



This is a repository copy of *Robust inference of genetic architecture in mapping studies.*

White Rose Research Online URL for this paper:
<http://eprints.whiterose.ac.uk/118996/>

Version: Accepted Version

Article:

Slate, J. orcid.org/0000-0003-3356-5123 (2017) Robust inference of genetic architecture in mapping studies. *Molecular Ecology*, 26 (6). pp. 1453-1455. ISSN 0962-1083

<https://doi.org/10.1111/mec.14052>

This is the peer reviewed version of the following article: Slate, J. (2017), Robust inference of genetic architecture in mapping studies. *Mol Ecol*, 26: 1453–1455, which has been published in final form at <https://doi.org/10.1111/mec.14052>. This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for Self-Archiving.

Reuse

Unless indicated otherwise, fulltext items are protected by copyright with all rights reserved. The copyright exception in section 29 of the Copyright, Designs and Patents Act 1988 allows the making of a single copy solely for the purpose of non-commercial research or private study within the limits of fair dealing. The publisher or other rights-holder may allow further reproduction and re-use of this version - refer to the White Rose Research Online record for this item. Where records identify the publisher as the copyright holder, users can verify any specific terms of use on the publisher's website.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

1 Robust inference of genetic architecture in mapping studies

2 Jon Slate

3 Department of Animal & Plant Sciences, University of Sheffield, Sheffield, S10 2TN, UK

4 j.slate@sheffield.ac.uk

5

6 Due to the development of next-generation sequencing and related tools, the feasibility of gene
7 mapping studies in wild populations has improved dramatically. However, phenotypic data collection
8 remains challenging and sample sizes are typically orders of magnitude smaller than are seen in
9 genome wide association studies (GWAS) of human populations, where hundreds of thousands of
10 people may be screened (Wood *et al.* 2014). Consequently, the power to detect quantitative trait loci
11 (QTL) remains modest, unless the focal trait is segregating for genes of major effect. Worse still, small
12 sample sizes result in crude estimates of effect sizes; those that are mostly severely overestimated, will
13 be the ones most likely to reach statistical significance - the well-known 'Beavis Effect' (Beavis 1994).
14 This makes inference from mapping studies very problematic. Does a significant peak with a large effect
15 on phenotypic variation represent a true hit, with the trait being determined by relatively few genes of
16 large effect (an 'oligogenic' architecture)? Or does the peak simply represent an upwardly biased
17 estimate and a false positive QTL? Without replication, it is very hard to know which scenario is true.
18 Similarly, interpretation of a null result (no significant QTL) can be problematic. Does this mean that the
19 study was underpowered to pick up any medium-large effect loci that are present? Or does the trait
20 have a genuinely polygenic architecture, caused by many loci of small effect, each of which is
21 undetectable in that particular experiment? It's not hard to see why there remains scepticism about the
22 value of gene mapping studies to evolutionary research (Rockman 2012; Travisano & Shaw 2013). In
23 this issue of *Molecular Ecology*, Li *et al.* (2017) describe a gene mapping experiment that goes some
24 way to addressing the issues of low power and inflated effect size that have plagued previous studies.
25 The authors have conducted a mapping study of brain traits in nine-spined sticklebacks *Pungitius*
26 *pungitius* to tackle two alternative hypotheses about the genetic architecture of brain morphology. At
27 face value, the study could be seen as a relatively standard GWAS, albeit with an impressive number of
28 markers for what is not a classical model organism. However, scratch beneath the surface a little, and it
29 becomes clear that some sophisticated analytical approaches have been used to try to understand trait
30 architecture in a more rigorous way than is typical.

31

32 In the study Li and colleagues used an F2 mapping population, derived by crossing a marine female
33 from the Baltic Sea with a freshwater male from a pond in Northern Finland. There were 239
34 phenotyped and genotyped F2 individuals, and a little over 15,000 SNPs obtained from genotyping-
35 by-sequencing, which were mapped to 21 linkage groups (the known number of chromosomes in this
36 species). The authors measured the volume of five different parts of the brain, and were interested in
37 comparing two alternative hypotheses about the genetic architecture of brain traits. Under the
38 *mosaic model* each brain component has a distinct genetic architecture, and it is free to evolve
39 without genetic constraint from other brain components. The alternative idea, the *concerted model*,
40 posits that brain component evolution is constrained, perhaps due to a common genetic architecture
41 influencing the different parts. Recognising that the experimental design was exactly the kind where
42 QTL of major effect could be identified spuriously in a standard linkage mapping (or GWAS)
43 experiment, the authors utilised an approach known as *de-biased Least Absolute Shrinkage and*
44 *Selection Operator (LASSO) mapping* (Van de Geer *et al.* 2014; Zhang & Zhang 2014). De-biased LASSO

45 has not been widely employed in gene-mapping studies but it has several advantageous properties.
46 Perhaps the most obviously different feature is that multiple markers are modelled simultaneously.
47 This reduces the risk of effect size overestimation, and can also lead to an increase in power. It also
48 facilitates the estimation of heritability by summing the effect of SNPs fitted in the model. Thus, Li *et*
49 *al.* (2017) had two main aims. The empirical goal was to understand the genetic architecture of brain
50 traits and evaluate whether the mosaic or concerted model was more plausible. The methodological
51 goal was to compare different mapping approaches using both real and simulated data, to establish
52 whether de-biased LASSO gave more reliable parameter estimates than approaches that fit single
53 markers consecutively.

