

permGWAS: Efficient Permutation-based Genome-wide Association Studies for Skewed Phenotypic Distributions

Maura John^{1,2}, Markus J. Ankenbrand³, Carolin Artmann³, Jan A. Freudenthal³, Arthur Korte³, Dominik G. Grimm^{1,2,4}

1 Technical University of Munich, TUM Campus Straubing for Biotechnology and Sustainability, Bioinformatics

2 Weihenstephan-Triesdorf University of Applied Sciences, Bioinformatics

3 Center for Computational and Theoretical Biology, University of Würzburg

4 Technical University of Munich, TUM Department of Informatics

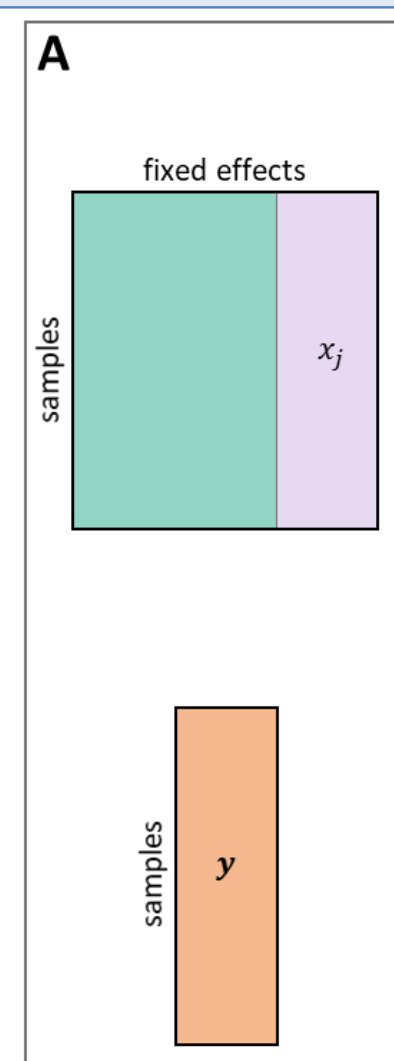
Motivation

- Genome-wide association studies (GWAS) are a key tool to analyze relationship between genotypes and phenotypes
- Linear mixed models (LMMs) test whether a single genetic marker is associated with a given phenotype via the model

$$y = X\beta + u + \epsilon \quad \epsilon \sim \mathcal{N}(0, \sigma_e^2 I), u \sim \mathcal{N}(0, \sigma_g^2 K)$$

phenotype fixed effects random effects noise variance components genetic similarity matrix
- Variance components are estimated by maximizing the likelihood function $L(\beta, \sigma_g^2, \sigma_e^2)$
- Assumptions of Gaussian distribution and independent genetic markers are often violated in real-world data
- Commonly used Bonferroni significance threshold [1] often too conservative for normally distributed phenotypes or not stringent enough for phenotypes with skewed distributions
- Idea: permutation-based threshold via *maxT* method [2]: Permute phenotype q times and compute minimal p-value for each permutation. The adjusted threshold is the α^{th} percentile
- Problem: enormous computational complexity
- Solution: permGWAS, an efficient reformulation of LMMs using 3D and 4D tensors that can provide permutation-based thresholds

