

RESEARCH NOTE

Open Access



Power calculator for detecting allelic imbalance using hierarchical Bayesian model

Katrina Sherbina¹, Luis G. León-Novelo² , Sergey V. Nuzhdin³, Lauren M. McIntyre⁴ and Fabio Marroni^{5*}

Abstract

Objective: Allelic imbalance (AI) is the differential expression of the two alleles in a diploid. AI can vary between tissues, treatments, and environments. Methods for testing AI exist, but methods are needed to estimate type I error and power for detecting AI and difference of AI between conditions. As the costs of the technology plummet, what is more important: reads or replicates?

Results: We find that a minimum of 2400, 480, and 240 allele specific reads divided equally among 12, 5, and 3 replicates is needed to detect a 10, 20, and 30%, respectively, deviation from allelic balance in a condition with power > 80%. A minimum of 960 and 240 allele specific reads divided equally among 8 replicates is needed to detect a 20 or 30% difference in AI between conditions with comparable power. Higher numbers of replicates increase power more than adding coverage without affecting type I error. We provide a Python package that enables simulation of AI scenarios and enables individuals to estimate type I error and power in detecting AI and differences in AI between conditions.

Keywords: Allelic imbalance, Type I error, Power, Simulation, Allele specific reads, Biological replicates

Introduction

Gene expression in a diploid individual is the result of the combined expression of both alleles. Allele Specific Expression (ASE) is the amount of mRNA transcribed at each allele. The two alleles of a diploid individual can show significantly different expression, a condition termed allelic imbalance (AI) [1]. AI is a result of genetic variation in regulation in *cis* (e.g. promoters, enhancers, and other noncoding sequences), *trans* (e.g. transcription factors) or resulting from *cis* by *trans* interactions [1–7]. AI has been observed as a consequence of imprinting [8–10] and nonsense mediated decay [11] and has been shown to contribute to heterosis [12] and hybrid incompatibility [13]. The extent of AI in human tissues can give information on the impact of heterozygous mutations

on the expression of the mutated allele in healthy [14] or cancerous human tissues [15]. Also, loss of heterozygosity can be detected using AI [16, 17].

Several methods have been proposed for the detection of AI, [5, 15, 18–21]. However, there is currently only one model developed to formally test for difference in AI across conditions [22]. Comparing AI between conditions or tissues can provide new insights into the mechanisms of gene expression regulation [6, 9, 22–27]. Most often, these comparisons are heuristic without a formal statistical test. However, statistical comparisons have been made of heterogeneity in AI between mated and virgin *Drosophila* female head tissue [22], human tissues types within an individual [11, 26, 28], and cell subpopulations in different developmental stages [29]. Some statistical tests have been performed to assess whether *cis* effects differ among alleles in a population [7] or in parent of origin effects in mice [5].

Type I error in AI studies has been well explored and is known to be high, particularly when failing to account

*Correspondence: fabio.marroni@uniud.it

⁵ Dipartimento di Scienze Agroalimentari, Ambientali e Animali, Università di Udine, 33100 Udine, Italy

Full list of author information is available at the end of the article



for map bias [30], and/or using the binomial test [5, 19, 20, 31–33]. What is currently absent from the literature is an understanding of the power for studies of AI and, in particular, what the best allocation of resources is for boosting power for detection of AI when the hypothesis of interest is a comparison of AI *between* conditions. What is more important: more reads or more replicates? Is there a minimum number of replicates needed? A minimum number of allele specific reads? As sequencing costs are dropping in price and the per sample cost of library preparation dramatically lower than a decade ago, it is time to stop relying on the magic number 3 and determine the necessary size and scope of such studies to control type I for a particular type II error/power. It is common practice to assess power before embarking on association studies [34, 35], but no tool is currently available for assessing power for detecting AI and differences in AI between conditions.

To address this need, we present here the package BayesASE_power. It consists of tools to enable the user to simulate RNA-seq read counts under a previously published Bayesian model of AI [7, 20] with any number of replicates, reads, and AI. The results are aggregated across multiple simulated datasets to estimate Type I error and power. We demonstrate how to use BayesASE_power to plan experiments to achieve the desired power in detecting AI within a condition and/or interactions of AI between conditions.

Main text

Methods

The model used for detection of AI in any condition and for comparing levels of AI between any two conditions has been described earlier and implemented in the package BayesASE [7, 22]. We give here the basic definitions and refer the reader to Additional file 1 for further details.

