

PAPER • OPEN ACCESS

GIFT: new method for the genetic analysis of small gene effects involving small sample sizes

To cite this article: Cyril Rauch *et al* 2023 *Phys. Biol.* **20** 016001

View the [article online](#) for updates and enhancements.

You may also like

- [The White-light Superflares from Cool Stars in GWAC Triggers](#)
Guang-Wei Li, , Liang Wang et al.
- [pygiftgenerator: a python module designed to prepare Moodle-based quizzes](#)
Jon Sáenz, Idoia G Gurtubay, Zunbeltz Izaola et al.
- [The Automatic Observation Management System of the GWAC Network. I. System Architecture and Workflow](#)
Xuhui Han, Yujie Xiao, PinPin Zhang et al.

OPEN ACCESS



CrossMark

RECEIVED

11 July 2022

REVISED

4 October 2022

ACCEPTED FOR PUBLICATION

12 October 2022

PUBLISHED

3 November 2022

Original content from this work may be used under the terms of the [Creative Commons Attribution 4.0 licence](#).

Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.



PAPER

GIFT: new method for the genetic analysis of small gene effects involving small sample sizes

Cyril Rauch^{1,*} , Panagiota Kyratzi², Sarah Blott¹, Sian Bray³ and Jonathan Wattis²¹ School of Veterinary Medicine and Science, University of Nottingham, College Road, Sutton Bonington, LE12 5RD, United Kingdom² School of Mathematical Sciences, University of Nottingham, University Park, NG7 2RD, United Kingdom³ School of Life Sciences, University of Nottingham, University Park, NG7 2RD, United Kingdom

* Author to whom any correspondence should be addressed.

E-mail: cyril.rauch@nottingham.ac.uk**Keywords:** phenotype–genotype mapping, complex traits, GWAS, GIFT, field theorySupplementary material for this article is available [online](#)

Abstract

Small gene effects involved in complex/omnigenic traits remain costly to analyse using current genome-wide association studies (GWAS) because of the number of individuals required to return meaningful association(s), a.k.a. study power. Inspired by field theory in physics, we provide a different method called genomic informational field theory (GIFT). In contrast to GWAS, GIFT assumes that the phenotype is measured precisely enough and/or the number of individuals in the population is too small to permit the creation of categories. To extract information, GIFT uses the information contained in the cumulative sums difference of gene microstates between two configurations: (i) when the individuals are taken at random without information on phenotype values, and (ii) when individuals are ranked as a function of their phenotypic value. The difference in the cumulative sum is then attributed to the emergence of phenotypic fields. We demonstrate that GIFT recovers GWAS, that is, Fisher's theory, when the phenotypic fields are linear (first order). However, unlike GWAS, GIFT demonstrates how the variance of microstate distribution density functions can also be involved in genotype–phenotype associations when the phenotypic fields are quadratic (second order). Using genotype–phenotype simulations based on Fisher's theory as a toy model, we illustrate the application of the method with a small sample size of 1000 individuals.

1. Introduction

Identifying the association between phenotypes and genotypes is the fundamental basis of genetic analyses. In the early days of genetic studies, beginning with Mendel's work at the end of the 19th century, genotypes were inferred by tracking the inheritance of phenotypes between individuals with known relationships (linkage analysis). In recent years, the development of molecular tools, culminating in high-density genotyping and whole genome sequencing, has enabled DNA variants to be directly identified and phenotypes to be associated with genotypes in large populations of unrelated individuals through association mapping. Genome-wide association studies (GWAS) have become the method of choice, largely replacing linkage analyses, because they are more powerful for mapping complex traits, that is, they can

be used to detect smaller gene effects, and they provide a greater mapping precision as they depend on population-level linkage disequilibrium rather than close family relationships. For example, the 2021 NHGRI-EBI GWAS catalogue currently lists 316 782 associations identified in 5149 publications describing GWAS results [1]. Additionally, extensive data collection has been initiated through efforts such as the UK Biobank [2], Generation Scotland [3] and NIH All of Us research programme (<https://allofus.nih.gov/>) with the expectation that large-scale GWAS will elucidate the basis of human health and disease and facilitate precision medicine.

While genomic technologies used to generate data have rapidly advanced within the last 20 years, the statistical models used in GWAS to analyse the data are still predominantly based on Fisher's method published than 100 years ago [4, 5]. Using probability

density functions (PDFs) and in particular the normal distribution, Fisher's method partitions genotypic values by performing a linear regression of the phenotype on marker allelic dosage [6]. The regression coefficient estimates the average allele effect size, and the regression variance is the additive genetic variance due to the locus [7]. While Fisher's method has been improved, for example using conditional probability linked to potential prior knowledge of genetic systems (Bayes' method) [8, 9], the overall determination of genotype–phenotype mapping is still grounded on PDFs. However, the use of PDFs become problematic in the case of complex/omnigenic traits as they require large scale-study or equivalently, large sample size.

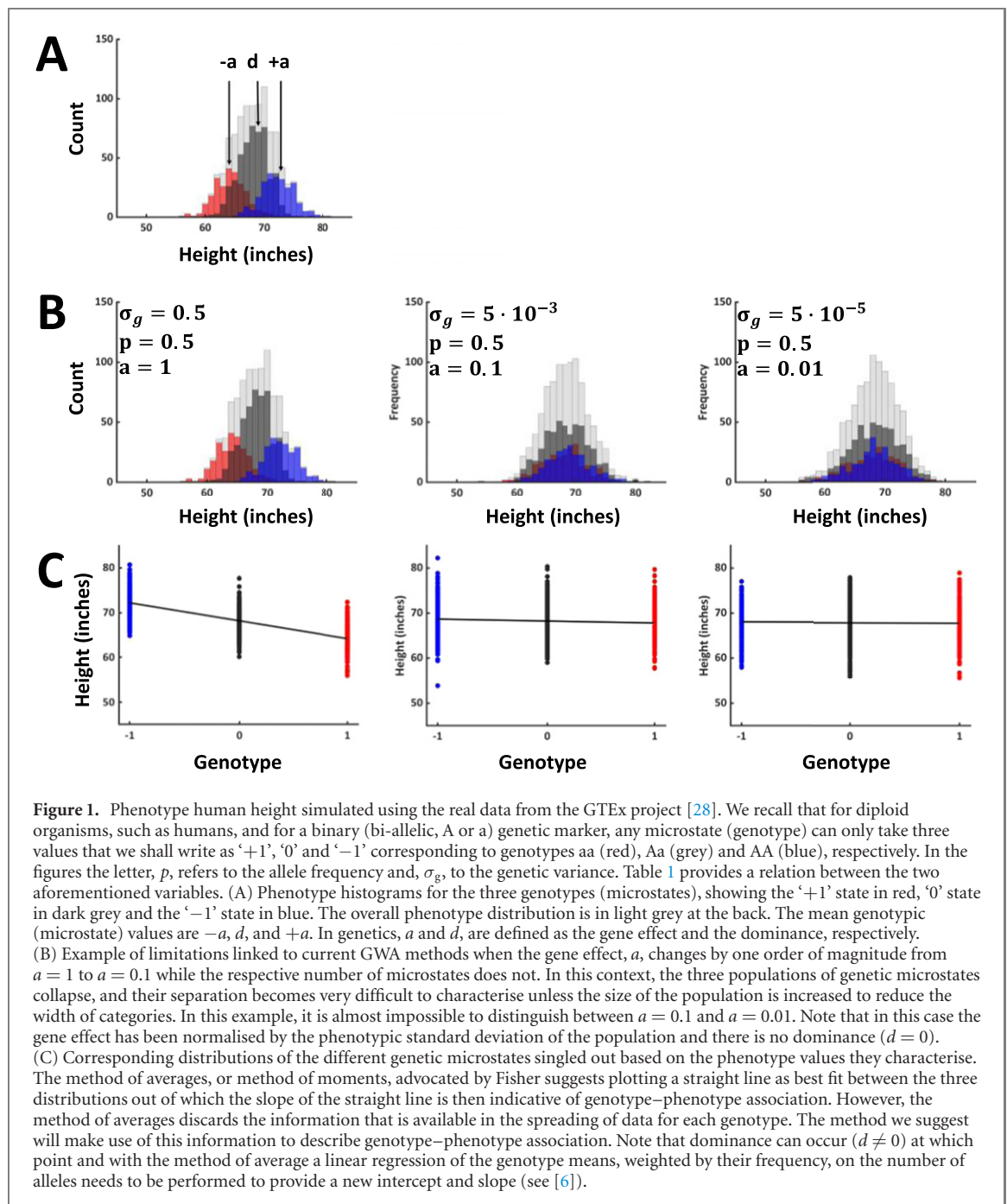
The results obtained by GWAS have demonstrated that complex traits are driven by a vast number of tiny-effect loci, namely a vast number of genes each with tiny-effect, and not by a handful of moderate-effect loci as initially thought. In turn, this has led to a re-conceptualisation of the genetic basis of complex traits from being polygenic (handful of loci/genes) to omnigenic (vast number of loci/genes) [10–16]. Although the omnigenic paradigm is central to further our understanding of biology, there is a practical issue concerning the extraction of information to relate genotype to phenotype in this case. Indeed, tiny-effect loci (i.e., very small gene effects) necessitate a remarkably large population to extract information. Figure 1 exemplifies the limit of GWAS with a restricted sample size of 1000 individuals. This issue regarding the need for large sample sizes was present, but dismissed, in Fisher's seminal work [4] as he assumed an 'infinite population' from the start to use the normal distribution density function in the continuum limit. This assumption allowed him to provide a method able to extract, in theory, the genetic information required to map any genotype to phenotype.

While one may assume an infinite population mathematically, in practice this comes at a huge cost. To give a 'real-life' example of the sample size needed to study complex traits the best is to turn to the phenotype 'height' in humans. The phenotype height in humans is a classical quantitative trait that has been studied for over a century as a model for investigating the genetic basis of complex traits [4, 17] and whose measured heritability is well known [4, 18, 19]. However, this phenotype has remained controversial [12] for a long time as current association methods were not been able to fully recover the heritability measured [21, 22]. While different reasons were put forward to explain this discrepancy including, for example, too restricted sample sizes, too stringent statistical tests or the involvement of the environment [6, 23]; this point seems to have been resolved only recently. Using a population containing a staggering 5.3 million individuals a recent

study claims to have captured nearly all of the common single nucleotide polymorphisms (SNPs)-based heritability [24].

This important study confirms that the precision of current quantitative genetic methods to determine omnigenic traits comes at an astronomical cost in line with the assumptions used, namely the need for a staggeringly large (near infinite) population. In this context one may wonder whether such large-scale study will ever be replicated in any other species and in particular those near extinction where small sample sizes need to be considered. Alternatively, one may try to understand where the need for large sample sizes comes from and determine whether it is possible to extract information linking genotype to phenotype in a different way.

There is a very good reason as to why large populations will always be required for omnigenic phenotypes when GWAS is used. As mentioned above, the reason is rooted in the fact that GWAS is mostly based on frequentist probabilities a.k.a. PDFs. Indeed, GWAS is based on statistics and, by definition, statistics deals with the measurement of uncertainties [25]. To draw inferences from the comparison of large datasets, a method that requires some understanding of its accuracy, including ways of measuring the uncertainty in data values, is needed. In this context, statistics is the science of collecting, analysing, and interpreting data, while PDFs defined through the notion of relative frequencies, is central to determining the validity of statistical inferences. In practice, the use of frequentist probabilities (or PDFs) and the resulting categorisation of data is justified when inaccuracy exists in experimental measurements. For example, measuring a continuous phenotype such as the height of individuals with a ruler with centimetre graduations, that is, to the nearest centimetre, warrants the use of frequentist probability (or PDFs). In this case, a frequency table of phenotypic values can be defined through 1 cm-width bins or categories, from which the PDFs of the phenotype height and of the genotypes can be deduced to address the statistical inferences. However, the precision available for the inferences will always be, at best, given by the width of categories created and linked to the experimental precision achieved (1 cm in this case). Consequently, if instead of using a ruler with centimetre graduations one was using a ruler with millimetre graduations to increase the precision in inferences, a larger sample size would be required such as to match the new 1 mm-width of categories to reform the PDFs. The trading between the sample size and the precision achieved by GWAS is known as 'study power' and its *raison d'être* is linked to the fact that the entire field of probability, and therefore the PDFs, has been conceptualised mathematically to represent the fact that information on a system is limited. It is for



this reason that the normal distribution was known before as the ‘error function’ or ‘law of errors’, where the term ‘error’ is defined experimentally (see rulers above). Accordingly, the creation of categories implies that a sort of ‘imprecision’ is necessary.

While the notion of imprecision can be genuine (see rulers above), the act of creating categories to use PDFs when precision in experimental measurements is available is not fully justified and can be seen as an act of ‘wilful ignorance’. This is so because information is lost by slotting different phenotypic values into the same category. To exemplify this point let us take an example, imagine a species close to extinction (very small population, say 50 individuals) and that it is possible to measure phenotypic values with very

high precision, for example, using highly advanced imaging techniques or biosensing technologies [26]. In this case, each measured individual could return a unique phenotypic value. Consequently, reforming categories to reform and use PDFs would mean embracing the relatively large width of categories created leading to the impression that a large imprecision is present, even so such imprecision did not exist in the first place. One may argue that PDFs, such as the normal distribution in Fisher’s theory, are not required since averages and variances can be mathematically calculated directly from data without the need to recreate PDFs. However, this argument is not valid as averages and variances (and any other moments) cannot be dissociated from PDFs, since

they are the ontological parameters that define PDFs, and that PDFs are used to determine statistical inferences in the field of quantitative genetics. Consequently, thinking in term of averages and variances necessitate to conceptualise and hold valid PDFs, i.e., categories, to describe any system.

Therefore, there is a need to formulate new methods using the full information generated through accurate and highly precise genotyping/phenotyping when sample sizes are small which does not require the categorisation of data. In fact, this problem is equivalent to finding a way to resolve genotype–phenotype mapping by assuming a finite-size population with phenotype values measured precisely enough to rule out the possibility that two phenotype values are in the same category. Taking this challenge as the starting point, a new and relatively simple method for extracting information for genotype–phenotype mapping can be defined. While this method is remarkably simple when explained in lay terms, its theoretical framework requires the introduction of a new concept called ‘phenotypic fields’. Phenotypic fields can also be defined within the context of Fisher’s theory.

