1    Classification: Biological Sciences; Evolution
2
3
4    Title: Detecting selection with a genetic cross
5
6
7    Hunter B. Fraser
8
9
10   Department of Biology, Stanford University, Stanford, CA 94305, USA. hbfraser@stanford.edu
11
12
13   Keywords: natural selection, genetic cross, variance

**Abstract**

Distinguishing which traits have evolved under natural selection, as opposed to neutral evolution, is a major goal of evolutionary biology. Several tests have been proposed to accomplish this, but these either rely on false assumptions or suffer from low power. Here, I introduce a new approach to detecting lineage-specific selection that makes minimal assumptions and only requires phenotypic data from ~10 individuals. The test compares the phenotypic difference between two populations to what would be expected by chance under neutral evolution, which can be estimated from the phenotypic distribution of an $F_2$ cross between those populations. Simulations show that the test is robust to parameters such as the number of loci affecting the trait, the distribution of locus effect sizes, heritability, dominance, and epistasis. Comparing its performance to the QTL sign test—an existing test of selection that requires both genotype and phenotype data—the new test achieves comparable power with 50- to 100-fold fewer individuals (and no genotype data). Applying the test to empirical data spanning over a century shows strong directional selection in many crops, as well as on naturally selected traits such as head shape in Hawaiian *Drosophila* and skin color in humans. Applied to gene expression data, the test reveals that the strength of stabilizing selection acting on mRNA levels in a species is strongly associated with that species' effective population size. In sum, this test is applicable to phenotypic data from almost any genetic cross, allowing selection to be detected more easily and powerfully than previously possible.

**Significance Statement**

Natural selection is the force that underlies the spectacular adaptations of all organisms to their environments. However, not all traits are under selection; a key question is which traits have been shaped by selection, as opposed to the random drift of neutral traits. Here, I develop a test of selection on quantitative traits that can be applied to almost any genetic cross between divergent populations or species. The test is robust to a wide range of potential confounders, and has greater power to detect selection than existing tests. Applied to empirical data, the test reveals widespread selection in both domesticated and wild species, allowing selection to be detected more easily and powerfully than previously possible.

**Introduction**

Trait-based tests of selection aim to distinguish the effects of two major forces of evolution: natural selection and neutral drift. Because many factors affect trait divergence—e.g. population size, divergence time, and genetic architecture—distinguishing these two forces is seldom straightforward. Several types of trait-based selection tests have been proposed, all of which view neutrality as a null model, but which differ in how they assess this null and in the type of data they require (reviewed in Chapter 12 of Walsh and Lynch (1)).

For example, time series tests use phenotypic measurements of a single species over time, typically from the fossil record (a stratophenetic series). If the trait shows departure from the neutral expectation of a random walk—e.g. many more time steps with trait increases than decreases—then neutrality is rejected. The key assumption is that environmental changes do not affect these phenotypic trends, which is difficult to justify considering how much environments can change over the millions of years typically covered in a stratophenetic series.

2

60    A more widely used approach is known as $Q_{ST}$, where the population structure of
61  phenotypic variance is compared to the analogous genetic metric $F_{ST}$. By utilizing genetic
62  crosses in common garden experiments, the confounding effects of environment can be
63  controlled, allowing selection to be assessed in a wide range of species (2). Limitations of this
64  approach include low power (requiring data from >10 populations (3)) and several assumptions
65  about epistasis and mutation rates (see Supplemental Note). However, an improved $Q_{ST}$-based
66  method has sufficient power to detect selection using only a few populations (4, 5).
67    Another widely used test is known as the quantitative trait locus (QTL) sign test (6, 7). In
68  this test, QTL are first mapped using genotype and phenotype data from a genetic cross between
69  two divergent parental lines. Under neutrality (and the absence of ascertainment bias), QTL
70  directionality—i.e. which parent's allele increases the trait at each QTL—is expected to be
71  binomially distributed around 50%, much like a series of coin flips (Fig. 1a left). In contrast,
72  under lineage-specific selection, QTL directions will be biased in one direction (Fig. 1a right).
73  Although this test is quite robust due to its minimal assumptions, it also suffers from low power:
74  a minimum of eight QTL (which is rarely reached in practice; see Supplemental Note) is
75  required to achieve a nominal $p < 0.01$.
76    The sign test's low power is largely due to the fact that it only uses QTL directionality
77  information, while ignoring the phenotypic divergence between the two parental lines. However,
78  the parental traits contain important information: if a trait evolves under directional selection, it
79  will diverge much faster than under neutrality (Fig. 1b). If it were possible to estimate the
80  divergence expected by chance under neutrality, then this could be used as a null hypothesis;
81  parental trait divergence significantly greater than this expectation would suggest lineage-
82  specific directional selection, whereas divergence less than this would suggest stabilizing
83  selection.
84    Indeed, this intuitive logic underlies another class of trait-based methods, "rate tests," that
85  ask whether the phenotypic divergence of multiple populations is consistent with neutral drift (1,
86  8). The neutral expectation is estimated from population genetic theory, using parameters such as
87  the effective population size, the mutational variance, and the number of generations since
88  population divergence. Since these parameters and their sampling variances can typically only be
89  roughly estimated (at best), and several strong assumptions must also be made, rate tests are
90  viewed as qualitative guides rather than quantitative tests of neutrality (1, 8) (Supplemental
91  Note).
92    In this work, I sought to develop a trait-based test of selection with the robustness of the
93  sign test, while utilizing the framework of rate tests to increase the power to detect selection.
94
95  **Results**
96    The logic underlying rate tests could lead to a more rigorous test of selection if the
97  expected divergence under neutrality could be more accurately estimated. This can be achieved
98  using the neutral model of the sign test, where the genetic variants underlying QTL (quantitative
99  trait nucleotides, or QTN) have no directionality bias (Fig 1a left and Fig 1b purple). How could
100  the distribution of phenotypes expected under this model of neutrality be estimated? One way
101  would be to measure the effect size of every QTN and then predict the phenotypes resulting from
102  random allelic combinations. However, there is a simpler solution: the $F_2$ trait distribution
103  represents exactly this null model. Regardless of the QTN directionalities in the parents, the $F_2$
104  phenotypes result from random combinations of the segregating alleles; this randomness mimics
105  the random directionality expected under neutral evolution (Fig 1a left, Fig 1b purple). Therefore

106    every $F_2$ individual can be thought of as a random draw from the distribution of potential
107    parental phenotypes resulting from neutral evolution. If the parental divergence is significantly
108    greater than expected based on the phenotypic variance observed in the $F_2$ population, then
109    neutrality can be rejected in favor of directional selection (Fig 1c top). If instead the parents are
110    significantly less diverged than expected—known as transgressive segregation—then stabilizing
111    selection is inferred (Fig 1c bottom).
112         I will begin with a simple model of a trait in a haploid species where all QTN are additive
113    with equal effect sizes, and there is no environmental variation or trait measurement error (i.e.
114    broad-sense heritability $H^2 = 1$). Let $n_q$ denote the number of QTN for this trait that differ
115    between two strains/populations. Under neutrality, traits diverge from the ancestor like a random
116    walk, proportional to $\sqrt{n_q}$ (Fig 1b purple). For two lineages evolving independently, the
117    expected absolute difference will be $\sqrt{2n_q}$. The difference between the two parental trait values
118    represents a single draw from a binomial distribution with $p = 0.5$ and $n = n_q$ (it is only one draw
119    since once one parent's allelic states are defined, the other's must be the opposite for any QTN
120    segregating in the cross); for $n_q > {\sim}20$ this approximates a normal distribution. The square of this
121    difference is proportional to the parental trait variance; dividing this variance by the variance
122    expected by chance under neutral evolution—i.e. the variance of the $F_2$ trait distribution, as
123    discussed above—results in the test statistic (denoted $v$ and illustrated in Fig 1c):
124

$$v = \frac{\sigma^2_{par}}{\sigma^2_{F2}} \tag{1}$$

125    where $\sigma^2_{F2}$ is the $F_2$ phenotypic variance and $\sigma^2_{par}$ is the between-strain variance of the two
126    parental strain (or population) means. This ratio is expected to be distributed as $F(1, n_{F2}\text{-}1)$ under
127    neutrality, where $n_{F2}$ is the number of $F_2$ individuals. This approximates a $\chi^2$ distribution with 1
128    degree of freedom for $n_{F2} > {\sim}20$.
129         We can now relax the simplifying assumptions above. In diploids, the expected variance
130    in the $F_2$ is half that between the parents; this is accounted for by multiplying the denominator by
131    a constant, denoted $c$ (more generally, any factors that affect the phenotypic variance in the
132    progeny—including dominance and other cross designs such as backcrosses or recombinant
133    inbred lines [RILs]—can be accommodated by adjusting the value of this constant; see
134    Supplemental Note). Allowing environmental variation and trait measurement error is equivalent
135    to adding random noise to both the numerator and denominator. Correcting for this yields:
136

$$v = \frac{\sigma^2_{par} - \dfrac{\sigma^2_{p1}}{n_{p1}} - \dfrac{\sigma^2_{p2}}{n_{p2}}}{\sigma^2_{F2} H^2 c} \tag{2}$$

