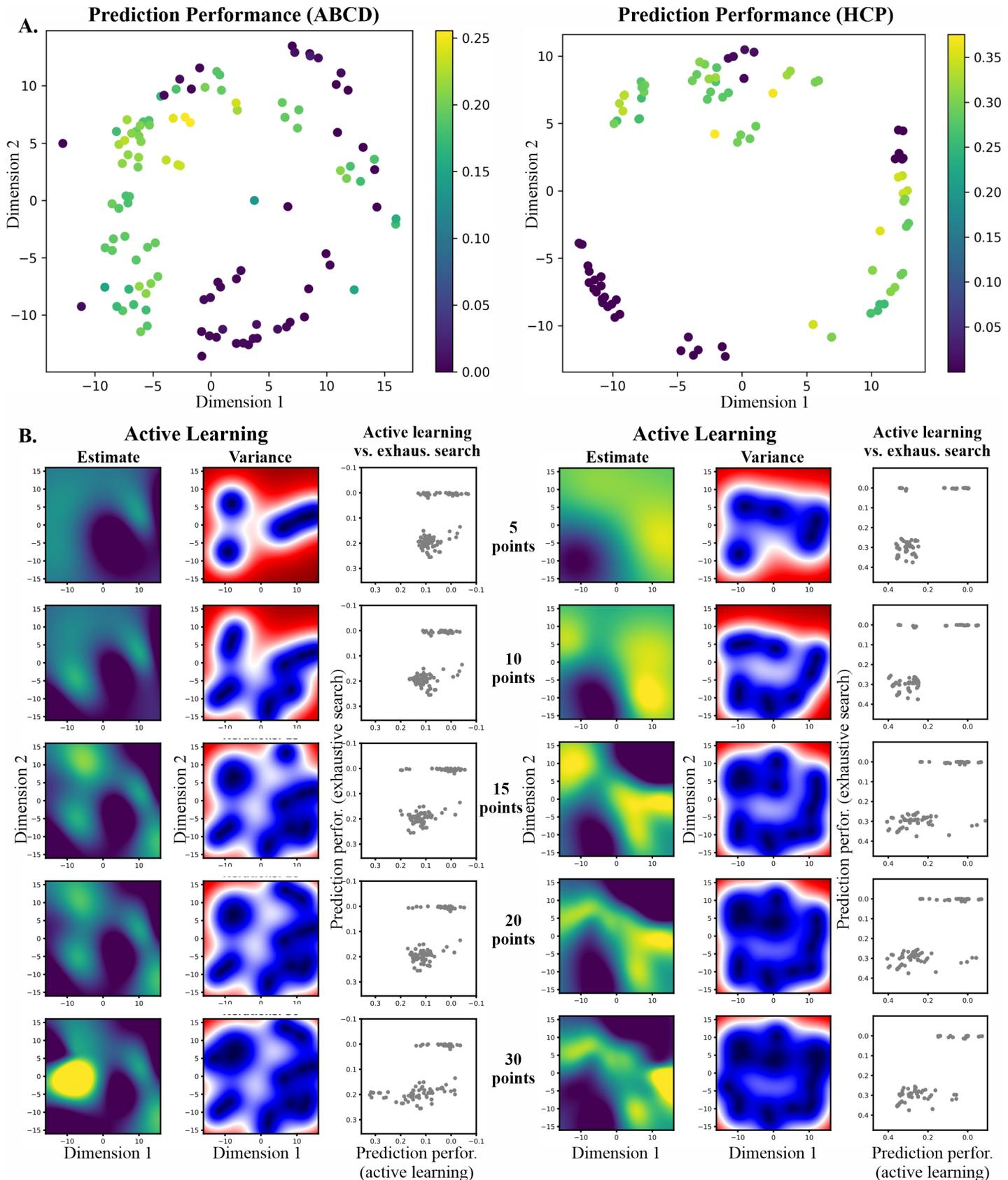


Figure 3: The prediction performance of forking paths estimated by exhaustive search and active learning; left panel illustrates results from the ABCD dataset and the right panel from the HCP dataset. (A) The prediction performances of all forking paths were also computed (exhaustive search) and served as the ground truth for the active learning. (B) The active learning process to estimate the prediction performance of all forking paths is visualized for different sampled points (5, 10, 15, 20, and 30 points). The first column is the estimated prediction performance of the space. After 30 iterations, the estimation is comparable with the ground truth (Fig. 3A). The second column indicates which points have been sampled. The third column is the prediction performance from the active learning versus exhaustive search for all forking paths.

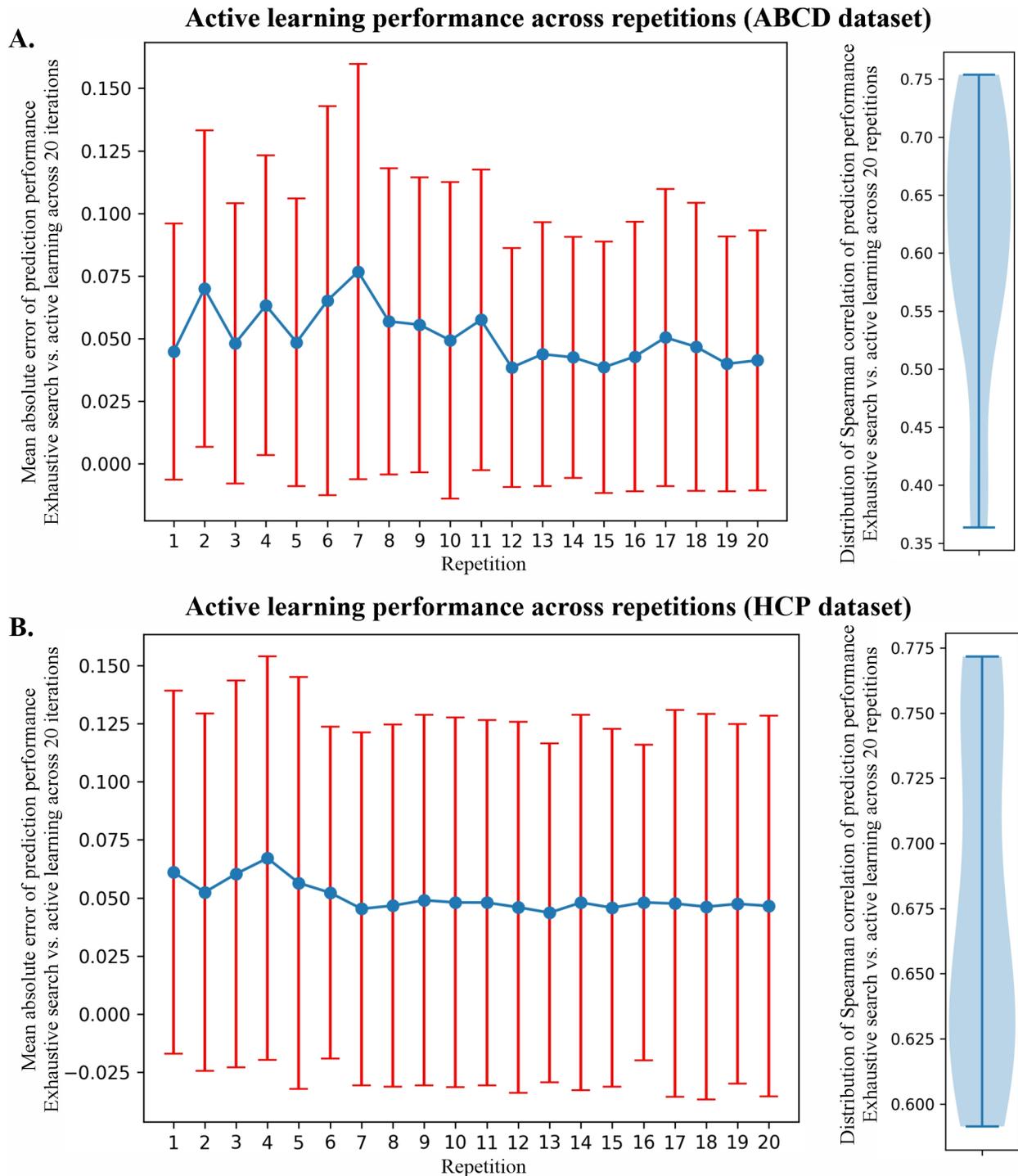


Figure 4: Prediction performance of active learning compared to exhaustive search across different repetitions for ABCD dataset (A) and HCP dataset (B). The left panel is the mean absolute error of prediction performance (explained variance of g) of all forking paths obtained by active learning and exhaustive search across repetitions. The right panel is the distribution of the Spearman correlations between the prediction performance of all forking paths obtained by exhaustive search and those from active learning across repetitions.

researcher's degree of freedom, as it provides an efficient method to explore the multiverse of possible decisions in a guided, well informed way. In line with this goal, we believe that this method can be further extended to target a larger user community by adding more flexibility and multiple features to it. Here, we added flexibility by allowing the combination of brain-wide and region-specific graph measures and added further modelling capabilities by integrating SEM as a predictive model to explain a latent outcome variable. Next, we also emphasize that the end product of this guided active learning method is not only to identify the best performing forking paths, but also to provide a full report of all possible forking paths and to gain knowledge about the sensitivity of the results with respect to the decisions of different researchers. In this sense, we introduced an interactive Shiny application to visualize not only the results of all forking paths, but also how the forking paths are related to each other in terms of similarity between the graph measures obtained.

Second, we now discuss possible future applications building upon the results of our study. Note that we predicted the latent variable using only a subset of brain areas related to cognitive control and risk-taking behaviors, and thus did not apply any regularization within the prediction model. A future study may include more predictors, such as areas from the whole brain, and apply regularization in the prediction model aiming to extract the important brain areas for prediction. We noticed that regularization options are also available in the *semopy* package (Meshcheryakov *et al.* 2021) which can be explored in the future. Relatedly, other features in *semopy* can also be explored to elaborate on more advanced SEM models which are potentially relevant to test brain-behavior associations. Furthermore, the forking paths we use in this study are also limited to data processing steps after functional connectivity definition. Follow-up studies may also consider additional forking paths related to fMRI data preprocessing in spatial or temporal domains, especially those dealing with noise removal. For these preprocessing decisions, a slightly different choice may contribute to significantly different results. In addition, other forking paths may be found beyond the data pre-processing domain or data multiverse. For example, one could consider different ways to compute the similarity matrix across forking paths and whether the Fisher's z-transformation is part of the method multiverse.

On a related note, the number of training data (selected points/-forking paths) can also be one important forking path for the method multiverse since the performance of the active learning model also depends on the selected training points. For an example, we performed a small analysis with two different scenarios to select 30 training points: (i) 5 points selected randomly and 25 points selected via Bayesian optimization and (ii) 30 points selected randomly. We found that after completing the training with 30 selected points, the Spearman correlations of the prediction performance of all the forking paths between the active learning and exhausting search are 0.48 and 0.58 for case (i) and case (ii) in ABCD dataset and 0.57 and 0.60 for case (i) and case (ii) for HCP dataset. For reference, the Spearman correlations from the original approach (10 points randomly selected and 20 points via Bayesian optimization) are 0.67 and 0.75 for ABCD and HCP datasets, respectively. The results suggest that variability of the outcomes may occur with different methods to select the learning points. Therefore, future studies may consider a more systematic multiverse analysis that considers the forking paths in the methodological steps, or the multiverse of the method (method multiverse). Moreover, the use of different behavioral measures to define a latent variable (e.g. g), or the multiverse of the

outcome variable (outcome multiverse), can also be explored.

Notably, there are steps/forking paths with a huge number of options (the one with continuous variables or the one with infinite discrete values) or with options that are mutually exclusive. In those cases, the number of forking paths exponentially grows and the implementation of multiverse analysis may not be computationally feasible. The approaches to define the garden of forking paths, or generally the question how to correctly and efficiently perform a multiverse analysis, are still being discussed. Del Giudice and Gangestad 2021, proposed that assessment of equivalence of the options/forking paths in terms of measurement, effect, and power/precision may help to reduce the number of forking paths. A sampling method across all the forking paths can also be conducted to reduce the number of forking paths (Paul *et al.* 2022).

Finally, the visualization application is still being improved, both visually and technically, e.g., by adding more features that allow users to interact more easily with the multiverse. Moreover, integrating the visualization app and the active learning approach into one toolbox will be an important contribution to studies on multiverse analysis. Especially, the implementation into a toolbox with graphical user interface (GUI) that can be generalized across different research domains will facilitate the application of multiverse analysis in different fields.

Conclusion

Guided multiverse analysis (Dafflon *et al.* 2022) is necessary in fields dealing with complex data structures, such as graph-theory fMRI based brain-behavior association research. Such associations are widely studied in computational psychiatry and neurocognitive psychology, where the behavioral variables of interest are inherently latent. Our extension of the guided multiverse analysis method (Dafflon *et al.* 2022) makes the approach suitable for a broader community interested in assessing the robustness of findings across a large number of possible analytical choices when predicting a latent variable with graph theory fMRI measures.

Acknowledgments

This work was supported by a grant from the German Research Foundation (DFG) to Andrea Hildebrandt (HI 1780/7-1), Carsten Gießing (GI 682/5-1), Stefan Debener (DE 779/8-1) and Christiane Thiel (TH 766/9-1) as part of the DFG priority program "META-REP: A Meta-scientific Programme to Analyse and Optimise Replicability in the Behavioural, Social, and Cognitive Sciences" (SPP 2317) and a fellowship awarded from the Hanse-Wissenschaftskolleg (Institute of Advanced Study) in Delmenhorst, Germany, and the School for Medicine and Health Sciences at the University of Oldenburg to Daniel Kristanto. Data used in the preparation of this article were obtained from the Adolescent Brain Cognitive Development (ABCD) Study <https://abcdstudy.org>, held in the NIMH Data Archive (NDA). This is a multisite, longitudinal study designed to recruit more than 10,000 children aged 9–10 and follow them over 10 years into early adulthood. The ABCD Study® is supported by the National Institutes of Health and additional federal partners under award numbers U01DA041048, U01DA050989, U01DA051016, U01DA041022, U01DA051018, U01DA051037, U01DA050987, U01DA041174, U01DA041106, U01DA041117, U01DA041028, U01DA041134, U01DA050988, U01DA051039, U01DA041156, U01DA041025, U01DA041120, U01DA051038, U01DA041148, U01DA041093, U01DA041089,

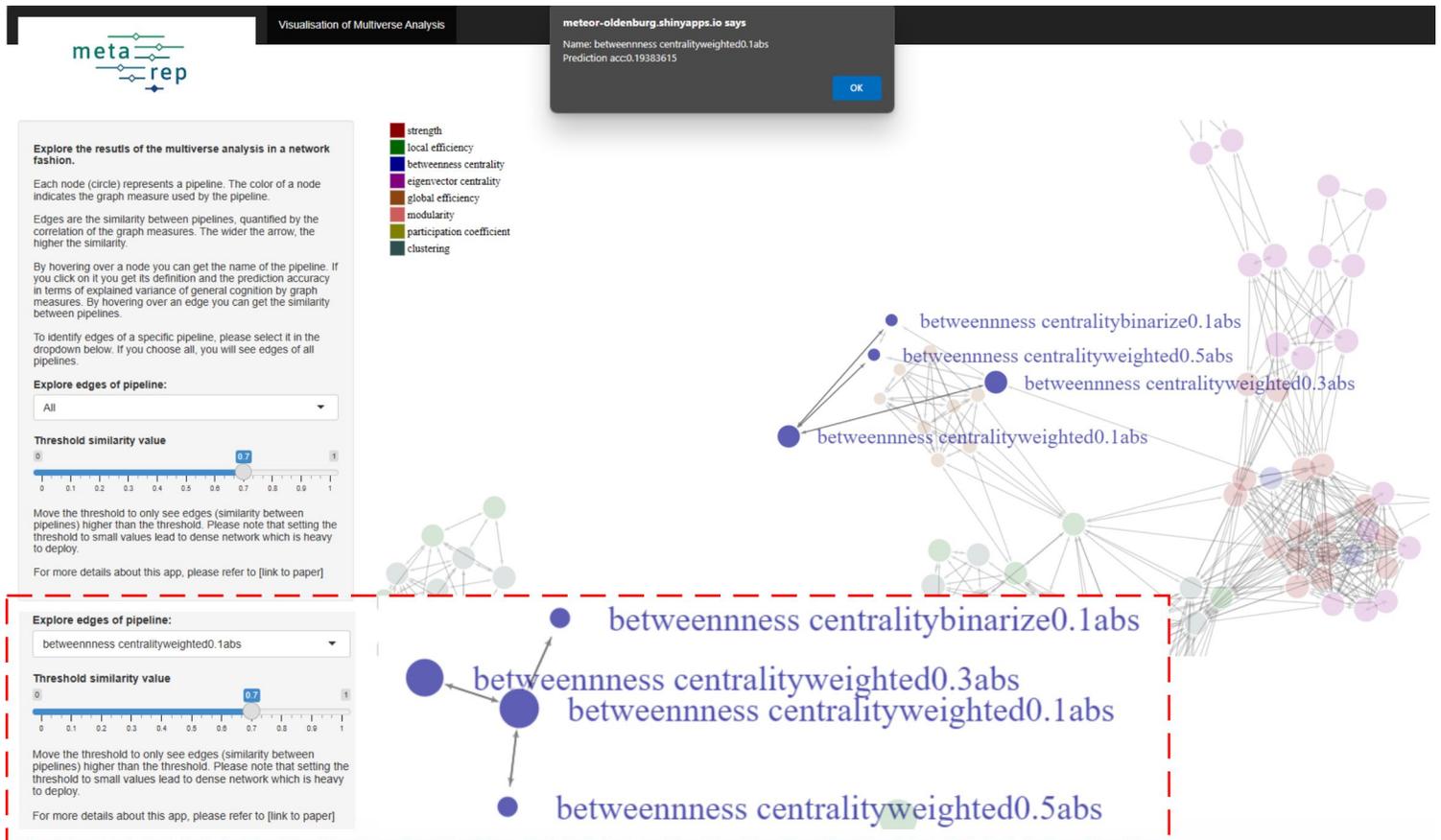


Figure 5: Screenshot of the interactive Shiny application for the visualization of the multiverse of analysis results. The multiverse is visualized as a network where the nodes are the forking paths and the edges indicate the similarity between the forking paths.

U24DAO41123, U24DAO41147. A full list of supporters is available at <https://abcdstudy.org/federal-partners.html>. A listing of participating sites and a complete listing of the study investigators can be found at https://abcdstudy.org/consortium_members/. ABCD consortium investigators designed and implemented the study and/or provided data, but did not necessarily participate in the analysis or writing of this report. This paper reflects the views of the authors and may not reflect the opinions or views of the NIH or ABCD consortium investigators. The ABCD data repository grows and changes over time. The ABCD data used in this report came from <https://doi.org/10.15154/1504041>.

Code and Data Availability

Codes to run the analysis of this study are available at <https://github.com/kristantodan12/ExtendedAL>. Computed graph measures, behavioral scores, and the list of pipeline for HCP dataset are available at <https://github.com/kristantodan12/ExtendedAL>, while for ABCD dataset are available at <https://nda.nih.gov/study.html?id=2428>.