The remainder of this paper is organised as follows. In the first part, an intuitive approach to the method genomic informational field theory (GIFT) is presented in which one shall see that association between datasets (i.e., genotype and phenotype) can be analysed in specific way that do not involve the use of means and variances necessarily, but phenotypic fields instead. This is followed by a second part stating and explaining the necessary ingredients from physics (entropy, energy and field) and how they must be combined, to provide GIFT. Since GIFT is not too difficult to model, we have relegated the theoretical development of GIFT in the appendix A. Finally (third part) one will demonstrate how Fisher’s seminal theory can be re-transcribed using GIFT. In particular one shall see that Fisher’s seminal intuition corresponds to the simplest form of GIFT. Finally (fourth part), one will compare GIFT to GWAS using simulated genotype and phenotype to demonstrate that GIFT outperforms GWAS.

2. Position of the problem and heuristic presentation of GIFT a method

The practical issue regarding genotype–phenotype mappings with current statistical methods concerns

the sample size needed to provide accurate/precise information when complex/omnigenic traits are involved. As stated in the introduction, this issue stems from the creation of categories historically linked to the notions of ‘imprecision’ or ‘error’ in measurements. At the dawn of the 21st century we are getting more precise in our measurements, and one may wonder what sort of scientific/mathematical tool we should be using if one were able to attain any level of precision wanted in cases where the population size studied is limited.

We recall here that for diploid organisms, such as humans, and for a binary (bi-allelic, A or a) genetic marker, any microstate (genotype) can only take three values that we shall write as ‘+1’, ‘0’ and ‘−1’ corresponding to genotypes aa (homozygote), Aa (heterozygote) and AA (homozygote), respectively.

One way to proceed to develop a method embracing precision is to start by looking at how density distribution functions are transformed when precision in phenotypic measurements increases. From figure 2(A), the conclusion is obvious, the bar charts are transformed into code bars, where each bar originates from a particular phenotype value representing one individual from the population studied. This result is expected since when the width of categories decreases due to an increase in precision in measurements, there will be a point where there can only be one individual per category. To extract information from the code bars represented by the bottom right chart in figure 2(A), let us now wonder what it means to have information on the phenotype as opposed to have none.

To answer this question the best thing is to further simplify the problem by considering the coloured bars only and not their spacing. Imagine, therefore, that a set of individuals has been genotyped and that those individuals are picked at random. That is, there is no information on any phenotype. Imagine also that one decides to concentrate, for example, on the genome position 1 000 000 on chromosome 4 for all the individuals since this genome position happens to display a biallelic SNPs across the set of individuals.

Thus, upon calling randomly but sequentially individuals, the genotypic information obtained in due course can be represented as a random string of genotypes including ‘+1’, ‘0’ and ‘−1’ microstates (representing homozygote-AA, heterozygote-Aa and homozygote-aa). An example of such random configuration is:

$$[-1, +1, 0, -1, -1, +1, +1, \dots, -1, +1, +1, 0, -1, 0, +1, \dots, 0, 0, -1, +1, 0, +1, -1].$$

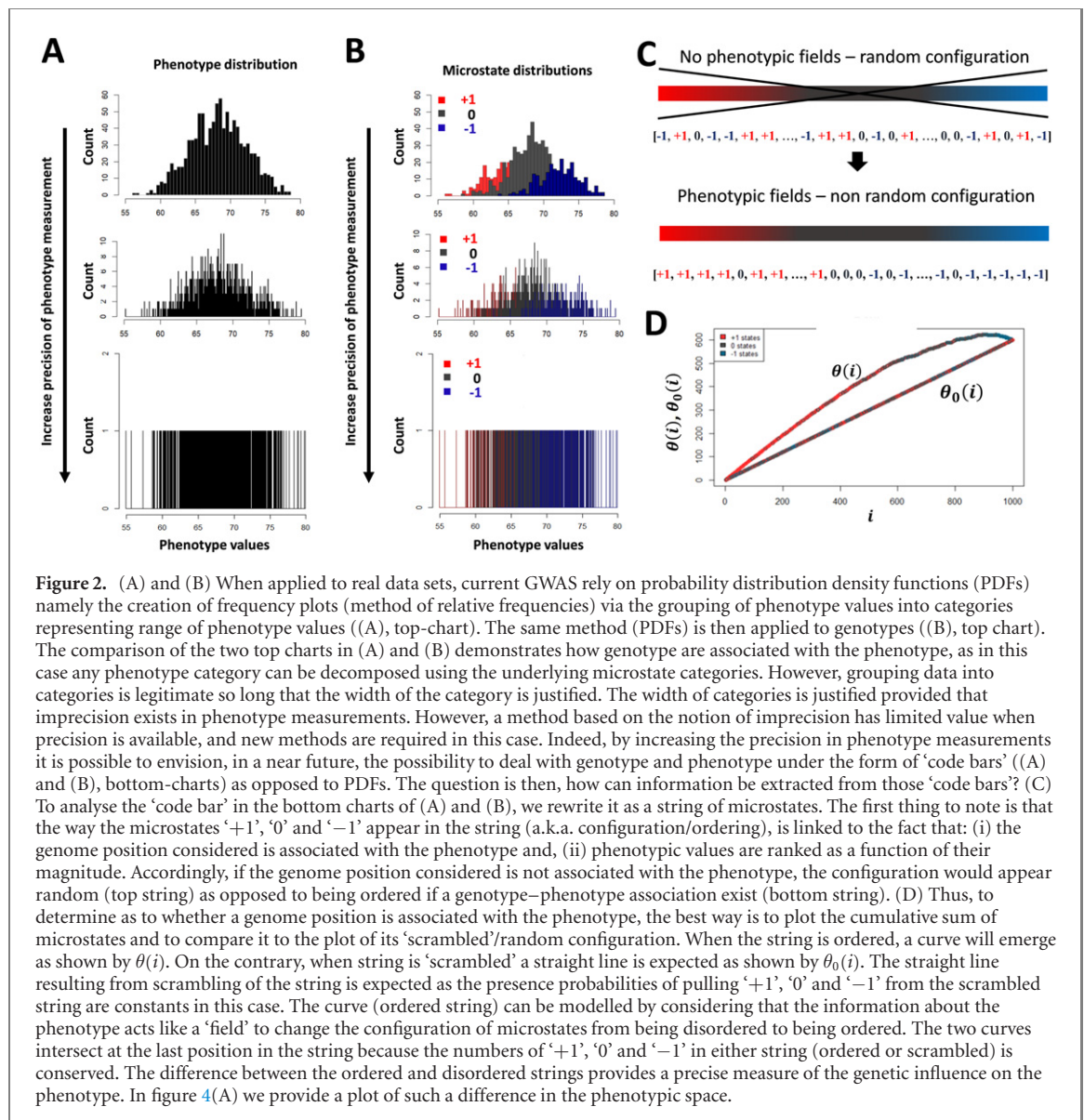


Figure 2. (A) and (B) When applied to real data sets, current GWAS rely on probability distribution density functions (PDFs) namely the creation of frequency plots (method of relative frequencies) via the grouping of phenotype values into categories representing range of phenotype values ((A), top-chart). The same method (PDFs) is then applied to genotypes ((B), top chart). The comparison of the two top charts in (A) and (B) demonstrates how genotype are associated with the phenotype, as in this case any phenotype category can be decomposed using the underlying microstate categories. However, grouping data into categories is legitimate so long that the width of the category is justified. The width of categories is justified provided that imprecision exists in phenotype measurements. However, a method based on the notion of imprecision has limited value when precision is available, and new methods are required in this case. Indeed, by increasing the precision in phenotype measurements it is possible to envision, in a near future, the possibility to deal with genotype and phenotype under the form of ‘code bars’ ((A) and (B), bottom-charts) as opposed to PDFs. The question is then, how can information be extracted from those ‘code bars’? (C) To analyse the ‘code bar’ in the bottom charts of (A) and (B), we rewrite it as a string of microstates. The first thing to note is that the way the microstates ‘+1’, ‘0’ and ‘−1’ appear in the string (a.k.a. configuration/ordering), is linked to the fact that: (i) the genome position considered is associated with the phenotype and, (ii) phenotypic values are ranked as a function of their magnitude. Accordingly, if the genome position considered is not associated with the phenotype, the configuration would appear random (top string) as opposed to being ordered if a genotype–phenotype association exist (bottom string). (D) Thus, to determine as to whether a genome position is associated with the phenotype, the best way is to plot the cumulative sum of microstates and to compare it to the plot of its ‘scrambled’/random configuration. When the string is ordered, a curve will emerge as shown by $\theta(i)$. On the contrary, when string is ‘scrambled’ a straight line is expected as shown by $\theta_0(i)$. The straight line resulting from scrambling of the string is expected as the presence probabilities of pulling ‘+1’, ‘0’ and ‘−1’ from the scrambled string are constants in this case. The curve (ordered string) can be modelled by considering that the information about the phenotype acts like a ‘field’ to change the configuration of microstates from being disordered to being ordered. The two curves intersect at the last position in the string because the numbers of ‘+1’, ‘0’ and ‘−1’ in either string (ordered or scrambled) is conserved. The difference between the ordered and disordered strings provides a precise measure of the genetic influence on the phenotype. In figure 4(A) we provide a plot of such a difference in the phenotypic space.

Note that the order in which the individuals were called is linked to the position in the string. Let us now repeat the same experiment using the same individuals in a context where accurate information on a chosen phenotype is available. That is, we call the individuals as a function of the magnitude of their phenotype we consider. For example, if the phenotype is height, one starts by calling the smallest individual and all subsequent individuals through successive increments in their phenotype height. Note again that because each

individual has a unique phenotype value there is no possibility for two individuals to be called at once.

If the genome position 1 000 000 on chromosome 4 is involved in the formation of the phenotype, then one would expect a change in the configuration of the string of microstates based on the fact that homozygotes would be found at the extremities of the string and heterozygotes towards the middle (see figure 2). An example of such a string would be, for example:

$$[+1, +1, +1, +1, 0, +1, +1, \dots, +1, 0, 0, 0, -1, 0, -1, \dots, -1, 0, -1, -1, -1, -1, -1].$$

Thus, the only thing that changes between the random and the phenotype-ordered configurations is the way the genetic microstates are allocated to positions in the string. However, as the genome position 1 000 000 on chromosome 4 is the only one that has been considered, the two configurations contain the same number of ‘+1’, ‘0’ and ‘−1’, since the same individuals were considered between the two configurations.

The *ansatz* is then to consider the cumulative sum of microstates as a function of the position in the string. Indeed, it is clear from the examples given above that if one starts by adding the microstates together, differences will be seen in the resulting cumulative sums. To give an example, let us consider the two strings above and note ‘ $\theta_0(i)$ ’ and ‘ $\theta(i)$ ’ the cumulative sums of microstates in the random and ordered configurations, respectively, where ‘ i ’ is the position in the string. Then adding the microstates starting from the left side of the strings one finds:

$$\begin{aligned}\theta_0(1) &= -1 = -1 & \theta(1) &= +1 = +1 \\ \theta_0(2) &= -1 + 1 = 0 & \theta(2) &= +1 + 1 = +2 \\ \theta_0(3) &= -1 + 1 + 0 = 0 & \theta(3) &= +1 + 1 + 1 = +3.\end{aligned}$$

As a result, the difference ‘ $\theta(i) - \theta_0(i)$ ’ is expected to be indicative of the importance of the phenotypic information and how gene microstates are related to the phenotype. The fact that the same individuals were considered in both configurations also imposes a conservation relation under the form: $\theta(N) - \theta_0(N) = 0$. One shall call the cumulative sums: ‘genetic paths’ whose mathematical definition will be précised below. To conclude, it is the information on the phenotypic values that provides a change in the configuration of microstates and one can start developing the formulations of, $\theta_0(i)$ and $\theta(i)$.

Noting, N_{+1} , N_0 and N_{-1} the number of genetic microstates ‘+1’, ‘0’ and ‘−1’, respectively. The genetic microstate frequencies for genome position 1 000 000 on chromosome 4 are defined by, $N_{+1}/N = \omega_+^0$, $N_0/N = \omega_0^0$ and $N_{-1}/N = \omega_-^0$.

When the positioning of the genetic microstates in the string is performed in a random fashion, the probabilities of finding ‘+1’, ‘0’ or ‘−1’ as genetic microstate at any position are ω_+^0 , ω_0^0 and ω_-^0 , respectively. The resulting cumulative sum is then: $\theta_0(i) = (+1 \cdot \omega_+^0 + 0 \cdot \omega_0^0 - 1 \cdot \omega_-^0)i$. Consequently, $\theta_0(i)$ is therefore a straight line. We shall call $\theta_0(i)$ the ‘default genetic path’ (figure 2(C)).

In the second configuration the microstates are ordered. Noting $\omega_+(i)$, $\omega_0(i)$ and $\omega_-(i)$ the occurrence probabilities of the genetic microstates ‘+1’, ‘0’ and ‘−1’ the cumulative sum is then: $\theta(i) = \sum_1^i (+1 \cdot \omega_+(j) + 0 \cdot \omega_0(j) - 1 \cdot \omega_-(j))$; where $\theta(i)$ is defined as the ‘phenotype-responding genetic path’ (figure 2(C)).

As a result, the signature of a gene interacting with the phenotype when considering the two aforementioned genetic paths is the difference: $\theta(i) - \theta_0(i) = \sum_1^i [(\omega_+(j) - \omega_-(j)) - (\omega_+^0 - \omega_-^0)]$. One can then be a little bit more prescriptive by introducing the notion of phenotypic fields.

3. A physics-inspired model for GIFT: notion of ‘phenotypic fields’ and resulting difference between the phenotype-responding and default genetic paths

The difference, $\theta(i) - \theta_0(i)$, can be described using field theory. Indeed, as the only difference between the two configurations is the information linked to the phenotypic values, the phenotypic information can be thought as an external field impacting the configuration of microstates. To provide a physics-inspired definition of genotype–phenotype mapping, let us reconsider the random string above and assume that the set of individuals in the string are particles and that the different genetic microstates ‘+1’, ‘0’ and ‘−1’ are their physical properties. One can then assume that it is those properties that interact with the field. Note that contrary to physics where a single field is defined, one needs in our case to define one field per microstate. One shall note by $u_+(\Omega)$, $u_0(\Omega)$ and $u_-(\Omega)$ the phenotypic fields acting on the microstates ‘+1’, ‘0’ and ‘−1’ respectively. Note that the variable Ω represents the phenotypic values measured precisely. By assuming further that the particles cannot interact together and that, when they are not forced into a specific configuration by the fields, they can hop and exchange positions when the field is null (similar to a diffusion/thermal process), one can then model the string of microstates as a closed system. Figure 2(C) provides an idea of how the ‘phenotypic fields’, when non null, impacts on the configuration of microstates by segregating them. With those assumptions and using basic principles from statistical physics it is then possible to model the presence probability of microstates at any position in the string.