137    where $n_{p1}$ and $n_{p2}$ are the number of replicate individuals measured for each parental strain, and
138    $\sigma^2_{p1}$ and $\sigma^2_{p2}$ are the within-strain variances of each parental strain (see Supplemental Note). All
139    of these terms can be estimated from phenotype data in a single genetic cross, provided that
140    multiple individuals of each parental type are included. Note that Equation 1 is a special case of
141    Equation 2, where $c = 1$ (for haploids) and $H^2 = 1$ (hence $\sigma^2_{p1} = \sigma^2_{p2} = 0$).
142         To explore the behavior of $v$ as a neutral null model, I conducted simulations of neutral
143    traits in parental strains and their $F_2$ progeny (see Methods). These simulations allowed the

4

144  precise manipulation of individual parameters to assess their effects on $v$. These are presented as
145  quantile-quantile (QQ) plots, where points following the Y=X line represent adherence to the
146  expected F distribution. Points above and below the line represent values of $v$ greater and less
147  than expected under the null, respectively.
148      The test followed the expected distribution of $v$ closely for a range of different QTL
149  effect size distributions (Fig 2a). This includes the exponential distribution, which is thought to
150  be a reasonable approximation for QTL (9, 10). However, extremely skewed distributions—e.g.
151  a monogenic trait where one QTL explains all trait variance—will lead to half of all $F_2$
152  individuals identical to one of the parents and $v$ values tightly distributed around 1, resulting in
153  little power to detect any deviation from the null. Therefore the test is only useful for polygenic
154  traits.
155      How polygenic must a trait be? More important than the number of genetic variants
156  affecting the trait—which is likely to be quite large for complex traits (11)—is the number of
157  independently segregating QTN, which is limited by the recombination index (12) (RI). The RI
158  is defined as the haploid number of chromosomes plus the number of recombinations per
159  meiosis; it represents the number of independent genomic regions segregating in a cross.
160  Adherence of $v$ to the null improves with greater RI (Fig 2b) since more shuffling of the two
161  parental genomes leads to more normally distributed $\sigma^2_{F2}$ values. For example, an extreme case
162  of a single non-recombining chromosome (RI = 1) would be equivalent to the monogenic
163  example above where $v$ cannot deviate far from 1. The $v$-test behaves conservatively in crosses
164  with RI < ~20 (Fig 2b). Fortunately, this is rarely an issue in practice since the mean RI for
165  plants is ~30 and for animals is ~40 (13). (For RILs, there are more generations for
166  recombination so the mean "effective RI" is ~45 for plants and ~58 for animals; see Methods.)
167      Sample size is another important consideration. Although more samples are always
168  preferable, $v$ follows the null distribution even with only three $F_2$ individuals (Fig 2c).
169      Other potential sources of noise include environmental variability and measurement
170  error. Although these affect both parental and $F_2$ phenotypic variances, they have less effect on
171  the parental estimates when parental replicates are included, because taking the mean of each
172  parental strain reduces this noise. Without correcting for this effect (i.e. setting $H^2 = 1$, and thus
173  $\sigma^2_{p1} = \sigma^2_{p2} = 0$, in Equation 2), low $H^2$ leads to severe underestimates of $v$ (Fig 2d left panels).
174  However including a correction for this (Equation 2) precisely accounts for this effect (Fig 2d
175  right panels).
176      The final effect explored via simulation is genetic interaction, which is the context-
177  dependence of phenotypic effects either between loci (epistasis) or between alleles at the same
178  locus (dominance). Epistasis can take many forms, but the large-scale pattern most often
179  observed is known as diminishing returns epistasis (14–16), where strains with the greatest trait
180  value have lower values than expected from additivity. To model this, I transformed the
181  simulated trait values as $\sqrt{t}$, where $t$ is each trait value. Since this affects both $\sigma^2_{F2}$ and $\sigma^2_{par}$
182  similarly, it has little impact on the distribution of $v$ (Fig 2e). I also modeled synergistic
183  (increasing returns) epistasis as $t^2$, and again found no effect (Fig 2e). However, epistasis in more
184  extreme forms could obscure any signal of selection (Supplemental Note). Dominance can be
185  accounted for by adjusting the value of $c$ to offset its effect on $\sigma^2_{F2}$ (Fig 2e, Supplemental Note).
186      In sum, simulations of neutral evolution show that $v$ is robust with respect to the number
187  of loci affecting the trait, the distribution of locus effect sizes, environmental variability,
188  measurement error, dominance, and epistasis.

189     Having established the behavior of $v$ under neutrality, I explored its power by simulating
190 directional selection. These were identical to the neutral simulations except that all QTN had
191 concordant parental directionality (Fig 1a right, Fig 1b blue). I compared the $v$-test to the sign
192 test to assess their relative performance in a range of parameter settings (other selection tests use
193 very different types of data, precluding a direct comparison). Both tests utilized the same
194 simulated $F_2$ phenotype data for each cross; the sign test was also provided with optimal
195 genotype data, meaning that each QTN was genotyped without error, had no linkage with other
196 QTN, and no other genetic markers were included. Both tests result in a probability of neutrality
197 ($p_{nut}$) value for each trait (see Methods). More extreme $p_{nut}$ values represent greater power to
198 detect selection, with points above the diagonal representing crosses with greater power for the
199 $v$-test and those below the diagonal indicating greater power for the sign test.
200     Even with optimal genotype data, at small sample sizes ($n_{F2}$ = 10 or 100) very few QTL
201 are mapped, resulting in the sign test's low power (Fig 3). In contrast, the $v$-test often rejects the
202 null with $p_{nut} < 10^{-4}$ even at $n_{F2}$ = 10. The $v$-test is generally more powerful than the sign test at
203 $n_{F2} < 10^3$. However, the sign test generally outperforms the $v$-test at very large sample sizes ($n_{F2}$
204 $> 10^3$ when $H^2 > 0.8$, and $n_{F2} > 10^4$ for all $H^2$).
205     To further explore the $v$-test's power at small sample sizes, I simulated crosses with a
206 total of 10, 20, or 30 phenotyped individuals (including the parents). For example at n = 10, 93%
207 of traits rejected the null at $p_{nut} < 0.05$ when $H^2 = 0.1$, and 99% when $H^2 = 0.8$ (Supp Fig 1a). The
208 $v$-test performed well with 10 individuals even when selection was weak (with up to 20% of
209 QTN acting in opposition to the majority) and heritability was low (Supp Fig 2). In contrast, the
210 sign test required 500-1000 phenotyped and genotyped individuals to reach the same power
211 (Supp Fig 1b). This difference in power of the two tests makes sense considering that although
212 they are both evaluating the same neutral null model, the sign test is doing so directly (by
213 comparing QTL directionalities to the neutral expectation; Fig 1a), whereas the $v$-test is doing so
214 indirectly. When enough QTL are mapped, the direct approach of the sign test is superior, but by
215 not needing to map QTL, the $v$-test requires fewer individuals, as well as no genotype data.
216     Due to its generality, the $v$-test can be applied to data from almost any genetic cross
217 where both parental strains and their $F_2$ (or other cross) progeny were phenotyped. This includes
218 most QTL studies, as well as genetic studies from before genotyping was possible. To explore
219 the test empirically, I collected published data for 126 traits from 21 species (Supp Table 1). The
220 $v$-test $p_{nut}$-values for artificially selected traits (in crops, livestock, and laboratory experiments)
221 revealed a strong skew towards low $p_{nut}$-values indicating directional selection, as expected[17]
222 (Fig 4a; for a discussion of trait ascertainment bias—a major caveat for all trait-based selection
223 tests—see Discussion and Supplemental Note). Three of the most significant traits were from
224 maize, including data for ear length and seed weight published in 1913, suggesting intense
225 artificial selection on these traits prior to this date [18].
226     In contrast, the $p_{nut}$-value distribution for traits naturally selected in the wild showed a
227 less extreme skew (Fig 4b; Supp Table 1; comparing distributions, Wilcoxon $p = 2 \times 10^{-5}$), with
228 peaks at both low and high p-values suggesting a wide range of selection pressures. The most
229 significant trait for directional selection was male head shape in a cross between two species of
230 Hawaiian *Drosophila* [19] (data from reciprocal $F_2$ crosses and an $F_6$ cross are all significant;
231 Supplemental Note). This is a well-known example of rapid morphological evolution, potentially
232 due to sexual selection [20], but whether the divergence could instead be explained by genetic
233 drift was not previously testable. The next most significant trait was human skin color, measured
234 in the "$F_2$" grandchildren of West Africans and Europeans [21]. Despite the small sample size