Citation

Brainiacs 2023 Volume 4 Issue 2 Edoc J962E0F53

Title: "An Extended Active Learning Approach to Multiverse Analysis: Predictions of Latent Variables from Graph Theory Measures of the Human Connectome and Their Direct Replication"

Authors: Daniel Kristanto, Carsten Gießing, Merle Marek, Changsong Zhou, Stefan Debener, Christiane Thiel, Andrea Hildebrandt

Dates: received 2023-09-22, presented 2023-10-09, published 2023-12-21, endorsed 2023-12-30

Copyright: © 2023 Brain Health Alliance

Contact: [D Kristanto at Univ Oldenburg](mailto:D.Kristanto@univ-oldenburg.de)

URL: [BrainiacsJournal.org/arc/pub/Kristanto2023MVMRIA](https://brainiacsjournal.org/arc/pub/Kristanto2023MVMRIA)

PDP: [/Nexus/Brainiacs/Kristanto2023MVMRIA](https://nexus.brainiacs.org/Kristanto2023MVMRIA)

DOI: [10.48085/J962E0F53](https://doi.org/10.48085/J962E0F53)

References

- [1] M. Alavash, C. C. Hilgetag, C. M. Thiel, and C. Gießing. "Persistence and flexibility of complex brain networks underlie dual-task interference." *Human Brain Mapping* 36.9 (2015), pp. 3542–3562. DOI: [10.1002/hbm.22861](https://doi.org/10.1002/hbm.22861) (cited p. 2).
- [2] J. J. Allaire, C. Gandrud, K. Russell, and C. J. Yetman. *networkD3: D3 JavaScript Network Graphs from R*. 2022. URL: <https://cran.r-project.org/web/packages/networkD3/networkD3.pdf> (cited p. 5).
- [3] D. L. Barabási, G. Bianconi, E. Bullmore, M. Burgess, et al. "Neuroscience Needs Network Science." *Journal of Neuroscience* 43.34 (2023), pp. 5989–5995. ISSN: 15292401. DOI: [10.1523/JNEUROSCI.1014-23.2023](https://doi.org/10.1523/JNEUROSCI.1014-23.2023) (cited p. 2).
- [4] A. K. Barbey. "Network Neuroscience Theory of Human Intelligence." *Trends in Cognitive Sciences* 22.1 (2018), pp. 8–20. ISSN: 1879307X. DOI: [10.1016/j.tics.2017.10.001](https://doi.org/10.1016/j.tics.2017.10.001) (cited p. 2).
- [5] D. M. Barch, G. C. Burgess, M. P. Harms, S. E. Petersen, et al. "Function in the human connectome." *Neuroimage* 19 (2009), pp. 389–399. ISSN: 08966273. DOI: [10.1016/j.asieco.2008.09.006](https://doi.org/10.1016/j.asieco.2008.09.006). EAST (cited p. 3).
- [6] R. Botvinik-Nezer, F. Holzmeister, C. F. Camerer, A. Dreber, et al. "Variability in the analysis of a single neuroimaging dataset by many teams." *Nature* 582.7810 (2020), pp. 84–88. ISSN: 14764687. DOI: [10.1038/s41586-020-2314-9](https://doi.org/10.1038/s41586-020-2314-9) (cited pp. 1, 2, 5, 6).
- [7] O. J. Bruton. "Is there a 'g-neuron'? Establishing a systematic link between general intelligence (g) and the von Economo neuron." *Intelligence* 86.November 2020 (2021), p. 101540. ISSN: 01602896. DOI: [10.1016/j.intell.2021.101540](https://doi.org/10.1016/j.intell.2021.101540) (cited p. 2).
- [8] B. J. Casey, T. Cannonier, M. I. Conley, A. O. Cohen, et al. "The Adolescent Brain Cognitive Development (ABCD) study : Imaging acquisition across 21 sites." *Developmental Cognitive Neuroscience* 32.March (2018), pp. 43–54. ISSN: 1878-9293. DOI: [10.1016/j.dcn.2018.03.001](https://doi.org/10.1016/j.dcn.2018.03.001) (cited p. 3).
- [9] W. Chang, J. Cheng, J. J. Allaire, C. Sievert, et al. *shiny: Web Application Framework for R*. 2023. URL: <https://github.com/rstudio/shiny> (cited p. 5).
- [10] G. Chen, G. Chen, C. Xie, and S. J. Li. "Negative Functional Connectivity and Its Dependence on the Shortest Path Length of Positive Network in the Resting-State Human Brain." *Brain Connectivity* 1.3 (2011), pp. 195–206. ISSN: 21580022. DOI: [10.1089/brain.2011.0025](https://doi.org/10.1089/brain.2011.0025) (cited p. 5).
- [11] J. Chen, A. Tam, V. Kebets, C. Orban, et al. "Shared and unique brain network features predict cognitive, personality, and mental health scores in the ABCD study." *Nature Communications* 13.1 (2022). ISSN: 2041-1723. DOI: [10.1038/s41467-022-29766-8](https://doi.org/10.1038/s41467-022-29766-8) (cited p. 2).
- [12] J. Dafflon, P. F. Da Costa, F. Váša, R. P. Monti, et al. "A guided multiverse study of neuroimaging analyses." *Nature Communications* 13.1 (2022). ISSN: 20411723. DOI: [10.1038/s41467-022-31347-8](https://doi.org/10.1038/s41467-022-31347-8) (cited pp. 2, 3, 6, 7, 10).
- [13] M. Del Giudice and S. W. Gangestad. "A Traveler's Guide to the Multiverse: Promises, Pitfalls, and a Framework for the Evaluation of Analytic Decisions." *Advances in Methods and Practices in Psychological Science* 4.1 (2021). ISSN: 25152467. DOI: [10.1177/2515245920954925](https://doi.org/10.1177/2515245920954925) (cited p. 10).
- [14] B. Fischl, D. H. Salat, E. Busa, M. Albert, et al. "Whole brain segmentation: Automated labeling of neuroanatomical structures in the human brain." *Neuron* 33.3 (2002), pp. 341–355. ISSN: 08966273. DOI: [10.1016/S0896-6273\(02\)00569-X](https://doi.org/10.1016/S0896-6273(02)00569-X) (cited p. 2).
- [15] A. R. Franco. "Spatial Stability of Functional Networks : A Measure to Assess the Robustness of Graph-Theoretical Metrics to Spatial Errors Related to Brain Parcellation." 15.February (2022), pp. 1–18. DOI: [10.3389/fnins.2021.736524](https://doi.org/10.3389/fnins.2021.736524) (cited p. 5).
- [16] M. F. Glasser, T. S. Coalson, E. C. Robinson, C. D. Hacker, et al. "A multi-modal parcellation of human cerebral cortex." *Nature* (2016). ISSN: 14764687. DOI: [10.1038/nature18933](https://doi.org/10.1038/nature18933) (cited p. 2).
- [17] M. F. Glasser, S. N. Sotiropoulos, J. A. Wilson, T. S. Coalson, et al. "The minimal preprocessing pipelines for the Human Connectome Project." *NeuroImage* 80 (Oct. 2013), pp. 105–124. ISSN: 1053-8119. DOI: [10.1016/J.NEUROIMAGE.2013.04.127](https://doi.org/10.1016/J.NEUROIMAGE.2013.04.127) (cited p. 2).
- [18] K. Hilger, M. Fukushima, O. Sporns, and C. J. Fiebach. "Temporal stability of functional brain modules associated with human intelligence." *Human Brain Mapping* 41.2 (2020), pp. 362–372. DOI: [10.1002/hbm.24807](https://doi.org/10.1002/hbm.24807) (cited p. 6).
- [19] L. Hubert and P. Arabie. "Comparing partitions." *Journal of Classification* 2.1 (1985), pp. 193–218. ISSN: 1432-1343. DOI: [10.1007/BF01908075](https://doi.org/10.1007/BF01908075) (cited p. 6).

- [20] A. A. Igolkina and G. Meshcheryakov. "semopy: A Python Package for Structural Equation Modeling." *Structural Equation Modeling: A Multidisciplinary Journal* 27.6 (Nov. 2020), pp. 952–963. ISSN: 1070-5511. DOI: [10.1080/10705511.2019.1704289](https://doi.org/10.1080/10705511.2019.1704289) (cited p. 5).
- [21] R. E. Jung and R. J. Haier. "The Parieto-Frontal Integration Theory (P-FIT) of intelligence: Converging neuroimaging evidence." *Behavioral and Brain Sciences* 30.2 (2007), pp. 135–154. ISSN: 0140525X. DOI: [10.1017/S0140525X07001185](https://doi.org/10.1017/S0140525X07001185) (cited pp. 2, 6).
- [22] R. B. Kline. *Principles and Practice of Structural Equation Modeling, Fourth Edition*. Methodology in the social sciences. New York, NY, US: Guilford Publications, 2015. ISBN: 9781462523009 (cited p. 2).
- [23] K. Kovacs and A. R. Conway. "Process Overlap Theory: A Unified Account of the General Factor of Intelligence." *Psychological Inquiry* 27.3 (2016), pp. 151–177. ISSN: 1047840X. DOI: [10.1080/1047840X.2016.1153946](https://doi.org/10.1080/1047840X.2016.1153946) (cited p. 2).
- [24] K. Kriegbaum, N. Becker, and B. Spinath. "The relative importance of intelligence and motivation as predictors of school achievement: A meta-analysis." *Educational Research Review* 25. February (2018), pp. 120–148. ISSN: 1747938X. DOI: [10.1016/j.edurev.2018.10.001](https://doi.org/10.1016/j.edurev.2018.10.001) (cited p. 2).
- [25] D. Kristanto, A. Hildebrandt, W. Sommer, and C. Zhou. "Cognitive abilities are associated with specific conjunctions of structural and functional neural subnetworks." *NeuroImage* 279. November 2022 (2023), p. 120304. ISSN: 1053-8119. DOI: [10.1016/j.neuroimage.2023.120304](https://doi.org/10.1016/j.neuroimage.2023.120304) (cited p. 2).
- [26] G. Meshcheryakov, A. A. Igolkina, and M. G. Samsonova. "semopy 2: A structural equation modeling package with random effects in python." *arXiv preprint arXiv:2106.01140* (2021) (cited pp. 5, 10).
- [27] K. Paul, C. A. Short, A. Beauducel, H. P. Carsten, et al. "The methodology and dataset of the coscience eeg-personality project – a large-scale, multi-laboratory project grounded in cooperative forking paths analysis." *Personality Science* 3 (2022), pp. 1–26. DOI: [10.5964/ps.7177](https://doi.org/10.5964/ps.7177) (cited pp. 1, 10).
- [28] R Core Team. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing, 2021. URL: <https://www.r-project.org/> (cited p. 5).
- [29] M. Rubinov, R. Kötter, P. Hagmann, and O. Sporns. "Brain connectivity toolbox: a collection of complex network measurements and brain connectivity datasets." *NeuroImage* 47 (2009), S169. ISSN: 1053-8119. DOI: [https://doi.org/10.1016/S1053-8119\(09\)71822-1](https://doi.org/10.1016/S1053-8119(09)71822-1) (cited p. 5).
- [30] A. Schaefer, R. Kong, E. M. Gordon, T. O. Laumann, X.-N. Zuo, A. J. Holmes, S. B. Eickhoff, and B. T. T. Yeo. "Local-Global Parcellation of the Human Cerebral Cortex from Intrinsic Functional Connectivity MRI." *Cerebral Cortex* 28.9 (2018), pp. 3095–3114. ISSN: 1047-3211. DOI: [10.1093/cercor/bhx179](https://doi.org/10.1093/cercor/bhx179) (cited pp. 2, 6).
- [31] B. Settles. *Active Learning Literature Survey*. Tech. rep. 1648. 2009. URL: <https://research.cs.wisc.edu/techreports/2009/TR1648.pdf> (cited p. 2).
- [32] N. C. Silver and W. P. Dunlap. "Averaging correlation coefficients: Should Fisher's z transformation be used?" *Journal of Applied Psychology* 72.1 (1987), pp. 146–148. ISSN: 1939-1854(Electronic), 0021-9010(Print). DOI: [10.1037/0021-9010.72.1.146](https://doi.org/10.1037/0021-9010.72.1.146) (cited p. 5).
- [33] C. Spearman. "General Intelligence" Objectively Determined and Measured." (1904) (cited p. 2).
- [34] S. Steegen, F. Tuerlinckx, A. Gelman, and W. Vanpaemel. "Increasing Transparency Through a Multiverse Analysis." *Perspectives on Psychological Science* 11.5 (2016), pp. 702–712. ISSN: 17456924. DOI: [10.1177/1745691616658637](https://doi.org/10.1177/1745691616658637) (cited p. 2).
- [35] R. L. Thorndike. "Who belongs in the family?" *Psychometrika* 18.4 (1953), pp. 267–276. ISSN: 1860-0980. DOI: [10.1007/BF02289263](https://doi.org/10.1007/BF02289263) (cited p. 6).
- [36] J. Xiang, J. Xue, H. Guo, D. Li, et al. "Graph-based network analysis of resting-state fMRI: test-retest reliability of binarized and weighted networks." *Brain Imaging and Behavior* 14.5 (2020), pp. 1361–1372. ISSN: 19317565. DOI: [10.1007/s11682-019-00042-6](https://doi.org/10.1007/s11682-019-00042-6) (cited p. 5).
- [37] B. T. Yeo, F. M. Krienen, J. Sepulcre, M. R. Sabuncu, et al. "The organization of the human cerebral cortex estimated by intrinsic functional connectivity." *Journal of Neurophysiology* 106.3 (2011), pp. 1125–1165. ISSN: 00223077. DOI: [10.1152/jn.00338.2011](https://doi.org/10.1152/jn.00338.2011) (cited p. 6).

BabbleNewt: A Reference Citation Format for Bibliographic Metadata

S. Koby Taswell, Aniruddh Anand, Max Montes-Soza, Carl Taswell



BabbleNewt: A Simplified, Consistent, and Interoperable Reference Citation Format for Bibliographic Metadata*

S. Koby Taswell, Aniruddh Anand, Max Montes-Soza, Carl Taswell†

Abstract

Of the diverse bibliographic metadata formats, BibTeX and BibLaTeX have been dominant across mathematics, computing, and engineering due to their use with the TeX and LaTeX typesetting compilers. Despite success in these fields as well as the publishing industry, both BibTeX and BibLaTeX have some deficiencies, notably inconsistencies in the format definitions and use of macros, pseudo-records, programs and processing methods across different software implementations and installations. These inconsistencies contribute to bibliography parsing and document typesetting errors especially problematic with difficult debugging for large bibliography files. A subproject within the PORTAL-DOORS Project (PDP), the BabbleNewt Project aims to address these concerns by designing a set of formats which iterate on the original BibTeX and BibLaTeX formats while enabling easy conversion between them and a newly designed simplified, consistent, and interoperable format called BabbleNewt. The set of related formats implemented for bibliography processors by PDP BabbleNewt includes two formats PdpBibtex and PdpBiblatex corresponding to the original BibTeX and BibLaTeX, two generalized transition formats PdpBibtexgen and PdpBiblatexgen, and the novel format PdpBabblenewt.

Keywords

Bibliographic metadata, interoperability, file formats, PdpBibtex, PdpBibtexgen, PdpBiblatex, PdpBiblatexgen, PdpBabblenewt.

Contents

Introduction

Format Description

| | |
|-------------------------------|---|
| ReferenceType and CitationKey | 2 |
| No Macros or Pseudo-Records | 3 |
| Formatted Record Examples | 3 |

Format Interoperability

Format Performance

| | |
|------------|---|
| Conclusion | 4 |
|------------|---|

| | |
|----------|---|
| Citation | 4 |
|----------|---|

| | |
|------------|---|
| References | 4 |
|------------|---|

Introduction

In February 1983, Oren Patashnik began work on BibTeX, “a tool for automating your list of references”, intended to accompany LaTeX document typesetting (Patashnik 1998; Patashnik 2003; Fenn 2006). Patashnik’s original format BibTeX was accompanied by a parsing utility of the same name, often written in lower case as the command name `bibtex` to run the parser. Since the original development of `bibtex`, various other tools for the format BibTeX have also been implemented including `bibtex8`, `biber`, `BibTeXu`, `CL-dfBibtex`, `MLBibTeX`, and `Bibulous`. Whereas the original parser `bibtex` supported only 7-bit ASCII characters, `bibtex8` supports 8-bit ASCII characters and `BibTeXu` supports the UTF-8 character set. Apart from differences in processing character sets, most of the BibTeX parser alternatives have not departed from the original `bibtex` parser intended for the original BibTeX format. In contrast, the parsing tool `biber` was developed for the BibLaTeX format, designed as an extended superset of the BibTeX format (Kime and Wemheuer 2023; Mittelbach 2023). The original BibTeX format has a fixed set of entry types where an entry type declares the type of reference (eg, article, book, etc.) described within the bibliographic metadata record that includes required and optional fields for that entry type such as author, title, publisher, etc. The extended format BibLaTeX improved the usefulness of the format with the addition of many more entry types and metadata fields.

These tools are used throughout mathematics, computing, and engineering fields where LaTeX document typesetting has become the standard expected for publication of manuscripts. Despite how widely these tools are now used throughout these communities, challenges still exist compromising both formats BibTeX and BibLaTeX (Markey 2009; Rees 2017; Mittelbach 2023) in a manner that derives from the original design which lacks the simplicity and consistency of a JSON-style format. Here is a sample record of a *.bib file in the BibTeX format with double quotes for the field values:

```
@article{Patashnik1998bibtex,
  author = "Oren Patashnik",
  journal = "TUGboat",
  number = "2",
```

*Presented 2023-10-09 with slides and and video at Guardians 2023

†Authors affiliated with Brain Health Alliance Virtual Institute, Ladera Ranch, CA 92694 USA; correspondence to CTaswell at Brain Health Alliance.

Table 1: Syntax for PDP BibCitRef Formats with Placeholder Symbols
 Etyp, Ekey, Anam, Aval for Entity Type and Key, Attribute Name and Value

| Format | File Extension | Entity Opener | Attribute Name-Value Pair | Entity Closer | Attribute List |
|----------------|----------------|---------------|---------------------------|---------------|----------------|
| PdpBibtex | *.pbtx | @Etyp{Ekey, | Anam= "Aval", | } | specified |
| PdpBibtexgen | *.pbtg | @Etyp{Ekey, | Anam= {Aval}, | } | unconstrained |
| PdpBiblatex | *.pblt | @Etyp{Ekey, | Anam= {Aval}, | } | specified |
| PdpBiblatexgen | *.pblg | @Etyp{Ekey, | Anam= [Aval], | } | unconstrained |
| PdpBabblenewt | *.pbbn | @{ | Anam= [Aval], | }@ | unconstrained |

```

pages = "204-207",
title = "BIBTEX 101",
volume = "19",
year = "1998",
}

```

Here is a sample record of a *.bib file in the BibLaTeX format with curly braces for the field values:

```

@article{Patashnik1998bibtex,
  author = {Oren Patashnik},
  date = {1998-03-22},
  journaltitle = {TUGboat},
  number = {2},
  pages = {204-207},
  title = {BIBTEX 101},
  volume = {19},
}

```

The PDP BabbleNewt Project maintains a set of five different but related PDP BibCitRef formats called PdpBibtex, PdpBibtexgen, PdpBiblatex, PdpBiblatexgen, and PdpBabblenewt, intended for use with bibliography file types denoted by the file extensions *.pbtx, *.pbtg, *.pblt, *.pblg, and *.pbbn, respectively. These related formats support both backward and forward compatibility and conversion between a collection of interoperable bibliographic metadata formats. The PdpBibtex and PdpBiblatex formats correspond to the original BibTeX and BibLaTeX formats. The PdpBibtexgen and PdpBiblatexgen formats serve as generalized variant formats for didactic, development, and test purposes. The PdpBabblenewt format provides a simplified, consistent, and interoperable format with a clean separation of data from code that should maximize parsing efficiency, minimize programming errors, and simplify debugging of both parsers and data.