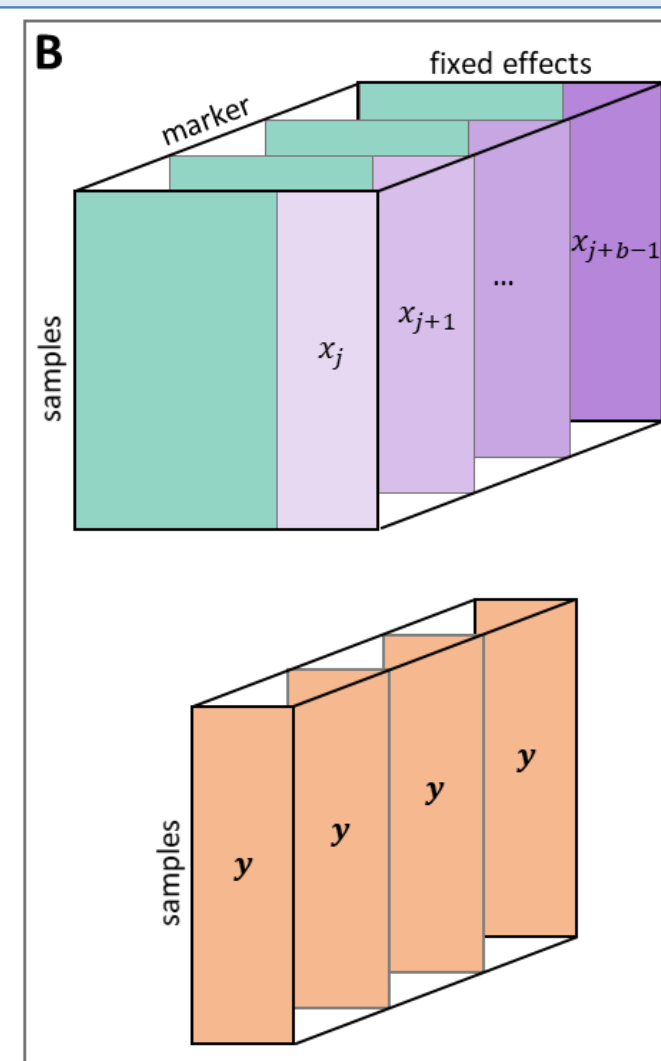
permGWAS Architecture



Basic Linear Mixed Model

Estimate the effect of a single genetic marker:

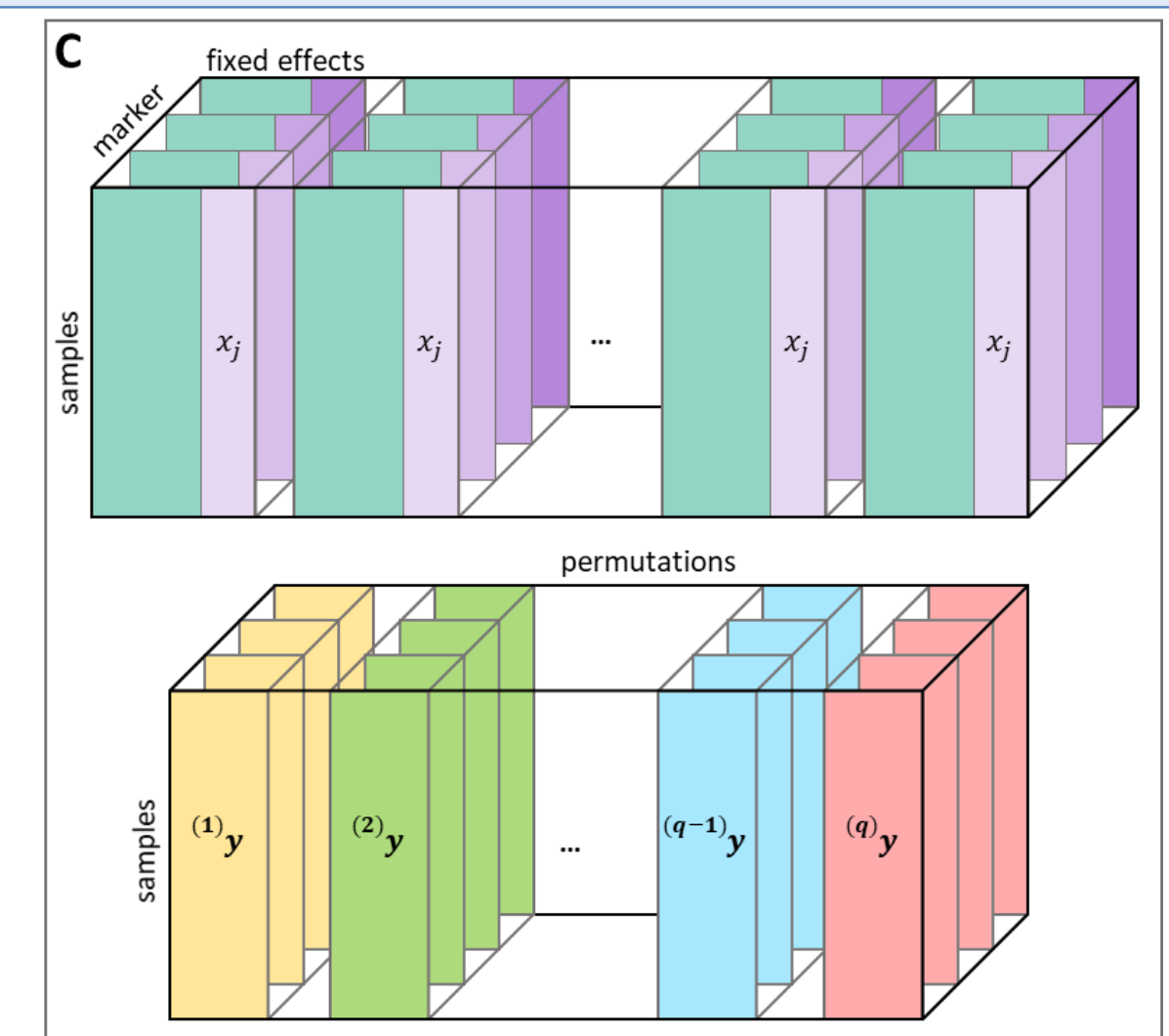
$$L(\beta, \sigma_g^2, \sigma_e^2) = \mathcal{N}(y | X\beta; \sigma_g^2 K + \sigma_e^2 I)$$