Thus, after re-expressing the genetic paths in the space of phenotypic values since the fields are function of the phenotypic values (appendix A.1), one can then construct functionals representing the entropy (appendix A.2) and the total interaction energy between the microstates and the subfields (appendix A.3). Finally, one can optimise a functional similar to the free energy to express how the fields are related to the asymmetry of states (appendix A.4). Consequently, one can demonstrate the familiar result concerning the presence probability of microstates

expressed in the phenotypic space as,

$$\hat{\omega}_+(\Omega) = \frac{\omega_+^0 e^{-\delta u_+(\Omega)}}{\omega_0^0 + \omega_+^0 e^{-\delta u_+(\Omega)} + \omega_-^0 e^{-\delta u_-(\Omega)}} \quad (1)$$

$$\hat{\omega}_0(\Omega) = \frac{\omega_0^0}{\omega_0^0 + \omega_+^0 e^{-\delta u_+(\Omega)} + \omega_-^0 e^{-\delta u_-(\Omega)}} \quad (2)$$

$$\hat{\omega}_-(\Omega) = \frac{\omega_-^0 e^{-\delta u_-(\Omega)}}{\omega_0^0 + \omega_+^0 e^{-\delta u_+(\Omega)} + \omega_-^0 e^{-\delta u_-(\Omega)}}. \quad (3)$$

Where the hat ‘ $\hat{\cdot}$ ’ is added to insist on the fact that the presence probabilities of microstates are expressed in the space of phenotypic values (and not positions) and, $\delta u_+(\Omega) \stackrel{\text{def}}{=} u_+(\Omega) - u_0(\Omega)$, $\delta u_-(\Omega) \stackrel{\text{def}}{=} u_-(\Omega) - u_0(\Omega)$. Equations (1)–(3) are familiar to physicists when dealing with Boltzmann’s weigh in statistical physics. Note that the default genetic path is defined when the fields are null. Noting, $\Delta\hat{\theta}(\Omega) \stackrel{\text{def}}{=} \hat{\theta}(\Omega) - \hat{\theta}_0(\Omega)$, the difference between the phenotype responding and default genetic paths expressed in the phenotypic space, $\Delta\hat{\theta}(\Omega)$ is therefore a function of the difference between equations (1) and (3). One can then make the symmetries of the problem more apparent by defining for the genetic microstates, $\Delta\omega_0 \stackrel{\text{def}}{=} \omega_+^0 - \omega_-^0$ and $\omega_0 \stackrel{\text{def}}{=} \omega_+^0 + \omega_-^0 = 1 - \omega_0^0$; and for the phenotypic fields, $2\bar{u}(\Omega) \stackrel{\text{def}}{=} \delta u_+(\Omega) + \delta u_-(\Omega) = u_+(\Omega) + u_-(\Omega) - u_0(\Omega)$ and $2\Delta u(\Omega) \stackrel{\text{def}}{=} \delta u_+(\Omega) - \delta u_-(\Omega) = u_+(\Omega) - u_-(\Omega)$. In this case, using hyperbolic functions one deduces (see appendix A.4 for development),

$$\hat{\omega}_+(\Omega) - \hat{\omega}_-(\Omega) = \frac{\text{sh}(\Delta u(\Omega_0) - \Delta u(\Omega))}{\alpha_0 e^{\bar{u}(\Omega)} + \text{ch}(\Delta u(\Omega_0) - \Delta u(\Omega))}. \quad (4)$$

Where, $\text{th}(\Delta u(\Omega_0)) \stackrel{\text{def}}{=} \frac{\Delta\omega_0}{\omega_0}$ and $\alpha_0 \stackrel{\text{def}}{=} \frac{1-\omega_0}{\omega_0} \text{ch}(\Delta u(\Omega_0)) = \frac{1}{2} \frac{\omega_0^0}{\sqrt{\omega_+^0 \omega_-^0}}$. The new variable ‘ Ω_0 ’ is the phenotype value corresponding to the condition $\hat{\omega}_+(\Omega_0) \sim \hat{\omega}_-(\Omega_0)$ and the meaning of the constant ‘ α_0 ’ can be related to the Hardy–Weinberg law from population genetic. Hardy–Weinberg law based on random mating in a population provides a relationship between the genetic microstate frequencies under the form: $p^2 + 2pq + q^2 = 1$, where p^2 and q^2 are the genotype frequencies of genetic microstates ‘+1’ and ‘−1’, i.e. homozygote genotypes aa and AA, respectively; and $2pq$ the genotype frequency for genetic microstate ‘0’, i.e. the heterozygote genotype Aa. In our case, this corresponds to replacing p^2 , q^2 and $2pq$ with, respectively, ω_+^0 , ω_-^0 and ω_0^0 . Consequently, the Hardy–Weinberg law imposes $\alpha_0 = 1$ with $\alpha_0 \neq 1$ corresponding to a deviation from the law. However, this term is expected to remain stable upon any changes of allele or genotype frequencies suggesting therefore that, genetically, any changes in ‘ $\Delta\omega_0$ ’ are to some extent compensated by corresponding changes in ‘ ω_0 ’. Finally, using equation (4) one deduces the difference between the phenotype responding and default genetic paths expressed in the phenotypic space, $\Delta\hat{\theta}(\Omega)$, as (appendix A.5):

$$\Delta\hat{\theta}(\Omega) = \int_{\Omega_{1/N}}^{\Omega} \left[\frac{\text{sh}(\Delta u(x_0) - \Delta u(x))}{\alpha_0 e^{\bar{u}(x)} + \text{ch}(\Delta u(x_0) - \Delta u(x))} - \frac{\text{sh}(\Delta u(x_0))}{\alpha_0 + \text{ch}(\Delta u(x_0))} \right] \frac{1}{\Delta(x)} dx. \quad (5)$$

Where $\Omega_{1/N}$ is the smallest phenotypic value measured and, $\Delta(x)$, is the spacing between individuals in the code bar figure 2(A) that can be related to the PDF of the phenotype when the population measured is dense (see appendix A.1 and SM1 in the supplementary materials). Finally, the conservation of genetic microstates needs to be added, that is, $\Delta\hat{\theta}(\Omega_1) = 0$, expressed as,

$$\int_{\Omega_{1/N}}^{\Omega_1} \frac{\text{sh}(\Delta u(x_0) - \Delta u(x))}{\alpha_0 e^{\bar{u}(x)} + \text{ch}(\Delta u(x_0) - \Delta u(x))} \frac{1}{\Delta(x)} dx = \frac{\text{sh}(\Delta u(x_0))}{\alpha_0 + \text{ch}(\Delta u(x_0))}. \quad (6)$$

Where Ω_1 is the largest phenotypic value measured. Therefore, as α_0 is constant since a single genome position is considered, the genetic paths difference can be re-expressed integrally using two independent reduced phenotypic fields, i.e., $\Delta u(\Omega)$ and $\bar{u}(\Omega)$, and equation (6) provides a coupling between those fields and $\Delta(\Omega)$. The advantage of using fields is the reduction of unknown parameters involved in the problem and the possibility of laying out genotype–phenotype associations based on fields’ symmetry. For example, and as a minimalist model, one may wonder what sort of expression would take the fields if the reference field $u_0(\Omega)$ was null and, δu_+ and δu_- , were acting anti-symmetrically and linearly on the microstates ‘+1’ and ‘−1’? This minimalist model can be developed (see SM2 in the supplementary materials) and is similar to Fisher’s seminal intuition concerning genotype–phenotype associations (see below).

Our aim is now to demonstrate that the idea of genetic paths mediated by phenotypic fields already exists in Fisher theory. This can be shown by coarse graining the paths.

4. Coarse graining GIFT: definition of fields in Fisher’s context, implication for small gene effects and definition of variance fields

4.1. Coarse graining GIFT

To derive a coarse-grained version of GIFT, that is, a genetic path difference for GWAS, we assume the existence of categories or bins and concentrate on the interval of phenotype values ranging between Ω and $\Omega + \delta\Omega$ defining one particular category or bin.

Based on frequentist probability, by noting by N the total number of individuals in the population

we can define by: $\delta N \sim N \cdot P_\Omega(\Omega)\delta\Omega$, the number of individuals in the phenotype category concerned namely with a phenotype value ranging between Ω and $\Omega + \delta\Omega$.

Similarly, concentrating on a single genome position, we can define by: $\delta N_+ \sim N_+^0 \cdot P_+(\Omega)\delta\Omega$, $\delta N_0 \sim N_0^0 \cdot P_0(\Omega)\delta\Omega$ and $\delta N_- \sim N_-^0 \cdot P_-(\Omega)\delta\Omega$ the respective number of ‘+1’, ‘0’ and ‘−1’ genetic microstates for the phenotype values ranging between Ω and $\Omega + \delta\Omega$, where N_+^0 , N_0^0 and N_-^0 correspond to the total number of ‘+1’, ‘0’ and ‘−1’ microstates in the population with respective presence PDFs, $P_+(\Omega)$, $P_0(\Omega)$ and $P_-(\Omega)$.

The design of categories generates two conservation relationships: the first one concerns the total number of individuals and microstates, namely, that for a given genome position, the sum of all possible microstates is also the sum of all individuals. This relationship is written as follows: $N = N_+^0 + N_0^0 + N_-^0$. The second conservation relationship is linked to the category considered. The number of individuals in the category concerned is also the sum of the microstates in this category: $\delta N = \delta N_+ + \delta N_0 + \delta N_-$. Consequently, the conservation relation concerning the number of individuals and microstates in the concerned category can be rewritten as,

$$\delta N/N = (N_+^0/N)(\delta N_+/N_+^0) + (N_0^0/N)(\delta N_0/N_0^0) + (N_-^0/N)(\delta N_-/N_-^0). \quad (7)$$

Using the PDF of both microstates and the phenotype defined above, the following is deduced:

$$P_\Omega(\Omega) = \omega_+^0 P_+(\Omega) + \omega_0^0 P_0(\Omega) + \omega_-^0 P_-(\Omega). \quad (8)$$

From equation (8) all the moments of microstate distributions are related to those of the phenotype distribution. Let us note by, $\langle\Omega\rangle$, and, a_+ , a_0 , a_- , the average values of the phenotype and microstates ‘+1’, ‘0’, ‘−1’ distribution functions, respectively; and by σ^2 and σ_+^2 , σ_0^2 , σ_-^2 the variances of the phenotype and microstates ‘+1’, ‘0’, ‘−1’ distribution functions, respectively. From equation (8), one deduces the conservation relations for the first two moments in the form:

$$\langle\Omega\rangle = \omega_+^0 a_+ + \omega_0^0 a_0 + \omega_-^0 a_- \quad (9a)$$

$$\sigma^2 = (\omega_+^0 \sigma_+^2 + \omega_0^0 \sigma_0^2 + \omega_-^0 \sigma_-^2) + \omega_+^0 (\langle\Omega\rangle - a_+)^2 + \omega_0^0 (\langle\Omega\rangle - a_0)^2 + \omega_-^0 (\langle\Omega\rangle - a_-)^2. \quad (9b)$$

The relations provided by equation (9) are valid by definition, namely, whatever PDFs are involved. While Fisher never formulated a conservation similar to the second from equation (9b), in his seminal paper [4] he used the notation α^2 to denote the genetic variance in the form: $\alpha^2 \stackrel{\text{def}}{=} \omega_+^0 (\langle\Omega\rangle - a_+)^2 + \omega_0^0 (\langle\Omega\rangle - a_0)^2 + \omega_-^0 (\langle\Omega\rangle - a_-)^2$.

Let us now define the coarse-grained version of equation (4) by noting, $\delta\hat{\omega}_+(\Omega) - \delta\hat{\omega}_-(\Omega)$, the difference in the presence probability of microstates ‘+1’ and ‘−1’ for the category of interest, it is then deduced:

$$\begin{aligned} \delta\hat{\omega}_+(\Omega) - \delta\hat{\omega}_-(\Omega) &= \frac{\delta N_+}{\delta N} - \frac{\delta N_-}{\delta N} \\ &= \frac{\omega_+^0 P_+(\Omega) - \omega_-^0 P_-(\Omega)}{\omega_+^0 P_+(\Omega) + \omega_0^0 P_0(\Omega) + \omega_-^0 P_-(\Omega)}. \end{aligned} \quad (10)$$

Note that equation (10) corresponds to a local relative difference, in the space of phenotypic values, of microstates ‘+1’ and ‘−1’. Also, it can be verified that the conservation of gene microstates holds by summing, $\delta N_+ - \delta N_-$, over all existing categories or bins, namely: $N_+ - N_- = \sum_{\text{bins}} (\delta N_+ - \delta N_-) = \sum_{\text{bins}} (\delta\hat{\omega}_+(\Omega) - \delta\hat{\omega}_-(\Omega))\delta N$. Using the definition $\delta N = NP_\Omega(\Omega)\delta\Omega$ together with equation (8), we deduce using the continuum limit, $\frac{N_+ - N_-}{N} \sim \int_{\Omega_{1/N}}^{\Omega_1} \frac{\omega_+^0 P_+(\Omega) - \omega_-^0 P_-(\Omega)}{\omega_+^0 P_+(\Omega) + \omega_0^0 P_0(\Omega) + \omega_-^0 P_-(\Omega)} P_\Omega(\Omega)\delta\Omega$ or equivalently, $\frac{N_+ - N_-}{N} \sim \int_{\Omega_{1/N}}^{\Omega_1} (\omega_+^0 P_+(\Omega) - \omega_-^0 P_-(\Omega))\delta\Omega = \omega_+^0 - \omega_-^0$.

Direct mapping of fields can then be performed between GWAS and GIFT (SM3 in the supplementary materials). As a result, we can define the coarse-grained versions of the difference in the genetic paths using the continuum limit as:

$$\begin{aligned} \Delta\hat{\theta}(\Omega) \sim \int_{\Omega_{1/N}}^{\Omega} \left[\frac{\omega_+^0 P_+(\Omega) - \omega_-^0 P_-(\Omega)}{\omega_+^0 P_+(\Omega) + \omega_0^0 P_0(\Omega) + \omega_-^0 P_-(\Omega)} \right. \\ \left. - (\omega_+^0 - \omega_-^0) \right] P_\Omega(\Omega)d\Omega. \end{aligned} \quad (11)$$

Equation (11) demonstrates that $\Delta\hat{\theta}(\Omega)$ is sensitive to the PDFs involved as a whole and not just to the average values. In other words, the variance of microstates and their average values will impact on genotype–phenotype associations. Note that in equation (11) the integration interval is unchanged. However, the convergence in probability of distributions allows some freedom, for example changing the integration interval from $[\Omega_{1/N}; \Omega_1]$ to $[0; +\infty]$, or from $[\Omega_{1/N} - \Omega_m; \Omega_1 - \Omega_m]$ to $[-\infty; +\infty]$, where ‘ Ω_m ’ is a median position.