Format Description

The BabbleNewt Project set of PDP BibCitRef formats remain related to each other in a progressive transition to facilitate migration and conversion of bibliography files from one format to another. In a bibliography file for any of these formats, a bibliographic citation record for a bibliographic reference entity consists of an entity opener, a list of attribute name-value pairs, and an entity closer. The related set of formats differ with respect to the syntax required for the entity opener/closer pair and the list of attribute name-value pairs, also importantly, whether the format specifies the list of attribute pairs or allows the list of attribute pairs to be unconstrained (see Table 1). The formats PdpBibtex and PdpBiblatex specify the lists of entity types and lists of attribute name-value pairs for each entity type, whereas the generalized formats PdpBibtexgen, PdpBiblatexgen, and PdpBabblenewt allow these lists to be unconstrained. Lists of attribute name-value pairs are separated by commas with each name and value in a pair separated by an equal sign.

For the progressive sequence of formats PdpBibtex, PdpBibtexgen, PdpBiblatex, and PdpBiblatexgen, the entity opener/closer pairs remain the same as in the PdpBibtex format for the entity type with entity key and curly braces. In contrast, the entity opener/closer has been simplified and made symmetric for the format PdpBabblenewt with the opener "@{" and closer "}" to create a more consistent JSON-style scheme for the bibliography file, while also allowing for unconstrained attribute name-value pairs in any order. The entity type and key from the PdpBibtex format have been mapped, respectively, to the attribute name-value pairs with names "referencetype" and "citationkey" in the PdpBabblenewt format.

Moreover for attribute values, the curly brace delimiters in the PdpBibtexgen and PdpBiblatex formats have been changed to square bracket delimiters in the PdpBiblatexgen and PdpBabblenewt formats. Switching from curly braces to square brackets as delimiters for attribute values improves human readability and also improves computer parsing by avoiding the nesting of curly braces, thus simplifying parsing with regular expressions, requirements for escape sequences, and reducing programming errors.

Inspired by a simplified JSON-style approach, the scheme found in the PdpBabblenewt format should permit development of faster parsers with fewer errors. While perhaps not important for small bibliography files with only a few dozen records, error-free efficiency becomes much more important for millions of records in large-scale databases. Indeed, Nurseitov et al. (2009) found that the data processing rates for JSON were much faster and less resource intensive than for XML with parsing of JSON up to 100 times faster than XML.

ReferenceType and CitationKey

The ReferenceType for PDP BibCitRef formats is defined as the type of cited reference such as article, book, report, etc. To describe the reference entity, the ReferenceType determines the list of allowed attribute name-value pairs for the reference entity in those bibliography formats (PdpBibtex and PdpBiblatex) that require them, and corresponds to what has been called the "entry type" in the past. ReferenceTypes, when used with bibliography styles that permissively allow both required and optional attribute name-value pairs for each ReferenceType can be better supported with the generalized and unconstrained bibliography formats (PdpBibtexgen, PdpBiblatexgen, and PdpBabblenewt). For more robust parsing, ReferenceTypes should be considered case-insensitive when processed in algorithms.

Each reference entity record in a bibliography file should always have both a ReferenceType and a CitationKey as a unique identifier to assure disambiguation of references. All PDP BibCitRef formats require, and generate if necessary, a unique CitationKey for each reference entity with a ReferenceType. In general, the CitationKey may be any arbitrary unique character string of arbitrary length. Long identifiers

for references quickly become inconvenient when typing the source for manuscripts, whereas use of a max-length patterned generator related to the bibliographic metadata for the reference entity provides a consistent mechanism that facilitates easier recognition of CitationKeys. PDP BibCitRef formats generate CitationKeys with a pattern comprised of 3 components with a 16-char-max identifier for provenance (from LastNameFirstAuthor, LastNameEditor, or OrganizationName), an 8-char-max identifier for date (from Date or Year), and an 8-char-max identifier for title (from AcronymFromTitle or WordFromTitle), yielding a CitationKey with a maximum length of 32 characters.

No Macros or Pseudo-Records

A pseudo-record in a bibliography file does not describe a bibliographic reference, but instead provides some other functionality. Pseudo-records can be macros to perform actions such as basic substitution or commands to trigger more complex actions, which interact with other inputs, outputs, or data files. Neither macros nor other kinds of pseudo-records are allowed in PDP BibCitRef bibliography files because the BabbleNewt Project maintains a guiding principle of imposing consistency on the set of related PDP BibCitRef formats in a simplified JSON-like style such that data and code do not mix. This guiding principle implies maintaining a clear boundary between data and code, ie, between the formatted and structured data in data files and the processing algorithms implemented in lexing, parsing, and other utilities of software libraries. Therefore, the PdpBabblenewt format will remain clean with only data and without any code, macros, or pseudo-records. Conversions to the PdpBabblenewt format from other formats requiring expansions of incomplete and/or abbreviated data should be pre-processed with the necessary macro substitutions.

Formatted Record Examples

PdpBibtex (*.pbtX)

```
@article{Patashnik2003BYTT,
  author = "Oren Patashnik",
  journal = "TUGboat",
  title = {BibTeX yesterday, today, and tomorrow},
  volume = "24",
  year = "2003",
}
```

PdpBibtexgen (*.pbtg)

```
@article{Patashnik2003BYTT,
  author = {Oren Patashnik},
  journal = {TUGboat},
  title = {BibTeX yesterday, today, and tomorrow},
  volume = {24},
  year = {2003},
}
```

PdpBiblatex (*.pbL)

```
@article{Patashnik2003BYTT,
  author = {Oren Patashnik},
  date = {2003},
  journaltitle = {TUGboat},
  title = {BibTeX yesterday, today, and tomorrow},
  volume = {24},
}
```

PdpBiblatexgen (*.pblg)

```
@article{Patashnik2003BYTT,
  author = [Oren Patashnik],
  date = [2003],
```

```
journaltitle = [TUGboat],
title = [BibTeX yesterday, today, and tomorrow],
volume = [24],
}
```

PdpBabblenewt (*.pbbn)

```
@{
  referencetype = [article],
  citationkey = [Patashnik2003BYTT],
  author = [Oren Patashnik],
  date = [2003],
  journaltitle = [TUGboat],
  title = [BibTeX yesterday, today, and tomorrow],
  volume = [24],
}@
```

Format Interoperability

Citation Style Language (CSL), developed by Zelle (2015), is an XML-based language for use with citations of references in bibliographies. Similar to BibTeX and BibLaTeX, CSL allows mixing of both code and data in the same file. The BabbleNewt format differs from CSL, BibTeX and BibLaTeX by requiring strict adherence to a data-only principle for the bibliography, thus disallowing macros, commands, styles, pseudo-records or other kinds of code mixed into the data file. As a JSON-like data format, BabbleNewt also differs from CSL implemented as an XML-based language. CSL uses a "CitationKey" but not a "ReferenceType". Differences between formats for entity-attribute names (aka record field names), such as "ReferenceType" and "CitationKey" regardless of punctuation use and letter casing in the names, can be accommodated by mappings for the related entity attributes when processing transforms from one format to another with import, export, and convert utilities. Thus, the BabbleNewt format is interoperable with CSL and any other bibliographic metadata format including both backward and forward compatibility with the BibTeX and BibLaTeX formats.

The BabbleNewt format maintains adherence to principles for simplifying the format design in order to reduce errors in both data and code, thereby improving reliability and efficiency of processing utilities. To be compatible with requirements for the PrincipalTags of NPDS resource entities (C. Taswell 2007; C. Taswell 2010), and to map a CitationKey for a BibCitRef record to the corresponding PrincipalTag for an NPDS record, use of punctuation symbols such as the hyphen must be avoided in both the attribute name and attribute value. The BabbleNewt format imposes this same requirement on all other attribute names (eg, "ReferenceType" and not "reference-type") but not on other attribute values for which it would be impractical. This no-punctuation rule for both value and name of an attribute only applies to the CitationKey.

This simplifying rule imposed on the CitationKey implements an important design principle: Avoid use of unnecessary escape symbols, punctuation, and characters that may complicate processing and contribute to additional requirements for more complexity in lexers and parsers. Unnecessary complexity only worsens the probability of coding errors in the software and faulty processing of the data. This simplified design of the BabbleNewt format with a consistent JSON-like style will support more robust lexing and parsing algorithms with greater portability across different programming languages.

Format Performance

Read-write accuracy and efficiency tests were performed on bibliography files in each of the 5 related formats BibTeX, PdpBibtexgen,

Table 2: Format Median Round Trip Timing Experiments per Number Records with Lossless Transfer in Seconds

| Timing Tests | PdpBibtex | PdpBibtexgen | PdpBiblatex | PdpBiblatexgen | PdpBabblenewt |
|----------------|-----------|--------------|-------------|----------------|---------------|
| Initialization | 0.37 | 0.36 | 0.33 | 0.40 | 0.32 |
| 8 records | 0.38 | 0.38 | 0.35 | 0.42 | 0.33 |
| 80 records | 0.61 | 0.55 | 0.57 | 0.61 | 0.52 |
| 800 records | 2.62 | 2.07 | 2.35 | 2.36 | 1.95 |
| 8000 records | 20.99 | 16.52 | 18.90 | 16.95 | 14.75 |

PdpBiblatex, PdpBiblatexgen and PdpBabblenewt corresponding to the same bibliography database of more than 8000 records. The experimental protocol involved several steps: 1) Initialize the BabbleNewt lexer for each format, 2) Measure time for a round-trip cycle of read from import file on disk to record list in memory then write back the records to export file on disk, and 3) Examine and compare export file to import file for any differences in lines or characters. These tests were repeated for varying counts of 8, 80, 800, and 8000 records for each of the 5 bibliography database formats. At all size counts from 8 to 8000, and for all 5 formats, the export files were observed to match the import files exactly with perfect reproducibility. Table 2 summarizes the processing times for these efficiency tests which shows that the PdpBabblenewt format was the most efficient.

Conclusion

The PDP BabbleNewt Project has developed a set of 5 related bibliography database formats called PdpBibtex, PdpBibtexgen, PdpBiblatex, PdpBiblatexgen, and PdpBabblenewt which iterate, extend, and generalize the original BibTeX and BibLaTeX formats while maintaining both backward and forward compatibility as well as supporting progressive transitional migrations between the formats. This set of related formats and the accompanying BabbleNewt lexer have been designed with adherence to the software engineering principle of separating the data files for the bibliographic data from the code files for algorithms implemented in utilities and programs that process the data. Guiding principles for design and implementation for both data formats and processing utilities in the BabbleNewt Project emphasize the concepts of simplicity, consistency, reproducibility, and interoperability with a JSON-like style. Whereas the BabbleNewt Project with its BabbleNewt lexer focuses on processing for interoperability between the set of 5 related formats presented herein, the BabbleBird Project with its BabbleBird parser (S. K. Taswell and C. Taswell 2024) focuses on processing for interoperability between other bibliography database repositories such as IEEE Xplore, NLM PubMed, and Unpaywall, as well as other bibliographic metadata formats such as BIBFRAME, MARC, and RIS (S. K. Taswell, Uhegbu, et al. 2020).

Citation

Brainiacs 2023 Volume 4 Issue 2 Edoc K562CB81C

Title: "BabbleNewt: A Simplified, Consistent, and Interoperable Reference Citation Format for Bibliographic Metadata"

Authors: S. Koby Taswell, Aniruddh Anand, Max Montes-Soza, Carl Taswell

Dates: created 2021-08-19, received 2023-06-28, presented 2023-10-09, updated 2023-12-18, published 2023-12-18, endorsed 2023-12-30

Copyright: © 2023 Brain Health Alliance

Contact: [CTaswell at Brain Health Alliance](mailto:CTaswell@BrainHealthAlliance.org)

URL: BrainiacsJournal.org/arc/pub/Taswell2023BBNewt

PDP: [/Nexus/Brainiacs/Taswell2023BBNewt](https://Nexus/Brainiacs/Taswell2023BBNewt)

DOI: [/10.48085/K562CB81C](https://doi.org/10.48085/K562CB81C)

References

- [1] J. Fenn. "Managing citations and your bibliography with BibTeX." *The PracTEX Journal* 4 (2006) (cited p. 1).
- [2] P. Kime and M. Wemheuer. *BibLaTeX – Sophisticated Bibliographies in LaTeX*. Developed and maintained 2006–2012 by Philipp Lehman; 2012–2017 by Philip Kime, Audrey Boruvka, Joseph Wright; 2018–2023 by Philip Kime, Moritz Wemheuer. 2023. URL: <https://ctan.org/pkg/biblatex> (cited p. 1).
- [3] N. Markey. *TameTheBeaST – A manual about bibliographies and especially BibTeX*. Oct. 11, 2009. URL: <https://ctan.org/pkg/tamethebeast> (visited on 02/27/2022) (cited p. 1).
- [4] F. Mittelbach. *The LaTeX companion*. Ed. by U. Fischer. Third edition. Tools and techniques for computer typesetting. Parts I & II. Boston: Addison-Wesley, 2023. ISBN: 013816648X (cited p. 1).
- [5] N. Nurseitov, M. Paulson, R. Reynolds, and C. Izurieta. "Comparison of JSON and XML data interchange formats: a case study." *Caine* 9 (2009), pp. 157–162 (cited p. 2).
- [6] O. Patashnik. "BIBTEX 101." *TUGboat* 19.2 (Mar. 22, 1998), pp. 204–207 (cited p. 1).
- [7] O. Patashnik. "BibTeX yesterday, today, and tomorrow." *TUGboat* 24.1 (2003), pp. 25–30 (cited p. 1).
- [8] C. F. Rees. *BibLaTeX/Biber Cheat Sheet*. CTAN, June 24, 2017. URL: <https://ctan.org/pkg/biblatex-cheatsheet> (visited on 02/27/2022) (cited p. 1).
- [9] C. Taswell. "DOORS to the Semantic Web and Grid with a PORTAL for Biomedical Computing." *IEEE Transactions on Information Technology in Biomedicine* 12.2 (2 Mar. 2007). In the Special Section on Bio-Grid published online 3 Aug. 2007, pp. 191–204. ISSN: 1089-7771. DOI: [10.1109/TITB.2007.905861](https://doi.org/10.1109/TITB.2007.905861) (cited p. 3).
- [10] C. Taswell. "A Distributed Infrastructure for Metadata about Metadata: The HDMM Architectural Style and PORTAL-DOORS System." *Future Internet* 2.2 (2010), pp. 156–189. ISSN: 1999-5903. DOI: [10.3390/FI2020156](https://doi.org/10.3390/FI2020156). URL: <https://www.mdpi.com/1999-5903/2/2/156> / (cited p. 3).
- [11] S. K. Taswell and C. Taswell. "BabbleBird: A Flexible Software Library for Converting Diverse Bibliographic Formats" (2024). Manuscript in preparation. (cited p. 4).
- [12] S. K. Taswell, K. Uhegbu, S. Mashkoor, S. Dutta, and C. Taswell. "Storing bibliographic data in multiple formats with the NPDS cyberinfrastructure." *Proceedings of the Association for Information Science and Technology* 57.1 (Oct. 2020). DOI: [10.1002/pra2.428](https://doi.org/10.1002/pra2.428) (cited p. 4).
- [13] R. M. Zelle. *Citation Style Language Primer – An Introduction to CSL*. 2015. URL: <https://docs.citationstyles.org/en/stable/primer.html> (cited p. 3).

Managing Lexical-Semantic Hybrid Records of FAIR Metrics Analyses

Adam Craig, Anousha Athreya, Carl Taswell



Managing Lexical-Semantic Hybrid Records of FAIR Metrics Analyses with the NPDS Cyberinfrastructure*

Adam Craig, Anousha Athreya, Carl Taswell†

Abstract

Current approaches to plagiarism detection often focus on finding lexical matches rather than semantic similarities in the text content that is compared. But the more important unanswered questions remain whether similar concepts expressed in related topical contexts are semantically equivalent as idea-laundering plagiarism by humans or algorithm-generated plagiarism by machines. Now publicly available and easily accessible, text-generating algorithms have automated the process of assembling a text derived from but not attributed to published content scraped from the web. The FAIR Metrics, with FAIR an acronym for Fair Attribution to Indexed Reports and Fair Acknowledgment of Information Records, measure how appropriately a document cites prior records based on whether they contain similar claims that are equivalent in meaning. We demonstrate herein a workflow with results for manual evaluation of the FAIR Metrics to quantify the extent of plagiarism in 8 articles retracted or reported for plagiarism. We also demonstrate use of the Nexus-PORTAL-DOORS-Scribe (NPDS) Cyberinfrastructure to manage semantic descriptions of the concept mappings and entity equivalence evaluations made using concepts and relationships from the PDP-DREAM Ontology.

Keywords

Plagiarism, bibliometrics, citation analysis, knowledge engineering, semantic web, equivalent entities, concept mapping, ontology.

Contents

[Introduction](#)

[Methods](#)

[Results](#)

[Discussion](#)

[Conclusion](#)

[Citation](#)

References

6

Introduction

With the rise of generative artificial intelligence (AI), scholastic institutions and scholarly publishers have recognized the need for tools to detect AI-generated documents, initiating an arms race with AI-assisted plagiarists. Earlier this year, the journal *Science* recently updated its editorial policies to clarify that use of artificial intelligence to produce papers is plagiarism (Thorp 2023). Tools such as [Copleaks 2023](#), [GPTZero 2023](#), and the OpenAI text classifier (Kirchner et al. 2023), attempt to detect the probability that a text document was produced by an AI algorithm instead of a living person (Orenstrakh et al. 2023). Manuscripts in Springer, Elsevier, IEEE, Wiley, and ProQuest utilize CrossCheck, a plagiarism detection tool by iThenticate that is available to publication editors within the journals (IEEE 2023). Per the IEEE webpage, CrossCheck compares manuscripts to a database of over 6 billion web pages of published technical papers and provides a report of the similarity to previously published work. [Copleaks 2023](#), which compares submitted documents against large datasets, also includes cross-language detection capabilities and may also detect image-based text plagiarism using optical character recognition technology. Scholarly publishers use Turnitin iThenticate (Young 2023) to detect plagiarism in publishing, while universities use Turnitin Similarity, another product from the same company, to check manuscripts in education. [Khalil and Er 2023](#) tested the ability of Turnitin iThenticate and Similarity to identify plagiarism in essays written by ChatGPT and found that similarity scores ranged from 0% to 68%, indicating the need for new approaches.

1 Several recent surveys have documented the search for new analytic
algorithms, especially for methods that look beyond superficial differ-
2 ences in wording to the meaning and structure of a work. [Vrbanc and
Meštrović 2017](#) evaluated plagiarism detection methods currently used
3 by Croatian higher education institutions for measuring the quality of
academic and scientific work. In a preliminary review, they discussed
4 the use of semantic similarity techniques as an alternative for plagiarism
5 detection by quantifying the similarity of meaning in texts. [Altheneyan
and Menai 2020](#) discussed use of paraphrase identification through
6 word overlap and structural representations for application to automatic
6 plagiarism detection. They compared existing methods, measuring Pre-
cision, Recall, and F-measure values. They found that the most optimal
results were obtained with SVM and deep learning classifiers while the
worst resulted from naive similarity-based methods. They found that
all methods have worse precision than recall due to the high overlap in

*Presented 2023-10-09 with [slides](#) and [video](#) at [Guardians 2023](#); preliminary version [Craig, Athreya, et al. 2023](#) presented as [poster](#) at [eScience 2023](#).

†Authors affiliated with Brain Health Alliance Virtual Institute, Ladera Ranch, CA 92694 USA; correspondence to [A Craig](#) at [BHAVI](#).

distributions of lexical similarity measures between false paraphrase pairs and true paraphrase pairs.

A recent survey of plagiarism detection tools by [Jiffriya et al. 2021](#) classified plagiarism detection methods as lexical, structural, semantic, stylometric, syntactic, citation, or cross-language based. For natural language plagiarism detection, style-based identification remains difficult because web-based software typically only analyses the authors' first submissions of manuscripts. Detection tools were found to have false-positive results and inability to detect copied content due to scope of detection, paraphrasing, and cross-language plagiarism. Some promising new methods of semantic plagiarism detection include those from [Javadi-Moghaddam et al. 2022](#) and [Eisa et al. 2020](#). [Javadi-Moghaddam et al. 2022](#) investigated semantic plagiarism detection methods using weighted values for matched instances within manuscript sections. The method utilizes the most frequent terms of the manuscript. They found that the model is more accurate depending on the number of surrounding terms, tested with 1-, 2-, and 3-term examples, with a larger window allowing the model to check for adjacent plagiarism. [Eisa et al. 2020](#) proposed a method for detection of image and figure plagiarism in scientific publications. Because image-based plagiarism detection is rooted in determining the meaning of the figure, the method obtains structural and textual features to check for a similarity score between the elements. It then uses semantic mapping to relate the associated concepts between figures.

Others began their fight against AI-assisted plagiarism before the present generative AI boom. In 2013, [C. Labbé and D. Labbé 2013](#) reported that they had identified 85 purportedly peer reviewed papers in 24 conference proceedings that were products of the SciGen text generation algorithm. As noted by [C. Labbé and D. Labbé 2013](#), even though SciGen produces grammatically correct, properly formatted documents, a human reader can easily discern them from actual reports of scientific research due to the lack of any coherent meaning behind the concatenations of technological buzzwords. Even though [C. Labbé and D. Labbé 2013](#) and [Xiong and Huang 2009](#) both provided effective methods for automatically detecting SciGen-derived text, as late as 2021, [Cabanac and C. Labbé 2021](#) identified 243 SciGen pseudo-articles, 192 of which remained in publication, neither retracted nor withdrawn.

Furthermore, the new wave of AI-assisted text generators represent a greater challenge. [Gao et al. 2022](#) found that even human reviewers could only identify ChatGPT-generated abstracts 68% of the time and that plagiarism detection software did not flag any of them as taken from other indexed online content. This automated remixing of content in which the plagiarizing author may be completely unaware of the existence of the original work (when the black-box intermediary of the AI generator hides the sources) represents a new level of social disconnection between plagiarist and victim that was not possible when taking words or ideas from a work required that one read it and manually copy or paraphrase its content. Bibliometric analysis from [Santos-d'Amorim et al. 2022](#) suggested a possible starting point for this trend with evidence of a rise in plagiarized work from paper mills beginning in 2015. But [Gaudino et al. 2021](#) showed the start of a meteoric rise in retractions for research misconduct beginning as early as the late 1990s. Although plagiarism certainly did not begin with the development of the internet and web, modern information technology has made it easier to discover literature for both proper citation and referencing of sources and for the illegitimate plagiarism of those sources.

The inability of both algorithms and human reviewers to reliably detect plagiarism and the slowness, dismissiveness and/or non-response

by some publishers to address reports of plagiarism shows that the scholarly publishing community needs a new approach. One such strategy proposed by [Craig, Lee, et al. 2022](#): Publishers should improve the quality and integrity of the peer review process to provide publicly accessible living documents which track, monitor, and record continued checking of the claims made, and sources cited, by a published document. As part of this more rigorous approach, the FAIR Metrics provide a framework for appraisals of how well a scholarly work adheres to community standards by accurately attributing ideas to their sources [Craig, Ambati, Dutta, Kowshik, et al. 2019](#). Different from approaches based solely on lexical similarity of texts, evaluation of FAIR Metrics depends on search of previously published literature for claims with equivalent meaning ([Athreya et al. 2020b](#)). Because this semantic analysis is more difficult to automate for machine algorithms than lexical analysis, and more labor intensive to perform by human persons, prior work has only demonstrated the properties of the FAIR metrics using hypothetical test cases ([Craig, Ambati, Dutta, Mehrotra, et al. 2019](#)).

However, in a recent report at [eScience 2019](#), we introduced a practical approach to evaluating FAIR Metrics by human analysts of semantic concepts for each test document with respect to similarities found in a limited pool of comparison texts, and summarized the results of this evaluation on a set of 5 different test examples ([Craig, Athreya, et al. 2023](#)). In the present report, we provide a more thorough account of the evaluation process and discuss how the FAIR Metrics scores relate to the shared social context of the evaluated test and comparison texts. Additionally, the present report provides more detail regarding use of the PDP-DREAM Ontology to represent the results of human-analyst FAIR Metric evaluations in machine-readable resource description format (RDF) knowledge graphs. These linked graphs can then serve as openly accessible and searchable records of the assessments with the FAIR Metrics, enabling transparency and discussion of both subjective and objective evaluations of the scientific claims contributed to the historical record of published literature ([S. K. Taswell et al. 2020](#); [Craig, Lee, et al. 2022](#)). For more about the FAIR Metrics and PDP-DREAM Ontology, see [Craig, Ambati, Dutta, Kowshik, et al. 2019](#), [Dutta, Uhegbu, et al. 2020](#), and [Craig and C. Taswell 2021](#).

Methods

[Craig, Ambati, Dutta, Mehrotra, et al. 2019](#) described 4 ratio metrics calculated from counts of 4 categories of claims: Quoted (Q) claims correctly attributed to prior work, Misquoted (M) claims misrepresenting prior work, Plagiarized (P) claims matching but not attributed to prior work, and Novel (N) claims not found in or reported as sourced from prior work. We now use subscripts with letters instead of numbers to clarify which ratio metric emphasizes which count with F_Q, F_M, F_P, F_N here corresponding respectively to F_1, F_2, F_3, F_4 in [Craig, Ambati, Dutta, Kowshik, et al. 2019](#). In the ideal automated use case described in [Craig, Ambati, Dutta, Mehrotra, et al. 2019](#), a semantic inference engine checks for equivalence relationships between the subject, verb, and object URIs of 2 RDF triples that reference appropriate formal ontologies. At present, creating sufficiently semantically rich descriptions of the scientific claims of a report to allow such automated comparison is a complex and labor-intensive task. We are not aware of an existing library of such descriptions extensive enough to permit a comprehensive search for equivalent statements.

As a practical interim approach to applying and using the FAIR Metrics that we can demonstrate now, we introduce limited-scope human-analyst evaluation of scientific claims for the FAIR Metric calculations.

Craig, Ambati, Dutta, Mehrotra, et al. 2019 described an earlier attempt at a pairwise comparison of scholarly articles, but the approach described there failed to produce usable results. Our new procedure differs in that we evaluate all claims with cited sources instead of discarding those that cite a source other than the comparison document, providing a more reliable and valid set of counts. In this current approach, a human evaluator compares the test document to any resources it cites and 1 or more specific references from which the authors have been proven or reported to have plagiarized previously published material. As a summary of the approach, we used the following procedure: 1) Access test T and comparison C documents and the set of references $\{R_j \mid j = 1, 2, \dots, J\}$ cited by T and/or C . 2) Relabel C as R_0 so that it can be analysed in the set of references $\{R_j \mid j = 0, 1, 2, \dots, J\}$. 3) List statements and select claims, ie, statements highlighted as novel or cited with a reference. 4) Initialise counts M, N, P, Q to 0 and iterate the comparison analysis over the claims. 5) If claim in T cites R_j , search R_j for equivalent claim. 6) If found, increment Q else increment M . 7) If claim in T does not cite a source, search R_j for equivalent claim. 8) If found, increment P else increment N . While this method (limited in scope to analysis of $J + 2$ documents) does not suffice to detect all cases of plagiarism, it can serve as a more objective method of assessing allegations of suspected plagiarism and/or of misrepresentations of previously published references and records in the literature when there exist known test T and comparison C documents.

The distinction between statements and claims reflects the practicality that not every phrase or sentence in a document represents a substantive and meaningful contribution to scholarly knowledge. The reports we examined as test cases contain general sentiments, reiterations of common knowledge, and technical details often found in what are considered materials and methods rather than scientific claims found in results, discussions, or conclusions. The selection of key claims found in a scientific, engineering, or medical report should also reflect the current state of knowledge regarding what community standards exist for the relevant field of scholarly inquiry and research. For example, in a genome-wide association study (GWAS) producing p-values for differential expression of many genes in the human genome, we would not consider the result of each statistical p-value test a meaningful claim in isolation. Instead, in this context, the final results, inferences, and conclusions drawn by the GWAS based on the lower-level intermediate results would be considered the key claims. Craig, Ambati, Dutta, Mehrotra, et al. 2019 did not clarify such a convention for this distinction between statements and claims for evaluation purposes when calculating FAIR metrics.

For the present analysis, we consider a claim to be any statement highlighted as an important concept in the abstract, statements implicitly or explicitly declared to be novel concepts, and/or statements corresponding to concepts otherwise attributed to a cited source. It is common practice for papers to reiterate their key claims in multiple sections, so it is important to take care to avoid double-counting claims. When evaluating texts organized into the standard set of 6 sections (Abstract, Introduction, Methods, Results, Discussion, and Conclusion), we found that counting claims from the Introduction and Discussion sections was most expedient. Claims in the Abstract and Conclusion typically lack citations, whereas the purpose of the Introduction and Discussion often strives to place the goals and results of the research in the proper context of the published literature and to cite relevant sources. Claims from the Methods and Results sections may also be

appropriate for consideration as Methods claims and Results claims. But most reports we have evaluated here restate the results that rise to the level of substantively meaningful claims in the Discussion. If a test T document cited multiple sources for the same claim, we considered it a quoted claim for incrementing the quoted Q count if at least 1 of the sources had an equivalent claim. Fundamentally, this comparison evaluation method depends on the ability to recognize equivalence between 2 claims, ie, when 2 claims are equivalent substantively in semantic meaning in context. For a detailed discussion of different interpretations of lexical and semantic equivalence, see Athreya et al. 2020a; Athreya et al. 2020b. For purposes of the current demonstration with analysis of 9 test cases (see Table 1), we used the following method to identify equivalent statements. For each claim indexed by i as $T_i \in T$, the human reader finds the corresponding claim indexed by k as $R_{j,k} \in R_j$ closest in meaning. The analyst then evaluates the pair of claims as either equivalent or not equivalent.

For this initial sample of 9 test cases in Table 1, representative examples of different real-world scenarios were chosen. As a negative control, we selected C. Taswell 2007 compared to Mons 2005 as an example pair on a related topic but with little overlap between T and C in terms of the concepts and ideas presented and discussed. As positive controls, 7 examples of journal articles known to have been retracted for different levels of plagiarism were selected. These reports with known plagiarism were found with a search of the Retraction Watch website and database *Retraction Watch Database User Guide 2023*. The example Uddin et al. 2022 plagiarized heavily from Foster et al. 2019, but also properly cited numerous sources. The example Gnat et al. 2022 cited Hoog et al. 2016, but also used content without citation.

Three of the examples illustrate a shared tactic for obfuscating plagiarism. Ullah et al. 2018 is a case of whole-text plagiarism with only cursory paraphrasing from Sansaniwal and Kumar 2015, a work describing a test of a solar-powered produce dryer, except that the plagiarists substituted their own home institution for the original authors' as the testing site and replaced ginger with asparagus as the vegetable being dried. The other 2, Yao et al. 2016 and Dai et al. 2015, applied the same tactics albeit with greater sophistication, replacing multiple content words from the original articles they plagiarized, G. Li et al. 2015 and Lv et al. 2015 respectively, and changed some background statements and references where simple substitution would have lead to factually incorrect statements or where the cited source referenced 1 of the replaced terms. These cases differ in that Dai et al. 2015 applied the latter tactic more thoroughly. By contrast, Guo et al. 2013 plagiarized almost all of Fischbach et al. 2009 without using this tactic of systematically swapping in meaningfully different content words. Instead, they paraphrased extensively, sometimes changing the meaning of a claim seemingly by accident.

As a positive control, we considered Su et al. 2005, an example of near-verbatim whole-text plagiarism with only minor edits of the content of Schwab et al. 2001. Finally, we examined Wilkinson et al. 2016, which has not yet been retracted for plagiarism. Previously reported in Craig, Ambati, Dutta, Kowshik, et al. 2019, all of the the Findable, Accessible, Interoperable, and Reusable (FAIR) Principles described in Wilkinson et al. 2016 plagiarized (as idea-laundering plagiarized versions) of some, but not all, of the design and practice principles described in C. Taswell 2007. The original PORTAL-DOORS Project Principles (C. Taswell 2007; C. Taswell 2010) have been renamed the PDP-DREAM Principles (Craig, Ambati, Dutta, Kowshik, et al. 2019).

For a demonstration of publishing these FAIR Metrics analyses in a

machine-readable manner, we have published records of the comparison documents at [PORTALDOORS.net](https://portaldoors.net) using the reference implementation of the Nexus-PORTAL-DOORS-Scribe (NPDS) Cyberinfrastructure. NPDS provides an online information management system for sharing and distributing data records about different kinds of online and offline resources grouped by problem domain (C. Taswell 2007; Dutta, Kowshik, et al. 2019). We have scoped the Fidentinus diristry for NPDS records with descriptions of known plagiarism cases, while other documents not suspected of plagiarism, such as C. Taswell 2007, have been described in NPDS records in other diristries appropriate to their problem domains, which for C. Taswell 2007 can be found in the DaVinci diristry for semantic web technologies. In addition to including the FAIR Metric values as metadata items in the NPDS records, we developed a FAIR Metrics sub-module of the PDP-DREAM Ontology, a formal OWL ontology for codifying the relationships among concepts relevant to the PORTAL-DOORS Project (Craig, Ambati, Dutta, Kowshik, et al. 2019). This sub-module features the classes and properties needed to record the key assertions an analyst makes when evaluating the FAIR Metrics: the identification of key claims in a document, attributions of claims to other documents, scoring of equivalence matches between claims across documents, classification of each claim in the test document, total counts for the 4 categories, and the FAIR Metrics ratios calculated from those counts.

Results

While we developed the most recent version of the PDP-DREAM Ontology formatted as N-Quads (Craig and C. Taswell 2021), we have also created a version of it formatted in standard Web Ontology Language (OWL) 2.0 XML using Stanford Protégé in order to support compatibility with a broad variety of consumer applications (Drummond et al. 2005). We found that Protégé dropped the fourth element of each quad, the graph label, and treated each quad as a triple. Subsequently, when developing a new FAIR Metrics sub-module, we organized all classes under the *FairMetricsRelatedEntity* class, all object properties under *hasFairMetricObjectProperty*, and all data properties under *hasFairMetricDataProperty*. We established 2 major classes: *Document* and *Statement*. We then assigned *Statements* to the subclasses *NonClaim*, *Claim*, or *FairMetricCategorizedClaim*, which in turn has 4 subclasses: *MisquotedClaim*, *NovelClaim*, *PlagiarizedClaim*, and *QuotedClaim*. We designated 2 object properties: *hasAttribution*, to indicate a reference from a *Claim* in 1 *Document* to another, and *hasFairMetricClaimCategory*, to indicate that a *Claim* belongs to 1 of the subclasses of *FairMetricCategorizedClaim*. We used 12 distinct data properties in each RDF record of a test text. For human readability, we embedded the title of a *Document* in its RDF description using *hasName*. Similarly, we used *hasText* to embed the original natural language representation of a *Claim* in its RDF description. We also used *hasEquivalenceScore* for the equivalence score and 4 data properties to represent the 4 FAIR Metric counts and another 4 to represent the 4 FAIR Metric ratios. While reviewing the results, we added the data property *hasEquivalentClaimText* with which to directly embed the text of a matching claim in the description of a claim being tested. We found that this procedure makes it easier for the reader to check equivalence of the claims.

We report the results of the FAIR Metrics analyses on the 9 example pairs in Table 1. The negative control, C. Taswell 2007 written as a literature review with integrated synthesis of a collection of design and practice principles, had no substantive overlap with Mons 2005 and cited all its sources adequately, resulting in F_M , F_P , and F_Q scores

of 1. The ratio of novel claims to cited claims was nearly even, leading to an F_N score close to 0. Different from the other FAIR Metric ratios that have increasing values of fairness and ideal values of 1, this novelty measure F_N does not necessarily have an ideal value which can vary according to the type of manuscript, eg, primary research report versus secondary literature review. Future work will establish what values of each of the FAIR metrics should be considered acceptable for that measure and what values should meet the standards of the scholarly community in a given research field.

Both Uddin et al. 2022 and Gnat et al. 2022 attained positive FAIR Metric scores, as both appropriately cite the sources of most concepts they presented. Although F_M , F_P , and F_Q are greater than 0, they are still well below 1, which would be sufficient to alert an editor to issues requiring further scrutiny. The negative F_P score of Ullah et al. 2018 demonstrates that the FAIR Metrics are immune to conventional paraphrasing. However, the non-zero N count shows that changing actual content words to those with different meanings can decrease the apparent extent of plagiarism. Nevertheless, this tactic of random word replacement did result in misrepresentations of the content of the cited sources, including such clearly erroneous statements as “About half of the total production of Asparagus is being consumed as white and red Asparagus, whereas the remaining 30% is converted into dry Asparagus for medicinal purposes, and 20% is used as seed material” (Deshmukh, Varma, et al. 2014). We excluded 5 claims from the analysis of Ullah et al. 2018 due to inability to locate any of the texts cited as their sources. Yao et al. 2016 copied most of the content of G. Li et al. 2015, but replaced several keywords with some additional paraphrasing. They replaced “chondrosarcoma” with “glioblastoma”, “Slug” with “Twist”, “CXCR7” with “CXCR4”, “CCL21” with “CXCL12”, “SW1353” with “U87”, and “transwell” with “wound healing”.

However, they completed a more deliberate substitution and paraphrasing than did the authors of Ullah et al. 2018. In particular, they rewrote the first few sentences of the introduction, because replacing “chondrosarcoma” with “glioblastoma” would have resulted in clearly false statements. Where the sources that G. Li et al. 2015 cited would not support the new statements, they found other, more relevant, sources to cite. However, they were not as deliberate in their paraphrasing throughout the text. G. Li et al. 2015 cited Nieto et al. 1994, titled “Control of cell behavior during vertebrate development by Slug, a zinc finger gene”, Haupt et al. 2006, titled “Clues from worms: a Slug at Puma promotes the survival of blood progenitors”, Y. Li et al. 2014, titled “Axl mediates tumor invasion and chemosensitivity through PI3K/Akt signaling pathway and is transcriptionally regulated by slug in breast carcinoma”, and He et al. 2012, titled “Ikaros inhibits proliferation and, through upregulation of Slug, increases metastatic ability of ovarian serous adenocarcinoma cells”. Instead of finding a new, more appropriate paper to cite, they changed the titles in the references to “Control of cell behavior during vertebrate development by twist, a zinc finger gene”, “Clues from worms: a twist at Puma promotes the survival of blood progenitors”, “Axl mediates tumor invasion and chemosensitivity through PI3K/Akt signaling pathway and is transcriptionally regulated by Twist in breast carcinoma”, and “Ikaros inhibits proliferation and, through upregulation of twist, increases metastatic ability of ovarian serous adenocarcinoma cells”. The lack of capitalization of “twist” is in the citations as presented in the text. Since these attributions misrepresented not only the key claims in the text, but also the claims made in the titles of the reports cited, the attributed claims count as misquoted.

If considered naively, the numerous substitutions would greatly in-

flate the number of novel claims. However, in this case and with Ullah et al. 2018, it is possible to mitigate this concern by abstracting out details of the sentences. For example, we can take the claims “However, to our knowledge, the potential mechanisms of the CXCL12/CXCR4 pathway in modulation of the EMT process have been largely unknown previously.” and “We were very interested in their relationships and investigated whether Slug signaling was up-regulated by CCL21/CXCR7 pathway to induce EMT in human chondrosarcoma tissues and cells.” in Yao et al. 2016 to be novel claims, as they identify the specific pathway, transcription factor, and type of cancer to be studied as different from those identified in the corresponding claims in G. Li et al. 2015: “However, to our knowledge, the potential mechanisms of the CXCR7 pathway in modulation of the EMT process have been largely unknown previously.” and “We were very interested in their relationships and investigated whether Slug signaling was up-regulated by CCL21/CXCR7 pathway to induce EMT in human chondrosarcoma tissues and cells.” However, in all subsequent claims, we can abstract out these details and consider each substituted word equivalent to the original. For example, we can abstract both “Twist” in Yao et al. 2016 and “Slug” in G. Li et al. 2015 to “the transcription factor of interest”.

When plagiarizing Lv et al. 2015, the authors Dai et al. 2015 applied the same tactics as did the authors Yao et al. 2016. Specifically, they substituted “glioma” for “glioblastoma”, “(SDF-1)/CXCR4” for “EGF”, and “U87” for “U251” and then rewrote some parts of the introduction to replace the resulting obvious misstatements with correctly sourced background information. We can apply the same method of abstraction in order to arrive at appropriate FAIR Metric counts. They were more deliberate about replacing references with those that included claims equivalence matching what they were asserting after the substitutions. But they still included 1 misquoted claim, that the “SDF-1 pathway mainly included the RAS/RAF/MEK/ERK and PI3K/AKT pathways.” While the sources they cited do refer to MEK, ERK, PI3K, and AKT as being part of pathways that include SDF-1, they did not mention RAS or RAF.

Guo et al. 2013 did not attempt any substantive substitutions of content words and instead relied only on paraphrasing and some slight abridgement to obfuscate their plagiarism of Fischbach et al. 2009. Every claim in Guo et al. 2013 had an apparent counterpart in Fischbach et al. 2009. The 1 novel claim found is a paraphrasing that is so garbled that it completely loses the meaning of its counterpart in the original text. Specifically, “Measurement of the SNR and CNR in the images does not allow for the assessment of aesthetic appearance, the depiction of tiny structural details, the distinction of different tissues, the impairment by artifacts, and, hence, the diagnostic value of the images.” in Fischbach et al. 2009 becomes “The diagnosing images can be influence by the artifacts and visualization ability of anatomical details by SNR and CNR at different tissues.” in Guo et al. 2013. One of the 2 misquoted claims is due to another instance of paraphrasing that altered the meaning of the sentence in a nonsensical way, taking an original sentence about eliminating a blood vessel from an image and altering it to be about eliminating the nerves that were originally the focus of the imaging. The other is due to the addition of citations to the plagiarized version of 1 of the novel claims in Fischbach et al. 2009, but attributing it to earlier works.

Among the 8 plagiarism examples analysed here, Su et al. 2005 represents the most overt and explicit plagiarism. Most of the text is a verbatim copy of Schwab et al. 2001 with only sparse rewordings. As such, all claims are either plagiarized or quoted.

In contrast, Wilkinson et al. 2016 did not copy text verbatim from C. Taswell 2007; C. Taswell 2010. Instead, Wilkinson et al. 2016 obfuscated their plagiarism of concepts and ideas by paraphrasing part of the Taswell 2007 collection without citation (Craig, Ambati, Dutta, Kowshik, et al. 2019) and the editors of *Nature Scientific Data* concealed this plagiarism by refusing to correct the omission of citation of the original sources — which constitutes both *idea-laundering plagiarism by authors* and *idea-bleaching censorship by editors* as defined by S. K. Taswell et al. 2020. Each of the 24 claims counted as plagiarized in Wilkinson et al. 2016 (the FAIR-named collection of principles) has a corresponding equivalent in C. Taswell 2007 (the PORTAL-DOORS Project collection of principles). Moreover, the 5 claims counted as novel in Wilkinson et al. 2016 focused merely on building consensus at workshops for their FAIR-named collection. The 6 claims counted as misquoted in Wilkinson et al. 2016 likely resulted from changes to the content of the websites cited as sources.

Discussion

The 8 cases of plagiarism in Table 1 illustrate the complexity and diversity of real-world plagiarism and demonstrate that the current version of the FAIR Metrics are useful in real-world peer reviews. The FAIR Metrics did not indicate any sign of plagiarism in the negative control case of the example pair C. Taswell 2007 and Mons 2005. Thus, the requirement of equivalence of meaning can assist in detecting plagiarism while not yielding false positives for plagiarism and possibly allegations of plagiarism in the scenario of different author groups writing about the same topics within the same field of study contemporaneously. The cases retracted for plagiarism show that the FAIR Metrics can positively identify cases of explicit plagiarism even with mild paraphrasing across problem domains as diverse as green technology (Ullah et al. 2018), dermatology (Gnat et al. 2022), and neuroscience (Uddin et al. 2022). Future work will more formally evaluate the sensitivity and specificity of the FAIR Metrics for the detection of plagiarism in various scenarios.

Although the FAIR Metrics provide helpful insights and alerts, the current version does not obviate the need for other forms of textual analysis, both lexical and semantic, to identify and understand the full nature and extent of plagiarism in research communications. In particular, the Q counts can be spuriously high in that many of the passages in the plagiarizing papers with correct attributions have nevertheless been plagiarized from the comparison papers. Since neither T nor C are the original references for the ideas presented, and since both attribute them to prior sources that do present such concepts, the FAIR Metric evaluation procedure as currently practiced deems the copies of such claims in both works to be valid quoted claims, even if they have identical wording. We originally designed the FAIR Metrics to evaluate the quality of primary research articles, which should present original results and analyses balanced with context from the existing literature (Craig and C. Taswell 2018). In their present form, they would not be suitable for a comparison of 2 pure reviews of the literature that summarize previously published content from the historical record devoid of any attempt in the literature reviews to provide commentary, analysis, or synthesis with new concepts, ideas, and claims. While we plan to develop FAIR Metrics customized for different kinds of scholarly research communications, current lexical and semantic comparison methods can still serve as complementary tools for use with the FAIR Metrics analyses. Regardless, when automated with machine algorithms these comparison evaluations for the detection of plagiarism should always be subject to final review by human analysts.

Table 1: FAIR Metrics for example comparison pairs listed in F_P descending order

| Pair | Test (T) text | Retracted? | Comparison (C) text | M | N | P | Q | F_M | F_N | F_P | F_Q |
|------|---------------------------------------|------------|-------------------------------------------|-----|-----|-----|-----|-------|-------|-------|-------|
| 1 | C. Taswell 2007 | no | Mons 2005 | 0 | 20 | 0 | 22 | 1.00 | 0.05 | 1.00 | 1.00 |
| 2 | Uddin et al. 2022 | yes | Foster et al. 2019 | 0 | 18 | 18 | 87 | 0.83 | 0.56 | 0.66 | 0.83 |
| 3 | Gnat et al. 2022 | yes | Hoog et al. 2016 | 0 | 3 | 10 | 30 | 0.75 | 0.63 | 0.50 | 0.75 |
| 4 | Wilkinson et al. 2016 | no | C. Taswell 2007 | 6 | 5 | 24 | 28 | 0.38 | 0.37 | 0.07 | 0.48 |
| 5 | Yao et al. 2016 | yes | G. Li et al. 2015 | 4 | 2 | 11 | 9 | 0.21 | 0.27 | -0.08 | 0.38 |
| 6 | Dai et al. 2015 | yes | Lv et al. 2015 | 1 | 2 | 18 | 14 | 0.39 | 0.34 | -0.12 | 0.42 |
| 7 | Guo et al. 2013 | yes | Fischbach et al. 2009 | 2 | 1 | 13 | 10 | 0.32 | 0.346 | -0.12 | 0.40 |
| 8 | Ullah et al. 2018 | yes | Sansaniwal and Kumar 2015 | 31 | 3 | 7 | 2 | -0.73 | -0.02 | -0.13 | 0.05 |
| 9 | Su et al. 2005 | yes | Schwab et al. 2001 | 0 | 0 | 20 | 12 | 0.38 | 0.38 | -0.25 | 0.38 |

M Misquoted, N Novel, P Plagiarized, Q Quoted Counts; F_M Misquoted, F_N Novel, F_P Plagiarized, F_Q Quoted FAIR Metrics.

The labor-intensive evaluation process required of human analysts, as demonstrated in this report, remains another current limitation on the practical utility of the FAIR Metrics for screening large numbers of documents. Even the pairwise comparison approach used in the present work requires that the reviewer list all statements in the test text, identify which ones are significant enough to be key claims, search the comparison text for equivalent claims, and perform at least a cursory search of every cited text for equivalents of the claims attributed to them. While more work than a typical peer review, this process can nevertheless be used as an important method for keeping the provenance and development of ideas traceable and verifiable when evaluating suspected cases of plagiarism. Publishers can make it worthwhile for reviewers by publishing the evaluation documents as citable works in their own right, thus featuring the scholarship and analytical skills of the reviewers ([Craig, Lee, et al. 2022](#)). Furthermore, making these records not only readable by humans but also by machines as subject-verb-object triples and linked quads will enhance their potential application and use in both scenarios of rapid screening of large numbers of documents as well as careful evaluation of a small number of documents suspected of plagiarism. The resulting linked knowledge graph can also be explored by semantic search and reasoning engines and provides a resource for the development and testing of tools to automate parts of the FAIR Metrics evaluation process. Maintaining a corpus of test cases known to contain matching entities can be useful for testing named entity recognition approaches such as [Taufiq et al. 2023](#) and [Khadilkar et al. 2018](#), which could be adapted to produce matched claim pairs for piping into an automated FAIR Metrics calculator.

Conclusion

We have demonstrated that the FAIR Metrics provide a quantitative method of evaluating the extent to which a scholarly communication adheres to a code of conduct of fairness when discussing and citing relevant research in the field, taking ideas from previously published literature, and properly crediting the original sources. We have shown that even a simple evaluation procedure against a limited pool of comparison texts yields differences in measures which can assist peer review to assess concerns about plagiarism, misrepresentation, citational justice, and fairness. We have created searchable online repositories of NPDS records with semantic representations of FAIR Metric analyses that serve as a prototype for a more reproducible, verifiable, and accountable approach to open and transparent peer review.

Citation

Brainiacs 2023 Volume 4 Issue 2 Edoc D5B2734F2

Title: "Managing Lexical-Semantic Hybrid Records of FAIR Metrics Analyses with the NPDS Cyberinfrastructure"

Authors: Adam Craig, Anousha Athreya, Carl Taswell

Dates: created 2023-06-16, received 2023-10-03, presented 2023-10-09, updated 2023-12-27, published 2023-12-27, endorsed 2023-12-31

Copyright: © 2023 Brain Health Alliance

Contact: [A Craig at BHAVI](#)

URL: [Brainiacsjournal.org/arc/pub/Craig2023MLSHRFMA](https://brainiacsjournal.org/arc/pub/Craig2023MLSHRFMA)

PDP: [/Nexus/Brainiacs/Craig2023MLSHRFMA](#)

DOI: [/10.48085/D5B2734F2](https://doi.org/10.48085/D5B2734F2)

References

- [1] A. Altheneyan and M. E. B. Menai. "Evaluation of state-of-the-art paraphrase identification and its application to automatic plagiarism detection." *International Journal of Pattern Recognition and Artificial Intelligence* 34.04 (2020), p. 2053004 (cited p. 1).
- [2] A. Athreya, S. K. Taswell, S. Mashkoor, and C. Taswell. "Essential Question: 'Equal or Equivalent Entities?' About Two Things as Same, Similar, or Different." In: *2020 Second International Conference on Transdisciplinary AI (TransAI)*. 2020, pp. 123–124. DOI: [10.1109/TransAI49837.2020.00028](https://doi.org/10.1109/TransAI49837.2020.00028) (cited p. 3).
- [3] A. Athreya, S. K. Taswell, S. Mashkoor, and C. Taswell. "The Essential Enquiry 'Equal or Equivalent Entities?' About Two Things as Same, Similar, Related, or Different." *Brainiacs Journal of Brain Imaging And Computing Sciences* 1.1, PEDADC885 (1 Dec. 30, 2020), pp. 1–7. DOI: [10.48085/PEDADC885](https://doi.org/10.48085/PEDADC885). URL: <https://brainiacsjournal.org/arc/pub/Athreya2020EEEEEE> (cited pp. 2, 3).
- [4] G. Cabanac and C. Labbé. "Prevalence of nonsensical algorithmically generated papers in the scientific literature." *Journal of the Association for Information Science and Technology* 72.12 (2021), pp. 1461–1476 (cited p. 2).
- [5] Copyleaks. *AI Content Detector Continues To Be Confirmed As Most Accurate By Third-Party Studies*. Oct. 30, 2023. URL: <https://copyleaks.com/blog/ai-detector-continues-top-accuracy-third-party> (cited p. 1).

- [6] A. Craig, A. Ambati, S. Dutta, P. Kowshik, S. Nori, S. K. Taswell, Q. Wu, and C. Taswell. "DREAM Principles and FAIR Metrics from the PORTAL-DOORS Project for the Semantic Web." In: *2019 IEEE 11th International Conference on Electronics, Computers and Artificial Intelligence (ECAI)* (June 28, 2019). Pitesti, Romania: IEEE, June 2019, pp. 1–8. DOI: [10.1109/ECAI46879.2019.9042003](https://portaldoors.org/pub/docs/ECAI2019DREAMFAIR0612.pdf). URL: <https://portaldoors.org/pub/docs/ECAI2019DREAMFAIR0612.pdf> (cited pp. 2–5).
- [7] A. Craig, A. Ambati, S. Dutta, A. Mehrotra, S. K. Taswell, and C. Taswell. "Definitions, Formulas, and Simulated Examples for Plagiarism Detection with FAIR Metrics." In: *2019 ASIS&T 82nd Annual Meeting* (Oct. 19, 2019). Vol. 56. Melbourne, Australia: Wiley, 2019, pp. 51–57. DOI: [10.1002/PRA2.6](https://portaldoors.org/pub/docs/ASIST2019FairMetrics0611.pdf). URL: <https://portaldoors.org/pub/docs/ASIST2019FairMetrics0611.pdf> (cited pp. 2, 3).
- [8] A. Craig, A. Athreya, and C. Taswell. "Example evaluations of plagiarism cases using FAIR Metrics and the PDP-DREAM Ontology." *IEEE eScience 2023* (Oct. 13, 2023) (cited pp. 1, 2).
- [9] A. Craig, C. Lee, N. Bala, and C. Taswell. "Motivating and Maintaining Ethics, Equity, Effectiveness, Efficiency, and Expertise in Peer Review." *Brainiacs Journal of Brain Imaging And Computing Sciences* 3.1, 15B147D9D (1 June 30, 2022), pp. 1–21. DOI: [10.48085/15B147D9D](https://BrainiacsJournal.org/arc/pub/Craig2022MMEPR). URL: <https://BrainiacsJournal.org/arc/pub/Craig2022MMEPR> (cited pp. 2, 6).
- [10] A. Craig and C. Taswell. "Formulation of FAIR Metrics for Primary Research Articles." In: *2018 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*. Madrid, Spain: IEEE, Dec. 3, 2018, pp. 1632–1635. DOI: [10.1109/BIBM.2018.8621399](https://doi.org/10.1109/BIBM.2018.8621399) (cited p. 5).
- [11] A. Craig and C. Taswell. "PDP-DREAM Software for Integrating Multimedia Data with Interoperable Repositories." *Brainiacs Journal of Brain Imaging And Computing Sciences* 2.1, HA46280EF (1 Dec. 31, 2021), pp. 1–6. DOI: [10.48085/HA46280EF](https://BrainiacsJournal.org/arc/pub/Craig2021SIMDIR). URL: <https://BrainiacsJournal.org/arc/pub/Craig2021SIMDIR> (cited pp. 2, 4).
- [12] C. Dai, S. Lv, R. Shi, J. Ding, et al. "RETRACTED ARTICLE: nuclear protein C23 on the cell surface plays an important role in activation of CXCR4 signaling in glioblastoma." *Molecular neurobiology* 52 (2015), pp. 1521–1526 (cited pp. 3, 5, 6).
- [13] A. W. Deshmukh, M. N. Varma, et al. "Investigation of solar drying of ginger (*Zingiber officinale*): Empirical modelling, drying characteristics, and quality study." *Chinese Journal of Engineering* 2014 (2014), pp. 1–7 (cited p. 4).
- [14] N. Drummond, M. Horridge, and H. Knublauch. "Protégé-OWL tutorial." In: *8th International Protégé Conference*. 2005 (cited p. 4).
- [15] S. Dutta, P. Kowshik, A. Ambati, S. Nori, S. K. Taswell, and C. Taswell. "Managing Scientific Literature with Software from the PORTAL-DOORS Project." In: *2019 IEEE 15th International Conference on eScience (eScience)* (Sept. 24, 2019). San Diego, California: IEEE, Sept. 2019. DOI: [10.1109/eScience.2019.00081](https://portaldoors.org/pub/docs/BCDC2019PdpDemo0806.pdf). URL: <https://portaldoors.org/pub/docs/BCDC2019PdpDemo0806.pdf> (cited p. 4).
- [16] S. Dutta, K. Uhegbu, S. Nori, S. Mashkoor, S. K. Taswell, and C. Taswell. "DREAM Principles from the PORTAL-DOORS Project and NPDS Cyberinfrastructure." In: *2020 IEEE 14th International Conference on Semantic Computing (ICSC)*. IEEE, Feb. 4, 2020, pp. 211–216. DOI: [10.1109/ICSC.2020.00044](https://portaldoors.org/pub/docs/ECAI2019DREAMFAIR0612.pdf). URL: <https://portaldoors.org/pub/docs/ECAI2019DREAMFAIR0612.pdf> (cited p. 2).
- [17] T. A. E. Eisa, N. Salim, and A. Abdelmaboud. "Content-based scientific figure plagiarism detection using semantic mapping." In: *Emerging Trends in Intelligent Computing and Informatics*. Ed. by F. Saeed, F. Mohammed, and N. Gazem. Springer International Publishing, 2020, pp. 420–427. ISBN: 978-3-030-33582-3 (cited p. 2).
- [18] F. Fischbach, M. Müller, and H. Bruhn. "High-resolution depiction of the cranial nerves in the posterior fossa (N III–N XII) with 2D fast spin echo and 3D gradient echo sequences at 3.0 T." *Clinical imaging* 33.3 (2009), pp. 169–174 (cited pp. 3, 5, 6).
- [19] E. M. Foster, A. Dangla-Valls, S. Lovestone, E. M. Ribe, and N. J. Buckley. "Clusterin in Alzheimer's Disease: Mechanisms, Genetics, and Lessons From Other Pathologies." *Frontiers in Neuroscience* 13, 164 (Feb. 2019). ISSN: 1662-453X. DOI: [10.3389/fnins.2019.00164](https://doi.org/10.3389/fnins.2019.00164) (cited pp. 3, 6).
- [20] C. A. Gao, F. M. Howard, N. S. Markov, E. C. Dyer, S. Ramesh, Y. Luo, and A. T. Pearson. "Comparing scientific abstracts generated by ChatGPT to original abstracts using an artificial intelligence output detector, plagiarism detector, and blinded human reviewers." *bioRxiv* (2022), pp. 2022–12 (cited p. 2).
- [21] M. Gaudino, N. B. Robinson, K. Audisio, M. Rahouma, U. Benedetto, P. Kurlansky, and S. E. Fremes. "Trends and characteristics of retracted articles in the biomedical literature, 1971 to 2020." *JAMA internal medicine* 181.8 (2021), pp. 1118–1121 (cited p. 2).
- [22] S. Gnat, A. Nowakiewicz, and P. Zięba. "Retraction of: Sebastian Gnat, Aneta Nowakiewicz, Przemysław Zięba: Taxonomy of dermatophytes – the classification systems may change but the identification problems remain the same." *Postępy Mikrobiologii - Advancements of Microbiology* 61.4 (2022), pp. 261–261. DOI: [10.2478/am-2022-013](https://doi.org/10.2478/am-2022-013) (cited pp. 3–6).
- [23] GPTZero. *Our detection technology*. 2023. URL: <https://gptzero.me/technology> (cited p. 1).
- [24] Z.-Y. Guo, J. Chen, H.-Y. Liao, Q.-Y. Cheng, S.-X. Fu, C.-X. Chen, D. Yu, et al. *RETRACTED: High-resolution MRI of cranial nerves in posterior fossa at 3.0 T*. 2013 (cited pp. 3, 5, 6).
- [25] S. Haupt, O. Alsheich-Bartok, and Y. Haupt. "Clues from worms: a Slug at Puma promotes the survival of blood progenitors." *Cell death and differentiation* 13.6 (2006), p. 913 (cited p. 4).
- [26] L.-C. He, F.-H. Gao, H.-Z. Xu, S. Zhao, C.-M. Ma, J. Li, S. Zhang, and Y.-L. Wu. "Ikaros inhibits proliferation and, through upregulation of Slug, increases metastatic ability of ovarian serous adenocarcinoma cells." *Oncology reports* 28.4 (2012), pp. 1399–1405 (cited p. 4).
- [27] G. S. de Hoog, K. Dukik, M. Monod, A. Packeu, et al. "Toward a Novel Multilocus Phylogenetic Taxonomy for the Dermatophytes." *Mycopathologia* 182.1–2 (Oct. 2016), pp. 5–31. ISSN: 1573-0832. DOI: [10.1007/s11046-016-0073-9](https://doi.org/10.1007/s11046-016-0073-9) (cited pp. 3, 6).
- [28] IEEE, ed. *CrossCheck Information*. 2023. URL: <https://www.ieee.org/publications/rights/cross-check-main.html> (cited p. 1).
- [29] S. Javadi-Moghaddam, F. Roosta, and A. Noroozi. "Weighted semantic plagiarism detection approach based on AHP decision model." *Accountability in Research* 29.4 (2022), pp. 203–223 (cited p. 2).
- [30] M. Jiffriya, M. Jahan, and R. Ragel. "Plagiarism detection tools and techniques: A comprehensive survey." *Journal of Science-FAS-SEUSL* 2.02 (2021), pp. 47–64 (cited p. 2).
- [31] K. Khadiikar, S. Kulkarni, and P. Bone. "Plagiarism detection using semantic knowledge graphs." In: *2018 Fourth International Conference on Computing Communication Control and Automation (ICCUBEA)*. IEEE, 2018, pp. 1–6 (cited p. 6).
- [32] M. Khalil and E. Er. "Will ChatGPT get you caught? Rethinking of plagiarism detection." *arXiv preprint arXiv:2302.04335* (2023) (cited p. 1).
- [33] J. H. Kirchner, L. Ahmad, S. Aaronson, and J. Leike. *New AI classifier for indicating AI-written text*. Jan. 31, 2023. URL: <https://openai.com/blog/new-ai-classifier-for-indicating-ai-written-text> (cited p. 1).

- [34] C. Labbé and D. Labbé. “Duplicate and fake publications in the scientific literature: how many SCiGen papers in computer science?” *Scientometrics* 94 (2013), pp. 379–396 (cited p. 2).
- [35] G. Li, Y. Yang, S. Xu, L. Ma, M. He, and Z. Zhang. “Slug signaling is up-regulated by CCL21/CXCR7 to induce EMT in human chondrosarcoma.” *Medical Oncology* 32.2 (2015), p. 2 (cited pp. 3–6).
- [36] Y. Li, L. Jia, D. Ren, C. Liu, Y. Gong, N. Wang, X. Zhang, and Y. Zhao. “Axl mediates tumor invasion and chemosensitivity through PI3K/Akt signaling pathway and is transcriptionally regulated by slug in breast carcinoma.” *IUBMB life* 66.7 (2014), pp. 507–518 (cited p. 4).
- [37] S. Lv, C. Dai, Y. Liu, B. Sun, R. Shi, M. Han, R. Bian, and R. Wang. “Cell surface protein C23 affects EGF-EGFR induced activation of ERK and PI3K-AKT pathways.” *Journal of Molecular Neuroscience* 55 (2015), pp. 519–524 (cited pp. 3, 5, 6).
- [38] B. Mons. “Which gene did you mean?” *BMC Bioinformatics* 6, 142 (2005). DOI: [10.1186/1471-2105-6-142](https://doi.org/10.1186/1471-2105-6-142) (cited pp. 3–6).
- [39] M. A. Nieto, M. G. Sargent, D. G. Wilkinson, and J. Cooke. “Control of cell behavior during vertebrate development by Slug, a zinc finger gene.” *Science* 264.5160 (1994), pp. 835–839 (cited p. 4).
- [40] M. S. Orenstrakh, O. Karnalim, C. A. Suarez, and M. Liut. *Detecting LLM-Generated Text in Computing Education: A Comparative Study for ChatGPT Cases*. 2023. DOI: [10.48550/ARXIV.2307.07411](https://doi.org/10.48550/ARXIV.2307.07411) (cited p. 1).
- [41] *Retraction Watch Database User Guide*. Apr. 2023. URL: <https://retractionwatch.com/retraction-watch-database-user-guide/> (cited p. 3).
- [42] S. K. Sansaniwal and M. Kumar. “Analysis of ginger drying.” *Journal of Mechanical Engineering and Sciences* 9 (2015), pp. 1671–1685. DOI: [10.15282/jmes.9.2015.13.0161](https://doi.org/10.15282/jmes.9.2015.13.0161) (cited pp. 3, 6).
- [43] K. Santos-d’Amorim, T. Wang, B. Lund, and R. N. Macedo Dos Santos. “From plagiarism to scientific paper mills: a profile of retracted articles within the SciELO Brazil collection.” *Ethics & Behavior* 34.1 (Nov. 2022), pp. 1–18. ISSN: 1532-7019. DOI: [10.1080/10508422.2022.2141747](https://doi.org/10.1080/10508422.2022.2141747) (cited p. 2).
- [44] S. Schwab, D. Georgiadis, J. Berroushot, P. D. Schellinger, C. Graffagnino, and S. A. Mayer. “Feasibility and safety of moderate hypothermia after massive hemispheric infarction.” *Stroke* 32.9 (2001), pp. 2033–2035 (cited pp. 3, 5, 6).
- [45] J. Su, Y. Qiou, Z. Chen, and Y. Chen. “Feasibility and safety of moderate hypothermia after acute ischemic stroke (vol 21, pg 353, 2003).” *INTERNATIONAL JOURNAL OF DEVELOPMENTAL NEUROSCIENCE* 23.4 (2005), pp. 411–411 (cited pp. 3, 5, 6).
- [46] C. Taswell. “DOORS to the Semantic Web and Grid with a PORTAL for Biomedical Computing.” *IEEE Transactions on Information Technology in Biomedicine* 12.2 (2 Mar. 2007). In the Special Section on Bio-Grid published online 3 Aug. 2007, pp. 191–204. ISSN: 1089-7771. DOI: [10.1109/TITB.2007.905861](https://doi.org/10.1109/TITB.2007.905861) (cited pp. 3–6).
- [47] C. Taswell. “A Distributed Infrastructure for Metadata about Metadata: The HDMM Architectural Style and PORTAL-DOORS System.” *Future Internet* 2.2 (2010), pp. 156–189. ISSN: 1999-5903. DOI: [10.3390/FI2020156](https://doi.org/10.3390/FI2020156). URL: <https://www.mdpi.com/1999-5903/2/2/156/> (cited pp. 3, 5).
- [48] S. K. Taswell, C. Triggler, J. Vayo, S. Dutta, and C. Taswell. “The Hitchhiker’s Guide to Scholarly Research Integrity.” In: *2020 ASIS&T 83rd Annual Meeting* (Oct. 22, 2020). Vol. 57. Wiley, 2020, e223. DOI: [10.1002/ppra2.223](https://doi.org/10.1002/ppra2.223). URL: <https://asistdl.onlinelibrary.wiley.com/doi/abs/10.1002/ppra2.223> (cited pp. 2, 5).
- [49] U. Taufiq, R. Pulungan, and Y. Suyanto. “Named entity recognition and dependency parsing for better concept extraction in summary obfuscation detection.” *Expert Systems with Applications* (2023), p. 119579 (cited p. 6).
- [50] H. H. Thorp. *ChatGPT is fun, but not an author*. 2023 (cited p. 1).
- [51] M. S. Uddin, M. T. Kabir, M. M. Begum, M. S. Islam, T. Behl, and G. M. Ashraf. *Retraction Note to: Exploring the Role of CLU in the Pathogenesis of Alzheimer’s Disease*. 2022. DOI: [10.1007/s12640-022-00519-1](https://doi.org/10.1007/s12640-022-00519-1) (cited pp. 3–6).
- [52] F. Ullah, M. Kang, M. K. Khattak, and S. Wahab. “Retracted: Experimentally investigated the asparagus (*Asparagus officinalis* L.) drying with flat-plate collector under the natural convection indirect solar dryer.” *Food Science & Nutrition* 6.6 (2018), pp. 1357–1357 (cited pp. 3–6).
- [53] T. Vrbancic and A. Meštrović. “The struggle with academic plagiarism: Approaches based on semantic similarity.” In: *2017 40th International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO)*. IEEE, 2017, pp. 870–875 (cited p. 1).
- [54] M. D. Wilkinson, M. Dumontier, I. J. Aalbersberg, G. Appleton, et al. “The FAIR Guiding Principles for scientific data management and stewardship.” *Scientific data* 3.1 (2016), pp. 1–9 (cited pp. 3, 5, 6).
- [55] J. Xiong and T. Huang. “An effective method to identify machine automatically generated paper.” In: *2009 Pacific-Asia Conference on Knowledge Engineering and Software Engineering*. IEEE, 2009, pp. 101–102 (cited p. 2).
- [56] C. Yao, P. Li, H. Song, F. Song, Y. Qu, X. Ma, R. Shi, and J. Wu. “RETRACTED ARTICLE: CXCL12/CXCR4 Axis Upregulates Twist to Induce EMT in Human Glioblastoma.” *Molecular neurobiology* 53 (2016), pp. 3948–3953 (cited pp. 3–6).
- [57] L. Young. *iThenticate 2.0: Advancing research integrity with AI writing detection*. Ed. by turnitin. Nov. 1, 2023. URL: <https://www.turnitin.com/blog/ithenticate-2-0-advancing-research-integrity-with-ai-writing-detection> (cited p. 1).

Reproducibility, Validity, and Integrity in Scholarly Research

Carl Taswell



Reproducibility, Validity, and Integrity in Scholarly Research: What Accountability for Willful Disregard?*

Carl Taswell†

Commentary

The [Aims and Scope](#) of the [Brainiacs Journal](#) encompass the *Information And Communication Sciences* in addition to the *Imaging And Computing Sciences*. These IACS fields find a nexus in the ethically and legally challenging realm of establishing credibility of witnesses in courts of law. How do we differentiate evidentiary fact from fiction and truth from lies when related as stories told by witnesses to juries in court? Should we trust scientific reports when researchers claim that brain imaging, physiologic monitoring devices, and psychological questionnaires ([US Congress Office of Technology Assessment 1983](#); [Langleben and Moriarty 2013](#); [Meltzer et al. 2013](#); [Walczyk et al. 2018](#); [Hsu et al. 2019](#); [Díaz Soto and Borbón 2022](#)) can or cannot be used as *lie detectors* to reveal when a person is not telling the truth? In the current era of information wars that have spread from politicians to scientists, how do we maintain reproducibility, validity, and integrity in scholarly research for science, engineering, and medicine?

Brain Health Alliance (BHA) has been studying these questions for the past several years ([Craig, Ambati, et al. 2019](#); [Taswell, Triggles, et al. 2020](#); [Athreya et al. 2020](#); [Taswell, Athreya, et al. 2021](#); [Craig, Lee, et al. 2022](#)). The BHA Virtual Institute (BHAVI) has now hosted two annual Guardians conferences ([Guardians 2022](#) and [Guardians 2023](#)) focused on truth and integrity in science ([Craig, Taswell, et al. 2022](#); [Taswell and Craig 2023](#)). BHAVI has answered the question “Who are the Guardians of Truth and Integrity?” each of the past two years by honoring Dr. Peter Wilmschurst and Dr. Anthony Fauci, respectively, as the [2022 Guardian](#) and [2023 Guardian](#). Next year for [Guardians 2024](#), we pose the following questions seeking answers:

- What accountability should be imposed on researchers for willful disregard of reproducibility, validity, and integrity in scholarly research? What should be the sanctions for those researchers who violate professional codes of conduct?
- Doctors must be licensed to practice medicine with patients. Lawyers must be licensed to practice law with clients. Teachers in schools must be licensed to teach students. Why do we not require researchers to be licensed to conduct research?
- Which organizations should respond to concerns and complaints

about fraud, falsification, plagiarism, and misconduct in research? Which organizations should impose and enforce the sanctions for those researchers who violate professional codes of conduct?

- To which independent and impartial venue (one which is devoid of conflicts of interest and mandated to disclose the legal and financial nature of the entity by clarifying the funding sources and controlling parties), should we submit complaints when a researcher with less power and money becomes a victim exploited by a researcher with more power and money?
- Educators at academic universities in our communities should be leaders and example role models who teach and promote moral, ethical, civil, courteous, tolerant, and respectful behavior between and amongst all members of our communities. How should we heal and cure the worsening triple-G problem in academia of *Grooming, Gaslighting, and Ghosting*?
- What is the meaning and purpose of policies and procedures in academic policy manuals if the rules are not imposed and enforced for all members of the university including both teachers and students? What is the meaning and purpose of ethical codes of conduct if the guidelines are not imposed and enforced by the professional societies that adopt and promote such codes of conduct for their members?
- Can truth and integrity exist in academic science, engineering, and medicine without accountability for willful disregard?

Report submissions for [Guardians 2024](#) will open on 9 January 2024. Author presentations will be held at an online event via GoToMeeting videoconference at meet.goto.com/965055533 on 9 October 2024. For further information, please contact [guardians at BHAVI.us](#).

Citation

Brainiacs 2023 Volume 4 Issue 2 Edoc L3570F30F

Title: “Reproducibility, Validity, and Integrity in Scholarly Research: What Accountability for Willful Disregard?”

Authors: Carl Taswell

Dates: created 2023-10-03, received 2023-10-03, presented 2023-10-09, updated 2023-12-31, published 2023-12-31,

Copyright: © 2023 Brain Health Alliance

Contact: [CTaswell at Brain Health Alliance](#)

*Presented 2023-10-09 with [slides](#) and [video](#) at [Guardians 2023](#)

† Authors affiliated with Brain Health Alliance Virtual Institute, Ladera Ranch, CA 92694 USA; correspondence to [CTaswell at Brain Health Alliance](#).

URL: [Brainiacsjournal.org/arc/pub/Taswell2023RVISR](https://brainiacsjournal.org/arc/pub/Taswell2023RVISR)

PDP: [/Nexus/Brainiacs/Taswell2023RVISR](https://Nexus/Brainiacs/Taswell2023RVISR)

DOI: [10.48085/L3570F30F](https://doi.org/10.48085/L3570F30F)

References

- [1] A. Athreya, S. K. Taswell, S. Mashkoo, and C. Taswell. "The Essential Enquiry 'Equal or Equivalent Entities?' About Two Things as Same, Similar, Related, or Different." *Brainiacs Journal of Brain Imaging And Computing Sciences* 1.1, PEDADC885 (1 Dec. 30, 2020), pp. 1–7. DOI: [10.48085/PEDADC885](https://doi.org/10.48085/PEDADC885). URL: <https://brainiacsjournal.org/arc/pub/Athreya2020EEEEEE> (cited p. 1).
- [2] A. Craig, A. Ambati, S. Dutta, P. Kowshik, S. Nori, S. K. Taswell, Q. Wu, and C. Taswell. "DREAM Principles and FAIR Metrics from the PORTAL-DOORS Project for the Semantic Web." In: *2019 IEEE 11th International Conference on Electronics, Computers and Artificial Intelligence (ECAI)* (June 28, 2019). Pitesti, Romania: IEEE, June 2019, pp. 1–8. DOI: [10.1109/ECAI46879.2019.9042003](https://doi.org/10.1109/ECAI46879.2019.9042003). URL: <https://portaldoors.org/pub/docs/ECAI2019DREAMFAIR0612.pdf> (cited p. 1).
- [3] A. Craig, C. Lee, N. Bala, and C. Taswell. "Motivating and Maintaining Ethics, Equity, Effectiveness, Efficiency, and Expertise in Peer Review." *Brainiacs Journal of Brain Imaging And Computing Sciences* 3.1, I5B147D9D (1 June 30, 2022), pp. 1–21. DOI: [10.48085/I5B147D9D](https://doi.org/10.48085/I5B147D9D). URL: <https://BrainiacsJournal.org/arc/pub/Craig2022MMEEPR> (cited p. 1).
- [4] A. Craig, S. K. Taswell, A. Athreya, and C. Taswell. "Who are the Guardians of Truth and Integrity?" *Brainiacs Journal of Brain Imaging And Computing Sciences* 3.2 (Dec. 29, 2022). ISSN: 2766-6883. DOI: [10.48085/y9f719aa4](https://doi.org/10.48085/y9f719aa4) (cited p. 1).
- [5] J. M. Díaz Soto and D. Borbón. "Neurorights vs. neuroprediction and lie detection: The imperative limits to criminal law." *Frontiers in Psychology* 13 (Dec. 2022). ISSN: 1664-1078. DOI: [10.3389/fpsyg.2022.1030439](https://doi.org/10.3389/fpsyg.2022.1030439) (cited p. 1).
- [6] C.-W. Hsu, C. Begliomini, T. Dall'Acqua, and G. Ganis. "The effect of mental countermeasures on neuroimaging-based concealed information tests." *Human Brain Mapping* 40.10 (Mar. 2019), pp. 2899–2916. ISSN: 1097-0193. DOI: [10.1002/hbm.24567](https://doi.org/10.1002/hbm.24567) (cited p. 1).
- [7] D. D. Langleben and J. C. Moriarty. "Using brain imaging for lie detection: Where science, law, and policy collide." *Psychology, Public Policy, and Law* 19.2 (May 2013), pp. 222–234. ISSN: 1076-8971. DOI: [10.1037/a0028841](https://doi.org/10.1037/a0028841) (cited p. 1).
- [8] C. C. Meltzer, G. Sze, K. S. Rommelfanger, K. Kinlaw, J. D. Banja, and P. R. Wolpe. "Guidelines for the Ethical Use of Neuroimages in Medical Testimony: Report of a Multidisciplinary Consensus Conference." *American Journal of Neuroradiology* 35.4 (Aug. 2013), pp. 632–637. ISSN: 1936-959X. DOI: [10.3174/ajnr.a3711](https://doi.org/10.3174/ajnr.a3711) (cited p. 1).
- [9] S. K. Taswell, A. Athreya, M. Akella, and C. Taswell. "Truth in Science." *Brainiacs Journal of Brain Imaging and Computing Sciences* 2.1 (1 Dec. 31, 2021), pp. 1–9. DOI: [10.48085/M85EC99EE](https://doi.org/10.48085/M85EC99EE). URL: <https://BrainiacsJournal.org/arc/pub/Taswell2021Truth> (cited p. 1).
- [10] S. K. Taswell and A. Craig. "Who are the Guardians of Truth and Integrity?" *Brainiacs Journal of Brain Imaging And Computing Sciences* 4.2 (Dec. 31, 2023). DOI: [10.48085/y331839fb](https://doi.org/10.48085/y331839fb) (cited p. 1).
- [11] S. K. Taswell, C. Triggler, J. Vayo, S. Dutta, and C. Taswell. "The Hitchhiker's Guide to Scholarly Research Integrity." In: *2020 ASIS&T 83rd Annual Meeting* (Oct. 22, 2020). Vol. 57. Wiley, 2020, e223. DOI: [10.1002/pra2.223](https://doi.org/10.1002/pra2.223). URL: <https://asistdl.onlinelibrary.wiley.com/doi/abs/10.1002/pra2.223> (cited p. 1).
- [12] US Congress Office of Technology Assessment. *Scientific Validity of Polygraph Testing: A Research Review and Evaluation*. Tech. rep. Washington, DC: United States Congress, Nov. 1983. URL: <https://sgp.fas.org/othergov/polygraph/ota/index.html> (cited p. 1).
- [13] J. J. Walczyk, N. Sewell, and M. B. DiBenedetto. "A Review of Approaches to Detecting Malingering in Forensic Contexts and Promising Cognitive Load-Inducing Lie Detection Techniques." *Frontiers in Psychiatry* 9 (Dec. 2018). ISSN: 1664-0640. DOI: [10.3389/fpsyg.2018.00700](https://doi.org/10.3389/fpsyg.2018.00700) (cited p. 1).

Who are the Guardians of Truth and Integrity?

S. Koby Taswell, Adam Craig



Who are the Guardians of Truth and Integrity?*

S. Koby Taswell and Adam Craig†

Abstract

On October 9th, Brain Health Alliance (BHA, a 501c3 not-for-profit organization) hosted Guardians 2023, our 2nd annual conference entitled “Who are the Guardians of Truth and Integrity?” The Guardians conferences focus on the global impact of information cyberwars on citizens of planet Earth. Internationally in media of many forms, information has been warped and twisted, resulting in disease, death, and destruction around the globe. To combat the spread of lies and extremist propaganda, the Guardians conferences strive to promote better understanding and awareness about the harm caused by information wars, and to advance learning and knowledge about how to support truth and integrity through technological and sociological research and education for communications in science, engineering, and medicine.

Keywords

Research integrity, citational justice, publishing ethics, scientific truth, GWAS, fake stuff, academic ghosting, FAIR Metrics.

Contents

| | |
|-------------------------------------------------|---|
| Guardians 2023 Program | 1 |
| 2023 Guardian: Anthony S. Fauci | 2 |
| Julian Hecker and Nan Laird | 2 |
| Walter Scheirer | 3 |
| Alicia Andrzejewski | 3 |
| Daniel Kristanto | 3 |
| Koby Taswell | 4 |
| Adam Craig | 4 |
| Carl Taswell | 4 |
| Citation | 4 |
| References | 5 |

Guardians 2023 Program

Guardians 2023 was held on October 9th as a half-day online event with 3 invited speakers:

- Dr. Nan Laird, Harvard University, Boston MA
- Dr. Walter Scheirer, University of Notre Dame, Notre Dame IN
- Dr. Alicia Andrzejewski, William & Mary, Williamsburg VA

who gave insightful presentations related to truth, integrity, information, and communication relevant to the current state of affairs for scientific research in today’s world. The workshop began with recognition of Dr. Anthony Fauci as our 2023 Guardian of Truth and Integrity.

Opening Remarks

- 09:00 Julie Neidich, BHAVI 2023 Guardian: Anthony S. Fauci (2023 Guardian [slides](#) and [video](#))

Invited Talks

- 09:15 Julian Hecker and Nan Laird, Fallacies and Pitfalls in Genome-Wide Association Studies ([JH slides](#), [NL slides](#), [JH+NL video](#))
- 10:15 Walter Scheirer, Photoshop Fantasies: Why is there so much fake stuff on the Internet? ([WS slides](#), [WS video](#))
- 11:15 Alicia Andrzejewski, Academic Ghosting: Towards an Academy of Truth-Telling ([AA slides](#), [AA video](#))

Technical Talks

- 12:30 Daniel Kristanto, Multiverse in Functional Magnetic Resonance Imaging Analysis ([DK slides](#), [DK video](#))
- 13:00 Koby Taswell, Consistent Bibliographic Data Formats with the BabbleNewt Project ([KT slides](#), [KT video](#))
- 13:30 Adam Craig, Managing Lexical-Semantic Hybrid Records of FAIR Metrics Analyses with the NPDS Cyberinfrastructure ([AC slides](#), [AC video](#))

Closing Remarks

- 14:00 Carl Taswell, Reproducibility, Reliability, and Integrity in Scholarly Research: What Accountability for Willful Disregard? ([CT slides](#), [CT video](#))

All slides and recordings of the talks are also available at www.BHAVI.us/BhaviHome/Symposia/202310.

* All talks at Guardians 2023 were presented virtually and are available online.

† Authors affiliated with Brain Health Alliance Virtual Institute, Ladera Ranch, CA 92694 USA; correspondence to [KTaswell at Brain Health Alliance](mailto:KTaswell@BrainHealthAlliance.org).

2023 Guardian: Anthony S. Fauci

2023 Guardian — BHA recognized Anthony S. Fauci, MD, as the BHA 2023 Guardian of Truth and Integrity. Throughout his career, Dr. Fauci has worked tirelessly to improve societal health through research in infectious diseases such as HIV, SARS, H1N1, and many more. During his time as the Director of the NIH National Institute of Allergy and Infectious Diseases, Dr. Fauci provided exemplary leadership in support of integrity in medical scientific research. His dedication to truth, honesty, and integrity proved crucial in the fight against the global COVID-19 pandemic during which Dr. Fauci and many other physician/scientists would have otherwise been swallowed by the sea of fake information that made it more difficult to save lives. For his lifetime of service and dedication to saving lives through integrity in clinical research, BHA honored Dr. Fauci as the 2023 Guardian of Truth and Integrity.

Julian Hecker and Nan Laird

Hecker, Craig, et al. 2023 — This talk and the associated slides, video, and article review common ways researchers misinterpret genome-wide association study (GWAS) results and how to avoid these fallacies and pitfalls. The authors review several relevant statistical methods, but one of the most important defenses against drawing wrong conclusions requires maintaining the appropriate mindset: Finding statistically significant associations between genetic variants and a phenotype of interest is not the end, but the beginning of a scientific journey.

Valid statistical methods do serve a purpose: Because a GWAS can consist of over a million statistical tests of association between the trait of interest and individual variants of all the genes in the genome, the likelihood is high that some tests will produce p-values less than 0.05 by random chance (**DerSimonian and Laird 2015**). Analysts need to address this bias by using statistical methods that correct the p-values with a stricter threshold of significance, such as the Bonferroni false discovery rate correction (**Tam et al. 2019**) or the versatile gene-based association study (VEGAS) methodology (**Hecker, Maaser, et al. 2017**). Even these corrections do not guarantee that subsequent studies will be able to replicate the results. The winner's curse is a statistical effect that leads to overestimation of the effect sizes of genetic variants that passed the significance threshold in GWAS (**Zhong and Prentice 2010**), which can lead researchers to underestimate the sample sizes they need for the next study to have the desired power. Several methods, including bootstrap resampling and empirical Bayesian estimation, can provide corrected estimates of effect sizes (**Forde et al. 2023**).

However strong the statistical association between the genetic variant and the trait variant, it does not tell us that the genotype causes the phenotype. As genes pass from one generation to the next, recombination, mutation, selection, and genetic drift act on them in complex ways, leading some pairs of genetic variants at different places in the genome to co-occur more or less often than one would expect, a situation known as linkage disequilibrium (LD, **Slatkin 2008**). Simple proximity of two genes to one another on the genome can lead to LD, creating statistical associations between genes without any causal relationship to the trait of interest simply because proximity causes them to be co-inherited more often (**Lappalainen and MacArthur 2021**). To account for this disconnect between correlation and causation, researchers use a category of methods called fine mapping to incorporate knowledge about the structure of the genome and rates of co-inheritance to identify which GWAS hit in a neighborhood with multiple such hits is most likely to be a causal variant (**Schaid et al. 2018**).

In addition to spatial proximity on the genome, natural selection and genetic drift can cause two genes to co-occur at higher or lower than overall average rates in a sector of the population with shared ancestry. This LD due to population stratification can lead the GWAS results to tag both genes as significantly associated with a trait that occurs at higher or lower rates in this same sector of the population, even when only one gene has any effect on it (**Derks et al. 2022**). While it is possible to correct for this association to some extent by including principal components of genetic ancestry as statistical covariates in the GWAS analysis (**Price et al. 2006**), it may be more effective to design the study to avoid introducing the effect in the first place. One way is with a family-based study design instead of a study that tests for associations throughout the general population (**Derks et al. 2022**). Examples include the transmission disequilibrium test (TDT, **Schaid 1998**) and family-based association test (FBAT, **Abecasis et al. 2000**).

The need for statistical corrections and additional analyses described above can make achieving the needed statistical power challenging. As an alternative to conducting a single sufficiently large study, researchers can use meta-analysis to combine results from multiple studies (**Mikolajewicz and Komarova 2019; Steel et al. 2021; Abdellaoui et al. 2023**). However, the predominance of study participants of European ancestry in past studies can bias results and limit the ability to generalize results to the rest of the population (**Derks et al. 2022**). Furthermore, the analysts need to check the metadata of the studies to ensure that the study designs are similar enough to allow comparison (**Mikolajewicz and Komarova 2019; Steel et al. 2021**).

Following the steps described above can help researchers to better identify a statistically promising association between a genetic variant and a phenotype, but they still cannot prove a causal association or reveal the mechanisms of cause and effect. GWAS hits often occur in non-coding regions of the genome with obscure regulatory functions (**Abdellaoui et al. 2023; Aguet et al. 2023**). Catalogues of known functional elements, such as the Encyclopedia of DNA Elements (ENCODE) (**Moore et al. 2020**) and GENCODE (**Frankish et al. 2020**), can help researchers leverage existing knowledge from other experimental methods (**Kichaev et al. 2019**). One important class of methods is quantitative trait locus (QTL) analyses, which involve searching for associations between locations on the genome, the QTLs, and various measurable features. These include associations between molecular QTLs and molecular phenotypes, including DNA methylation and production of specific metabolites (**Lappalainen and MacArthur 2021; Aguet et al. 2023**). They also include expression QTLs, which associate with downstream differences in the expression levels of other genes, as well as loci where both kinds of effects colocalize (**Rheenen et al. 2021**). These methods can extend to comprehensive post-GWAS analyses, particularly transcriptome-wide association studies (TWAS) and proteome-wide association studies (PWAS) (**Gedik et al. 2023**). Combining these locus-focused methods with information indicating similarity of annotations from single-cell gene expression, protein-protein interaction, and pathway participation features can lead to even more accurate identification of causal variants (**Weeks et al. 2023**). In this way, statistical data analysis and biological activity functional testing achieve a kind of synergy wherein statistical methods like GWAS identify candidates for study through lower-throughput laboratory experiments, which in turn provide knowledge of mechanisms of interaction that advanced statistical methods can use to more effectively find and prioritize subsequent candidate variants (**Gallagher and Chen-Plotkin 2018**).

Walter Scheirer

Photoshop Fantasies — In the past decade, the world has witnessed an increasing trend in fake posts, news, and online information. Recent estimates approximate that less than 60% of all web traffic is human with a majority of social media accounts operated by bots driven by automated algorithms (Read 2018). Elections and politics have been adversely impacted by fake information touted as fact with debates over what is fact or fiction becoming increasingly prevalent. When everything from deep fakes to memes recasts real life into fictional fantasies, how much can we trust the information we consume online?

Though fake information is rampant these days, modified images have been central to the history of propaganda, art, and entertainment. Authoritarian regimes such as the Chinese Communist Party under Mao Zedong used photo editing to rewrite history (Jaubert 1989). One significant example of this was removing the Gang of Four from images of Mao Zedong's funeral in an attempt to erase them from history. Although actions to rewrite history are often malicious at worst, misguided at best, picture editing in and of itself is not a purely evil act.

Since the first edited pictures in the 1840s, various methods have been employed to change the original image to something different, such as removing figures from a picture, cropping the edges to remove context, face swapping with cutouts, adding props, and more. Early photographers and artists used these techniques to improve the subject's appearance or add a sense of whimsy to the image by including fantastical figures. Some edited photography was obviously intended as a joke, such as an image of farmers cutting corn the size of logs. However, this intentional humor has not been the case for all edited images. Disputes only arise when an obviously impossible image is portrayed and contextualized deceptively as capturing and conveying truth, rather than recognizing the edited image as just that, ie, an image that has been altered no longer representing the truth.

Moving forward into the future of the digital age, photo editing techniques naturally lent themselves to digital picture editing with the creation of various image filters that would later become the foundations of software programs such as Adobe Photoshop. This change paved the way for cleaner removal or duplication of objects within an image as well as face swaps and other digital effects. Building from computer based signal processing for picture editing, AI became the next major step, enabling an interested layperson to engage in the art, without themselves having practiced the skills needed for photo-editing.

Today, with the internet as the 'frontier of the imagination', photoshop battles, AI-generated art, and edited memes have become the norm. By understanding the history, as a community, we can better learn to handle debates over truth and reality in the present and prepare for new futures as generated images become more widespread than in the past. For further discussion of this topic, see the book entitled "A History of Fake Things on the Internet" by Walter Scheirer 2023.

Alicia Andrzejewski

Academic Ghosting — Despite the respected image and prestige of those persons participating in academia, there exists a much darker underbelly to the institution. Universities may refuse or otherwise fail to protect their faculty professors and teachers from harassment by students and vice versa. Discussion of neurodivergence and/or mental illness has been heavily stigmatized, leaving faculty members with behavioral health challenges without support by their colleagues, mentors, and supporting staff. Ghosting whether by the institution during

the hiring process and/or by colleagues and mentors while at work has become endemic to the academic system.

Ghosting can be defined as the act of disappearing from someone's life without a word. A pattern once rampant within online dating, it has now reared its ugly head within academia. During hiring, many aspiring faculty send in applications, yet never hear back from the hiring committee, the HR department, or administrative staff. Although being ghosted during hiring can be upsetting and leave highly talented individuals in the dark, some members of HR departments claim that they do so to not sour the individual on the institution should there be a later attempt to hire that individual. This perspective seems illogical since an interested academic on the job market could also be soured on the institution by not hearing back with an update in the first place.

Ghosting can be even more devastating when the ghoster is a mentor or colleague of the ghostee. Outside of a student's own skills, successful completion of Ph.D. degrees are dependent on the faculty mentor's timely participation in advising the student and helping organize the thesis review committee. Despite a mentor's crucial role, it is not uncommon for mentors to ghost a student, while still clearly being a part of the institution, often leading to great emotional, career, and financial damage to the student. Sadly, those who have successfully made it to the position of faculty professor are also not necessarily safe from being targeted because colleagues ghost others for seemingly no apparent reason. It can be completely insidious with a slow decay in communication over time, or more obvious with an abrupt shift from constant and friendly communication to absolutely nothing at all.

Ghosting within academia may represent a lack of motivation to face tough conversations which simply must happen for the health of both individuals and the organization as a whole. Without addressing these problems within academia, the community will only become more unstable and the problems will likely worsen. Unfortunately, clear methods to fix the concerns are not immediately apparent. For ghosting during hiring, some hiring committee leaders have taken it upon themselves to personally email each and every applicant. But this task may impose a great cost in time and cannot be applied similarly to some of the other ghosting problems within academia. However, the best way forward to start is at least to begin the conversation and to spread awareness of the problems caused by academic ghosting. For more information and discussion, see articles in the Chronicle of Higher Education (Andrzejewski 2022; Andrzejewski 2023a; Andrzejewski 2023b) as well as an episode of the Academic Life podcast hosted by Dr. Christina Gessler featuring Dr. Andrzejewski (Gessler and Andrzejewski 2023).

Daniel Kristanto

Kristanto et al. 2023 — When conducting research, each choice made about methodology can impact to the results. These selections range from the actual experimental methods to the data processing methods performed, including statistical analysis of the results. To optimize across these various methods, a researcher can perform 'multi-universe analysis', which considers the various branching paths of possible methodologies, ie, of different methodologic data processing pipelines.

This paradigm can be applied to functional magnetic resonance imaging (fMRI) assessment of human brain networks. To begin, a systematic literature review of 252 papers was conducted to determine the possible forking methodological paths. Some common methods used were structural pre-processing, functional pre-processing, noise removal, functional connectivity definition, and graph analysis. Then using active machine learning, a smaller set of optimal paths can be deduced. Both

study results and an online interactive web application where viewers can see the various possible pipelines are discussed.

Koby Taswell

[S. K. Taswell, Anand, et al. 2023](#) — Appropriate reference citation serves as the foundation for ensuring research integrity, but managing a large number of resources can become burdensome. To address this concern, numerous organizations and research groups have developed a variety of automated tools for reference citation management associated with metadata formats to store the bibliographic data. Once properly stored in bibliography data files, citations can be used for references in documents or in other operations such as automated citation analysis for plagiarism detection.

BibTeX and BibLaTeX are widely used reference citation formats, common to the mathematics, computer science, and engineering communities for use with TeX and LaTeX document typesetting. Despite their decades-long history and wide recognition within these communities, they remain error prone due to format inconsistencies combined with various issues of instability and difficulty when debugging large-scale bibliographies. The BabbleNewt Project aims to address these deficiencies by providing a new format that can be easily converted to and from past versions of BibTeX and BibLaTeX while supporting migration to a more robust, fast, simple, and consistent JSON-like interoperable format.

Adam Craig

[Craig, Athreya, et al. 2023b](#) — Citation metrics that rate a publication more highly based on how many other works cite it create a perverse incentive to avoid citing potential rivals ([S. K. Taswell, Triggler, et al. 2020](#)). The FAIR Metrics, with FAIR an acronym for Fair Attribution to Indexed Reports and Fair Acknowledgment of Information Records, as defined in [Craig, Ambati, Dutta, Mehrotra, et al. 2019](#), solves this problem directly by quantifying how fairly a publication cites previously published work, thus providing alternative metrics to incentivize fairness with citational justice ([C. Taswell 2022](#)). The 4 FAIR Metric counts measure the numbers of claims misquoted from or misattributed to prior work, quoted from prior work, presented as novel, or plagiarized from other sources. These counts are used to calculate the corresponding 4 FAIR Metric ratios which provide summary scores, each emphasizing a different aspect of citation practice. Unlike commonly used lexical plagiarism detection tools, the FAIR Metrics depend on entity equivalences between the concepts and ideas expressed in documents, not just lexical similarity between documents. Demonstration with a human analyst evaluating the FAIR Metrics on example texts provides a prototype workflow for use of the FAIR Metrics that enables human-performed peer review to be more objective and that serves as a standard for comparison of results from future automated algorithms.

This work extends the preliminary version of an analysis presented at eScience 2023 ([Craig, Athreya, et al. 2023a](#)), which described the successful application of a human-performed FAIR Metrics evaluation workflow to 5 reports and a brief description of methods for publishing semantic descriptions of the evaluation with the PDP-DREAM Ontology. Both the original workflow proposed in [Craig, Ambati, Dutta, Mehrotra, et al. 2019](#) and that of [Craig, Athreya, et al. 2023a](#) focus on comparison of claims between one test document and one comparison document. This extended version [Craig, Athreya, et al. 2023b](#) of [Craig, Athreya, et al. 2023a](#) elaborates on the structure of these RDF description records and

analyses 4 more example pairs. In a change from the earlier procedures, analysts were required to evaluate scores in comparison to all references cited by the test and comparison documents. This approach provides a more robust way to evaluate allegations of plagiarism via the creation of a RDF document clarifying which claims from the test document match claims that may or may not be referenced in other documents.

The authors selected 9 example test-comparison pairs for evaluation. One case was selected as a negative control representing a known pair of documents without plagiarism. Seven cases were selected as known plagiarism from the Retraction Watch database with differing forms and extents of plagiarism. The last case was selected as reported plagiarism based on the comparison documented in detail in [Craig, Ambati, Dutta, Kowshik, et al. 2019](#). For each test document, a designated comparison document was chosen for evaluation. In general, FAIR Metrics ratio scores for test-comparison pairs of known plagiarism were lower than the negative control case, with the most extreme instances of known plagiarism having the lowest scores.

In the case of reported but not-yet-retracted plagiarism, which compared [Wilkinson et al. 2016](#) to [C. Taswell 2007](#), the FAIR Metrics analysis confirms this reported plagiarism as paraphrasing plagiarism, classifying the FAIR Principles claims wrongly misrepresented as novel by [Wilkinson et al. 2016](#) instead as plagiarized from [C. Taswell 2007](#). Overall, the FAIR Metrics scores found for [Wilkinson et al. 2016](#) align with those of the extreme examples of plagiarism, thus confirming the plagiarism by [Wilkinson et al. 2016](#) of [C. Taswell 2007](#).

Carl Taswell

[C. Taswell 2023](#) — As each year passes in the current era of information wars, the importance of maintaining reproducibility, reliability, validity, and integrity in scholarly research only grows greater, but there are not yet enforceable safeguards that have been adopted. In the long term, as a community of researchers, we should consider licensing analogous to that required in the professions of medicine, law and education. In the short term, the current situation leaves many questions unanswered. What steps can be taken now to respond to complaints from victims of plagiarism, misconduct, and fraud? Which organizations will remain committed not only to talking about preventing plagiarism, misconduct, and fraud but also sanctioning these violations of professional conduct when they occur? How can we not only heal from but also cure and prevent the problems of grooming, gaslighting, and ghosting in academia? [Guardians 2024](#) will continue this conversation with the website open for submissions beginning on 9 January 2024.

Citation

Brainiacs 2023 Volume 4 Issue 2 Edoc Y331839FB

Title: "Who are the Guardians of Truth and Integrity?"

Authors: S. Koby Taswell and Adam Craig

Dates: created 2023-10-09, received 2023-10-09, updated 2023-12-23, published 2023-12-23, endorsed 2023-12-31

Copyright: © 2023 Brain Health Alliance

Contact: [KTaswell at Brain Health Alliance](mailto:KTaswell@BrainHealthAlliance.org)

URL: BrainiacsJournal.org/arc/pub/Taswell2023WAGTI

PDP: [/Nexus/Brainiacs/Taswell2023WAGTI](https://Nexus/Brainiacs/Taswell2023WAGTI)

DOI: [/10.48085/Y331839FB](https://doi.org/10.48085/Y331839FB)

References

- [1] A. Abdellaoui, L. Yengo, K. J. Verweij, and P. M. Visscher. "15 years of GWAS discovery: Realizing the promise." *The American Journal of Human Genetics* 110.2 (Feb. 2023), pp. 179–194. ISSN: 0002-9297. DOI: [10.1016/j.ajhg.2022.12.011](https://doi.org/10.1016/j.ajhg.2022.12.011) (cited p. 2).
- [2] G. R. Abecasis, W. O. Cookson, and L. R. Cardon. "Pedigree tests of transmission disequilibrium." *European Journal of Human Genetics* 8.7 (2000), pp. 545–551 (cited p. 2).
- [3] F. Aguet, K. Alasoo, Y. I. Li, A. Battle, H. K. Im, S. B. Montgomery, and T. Lappalainen. "Molecular quantitative trait loci." *Nature Reviews Methods Primers* 3.1 (Jan. 2023). ISSN: 2662-8449. DOI: [10.1038/s43586-022-00188-6](https://doi.org/10.1038/s43586-022-00188-6) (cited p. 2).
- [4] A. Andrzejewski. "When Students Harass Professors." *Chronicle of Higher Education* (Aug. 8, 2022). URL: <https://www.chronicle.com/article/when-students-harass-professors> (cited p. 3).
- [5] A. Andrzejewski. "Academics Don't Talk About Our Mental Illnesses. We Should." *Chronicle of Higher Education* (July 5, 2023). URL: <https://www.chronicle.com/article/academics-dont-talk-about-our-mental-illnesses-we-should> (cited p. 3).
- [6] A. Andrzejewski. "The Sad Humiliations of Academic Ghosting." *Chronicle of Higher Education* (Jan. 30, 2023). URL: <https://www.chronicle.com/article/the-sad-humiliations-of-academic-ghosting> (cited p. 3).
- [7] A. Craig, A. Ambati, S. Dutta, P. Kowshik, S. Nori, S. K. Taswell, Q. Wu, and C. Taswell. "DREAM Principles and FAIR Metrics from the PORTALDOORS Project for the Semantic Web." In: *2019 IEEE 11th International Conference on Electronics, Computers and Artificial Intelligence (ECAI)* (June 28, 2019). Pitesti, Romania: IEEE, June 2019, pp. 1–8. DOI: [10.1109/ECAI46879.2019.9042003](https://doi.org/10.1109/ECAI46879.2019.9042003). URL: <https://portaldoors.org/pub/docs/ECAI2019DREAMFAIRO612.pdf> (cited p. 4).
- [8] A. Craig, A. Ambati, S. Dutta, A. Mehrotra, S. K. Taswell, and C. Taswell. "Definitions, Formulas, and Simulated Examples for Plagiarism Detection with FAIR Metrics." In: *2019 ASIS&T 82nd Annual Meeting* (Oct. 19, 2019). Vol. 56. Melbourne, Australia: Wiley, 2019, pp. 51–57. DOI: [10.1002/PRA2.6](https://doi.org/10.1002/PRA2.6). URL: <https://portaldoors.org/pub/docs/ASIST2019FairMetrics0611.pdf> (cited p. 4).
- [9] A. Craig, A. Athreya, and C. Taswell. "Example evaluations of plagiarism cases using FAIR Metrics and the PDP-DREAM Ontology." *IEEE eScience 2023* (Oct. 13, 2023) (cited p. 4).
- [10] A. Craig, A. Athreya, and C. Taswell. "Managing Lexical-Semantic Hybrid Records of FAIR Metrics Analyses with the NPDS Cyberinfrastructure." *Brainiacs Journal of Brain Imaging And Computing Sciences* 4.2 (Dec. 27, 2023). DOI: [10.48085/d5b2734f2](https://doi.org/10.48085/d5b2734f2) (cited p. 4).
- [11] E. M. Derks, J. G. Thorp, and Z. F. Gerring. "Ten challenges for clinical translation in psychiatric genetics." *Nature Genetics* 54.10 (Sept. 2022), pp. 1457–1465. ISSN: 1546-1718. DOI: [10.1038/s41588-022-01174-0](https://doi.org/10.1038/s41588-022-01174-0) (cited p. 2).
- [12] R. DerSimonian and N. Laird. "Meta-analysis in clinical trials revisited." *Contemporary Clinical Trials* 45 (Nov. 2015), pp. 139–145. ISSN: 1551-7144. DOI: [10.1016/j.cct.2015.09.002](https://doi.org/10.1016/j.cct.2015.09.002) (cited p. 2).
- [13] A. Forde, G. Hemani, and J. Ferguson. "Review and further developments in statistical corrections for Winner's Curse in genetic association studies." *PLoS Genetics* 19.9 (2023), e1010546 (cited p. 2).
- [14] A. Frankish, M. Diekhans, I. Jungreis, J. Lagarde, et al. "GENCODE 2021." *Nucleic Acids Research* 49.D1 (Dec. 2020), pp. D916–D923. ISSN: 1362-4962. DOI: [10.1093/nar/gkaa1087](https://doi.org/10.1093/nar/gkaa1087) (cited p. 2).
- [15] M. D. Gallagher and A. S. Chen-Plotkin. "The Post-GWAS Era: From Association to Function." *The American Journal of Human Genetics* 102.5 (May 2018), pp. 717–730. ISSN: 0002-9297. DOI: [10.1016/j.ajhg.2018.04.002](https://doi.org/10.1016/j.ajhg.2018.04.002) (cited p. 2).
- [16] H. Gedik, R. E. Peterson, B. P. Riley, V. I. Vladimirov, and S.-A. Bacanu. "Integrative Post-Genome-Wide Association Study Analyses Relevant to Psychiatric Disorders: Imputing Transcriptome and Proteome Signals." *Complex Psychiatry* 9.1-4 (2023), pp. 130–144 (cited p. 2).
- [17] C. Gessler and A. Andrzejewski. *Academic Ghosting*. June 8, 2023. URL: <https://newbooksnetwork.com/academic-ghosting> (cited p. 3).
- [18] J. Hecker, A. Craig, A. Hughes, J. Neidich, C. Taswell, and N. Laird. "Fallacies and Pitfalls in Genome-Wide Association Studies." *Brainiacs Journal of Brain Imaging And Computing Sciences* 4.2 (Dec. 21, 2023). DOI: [10.48085/gfa4e8812](https://doi.org/10.48085/gfa4e8812) (cited p. 2).
- [19] J. Hecker, A. Maaser, D. Prokopenko, H. L. Fier, and C. Lange. "Reporting Correct p Values in VEGAS Analyses." *Twin Research and Human Genetics* 20.3 (Mar. 2017), pp. 257–259. ISSN: 1839-2628. DOI: [10.1017/tbg.2017.16](https://doi.org/10.1017/tbg.2017.16) (cited p. 2).
- [20] A. Jaubert. *Making people disappear. An amazing chronicle of photographic deception*. Intelligence & National Security Library. Bibliography: p. 187-190. Pergamon-Brassey's International Defense Publishers, 1989. 190 pp. ISBN: 0080374301 (cited p. 3).
- [21] G. Kichaev, G. Bhatia, P.-R. Loh, S. Gazal, et al. "Leveraging Polygenic Functional Enrichment to Improve GWAS Power." *The American Journal of Human Genetics* 104.1 (Jan. 2019), pp. 65–75. ISSN: 0002-9297. DOI: [10.1016/j.ajhg.2018.11.008](https://doi.org/10.1016/j.ajhg.2018.11.008) (cited p. 2).
- [22] D. Kristanto, C. Gießing, M. Marek, C. Zhou, S. Debener, C. Thiel, and A. Hildebrandt. "An Extended Active Learning Approach to Multiverse Analysis: Predictions of Latent Variables from Graph Theory Measures of the Human Connectome and Their Direct Replication." *Brainiacs Journal of Brain Imaging And Computing Sciences* 4.2 (Dec. 21, 2023). DOI: [10.48085/j962e0f53](https://doi.org/10.48085/j962e0f53) (cited p. 3).
- [23] T. Lappalainen and D. G. MacArthur. "From variant to function in human disease genetics." *Science* 373.6562 (Sept. 2021), pp. 1464–1468. ISSN: 1095-9203. DOI: [10.1126/science.abi8207](https://doi.org/10.1126/science.abi8207) (cited p. 2).
- [24] N. Mikolajewicz and S. V. Komarova. "Meta-Analytic Methodology for Basic Research: A Practical Guide." *Frontiers in Physiology* 10 (Mar. 2019). ISSN: 1664-042X. DOI: [10.3389/fphys.2019.00203](https://doi.org/10.3389/fphys.2019.00203) (cited p. 2).
- [25] J. E. Moore, M. J. Purcaro, H. E. Pratt, C. B. Epstein, et al. "Expanded encyclopaedias of DNA elements in the human and mouse genomes." *Nature* 583.7818 (2020), pp. 699–710 (cited p. 2).
- [26] A. L. Price, N. J. Patterson, R. M. Plenge, M. E. Weinblatt, N. A. Shadick, and D. Reich. "Principal components analysis corrects for stratification in genome-wide association studies." *Nature genetics* 38.8 (2006), pp. 904–909 (cited p. 2).
- [27] M. Read. "How Much of the Internet Is Fake? Turns Out, a Lot of It, Actually." *New York Intelligencer* (Dec. 26, 2018). URL: <https://nymag.com/intelligencer/2018/12/how-much-of-the-internet-is-fake.html> (cited p. 3).
- [28] W. van Rheenen, R. A. A. van der Spek, M. K. Bakker, J. J. F. A. van Vugt, et al. "Common and rare variant association analyses in amyotrophic lateral sclerosis identify 15 risk loci with distinct genetic architectures and neuron-specific biology." *Nature Genetics* 53.12 (Dec. 2021), pp. 1636–1648. ISSN: 1546-1718. DOI: [10.1038/s41588-021-00973-1](https://doi.org/10.1038/s41588-021-00973-1) (cited p. 2).
- [29] D. J. Schaid. "Transmission disequilibrium, family controls, and great expectations." *The American Journal of Human Genetics* 63.4 (1998), pp. 935–941 (cited p. 2).

- [30] D. J. Schaid, W. Chen, and N. B. Larson. "From genome-wide associations to candidate causal variants by statistical fine-mapping." *Nature Reviews Genetics* 19.8 (May 2018), pp. 491–504. ISSN: 1471-0064. DOI: [10.1038/s41576-018-0016-z](https://doi.org/10.1038/s41576-018-0016-z) (cited p. 2).
- [31] W. J. Scheirer. *A history of fake things on the internet*. Includes bibliographical references (pages 193–228) and index. Stanford, California: Stanford University Press, Dec. 5, 2023. 241 pp. ISBN: 978-1503632882 (cited p. 3).
- [32] M. Slatkin. "Linkage disequilibrium – understanding the evolutionary past and mapping the medical future." *Nature Reviews Genetics* 9.6 (June 2008), pp. 477–485. ISSN: 1471-0064. DOI: [10.1038/nrg2361](https://doi.org/10.1038/nrg2361) (cited p. 2).
- [33] P. Steel, S. Beugelsdijk, and H. Aguinis. "The anatomy of an award-winning meta-analysis: Recommendations for authors, reviewers, and readers of meta-analytic reviews." *Journal of International Business Studies* 52.1 (Jan. 2021), pp. 23–44. ISSN: 1478-6990. DOI: [10.1057/s41267-020-00385-z](https://doi.org/10.1057/s41267-020-00385-z) (cited p. 2).
- [34] V. Tam, N. Patel, M. Turcotte, Y. Bossé, G. Paré, and D. Meyre. "Benefits and limitations of genome-wide association studies." *Nature Reviews Genetics* 20.8 (May 2019), pp. 467–484. ISSN: 1471-0064. DOI: [10.1038/s41576-019-0127-1](https://doi.org/10.1038/s41576-019-0127-1) (cited p. 2).
- [35] C. Taswell. "DOORS to the Semantic Web and Grid with a PORTAL for Biomedical Computing." *IEEE Transactions on Information Technology in Biomedicine* 12.2 (2 Mar. 2007). In the Special Section on Bio-Grid published online 3 Aug. 2007, pp. 191–204. ISSN: 1089-7771. DOI: [10.1109/TITB.2007.905861](https://doi.org/10.1109/TITB.2007.905861) (cited p. 4).
- [36] C. Taswell. "Epistemic Injustice, Open Access, and Citational Justice." *Brainiacs Journal of Brain Imaging And Computing Sciences* 3.2 (Dec. 30, 2022). ISSN: 2766-6883. DOI: [10.48085/x3b678b7a](https://doi.org/10.48085/x3b678b7a) (cited p. 4).
- [37] C. Taswell. "Reproducibility, Validity, and Integrity in Scholarly Research: What Accountability for Willful Disregard?" *Brainiacs Journal of Brain Imaging And Computing Sciences* 4.2 (Dec. 31, 2023). DOI: [10.48085/13570f30f](https://doi.org/10.48085/13570f30f) (cited p. 4).
- [38] S. K. Taswell, A. Anand, M. Montes-Soza, and C. Taswell. "BabbleNewt: A Simplified, Consistent, and Interoperable Citation Format for Bibliographic Metadata." *Brainiacs Journal of Brain Imaging And Computing Sciences* 4.2 (Dec. 18, 2023). DOI: [10.48085/k562cb81c](https://doi.org/10.48085/k562cb81c) (cited p. 4).
- [39] S. K. Taswell, C. Triggler, J. Vayo, S. Dutta, and C. Taswell. "The Hitchhiker's Guide to Scholarly Research Integrity." In: *2020 ASIS&T 83rd Annual Meeting* (Oct. 22, 2020). Vol. 57. Wiley, 2020, e223. DOI: [10.1002/ppra2.223](https://doi.org/10.1002/ppra2.223). URL: <https://asistdl.onlinelibrary.wiley.com/doi/abs/10.1002/ppra2.223> (cited p. 4).
- [40] E. M. Weeks, J. C. Ulirsch, N. Y. Cheng, B. L. Trippe, et al. "Leveraging polygenic enrichments of gene features to predict genes underlying complex traits and diseases." *Nature Genetics* 55.8 (July 2023), pp. 1267–1276. ISSN: 1546-1718. DOI: [10.1038/s41588-023-01443-6](https://doi.org/10.1038/s41588-023-01443-6) (cited p. 2).
- [41] M. D. Wilkinson, M. Dumontier, I. J. Aalbersberg, G. Appleton, et al. "The FAIR Guiding Principles for scientific data management and stewardship." *Scientific data* 3.1 (2016), pp. 1–9 (cited p. 4).
- [42] H. Zhong and R. L. Prentice. "Correcting "winner's curse" in odds ratios from genomewide association findings for major complex human diseases." *Genetic Epidemiology: The Official Publication of the International Genetic Epidemiology Society* 34.1 (2010), pp. 78–91 (cited p. 2).