Batch-wise Linear Mixed Model

Estimate the effects of several genetic markers simultaneously using 3D-tensors:

$$L(\beta_j^b, \sigma_g^2, \sigma_e^2) = \mathcal{N}(y | X_j^b \beta_j^b; \sigma_g^2 K + \sigma_e^2 I)$$



Permutation-based Linear Mixed Model

Estimate the effects of several genetic markers and permutations simultaneously using 4D-tensors:

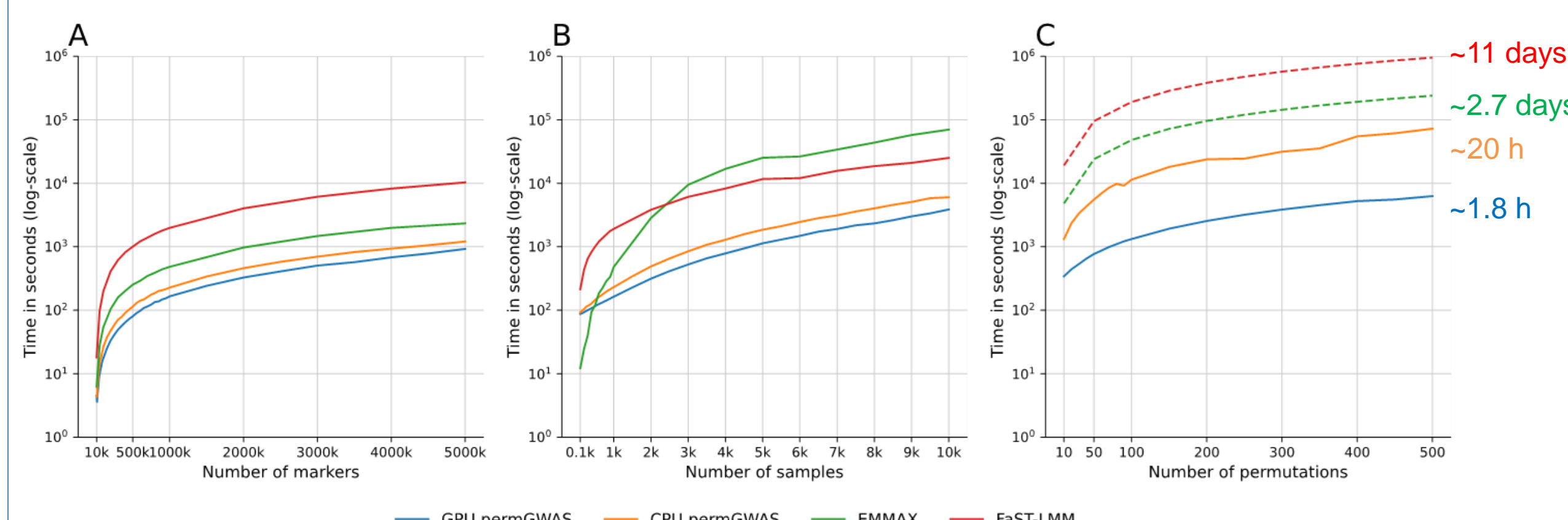
$$L(q\beta_j^b, \sigma_g^2, \sigma_e^2) = \mathcal{N}(y | X_j^b q\beta_j^b; \sigma_g^2 K + \sigma_e^2 I)$$