The fields can then be defined in Fisher’s context setting: $-\delta u_+(\Omega) \stackrel{\text{def}}{=} \ln[P_+(\Omega)/P_0(\Omega)]$ and $-\delta u_-(\Omega) \stackrel{\text{def}}{=} \ln[P_-(\Omega)/P_0(\Omega)]$; and from those relations it is deduced: $\bar{u}(\Omega) = \ln[P_0(\Omega)/\sqrt{P_+(\Omega) \cdot P_-(\Omega)}]$ and $\Delta u(\Omega) = \frac{1}{2} \ln[P_-(\Omega)/P_+(\Omega)]$.

Finally using equation (5)’s notations one deduces:

$$\Delta\hat{\theta}(\Omega) = \int_{\Omega_{1/N}}^{\Omega} \frac{\text{sh}\left(\frac{1}{2}\ln\left(\frac{P_+(\Omega)}{P_-(\Omega)} \cdot \frac{P_-(\Omega_0)}{P_+(\Omega_0)}\right)\right)}{\frac{1}{2}\frac{\omega_0^0 P_0(\Omega)}{\sqrt{\omega_+^0 P_+(\Omega)\omega_-^0 P_-(\Omega)}} + \text{ch}\left(\frac{1}{2}\ln\left(\frac{P_+(\Omega)}{P_-(\Omega)} \cdot \frac{P_-(\Omega_0)}{P_+(\Omega_0)}\right)\right)} P_{\Omega}(\Omega) d\Omega$$

$$- \frac{\text{sh}\left(\frac{1}{2}\ln\left(\frac{P_-(\Omega_0)}{P_+(\Omega_0)}\right)\right)}{\frac{1}{2}\frac{\omega_0^0}{\sqrt{\omega_+^0 \omega_-^0}} + \text{ch}\left(\frac{1}{2}\ln\left(\frac{P_-(\Omega_0)}{P_+(\Omega_0)}\right)\right)} \int_{\Omega_{1/N}}^{\Omega} P_{\Omega}(\Omega) d\Omega. \quad (12)$$

The significance of these fields can now be addressed. The field $\bar{u}(\Omega)$ describes local deviations from the Hardy–Weinberg law, valid for each bin or category of phenotype values. For example, if a population was under no selection and random mating occurred, then the whole population would follow the Hardy–Weinberg equilibrium law, that is, $\frac{1}{2}\omega_0^0/\sqrt{\omega_+^0\omega_-^0} \sim 1$. However, selecting a particular bin or category of phenotype values would demonstrate a local deviation of this law, given by the term $e^{\bar{u}(\Omega)} = P_0(\Omega)/\sqrt{P_+(\Omega)P_-(\Omega)}$.

The signification of $\Delta u(\Omega)$ can be addressed using Fisher's approach.

4.2. Definition of fields using Fisher's theory

In his seminal paper [5], Fisher hypothesised that in a context where the population is infinite to use the normal distribution, the genetic variance ' α^2 ' is much smaller than the phenotype variance and that the variances of microstate distribution density functions for each gene are similar to that of the variance of the phenotype. While his hypothesis can be understood intuitively when all distribution density functions nearly overlap, it can also be demonstrated using equation (9b). Indeed, assuming $\alpha^2 \ll \sigma^2$ implies $\sigma^2 - \alpha^2 \sim \sigma^2$ and therefore $\sigma^2 \sim \omega_+^0 \sigma_+^2 + \omega_0^0 \sigma_0^2 + \omega_-^0 \sigma_-^2$. As $\omega_+^0 + \omega_0^0 + \omega_-^0 = 1$, posing $\sigma_+^2 \sim \sigma_0^2 \sim \sigma_-^2 \sim \sigma^2$ as Fisher did, is one valid solution. However, the relation $\sigma^2 \sim \omega_+^0 \sigma_+^2 + \omega_0^0 \sigma_0^2 + \omega_-^0 \sigma_-^2$, is the equation of an ellipse, and an infinite number of solutions are, in theory, possible. Those solutions will depend on the variances of microstates (see section 4.4. below, that is, the definition of fields linked to the variance of microstates).

Following Fisher's assumption, the probability of finding the microstates '+1', '0' and '-1' as a function of phenotype values are expressed as [4] (figure 1): $P_+(\Omega) \sim \frac{K}{\sigma} \exp\left(-\frac{1}{2}\frac{(\Omega-a_+)^2}{\sigma^2}\right)$, $P_0(\Omega) \sim \frac{K}{\sigma} \exp\left(-\frac{1}{2}\frac{(\Omega-a_0)^2}{\sigma^2}\right)$ and $P_-(\Omega) \sim \frac{K}{\sigma} \exp\left(-\frac{1}{2}\frac{(\Omega-a_-)^2}{\sigma^2}\right)$. Where 'K' is the normalisation constant identical for all microstates distribution

functions since they nearly overlap. Noting, $\delta x \stackrel{\text{def}}{=} x - \langle \Omega \rangle$, the difference between the variable, x , and the average of phenotype values to simplify the notations, as the second orders defined by $(\delta a_+/\sigma)^2$, $(\delta a_0/\sigma)^2$ and $(\delta a_-/\sigma)^2$ are neglected in Fisher's context, one deduces then:

$$\Delta u(\Omega_0) = \frac{1}{2} \ln\left(\frac{P_-(\Omega_0)}{P_+(\Omega_0)}\right)$$

$$\sim \frac{1}{2} \left(\frac{\delta a_- - \delta a_+}{\sigma}\right) \left(\frac{\delta \Omega_0}{\sigma}\right) \quad (13a)$$

$$\Delta u(\Omega_0) - s\Delta u(\Omega) = \frac{1}{2} \ln\left(\frac{P_+(\Omega)}{P_-(\Omega)} \cdot \frac{P_-(\Omega_0)}{P_+(\Omega_0)}\right)$$

$$= \frac{1}{2} \left(\frac{\delta a_- - \delta a_+}{\sigma}\right) \left(\frac{\delta \Omega_0}{\sigma} - \frac{\delta \Omega}{\sigma}\right) \quad (13b)$$

$$\bar{u}(\Omega) = \ln\left(\frac{P_0(\Omega)}{\sqrt{P_+(\Omega) \cdot P_-(\Omega)}}\right)$$

$$\sim \frac{1}{2} \left(\frac{2\delta a_0 - \delta a_+ - \delta a_-}{\sigma}\right) \frac{\delta \Omega}{\sigma}. \quad (13c)$$

In equations (13a) and (13b) the term $\delta a_- - \delta a_+ = a_- - a_+ = 2a$ (see figure 1(A)) is known as the 'gene effect' in GWAS. In equation (13c), the term $2\delta a_0 - \delta a_+ - \delta a_- = 2a_0 - a_+ - a_- = d$ (see figure 1(A)) is the dominance as defined in the GWAS. In his seminal paper, Fisher considered: $d \sim 0$.

Altogether, these results demonstrate that Fisher's theory can be described by phenotypic fields and genetic paths. As it turns out Fisher's model corresponds to the minimalist model aforementioned (SM2 in the supplementary materials). Using these fields, it is also possible to determine a generic solution to equation (8) (see SM4 in the supplementary materials).

4.3. Implication for small gene effects

Complex traits involve genes with very small effects that are difficult to characterise. The aim is to determine the resulting difference in the genetic paths in this case, that is, when the gene effect, $a =$

$(a_- - a_+)/2$ (see definition above), tends towards zero: $a \rightarrow 0$. Because PDFs are used, the integration interval can be altered using the convergence property of the distributions. In this context, the conservation of genetic microstates (equation (6)) can be written using Fisher's fields as,

$$\int_0^{+\infty} \frac{\text{sh}\left(\frac{1}{2}\left(\frac{a}{\sigma}\right)\left(\frac{\delta\Omega_0}{\sigma} - \frac{\delta\Omega}{\sigma}\right)\right)}{\alpha_0 e^{\frac{1}{2}\left(\frac{a}{\sigma}\right)\left(\frac{\delta\Omega}{\sigma}\right)} + \text{ch}\left(\frac{1}{2}\left(\frac{a}{\sigma}\right)\left(\frac{\delta\Omega_0}{\sigma} - \frac{\delta\Omega}{\sigma}\right)\right)} P_{\Omega}(\Omega) d\Omega \sim \frac{\text{sh}\left(\frac{1}{2}\left(\frac{a}{\sigma}\right)\left(\frac{\delta\Omega_0}{\sigma}\right)\right)}{\alpha_0 + \text{ch}\left(\frac{1}{2}\left(\frac{a}{\sigma}\right)\left(\frac{\delta\Omega_0}{\sigma}\right)\right)}. \quad (14)$$

Consider, as Fisher did in his seminal paper [4], a phenotype distribution of the form, $P_{\Omega}(\Omega) \sim \frac{K}{\sigma} e^{-\frac{1}{2}\left(\frac{\Omega - \langle\Omega\rangle}{\sigma}\right)^2}$, and rescale the phenotype values in the integral using, a/σ , as a scaling parameter as follows:

$$\int_0^{+\infty} \frac{\text{sh}\left(\frac{1}{2}\left(\frac{\delta\Omega_0}{\sigma} - \frac{\delta\Omega}{\sigma}\right)\right)}{\alpha_0 e^{\frac{1}{2}\left(\frac{a}{\sigma}\right)\left(\frac{\delta\Omega}{\sigma}\right)} + \text{ch}\left(\frac{1}{2}\left(\frac{\delta\Omega_0}{\sigma} - \frac{\delta\Omega}{\sigma}\right)\right)} K e^{-\frac{1}{2}\left(\frac{\delta\Omega}{\sigma}\right)^2} \times d\left(\frac{\Omega}{a}\right) \sim \frac{\text{sh}\left(\frac{1}{2}\left(\frac{a}{\sigma}\right)\left(\frac{\delta\Omega_0}{\sigma}\right)\right)}{\alpha_0 + \text{ch}\left(\frac{1}{2}\left(\frac{a}{\sigma}\right)\left(\frac{\delta\Omega_0}{\sigma}\right)\right)}. \quad (15)$$

By taking the limit $a \rightarrow 0$, the rescaled phenotype distribution becomes a Dirac distribution, dominating any convergences; thus, the left-hand side can be transformed as:

$$\lim_{a \rightarrow 0} \int_0^{+\infty} \frac{\text{sh}\left(\frac{1}{2}\left(\frac{\delta\Omega_0}{\sigma} - \frac{\delta\Omega}{\sigma}\right)\right)}{\alpha_0 e^{\frac{1}{2}\left(\frac{a}{\sigma}\right)\left(\frac{\delta\Omega}{\sigma}\right)} + \text{ch}\left(\frac{1}{2}\left(\frac{\delta\Omega_0}{\sigma} - \frac{\delta\Omega}{\sigma}\right)\right)} K e^{-\frac{1}{2}\left(\frac{\delta\Omega}{\sigma}\right)^2} d\left(\frac{\Omega}{a}\right) \sim \frac{\text{sh}\left(\frac{1}{2}\left(\frac{\delta\Omega_0}{\sigma}\right)\right)}{\alpha_0 + \text{ch}\left(\frac{1}{2}\left(\frac{\delta\Omega_0}{\sigma}\right)\right)} \sim \lim_{a \rightarrow 0} \frac{\text{sh}\left(\frac{1}{2}\left(\frac{a}{\sigma}\right)\left(\frac{\delta\Omega_0}{\sigma}\right)\right)}{\alpha_0 + \text{ch}\left(\frac{1}{2}\left(\frac{a}{\sigma}\right)\left(\frac{\delta\Omega_0}{\sigma}\right)\right)} \sim 0. \quad (16)$$

Therefore, small gene effects imply: $\delta\Omega_0/\sigma \ll 1$. Recalling that $\Delta\omega_0/\omega_0 = \text{th}\left(\left(\frac{a}{\sigma}\right)\left(\frac{\delta\Omega_0}{\sigma}\right)\right)$, one also deduces $\lim_{a \rightarrow 0} \text{th}\left(\left(\frac{a}{\sigma}\right)\left(\frac{\delta\Omega_0}{\sigma}\right)\right) \sim 0$, that is, small gene effects always involve common allele frequencies, namely $\Delta\omega_0/\omega_0 \rightarrow 0$. Using this result, the genetic paths difference can then be developed when gene effects are small and by assuming that $P_{\Omega}(\Omega_{1/N}) \ll 1$ and that ' $P_{\Omega}(\Omega)$ ' is normally distributed one obtains at the leading order:

$$\Delta\hat{\theta}(\Omega) \sim \frac{\omega_0}{2}\left(\frac{a}{\sigma}\right) \int_{\Omega_{1/N}}^{\Omega} -\left(\frac{\delta\Omega}{\sigma}\right) P_{\Omega}(\Omega) d\Omega + O(a^2) \sim \omega_0 \frac{a}{2} P_{\Omega}(\Omega) + O(a^2). \quad (17)$$

Equation (17) shows that in the context of Fisher's theory, a small gene effect corresponds to an overlapping symmetry between the genetic microstates and the phenotype distribution, with an amplitude proportional to the gene effect.

4.4. Fields linked to the variance of microstates

The involvement of variances in microstate distribution functions in genotype–phenotype associations is a highly debated matter (see [20] and references within). As mentioned above, the expression of the difference of the genetic paths considers the distribution density function as a whole, including the role of the microstate variances. In this context, we saw that equation (9b) provides a relation between variances in the form of an ellipse. Assuming a single variance for all microstate and phenotype distributions, as Fisher did, is plausible; but other solutions exist that would not violate equation (9b). In this context, let us imagine that the gene effect and dominance are nulls but that the distribution density function of microstates '+1', '0', '-1' and of the phenotype have distinct variances; written, respectively, as: $P_+(\Omega) \sim \frac{K}{\sigma_+} \exp\left(-\frac{1}{2}\left(\frac{\delta\Omega}{\sigma_+}\right)^2\right)$, $P_0(\Omega) \sim \frac{K}{\sigma_0} \exp\left(-\frac{1}{2}\left(\frac{\delta\Omega}{\sigma_0}\right)^2\right)$, $P_-(\Omega) \sim \frac{K}{\sigma_-} \exp\left(-\frac{1}{2}\left(\frac{\delta\Omega}{\sigma_-}\right)^2\right)$ and $P_{\Omega}(\Omega) \sim \frac{K}{\sigma} \exp\left(-\frac{1}{2}\left(\frac{\delta\Omega}{\sigma}\right)^2\right)$.

By noting $\lambda_+ = \sigma/\sigma_+$, $\lambda_0 = \sigma/\sigma_0$ and $\lambda_- = \sigma/\sigma_-$, the fields can be mapped under the form:

$$\Delta u(\Omega_0) = \frac{1}{2} \ln\left(\frac{P_-(\Omega_0)}{P_+(\Omega_0)}\right) = +\frac{1}{2}(\lambda_+^2 - \lambda_-^2) \left(\frac{\delta\Omega_0}{\sigma}\right)^2 + \ln\left(\frac{\lambda_+}{\lambda_-}\right) \quad (18a)$$

$$\Delta u(\Omega_0) - s\Delta u(\Omega) = \frac{1}{2} \ln\left(\frac{P_+(\Omega)}{P_-(\Omega)} \cdot \frac{P_-(\Omega_0)}{P_+(\Omega_0)}\right) = +\frac{1}{2}(\lambda_+^2 - \lambda_-^2) \left(\left(\frac{\delta\Omega_0}{\sigma}\right)^2 - \left(\frac{\delta\Omega}{\sigma}\right)^2\right) \quad (18b)$$

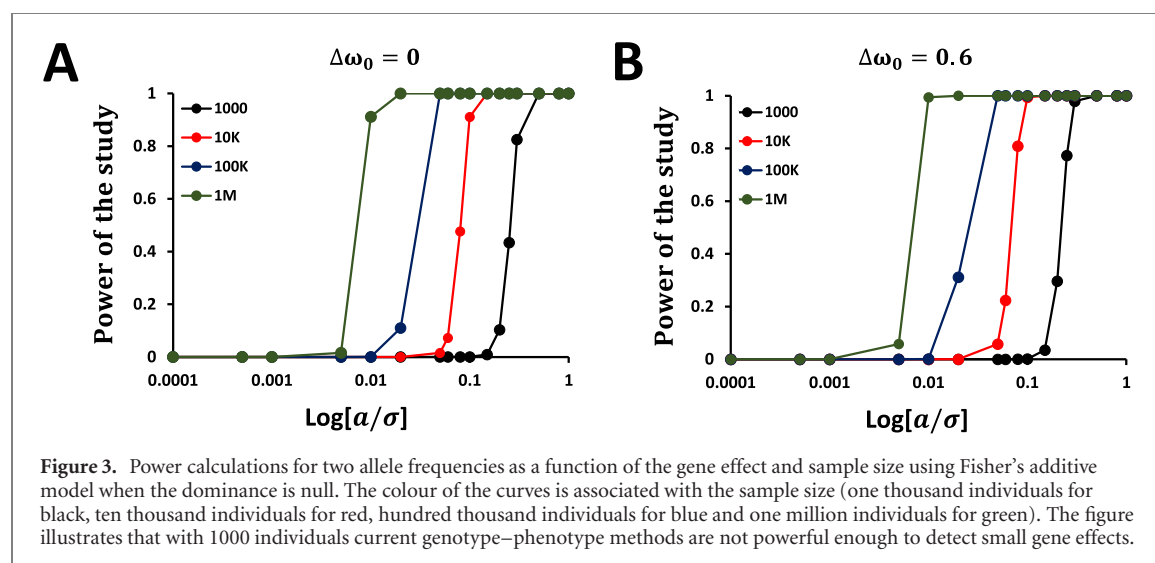
$$\bar{u}(\Omega) = \ln\left(\frac{P_0(\Omega)}{\sqrt{P_+(\Omega) \cdot P_-(\Omega)}}\right) \sim -\frac{1}{2}\left(\lambda_0^2 - \frac{1}{2}\lambda_+^2 - \frac{1}{2}\lambda_-^2\right) \left(\frac{\delta\Omega}{\sigma}\right)^2 + \ln\left(\frac{\lambda_+ \lambda_-}{\lambda_0}\right). \quad (18c)$$

Consequently, pseudo-gene effect and pseudo-dominance linked to the variances of genetic microstates can be defined, respectively, as: $a' = \frac{1}{2}(\lambda_+^2 - \lambda_-^2)$ and $d' = \frac{1}{2}(\lambda_0^2 - \frac{1}{2}\lambda_+^2 - \frac{1}{2}\lambda_-^2)$.

To conclude, as this new method does not only concentrate on average values, it captures more information as far as genotype–phenotype associations are involved.

Table 1. Genetic variances as a function of the gene effect and allele frequency.

		Allele frequency (p)					Genetic variances
		0.1	0.2	0.3	0.4	0.5	
Gene effect in unit of phenotypic standard deviations (a/σ)	1	0.18	0.32	0.42	0.48	0.5	
	0.8	0.1152	0.2048	0.2688	0.3072	0.32	
	0.5	0.045	0.08	0.105	0.12	0.125	
	0.3	0.0162	0.0288	0.0378	0.0432	0.045	
	0.15	0.00405	0.0072	0.00945	0.0108	0.01125	
	0.1	0.0018	0.0032	0.0042	0.0048	0.005	
	0.05	0.00045	0.0008	0.00105	0.0012	0.00125	
	0.01	0.000018	0.000032	0.000042	0.000048	0.00005	
	0.005	4.5×10^{-6}	0.000008	1.05×10^{-5}	0.000012	1.25×10^{-5}	
	0.001	1.8×10^{-7}	3.2×10^{-7}	4.2×10^{-7}	4.8×10^{-7}	5×10^{-7}	
	0.0005	4.5×10^{-8}	8×10^{-8}	1.05×10^{-7}	1.2×10^{-7}	1.25×10^{-7}	



5. Illustration of the application of GIFT using simulated data

We intend to illustrate how GIFT can be applied using simulated data and qualitatively assess its sensitivity to extract information.

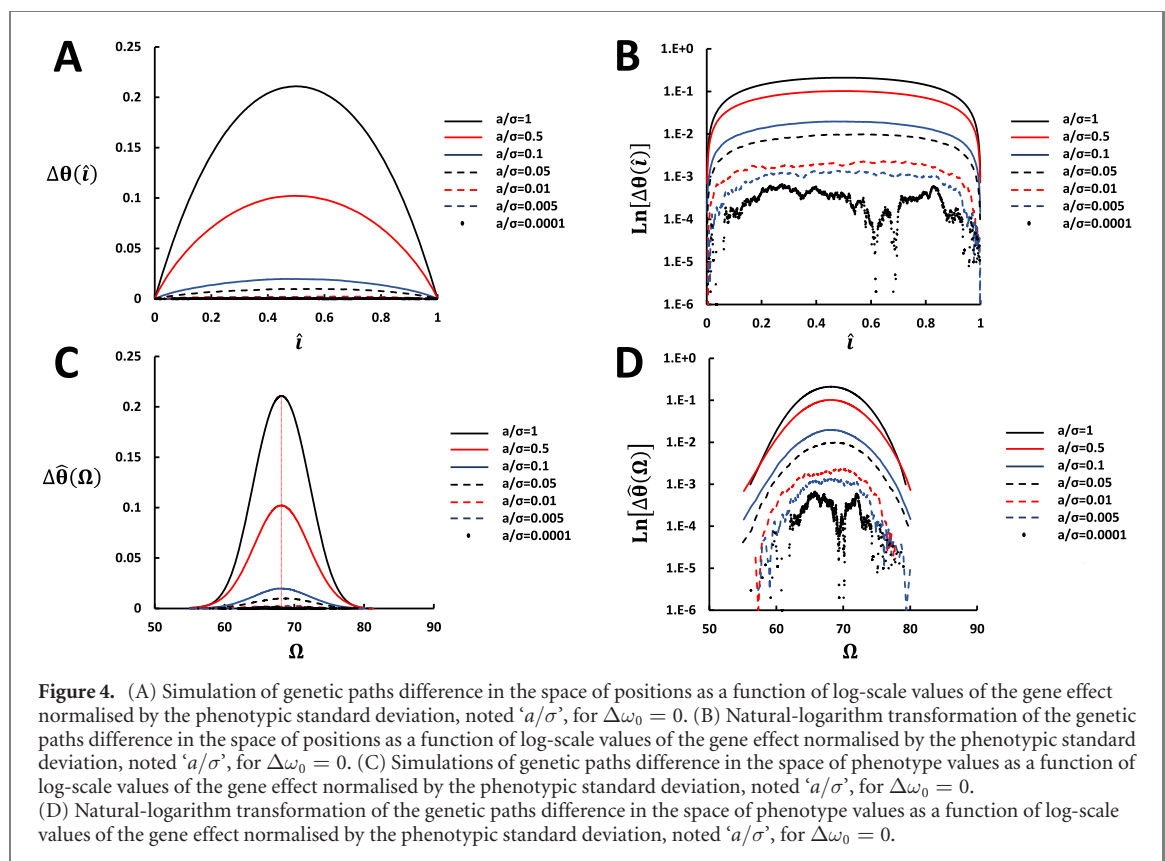
5.1. Data simulations

The codes used are provided in SM5, see the supplementary materials.

Data was simulated according to quantitative genetic models defined by Falconer and Mackay (1996) [27]. A single bi-allelic quantitative trait locus associated with a continuous phenotype was modelled, with an additive allele effect, a , and allele frequencies, p and q , where $p + q = 1$. The simulation parameters were set as the number of individuals sampled, $N = 1000$; number of simulation replicates, $n = 1000$; allele frequency, p ; additive allele effect, a , and dominance, d ; note that the number of simulation replicates allows one to determine the best outcomes. While the theory provided in this paper is general, the simulation of data will be restricted to individuals' genotypes allocated according to Hardy–Weinberg proportions. For N individuals, Np^2 had genotype

AA (corresponding to microstate -1), $2pqN$ had genotype Aa (microstate 0), and Nq^2 had genotype aa (microstate $+1$). The allele effect, a , is defined as half the difference between the $+1$ and -1 genotype (microstate) means, and d , is the position of the 0 genotype (microstate) mean (figure 1). Dominance is measured as the deviation of the mean of microstate 0 from the midpoint between the means of the $+1$ and -1 microstates. For the purposes of the simulation dominance, d , was 0 , that is, the mean of microstate 0 was mid-way between the mean of microstates $+1$ and -1 .

The additive genetic variance due to the quantitative trait loci (σ_{QTL}^2) was defined as [27]: $\sigma_{QTL}^2 = 2pq[a + d(q - p)]^2$. Each individual was assigned a genotypic value, depending on their microstate: $-a$ for the $+1$ microstate, 0 for the 0 microstate, and $+a$ for the -1 microstate. Individual phenotypes were generated by adding a random environmental effect to the genotypic value of each individual. The added environmental effect was a random variate drawn from a normal distribution with a mean of 0 and variance of $1 - \sigma_{QTL}^2$. The phenotype was then rescaled to a value representing a realistic dataset: phenotype = (simulated phenotype \times standard deviation of real



data) + mean of real data. In this case, the real dataset modelled was a summary of the genotype-tissue expression (GTEx) project [28]. In particular, the phenotype was height with a mean of 68.17206 inches and standard deviation of 4.030 07 inches. For each simulated replicate of N individuals, the difference between the cumulative sums of microstates ordered by phenotype value and genotypes in a randomised order with respect to phenotype was determined to create the difference in the genetic paths difference. The maximum value of this difference was identified and its position and phenotypic value in the ordered string of microstates were recorded. Where the maximum value extended over several positions, the mean position and phenotypic value were recorded. Finally, to simplify representation, the amplitudes of the genetic path differences were normalised by population size ($N = 1000$ in this case).

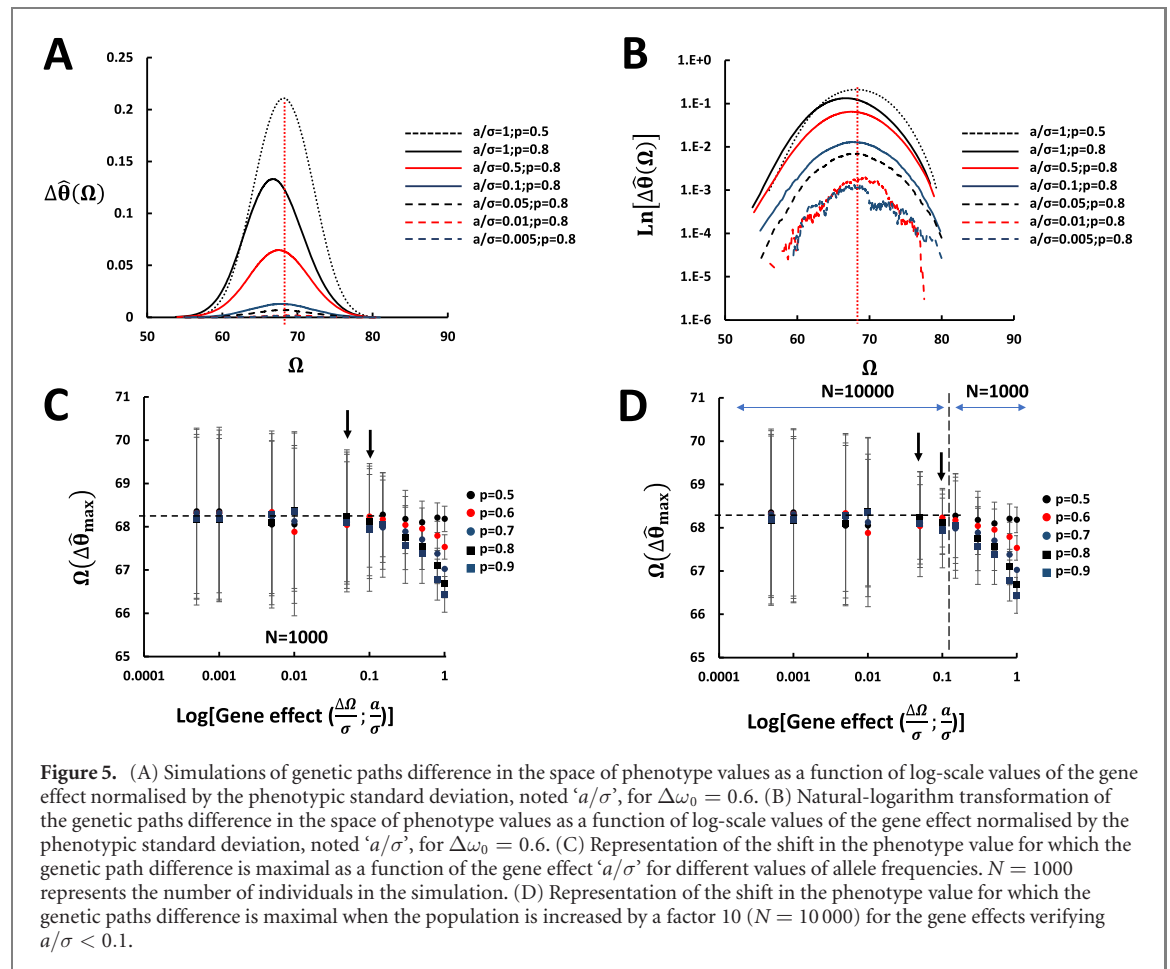
Note that the standard deviation(s) arising from genotype-phenotype simulations were not considered in the analysis that follows. Instead, we report a theoretical analysis of the convergence of the genetic path difference method, and its self-consistency, as well as its sensitivity to detect genotype-phenotype associations using simulations, in SM6 and SM7, respectively, see supplementary materials.

5.2. Analysis of simulated results

For information, table 1 shows how genetic variance, gene effect, i.e., a/σ , and allele frequency are numerically related using GWAS method. Similarly

figures 3(A) and (B) represent, in the context of GWAS and for the allele frequencies $p = 0.5$ ($\Delta\omega_0 = 0$) and $p = 0.8$ ($\Delta\omega_0 = 0.6$) that will be used below as examples, the relationship between the power of the study, the gene effect and the sample size as described in [29]. Briefly, the power of a study is related to the concepts of type I and type II errors. A type I error (a.k.a. α) is rejecting the null hypothesis in favour of a false alternative hypothesis, and a type II error (a.k.a. β) is failing to reject a false null hypothesis in favour of a true alternative hypothesis. The power of a study is then the probability of avoiding a type II error. Mathematically, the power is defined by, $1 - \beta$, where $0 < \beta < 1$. If the power is close to 1, i.e., $\beta \sim 0$, the hypothesis test is very good at detecting a false null hypothesis. β is commonly set at 0.2, to provide a power ~ 0.8 (or 80%). Powers lower than 0.8, while not impossible, would typically be considered too low for GWAS. The four primary factors affecting power are, the sample size, the significance level (or α), the variance/variability in the measured response and the magnitude of the effect of the variable. Only the first variable can be altered in a study since all the others are fixed by the genes. To conclude, power is increased when the sample size or effect sizes (gene effect) are increased. Accordingly, figure 3 demonstrates that 1000 individuals would not allow 80% power to be achieved unless the gene effect is sufficiently large, that is for $a/\sigma \geq 0.5$.

Using simulated data, we can now represent the genetic paths difference and its log transformation for



$\Delta\omega_0 = 0$, either in the space of positions (figures 4(A) and (B)) or of phenotypic values (figures 4(C) and (D)) for different log-scale values of the normalised gene effect a/σ (see inset).

As shown in figure 4(C), the profile of the phenotype distribution density function is recovered with an amplitude that decreases as a/σ decreases. The red vertical dashed line in figure 4(C) represents the mean phenotypic value. Using the natural logarithm to transform $\Delta\hat{\theta}(\Omega)$ (figure 4(C)) to $\ln[\Delta\hat{\theta}(\Omega)]$ (figure 4(D)) demonstrates that a difference between genetic paths can be seen for small gene effects.

One may then compare how perceptible the associations are using the new method by comparing figures 1(B) and (C) (method of averages based on Fisher's theory) and figure 4(D) for identical allele frequencies and similar gene effects. Recall that GWAS rely on determining difference in averages (see figure 1(B) or figure 1(C)). However, the determination of a difference in the microstate averages rely on a strong gene effect (figure 1(B)) or a very large population (see figure 3) as otherwise the density functions of microstates collapse onto one. This is particularly visible when one compares the right-hand and left-hand graphs in figure 1(B) or figure 1(C). Thus, the results provided by figure 4 suggest that GIFT can

be applied to 1000 individuals to return information regarding potential genotype–phenotype associations that would not otherwise be possible, or extremely difficult, with current association studies.

Concentrating on different allele frequencies given by $p = 0.8$ ($\Delta\omega_0 = 0.6$) as an example. Figures 5(A) and (B) are representations of $\Delta\hat{\theta}(\Omega)$ and its natural-log transformation for log-scale values of a/σ .

Differences are clearly visible between $\Delta\omega_0 = 0$ and $\Delta\omega_0 = 0.6$, since the phenotype values for which $\Delta\hat{\theta}(\Omega)$'s are maximal have been shifted from the average value of the phenotype (indicated by the vertical red dashed line). This is not surprising because the simulation only imposed a set of genetic variances, without any constraint on the conservation of the average phenotypic value.

However, the shift of the phenotype value for which $\Delta\hat{\theta}(\Omega)$ is maximal is of interest. As equation (17) demonstrates that for small gene effects, the genetic path difference should be proportional to the phenotype distribution, that is, the phenotype value for which $\Delta\hat{\theta}(\Omega)$ is extreme should be the average value of the phenotype.

Thus, to obtain a better visualisation of the impact of the gene effect on the positioning of the phenotype value $\Omega(\Delta\hat{\theta}_{\max})$ for which the genetic path difference is maximal, a set of simulations were also performed based on allele frequencies

Table 2. Nonlinear fit results for figure 4(C) using $\text{Ln}[\Delta\hat{\theta}(\Omega)] = A\Omega^2 + B\Omega + C$.

	$a/\sigma = 1$ $p = 0.5$	$a/\sigma = 0.5$ $p = 0.5$	$a/\sigma = 0.05$ $p = 0.5$	$a/\sigma = 0.005$ $p = 0.5$
A	−0.0371	−0.0321	−0.0311	−0.0435
B	5.0539	4.3799	4.2425	5.9241
C	−173.69	−151.48	−149.5	−208.35
r^2	0.9960	0.9990	0.9937	0.8329
$\bar{\Omega}$	68.11	68.22	68.21	68.09
σ^2	13.48	15.58	16.08	11.49

defined by $p \in \{0.5; 0.4; 0.3; 0.2; 0.1\}$ for log-ranging values of gene effects (figure 5(C)). Note that $p \in \{0.5; 0.6; 0.7; 0.8; 0.9\}$ can be deduced from the symmetry around the average value of the phenotype.

While the standard deviations obtained were not always negligible concerning $\Omega(\Delta\hat{\theta}_{\max})$, typically between 0.5 and 1 phenotypic standard deviation for small gene effects; figure 5(A) demonstrated trends toward the average value of the phenotype with small gene effects. Indeed, below the simulated gene effect of $a/\sigma \sim 10^{-1}$, the average value of $\Omega(\Delta\hat{\theta}_{\max})$ was remarkably similar to that of the average value of the phenotype, marked by the horizontal black dashed line.

To confirm this trend for small gene effects, that is, $a/\sigma \leq 0.1$, we varied the population size from $N = 10^3$ to $N = 10^4$ to determine the presence of potential variations in $\Omega(\Delta\hat{\theta}_{\max})$ linked to the simulations. Results summarised in figure 5(B) demonstrate that the only difference was a reduction in the standard deviations obtained for $\Omega(\Delta\hat{\theta}_{\max})$ for the simulated gene effects comprised between 0.01 and 0.1 (see arrows figures 5(C) and (D) pointing to different magnitude of the standard deviations). Namely, the initial symmetry of the phenotype distribution density function reappears, as expected (equation (17)).

Finally, equation (17) suggests that for small gene effects $\Delta\hat{\theta}(\Omega)$ is proportional to the gene effect a/σ in the form, $\Delta\hat{\theta}(\Omega) \sim \omega_0 \frac{a}{\sigma} K e^{-\frac{1}{2\sigma^2}(\Omega - \langle\Omega\rangle)^2}$. As a consequence it was decided to fit the all the curves $\Delta\hat{\theta}(\Omega)$ in figures 4(D) and 5(B) with quadratic equations of the form $\text{Ln}[\Delta\hat{\theta}(\Omega)] \sim A\Omega^2 + B\Omega + C$; see tables 2 and 3. Then, as $\Delta\hat{\theta}(\Omega) \sim e^{(C - \frac{B^2}{2A})} e^{A(\Omega + \frac{B}{2A})^2}$ with such a fit, we expect by identification of equation (17) that for small gene effects: $\bar{\Omega} \sim -\frac{B}{2A}$, $-\frac{1}{2\sigma^2} \sim A$ and $\frac{a}{\sigma} \sim \frac{1}{\omega_0 K} e^{(C - \frac{B^2}{2A})}$. Tables 2 and 3 provide the estimations for both $\bar{\Omega}$ and σ and setting $K \sim 1/\sqrt{2\pi}$, figure 6 provides a comparison between the gene effect from the simulations, $(\frac{a}{\sigma})_{\text{sim}}$, and the gene effect deduced from equation (17), $(\frac{a}{\sigma})_{\text{theo}}$.

Thus, recalling that the phenotype average and variance of the population modelled are, respectively, 68.17 inch and 16.24 inch²; tables 2 and 3 demonstrate that fitting the genetic paths difference as a function

Table 3. Nonlinear fit results for figure 4(F) using $\text{Ln}[\Delta\hat{\theta}(\Omega)] = A\Omega^2 + B\Omega + C$.

	$a/\sigma = 1$ $p = 0.2$	$a/\sigma = 0.5$ $p = 0.8$	$a/\sigma = 0.05$ $p = 0.2$	$a/\sigma = 0.005$ $p = 0.2$
A	−0.0318	−0.0292	−0.0282	−0.0300
B	4.2666	3.9583	3.8514	4.0741
C	−145.14	−136.71	−136.45	−145.40
r^2	0.9135	0.8892	0.8004	0.7554
$\bar{\Omega}$	67.08	67.77	68.28	67.90
σ^2	15.72	17.12	17.73	16.67

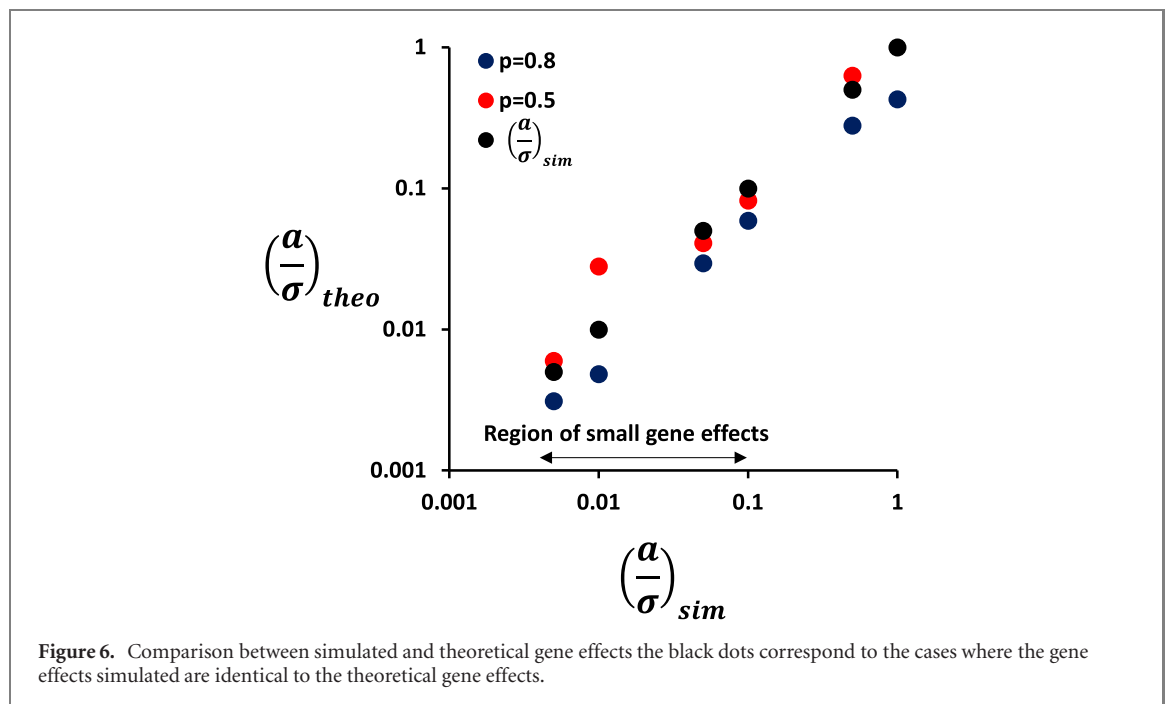
of phenotype values with a quadratic curve recovers the magnitude of the average and variance of the phenotype used for the simulations for most log-scale values of the gene effect. Furthermore, the amplitude of $\Delta\hat{\theta}(\Omega)$ is also indicative of the gene effects involved.

6. Discussion

In his seminal paper [4], Fisher provided a synthesis between the genetic inheritance of continuous traits and the Mendelian scheme of inheritance using statistics and probability. His theory has become a landmark in genetics and heredity and its conceptual framework is still used today. While statistics is a natural field to employ when dealing with large datasets, the interpretation of data as well as the inferences that can be drawn from it rely fundamentally on PDFs. As the act of creating categories to work with density functions is acknowledging imprecisions, our aim was to devise a different method ruling out the need for category. This new theoretical method, inspired by physics and named GIFT, uses the concept of phenotypic fields, and concentrate on ‘genetic paths’ to extract information on genotype phenotype mappings. It is then important to discuss the conceptual similarities and differences between GWAS and GIFT.

In term of conceptual similarities, we saw that the theory underscoring GIFT developed using Fisher’s assumptions recovers key concepts from quantitative genetics, including: (i) the Hardy–Weinberg coefficient locally, (ii) the Hardy–Weinberg coefficient at the population level, (iii) the gene effect, (iv) dominance, and (v) small gene effects involving common allele frequencies [30]. In this context, GIFT and GWAS are similar. Finally, applying GIFT to simulated data based on Fisher’s assumption proved its sensitivity for extracting information on genotype–phenotype associations when sample sizes and gene effects were small. The reason for not considering the dominance in the simulations is linked to the fact that realistic GWAS have shown that with small effect sizes/small gene effect (which is the main area of concern of the current paper), dominance effects are often too small, and an additive model as suggested by Fisher works well enough [31].

In term of conceptual differences, three essential points can be discussed.



First of all, GIFT is more general than GWAS in the sense that the phenotypic fields can be any, namely do not have to be linear. The prescription of linear phenotypic fields in Fisher's context comes from the symmetry associated with using the normal distribution as a template for any distribution density function, together with the assumption that the phenotypic and microstate variances are identical [4]. When the constraint on the variances is released, the phenotypic fields become quadratic involving the variances as well as the averages. In this context GIFT has enabled us to define new parameters linked to microstate variances, that are, the pseudo-gene effect and the pseudo-dominance, which will probably help resolve controversies [20].

Secondly, in terms of genetics what has been achieved so far is rather at odds with traditional ways of thinking about the notion of gene. Indeed, by defining the difference in genetic paths, $\Delta\theta(\Omega)$, one can say that it is the phenotype, i.e., phenotypic fields or information, that organises the configuration of genotypes and not the converse. In genetics the tradition is to think of genes as causing phenotypes. Here, a different way of thinking is suggested since it is the variation in phenotype values, resulting in our ability to generate a ranking process, that interacts with the microstates. Therefore, the phenotype is able to 'select' a set of genetic microstates. Recall that microstates 'respond' to, or interact with, the phenotypic fields only if they are associated with the phenotype. Consequently, this model suggests considering a genotype–phenotype 'loop', a.k.a. self-consistency. That is to say that if genes cause phenotypes (traditional view) and that phenotypes select gene microstates (present view), then an equivalence exists between phenotype and genotype.

Supplementary materials contain more information and development concerning the convergence and self-consistency of GIFT (SM6).

Finally, Fisher's theory/GWAS has been built on considering the normal distribution. In general, 'real' density functions never come as normally distributed. Given that Fisher's theory gives biological meanings to average and variance only, to define the 'gene effect' and 'genetic/phenotypic variance' linked to heredity, respectively; there is no biological meaning to any other statistical/mathematical parameters describing real density functions, such as for example the 'skewness'. As GIFT uses curves, namely does not use average and variance as central parameters, this issue does not exist with GIFT. Said differently, GIFT frees GWAS from any preconceived idea of what statistics and probability applied to biology should be.

Taken as a whole, the work presented here is a first step suggesting that GIFT can be considered as a potential method for genotype–phenotype mappings. Supplementary material SM7 contains more information concerning the signal-to-noise ratio when GIFT is used and SM8 (see supplementary materials) provides an initial illustration of the application of GIFT using real data based on GWAS results.

However the authors agree with the fact that more work needs to be done to compare GIFT to the vast literature concerning GWASes. For example, at present the model is quite simplistic in the way that, by construction, it does not allow the easy incorporation of covariates. Future works will relate covariates, such as age or sex, and for the case of human populations, genetic principal components (to account for population structure). In addition, GIFT will be compared against well-known statistical tests as used

in GWAS (e.g., z-test/chi-square as parametric test, or the Kolmogorov–Smirnov test as non-parametric).

7. Conclusion

A century ago, Fisher devised a statistical method to map genotypes and phenotypes, which was essentially based on the measure of uncertainty. We present here a method taking as a paradigm the fact that certainty can exist with the possibility to measure phenotype and genotype with very high precision. In an associated paper, we present a theoretical methodology based on Shannon's information enabling the significance of correlation using real genotype–phenotype data to be quantified [32]. To conclude, this new method (GIFT) opens a new way to analyse genotype–phenotype mapping.

Authors' contributions

Conceptualisation of GIFT as a new method based on physics field theory: CR; theoretical developments: CR, JW; simulation of genotype and phenotype data based on Fisher's theory: CR, PK, SaBl; analysis of simulated or real data using GIFT: CR, SB, PK.

Conflict of interest

We have no competing interests.

Funding

This work was supported by the Physics of Life Networks and The University of Nottingham.

Acknowledgments

We thank Professors Andras Paldi, James Leigh, Patricia Harris, and Darren Logan for fruitful discussions concerning the urgent need to improve genotype–phenotype association methods for small sample sizes.

Data availability statement

Code for the simulations presented in the manuscript can be found in the supplementary materials (see SM5). Genomic and phenotypic data analysed as part of this study in supplementary materials (see SM8) is already published and freely available online. The genomic data is part of the well-known *Arabidopsis thaliana* '1001 genomes' project (<https://doi.org/10.1186/gb-2009-10-5-107>) and can be found online here: [https://1001genomes.org/data-](https://1001genomes.org/data-center.html)

[center.html](https://1001genomes.org/data-center.html). The phenotype data was published in The Plant Journal (DOI: [10.1111/tjp.15177](https://doi.org/10.1111/tjp.15177)) and is freely available online via Ion Explorer here: https://bitbucket.org/ADAC_UoN/dr000081-web-service-ionome-seed-and-leaf-map/. The code used for the real-data analysis in SM8 in the supplementary materials can be found on Dryad, DOI: <https://doi.org/10.5061/dryad.vx0k6djtp>.

Appendix A

A.1. Differential expression of the genetic path in the space of phenotypic values

We assume that the phenotype values are measured precisely enough such that each individual has a unique phenotype value noted Ω_i . As the population is composed of 'N' individuals one defines, $\hat{i} \stackrel{\text{def}}{=} i/N$ and $\Omega_{\hat{i}}$, as the new position and its corresponding phenotype value, respectively. Thus $\Omega_{1/N}$ and Ω_1 are the smallest and largest phenotype values, respectively.

As a result, the cumulative sum of presence probability of genetic microstates as a function of ' \hat{i} ' can be written as: $\theta(\hat{i}) = \sum_{j=1}^{\hat{i}} (\omega_+(j) - \omega_-(j))$, where $\omega_+(j)$ and $\omega_-(j)$ are the presence probabilities of microstates '+1' and '-1' at the position j . Using the continuum limit, one deduces also, $\theta(\hat{i}) \sim \int_{1/N}^{\hat{i}} (\omega_+(j) - \omega_-(j)) dj$.

As the ranking of phenotype values was introduced to define, $\theta(\hat{i})$, the genetic path can also be expressed as a function of phenotype values under the form: $\hat{\theta}(\Omega_{\hat{i}}) = \sum_{j=1/N}^{\hat{i}} (\hat{\omega}_+(\Omega_j) - \hat{\omega}_-(\Omega_j))$ where the hat $\hat{\cdot}$ is added to describe new functions linked to phenotypic values. Accordingly, by definition, $\hat{\theta}(\Omega_{\hat{i}}) - \hat{\theta}(\Omega_{\hat{i}-1}) = \hat{\omega}_+(\Omega_{\hat{i}}) - \hat{\omega}_-(\Omega_{\hat{i}})$. Assuming that the measured phenotype values are sufficiently close, the continuum limit can be used to determine $\hat{\theta}(\Omega_{\hat{i}}) - \hat{\theta}(\Omega_{\hat{i}-1}) \sim \frac{d\hat{\theta}}{d\Omega} \cdot \Delta(\Omega)$, where $\Delta(\Omega)$ is the spacing between two consecutive individuals in the space of phenotype values. One deduces then: $\frac{d\hat{\theta}}{d\Omega} \sim (\hat{\omega}_+(\Omega) - \hat{\omega}_-(\Omega)) \frac{1}{\Delta(\Omega)}$. As a result, the genetic path can be expressed in the space of the phenotype values as:

$$\hat{\theta}(\Omega) \sim \int_{\Omega_{1/N}}^{\Omega} (\hat{\omega}_+(x) - \hat{\omega}_-(x)) \frac{1}{\Delta(x)} dx. \quad (\text{A.1})$$

Where $\int_{\Omega_{1/N}}^{\Omega_1} \Delta(\Omega) d\Omega = \Omega_1 - \Omega_{1/N} = \Delta\Omega$. As $\Delta(\Omega)$ is the phenotypic space between individual, provided that the population is large and dense enough one can relate ' $\Delta(\Omega)$ ' to the phenotype distribution density function, $P_{\Omega}(\Omega)$, under the form $1/\Delta(\Omega) \sim P_{\Omega}(\Omega)$ (SM1 in the supplementary materials). The different

elements that lead to the formation of a genetic path can now be addressed.

A.2. Entropy of the string of microstates

Keeping the notations ω_+^0 , ω_0^0 and ω_-^0 , for the genetic microstate frequencies of a given genome position across the population of individuals, we aim to determine the expressions of $\omega_+(\hat{i})$, $\omega_0(\hat{i})$ and $\omega_-(\hat{i})$ given the information obtained upon ordering the genotypes as a function of phenotype values along the \hat{i} -axis.

The default genetic path, θ_0 , as a straight line is defined by the absence of information on phenotype values, which is similar to an absence of association between the genetic microstates and the phenotype values, leading to an apparent disordering of genetic microstates. One way to measure this disordering is by using the ‘entropy’ of the string of genetic microstates for the genome position considered. In this context, the entropy is given by calculating the number of possible combinations of placing $N_+ = N\omega_+^0$, $N_0 = N\omega_0^0$ and $N_- = N\omega_-^0$ genetic microstates over ‘ N ’ possible positions. Consequently, the entropy of the default genetic path is $S_0 = N!/N_+!N_0!N_-!$; and for ‘ N ’ large enough using Stirling’s formula one deduces: $S_0/N \sim -\omega_+^0 \ln(\omega_+^0) - \omega_0^0 \ln(\omega_0^0) - \omega_-^0 \ln(\omega_-^0)$, that can be rewritten in the continuum limit as: $S_0 \sim -N \int_{\Omega_{1/N}}^{\Omega_1} (\omega_+^0 \ln(\omega_+^0) + \omega_0^0 \ln(\omega_0^0) + \omega_-^0 \ln(\omega_-^0)) d\Omega$. Note that, as the genetic microstate frequencies are constant in this case, the entropy can be rewritten using the phenotype values as,

$$S_0 \sim -N \int_{\Omega_{1/N}}^{\Omega_1} (\omega_+^0 \ln(\omega_+^0) + \omega_0^0 \ln(\omega_0^0) + \omega_-^0 \ln(\omega_-^0)) \frac{1}{\Delta(\Omega)} d\Omega. \quad (\text{A.2})$$

When information about phenotype values and their ranking is given and when the genome position considered is associated with the phenotype, S_0 is transformed to S where:

$$S = -N \int_{\Omega_{1/N}}^{\Omega_1} (\hat{\omega}_+(\Omega) \ln(\hat{\omega}_+(\Omega)) + \hat{\omega}_0(\Omega) \ln(\hat{\omega}_0(\Omega)) + \hat{\omega}_-(\Omega) \ln(\hat{\omega}_-(\Omega))) \frac{1}{\Delta(\Omega)} d\Omega. \quad (\text{A.3})$$

As a result, the entropy difference, $S - S_0$ when non-null provides information on whether the genome position is associated with phenotype values. Thus the difference, $S - S_0$, can be thought of as a ‘transformation’ in a physical/thermodynamic sense. That is, the difference in entropies must be balanced by a term that is linked to the association (or interaction) between the genetic microstates and the phenotype values.

A.3. Interaction energy between microstates and subfields

As the difference, $S - S_0$, is linked to the information gained from knowing phenotypic values and ranking them (appendix A.2), given the existence of three distinct genetic microstates, one can define three distinct functions a.k.a. phenotypic fields ‘ $u_+(\Omega)$ ’, ‘ $u_0(\Omega)$ ’ and ‘ $u_-(\Omega)$ ’ that are fundamentally related to changes in the phenotype-associated genetic path. In this context, the entire genetic path can be defined with a function representing the sum of interactions between each of the genetic microstates and phenotypic fields under the form:

$$E \sim N \int_{\Omega_{1/N}}^{\Omega_1} (u_+(\Omega) \hat{\omega}_+(\Omega) + u_0(\Omega) \hat{\omega}_0(\Omega) + u_-(\Omega) \hat{\omega}_-(\Omega)) \frac{1}{\Delta(\Omega)} d\Omega. \quad (\text{A.4})$$

In this context, one may consider that the set of microstates changes the configuration because the fields are ‘switch on’. This implies that for the genome positions that are not involved in the formation of the phenotype considered, the switch does not work, that is, the fields are null. In this context, one can consider the equivalence, $S - S_0 \sim E$. As a result, the relationship to optimise is: $\Delta S - E = 0$.

A.4. Optimisation of $\Delta S - E$

Recalling the conservation of genetic microstates for the genome position considered: $\int_{\Omega_{1/N}}^{\Omega_1} \hat{\omega}_+(\Omega) \frac{1}{\Delta(\Omega)} d\Omega = \omega_+^0$, $\int_{\Omega_{1/N}}^{\Omega_1} \hat{\omega}_0(\Omega) \frac{1}{\Delta(\Omega)} d\Omega = \omega_0^0$ and $\int_{\Omega_{1/N}}^{\Omega_1} \hat{\omega}_-(\Omega) \frac{1}{\Delta(\Omega)} d\Omega = \omega_-^0$ together with the conservation of probability, $\hat{\omega}_+(\Omega) + \hat{\omega}_0(\Omega) + \hat{\omega}_-(\Omega) = 1$, Euler–Lagrange’s method can then be used to determine the optimal configuration for $\hat{\omega}_+(\Omega)$, $\hat{\omega}_0(\Omega)$ and $\hat{\omega}_-(\Omega)$ in a context where the phenotypic fields are imposed. By defining α_+ , α_0 and α_- , the Lagrange multipliers for the conservation of genetic microstates, the relation to optimise with regard to the genetic microstate frequencies $\hat{\omega}_+(\Omega)$, $\hat{\omega}_0(\Omega)$ and $\hat{\omega}_-(\Omega)$ is then,

$$\begin{aligned} \Delta S/N - E/N - \alpha_+ \left(\omega_+^0 - \int_{\Omega_{1/N}}^{\Omega_1} \hat{\omega}_+(\Omega) \frac{1}{\Delta(\Omega)} d\Omega \right) \\ - \alpha_0 \left(\omega_0^0 - \int_{\Omega_{1/N}}^{\Omega_1} \hat{\omega}_0(\Omega) \frac{1}{\Delta(\Omega)} d\Omega \right) \\ - \alpha_- \left(\omega_-^0 - \int_{\Omega_{1/N}}^{\Omega_1} \hat{\omega}_-(\Omega) \frac{1}{\Delta(\Omega)} d\Omega \right) = 0. \quad (\text{A.5}) \end{aligned}$$

Using the conservation of genetic microstate frequencies, $\hat{\omega}_0(\Omega)$, can be replaced by $1 - \hat{\omega}_+(\Omega) - \hat{\omega}_-(\Omega)$, and a variational calculus can be performed on the genetic microstate frequencies, leading to two conditions:

$$\delta\hat{\omega}_+(\Omega) \left[\ln \left(\frac{\hat{\omega}_+(\Omega)}{1 - \hat{\omega}_+(\Omega) - \hat{\omega}_-(\Omega)} \right) - \delta u_+(\Omega) + (\alpha_+ - \alpha_0) \right] = 0 \quad (\text{A.6})$$

$$\delta\hat{\omega}_-(\Omega) \left[\ln \left(\frac{\hat{\omega}_-(\Omega)}{1 - \hat{\omega}_+(\Omega) - \hat{\omega}_-(\Omega)} \right) - \delta u_-(\Omega) + (\alpha_- - \alpha_0) \right] = 0. \quad (\text{A.7})$$

Where $\delta\hat{\omega}_+(\Omega)$ and $\delta\hat{\omega}_-(\Omega)$ are small variations in the presence probabilities of microstates '+1' and '-1', and $\delta u_+(\Omega) \stackrel{\text{def}}{=} u_+(\Omega) - u_0(\Omega)$, $\delta u_-(\Omega) \stackrel{\text{def}}{=} u_-(\Omega) - u_0(\Omega)$. Finally, $\hat{\omega}_0(\Omega)$ can be deduced using $\hat{\omega}_0(\Omega) = 1 - \hat{\omega}_+(\Omega) - \hat{\omega}_-(\Omega)$. Using the conditions $\hat{\omega}_+(\Omega) = \omega_+^0$, $\hat{\omega}_0(\Omega) = \omega_0^0$ and $\hat{\omega}_-(\Omega) = \omega_-^0$ when the fields are null, one obtains,

$$\hat{\omega}_+(\Omega) = \frac{\omega_+^0 e^{-\delta u_+(\Omega)}}{\omega_+^0 + \omega_+^0 e^{-\delta u_+(\Omega)} + \omega_-^0 e^{-\delta u_-(\Omega)}} \quad (\text{A.8})$$

$$\hat{\omega}_-(\Omega) = \frac{\omega_-^0 e^{-\delta u_-(\Omega)}}{\omega_+^0 + \omega_+^0 e^{-\delta u_+(\Omega)} + \omega_-^0 e^{-\delta u_-(\Omega)}}. \quad (\text{A.9})$$

To make the asymmetries of the problem more apparent, the following are defined for genetic microstates: $\Delta\omega_0 \stackrel{\text{def}}{=} \omega_+^0 - \omega_-^0$ and $\omega_0 \stackrel{\text{def}}{=} \omega_+^0 + \omega_-^0 = 1 - \omega_0^0$; and for the phenotypic fields:

$2\bar{u}(\Omega) \stackrel{\text{def}}{=} \delta u_+(\Omega) + \delta u_-(\Omega) = u_+(\Omega) + u_-(\Omega) - u_0(\Omega)$ and $2\Delta u(\Omega) \stackrel{\text{def}}{=} \delta u_+(\Omega) - \delta u_-(\Omega) = u_+(\Omega) - u_-(\Omega)$. Then, the difference and sum of $\hat{\omega}_+(\Omega)$ and $\hat{\omega}_-(\Omega)$ can be rewritten as follow:

$$\begin{aligned} \hat{\omega}_+(\Omega) - \hat{\omega}_-(\Omega) &= \frac{\omega_0 e^{-\bar{u}(\Omega)} \left[\frac{\Delta\omega_0}{\omega_0} \text{ch}(\Delta u(\Omega)) - \text{sh}(\Delta u(\Omega)) \right]}{1 - \omega_0 + \omega_0 e^{-\bar{u}(\Omega)} \left[\text{ch}(\Delta u(\Omega)) - \frac{\Delta\omega_0}{\omega_0} \text{sh}(\Delta u(\Omega)) \right]} \end{aligned} \quad (\text{A.10})$$

$$\begin{aligned} \hat{\omega}_+(\Omega) + \hat{\omega}_-(\Omega) &= \frac{\omega_0 e^{-\bar{u}(\Omega)} \left[\text{ch}(\Delta u(\Omega)) - \frac{\Delta\omega_0}{\omega_0} \text{sh}(\Delta u(\Omega)) \right]}{1 - \omega_0 + \omega_0 e^{-\bar{u}(\Omega)} \left[\text{ch}(\Delta u(\Omega)) - \frac{\Delta\omega_0}{\omega_0} \text{sh}(\Delta u(\Omega)) \right]} \end{aligned} \quad (\text{A.11})$$

Noting that: $-1 \leq \Delta\omega_0/\omega_0 \leq +1$, a new phenotype value is defined and noted ' Ω_0 ' by setting $\text{th}(\Delta u(\Omega_0)) \stackrel{\text{def}}{=} \frac{\Delta\omega_0}{\omega_0}$. Then, the difference and sum of $\hat{\omega}_+(\Omega)$ and $\hat{\omega}_-(\Omega)$ can be rewritten as:

$$\begin{aligned} \hat{\omega}_+(\Omega) - \hat{\omega}_-(\Omega) &= \frac{\text{sh}(\Delta u(\Omega_0) - \Delta u(\Omega))}{\alpha_0 e^{\bar{u}(\Omega)} + \text{ch}(\Delta u(\Omega_0) - \Delta u(\Omega))} \end{aligned} \quad (\text{A.12})$$

$$\begin{aligned} \hat{\omega}_+(\Omega) + \hat{\omega}_-(\Omega) &= \frac{\text{ch}(\Delta u(\Omega_0) - \Delta u(\Omega))}{\alpha_0 e^{\bar{u}(\Omega)} + \text{ch}(\Delta u(\Omega_0) - \Delta u(\Omega))} \end{aligned} \quad (\text{A.13})$$

$$\begin{aligned} \alpha_0 &\stackrel{\text{def}}{=} \frac{1 - \omega_0}{\omega_0} \text{ch}(s\Delta u(\Omega_0)) \\ &= \frac{1 - \omega_0}{\omega_0} \frac{1}{\sqrt{1 - (\text{th}(s\Delta u(\Omega_0)))^2}} \\ &= \frac{1 - \omega_0}{\sqrt{\omega_0^2 - \Delta\omega_0^2}} \\ &= \frac{1}{2} \frac{\omega_0^0}{\sqrt{\omega_+^0 \omega_-^0}}. \end{aligned} \quad (\text{A.14})$$

The new variable ' Ω_0 ' is the phenotype value corresponding to the condition $\hat{\omega}_+(\Omega_0) \sim \hat{\omega}_-(\Omega_0)$. The meaning of the constant ' α_0 ' can be related to the Hardy–Weinberg law. Hardy–Weinberg law based on random mating in a population provides a relationship between the genetic microstate frequencies under the form: $p^2 + 2pq + q^2 = 1$, where p^2 and q^2 are the genotype frequencies of genetic microstates '+1' and '-1', i.e. homozygote genotypes aa and AA, respectively; and $2pq$ the genotype frequency for genetic microstate '0', i.e. the heterozygote genotype Aa. In our case, this corresponds to replacing p^2 , q^2 and $2pq$ with, respectively, ω_+^0 , ω_-^0 and ω_0^0 . Consequently, the Hardy–Weinberg law imposes $\alpha_0 = 1$ with $\alpha_0 \neq 1$ corresponding to a deviation from the law. However, this term is expected to remain stable upon any changes of allele or genotype frequencies suggesting therefore that, genetically, any changes in ' $\Delta\omega_0$ ' are to some extent compensated by corresponding changes in ' ω_0 '.

We can now turn to the full expression of the genetic path difference in the space of the phenotype value:

A.5. Expression of the difference between the phenotype responding and default genetic paths expressed in the phenotypic space and conservation of genetic microstate frequencies

The phenotype-associated genetic path is simply the integration of (equation (A.12)) over the phenotype values that is given, as seen above (equation (A.1) in appendix A.1), by: $\int_{\Omega_{1/N}}^{\Omega} [\hat{\omega}_+(\Omega) - \hat{\omega}_-(\Omega)] \frac{1}{\Delta(\Omega)} d\Omega$. The default genetic path is deduced from considering that the difference in the presence probabilities between the genetic microstates '+1' and '-1' is constant, i.e.: $\int_{\Omega_{1/N}}^{\Omega} \Delta\omega_0 \frac{1}{\Delta(\Omega)} d\Omega$. By rewriting ' $\Delta\omega_0$ ' as ' $\omega_0(\Delta\omega_0/\omega_0)$ ', where $\Delta\omega_0/\omega_0 = \text{th}(\Delta u(\Omega_0))$ and deducing ' ω_0 ' from $\alpha_0 = (1 - \omega_0) \text{ch}(\Delta u(\Omega_0))/\omega_0$; it follows: $\Delta\omega_0 = \frac{\text{sh}(\Delta u(\Omega_0))}{\alpha_0 + \text{ch}(\Delta u(\Omega_0))}$. As a result, the difference between the phenotype-associated and default genetic paths expressed as a function of phenotypic fields within the continuum limit is:

$$\begin{aligned} \Delta\hat{\theta}(\Omega) &= \int_{\Omega_{1/N}}^{\Omega} \frac{\text{sh}(\Delta u(\Omega_0) - \Delta u(\Omega))}{\alpha_0 e^{\bar{u}(\Omega)} + \text{ch}(\Delta u(\Omega_0) - \Delta u(\Omega))} \frac{1}{\Delta(\Omega)} d\Omega \\ &\quad - \frac{\text{sh}(\Delta u(\Omega_0))}{\alpha_0 + \text{ch}(\Delta u(\Omega_0))} \int_{\Omega_{1/N}}^{\Omega} \frac{1}{\Delta(\Omega)} d\Omega. \end{aligned} \quad (\text{A.15})$$

The conservation of genetic microstates needs to be added regardless of the genetic path taken, that is, $\Delta\hat{\theta}(\Omega_1) = 0$, expressed as:

$$\int_{\Omega_{1/N}}^{\Omega_1} \frac{\text{sh}(\Delta u(\Omega_0) - \Delta u(\Omega))}{\alpha_0 e^{\overline{u}(\Omega)} + \text{ch}(\Delta u(\Omega_0) - \Delta u(\Omega))} \frac{1}{\Delta(\Omega)} d\Omega = \frac{\text{sh}(\Delta u(\Omega_0))}{\alpha_0 + \text{ch}(\Delta u(\Omega_0))}. \quad (\text{A.16})$$

ORCID iDs

Cyril Rauch  <https://orcid.org/0000-0001-8584-420X>

References

- [1] Buniello A *et al* 2019 The NHGRI-EBI GWAS catalog of published genome-wide association studies, targeted arrays and summary statistics 2019 *Nucleic Acids Res.* **47** D1005–12
- [2] Sudlow C *et al* 2015 UK Biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age *PLoS Med.* **12** e1001779
- [3] Smith B H *et al* 2013 Cohort profile: generation Scotland: Scottish family health study (GS:SFHS). The study, its participants and their potential for genetic research on health and illness *Int. J. Epidemiol.* **42** 689–700
- [4] Fisher R A 1919 XV. The correlation between relatives on the supposition of Mendelian inheritance *Trans. R. Soc. Edinburgh* **52** 399–433
- [5] Fisher R A 1923 XXI. On the dominance ratio *Proc. R. Soc. Edinburgh* **42** 321–41
- [6] Visscher P M and Goddard M E 2019 From R A Fisher's 1918 paper to GWAS a century later *Genetics* **211** 1125–30
- [7] Hivert V, Wray N R and Visscher P M 2021 Gene action, genetic variation, and GWAS: a user-friendly web tool *PLoS Genet.* **17** e1009548
- [8] Stephens M and Balding D J 2009 Bayesian statistical methods for genetic association studies *Nat. Rev. Genet.* **10** 681–90
- [9] Beaumont M A and Rannala B 2004 The Bayesian revolution in genetics *Nat. Rev. Genet.* **5** 251–61
- [10] Boyle E A, Li Y I and Pritchard J K 2017 An expanded view of complex traits: from polygenic to omnigenic *Cell* **169** 1177–86
- [11] van der Harst P and Verweij N 2018 Identification of 64 novel genetic loci provides an expanded view on the genetic architecture of coronary artery disease *Circ. Res.* **122** 433–43
- [12] Herrmann M and Yampolsky L Y 2021 False and true positives in arthropod thermal adaptation candidate gene lists *Genetica* **149** 143–53
- [13] Mähler N *et al* 2020 Leaf shape in *Populus tremula* is a complex, omnigenic trait *Ecol. Evol.* **10** 11922–40
- [14] Zhang W, Reeves G R and Tautz D 2021 Testing implications of the omnigenic model for the genetic analysis of loci identified through genome-wide association *Curr. Biol.* **31** 1092–8
- [15] Vuckovic D *et al* 2020 The polygenic and monogenic basis of blood traits and diseases *Cell* **182** 1214–31
- [16] Mathieson I 2021 The omnigenic model and polygenic prediction of complex traits *Am. J. Hum. Genet.* **108** 1558–63
- [17] Galton F 1886 Regression towards mediocrity in hereditary stature *J. R. Anthropol. Inst. GB Irel.* **15** 246–63
- [18] Visscher P M, Medland S E, Ferreira M A R, Morley K I, Zhu G, Cornes B K, Montgomery G W and Martin N G 2006 Assumption-free estimation of heritability from genome-wide identity-by-descent sharing between full siblings *PLoS Genet.* **2** e41
- [19] Silventoinen K *et al* 2003 Heritability of adult body height: a comparative study of twin cohorts in eight countries *Twin Res.* **6** 399–408
- [20] Nelson R M, Pettersson M E and Carlborg C 2013 A century after Fisher: time for a new paradigm in quantitative genetics *Trends Genet.* **29** 669–76
- [21] Yang J *et al* 2010 Common SNPs explain a large proportion of the heritability for human height *Nat. Genet.* **42** 565–9
- [22] Wood A R *et al* 2014 Defining the role of common variation in the genomic and biological architecture of adult human height *Nat. Genet.* **46** 1173–86
- [23] Visscher P M, Hill W G and Wray N R 2008 Heritability in the genomics era—concepts and misconceptions *Nat. Rev. Genet.* **9** 255–66
- [24] Yengo L *et al* 2022 A saturated map of common genetic variants associated with human height *Nature* <https://doi.org/10.1038/s41586-022-05275-y>
- [25] Stigler S M 1990 *The History of Statistics: The Measurement of Uncertainty before 1900* (Cambridge, MA: Harvard University Press)
- [26] Macdonald A, Hawkes L A and Corrigan D K 2021 Recent advances in biomedical, biosensor and clinical measurement devices for use in humans and the potential application of these technologies for the study of physiology and disease in wild animals *Phil. Trans. R. Soc. B* **376** 20200228
- [27] Falconer D S 1996 *Introduction to Quantitative Genetics* (Englewood Cliffs, NJ: Prentice-Hall)
- [28] Lonsdale J *et al* 2013 The genotype-tissue expression (GTEx) project *Nat. Genet.* **45** 580–5
- [29] Delongchamp R, Faramawi M F, Feingold E, Chung D and Abouelenein S 2018 The association between SNPs and a quantitative trait: power calculation *Eur. J. Environ. Public Health* **2** 10
- [30] Park J-H, Gail M H, Weinberg C R, Carroll R J, Chung C C, Wang Z, Chanock S J, Fraumeni J F and Chatterjee N 2011 Distribution of allele frequencies and effect sizes and their interrelationships for common genetic susceptibility variants *Proc. Natl Acad. Sci. USA* **108** 18026–31
- [31] Sham P C, Cherny S S, Purcell S and Hewitt J K 2000 Power of linkage versus association analysis of quantitative traits, by use of variance-components models, for sibship data *Am. J. Hum. Genet.* **66** 1616–30
- [32] Wattis J, Bray S, Kyrtzi P and Rauch C 2022 Analysis of genotype–phenotype association using fields and information theory (arXiv:2202.11989)