235    ($n_{F2} = 14$), the $v$-test is significant at all three reflectance wavelengths tested (Supp Table 1).
236    Human skin color has long been thought to be adaptive based on its correlation with local UV
237    radiation (22); these results provide independent confirmation for the role of selection. The third
238    most significant wild trait was burrowing behavior in *Peromyscus* mice (23), as measured by
239    burrow length in an interspecies backcross.
240            The $v$-test can also be applied to molecular-level traits such as gene expression levels,
241    which avoids the effects of trait ascertainment bias (Supp Note). The BXD collection of mouse
242    RILs is an excellent test case, with gene expression data available for 16 tissues. Performing the
243    $v$-test revealed that the $p_{nut}$-value distributions of all tissues were shifted towards stabilizing
244    selection to varying degrees (Fig 4c). To estimate the strength of stabilizing selection in each
245    tissue, I calculated the fraction of genes with $v$-test $p_{nut} > 0.99$. All 16 tissues had values between
246    1.0%-1.8% (note this should not be interpreted as the fraction of genes under stabilizing
247    selection, which is likely to be far higher). Interestingly, gene expression from six different
248    regions of the brain had significantly stronger stabilizing selection than in the ten non-brain
249    tissues (Fig 4c). This is consistent with previous reports of slower evolution of gene expression
250    in the mammalian brain compared to other tissues(24, 25), and suggests that this slower
251    evolution is at least partially due to greater selective constraint (as opposed to a lower mutational
252    variance, which can also lead to a slower evolutionary rate (1)).
253            Another type of molecular trait measured in the BXD cross is metabolite levels in liver
254    (26). Applying the $v$-test to these metabolomic data, cohorts fed two diets (high-fat and normal
255    chow) showed $p_{nut}$-values strongly skewed towards one (Fig 4c). Therefore the measured
256    metabolites appear to be under stronger stabilizing selection than mRNA levels.
257            Despite the variation in stabilizing selection on gene expression across the 16 tissues, all
258    tissues were relatively close to the neutral expectation (1% of genes at $p_{nut} > 0.99$). To compare
259    this to other species, I collected gene expression data from genetic crosses of five additional
260    species (*Saccharomyces cerevisiae*, *Oryza sativa*, *Arabidopsis thaliana, Brassica rapa*, and
261    *Caenorhabditis elegans*; Supp File 1). In contrast to mouse, some species had much stronger
262    stabilizing selection (e.g. *B. rapa* with 10.7% of genes under stabilizing selection). One
263    hypothesis to explain this wide range of values is that natural selection is expected to be stronger
264    in species with larger effective population sizes ($N_e$). Direct measurements of $N_e$ for these
265    species are not possible, but an indirect indicator is the fraction of neutral genomic positions that
266    are heterozygous (known as $\pi$), which is expected to increase with $N_e$ (27). Plotting the strength
267    of stabilizing selection against published values of $\pi$, I observed a strong correlation (Fig 4d).
268    This suggests that $N_e$ (or another factor correlated with $N_e$; see Supp Note) may be a major
269    determinant of stabilizing selection on gene expression levels, as has been previously proposed
270    (24, 25).
271            The peaks of gene expression $p_{nut}$-values near zero (Fig 4d) suggest that directional
272    selection may also be detectable from these data. For example, in *S. cerevisiae* genes with low
273    $p_{nut}$ are highly enriched for roles in mitochondrial translation (FDR = $6 \times 10^{-23}$ among genes
274    down-regulated in the lab strain BY; no enrichment among genes up-regulated in BY). In *B.*
275    *rapa*, the defense response to other organisms is the most enriched function among genes with
276    low $p_{nut}$ (FDR = 0.02). Similarly in *C. elegans*, immune response was the most enriched (FDR =
277    0.04).
278
279    **Discussion**
280

7

281     In this work, I introduced a new test of selection that combines the logic of rate tests with
282     the neutral null model of the sign test. The result is a test that is simple, robust, and more
283     powerful than existing tests. This will allow researchers to assess selection on traits of interest
284     any time they perform a genetic cross. Moreover, if multiple traits are measured in the same
285     cross, results from the test can be directly compared to assess the strength of selection on diverse
286     phenotypes (as in Fig 4c).
287     There are many potential extensions to this test. For example, the *v*-test framework could
288     be applied any time the genomes of two divergent populations are mixed, including naturally
289     admixed populations. In this case, the only modification needed would be in the calculation of c,
290     representing the expected ratio of parental variance to admixed progeny variance (Supplemental
291     Note). Other extensions could include testing multiple correlated traits simultaneously, focusing
292     only on additive effects, or estimating confidence intervals (Supplemental Note).
293     It is important to note that results from all trait-based tests of selection must be treated
294     with caution when trait ascertainment bias is present. If traits are chosen for study based on the
295     extent of their divergence between populations, then the neutral model no longer holds. For
296     example, imagine 100 neutral traits; in any properly calibrated selection test, we would expect ~5
297     of these to reach a nominal $p < 0.05$. If these same five are the only traits included in a study
298     (e.g. because they have the strongest phenotypic divergence), then they will appear to be
299     inconsistent with whichever null model they are tested against. In some cases it is possible to
300     correct for ascertainment bias, either by modifying the test itself (6) or by using a more
301     conservative p-value threshold (see Supplemental Note). However, the ideal solution is to
302     analyze traits that were selected for study independently of the parental trait values, which by
303     definition lack any ascertainment bias. The most widespread examples of this are molecular-
304     level traits such as the levels of mRNAs, proteins, metabolites, etc. Similarly, standardized
305     phenotyping (28) can be free of bias.
306     Notably, Equation 2 is identical to the widely used Castle-Wright (CW) estimator for the
307     number of loci underlying divergence in a quantitative trait (29–31) (though it was derived
308     independently). The maximum possible value of this estimator is the RI of the species being
309     crossed, resulting in a strong downward bias for most complex traits, which can have thousands
310     of variants contributing (11). It is therefore rather fortuitous that this severely biased estimator is
311     also precisely F-distributed under the null hypothesis of neutrality, even though neutrality had no
312     role in its original derivation and the F distribution has no role in its traditional interpretation (29,
313     30). Furthermore, the true number of loci underlying a trait (what the CW estimator aims to
314     estimate) is not indicative of selection; a neutral trait could have any number of underlying loci,
315     so this cannot be used to assess a trait's neutrality.
316     How can this one equation have two seemingly unrelated interpretations? The CW
317     estimator requires a number of restrictive assumptions, including that all QTL must act in the
318     same direction with respect to their parent of origin (as in Fig 1a right panel). Rather than being
319     an assumption of the *v*-test, this concordant QTL directionality is exactly what the *v*-test was
320     designed to detect. Therefore the connection between the two interpretations rests on the fact that
321     the strength of the signal detected by the CW estimator—the number of reinforcing QTL acting
322     in the same parental direction—is indicative not only of the number of loci, but also of any
323     selection that has acted on those loci since the divergence of the two parental strains.
324     One consequence of this mathematical homology is that the hundreds of published CW
325     estimator values dating back to 1921 (29) can now be immediately reinterpreted as tests of

326    neutral evolution (Supplemental Note), even when the phenotype data from these studies are not
327    available.
328
329
330
331

## Methods

### Neutral simulations

Parameters in neutral simulations are shown in Fig 2: QTN effect size distribution, RI (the number of independently segregating QTN), $n_{F2}$, $H^2$, epistasis, and dominance (c = 2 for additive, c = 1 for bidirectional dominant, and c = 4/3 for unidirectional dominant; see Supplemental Note). Parameter values were chosen to reflect typical published data sets rather than optimal parameters for test performance. First I generated effect sizes for the specified number of QTN by sampling from the specified distribution. In neutral evolution, QTN directions in each parent are random (Fig 1a), so the first parent's traits were determined by flipping the sign of each effect size with 50% chance and then summing the values. The second parent's traits were calculated the same way but with all signs flipped from the first parent (each QTN increases the parental difference by twice its effect size, assuming parents are homozygous). $F_2$ traits were determined by multiplying each effect size by one number chosen randomly from a set of four numbers that represent the four possible $F_2$ genotypes at each locus and their resulting phenotypic effects ([0 1 1 2] for additive, [0 2 2 2] for unidirectional dominant, or [0 0 2 2] for bidirectional dominant; see Supplemental Note), separately for every individual, and summing across all QTN. When $H^2 < 1$, random noise was also added to each parental and $F_2$ individual as a normally distributed variable with zero mean and standard deviation $\sigma = \sqrt{\frac{\sigma_S^2 * (1 - H^2)}{H^2}}$, where $\sigma_S^2 = \sigma_{F2}^2 - \sigma_{env}^2$ (which in practice is calculated as the variance of the $F_2$ trait values before noise is added, since $\sigma_{F2}^2$ and $\sigma_{env}^2$ are not known until the noise is calculated). In epistasis simulations, traits were additionally transformed either by $\sqrt{t}$ or $t^2$, where $t$ is the trait value. Parental within-strain variances ($\sigma_{p1}^2$ and $\sigma_{p2}^2$) were then estimated from 10 replicates per parent, and $\sigma_{F2}^2$ was estimated from the $F_2$ population. From these variables, $v$ was calculated using Equation 2 and converted into a $p$-value based on the F(1, $n_{F2}$-1) cumulative distribution.

### Selection simulations

Selection was simulated using the framework described above, with one difference: omitting the step of flipping the sign of each effect size with 50% chance. This meant that all QTN were reinforcing in their directionalities. $v$ was then calculated as described above and converted to a p-value based on the F(1, $n_{F2}$-1) cumulative distribution. Parameter values are listed in Fig 3 legend.

To calculate the sign test p-value, QTL must first be mapped. The genotype of each QTN variant (randomly generated as described above) was provided in a genotype matrix, with no genotyping error and no additional genetic markers (this represents an unrealistic best-case scenario for QTL mapping). Pearson's correlation coefficient between each marker and the $F_2$ phenotypes were then converted into LOD scores (32) as $LOD = \frac{-n * ln(1 - r^2)}{2 ln(10)}$. LOD > 3 was required to call a QTL. The directionalities for the full set of QTL for each cross were then assessed for their fit to the binomial distribution cumulative distribution with expected frequency = ½. The resulting two-sided p-value was the sign test p-value.