Results

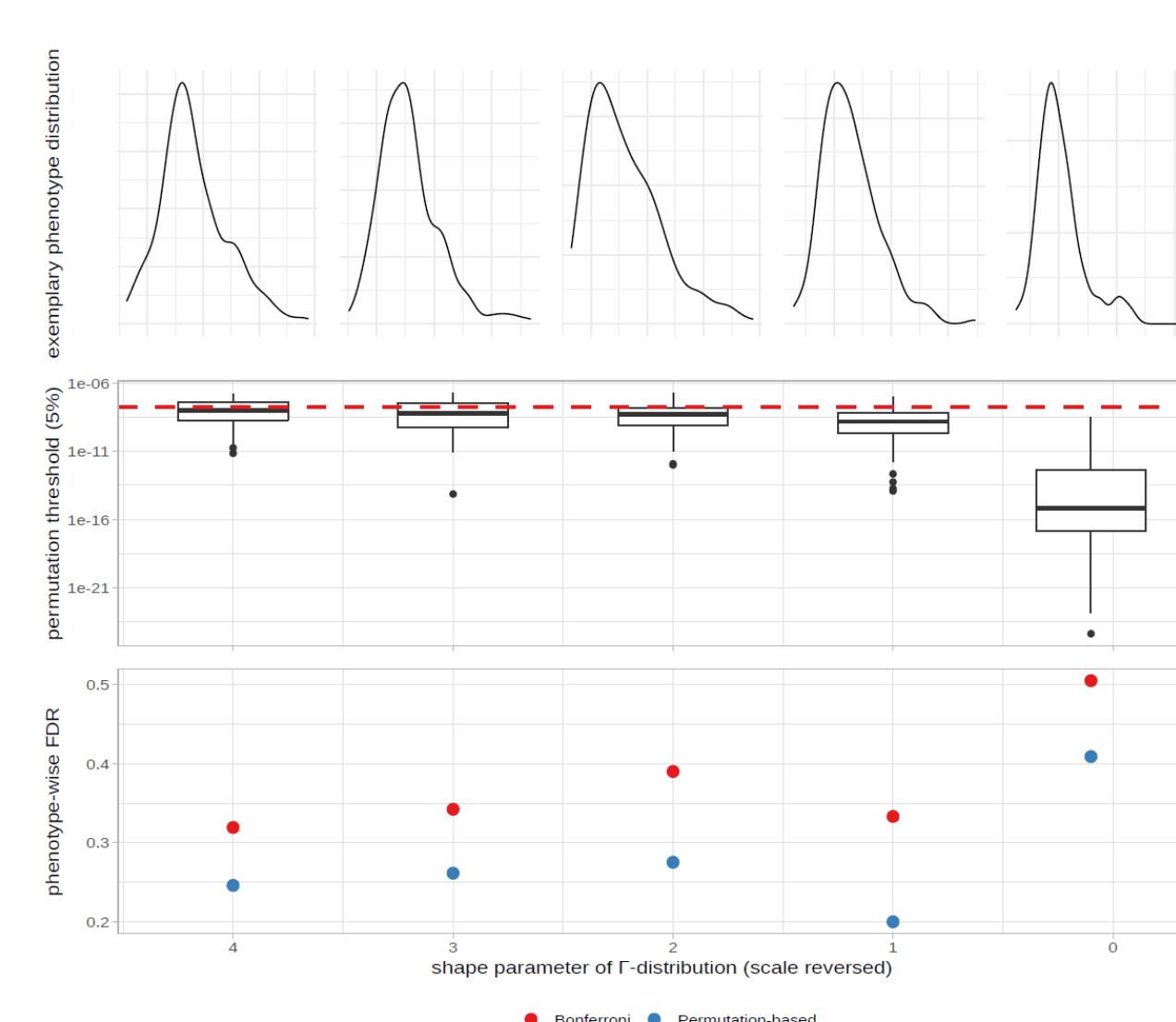
Runtime comparison to state-of-the-art

Comparison of the runtime of permGWAS to EMMAX [3] and FaST-LMM [4] with respect to

- Number of **markers** with fixed number of 1000 samples