54

55 For many of the traits that Li and colleagues studied, they found genomic regions that explained
56 significant genetic variation, even at a stringent genomewide significance threshold. There was very
57 little between-trait overlap in QTL locations, so the data were consistent with the mosaic model
58 (different genomic regions affecting different brain components). When running single SNP (i.e.
59 conventional) analyses, there were frequently numerous, tightly-linked significant SNPs, and their
60 estimated effect sizes were frequently 5-10% of the overall trait variation. These would be regarded
61 as genes of reasonably large effect. The multi-marker (i.e. de-biased LASSO) approach usually
62 identified the same genomic regions. However, the effect sizes were typically much smaller - 1% or
63 less of the phenotypic variation. The multimarker analyses could also measure the effect of each
64 linkage group on trait variation. Summing these effects across linkage groups provided estimates of
65 trait heritability; depending on the trait these ranged from ~0.10 - ~0.45. Some linkage groups
66 contributed disproportionately to additive genetic variation, but not in a way that supported the
67 concerted model (which would have predicted that the same linkage group would contribute to
68 different traits). Overall, the data suggest that brain traits are moderately heritable in this cross, and
69 that the trait architecture is consistent with the mosaic model of brain evolution. One slight caveat is
70 that the experimental cross was derived from just a single pair of fish, and we are largely ignorant of
71 how much genetic variation is segregating within *versus* between marine and freshwater populations.
72 Other experimental designs may have yielded quite different conclusions.

73

74 How well does the multi-marker de-biased LASSO method perform? The simulations, which are
75 presented in the supplementary material, explored both an oligogenic and a polygenic scenario.
76 Unsurprisingly, the single-locus approach had a high false positive rate, and a lower power to detect
77 true positives than de-biased LASSO (although neither approach had high power when the simulated
78 QTL were small). Effect size estimation of individual QTL was actually greatly downward-biased with
79 de-biased LASSO and tended to be upwardly biased with single-locus estimates. Encouragingly
80 though, de-biased LASSO provided accurate estimates of overall trait heritability, regardless of trait
81 architecture. Single-locus approaches fail in this regard, especially when the true architecture is
82 polygenic. It may well be the case that the downward bias of de-biased LASSO effect size estimates
83 can be rectified by summing the effects of linked SNPs in the region, or by first pruning the marker
84 data, so that retained SNPs are not in high linkage disequilibrium. The F2 design used by Li and
85 colleagues, probably causes linked SNPs to be in strong linkage disequilibrium. More generally, de-
86 biased LASSO is one of several recently introduced approaches that fit multiple markers
87 simultaneously (Moser *et al.* 2015; Zhou *et al.* 2013). These methods are beginning to be adopted in
88 evolutionary / ecological studies (Comeault *et al.* 2015) and they are attractive for several reasons.
89 First, they facilitate a more holistic approach to studying trait architecture, where instead of paying
90 slavish attention to 'significant peaks', QTL effect sizes, heritabilities, individual breeding values and

91 whole-chromosome contributions (*sensu* Yang et al. (2011) to genetic variation can be estimated
92 simultaneously. Second, they go a long way towards avoiding incorrect inference of an oligogenic trait
93 architecture that almost inevitably comes about from upwardly-biased single-marker effect size
94 estimates of true or false QTLs. It remains to be seen which multimarker method performs best,
95 although benchmarking studies of some approaches do exist (Moser *et al.* 2015). Perhaps, the most
96 heartening thing is that ecological genetic mapping studies such as the one by Li and colleagues, are
97 beginning to mature in a way that hypothesis-driven questions can be addressed without the reliance
98 on the detection and identification of specific QTL that may or may not true positives.

99

100

101 **References**

102

103 Beavis WD (1994) The power and deceit of QTL experiments: lessons from comparative QTL studies,
104 250-266.

105 Comeault AA, Flaxman SM, Riesch R, *et al.* (2015) Selection on a Genetic Polymorphism Counteracts
106 Ecological Speciation in a Stick Insect. *Current Biology* **25**, 1975-1981.

107 Li Z, Guo B, Yang J, *et al.* (2017) Deciphering the genomic architecture of the stickleback brain with a
108 novel multi-locus gene-mapping approach. *Molecular Ecology* **26**, xxx-xxx.

109 Moser G, Lee SH, Hayes BJ, *et al.* (2015) Simultaneous Discovery, Estimation and Prediction Analysis
110 of Complex Traits Using a Bayesian Mixture Model. *Plos Genetics* **11**.

111 Rockman MV (2012) THE QTN PROGRAM AND THE ALLELES THAT MATTER FOR EVOLUTION: ALL
112 THAT'S GOLD DOES NOT GLITTER. *Evolution* **66**, 1-17.

113 Travisano M, Shaw RG (2013) LOST IN THE MAP. *Evolution* **67**, 305-314.

114 Van de Geer S, Buhlmann P, Ritov Y, Dezeure R (2014) On Asymptotically Optimal Confidence
115 Regions and Tests for High-Dimensional Models. *Annals of Statistics* **42**, 1166-1202.

116 Wood AR, Esko T, Yang J, *et al.* (2014) Defining the role of common variation in the genomic and
117 biological architecture of adult human height. *Nat Genet* **46**, 1173-1186.

118 Yang J, Manolio TA, Pasquale LR, *et al.* (2011) Genome partitioning of genetic variation for complex
119 traits using common SNPs. *Nature Genetics* **43**, 519-525.

120 Zhang CH, Zhang SS (2014) Confidence intervals for low dimensional parameters in high dimensional
121 linear models. *Journal of the Royal Statistical Society Series B-Statistical Methodology* **76**,
122 217-242.

123 Zhou X, Carbonetto P, Stephens M (2013) Polygenic Modeling with Bayesian Sparse Linear Mixed
124 Models. *Plos Genetics* **9**.

125