The simulations in Supp Figs 1-2 were identical to those in Fig 3, but with different parameter values and different visualizations. 1000 simulations were performed for each combination of parameter values shown. In Supp Fig 2, true positives were defined as crosses

10

376 simulating selection where the *v*-test p-value was below a given cutoff; false positives were
377 crosses simulating neutral evolution where the p-value was below the same cutoff (the cutoff
378 varied from 0 to 1 to generate each ROC curve). Selection strength was represented by the
379 fraction of QTN with reinforcing directionality. This cannot be translated into a selection
380 coefficient since it would depend on how the trait relates to fitness, the population size, time
381 since population divergence, and many other parameters.
382
383 **Empirical analysis**
384     Traits were collected for Fig 4a from two sources: Wright (30) Tables 15.1-15.8, and
385 Lynch and Walsh (32) Table 9.6. Traits were collected for Fig 4b from Lynch and Walsh (32)
386 Table 9.6, Rieseberg (17) Table 5, and literature searches for cichlid and *Peromyscus* data.
387     For traits in Fig 4a-b and Supp Table 1, $H^2$ values were estimated as follows. If values
388 were provided by the authors of the original study, these were used. If not, the environmental
389 variance was estimated using Wright's preferred method (30), a weighted average of within-
390 strain variances: $\sigma_{env}^2 = \sigma_{p1}^2/4 + \sigma_{p2}^2/4 + \sigma_{F1}^2/2$. If data from the $F_1$ were not available, then I
391 used the sample size of each parental type as weights: $\sigma_{env}^2 = (n_{p1}\sigma_{p1}^2 + n_{p2}\sigma_{p2}^2)/(n_{p1} + n_{p2})$.
392 This variance was then used to calculate $H^2 = (\sigma_{F2}^2 - \sigma_{env}^2)/\sigma_{F2}^2$. In some cases this can lead to a
393 negative $H^2$ (likely due to overestimation of $\sigma_{env}^2$ since this was always based on fewer replicates
394 than $\sigma_{F2}^2$); therefore values of $H^2 < 0.1$ were set to 0.1. For the two cases of traits with outbred
395 parents (burrowing and parenting behavior in *Peromyscus*), $H^2$ may be underestimated due to
396 within-strain genetic variation contributing to $\sigma_{env}^2$; therefore I conservatively set $H^2 = 0.4$ for
397 these traits (which is higher than the heritability of most behavioral traits (33)), resulting in less
398 significant values of $p_{nut}$ (see Supplemental Note).
399     Data in Fig 4c were collected from http://www.genenetwork.org/ (selecting species:
400 mouse; group: BXD family; type: any tissue with parental data). Since most mouse gene
401 expression data sets in Fig 4c had only one sample per parental strain, $H^2$, $\sigma_{p1}^2$, and $\sigma_{p2}^2$ could not
402 be accurately estimated. To allow a direct comparison between data sets with parental replicates
403 vs. those without, I made two modifications: 1) for all tissues, I assumed half of the parental
404 variance was genetic, and half environmental (i.e. the numerator of Equation 2 was set to $\sigma_{par}^2/$
405 2). 2) I set $H^2=0.31$ for all genes, this being the median value estimated for yeast (34)
406 (specifically the median of 1-*e*, where *e* is the residual gene expression variance not explained by
407 either additive or pairwise epistatic effects). I selected yeast because it had the largest number of
408 gene expression profiles from a genetic cross of any species, a comprehensive heritability
409 analysis performed by the original authors, and the closest $\pi$ to mouse. Note that these
410 modifications affect the Y-axis values in Fig 4c, but not the relative relationship between points;
411 any values could have been used without affecting the trend shown. The complete list of tissues
412 and stabilizing selection scores are in Supp File 1.
413     Gene expression data for the six species in Fig 4d were from the following sources: *S.*
414 *cerevisiae* (34), *A. thaliana* (35), *B. rapa* (36) (normal phosphorus condition), *C. elegans* (37,
415 38) (control condition), *O. sativa* (39), and *M. musculus* (see above). Published expression data
416 from other species' crosses were not usable (e.g. no parental data). For mouse, the median
417 stabilizing selection level across all 16 tissues was used. To avoid spuriously low or negative
418 estimates of $H^2$, for all six species any genes with $H^2 < 0.1$ were set to 0.1 (as described above).
419 As above, to facilitate comparison across data sets I assumed half of the parental variance was
420 genetic, and half environmental. All $p_{nut}$ values are listed in Supp File 1. $\pi$ values were taken
421 from Leffler et al. (40) as the median of autosomal $\pi$ estimates for each species. No values were

422    listed for *O. sativa* or *B. rapa*, so other published estimates were used (41, 42). For *O. sativa*,
423    both parents of the genetic cross (Zhengshan 97 and Minghui 63) were in population group XI,
424    so the π for this group was used. Using a more recently published π estimate for *S. cerevisiae*
425    (0.18%, which is the median π across all 14 non-mosaic populations (43)) yielded a slightly
426    stronger Pearson correlation ($r = 0.935$). Omitting *B. rapa* as an outlier also strengthened the
427    correlation ($r = 0.960$). Gene Ontology process enrichments were calculated with GOrilla (44).
428

429    **Estimating recombination index**
430        RI was estimated as [mean cM/100] + [mean haploid chromosome number] for 189
431    plants and 140 animals (13). RILs experience around twice as many detectable recombinations
432    (defined as those occurring in a heterozygous genomic region) as an $F_2$ cross, regardless of the
433    exact number of generations of inbreeding, so for RILs the recombination values were doubled.
434    Backcrosses experience about half as many detectable recombinations as an $F_2$ cross.

## References

1. B. Walsh, M. Lynch, *Evolution and Selection of Quantitative Traits* (2018).

2. T. Leinonen, R. B. O'Hara, J. M. Cano, J. Merilä, Comparative studies of quantitative trait and neutral marker divergence: A meta-analysis. *J. Evol. Biol.* **21**, 1–17 (2008).

3. R. B. O'Hara, J. Merilä, Bias and precision in QST estimates: Problems and some solutions. *Genetics* **171**, 1331–1339 (2005).

4. R. B. O'Hara, *et al.*, driftsel: An R package for detecting signals of natural selection in quantitative traits. *Mol. Ecol. Resour.* **13**, 746–754 (2013).

5. O. Ovaskainen, M. Karhunen, C. Zheng, J. M. C. Arias, J. Merilä, A new method to uncover signatures of divergent and stabilizing selection in quantitative traits. *Genetics* **189**, 621–632 (2011).

6. H. A. Orr, Testing natural selection vs. genetic drift in phenotypic evolution using quantitative trait locus data. *Genetics* **149**, 2099–104 (1998).

7. H. B. Fraser, Genome-wide approaches to the study of adaptive gene expression evolution. *BioEssays* **33**, 469–477 (2011).

8. M. Turelli, J. H. Gillespie, R. Lande, RATE TESTS FOR SELECTION ON QUANTITATIVE CHARACTERS DURING MACROEVOLUTION AND MICROEVOLUTION. *Evolution (N. Y).* **42**, 1085–1089 (1988).

9. S. P. Otto, C. D. Jones, Detecting the undetected: estimating the total number of loci underlying a quantitative trait. *Genetics* **156**, 2093–107 (2000).

10. H. A. Orr, THE POPULATION GENETICS OF ADAPTATION: THE DISTRIBUTION OF FACTORS FIXED DURING ADAPTIVE EVOLUTION. *Evolution (N. Y).* **52**, 935–949 (1998).

11. E. A. Boyle, Y. I. Li, J. K. Pritchard, An Expanded View of Complex Traits: From Polygenic to Omnigenic. *Cell* **169**, 1177–1186 (2017).

12. C. D. Darlington, The Biology of Crossing-over. *Nature* **140**, 759–761 (1937).

13. J. Stapley, P. G. D. Feulner, S. E. Johnston, A. W. Santure, C. M. Smadja, Variation in recombination frequency and distribution across eukaryotes: patterns and processes. *Philos. Trans. R. Soc. B Biol. Sci.* **372**, 20160455 (2017).

14. S. Kryazhimskiy, D. P. Rice, E. R. Jerison, M. M. Desai, Microbial evolution. Global epistasis makes adaptation predictable despite sequence-level stochasticity. *Science* **344**, 1519–1522 (2014).

15. H.-H. Chou, H.-C. Chiu, N. F. Delaney, D. Segrè, C. J. Marx, Diminishing returns epistasis among beneficial mutations decelerates adaptation. *Science* **332**, 1190–2 (2011).

16. A. I. Khan, D. M. Dinh, D. Schneider, R. E. Lenski, T. F. Cooper, Negative epistasis between beneficial mutations in an evolving bacterial population. *Science* **332**, 1193–6 (2011).

17. L. H. Rieseberg, A. Widmer, A. M. Arntz, J. M. Burke, Directional selection is the primary cause of phenotypic diversification. *Proc. Natl. Acad. Sci.* **99**, 12242–12245 (2002).

18. R. A. Emerson, E. M. East, The inheritance of quantitative characters in maize. *Res. Bull. Bull. Agric. Exp. Stn. Nebraska No.2* (1913).

19. A. R. Templeton, Analysis of Head Shape Differences Between Two Interfertile Species of Hawaiian Drosophila. *Evolution (N. Y).* **31**, 630 (1977).

20. C. R. B. Boake, Sexual selection and speciation in Hawaiian Drosophila. *Behav. Genet.*

13

481     **35**, 297–303 (2005).
482  21.  G. A. HARRISON, J. J. OWEN, STUDIES ON THE INHERITANCE OF HUMAN
483     SKIN COLOUR. *Ann. Hum. Genet.* **28**, 27–37 (1964).
484  22.  N. G. Jablonski, G. Chaplin, The colours of humanity: the evolution of pigmentation in
485     the human lineage. *Philos. Trans. R. Soc. B Biol. Sci.* **372**, 20160349 (2017).
486  23.  J. N. Weber, B. K. Peterson, H. E. Hoekstra, Discrete genetic modules are responsible for
487     complex burrow evolution in Peromyscus mice. *Nature* **493**, 402–405 (2013).
488  24.  D. Brawand, *et al.*, The evolution of gene expression levels in mammalian organs. *Nature*
489     **478**, 343–348 (2011).
490  25.  P. Khaitovich, W. Enard, M. Lachmann, S. Pääbo, Evolution of primate gene expression.
491     *Nat. Rev. Genet.* **7**, 693–702 (2006).
492  26.  E. G. Williams, *et al.*, Systems proteomics of liver mitochondria function. *Science (80-. ).*
493     **352** (2016).
494  27.  M. Kimura, *The neutral theory of molecular evolution* (1983).
495  28.  T. F. Meehan, *et al.*, Disease model discovery from 3,328 gene knockouts by The
496     International Mouse Phenotyping Consortium. *Nat. Genet.* **49**, 1231–1238 (2017).
497  29.  W. E. Castle, AN IMPROVED METHOD OF ESTIMATING THE NUMBER OF
498     GENETIC FACTORS CONCERNED IN CASES OF BLENDING INHERITANCE.
499     *Science (80-. ).* **54**, 223–223 (1921).
500  30.  S. Wright, *Evolution and the Genetics of Populations, Volume 1* (1968).
501  31.  C. C. Cockerham, Modifications in estimating the number of genes for a quantitative
502     character. *Genetics* **114**, 659–664 (1986).
503  32.  M. Lynch, B. Walsh, *Genetics and analysis of quantitative traits* (1998).
504  33.  D. G. Stirling, D. Réale, D. A. Roff, Selection, structure and the heritability of behaviour.
505     *J. Evol. Biol.* **15**, 277–289 (2002).
506  34.  F. W. Albert, J. S. Bloom, J. Siegel, L. Day, L. Kruglyak, Genetics of trans-regulatory
507     variation in gene expression. *Elife* **7**, 1–39 (2018).
508  35.  M. A. L. West, *et al.*, Global eQTL Mapping Reveals the Complex Genetic Architecture
509     of Transcript-Level Variation in Arabidopsis. *Genetics* **175**, 1441–1450 (2007).
510  36.  J. P. Hammond, *et al.*, Regulatory Hotspots Are Associated with Plant Gene Expression
511     under Varying Soil Phosphorus Supply in Brassica rapa. *Plant Physiol.* **156**, 1230–1241
512     (2011).
513  37.  M. G. Sterken, *et al.*, Dissecting the eQTL micro-architecture in *Caenorhabditis elegans*.
514     *bioRxiv*, 651885 (2019).
515  38.  B. L. Snoek, *et al.*, Contribution of trans regulatory eQTL to cryptic genetic variation in
516     C. elegans. *BMC Genomics* **18**, 500 (2017).
517  39.  J. Wang, *et al.*, An expression quantitative trait loci-guided co-expression analysis for
518     constructing regulatory network using a rice recombinant inbred line population. *J. Exp.*
519     *Bot.* **65**, 1069–1079 (2014).
520  40.  E. M. Leffler, *et al.*, Revisiting an Old Riddle: What Determines Genetic Diversity Levels
521     within Species? *PLoS Biol.* **10**, e1001388 (2012).
522  41.  S. Park, H.-J. Yu, J.-H. Mun, S.-C. Lee, Genome-wide discovery of DNA polymorphism
523     in Brassica rapa. *Mol. Genet. Genomics* **283**, 135–145 (2010).
524  42.  W. Wang, *et al.*, Genomic variation in 3,010 diverse accessions of Asian cultivated rice.
525     *Nature* **557**, 43–49 (2018).
526  43.  J. Peter, *et al.*, Genome evolution across 1,011 Saccharomyces cerevisiae isolates. *Nature*

527           **556**, 339–344 (2018).

528   44.   E. Eden, R. Navon, I. Steinfeld, D. Lipson, Z. Yakhini, GOrilla: a tool for discovery and
529         visualization of enriched GO terms in ranked gene lists. *BMC Bioinformatics* **10**, 48
530         (2009).

531

532

## Acknowledgements

**Figure Legends**

**Figure 1. The sign test and the *v*-test. a.** Illustration of the sign test applied to the trait of mouse size. Left panel: two mice from separate populations that have had no selection acting on size are expected to have approximately equal numbers of QTL (or QTN) alleles increasing size (binomially distributed with expected frequency = ½; stabilizing selection on size would result in a similar pattern, but with a smaller expected parental trait divergence). Right panel: In contrast, two populations that have experienced lineage-specific directional selection on size will show greater phenotypic divergence and a preponderance of QTL alleles increasing size in the larger strain. A significant deviation from the binomial expectation indicates rejection of the null hypothesis of neutral evolution.  **b.** Simulation of trait divergence under a simple model of three selection regimes. One exponentially distributed QTL (or QTN) is added per time step and the number and effect sizes of QTL are identical in each selection regime; the only difference is their directionality. Under directional selection all QTL increase the trait value (as in Fig. 1a right panel); under neutral evolution their directionalities are random; and under stabilizing selection, their directionalities are whatever will bring the trait closer to the optimum (e.g. if the trait is above the optimum, the next QTL will be negative). Each selection regime has 100 lineages simulated for 100 time steps. **c.** Illustration of the *v*-test. Under a simple model, the variance of a neutral trait in two populations is expected to be approximately equal to that of their $F_2$ progeny (Equation 1). Lineage-specific directional selection will result in higher parental variance, whereas stabilizing selection will lead to lower parental variance (transgressive segregation).

**Figure 2. Neutral simulations.** Each panel shows 20 quantile-quantile (QQ) plots where every point is an independent simulation of a genetic cross between two lineages where the trait in question has been evolving neutrally (i.e. QTL directions in each parent are random; Fig 1a-b). The X-axis shows expected p-value quantiles (uniform between zero and one), and the Y axis shows observed values of *v* (Equation 2) in 20 QQ plots each with $10^4$ simulations. For each panel, one parameter is varied from the baseline model (Exponential distribution of QTL effect sizes, recombination index = 50, number of parental replicates = 10, $n_{F2}$= 100, $H^2$=1, diploid, no epistasis or dominance), except for the upper right panel which is the baseline model. **a.** Effects of varying QTL effect sizes. **b.** Effects of varying recombination index. **c.** Effects of varying the number of $F_2$ individuals. **d.** Effects of varying $H^2$, with or without the correction in Equation 2. **e.** Effects of varying epistasis and dominance. Bidirectional dominance means all loci are fully dominant but with ~50% of loci being dominant towards one parent, and ~50% towards the other. Unidirectional means all loci are fully dominant in the same direction (i.e. the $F_1$ phenotype is identical to one of the parents).

**Figure 3. Directional selection simulations.** All panels show scatter plots where every point is an independent simulation of a genetic cross between two lineages where the trait in question has been evolving under directional selection (i.e. all QTL are in the same direction; Fig 1a-b). The X-axis shows sign test log p-values, and the Y axis shows *v* test log p-values. For each panel, two key parameters ($H^2$ and $n_{F2}$) are set to the values shown and a third is varied within the panel (RI, which takes on all integer values from 5 to 100). For each value of RI, 10 simulations are shown, each with an independent set of QTL effect sizes; this results in 960 simulations (data points) per panel. All other parameters are kept constant throughout the figure (Exponential distribution of QTL effect sizes, $n_{par} = 10$, diploid, no epistasis or dominance).

16

583

584   **Figure 4. Empirical analysis. a.** Results for artificially selected traits in crops, livestock, and

585   laboratory selection experiments. Inset shows the six most significant traits. **b.** Results for

586   naturally selected traits in plants and animals. Inset shows the three most significant traits.  **c.**

587   Results for gene expression (mRNA levels) and metabolite levels measured in the same mouse

588   RIL panel (BXD). Note that any selection detected between the two parental lineages could

589   involve divergence of their wild ancestors (mostly *M. musculus domesticus*) and/or artificial

590   selection during their inbreeding in the lab. T-test p-values shown for each comparison. **d.** Center

591   panel: The strength of stabilizing selection vs. heterozygosity ($\pi$) in six species (in order of

592   decreasing $\pi$: *Brassica rapa, Arabidopsis thaliana, Caenorhabditis elegans, Oryza sativa, Mus*

593   *musculus, Saccharomyces cerevisiae*). Side panels: the full distribution of $p_{nut}$ values for the

594   species with the highest (right) and lowest (left) $\pi$.

595

596   **Supplemental Figure 1. a.** Directional selection was simulated as in Figure 3. Two replicates of

597   each parental line were used together with the indicated number of $F_2$ individuals. Therefore the

598   total number of individuals n = $n_{F2}$ + 4.  **b.** The *v*-test simulations from panel (a) were compared

599   to sign test results using 50-fold (left) or 100-fold (right) more individuals (all $F_2$ since the sign

600   test does not require parental data). Negative values indicate the *v*-test had lower median p-value

601   than the sign test in that comparison. The sign test generally requires >50-fold more phenotyped

602   individuals (as well as genotypes) to reach the same median p-value as the *v*-test; for traits with

603   low $H^2$ it requires >100-fold more individuals.

604

605   **Supplemental Figure 2.** Receiver-operator characteristic curves are shown for the *v*-test with a

606   range of $H^2$ and selection strengths. In these plots, a perfect classifier would have 100% true

607   positives and 0% false positives; a random classifier would be on the diagonal X=Y line. All

608   simulations used 10 individuals (two of each parental strain and six $F_2$), exponential distribution

609   of QTN effect sizes, RI = 50, diploid, and no epistasis or dominance. Neutral evolution was

610   simulated as in Figure 2 (defined as having QTN directionality binomially distributed around

611   50%). Directional selection was simulated as in Figure 3, except that the fraction of reinforcing

612   QTN (reflecting the strength of selection) was allowed to vary. Note that in these simulations

613   QTN directionality was independent of effect size; a more realistic case is that QTN of larger

614   effect would be less likely to oppose the direction of selection, in which case the test

615   performance would increase for any given % of reinforcing QTN.
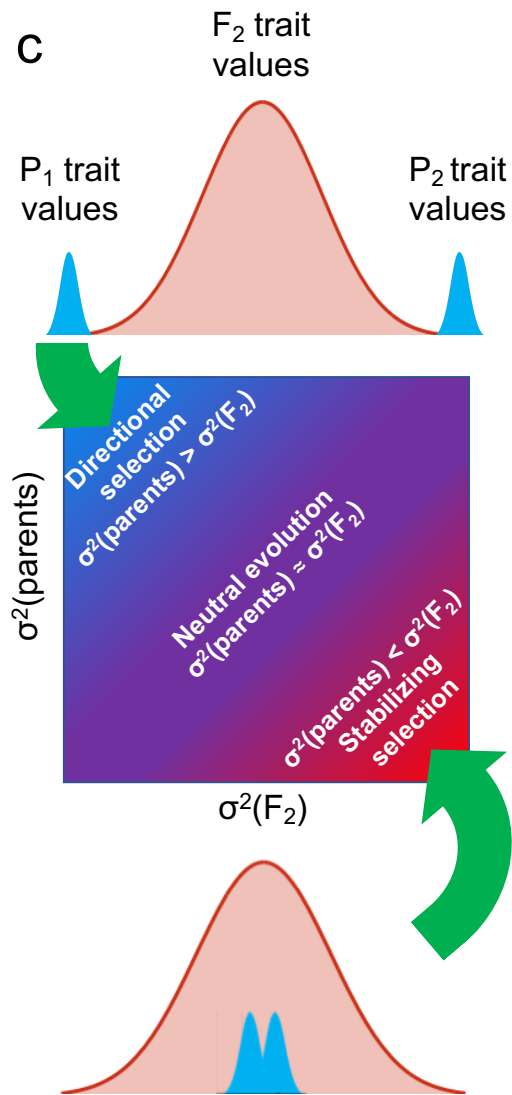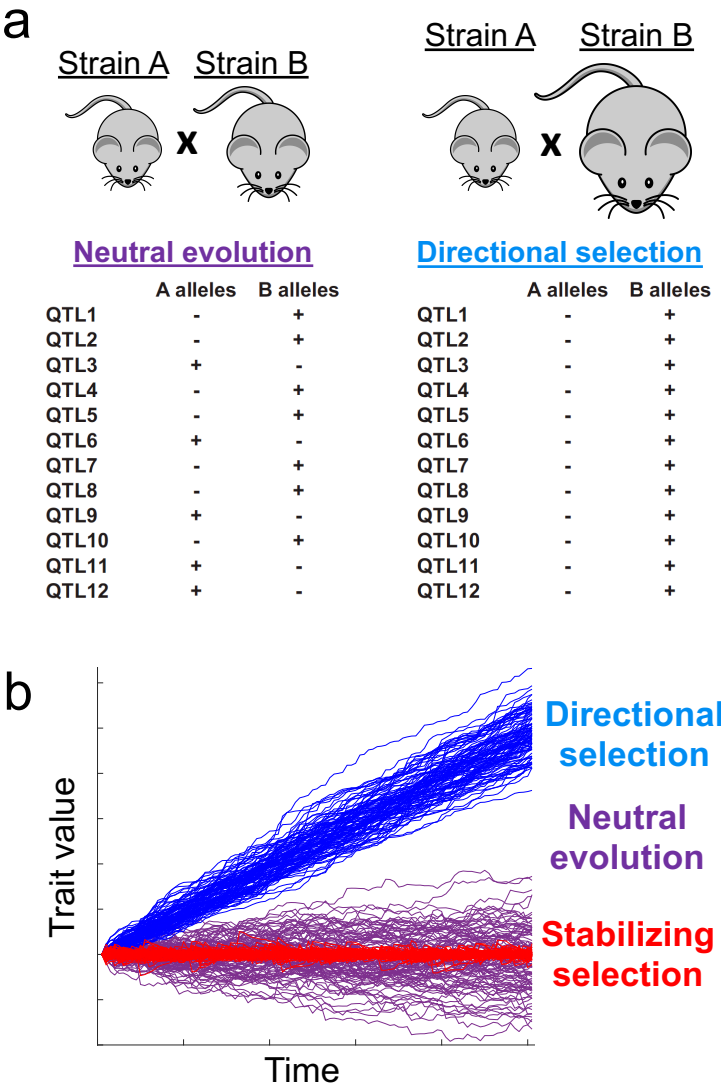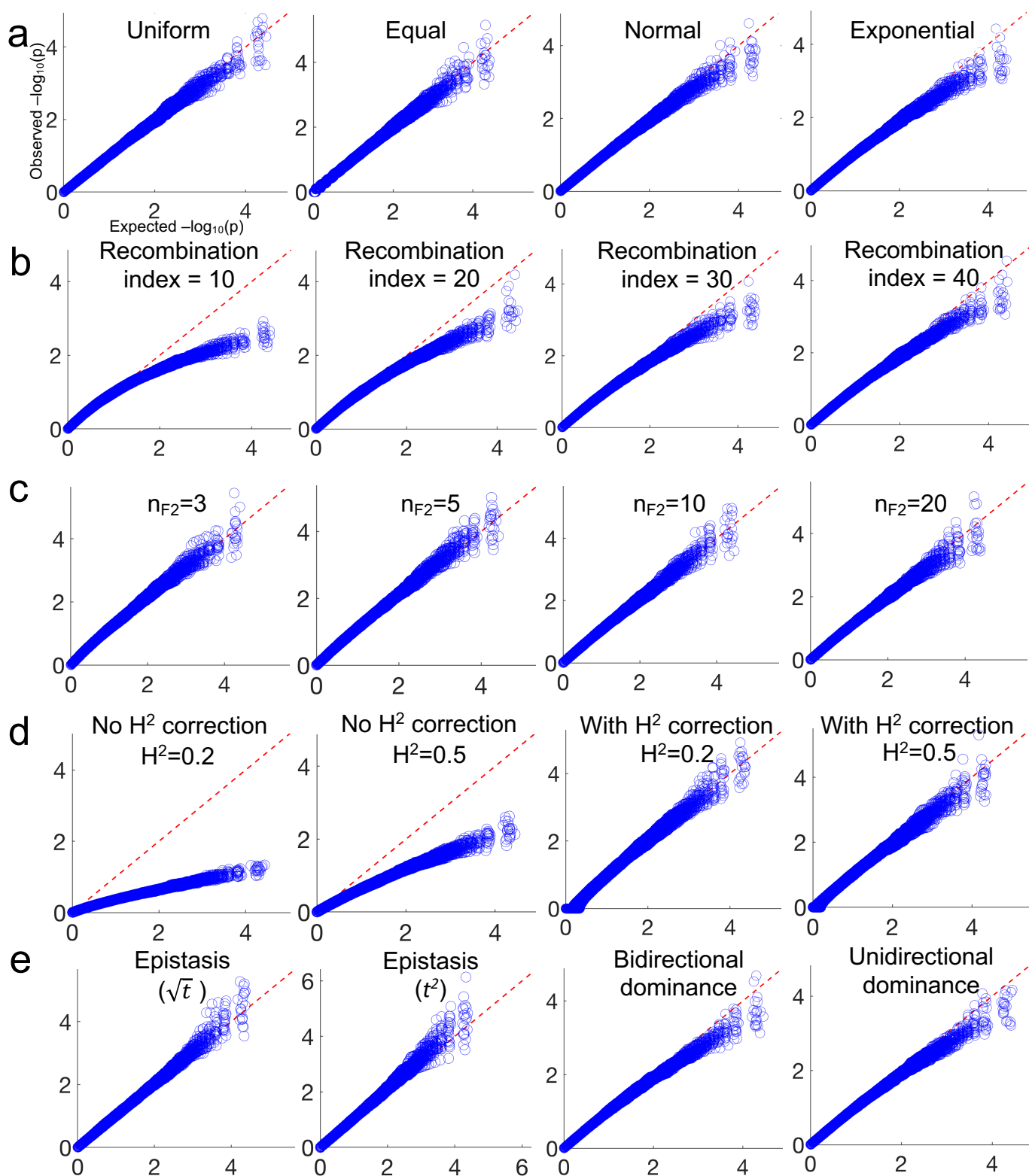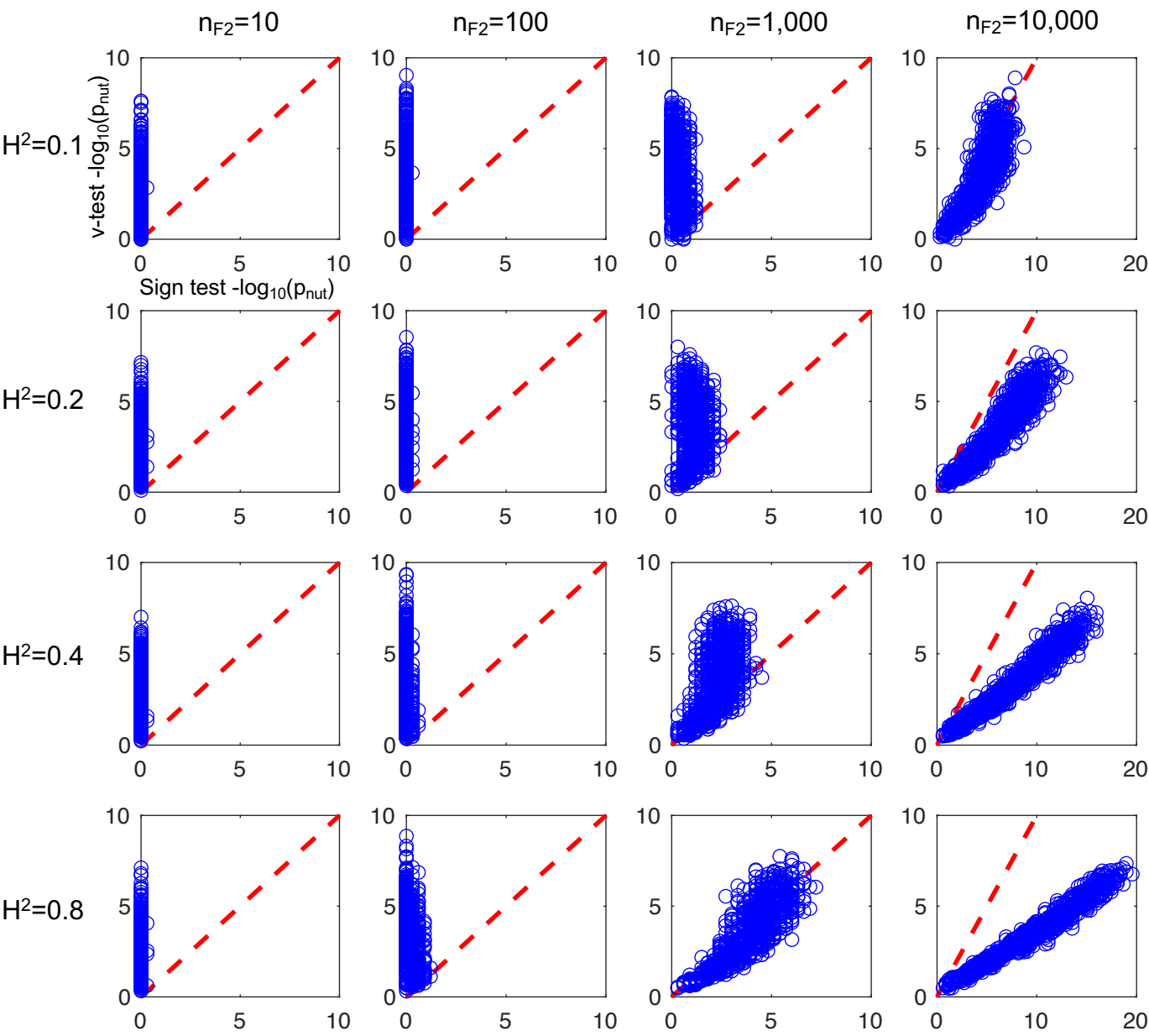
17

Fig 1



a

Strain A    Strain B        Strain A    Strain B

**Neutral evolution**          **Directional selection**

|       | A alleles | B alleles |
|-------|-----------|-----------|
| QTL1  | -         | +         |
| QTL2  | -         | +         |
| QTL3  | +         | -         |
| QTL4  | -         | +         |
| QTL5  | -         | +         |
| QTL6  | +         | -         |
| QTL7  | -         | +         |
| QTL8  | -         | +         |
| QTL9  | +         | -         |
| QTL10 | -         | +         |
| QTL11 | +         | -         |
| QTL12 | +         | -         |

|       | A alleles | B alleles |
|-------|-----------|-----------|
| QTL1  | -         | +         |
| QTL2  | -         | +         |
| QTL3  | -         | +         |
| QTL4  | -         | +         |
| QTL5  | -         | +         |
| QTL6  | -         | +         |
| QTL7  | -         | +         |
| QTL8  | -         | +         |
| QTL9  | -         | +         |
| QTL10 | -         | +         |
| QTL11 | -         | +         |
| QTL12 | -         | +         |

b

**Directional selection**

**Neutral evolution**

**Stabilizing selection**

Trait value

Time

c

$F_2$ trait values

$P_1$ trait values          $P_2$ trait values

Directional selection $\sigma^2(parents) > \sigma^2(F_2)$

Neutral evolution $\sigma^2(parents) = \sigma^2(F_2)$

$\sigma^2(parents) < \sigma^2(F_2)$ Stabilizing selection

$\sigma^2(parents)$

$\sigma^2(F_2)$

Fig 2

# Fig 3

# Fig 4



a **Artificially selected traits**

| Species/trait | p-value |
|---|---|
| Maize % oil | $2.0 \times 10^{-5}$ |
| Lab mouse weight | $2.2 \times 10^{-4}$ |
| Red pepper fruit weight | $3.3 \times 10^{-4}$ |
| Tobacco corolla length | $3.8 \times 10^{-4}$ |
| Maize seed weight | $4.5 \times 10^{-4}$ |
| Maize ear length | $5.5 \times 10^{-4}$ |

Directional selection ← → Stabilizing selection
v-test p-value

b **Naturally selected traits**

| Species/trait | p-value |
|---|---|
| Fruit fly head shape | $5.3 \times 10^{-6}$ |
| Human skin color | $1.2 \times 10^{-4}$ |
| Deer mouse burrow length | $2.1 \times 10^{-3}$ |