One important parameter in determining AI is the probability r of a read aligning to allele $g1$ ($g2$) given that it came from that allele, that we define as $r_{i,g1}$ ($r_{i,g2}$). Low values of these probabilities correspond to a high degree of ambiguously mapped reads, which occurs when there is little sequence divergence between the two alleles. Reads that do not map ambiguously are termed allele specific reads or informative reads.

AI in condition i is measured by the parameter θ_i representing the proportion of reads originating from the allele $g1$, which that can be written as follows:

$$\theta_i = \frac{\mathbb{E}(x_{i,k}/r_{i,g1})}{\mathbb{E}(x_{i,k}/r_{i,g1} + y_{i,k}/r_{i,g2})} = \frac{1/\alpha_i}{\alpha_i + 1/\alpha_i}$$

When θ_i is close to 0, we have one extreme case of AI with almost all the reads originating from $g2$. When $\theta_i = 0.5$, we have perfect allelic balance with 50% of the reads from each allele. With $\theta_i = 1$, we are in the opposite direction of extreme AI with all the reads originating from $g1$.

The following null hypotheses are defined:

1. Allelic balance in condition 1, i.e. null H1: $\theta_1 = 0.5$ or equivalently $\alpha_1 = 1$.
2. Allelic balance in condition 2, i.e. null H2: $\theta_2 = 0.5$ or equivalently $\alpha_2 = 1$.
3. Level of AI is the same in both conditions, i.e. null H3: $\theta_1 = \theta_2$ or equivalently $\alpha_1 = \alpha_2$.

To test these hypotheses, three cases are defined (Fig. 1):

1. H1, H2 and H3 are satisfied
2. H1 is satisfied, H2 and H3 are violated
3. H1 and H2 are violated, H3 is satisfied

In our simulation, magnitudes of deviation from the null are reported as ΔAI . Given $\theta_0 = 0.5, \Delta AI_1 = \frac{|\theta_1 - \theta_0|}{\theta_0}$ for H1, $\Delta AI_2 = \frac{|\theta_2 - \theta_0|}{\theta_0}$ for H2, and $\Delta AI_3 = \frac{|\theta_2 - \theta_1|}{\theta_1}$ for H3. Simulated deviations of ΔAI from the null are moderate, generally between 0.1 and 0.3, with a maximum of 0.5.

Scenarios that vary the number of allele specific reads, number of replicates, and AI in the different cases were simulated (Additional file 2).

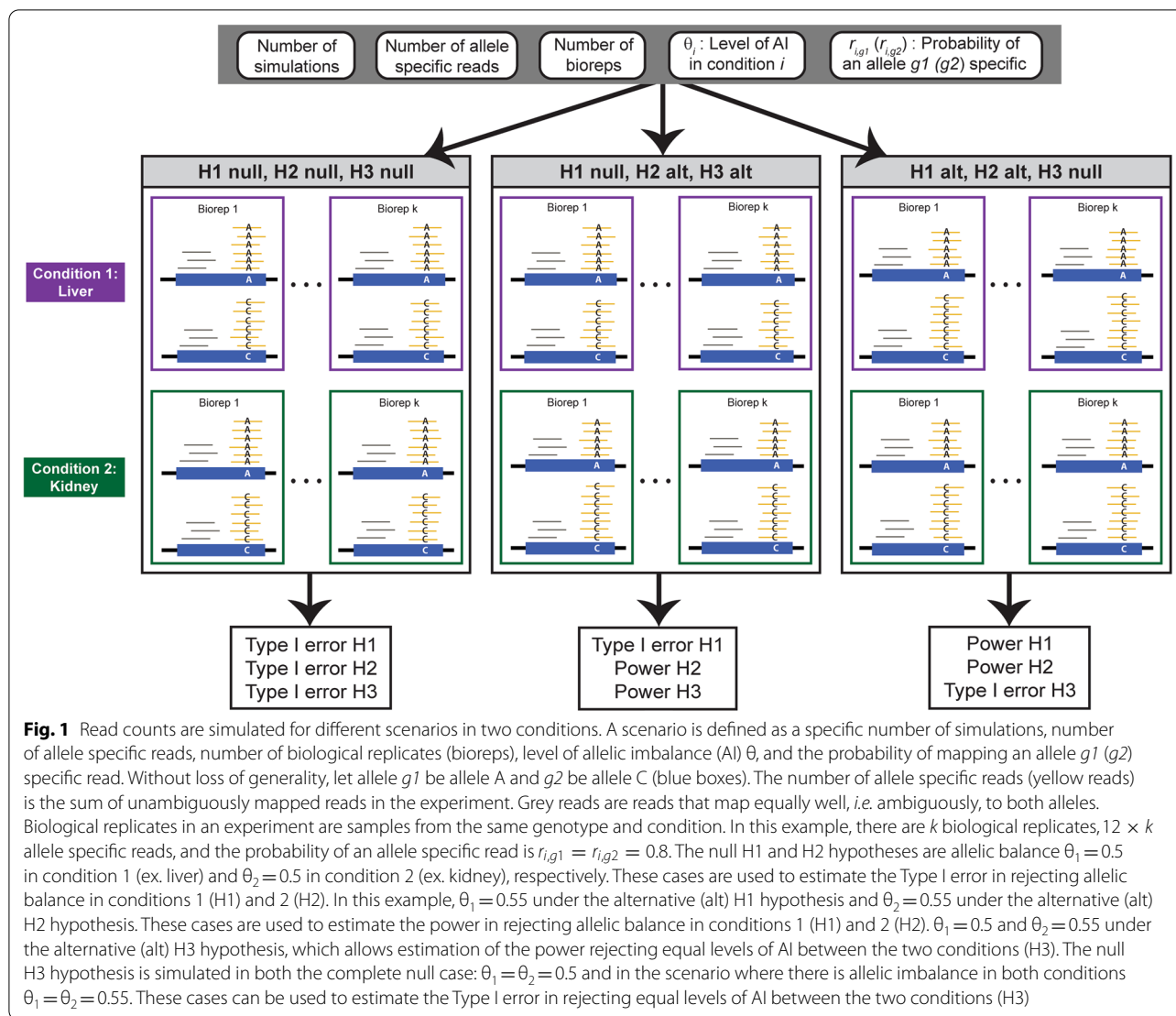
The model used for the detection of AI in any condition and for comparing levels of AI between any two conditions has been described earlier and implemented in the package BayesASE (<https://github.com/McIntyre-Lab/BayesASE>) [7, 22].

Results

Type I error is controlled except when total allele specific reads exceed 2400 allele specific reads dispersed across 8 or more biological replicates (Fig. 2a, b). However, type I error is always less than 0.08.

Under all conditions Type I error is low (Fig. 2, Additional file 3), and only exceeds the nominal value of 5% in scenarios with very high numbers of allele specific reads ($n > 2400$) and biological replicates.

Power (Fig. 3) for detecting small deviations from the null (0.1) is less than 0.4 when the number of bioreps is 3 and only exceeds 0.6 when the number of bioreps is greater than 6 and the total number of allele specific reads is large. H1 is rejected with power $> 80\%$ when the total number of reads is at least 2400, and the number of independent biological replicates is at least 12 (an average of 200 allele specific reads per biological replicate). Power

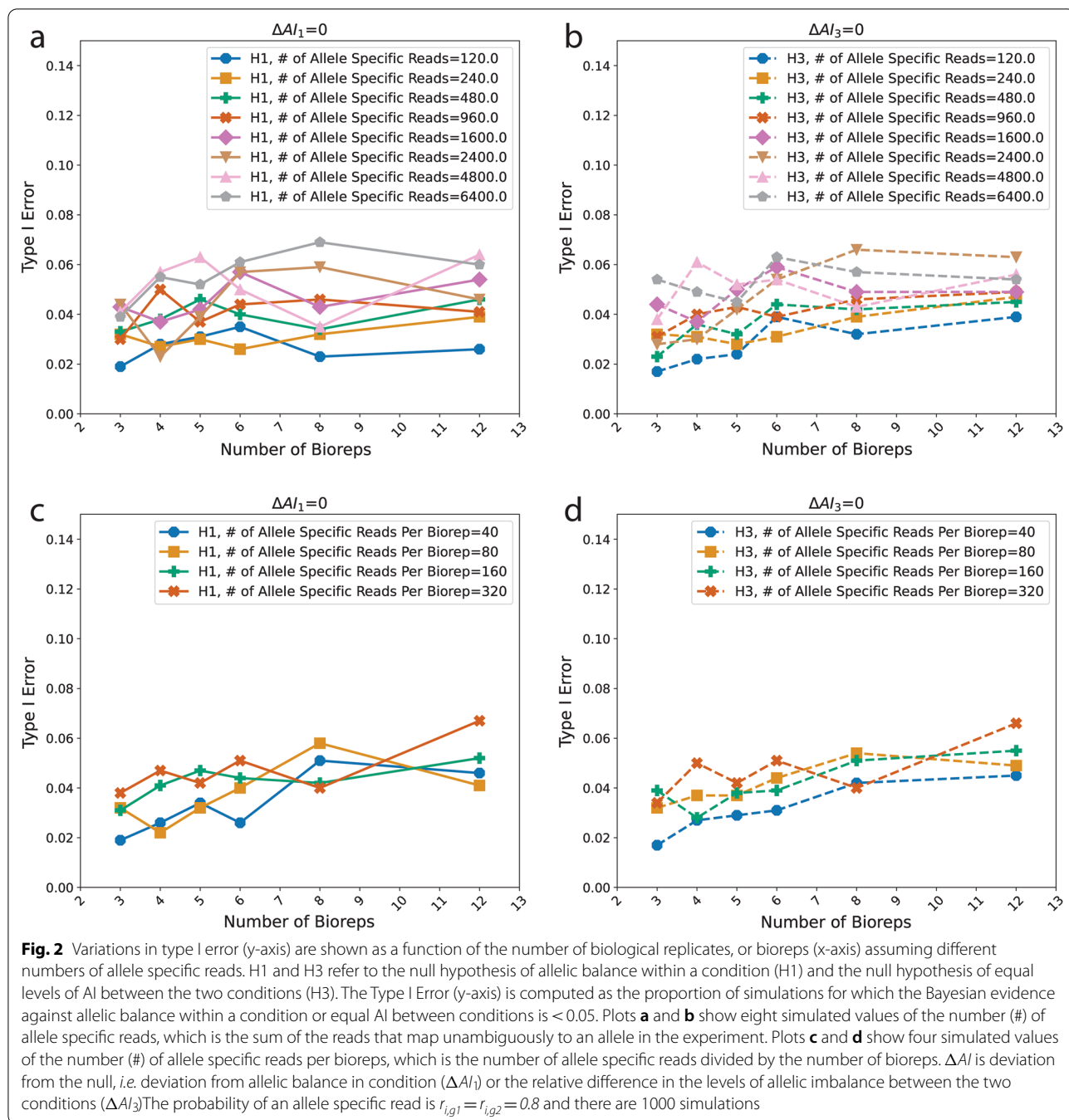


for rejecting H3 is, as expected, lower than H1 (Fig. 3, Additional file 4). For $\Delta AI = 0.2$ (central panels) and 960 informative reads, H1 is rejected with power > 80% with 3 biological replicates (average of 320 reads per replicate). H3 is rejected with power > 80% with 960 informative reads in 8 replicates (average of 120 reads per replicate). When $\Delta AI = 0.3$ (bottom panels), power approaches 100% except when the number of informative reads is low (120). As expected, the number of simulations does not affect estimates of power (Additional file 5). Power for the test of H3 for 3 biological replicates is maximal at 640 informative reads (Additional file 6). When $\Delta AI = 0.3$, most scenarios have power greater than 80% for both H1 and H3 (Fig. 3, Additional file 4). When ΔAI is 0.5, power for both H1 and H3 is ~ 100% even when the total number of informative reads is low. This represents an