- Number of **samples** with fixed number of 10^6 markers
- Number of **permutations** with 1000 samples and 10^6 markers



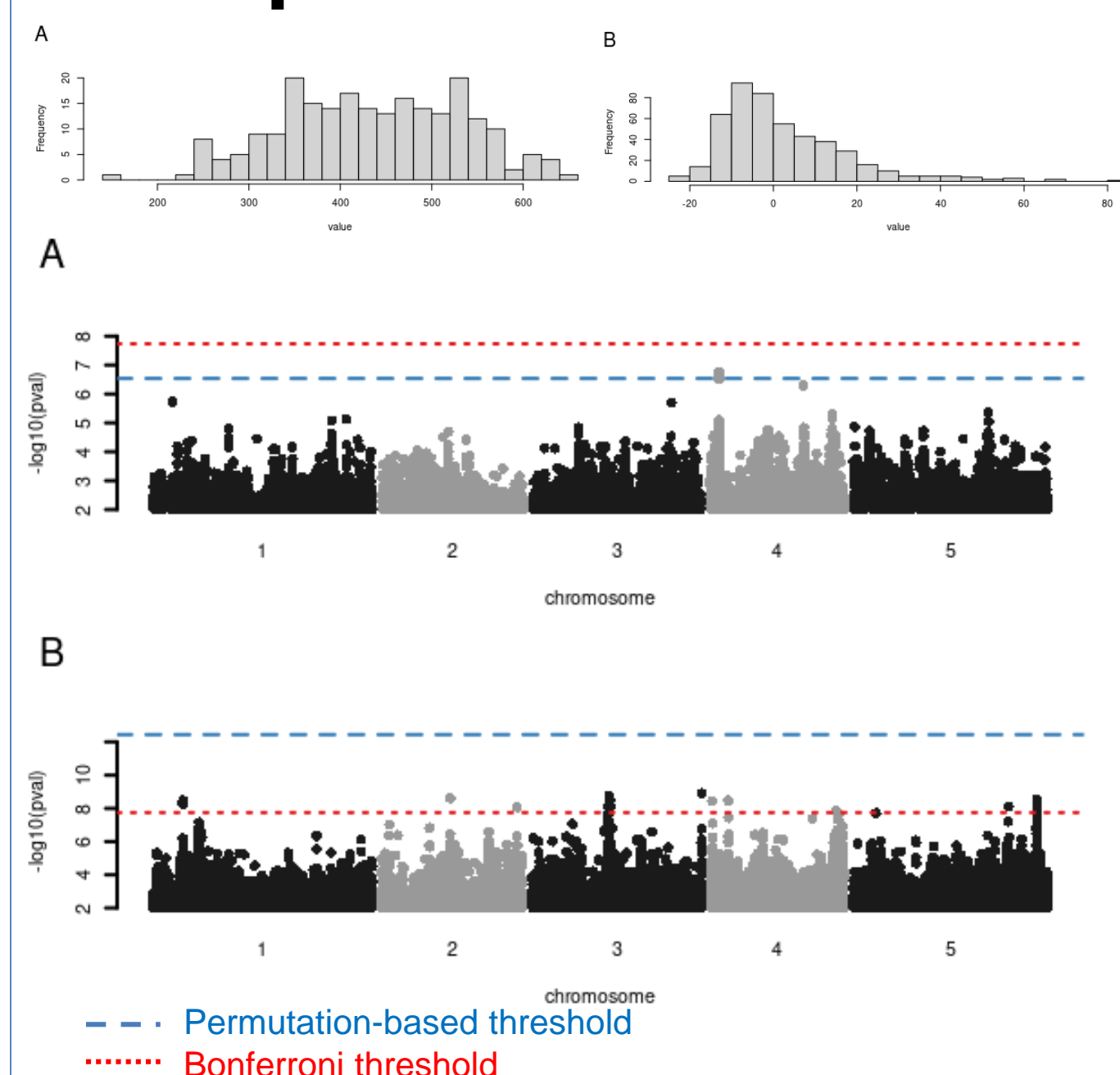
False Discovery Rate for skewed phenotypes



Comparison of False Discovery Rate (FDR) of permutation-based thresholds and Bonferroni threshold for synthetic phenotypes with gamma-distributed noise:

- FDR increases when phenotypes become more skewed
- Permutation-based thresholds become more stringent for skewed phenotypes
- Permutation-based thresholds have lower FDR for skewed phenotypes than Bonferroni threshold

Examples of GWAS in *Arabidopsis thaliana*



Analysis of GWAS results for real phenotypes from model plant *Arabidopsis thaliana* [5,6]:

- For normally distributed phenotypes the permutation-based threshold is less conservative
- For non-normally distributed phenotypes the permutation-based threshold is more stringent

References

Publication:

John, M., Ankenbrand, M. J., Artmann, C., Freudenthal, J. A., Korte, A., & Grimm, D. G. (2022). Efficient Permutation-based Genome-wide Association Studies for Normal and Skewed Phenotypic Distributions. *bioRxiv*. **Accepted at ECCB 2022**

References:

- Bonferroni, C. (1936). Teoria statistica delle classi e calcolo delle probabilit . Pubblicazioni del R Istituto Superiore di Scienze Economiche e Commerciali di Firenze, 8, 3–62.
- Westfall, P. H. and Young, S. S. (1993). Resampling-based multiple testing: Examples and methods for p-value adjustment, volume 279. John Wiley & Sons.
- Kang, H. M., Sul, J. H., Service, S. K., Zaitlen, N. A., Kong, S.-y., Freimer, N. B., Sabatti, C., and Eskin, E. (2010). Variance component model to account for sample structure in genome-wide association studies. *Nature genetics*, 42(4), 348–354.
- Lippert, C., Listgarten, J., Liu, Y., Kadie, C. M., Davidson, R. I., and Heckerman, D. (2011). Fast linear mixed models for genome-wide association studies. *Nature methods*, 8(10), 833–835.
- Seren,  ., Grimm, D., Fitz, J., Weigel, D., Nordborg, M., et al. (2016). Arapheno: a public database for arabidopsis thaliana phenotypes. *Nucleic Acids Research*, page gkw986.
- Togninalli, M., Seren,  ., Freudenthal, J. A., Monroe, J. G., Meng, D., et al. (2020). Arapheno and the aragwas catalog 2020: a major database update including rna-seq and knockout mutation data for arabidopsis thaliana. *Nucleic acids research*, 48(D1), D1063–D1068.

Website:

<http://bit.cs.tum.de>

GitHub:

<https://github.com/grimmlab/permGWAS>



Funding:

