

1 RESEARCH ARTICLE

2 RUNNING HEAD: Upgrading the investigative power of biological/medical datasets

3 Investigative power of Genomic Informational Field  
4 Theory (GIFT) relative to GWAS  
5 for genotype-phenotype mapping  
6

7 Panagiota Kyratzi<sup>1,2</sup>, Oswald Matika<sup>3</sup>, Amey H Brassington<sup>4</sup>, Connie E Clare<sup>5</sup>, Juan Xu<sup>6</sup>, David A  
8 Barrett<sup>7</sup>, Richard D Emes<sup>8</sup>, Alan L Archibald<sup>3</sup>, Andras Paldi<sup>2</sup>, Kevin D Sinclair<sup>4</sup>, Jonathan Wattis<sup>9</sup>  
9 and, Cyril Rauch<sup>1</sup>

10

11 <sup>1</sup>School of Veterinary Medicine and Science, University of Nottingham, College Road, Sutton Bonington,  
12 LE12 5RD, UK.

13 <sup>2</sup>École Pratique des Hautes Études, PSL Research University, St-Antoine Research Center, Inserm U938, 34  
14 rue Crozatier, 75012 Paris, France.

15 <sup>3</sup>Div. Genetics and Genomics, The Roslin Institute and Royal (Dick) School of Veterinary Studies, University  
16 of Edinburgh, Easter Bush, Midlothian EH25 9RG, Scotland, UK

17 <sup>4</sup>Agriculture and Horticulture Development Board, Middlemarch Business Park Siskin Parkway, East  
18 Coventry CV3 4PE, UK.

19 <sup>5</sup>School of Biosciences, University of Nottingham, College Road, Sutton Bonington, LE12 5RD, UK.

20 <sup>6</sup>Shanghai Leadingtac Pharmaceutical Co., Ltd, 781 Cailun Road, China (Shanghai) Pilot Free Trade Zone,  
21 Pudong, Shanghai 201203, China

22 <sup>7</sup>Centre for Analytical Bioscience, School of Pharmacy, University of Nottingham, Nottingham NG7 2RD,  
23 UK.

24 <sup>8</sup>Nottingham Trent University, 50 Shakespeare Street, Nottingham NG1 4FQ, UK.

25 <sup>9</sup>Centre for Mathematical Medicine and Biology, School of Mathematical Sciences, University of  
26 Nottingham, University Park, Nottingham NG7 2RD, UK.

27

28 Correspondence: Cyril Rauch ([cyril.rauch@nottingham.ac.uk](mailto:cyril.rauch@nottingham.ac.uk)).

29

---

30 **ABSTRACT**

31 Identifying associations between phenotype and genotype is the fundamental basis of genetic analyses.  
32 Inspired by frequentist probability and the work of R.A. Fisher, genome-wide association studies (GWAS)  
33 extract information using averages and variances from genotype-phenotype datasets. Averages and  
34 variances are legitimated upon creating distribution density functions obtained through the grouping of  
35 data into categories. However, as data from within a given category cannot be differentiated, the  
36 investigative power of such methodologies is limited. Genomic Informational Field Theory (GIFT) is a  
37 method specifically designed to circumvent this issue. The way GIFT proceeds is opposite to that of GWAS.  
38 Whilst GWAS determines the extent to which genes are involved in phenotype formation (bottom-up  
39 approach), GIFT determines the degree to which the phenotype can select microstates (genes) for its

40 subsistence (top-down approach). Doing so requires dealing with new genetic concepts, a.k.a. genetic  
41 paths, upon which significance levels for genotype-phenotype associations can be determined. By using  
42 different datasets obtained in *ovis aries* related to bone growth (Dataset-1) and to a series of linked  
43 metabolic and epigenetic pathways (Dataset-2), we demonstrate that removing the informational barrier  
44 linked to categories enhances the investigative and discriminative powers of GIFT, namely that GIFT  
45 extracts more information than GWAS. We conclude by suggesting that GIFT is an adequate tool to study  
46 how phenotypic plasticity and genetic assimilation are linked.

47

## 48 **NEW & NOTEWORTHY**

49 The genetic basis of complex traits remains challenging to investigate using classic GWASs. Given the  
50 success of gene editing technologies this point needs to be addressed urgently since there can only be  
51 useful editing technologies if precise genotype-phenotype mapping information is available initially. GIFT  
52 is a new mapping method designed to increase the investigative power of biological/medical datasets  
53 suggesting, in turn, the need to rethink the conceptual bases of quantitative genetics.

54 **Keywords:** Complex traits; GIFT; genotype-phenotype mapping studies; GWAS

55

---

## 56 **INTRODUCTION**

57 Identifying associations between phenotype and genotype is the fundamental basis of genetic analysis.  
58 The development of high-density genotyping and whole genome sequencing has enabled DNA variants to  
59 be directly identified and Genome-Wide Association Studies (GWASs) have become the method of choice  
60 for mapping genotype to phenotype in large populations of unrelated individuals. GWAS have been  
61 employed in many species, and especially in the study of human disease (1). By 2021 the NHGRI-EBI GWAS  
62 Catalog listed 316,782 associations identified in 5149 publications describing GWAS results (2).  
63 Additionally, extensive collection of data has been initiated through efforts such as the UK Biobank (3),  
64 Generation Scotland (4) and NIH *All of Us* research program (<https://allofus.nih.gov/>) in the expectation  
65 that large-scale GWAS will elucidate the basis of human health and disease and facilitate precision  
66 medicine.

67

68 While genomic technologies have advanced rapidly, statistical models used to analyze genetic data are  
69 still based on the models developed by Fisher more than 100 years ago (5, 6). GWASs essentially make  
70 use of the Fisher method of partitioning genotypic values by performing a linear regression of phenotype  
71 on marker allelic dosage (7). Regression coefficients estimate the average allele effect size, and the  
72 regression variance is the additive genetic variance due to the locus (8). However, an ongoing debate exists  
73 over whether the present analysis paradigm in quantitative genetics is at its limits for truly understanding  
74 complex traits, namely traits resulting from many genes each with very small effect size (9). As a result,  
75 one may wonder whether alternative statistical model(s) could be invented and used to determine  
76 genotype-phenotype mappings.

77

78 GWASs are fundamentally linked to frequentist probabilities that, defined through relative frequencies,  
79 determines the validity of statistical inferences. In practice, frequentist probabilities are generated  
80 through the grouping of data into bins or categories to generate a bar chart, that is then interpolated to  
81 create a distribution density function (DDF) in the continuum limit. The DDF is, in turn, used to determine  
82 statistical inferences including average, variance, p-value and so on. However, since the DDF approximates

83 the bar chart (and not the converse), and that it is not possible to differentiate data from within any given  
84 group/category, the DDF is constructed mathematically on the implicit assumption that information is  
85 missing to differentiate data from within any given group/category.

86  
87 The notion of ‘missing information’ can be legitimate and defined experimentally. For example, measuring  
88 the phenotype human height with a ruler with centimetre graduations implies that any height can be  
89 measured to the nearest centimetre. Consequently, one centimetre-width bins/categories need to be  
90 used to generate a frequency table of range of phenotype values upon which the phenotype and genotype  
91 DDFs are defined. In this case, all the resulting statistical inferences are defined with a precision  
92 corresponding to the nearest centimetre. The ‘missing information’ (i.e., that what cannot be measured  
93 by the ruler) corresponds then to sub-centimetric scales (i.e., distances to the nearest millimetre for this  
94 example). In practice the ‘missing information’ is therefore linked to the one of ‘imprecision’ and deciding  
95 to provide more precise statistical inferences implies that the width of categories be reduced, which can  
96 only be achieved by increasing the sample size. It is not by chance that the ‘normal distribution’ created  
97 by mathematicians and physicists was initially called the ‘law of errors’, where the notion of error  
98 (misinformation) results from imprecisions in experimental measurements. As a result, GWAS is faced  
99 with a fundamental issue involving the extraction of precise information using a method that,  
100 conceptually, assumes that information is missing or that data is mis-(in)formed.

101  
102 In general, the problem concerning the ‘missing information’ is never mentioned since the DDF in the  
103 continuum limit is never considered as an approximation but as something that has its own reality. Namely  
104 a DDF must exist independently of data measured (i.e., data must fit the DDF and not the converse). The  
105 latter remark leads to an interesting conceptual territory where the notions of average and variance, and  
106 their usage, may be questioned. If one considers the normal distribution (or any other DDFs) is inherent  
107 to life and that data must fit it (them), then the moments of the distribution (e.g., average and variance)  
108 are also essential parameters to describe life, and the variance often interpreted as noise in the data is  
109 then a nuisance. If, on the contrary, data is the important thing, and that the DDF is considered solely as  
110 a tool to interpolate data based on missing information, then average and variance are parameters  
111 derived from a lack of information and are, as a result, poorly informative. The latter point should not  
112 come as a surprise as reducing the huge diversity of populations to a handful of parameters (i.e., average  
113 and variance) is highly reductionist and likely to be poorly descriptive. Thus, while the notions of average  
114 and variance may help representing datasets, they are inventions nonetheless, i.e., thought constructions  
115 akin to the field of frequentist probability. Thus, using average and variance as a starting point to map  
116 genotype-phenotype (GWAS) is a matter of choice. Accordingly, different statistical methods can be  
117 suggested.

118  
119 To avoid those conceptual and practical issues a new method called GIFT (Genomic Informational Field  
120 Theory) has been designed and applied to simulated genotype-phenotype data in (10, 11), reviewed in  
121 (12). In short, to associate genotype to phenotype GIFT does not presume that the only important  
122 information concerning the gene effect is found in averages or variances, nor does it presume that DDFs  
123 are central. On the contrary, GIFT starts with the pre-requisite that phenotypic values, or phenotypic  
124 residuals after considering the environment/fixed effects, may be measured with sufficient precision to  
125 be unique in a population. Then, by avoiding grouping data into bins/categories, which would otherwise  
126 create an artificial imprecision, GIFT considers the entire information contained in the data, (i.e., variance  
127 is not a nuisance anymore) making use of the cumulative sum of microstates. Figure 1 provides the  
128 intuition underscoring GIFT as a method.

129

130 The current article extends our previous theoretic studies using simulated data to analyse for the first  
131 time two real datasets:

- 132 i. Dataset-1 is derived from a study concerned with the genetic background of carcass composition in  
133 sheep (*ovis aries*) (13). Using GWAS this study demonstrated a strong association between  
134 chromosome 6 and the carcass composition trait 'bone area at the ischium'. We now apply GIFT to  
135 reanalyse this dataset to benchmark it against GWAS. Since GWAS previously identified a QTL in  
136 chromosome 6, our hypothesis was that GIFT would at least replicate GWAS results and identify  
137 additional putative QTLs.
- 138 ii. Dataset-2 comprises biochemical data arising from an ongoing study in sheep which seeks to identify  
139 risk allele variants in genes whose products direct a series of metabolic pathways, collectively referred  
140 to as one carbon (1C) metabolism and associated epigenetic regulators. The gene array was designed  
141 to include all single nucleotide polymorphisms (SNPs) linked to known biochemical enzymes involved  
142 in these pathways. Given that Dataset-2 preselected genes for a targeted analysis of enzymes involved  
143 in these metabolic/epigenetic pathways, it can be considered more specific.

144  
145 The present article initially introduces the reader to the way data may be used and analysed differently  
146 using GIFT, contrasting to more conventional methods mostly based on an analysis of averages and  
147 variances. More specifically in Part 1, the null hypothesis defined by GIFT will be established. Using  
148 Dataset-1 the concept of genetic path pertaining to GIFT will be introduced (Part 2) out of which a p-value  
149 for GIFT will be defined (Part 3). Then Dataset-1 (Part 4) and Dataset-2 (Part 5) will be analysed comparing  
150 the informational/investigative power of GIFT relative to GWAS using Manhattan plots prior to performing  
151 enrichment analyses.

## 152 MATERIALS AND METHODS

### 153 Biological datasets.

154 The first dataset (Dataset-1) analysed 600 pedigree-recorded Scottish Blackface lambs using CT scans to  
155 determine *in vivo* carcasses composition (13). The trait selected for the present study is the bone areas of  
156 the ischium (BAI) measured in mm<sup>2</sup> from cross-sectional CT scans. The ischium is one of the three bones  
157 that make up the pelvis. It is located beneath the ilium and behind the pubis. The upper portion of the  
158 ischium forms a major part of the concave portion of the pelvis that forms the hip. The BAI crossed a  
159 genome-wide significance threshold on Chromosome 6 (OAR6). The pre-corrected phenotype values were  
160 obtained fitting fixed effects of age of dam, year of birth, the effect of management group (as sheep were  
161 from different farms), sex (males or females) and litter size (singles or twins) and as covariate the day of  
162 birth. Further information can be found in Matika et al. (2016) (13). Supplemental S1 provides the raw  
163 data used (Dataset-1).

164 The second dataset (Dataset-2) was from previously unpublished data extracted from a large ongoing  
165 programme of research to investigate genome regions (Quantitative trait loci, (QTL)) that determine  
166 metabolic and epigenetic responses to nutritionally induced deficiencies in one carbon metabolism (14,  
167 15). For this study sheep were used as an experimental model. All animal procedures relating to this study  
168 adhered to the Animals (Scientific Procedures) Act, 1986. Associated protocols complied with the ARRIVE  
169 guidelines and were approved by the University of Nottingham Animal Welfare and Ethical Review Body  
170 (AWERB) with Home-Office project licensed authority (30/3376;10<sup>th</sup> February 2016). Supplemental S2  
171 provides the raw data used (Dataset-2).

172

## 173 **Dataset-2: Sheep genome resequencing, custom array design and SNP profiling on test subjects.**

174 Twenty-four unrelated Texel ewes were sequenced to a depth of 30x in 2 pools at Edinburgh Genomics.  
175 DNA samples were prepared using Illumina's TruSeq PCR free kits and sequenced on an Illumina HiSeq  
176 2500 Rapid Mode (serial no. D00125), read length of 150PE. Reads were trimmed to remove adapter  
177 sequences and low-quality bases using skewer with commands (-Q 20, -q 3) (16) and mapped to the  
178 reference sheep genome assembly (Oar\_v3.1) using bwa mem (options -M -t 4) (17). Following  
179 deduplication using Picard-tools version 1.92, variants were called using GATK pipeline (18) including  
180 realignment around known indels and recalibration of bases, and FreeBayes (--use-best-n-alleles 4 --  
181 pooled-discrete --min-alternate-count 4). Annotation of SNPs was performed using Ensembl variant effect  
182 predictor VEP version ensembl tools release 79 (19). 15,347,831 variants were identified. Of these, ~3  
183 million were novel SNPs and ~12 million were already present in the Ensembl genome database. SNPs  
184 within annotated coding regions (VEP annotated "downstream gene variant" or "intron variant" removed)  
185 and within 3Kb upstream of a gene were retained. SNPs with a minor allele frequency of greater than 0.5  
186 were used to design an Illumina Infinium® iSelect® Custom Array consisting of 4,576 probes. This captured  
187 SNPs in 115 1C metabolism and related genes, and 108 related epigenetic regulators as well as 33 control  
188 SNPs (Supplemental S1).

189  
190 Liver samples were next collected post-mortem from 360 male and female Texel lambs (6 to 11 months  
191 of age) representing 11 farms dispersed regionally across the UK. Collections took place at regional  
192 abattoirs and samples immediately snap frozen in liquid N and stored at -80°C until analyses. DNA was  
193 then extracted using AllPrep DNA/RNA Mini kit (Qiagen, Manchester UK). Briefly approximately 20 mg of  
194 liver were mechanically disrupted using a TissueLyser (Qiagen, Manchester, UK) in 600 RLT plus buffer  
195 containing  $\beta$ -mercaptoethanol. Tissue lysates were then used to extract RNA and DNA according to the  
196 manufacturer instructions. The custom designed array was then used to SNP profile DNA from these Texel-  
197 sheep. For this purpose, liver samples were collected post-mortem from lambs (aged 6 to 11 months)  
198 representing 11 farms dispersed regionally across the UK. Collections took place at regional abattoirs and  
199 samples immediately snap frozen in liquid N and stored at -80°C until analyses. DNA was then extracted  
200 using AllPrep DNA/RNA Mini kit (Qiagen, Manchester UK). Briefly approximately 20 mg of liver were  
201 mechanically disrupted using a TissueLyser (Qiagen, Manchester, UK) in 600 RLT plus buffer containing  $\beta$ -  
202 mercaptoethanol. Tissue lysates were then used to extract RNA and DNA according to the manufacturer  
203 instructions.

204

## 205 **Dataset-2: Metabolic profiling.**

206 For the purposes of the current study the following seven liver metabolites were selected from a larger  
207 pool of 1C metabolites: S-adenosyl methionine (SAM), methylcobalamin (mB12), adenosylcobalamin  
208 (aB12), trimethylglycine (TMG), dimethylglycine (DMG), propionate (PPA) and methylmalonic acid (MMA).  
209 The first four metabolites were selected as representative intermediates of the methionine cycle whilst  
210 the latter two are intermediates in the hepatic synthesis of succinate (15) (Fig.2 & Supplemental S1).

211  
212 Hepatic concentrations of four metabolites (i.e., mB12, aB12, TMG and DMG) were determined by  
213 hydrophilic interaction chromatography (HILIC) coupled to electrospray ionization tandem mass  
214 spectrometry (MS/MS) as reported previously (20). For the analysis of SAM (determined separately by  
215 HILIC), the standard was purchased from Sigma-Aldrich (Poole, Dorset, UK). Stock solutions of this  
216 standard were prepared in potassium phosphate extraction buffer ( $\text{KH}_2\text{PO}_4$  and  $\text{K}_2\text{HPO}_4$ ; 40 mmol/L)  
217 containing 0.1% L-ascorbic acid, 0.15% citric acid and 0.1% MCE (adjusted to pH 7 with NaOH), each at a  
218 final concentration of 100  $\mu\text{mol/L}$ . Also, for SAM the mobile phase was modified from that used for the  
219 three other reported metabolites by adjusting the pH of the aqueous ammonium carbonate buffer



220 solution from 3.5 to 9.1. Mass spectrometer parameters for SAM were as follows: retention time = 7.69  
221 min; Q1mass = 399.1 amu; Q3 mass = 250.1 amu; declustering potential = 56; collision energy = 25;  
222 collision cell exit potential = 16.

223  
224 Hepatic concentrations of PPA and MMA were determined by gas chromatography coupled to mass  
225 spectroscopic-detection (GC-MS). Briefly, for PPA, 750  $\mu$ L 5-Sulfosalicylic acid (SSA, 0.04 mg/ml) was  
226 added to 150mg frozen liver, homogenised for 2 min and cooled on ice for 10 min. The sample was  
227 centrifuged for 15 min at 14,500 x g and 200  $\mu$ L liver homogenate transferred to a 2.5 mL screw capped  
228 glass vial. To this, 20  $\mu$ L internal standard (MBA, 400  $\mu$ M), 3.5  $\mu$ L HCl (37%) and 1 mL diethylether were  
229 added, vortexed for 2 min and centrifuged for 10 min at 14,500 x g. 600  $\mu$ L of the upper layer was  
230 transferred to a screw capped glass vial containing 3.5  $\mu$ L 1-(tert-butyldimethylsilyl)imidazole (TMDMSIM,  
231 97%), vortexed for 2 min and heated at 60°C for 30 min. GC-MS analysis proceeded after cooling. The  
232 method used a DB-5MS column (J&W Scientific Agilent technology, 30 m x 0.25 mm; 0.25  $\mu$ m film  
233 thickness). The carrier gas (He) was set at a constant flow rate of 1.3 ml/min. The injection volume was 5  
234  $\mu$ L for SCAN mode (for qualification) and SIM (selected ion monitoring) mode (for quantification), both  
235 using splitless mode. The injection port and MS selective detector interference temperatures were 260°C  
236 and 250°C respectively. The chromatograph was programmed for an initial temperature of 40°C for 1 min,  
237 increased to 60°C at 70°C min<sup>-1</sup>, then to 110°C at 15°C min<sup>-1</sup>, and finally 250°C at 70°C min<sup>-1</sup>. MS was  
238 tuned regularly and operated in electron impact (EI) ionization mode with the ionization energy of 70eV.  
239 SCAN mode measured at m/z: 30-300 and SIM ions were set at 159 (for MBA) and 131 (for PPA). The same  
240 method was used to produce a calibration curve for PPA using standards at concentrations ranging from  
241 19.5 nmol/g to 5 $\mu$ mol/g. The limit of detection was 19.5 nmol/g. CVs for low, medium and high QCs were  
242 10.4, 6.3 and 6.5% and the inter-assay CV was 4.7%.

243  
244 For MMA, 250  $\mu$ L 80% MeOH was added to 50 mg frozen liver, homogenised for 2 min and cooled on ice  
245 for 10 min. The sample was ten centrifuged for 15 min at 14,500 x g and 200  $\mu$ L liver homogenate  
246 transferred to a 2.5 mL screw capped glass vial. To this, 4  $\mu$ L internal standard (1 mM 4-chlorobutyric acid  
247 (CBA) in 1 mM HCl) followed by 250  $\mu$ L 12% BF<sub>3</sub>-Methanol were added, vortexed for 1 min and heated at  
248 95°C for 15 min. After cooling, 250  $\mu$ L cold distilled water and 250  $\mu$ L cold dichloromethane (CH<sub>2</sub>Cl<sub>2</sub>) were  
249 added to the vial, vortexed for 30s and centrifuged for 10 min at 14,500 x g. The lower dichloromethane  
250 layer was transferred to a screw capped glass auto-sampler vial with insert for GC-MS analysis. The  
251 method used a DB-WAX column (cross-linked polyethylene glycol; J&W Scientific Agilent technology) (30  
252 mm x 0.25 mm; 0.15  $\mu$ m film thickness). The carrier gas (He) was set at a constant flow rate of 1.0 ml/min.  
253 The injection volume was 1  $\mu$ L for SCAN mode (for qualification) and SIM mode (for quantification), both  
254 using splitless mode. The injection port and MS selective detector interference temperatures were 260°C  
255 and 280°C respectively. The chromatograph was programmed for an initial temperature of 50°C for 2 min,  
256 increasing to 150°C at 8°C min<sup>-1</sup>, then to 220°C at 100°C min<sup>-1</sup> and held for 5 min at the final temperature.  
257 MS was tuned regularly and operated in EI ionization mode with the ionization energy of 70eV. The limit  
258 of detection was 0.75 nmol/g for both MMA and SA and inter-assay CVs were 8.4% for MMA and 11.0%  
259 for SA.

## 260 **Dataset-2: Determination of GWAS for 1C-metabolites.**

261 Preliminary data analysis indicated the need to log-transform using the natural logarithm (Supplemental  
262 S3) to approximate normality. Transformed data were then pre-corrected for the fixed effects of farm (F)  
263 and sex (S) in ASReml using the following model,  $y_{ij} = \mu + F_i + S_j + e_{ij}$ , where  $y_{ij}$  is the log-transformed  
264 phenotype, that is the log-transformed metabolite concentration studied;  $\mu$  is the overall mean for the  
265 log-transformed metabolite concentration;  $F_i$  is the effect of the  $i^{\text{th}}$  farm ( $i = 1, \dots, 11$ );  $S_j$  the effect of  $j^{\text{th}}$  Sex

266 (Male vs Female) and,  $e_{ij}$  is the residual. The genotype dataset was filtered using PLINK (HWE p-value  
267 threshold of  $10^{-6}$ , call rate for genotypes of 10% and a MAF of 5%), the number of independent SNPs was  
268 determined using BCFTOOLS ( $r^2$ -threshold=0.1) and the GWAS Manhattan plots, linked to the  
269 determination of  $p_{GWAS}$ , were obtained using GEMMA. The same genotype and residual phenotypes as  
270 filtered by GWAS were used by GIFT.

## 271 Data representation using GIFT.

272 Adjusted phenotypic data (i.e., residuals, from Dataset-1 and Dataset-2) were used for this study.  
273 Regarding the representation of GIFT, upon selecting a SNP for all individuals, the different corresponding  
274 genotypes, aa, aA/Aa and AA, were assigned the arbitrary values +1, 0 and -1, respectively. With this  
275 convention any barcode can be represented by a string of numbers from which a GIFT analysis can be  
276 inferred. More specifically, the assignment of values +1, 0 and -1 were done as a function of the base pairs  
277 as follow: AA=TT=+1, GG=CC=-1 and 0 otherwise. As shown schematically in Fig.1, the residuals obtained  
278 were ranked by order of magnitude and the cumulative sum of their corresponding genotypic values  
279 performed to obtain the 'genetic path' for the SNP considered. The genetic path of a SNP is noted  $\theta(i)$  in  
280 the text (Fig.1). The null hypothesis for GIFT as well as the notion of significance when GIFT is used will be  
281 introduced and fully explained in the RESULTS section.  
282

## 283 RESULTS

### 284 Analyze of the null hypothesis $\theta_0(i)$ for GIFT

285 While  $\theta(i)$  is obtained using phenotypic information ( configuration ① in Fig.1 and 'Data representation  
286 using GIFT' in MATERIALS AND METHODS), it is also possible to plot the cumulative sum of microstates  
287 when no phenotypic information is present that is equivalent to 'scrambling' or permutating the string of  
288 microstates in Fig.1A also corresponding to the configuration ② in Fig.1B. Recall that since our focus is  
289 on a given SNP, then the number of microstates,  $N_+$ ,  $N_0$  and  $N_-$ , are identical between the configurations  
290 ① and ②. This new cumulative sum noted  $\theta_0(i)$  is expected to be a sort of null hypothesis solely  
291 dependent on the bulk microstate frequencies  $N_+/N$ ,  $N_0/N$  and  $N_-/N$ , where  $N_q$   $q \in \{+,0,-\}$  is the  
292 number of microstates of type  $q$ . This is so because there is no further information that could inform on  
293 the positioning of microstates in their list when the scrambled state is considered. However, while  $\theta(i)$  is  
294 unique since phenotypic information is used to generate it,  $\theta_0(i)$  is not as each time the string of  
295 microstates from Fig.1A is scrambled, a new  $\theta_0(i)$  appears. Accordingly, one needs to consider the set of  
296 possible  $\theta_0(i)$ s generated bounded to the microstate frequencies  $N_+/N$ ,  $N_0/N$  and  $N_-/N$ .

297 Using a selection of theoretic SNPs defined by different microstate frequencies (Table-1). Fig.3A illustrates  
298 the global shape resulting from simulating 1000  $\theta_0(i)$ s. The results demonstrate that the global shape of  
299 the  $\theta_0(i)$ s plotted as a function of the position in the string is ellipsoidal with short and long axes changing  
300 as a function of microstate frequencies involved, and where the different averages of  $\theta_0(i)$ s represented  
301 by black lines in Fig.3A, are straight lines with slopes linked to the difference,  $\Delta N/N = (N_+ - N_-)/N$ . The  
302 fact that the averages of  $\theta_0(i)$ s for a given set of microstates,  $N_+$ ,  $N_0$  and  $N_-$ , is always a straight line  
303 linked to microstate frequencies,  $N_+/N$ ,  $N_0/N$  and  $N_-/N$ , can be understood intuitively by the fact that  
304 scrambling or permutating an infinite number of times the string of microstates is equivalent to  
305 determining, for any position  $i$ , the presence probability,  $N_q/N$ , of each microstate in the string.  
306 Accordingly, for a given set of microstates,  $N_+$ ,  $N_0$  and  $N_-$ , the average of  $\theta_0(i)$ s, noted  $\langle \theta_0(i) \rangle$ , is  
307  $\langle \theta_0(i) \rangle = \frac{(N_+ - N_-)}{N} i$ . Further theoretic details can be found in (10, 11). Using  $\langle \theta_0(i) \rangle$  as a reference for the

308 null hypothesis, Fig.3B show the sur-imposition of the differences,  $\Delta\theta_0(i) = \theta_0(i) - \langle\theta_0(i)\rangle$ , obtained  
309 from simulations using SNPs from Table-1.

310 Finally, to assess the impact of the sample size (population size) on the null hypothesis the initial size  
311 ( $N=565$ , Table-1) was divided ( $N=280$ ) and multiplied ( $N=1130$ ) by a factor  $\sim 2$  while keeping constant the  
312 microstate frequencies  $N_+/N$ ,  $N_0/N$  and  $N_-/N$  from Table-1. The simulations Fig.3A show that the  
313 appearance of ellipsoids is affected when the sample size changes, becoming thinner as the population  
314 size increases. Plotting the standard deviation,  $\sigma(i/N)$ , as a function of the position once normalized by

315 the sample size,  $\sigma(i/N) = \sqrt{[\langle(\theta_0(i/N))^2\rangle - \langle\theta_0(i/N)\rangle^2]}/N$ , resulting from the different simulations in  
316 Fig.3C demonstrates that the standard deviation from GIFT is quadratic, and independent of the sample  
317 size, as expected from a random allocation of different microstates in the string of positions.

318 At first sight and with this primary analysis one could suggest that any genetic path departing from the  
319 cloud of genetic paths formed by the set of  $\theta_0(i)$ s upon the permutation of microstates (grey surface in  
320 Fig.3A or black surface in Fig.3B) would likely result in an association between the genotype and the  
321 phenotype. While true this assumption needs to be handed out carefully as it is not exhaustive. Indeed,  
322 some genetic paths may be highly structured and of relatively small amplitude. Examples of genetic path  
323 using real data from Dataset-1 will demonstrate this point.

324

#### 325 **Examples of genetic path using the bone area of the ischium (BAI) as phenotype (Dataset-1)**

326 The resulting average,  $\langle\theta_0(i)\rangle$ , and variance,  $\sigma(i)$ , can be used to inform the null hypothesis of a particular  
327 SNP from 'real' datasets. However, since there are as many different sets of  $\theta_0(i)$ s as number of SNPs,  
328 each SNP will return its own  $\langle\theta_0(i)\rangle$  (null hypothesis) upon scrambling. A comparison between SNPs using  
329 GIFT/genetic paths requires then to concentrate on the differences,  $\Delta\theta(i) = \theta(i) - \langle\theta_0(i)\rangle$ . In the  
330 remaining text one shall rewrite  $\langle\theta_0(i)\rangle$  as  $\theta_0(i)$  to simplify notations.

331 Concentrating now on 'real' dataset, the genetic paths were obtained further to ranking BAI residual  
332 values (Dataset 1) using an incremental rank from small to large values. As an example, Fig.4 shows the  
333 two genetic paths  $\theta(i)$  and  $\theta_0(i)$  for six SNPs, renamed SNP1-6 (see Table-2 for accurate genetic  
334 information) enabling us to appreciate the qualitative difference between the genetic paths. While the  
335 null hypothesis, i.e.,  $\theta_0(i)$ , resulting from the scrambling of phenotypic values many times always returns  
336 a straight line with a different slope for each SNP as seen above, the  $\theta(i)$ s have different shape. To  
337 represent the set of  $\theta(i)$ s in relation to the different microstates involved, each datapoint of the  $\theta(i)$ s is  
338 colour coded as in Fig.1C.

339 Since  $\theta_0(i)$  is linked to the difference between the genetic microstate frequencies of homozygotes,  $\Delta N =$   
340  $N_+ - N_-$ , in Fig.4 we represent by the angle  $\alpha$  such difference. Since  $\tan(\alpha) = +N_+/N - N_-/N$  where  $N$   
341 is the total number of positions ( $i = 1, 2, \dots, N$ )  $\theta_0(i)$  can be rewritten as,  $\theta_0(i) = \tan(\alpha)i$ . As any analysis  
342 must concentrate on the difference,  $\Delta\theta(i) = \theta(i) - \theta_0(i)$ , such as to cancel the apparent variability in  
343 the null hypothesis across SNPs, we represent the plots of the different  $\Delta\theta(i)$ s obtained in the right panel  
344 of Figs.4A-4F.

345 Figs.4A-4B display two distinct genetic paths that are globally similar. While they have different number  
346 of microstates of each type (see Table-2) the  $\Delta\theta(i)$ s of SNP1 and SNP2 are characterized by their small  
347 amplitudes and the fact that they are erratic crossing several times the axis of position corresponding to  
348 the null hypothesis. In those cases, using the information contained in the phenotypic residuals, namely  
349 ranking the phenotypic residuals from small to large values, does not permit to fully differentiate  $\theta(i)$   
350 from  $\theta_0(i)$ . On the other hand, the right panel in Figs.4C-4D for SNP3 and SNP4 demonstrates, in a more



351 noticeable way, a paraboloid shape for the  $\Delta\theta(i)$ s resulting from a segregation of microstates upon  
352 ordering the phenotypic residuals. The segregation of microstates +1 and -1 in opposite direction is  
353 reminiscent of Fisher theoretic works (Fig.1). As it turns out Figs.4C-4D show some similarities with Fig.1C  
354 based on a simulation inspired by Fisher's seminal works. Importantly the  $\Delta N$ -values of SNP1 and SNP4  
355 while of opposite sign are similar in absolute value, are as those of SNP2 and SNP3, suggesting, in turn,  
356 the  $\Delta N$ -values do not impact on the ability to differentiate  $\theta(i)$  from  $\theta_0(i)$ . Namely that a segregation of  
357 microstates can be inferred also with relatively large and opposed  $\Delta N$ -values.

358 Envisaging the migration of microstates +1 and -1 in opposite direction as initially postulated by Fisher as  
359 the sole framework to associate genotype and phenotype is not always valid. This is demonstrated by  
360 SNP5 and SNP6 and the appearance of structured genetic paths displaying clear sigmoidal shapes for the  
361  $\Delta\theta(i)$ s as shown in Figs.4E-4F. Theoretically this phenomenon can be understood and explained by the  
362 presence of non-linear phenotypic fields, see (11) also reviewed in (12), in turn breaking the symmetry  
363 postulated by Fisher assuming the sole presence of linear phenotypic fields. This type of sigmoidal shapes  
364 is of interest since they inform on potential regulation mechanisms involving very probably 'regulatory  
365 variants' (21). Indeed, the right panels in Figs.4E-4F can be envisioned as representing the genetic  
366 organization of two distinct subpopulations of phenotypic residual values, one above the dashed line and  
367 the other one underneath it. Taken separately those two subpopulations draw curves like Figs.4C-4D or  
368 Fig.1C. In this context it is tempting to suggest that sigmoid genetic paths reveal a type of genotype-  
369 phenotype association that is inherently 'scale-dependent', namely function of the magnitude of  
370 phenotypic residuals. Because traditional GWAS concentrates on averages and variances, these sigmoid  
371 paths would be remarkably difficult to characterize with traditional methods. This is so because there is  
372 no clear antisymmetric segregation of microstates. As an example, using SNPs1-6 (from Fig.4) we have  
373 plotted, in Fig.5, the average values of phenotypic residuals for each microstate, and in Table-2 we provide  
374 the resulting gene/size effects and the dominances associated with those. Fig.5 and Table-2 demonstrate  
375 that sigmoid genetic paths (SNP5 and SNP6) are much less detectable with traditional methods while  
376 paraboloid genetic paths (SNP3 and SNP4) are. Note that the numerical determination of ' $-\text{Log}_{10}(p_{\text{GIFT}})$ '  
377 in Table-2, that is the significance for GIFT, is explained in the next part below.

378 To conclude, based on Fisher's theoretic works, the traditional GWAS method has been optimized to map  
379 SNPs that, using GIFT, would draw paraboloid genetic paths (see Fig.1C). The potential novelty using GIFT  
380 resides in its ability to provide new information and detect relatively regular/structured sigmoid genetic  
381 paths that would otherwise not be detected by traditional methods.

382

### 383 **$p_{\text{GIFT}}$ : p-value for GIFT**

384 GIFT and GWAS extract information on genotype-phenotype associations in totally different ways. While  
385 GIFT concentrates on the significance of curves drawn using  $\Delta\theta(i) = \theta(i) - \theta_0(i)$ , GWAS focuses solely  
386 on the significance of difference of averages. However, to compare GIFT to GWAS it is essential to  
387 determine a p-value for GIFT that is exhaustive enough such as to also capture the information that GWAS  
388 provides. To this end a p-value was derived that concentrates on the maximal amplitudes difference of  
389 genetic paths (see Figs.6A-6B).

390 The p-value for GIFT can be understood as follows. Since the number of possible paths is linked to the  
391 number of configuration possible resulting from lodging  $N_+$ ,  $N_0$  and  $N_-$  microstates into a list composed  
392 of  $N = N_+ + N_0 + N_-$  components, the number of possible paths is,  $N_{\text{path}}^0 = \frac{N!}{N_+!N_0!N_-!}$ . Let us now divide  
393 the genetic paths into regions,  $\Delta i_1$ ,  $\Delta i_2$  and  $\Delta i_3$  as shown in Figs.6A-6B. As the number of microstates of  
394 each sort can be determined in each region using an adequate algorithm, then the total number of

395 possible genetic paths in this first, second and third regions are, respectively,  $N_1 = \frac{\Delta i_1!}{(n_+)_{i_1}!(n_0)_{i_1}!(n_-)_{i_1}!}$ ,  $N_2 =$   
396  $\frac{\Delta i_2!}{(n_+)_{i_2}!(n_0)_{i_2}!(n_-)_{i_2}!}$  and,  $N_3 = \frac{\Delta i_3!}{(n_+)_{i_3}!(n_0)_{i_3}!(n_-)_{i_3}!}$ , where  $(n_q)_p$  is the number of microstate of type  $q$  in the  $p^{\text{th}}$   
397 region,  $q \in \{+, 0, -\}$  and  $p \in \{1, 2, 3\}$ . Consequently, the probability of a genetic path in this context is,  
398  $\hat{p}_{\text{GIFT}} = N_1 N_2 N_3 / N_{\text{path}}^0$ . Using the null hypothesis simulations shown in Fig.3 based on the theoretic SNPs  
399 given in Table-1,  $\hat{p}_{\text{GIFT}}$  may be determined for each genetic path simulated. Its statistic plotted in Fig.6C  
400 for each SNP demonstrates very little variations across SNPs or when the sample size changes by a factor  
401 two. Based on this observation confidence intervals were determined for all SNPs by averaging the  $\hat{p}_{\text{GIFT}}$   
402 values obtained. The upper and lower red dashed lines represent the 99% and 95% confidence intervals.  
403 To consider the false discovery rate (FDR) and adjust p-values to remove type-I errors,  $\hat{p}_{\text{GIFT}}$ -values in  
404 Fig.6C were corrected using the Benjamini-Hochberg procedure leading to a new set of adjusted, i.e.,  
405 reduced, p-values, noted  $p_{\text{GIFT}}$  (see Fig.6D), that may be used to determine the true significance of DNA  
406 variants (SNPs). Returning to Table-2 the numerical value of  $p_{\text{GIFT}}$  was determined for the genetic paths  
407 shown in Fig.4 demonstrating that GIFT can extract information when sigmoid genetic paths are involved  
408 while traditional GWAS is unable to do so.  
409 Armed with  $p_{\text{GIFT}}$  an analysis of datasets can now be performed.

#### 410 **Comparison between GWAS and GIFT considering the bone area of the ischium (BAI) as phenotype** 411 **(Dataset-1)**

412 The first dataset (Dataset-1) analysed 567 pedigree-recorded Scottish Blackface lambs concentrating on  
413 the bone areas of the ischium measured in  $\text{mm}^2$  from cross-sectional CT scans (13). After adjusting  
414 phenotypic values, the work demonstrated a clear involvement of chromosome 6 as shown in Fig.7A. The  
415 genome-wide significant thresholds applied for GWAS in Fig.7A correspond to Bonferroni corrections at  
416 1% (upper red dashed line) and 5% (lower dashed red line) determined by using independent SNPs only.  
417 Formally a 1% (resp. 5%) Bonferroni correction is given by,  $-\text{Log}_{10}(0.01/N_{\text{ind-SNPs}})$  (resp.  
418  $-\text{Log}_{10}(0.05/N_{\text{ind-SNPs}})$ ) where  $N_{\text{ind-SNPs}} = 10433$  is the number of independent SNPs. Using its own  
419 thresholds (Fig.6D) GIFT was applied using the same set of phenotypic residuals. Figs.7A-7B demonstrate  
420 the results obtained by GWAS and GIFT using Manhattan plots.

421 The significance threshold by GIFT was defined by a null hypothesis using theoretic SNPs. To demonstrate  
422 that the theoretic results obtained from Fig.6D are transferrable to 'real' SNPs (Fig.7B), namely that the  
423 significant SNPs obtained in Fig.7B have null hypotheses with similar properties like those shown in Fig.6D,  
424 each significant SNP (Fig.7B) had its genetic path randomly permuted a thousand times to determine  
425 the distribution of  $-\text{Log}_{10}(p_{\text{GIFT}})$ -values corresponding to their null hypothesis. Results show that the  
426 null hypotheses are remarkably similar across SNPs and that the threshold determined using theoretic  
427 SNPs (Fig.6D) holds when 'real' SNPs are used (Supplemental S5).

428 Overall, Fig.7A and Fig.7B demonstrate that there is an agreement between GWAS and GIFT that  
429 chromosome 6 is involved. However, differences exist that are shown through the involvement of several  
430 chromosomes when GIFT is used. Considering the thresholds involved, for GWAS the phenotype studied  
431 may be considered as a sort of 'single gene trait' while for GIFT, the phenotype looks very much like a  
432 'complex trait' involving more chromosomes than chromosome 6. Detailed information of all significant  
433 SNPs by GWAS or GIFT is given in Supplemental S6.

434 Concentrating on Chromosome 6 to address the overlap of information provided by GIFT and GWAS, a  
435 Venn-diagram including highly significant SNPs only, namely SNPs beyond the upper red dashed-line in  
436 Figs.7A-7B, was plotted. The Venn-diagram (Fig.7C) reveals that most SNPs deemed significant by GWAS  
437 were also deemed significant by GIFT. Curiously, only one SNP seemed highly significant by GWAS but  
438 irrelevant for GIFT. As  $p_{\text{GIFT}}$  was designed to collect exhaustive information from GWAS, the SNP was  
439 identified (OAR6\_40311379) and its genetic path, i.e., its  $\Delta\theta(i)$ , plotted (Fig.7D-left) together with its

440 GWAS-representations (Fig.7D-right). The genetic path, being erratic of relatively small amplitude and  
441 crossing several times the axis of positions, did not display any obvious ‘parabolic or sigmoidal’  
442 associations at first sight, in turn justifying its small  $p_{\text{GIFT}}$ -value. The GWAS-representation of  
443 OAR6\_40311379 however, demonstrated the absence of microstate ‘-1’ as well as a near overlap of  
444 microstates ‘0’ and ‘+1’ further demonstrated by the similarities between their boxplots, suggesting the  
445 occurrence of a false-positive. To confirm this a comparison of phenotypic means for the microstates ‘0’  
446 and ‘+1’ was performed returning a t-test value of 1.1485 ( $p$ -value of 0.2512), confirming the presence of  
447 a false-positive.

448 In order to assess the overlap of information between GWAS and GIFT we plotted in Fig.7E the first 100  
449 more significant SNPs detected by GIFT and GWAS. Results confirm an overlap of SNPs associated with  
450 the phenotypic residuals for large values of  $p_{\text{GIFT}}$  and  $p_{\text{GWAS}}$  (see purple dots in  $Q_2$  in Fig.7E).  
451 Interestingly, two SNPs considered as significant by GWAS (two blue dots in  $Q_2$ ) were not by GIFT. That is  
452 because the  $p_{\text{GIFT}}$ -values for these dots were less than other SNPs detected by GIFT. As already stated  
453 above many SNPs from other chromosomes were considered significant by GIFT that were not by GWAS  
454 (see red dots in  $Q_4$ ). Finally, the quadrant  $Q_1$  in Fig.7E confirms that OAR6\_40311379, i.e., the false  
455 positive detected by GWAS, is a standalone SNP among the 100 SNPs for which  $p_{\text{GWAS}} > p_{\text{GIFT}}$ . Finally,  
456 the biotype of significant SNPs on Chromosome 6 for GIFT and GWAS are also presented in Fig.7F.

457 The primary conclusion provided by Figs.7A-F is that, when compared to GWAS, GIFT returns substantially  
458 more genetic information.

459 However, a central question concerns the genetic pertinence of the significant SNPs obtained by GIFT. As  
460 GIFT has been designed with the aim to increase the investigative power of biological datasets, we may  
461 assume that the significant SNPs obtained by GIFT once translated into gene names should underline some  
462 level of non-random gene-gene interactions. The latter point is particularly relevant since GIFT is expected  
463 to detect regulatory variants (c.f. sigmoidal genetic paths). To assess this point we performed an  
464 enrichment analysis based on gene names using the String database, which helps determine known and  
465 predicted protein-protein interactions. In order to apply String the significant SNPs obtained using GWAS  
466 and GIFT were mapped to the reference sheep genome assembly from ensembl (Oar\_v3.1) to obtain the  
467 gene names. Using those gene names String analyses were performed for GWAS and GIFT using a  
468 minimum required interaction score of 0.4. Fig.7G and Fig.7H show the networks obtained. With  
469 enrichment  $p$ -values for GWAS and GIFT of 0.176 and 0.00008, respectively, these results confirm that the  
470 set of genes determined by GIFT have more interactions among themselves than what would be expected  
471 for a random set of genes of the same size and degree distribution drawn from the genome. Namely that  
472 GIFT increases the investigative power of biological datasets.

473 At present, we do not know how the whole information provided by GIFT may inform on the putative  
474 biology of the phenotype studied (BAI). As it turns out, a full validation of the information provided by  
475 GIFT on Dataset-1 would require an in-depth mutational/deletion/insertion/gene-editing analyses in live  
476 animals, extending beyond the scope of this present article.

477 To demonstrate the relevance of the information provided by GIFT we decided to challenge GIFT using a  
478 different dataset (Dataset-2) concentrating on a complex trait related to 1C-metabolism.

479

#### 480 **Comparison between GWAS and GIFT considering 1C-metabolites as phenotype (Dataset-2)**

481 Dataset-2 concerns biochemical data which seeks to identify risk allele variants in genes whose products  
482 direct a specific series of metabolic pathways, known as one carbon (1C) metabolism (Fig.2). The  
483 significance of 1C metabolism is that it is a complex trait involving a series of interlinking metabolic  
484 pathways that provide 1C units (methyl groups) for the synthesis and methylation of biological molecules.  
485 After 1% and 5% Bonferroni corrections for GWAS and the Benjamini-Hochberg procedure applied to GIFT,  
486 the Manhattan plots were obtained (Fig.8A). Note that the number of independent SNPs in this case is  
487 624 (out of 3923 SNPs from the gene array). Fig.8A demonstrates clearly that the informational power of

488 GWAS is less than that of GIFT. Finally, in Fig.8B we provide the biotypes of the most significant SNPs  
489 shown by the upper red dashed lines obtained using GIFT. Detailed genetic information of the most  
490 significant SNPs obtained using GIFT is provided in Supplemental S8.

491 Since the gene array was synthesized using SNPs from known genes involved in 1C metabolism, the  
492 relevance of String analyses (i.e., enrichment p-values) would be minimal and of little interest.

493 Besides validating that GIFT may extract more information from genotype-phenotype datasets, it is worth  
494 underlying the biological importance and novelty of results obtained. One carbon metabolism in sheep is  
495 comparable to that in humans. The significance of 1C metabolism is that it is a complex trait involving a  
496 series of interlinking metabolic pathways that provide 1C units (methyl groups) for the synthesis and  
497 methylation of chromatin among other molecules (15). S-adenosylmethionine (SAM) is a potent methyl  
498 donor within these cycles and serves as the principal substrate for methylation of DNA, associated  
499 proteins, and RNA. It was previously demonstrated in sheep, cattle, rodent and human studies that  
500 disrupting these cycles during early pregnancy, by either dietary means (i.e., reducing dietary vitamin B12,  
501 folate, choline and/or methionine), or through exposure to environmental chemicals such as cigarette  
502 smoking, can lead to epigenetic dysregulation and impaired foetal development with long-term  
503 consequences for offspring cardiometabolic health (22–25). It was also advocated that interindividual and  
504 ethnic variability in epigenetic gene regulation arises because of single-nucleotide polymorphisms (SNPs)  
505 within 1C genes, associated epigenetic regulators, and differentially methylated target DNA sequences  
506 (15). However, information concerning the nature and extent of interactions between parental genotype,  
507 diet and EC exposure was, until now, limited to just a few 1C genes in humans (15). Consequently, data  
508 obtained by the current study provide new evidence concerning significant genetic variants in 1C-  
509 metabolism and directly associated metabolic genes and epigenetic regulators that rely on SAM as the  
510 methyl donor, potentially applicable to the human species.

## 511 DISCUSSION

512 While statistical association methods should not favor any biases when analyzing datasets, the way they  
513 are built mathematically is often indicative of a particular way of thinking. For example, with GWAS the  
514 phenotype is decomposed onto more fundamental sub-distributions characterized by the distribution of  
515 microstates (see Fig.1A). This approach underlines a sort of bottom-up approach that, within a  
516 reductionist framework, defines genes as biological agents controlling the phenotype aligned with the  
517 ‘Neo-Darwinian synthesis’. However, nothing prevents considering the opposite as far as statistical  
518 association methods are involved, and GIFT uses this degree of freedom. By using the full range of  
519 phenotypic information, GIFT transforms a random or disordered string of microstates (the straight line  
520 in the asymptotic limit seen in Fig.1C or Fig.3A) into an ‘ordered’ configuration of microstates (see Fig1C  
521 or Figs.4C-4F), in turn providing the signature of a genotype-phenotype association. Accordingly, since the  
522 phenotypic information controls the configuration of microstates it is a top-down approach, which turns  
523 out to be remarkably sensitive. GIFT has been estimated to be ~1000 more sensitive than GWAS (11).

524 There are three main reasons as to why GIFT is more sensitive. The first is that GIFT determines the  
525 significance of curves composed of an entire population of datapoints. As curves provide a greater level  
526 of significance than considering differences between microstate/phenotypic averages/variances as  
527 advocated by GWAS, hence GIFT is statistically more powerful. The second reason is that the null  
528 hypothesis for GIFT, namely  $\theta_0(i)$ , is contained in the definition of  $\Delta\theta(i)$  and is therefore specific to the  
529 genome position, or SNP, studied. With GIFT there are as many null hypotheses as SNPs. This contrasts  
530 with GWAS defining a null-hypothesis valid for all SNPs at the population level when the average of  
531 microstate distributions overlap. Consequently, the discriminative power of GIFT is amplified. The third  
532 reason is that GIFT is simpler than GWAS. Indeed, based on R.A. Fisher’s seminal work, GWAS is based on  
533 a complex theory that seeks to determine genotype-phenotype associations on one hand (aim 1), and the

534 heritability of phenotypes/traits studied on the other (aim 2). To achieve those two aims, the GWAS  
535 approach relies on frequentist probability to determine the validity of statistical inferences giving the  
536 notions of average and variance fundamental meanings related to aim 1 and 2, respectively. However,  
537 because average and variance are antinomic it is nearly impossible to have a clear picture of associations  
538 (size effects) since the noise (variance/heredity) blurs the average(s). On the other hand, by concentrating  
539 on genetic paths (curves) GIFT determines a global association. This does not mean that GIFT rules out the  
540 notions of size effect, dominance, and heritability, on the contrary, it encapsulates them under the generic  
541 notion of phenotypic field, i.e., size effect, dominance and heritability can be rederived from the  
542 phenotypic field. The term 'field' in the acronym GIFT is used to explain the disorder-order transition in  
543 the string of microstates using an analogy related to physics field theory, see (11, 12) for more details.

544 Finally, it is important to reframe GIFT within current debates in the field of biology. With GIFT it is the  
545 (information on the) phenotype that selects which SNP is required for its subsistence and it is interesting  
546 to note that, at the conceptual level and as a top-down approach, GIFT has some familiarity with the  
547 notion of phenotypic plasticity. Phenotypic plasticity refers to the ability of phenotypes to respond to a  
548 change in the environment favoring a divergence from the ancestor phenotype. As the phenotype relies  
549 on traits (modules), the responsiveness to any new input(s) must involve a re-organization of the  
550 phenotype architecture by allowing phenotypic sub-components (modular traits) to adapt the changes  
551 (26). Namely that genetic accommodation linked to a standing pool of genetic variations characterizing  
552 any trait is central to phenotypic plasticity that, through persistence, may genetically assimilate the new  
553 architecture (selection) (26, 27). In this context the top-down method GIFT, which is essentially a  
554 phenotype-genotype (and not genotype-phenotype) association method, can pull out any standing genes  
555 awaiting to be used by phenotypes.

556 To conclude, we provide evidence that GIFT enhances the investigative power of biological datasets.  
557 Additionally, we provide evidence also for the need to rethink the conceptual bases of genotype-  
558 phenotype association methods, such as use more information from the whole biodiversity of data.

559

## 560 **DATA AVAILABILITY**

561 Data including supplementary materials are available using the link:

## 562 **SUPPLEMENTAL MATERIAL**

563 S1 provides the raw data for Dataset-1; S2 provides the raw data for Dataset-2. S3 provides the  
564 statistical summary for the phenotypic adjustment prior to running GWAS on Dataset-2; S4 provides the  
565 code to obtain Fig.3 and Fig.6C. S5 represents the permutation analysis of significant SNPs obtained by  
566 GIFT from dataset-1. S6 represents the list of significant SNPs obtained by GWAS and GIFT when applied  
567 on Dataset-1; S7 provides the code to obtain Fig.4, Fig5 and Fig.7. S8 provides the list of significant SNPs  
568 by GIFT for Dataset-2. S9 provides the code to obtain Fig.8.

## 569 **ACKNOWLEDGMENTS**

570 The authors would like to thank Dr Barbara Bravi, Dr Wing-Yee Kwong and Dr Dongfang Li for useful  
571 discussions and/or technical assistance.

572 Present address of B.B.: Imperial College, London, UK.

573 Present address of W-Y.K. & D.L.: University of Nottingham, Sutton Bonington, UK.



574

## 575 GRANTS

576 This work was supported by the Biotechnology and Biological Sciences Research Council (BBSRC) Industrial  
577 Partnership Award with the Agriculture and Horticulture Development Board, Meat Promotion Wales and  
578 Agrisearch [BB/K017810/1; BB/K017993/1], and National Institutes of Health (R01 ES030374/ES/NIEHS  
579 NIH HHS/United States). P.K. is currently supported by a Doctoral Scholarship from the EPHE, Sorbonne  
580 University in collaboration with the University of Nottingham. C.E.C. was in receipt of a BBSRC Doctoral  
581 Training Partnership scholarship (1796056) and A.H.B. was in receipt of a scholarship from The Perry  
582 Foundation.

## 583 DISCLOSURES

584 Authors declare no conflict of interest, financial or otherwise.

## 585 AUTHOR CONTRIBUTIONS

586 CR conceptualized GIFT; CR & JW formalized GIFT; PK coded GIFT simulations; Dataset-2 was designed  
587 and obtained by KDS, OM, AHB, CEC, JX, DAB, RDE; Dataset-1 and Dataset-2 were analysed by PK, CR,  
588 JW, KDS, OM, AP; paper was written by CR and KDS, and proofread by CR and KDS.

## 589 REFERENCES

590 1. Burton PR, Clayton DG, Cardon LR, Craddock N, Deloukas P, Duncanson A, Kwiatkowski DP,  
591 McCarthy MI, Ouwehand WH, Samani NJ, Todd JA, Donnelly P, Barrett JC, Burton PR, Davison D,  
592 Donnelly P, Easton D, Evans D, Leung H-T, Marchini JL, Morris AP, Spencer CCA, Tobin MD, Cardon  
593 LR, Clayton DG, Attwood AP, Boorman JP, Cant B, Everson U, Hussey JM, Jolley JD, Knight AS, Koch  
594 K, Meech E, Nutland S, Prowse CV, Stevens HE, Taylor NC, Walters GR, Walker NM, Watkins NA,  
595 Winzer T, Todd JA, Ouwehand WH, Jones RW, McArdle WL, Ring SM, Strachan DP, Pembrey M,  
596 Breen G, St Clair D, Caesar S, Gordon-Smith K, Jones L, Fraser C, Green EK, Grozeva D, Hamshere  
597 ML, Holmans PA, Jones IR, Kirov G, Moskvina V, Nikolov I, O'Donovan MC, Owen MJ, Craddock N,  
598 Collier DA, Elkin A, Farmer A, Williamson R, McGuffin P, Young AH, Ferrier IN, Ball SG, Balmforth AJ,  
599 Barrett JH, Bishop DT, Iles MM, Maqbool A, Yuldasheva N, Hall AS, Braund PS, Burton PR, Dixon RJ,  
600 Mangino M, Stevens S, Tobin MD, Thompson JR, Samani NJ, Bredin F, Tremelling M, Parkes M,  
601 Drummond H, Lees CW, Nimmo ER, Satsangi J, Fisher SA, Forbes A, Lewis CM, Onnie CM, Prescott  
602 NJ, Sanderson J, Mathew CG, Barbour J, Mohiuddin MK, Todhunter CE, Mansfield JC, Ahmad T,  
603 Cummings FR, Jewell DP, Webster J, Brown MJ, Clayton DG, Lathrop GM, Connell J, Dominiczak A,  
604 Samani NJ, Marcano CAB, Burke B, Dobson R, Gungadoo J, Lee KL, Munroe PB, Newhouse SJ,  
605 Onipinla A, Wallace C, Xue M, Caulfield M, Farrall M, Barton A, and Genomics (BRAGGS) TB in RG,  
606 Bruce IN, Donovan H, Eyre S, Gilbert PD, Hider SL, Hinks AM, John SL, Potter C, Silman AJ, Symmons  
607 DPM, Thomson W, Worthington J, Clayton DG, Dunger DB, Nutland S, Stevens HE, Walker NM,  
608 Widmer B, Todd JA, Frayling TM, Freathy RM, Lango H, Perry JRB, Shields BM, Weedon MN,  
609 Hattersley AT, Hitman GA, Walker M, Elliott KS, Groves CJ, Lindgren CM, Rayner NW, Timpson NJ,  
610 Zeggini E, McCarthy MI, Newport M, Sirugo G, Lyons E, Vannberg F, Hill AVS, Bradbury LA, Farrar C,  
611 Pointon JJ, Wordsworth P, Brown MA, Franklyn JA, Heward JM, Simmonds MJ, Gough SCL, Seal S,  
612 Susceptibility Collaboration (UK) BC, Stratton MR, Rahman N, Ban M, Goris A, Sawcer SJ, Compston  
613 A, Conway D, Jallow M, Newport M, Sirugo G, Rockett KA, Kwiatkowski DP, Bumpstead SJ, Chaney  
614 A, Downes K, Ghorri MJR, Gwilliam R, Hunt SE, Inouye M, Keniry A, King E, McGinnis R, Potter S,

- 615 Ravindrarajah R, Whittaker P, Widden C, Withers D, Deloukas P, Leung H-T, Nutland S, Stevens HE,  
616 Walker NM, Todd JA, Easton D, Clayton DG, Burton PR, Tobin MD, Barrett JC, Evans D, Morris AP,  
617 Cardon LR, Cardin NJ, Davison D, Ferreira T, Pereira-Gale J, Hallgrimsdóttir IB, Howie BN, Marchini  
618 JL, Spencer CCA, Su Z, Teo YY, Vukcevic D, Donnelly P, Bentley D, Brown MA, Cardon LR, Caulfield  
619 M, Clayton DG, Compston A, Craddock N, Deloukas P, Donnelly P, Farrall M, Gough SCL, Hall AS,  
620 Hattersley AT, Hill AVS, Kwiatkowski DP, Mathew CG, McCarthy MI, Ouwehand WH, Parkes M,  
621 Pembrey M, Rahman N, Samani NJ, Stratton MR, Todd JA, Worthington J, The Wellcome Trust Case  
622 Control Consortium, Management Committee, Data and Analysis Committee, UK Blood Services  
623 and University of Cambridge Controls, 1958 Birth Cohort Controls, Bipolar Disorder, Coronary  
624 Artery Disease, Crohn's Disease, Hypertension, Rheumatoid Arthritis, Type 1 Diabetes, Type 2  
625 Diabetes, Tuberculosis, Ankylosing Spondylitis, Autoimmune Thyroid Disease, Breast Cancer,  
626 Multiple Sclerosis, Gambian Controls, DNA G Data QC and Informatics, Statistics, Primary  
627 Investigators. Genome-wide association study of 14,000 cases of seven common diseases and  
628 3,000 shared controls. *Nature* 447: 661–678, 2007. doi: 10.1038/nature05911.
- 629 2. Buniello A, MacArthur JAL, Cerezo M, Harris LW, Hayhurst J, Malangone C, McMahon A, Morales J,  
630 Mountjoy E, Sollis E, Suveges D, Vrousseau O, Whetzel PL, Amode R, Guillen JA, Riat HS, Trevanion  
631 SJ, Hall P, Junkins H, Flicek P, Burdett T, Hindorf LA, Cunningham F, Parkinson H. The NHGRI-EBI  
632 GWAS Catalog of published genome-wide association studies, targeted arrays and summary  
633 statistics 2019. *Nucleic Acids Res* 47: D1005–D1012, 2019. doi: 10.1093/nar/gky1120.
- 634 3. Sudlow C, Gallacher J, Allen N, Beral V, Burton P, Danesh J, Downey P, Elliott P, Green J, Landray M,  
635 Liu B, Matthews P, Ong G, Pell J, Silman A, Young A, Sprosen T, Peakman T, Collins R. UK biobank:  
636 an open access resource for identifying the causes of a wide range of complex diseases of middle  
637 and old age. *PLoS Med* 12: e1001779, 2015. doi: 10.1371/journal.pmed.1001779.
- 638 4. Smith BH, Campbell A, Linksted P, Fitzpatrick B, Jackson C, Kerr SM, Deary IJ, Macintyre DJ,  
639 Campbell H, McGilchrist M, Hocking LJ, Wisely L, Ford I, Lindsay RS, Morton R, Palmer CNA,  
640 Dominiczak AF, Porteous DJ, Morris AD. Cohort Profile: Generation Scotland: Scottish Family  
641 Health Study (GS:SFHS). The study, its participants and their potential for genetic research on  
642 health and illness. *Int J Epidemiol* 42: 689–700, 2013. doi: 10.1093/ije/dys084.
- 643 5. Fisher RA. XXI.—On the Dominance Ratio. *Proceedings of the Royal Society of Edinburgh* 42: 321–  
644 341, 1923. doi: 10.1017/S0370164600023993.
- 645 6. Fisher RA. XV.—The Correlation between Relatives on the Supposition of Mendelian Inheritance.  
646 *Transactions of the Royal Society of Edinburgh* 52: 399–433, 1919. doi:  
647 10.1017/S0080456800012163.
- 648 7. Visscher PM, Goddard ME. From R.A. Fisher's 1918 Paper to GWAS a Century Later. *Genetics* 211:  
649 1125–1130, 2019. doi: 10.1534/genetics.118.301594.
- 650 8. Hivert V, Wray NR, Visscher PM. Gene action, genetic variation, and GWAS: A user-friendly web  
651 tool. *PLoS Genet* 17: e1009548, 2021. doi: 10.1371/journal.pgen.1009548.
- 652 9. Nelson RM, Pettersson ME, Carlborg Ö. A century after Fisher: time for a new paradigm in  
653 quantitative genetics. *Trends Genet* 29: 669–676, 2013. doi: 10.1016/j.tig.2013.09.006.

- 654 10. Wattis JAD, Bray SM, Kyratzi P, Rauch C. Analysis of phenotype-genotype associations using  
655 genomic informational field theory (GIFT). *J Theor Biol* 548: 111198, 2022. doi:  
656 10.1016/j.jtbi.2022.111198.
- 657 11. Rauch C, Kyratzi P, Blott S, Bray S, Wattis J. GIFT: new method for the genetic analysis of small gene  
658 effects involving small sample sizes. *Phys Biol* 20, 2022. doi: 10.1088/1478-3975/ac99b3.
- 659 12. Rauch C, Wattis J, Bray S. On the Meaning of Averages in Genome-wide Association Studies: What  
660 Should Come Next? *Organisms Journal of Biological Sciences* 6: 7–22, 2023. doi: 10.13133/2532-  
661 5876/17811.
- 662 13. Matika O, Riggio V, Anselme-Moizan M, Law AS, Pong-Wong R, Archibald AL, Bishop SC. Genome-  
663 wide association reveals QTL for growth, bone and in vivo carcass traits as assessed by computed  
664 tomography in Scottish Blackface lambs. *Genet Sel Evol* 48: 11, 2016. doi: 10.1186/s12711-016-  
665 0191-3.
- 666 14. Clare CE, Pestinger V, Kwong WY, Tutt DAR, Xu J, Byrne HM, Barrett DA, Emes RD, Sinclair KD.  
667 Interspecific Variation in One-Carbon Metabolism within the Ovarian Follicle, Oocyte, and  
668 Preimplantation Embryo: Consequences for Epigenetic Programming of DNA Methylation. *Int J*  
669 *Mol Sci* 22, 2021. doi: 10.3390/ijms22041838.
- 670 15. Clare CE, Brassington AH, Kwong WY, Sinclair KD. One-Carbon Metabolism: Linking Nutritional  
671 Biochemistry to Epigenetic Programming of Long-Term Development. *Annu Rev Anim Biosci* 7:  
672 263–287, 2019. doi: 10.1146/annurev-animal-020518-115206.
- 673 16. Jiang H, Lei R, Ding S-W, Zhu S. Skewer: a fast and accurate adapter trimmer for next-generation  
674 sequencing paired-end reads. *BMC Bioinformatics* 15: 182, 2014. doi: 10.1186/1471-2105-15-182.
- 675 17. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.  
676 *Bioinformatics* 25: 1754–1760, 2009. doi: 10.1093/bioinformatics/btp324.
- 677 18. Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, Del Angel G, Levy-Moonshine A, Jordan T,  
678 Shakir K, Roazen D, Thibault J, Banks E, Garimella KV, Altshuler D, Gabriel S, DePristo MA. From  
679 FastQ data to high confidence variant calls: the Genome Analysis Toolkit best practices pipeline.  
680 *Curr Protoc Bioinformatics* 43: 11.10.1-11.10.33, 2013. doi: 10.1002/0471250953.bi1110s43.
- 681 19. McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GRS, Thormann A, Flicek P, Cunningham F. The Ensembl  
682 Variant Effect Predictor. *Genome Biol* 17: 122, 2016. doi: 10.1186/s13059-016-0974-4.
- 683 20. Xu J, Clare CE, Brassington AH, Sinclair KD, Barrett DA. Comprehensive and quantitative profiling of  
684 B vitamins and related compounds in the mammalian liver. *J Chromatogr B Analyt Technol Biomed*  
685 *Life Sci* 1136: 121884, 2020. doi: 10.1016/j.jchromb.2019.121884.
- 686 21. Boyle EA, Li YI, Pritchard JK. An Expanded View of Complex Traits: From Polygenic to Omnigenic.  
687 *Cell* 169: 1177–1186, 2017. doi: 10.1016/j.cell.2017.05.038.
- 688 22. Maloney CA, Hay SM, Young LE, Sinclair KD, Rees WD. A methyl-deficient diet fed to rat dams  
689 during the peri-conception period programs glucose homeostasis in adult male but not female  
690 offspring. *J Nutr* 141: 95–100, 2011. doi: 10.3945/jn.109.119453.

- 691 23. Sinclair KD, Allegrucci C, Singh R, Gardner DS, Sebastian S, Bispham J, Thurston A, Huntley JF, Rees  
692 WD, Maloney CA, Lea RG, Craigon J, McEvoy TG, Young LE. DNA methylation, insulin resistance,  
693 and blood pressure in offspring determined by maternal periconceptional B vitamin and  
694 methionine status. *Proc Natl Acad Sci U S A* 104: 19351–19356, 2007. doi:  
695 10.1073/pnas.0707258104.
- 696 24. Drake AJ, O’Shaughnessy PJ, Bhattacharya S, Monteiro A, Kerrigan D, Goetz S, Raab A, Rhind SM,  
697 Sinclair KD, Meharg AA, Feldmann J, Fowler PA. In utero exposure to cigarette chemicals induces  
698 sex-specific disruption of one-carbon metabolism and DNA methylation in the human fetal liver.  
699 *BMC Med* 13: 18, 2015. doi: 10.1186/s12916-014-0251-x.
- 700 25. Rubini E, Snoek KM, Schoenmakers S, Willemsen SP, Sinclair KD, Rousian M, Steegers-Theunissen  
701 RPM. First Trimester Maternal Homocysteine and Embryonic and Fetal Growth: The Rotterdam  
702 Periconception Cohort. *Nutrients* 14, 2022. doi: 10.3390/nu14061129.
- 703 26. West-Eberhard MJ. Developmental plasticity and the origin of species differences. *Proc Natl Acad*  
704 *Sci U S A* 102 Suppl 1: 6543–6549, 2005. doi: 10.1073/pnas.0501844102.
- 705 27. Palmer AR. Symmetry Breaking and the Evolution of Development. *Science* 306: 828–833, 2004.  
706 doi: 10.1126/science.1103707.

707

## 708 **FIGURE LEGENDS**

709 **Figure 1: (A)** For diploid organisms and for a binary (bi-allelic, A or a) genetic marker, any microstate  
710 (genotype) can only take three values that we shall write as ‘+1’, ‘0’ and ‘-1’ corresponding to genotypes  
711 aa, Aa and AA, respectively. The genotypes are color-coded to facilitate the representation of GIFT (+1: aa  
712 (red), 0: aA/Aa (black) and -1: AA (blue)). GWASs rely on probability density functions formed through the  
713 grouping of data into bins/categories. The phenotype distribution density function (A-top left) is then  
714 decomposed onto the distribution density function of genetic microstates (A-top right) for every single  
715 nucleotide polymorphism (SNP). Using an analysis of averages and variances such decomposition  
716 determines whether the SNP studied is associated with the phenotype by comparing the average and  
717 variances of distributions. Repeating the same operation for every SNP in the genome permits to map  
718 genotype to phenotype. However, as more precise inferences can only come with, and are only legitimized  
719 by, a reduction in the width of categories, larger sample sizes are needed. To overcome this issue one way  
720 to proceed is to deconstruct density functions and wonder what would happen if one were able to reduce  
721 the width of categories, that is increasing the precision in the measurement of the phenotype or  
722 equivalently getting access to the whole information of datasets, without changing the sample sizes (**A**  
723 from top-to-bottom). The mathematical object that emerges is then a coloured barcode that is a list of  
724 microstates that can be analysed precisely by GIFT. (**B**) Such barcode can be obtained simply at the  
725 practical level through field studies. Assume a flock of sheep has been genotyped and that their phenotype  
726 has been measured sufficiently precisely such as to exclude the possibility that any two phenotypic values  
727 are identical. In the figure the magnitude of the phenotypic value for each sheep is characterised by the  
728 (unique) ‘size’ of the sheep. The barcode is obtained by ranking animals as a function of the magnitude of  
729 their phenotypic values (configuration ① in Fig.1B). The null hypothesis is obtained via the random  
730 ranking of sheep that is equivalent to a lack of information on phenotypic values (configuration ② in  
731 Fig.1B). As GWAS works on phenotypic residual values after adjusting for fixed/environmental effects a  
732 similar barcode can be generated considering the magnitude of residual phenotypic values. (**C**) GIFT

733 proceeds by plotting the cumulative sum of microstates as a function of the position in the list generating  
734 a curve called genetic path that is represented by  $\theta(i)$  in Fig.1C and is unique to the SNP considered. While  
735 the curve  $\theta(i)$  does not provide any significant information on its own, one may generate, for the same  
736 SNP, a curve (genetic path) corresponding to a sort of null hypothesis when ranking the phenotype does  
737 not bring any informational value. This is possible by scrambling (permutating) the string of microstates  
738 an infinite number of times. It is then possible to show that, in the asymptotic limit, the null hypothesis  
739 returns a straight line, noted  $\theta_0(i)$  (Fig.1C) out of which inferences may be suggested regarding potential  
740 association between the genotype and the phenotype by comparing  $\theta_0(i)$  to  $\theta(i)$ . Note, the simulation  
741 shown in **(A)** adhering to Fisher seminal model is based on a constant sample size of 1000 involving an  
742 arbitrary normally distributed phenotype of mean and variance 68 and 4 units, respectively. Each  
743 microstate is normally distributed with a gene effect identical to the standard deviation of the phenotype  
744 but without dominance. The frequency of the genotypes aa (red), Aa/aA (Grey) and AA (blue) are 64%,  
745 32% and 4%, respectively and within Hardy-Weinberg ratio.

746

747 **Figure 2:** Linked methionine and propionate metabolism adapted from Clare et al. (2019) where all  
748 metabolites studied for this study are in red. The methionine cycle facilitates the re-methylation of  
749 homocysteine (Hcy) to methionine (Met) and ultimately S-adenosylmethionine (**SAM**) with methyl ( $\text{CH}_3$ )  
750 groups donated either from folate (5-mTHF) or betaine (trimethylglycine; **TMG**), thus leading to the  
751 formation of dimethylglycine (**DMG**). Methylcobalamin (**mb12**) serves as a cofactor for the reduction of  
752 the inactive form of methionine synthase to its active state (MTR), which then transfers a methyl group  
753 from 5-mTHF to Hcy. The linked metabolism of propionate (**PPA**) to succinate (an intermediary metabolite  
754 in the tricarboxylic cycle) requires adenosylcobalamin (**ab12**), which serves as a cofactor for  
755 methylmalonyl-CoA-mutase (MUT) leading to the generation of succinyl-CoA and methylmalonic acid  
756 (**MMA**) in this pathway. Other intermediary metabolites and enzymes listed: glycine (Gly), sarcosine (Sar),  
757 S-adenosylhomocysteine (SAH), tetrahydrofolate (THF), serine (Ser), cystathionine (Cth), cysteine (Cys),  
758 alpha-ketobutyrate ( $\alpha$ -KB), methylmalonic acid (MMA); Betaine homocysteine methyltransferase (BHMT),  
759 Methionine adenosyl-transferase (MAT), Glycine methyl-transferase (GNMT), Adenosyl-homocysteinase  
760 (AHCY), Cystathionine beta-synthase (CBS), cystathionine gamma-lyase (Cth).

761

762 **Figure 3: (A-left panel)** Simulations of genetic paths corresponding to null hypotheses using GIFT as a  
763 method. The data used for the simulation are given in Table-1. **(A-right panel)** Simulations of genetic paths  
764 corresponding to null hypotheses when the sample size is divided or multiplied by a factor two. **(B)**  
765 Representation of  $\Delta\theta_0(i) = \theta_0(i) - \langle\theta_0(i)\rangle$  for the microstates data as given in Table-1. **(C)** Plots of the  
766 standard deviation normalised by the square root of the sample size and where the position is also  
767 normalised by the sample size. The code for the simulations is given in Supplemental S4.

768

769 **Figure 4: A** sample of genetic paths selected from Dataset-1. The details of the different SNPs displayed  
770 are given in Table 2.

771 **Figure 5:** Analysis of averages (GWAS) for SNP1-6 (see Fig.4 and Table-2). Values for the size/gene effects  
772 (a) and dominances (d) are given in Table 2.

773 **Figure 6:** To provide a p-value extracting genotype-phenotype associations in an exhaustive manner for  
774 both GWAS and GIFT a method concentrating on the largest and smallest extreme values of the genetic  
775 path was focused upon. This method can be applied to paraboloid (GWAS or GIFT-like) **(A)** and sigmoid  
776 (GIFT-like) **(B)** genetic paths. The overall idea consists in determining how many paths  $N_1$ ,  $N_2$  and  $N_3$  can  
777 be generated from the respective interval of positions  $\Delta i_1$ ,  $\Delta i_2$  and  $\Delta i_3$  given that the constraints for the



778 extrema are  $\Phi_1$  and  $\Phi_2$ . Then a p-value ( $\hat{p}_{GIFT}$ ) can be determined as seen in the text. **(C)** Using  
779 simulations (K=1000 replicates) a statistic of  $\hat{p}_{GIFT}$  for the null hypothesis can be generated using theoretic  
780 SNPs (Table-1). Simulations demonstrate that  $\hat{p}_{GIFT}$  is relatively independent of the microstate's  
781 frequencies upon which a 99% (upper dashed line) and 95% (lower dashed line) interval confidences can  
782 be generated. **(D)**  $\hat{p}_{GIFT}$ -values were adjusted to consider FDR using Benjamini-Hochberg procedure  
783 leading to a new set of  $p_{GIFT}$ -values. The code for the simulations is given in Supplemental S4.

784 **Figure 7:** Manhattan plots based on p-values obtained by GWAS **(A)** and GIFT **(B)** demonstrating significant  
785 differences between the methods concerning potential genotype-phenotype associations. Note that the  
786 presence of a chromosome '0' results from the fact that some SNPs identified by (Matika et al., 2016)  
787 were not allocated to specific chromosomes/genomic positions due to lack of information at the time. A  
788 fathom chromosome (chromosome zero) was created to allocate those SNPs. **(C)** Venn-diagram  
789 representing the most significant SNPs by GWAS and GIFT. One SNP (OAR6\_40311379) demonstrated a  
790 large p-value for GWAS and a small p-value for GIFT. A representation of its genetic path **(D-left)** did not  
791 underscore any 'parabolic' or 'sigmoidal' associations. As it turned out this SNP was a false-positive by  
792 GWAS since the difference between the phenotypic means was not significant **(D-right)**. **(E)** The 100  
793 most significant SNPs by GWAS and GIFT were extracted, and their p-values plotted against each other.  
794 The dashed lines represent the threshold applied for GWAS (blue dashed line) and GIFT (red dashed line).  
795 The SNP OAR6\_40311379 pointed by the black arrow is the single one standing out in  $Q_1$  confirming its  
796 false-positive status. **(F)** Biotypes of the most significant SNPs by GIFT and GWAS. **(G)** String analysis  
797 performed to determine gene networks using significant SNPs by GWAS. **(H)** String analysis performed to  
798 determine gene networks using significant SNPs by GIFT, note that the dashed square underlines mTOR  
799 and FOXO3 determined by GWAS. The code for obtaining Figs.7B, 7C, 7D, 7E is given in Supplemental S7.

800 **Figure 8: (A)** Comparison of the information extracted by GWAS and GIFT using Manhattan plots for the  
801 metabolites presented in red in Fig.2. We recall the acronyms, S-adenosyl methionine (SAM),  
802 methylcobalamin (mB12), adenosylcobalamin (aB12), trimethylglycine (TMG), dimethylglycine (DMG),  
803 propionate (PPA) and methylmalonic acid (MMA). It should be noted that due to inherent difficulty linked  
804 to the measure of metabolite the sample sizes were not similar across metabolites, that is the values for  
805 N differ between the Manhattan plots (SAM: N=344; mB12: N=183; aB12: N=338; DMG: N=338; TMG:  
806 N=340; MMA: N=348; PPA: N=345). **(B)** Biotypes corresponding to the most significant SNPs for each  
807 metabolite determined by GIFT (a detailed list of information concerning those SNPs is given in  
808 supplemental S8). The code for the Manhattan plots and the determination of biotypes is given in  
809 Supplemental S9.

810

811

812 **TABLES**

813 **Table 1:** Theoretic SNPs used to capture the null hypothesis associated with GIFT upon 1000 simulations  
 814 of microstates permutation\*.

SNP NAME	N <sub>+</sub>	N <sub>0</sub>	N <sub>-</sub>	N
SNP1	25	25	515	565
SNP2	25	125	415	565
SNP3	25	225	315	565
SNP4	25	325	215	565
SNP5	25	425	115	565
SNP6	25	525	15	565

815 (\*): The difference between consecutive SNPs in the table is linked to the transfer of 100 microstates from  
 816 the microstates '-1' to the microstate '0' leaving the number of microstates '+' invariant. By permutating  
 817 the microstates '+' and '-' in the table similar plots as those obtained in Fig.3A could have been obtained,  
 818 the only difference would have been the slopes of the average  $\langle \theta_0(i) \rangle$  changing sign.

819

820 **Table 2:** Determination of gene/size effect (a) and dominance (d) for SNP1-6 from dataset-1. The level of  
 821 significance for GIFT and GWAS is colour coded: red=not-significant, green=significant.

CHR	NAME	POSITION	-Log <sub>10</sub> (p <sub>GIFT</sub> )	-Log <sub>10</sub> (p <sub>GWAS</sub> )	N <sub>+</sub>	N <sub>0</sub>	N <sub>-</sub>	a*	d**
9	OAR9_58767921 (SNP1)	56039025	2.7895	0.2735	391	160	16	N/A	N/A
3	s02120 (SNP2)	213625709	2.8893	0.0018	198	291	78	N/A	N/A
6	OAR6_40855809 (SNP3)	36655091	28.5105	9.8639	229	262	76	96.85	-13.01
6	OAR6_38315830 (SNP4)	34256151	20.7541	3.7366	24	222	321	-70.02	-0.05
23	OAR23_35510473 (SNP5)	33556377	19.7239	0.2301	254	260	53	N/A	N/A
25	OAR25_30372586 (SNP6)	29046746	18.5806	1.0692	90	266	211	N/A	N/A

822 (a\*): The gene/size effect is calculated considering the mid-distance between the average values of  
 823 phenotypic residuals of microstates '-1' and '+1'. (d\*\*): The dominance is calculated considering the  
 824 difference between the gene/size effect (a) and the position of the average value of phenotypic residuals  
 825 for the microstate '0'.

Fig.1

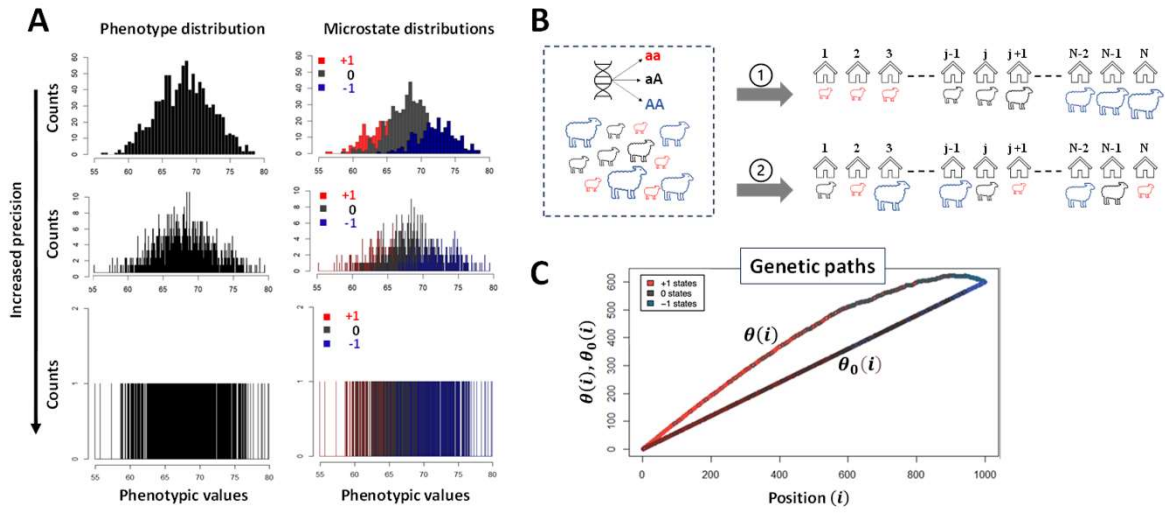


Fig.2

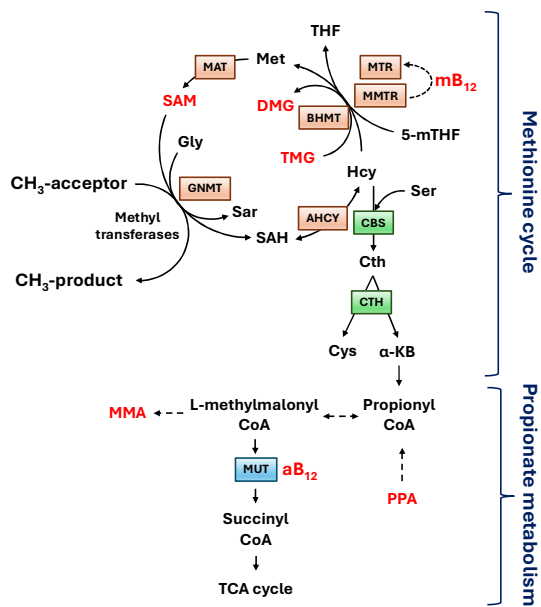


Fig.3

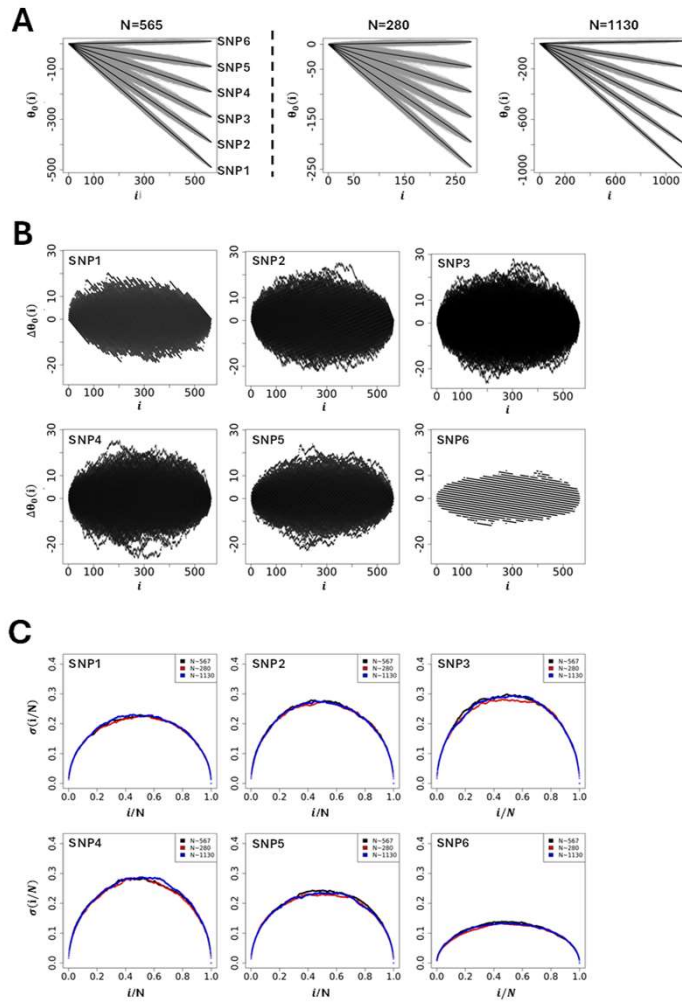


Fig.4

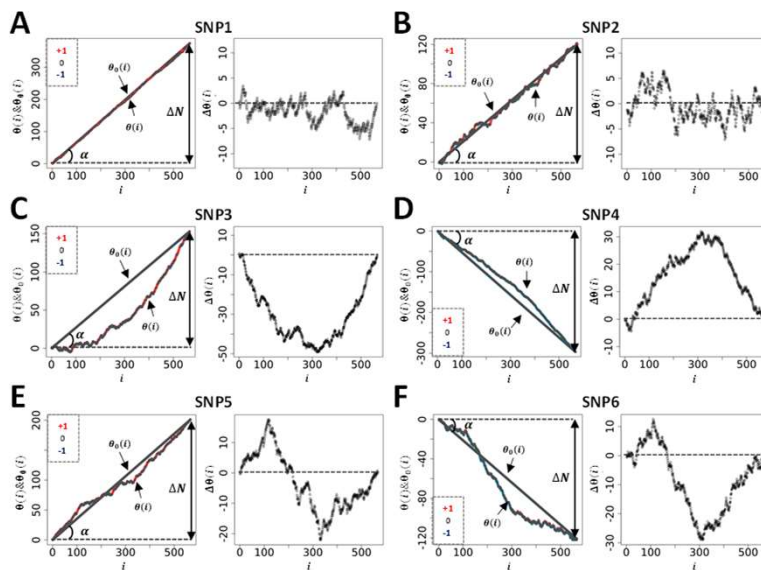


Fig.5

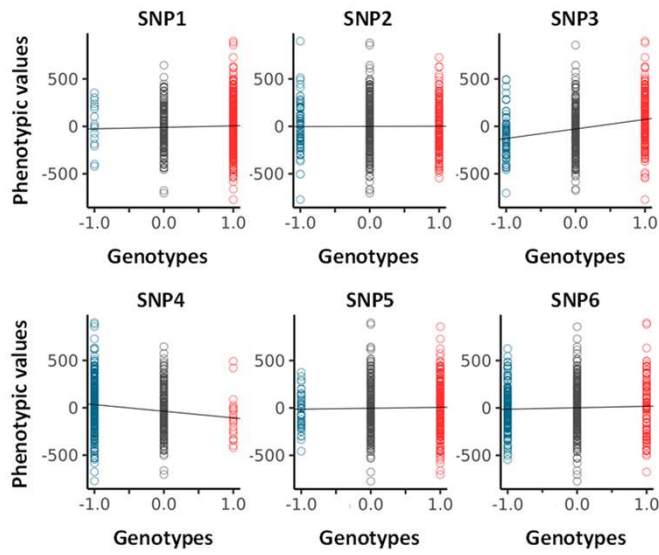


Fig.6

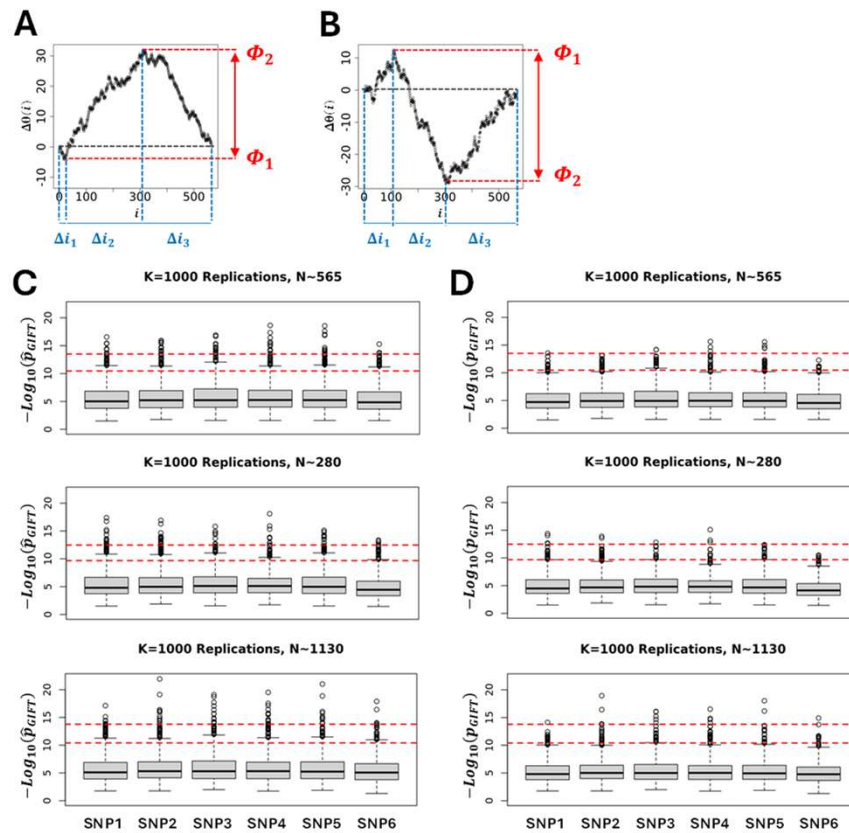




Fig.7

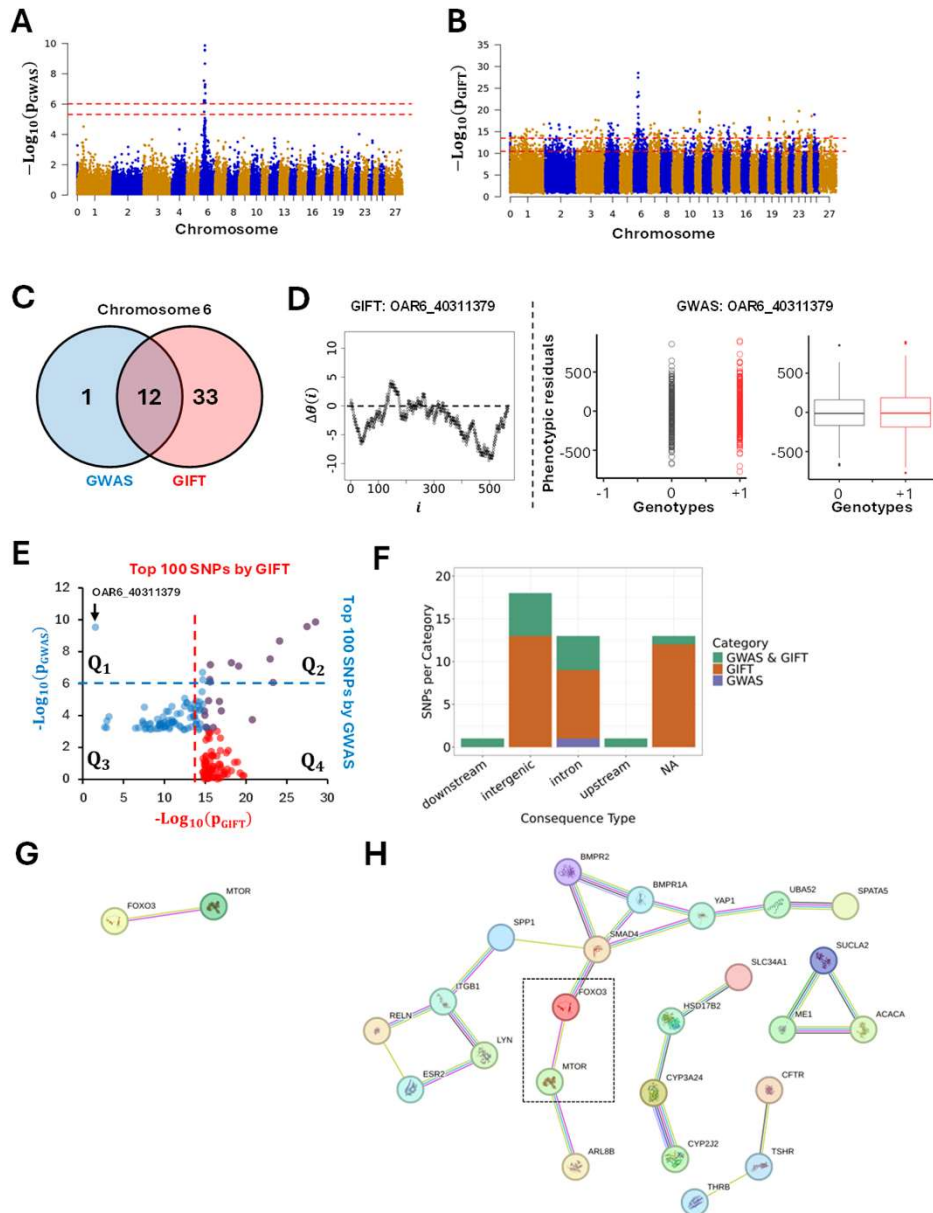


Fig.8

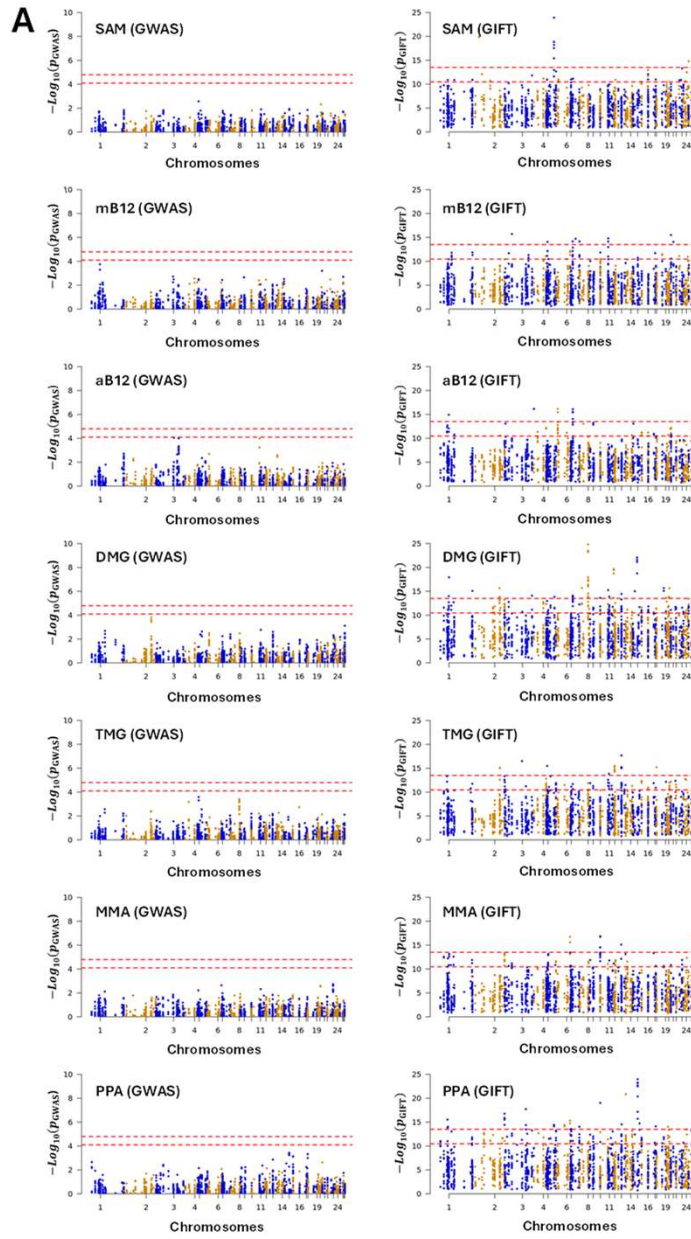


Fig.8 (continued)