extreme scenario, but one that is often observed in situations with loss of heterozygosity, indicating that in these scenarios relatively few reads are needed to detect AI with confidence (Additional file 7).

Discussion

Type I error rarely exceeds the nominal value of 5% even for very high numbers of allele specific reads, while increasing the number of allele specific reads substantially increases power. These observations are in agreement with other approaches [5, 15, 21, 36], including with simulations performed using a previous version of this model [22]. However, except for Zou et al. [5], the power of these approaches was not assessed for jointly changing the number of allele specific reads and biological replicates. BayesASE can directly test for a difference in AI



(See figure on next page.)

Fig. 3 H1 refers to simulations under the alternative hypothesis of allelic imbalance within a condition and H3 refers to unequal levels of AI between the two conditions. For H1, the x-axis is the effect size, which is the relative deviation from allelic balance $\Delta AI_1 = \frac{|\theta - \theta_0|}{\theta_0}$, where $\theta_0 = 0.5$. For H3, the x-axis is the relative difference in levels of AI between the two conditions $\Delta AI_3 = \frac{|\theta_2 - \theta_1|}{\theta_1}$ where the first condition is simulated under the null hypothesis and the second under the alternative hypothesis $\theta \neq 0.5$. The power (y-axis) is computed as the proportion of simulations for which the Bayesian evidence against allelic balance within a condition or against equal levels of AI between conditions is < 0.05. There are 1000 simulations and the probability of an allele specific read is $r_{i,g1} = r_{i,g2} = 0.8$. Simulations for 3, 4, 5, 6, 8, and 12 biological replicates (bioreps, x-axis) for varying numbers (#) of allele specific reads are reported

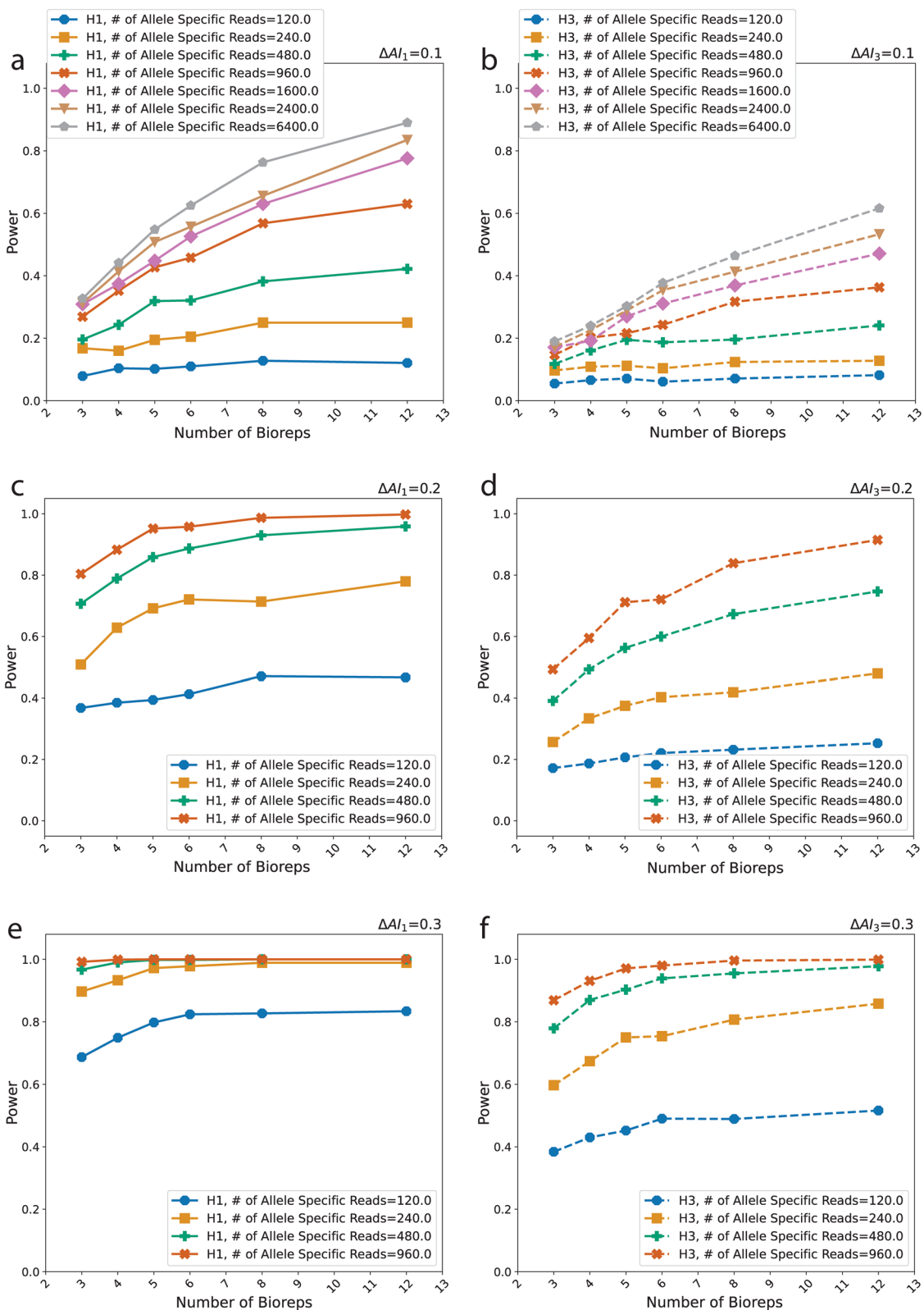


Fig. 3 (See legend on previous page.)

between two conditions or genotypes and, accordingly, we can assess how variation in both the number of replicates and reads affects the power to not only detect AI in a condition but differences in AI between conditions.

BayesASE has adequate power to detect moderate deviations from the null hypothesis. However, the minimum number of reads and biological replicates to achieve this power is greater for smaller deviations from the null. Our simulations suggest that a minimum of 2400 informative reads across 12 replicates, 480 informative reads across 5 replicates (or a minimum of 3 replicates with a total of 960 informative reads), and 240 informative reads across 3 replicates results in >80% power to detect ΔAI_1 (or ΔAI_2) = 0.1, 0.2, and 0.3, respectively. While the power to detect $\Delta AI_3 = 0.1$ does not surpass 60% in our simulations, we can detect a difference in AI between conditions (ΔAI_3) of 0.2 and 0.3 with comparable power for the same deviation from the null within a condition with the same number of informative reads but only when spread over more replicates (*i.e.* 8). A deviation from the null of $\Delta AI = 0.3$ has power >80% in most scenarios and even higher deviations can be detected with almost 100% power. Such large differences are indicative of loss of heterozygosity as observed in cancers [17] and imprinting [9, 10].

The results presented here describe general trends. In order to estimate power (and type I error) for a particular scenario of interest, the simulator developed as a part of this work can be used. The simulator explicitly enables the exploration of the total number of reads relative to the number of informative reads. The number of informative reads depends on the length of the feature (exon, gene), and on the density of polymorphisms. While, it is possible to analyze individual SNPs, investigators should take care to ensure that individual reads are not used in support of multiple SNPs. In addition, both the number allele specific reads (dependent on the distribution of polymorphisms) and the overall number of reads (dependent on library size and expression levels) affect power and should be accounted for in any modeling approach comparing AI between conditions.

One aspect that will be interesting to test in future studies is the behavior of nearby genes. In organisms, such as *D. melanogaster*, in which topologically associated domains (TADs) aggregate genes with similar expression patterns [37], we could expect that TADs also discriminate different patterns of AI, if the imbalance is due to shared sequence (*i.e.* polymorphic enhancer), but not if the imbalance is due to gene-private sequence (*i.e.* polymorphic gene or promoter).

We present results of an extensive simulation study to quantify type I error and power in detecting AI using the model implemented in the BayesASE pipeline [7]. Both

number of reads and number of replicates are important, and they both should be maximized. However, for any given number of reads, the best idea is to maximize the number of replicates. This is in agreement with previous studies that suggested that increased biological replication should be favored over increased depth of coverage [5, 38]. This of course should be balanced against the fact that having several replicates is more expensive. This said, we do not recommend designing any biological experiment with less than three biological replicates.

Limitations

This simulation study, like most such studies, makes simplifying assumptions for computational ease and efficiency. It is performed under optimal scenarios for a single gene and, thus, may not account for all limitations that are inherent to real data. Thus, the recommendations based on the simulation results should be considered a minimum threshold for study size planning. However, despite their drawbacks, simulations are necessary because it is not possible to estimate power without them.

Abbreviations

AI: Allelic Imbalance; ASE: Allele Specific Expression; TAD: Topologically associated domain.

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13104-021-05851-x>.

Additional file 1. Additional Methods. This file contains the section additional methods, in which we summarize the definition of the Bayesian model used in this work. The model has been previously published and described, and we provide in additional methods a brief summary just to facilitate readers.

Additional file 2. List of simulation parameters. Excel file consisting of two worksheet. Worksheet "Data" contains the simulation parameters used in the various simulations performed for this work. Worksheet "Legend" contains the description of the parameters.

Additional file 3. Variation of type I error as a function of number of simulations, number of allele specific reads per bioreps and extent of deviation from allelic balance.

Additional file 4. Variation of power as a function of number of bioreps.

Additional file 5. Variation of power as a function of number of simulations.

Additional file 6. Variation of power as a function of number of allele specific reads per biorep.

Additional file 7. Variation of power as a function of the extent of deviation from allelic balance.

Acknowledgements

Authors acknowledge the HiPerGator High Performance Super Computer at the University of Florida, and are grateful to UFRC (University of Florida Research Computing) for valuable assistance.

Authors' contributions

KS performed analysis and wrote the code, LGL-N developed the statistical model and wrote the code, SN interpreted the data, LMM developed the statistical model and interpreted the data, FM wrote the code and interpreted the data. All authors contributed to writing the manuscript. All authors read and approved the final manuscript

Funding

This work was financed by the National Institutes of Health (NIH) grant, National Institute of General Medical Sciences Grant Number GM128193/GM/NIGMS), Lauren M. McIntyre.

Availability of data and materials

This study was performed using programs written in Python and R that are available using the MIT license as the package BayesASE_power: https://github.com/McIntyre-Lab/BayesASE_power. The package requires the installation of BayesASE available on PyPI: <https://pypi.org/project/BayesASE/>. All the additional files are available at <https://osf.io/sw3r2/>, with the following <https://doi.org/10.17605/OSF.IO/SW3R2>

Declarations**Ethics approval and consent to participate**

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Author details

¹Quantitative and Computational Biology Section, University of Southern California, Los Angeles, CA 90046, USA. ²Department of Biostatistics and Data Science, The University of Texas Health Science Center at Houston-School of Public Health, Houston, TX 77030, USA. ³Molecular and Computational Biology Section, University of Southern California, Los Angeles, CA 90046, USA. ⁴Genetics Institute and Department of Molecular Genetics and Microbiology, University of Florida, Gainesville, FL 32603, USA. ⁵Dipartimento di Scienze Agroalimentari, Ambientali e Animal, Università di Udine, 33100 Udine, Italy.

Received: 16 July 2021 Accepted: 15 November 2021

Published online: 27 November 2021

References

- Wittkopp PJ, Haerum BK, Clark AG. Evolutionary changes in cis and trans gene regulation. *Nature*. 2004;430:85–8.
- Genissel A, McIntyre LM, Wayne ML, Nuzhdin SV. Cis and trans regulatory effects contribute to natural variation in transcriptome of *Drosophila melanogaster*. *Mol Biol Evol*. 2007;25:101–10.
- Graze RM, McIntyre LM, Main BJ, Wayne ML, Nuzhdin SV. Regulatory divergence in *Drosophila melanogaster* and *D. simulans*, a genomewide analysis of allele-specific expression. *Genetics*. 2009;183:547–61.
- Graze RM, Novelo LL, Amin V, Fear JM, Casella G, Nuzhdin SV, et al. Allelic imbalance in *Drosophila* hybrid heads: exons, isoforms, and evolution. *Mol Biol Evol*. 2012;29:1521–32.
- Zou F, Sun W, Crowley JJ, Zhabotynsky V, Sullivan PF, de Pardo-Manuel Villena F. A novel statistical approach for jointly analyzing RNA-Seq data from F1 reciprocal crosses and inbred lines. *Genetics*. 2014;197:389–99.
- Fear JM, León-Novelo LG, Morse AM, Gerken AR, Van Lehmann K, Tower J, et al. Buffering of genetic regulatory networks in *Drosophila melanogaster*. *Genetics*. 2016;203:1177–90.
- Miller BR, Morse AM, Borgert JE, Liu Z, Sinclair K, Gamble G, et al. Test-crosses are an efficient strategy for identifying cis-regulatory variation: bayesian analysis of allele-specific expression (BayesASE). *G3*. 2021. <https://doi.org/10.1093/g3journal/jkab096>.
- Gregg C, Zhang J, Weissbourd B, Luo S, Schroth GP, Haig D, et al. High-resolution analysis of parent-of-origin allelic expression in the mouse brain. *Science*. 2010;329:643–8.
- Baran Y, Subramaniam M, Biton A, Tukiainen T, Tsang EK, Rivas MA, et al. The landscape of genomic imprinting across diverse adult human tissues. *Genome Res*. 2015;25:927–36.
- Crowley JJ, Zhabotynsky V, Sun W, Huang S, Pakatci IK, Kim Y, et al. Analyses of allele-specific gene expression in highly divergent mouse crosses identifies pervasive allelic imbalance. *Nat Genet*. 2015;47:353–60.
- Rivas MA, Pirinen M, Conrad DF, Lek M, Tsang EK, Karczewski KJ, et al. Impact of predicted protein-truncating genetic variants on the human transcriptome. *Science*. 2015;348:666–9.
- Shao L, Xing F, Xu C, Zhang Q, Che J, Wang X, et al. Patterns of genome-wide allele-specific expression in hybrid rice and the implications on the genetic basis of heterosis. *PNAS*. 2019;116:5653–8.
- Mugal CF, Wang M, Backström N, Wheatcroft D, Ålund M, Sémon M, et al. Tissue-specific patterns of regulatory changes underlying gene expression differences among *Ficedula flycatchers* and their naturally occurring F1 hybrids. *Genome Res*. 2020. <https://doi.org/10.1101/gr.254508.119>.
- Kukurba KR, Zhang R, Li X, Smith KS, Knowles DA, Tan MH, et al. Allelic expression of deleterious protein-coding variants across human tissues. *PLoS Genet*. 2014;10:e1004304.
- Mayba O, Gilbert HN, Liu J, Haverly PM, Jhunjhunwala S, Jiang Z, et al. MBASED: allele-specific expression detection in cancer tissues and cell lines. *Genome Biol*. 2014;15:405.
- Tuch BB, Laborde RR, Xu X, Gu J, Chung CB, Monighetti CK, et al. Tumor transcriptome sequencing reveals allelic expression imbalances associated with copy number alterations. *PLoS ONE*. 2010;5:e9317.
- Liu Z, Dong X, Li Y. A genome-wide study of allele-specific expression in colorectal cancer. *Front Genet*. 2018;9:570.
- Pandey RV, Franssen SU, Futschik A, Schlötterer C. Allelic imbalance metre (Allim), a new tool for measuring allele-specific gene expression with RNA-seq data. *Mol Ecol Resour*. 2013;13:740–5.
- Skelly DA, Johansson M, Madeoy J, Wakefield J, Akey JM. A powerful and flexible statistical framework for testing hypotheses of allele-specific gene expression from RNA-seq data. *Genome Res*. 2011;21:1728–37.
- León-Novelo LG, McIntyre LM, Fear JM, Graze RM. A flexible Bayesian method for detecting allelic imbalance in RNA-seq data. *BMC Genomics*. 2014;15:920.
- Edsgård D, Iglesias MJ, Reilly S-J, Hamsten A, Tornvall P, Odeberg J, et al. GeneiASE: Detection of condition-dependent and static allele-specific expression from RNA-seq data without haplotype information. *Sci Rep*. 2016;6:21134.
- León-Novelo L, Gerken AR, Graze RM, McIntyre LM, Marroni F. Direct testing for allele-specific expression differences between conditions. *G3*. 2018;8:447–60.
- Guo M, Yang S, Rupe M, Hu B, Bickel DR, Arthur L, et al. Genome-wide allele-specific expression analysis using Massively Parallel Signature Sequencing (MPSS™) Reveals cis- and trans-effects on gene expression in maize hybrid meristem tissue. *Plant Mol Biol*. 2008;66:551–63.
- Castel SE, Aguet F, Mohammadi P, Aguet F, Anand S, Ardlie KG, et al. A vast resource of allelic expression data spanning human tissues. *Genome Biol*. 2020;21:234.
- Springer NM, Stupar RM. Allele-specific expression patterns reveal biases and embryo-specific parent-of-origin effects in hybrid maize. *Plant Cell*. 2007;19:2391–402.
- Tukiainen T, Villani A-C, Yen A, Rivas MA, Marshall JL, Satija R, et al. Landscape of X chromosome inactivation across human tissues. *Nature*. 2017;550:244–8.
- Ardlie KG, Deluca DS, Segre AV, Sullivan TJ, Young TR, The GTEx Consortium, et al. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science*. 2015;348:648–60.
- Pirinen M, Lappalainen T, Zaitlen NA, Dermitzakis ET, Donnelly P, GTEx Consortium, et al. Assessing allele-specific expression across multiple tissues from RNA-seq read data. *Bioinformatics*. 2015;31:2497–504.
- Choi K, Raghupathy N, Churchill GA. A Bayesian mixture model for the analysis of allelic expression in single cells. *Nat Commun*. 2019;10:5188.
- Degner JF, Marioni JC, Pai AA, Pickrell JK, Nkadori E, Gilad Y, et al. Effect of read-mapping biases on detecting allele-specific expression from RNA-sequencing data. *Bioinformatics*. 2009;25:3207–12.

31. Wang M, Uebbing S, Ellegren H. Bayesian inference of allele-specific gene expression indicates abundant cis-regulatory variation in natural flycatcher populations. *Genome Biol Evol.* 2017;9:1266–79.
32. van de Geijn B, McVicker G, Gilad Y, Pritchard JK. WASP: allele-specific software for robust molecular quantitative trait locus discovery. *Nat Methods.* 2015;12:1061–3.
33. Castel SE, Levy-Moonshine A, Mohammadi P, Banks E, Lappalainen T. Tools and best practices for data processing in allelic expression analysis. *Genome Biol.* 2015;16:195.
34. Sham PC, Purcell SM. Statistical power and significance testing in large-scale genetic studies. *Nat Rev Genet.* 2014;15:335–46.
35. Hong EP, Park JW. Sample size and statistical power calculation in genetic association studies. *Genomics Inf.* 2012;10:117–22.
36. Fontanillas P, Landry CR, Wittkopp PJ, Russ C, Gruber JD, Nusbaum C, et al. Key considerations for measuring allelic expression on a genomic scale using high-throughput sequencing. *Mol Ecol.* 2010;19(Suppl 1):212–27.
37. Torosin NS, Anand A, Golla TR, Cao W, Ellison CE. 3D genome evolution and reorganization in the *Drosophila melanogaster* species group. *PLoS Genet.* 2020;16:e1009229.
38. Main BJ, Bickel RD, McIntyre LM, Graze RM, Calabrese PP, Nuzhdin SV. Allele-specific expression assays using Solexa. *BMC Genomics.* 2009;10:422.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