Directional selection ← → Stabilizing selection
v-test p-value

c **Mouse (*Mus musculus*)**

$p = 4 \times 10^{-9}$
$p = 4 \times 10^{-3}$
$p = 1 \times 10^{-4}$

Liver mRNA

Nucleus accumbens mRNA

Liver metabolites

d *Saccharomyces cerevisiae* 1.1% stabilizing selection $\pi = 0.14\%$

Pearson's $r = 0.93$, p = 0.0035
Spearman's $r = 1$, p = 0.0014

*Brassica rapa* 10.7% stabilizing selection $\pi = 0.92\%$

# Supp Fig 1

a

### Power to detect directional selection (fraction of simulations at $p_{nut} < 0.05$)

### Power to detect directional selection (median $p_{nut}$)



b

### Power to detect directional selection ($\log_2$ ratio v-test/sign test median $p_{nut}$-value)



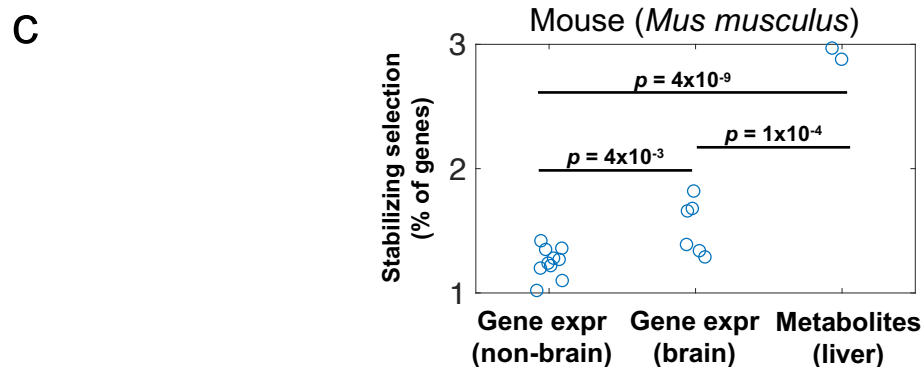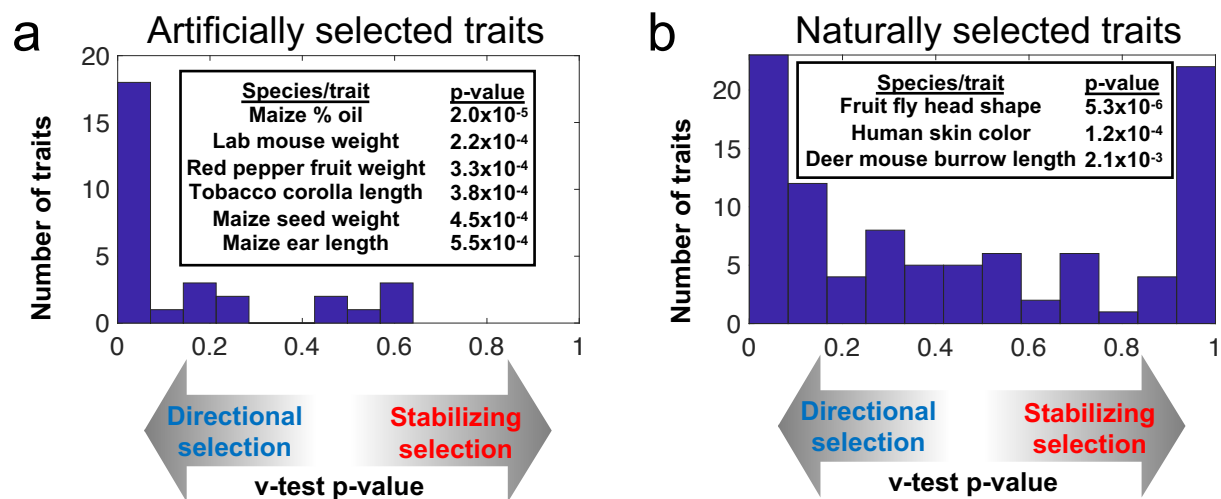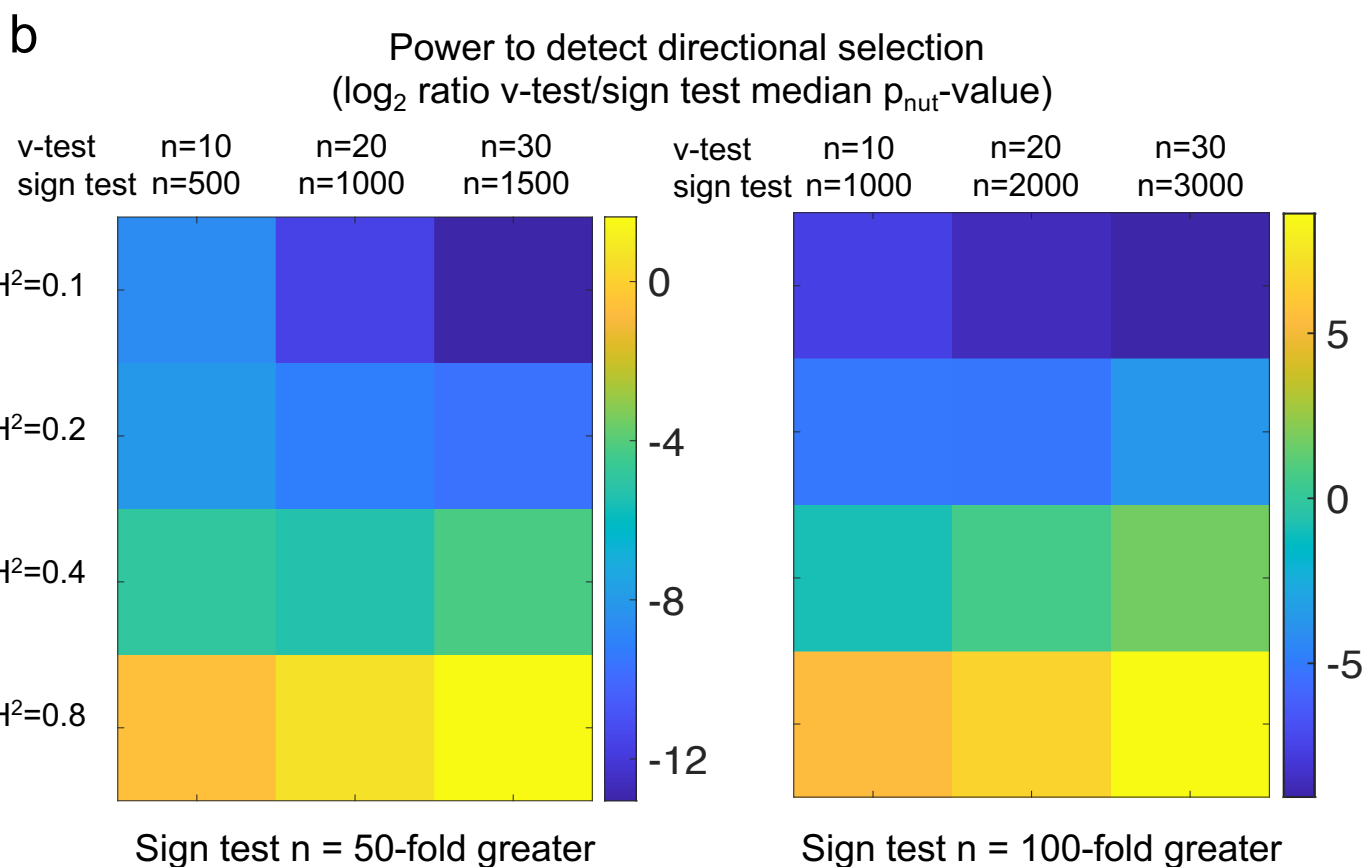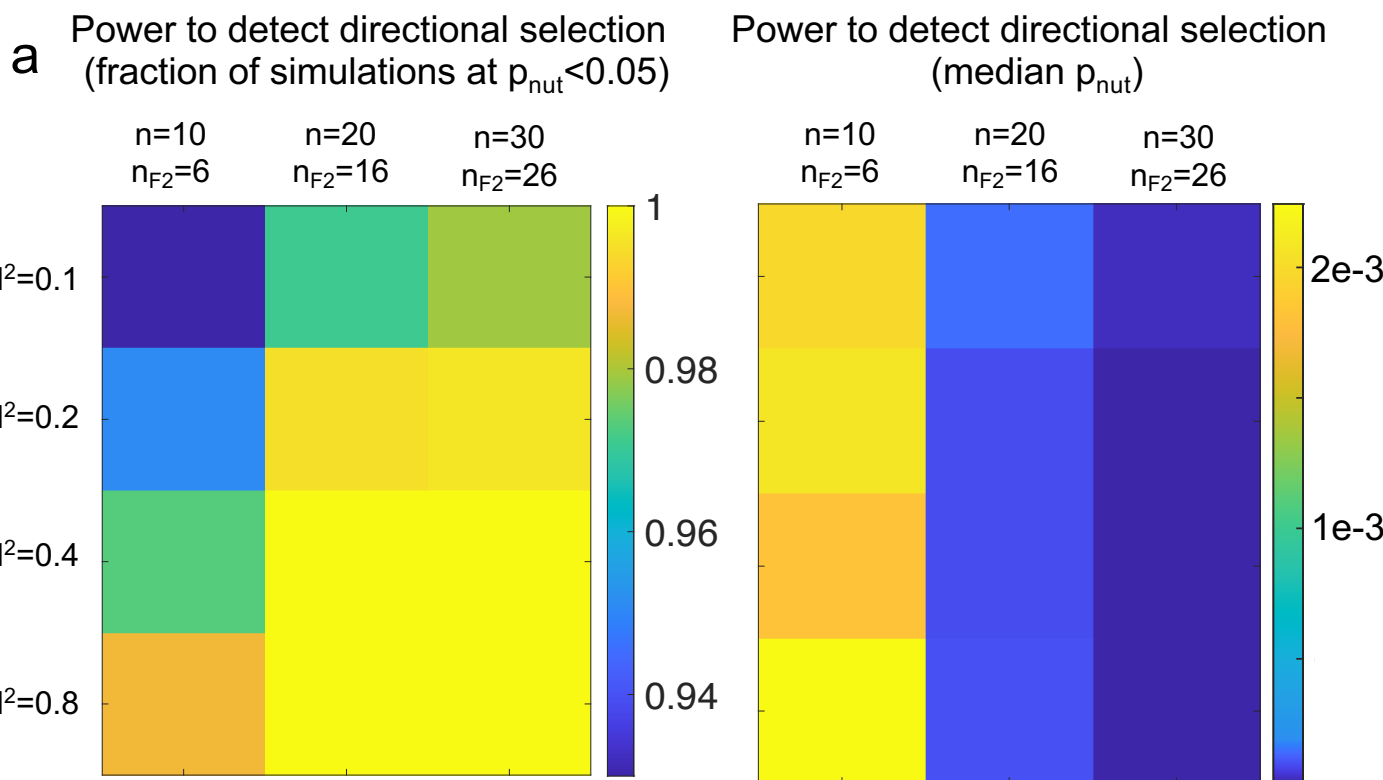Sign test n = 50-fold greater

Sign test n = 100-fold greater

# Supp Fig 2

Fraction of QTNs in same parental direction: