

1 **TITLE:**

2 A genotype-phenotype-fitness map reveals local modularity and global pleiotropy of adaptation

3

4 **AUTHORS:**

5 Grant Kinsler^{1*}, Kerry Geiler-Samerotte^{1,2*}, Dmitri Petrov¹

6

7 **AFFILIATIONS:**

8 ¹Department of Biology, Stanford University, Stanford, CA

9 ²Center of Mechanisms of Evolution, School of Life Sciences, Arizona State University, Tempe,
10 AZ

11 *Equal contribution

12

13 **SUMMARY:**

14 Building a genotype to phenotype to fitness map of adaptation is a central goal in evolutionary
15 biology. It is also notoriously difficult even when the adaptive mutations are known because it is
16 hard to identify which phenotypes make these mutations adaptive. We solve this problem by
17 first quantifying how the fitness of hundreds of adaptive mutants responds to subtle
18 environmental shifts and then inferring the existence of fitness-relevant phenotypes implicit in
19 these patterns of fitness variation. We find that a small number of phenotypes predicts the
20 fitness of the adaptive mutations near their original glucose-limited evolution condition.
21 Importantly, phenotypes that matter little to fitness at or near the evolution condition can matter
22 strongly in distant environments. This suggests that adaptive mutants are locally modular —
23 affecting a small number of phenotypes that matter to fitness in the environment where they
24 evolve — yet globally pleiotropic — affecting many phenotypes that contribute to fitness in new
25 environments.

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43 INTRODUCTION
44

45 High-replicate laboratory evolution experiments are opening an unprecedented window into the
46 dynamics and genetic basis of adaptive change by *de novo* mutation (Crozat et al., 2010; Good
47 et al., 2017; Lang et al., 2013; Levy et al., 2015; Tenaillon et al., 2012; Venkataram et al., 2016).
48 One of the key insights revealed by these studies is that in many systems, evolution can
49 proceed rapidly via many large-effect single mutations (Venkataram et al., 2016)(Crozat et al.,
50 2010; Good et al., 2017; Lang et al., 2013; Levy et al., 2015; Tenaillon et al., 2012)Venkataram
51 et al., 2016). While the identities of these adaptive mutations are often unique to a specific
52 replicate of the evolutionary experiment, across many replicates they tend to occur in similar
53 genes and pathways. Thus, while the diversity of mutations suggests that there might be many
54 ways to adapt, the much smaller number of apparent functional units implies, in contrast, that
55 most adaptive mutations affect a small set of key phenotypes (Fig 1A).

56 Consider the seminal study by Tenaillon et al. (Tenaillon et al., 2012) in which 115 populations
57 were evolved at high temperature for ~2000 generations. While the authors identified over a
58 thousand mutations that were largely unique to each population, the number of affected genes
59 was much smaller with 12 genes being hit over 25 times each. Even greater convergence was
60 seen at higher levels of organization such as operons. Similarly, Venkataram et al (Venkataram
61 et al., 2016) find that, of the hundreds of unique genetic mutations that occur during adaptation
62 to glucose-limitation, the vast majority fall into a relatively small number of genes (mostly *IRA1*,
63 *IRA2*, *GPB2*, *PDE2*) and primarily two pathways - Ras/PKA and TOR/Sch9. Thus despite the
64 diversity of mutations, it is possible that all of their effects can be mapped in one or few
65 dimensions required to describe their effects on the Ras/PKA or TOR/Sch9 pathways. These
66 are just two examples, but the pattern has been seen repeatedly (Barghi et al., 2019; Crozat et
67 al., 2010; Good et al., 2017; Lang et al., 2013; Lind et al., 2015). Note that this pattern is seen
68 not only in experimental evolution but also in cancer evolution. Individual tumors are largely
69 unique in terms of specific mutations, but these mutations affect a much smaller set of driver
70 genes and an even smaller number of higher functional units such as signalling pathways
71 (Hanahan and Weinberg, 2011, 2000).

72 The mapping of adaptive mutations to a smaller number of functional units and thus a low-
73 dimensional space representing their phenotypic effects (Fig 1A) is consistent with theoretical
74 models of adaptation. These theoretical models argue that adaptive mutations, especially those
75 of substantial fitness benefit, cannot affect too many phenotypes at once as most such effects
76 should be deleterious and thus inconsistent with the overall positive effect on fitness (Fisher,
77 1930; Orr, 2000). More recent studies likewise suggest that selection against mutations with
78 high pleiotropy, *i.e.* mutations that affect many phenotypes, has resulted in a modular
79 architecture of the genotype-phenotype map, in which genetic changes can influence some
80 phenotypes without affecting others (Altenberg, 2005; Collet et al., 2018; Melo et al., 2016;
81 Wagner et al., 2007; Wagner and Altenberg, 1996; Wagner and Zhang, 2011; Welch and
82 Waxman, 2003). This architecture would allow single mutations to drive large low-dimensional
83 phenotypic shifts. It would also explain the observations that very large collections of adaptive
84 mutations are not diverse in terms of affected genes, pathways, and phenotypes. The reason for
85 this is that only mutations that affect the genes, pathways, and phenotypes that represent the
86 correct module most relevant to adaptation in the study environment will be adaptive.

87 While theoretically appealing, the possibility that observed adaptive mutations indeed affect only
88 a very small number of phenotypes is difficult to reconcile with the notion that organisms are
89 tightly integrated (Kacser and Burns, 1981; Paaby and Rockman, 2013; Rockman, 2012).
90 Further, there is experimental evidence of widespread pleiotropy, for example, from genome
91 wide association studies that suggest every gene can influence every trait, at least to some
92 extent (Boyle et al., 2017; Chesmore et al., 2018; Sella and Barton, 2019; Sivakumaran et al.,

93 2011; Visscher and Yang, 2016). It is possible that pleiotropy is common, but strongly adaptive
94 mutations observed in experimental evolution are unusual in that they have few pleiotropic
95 phenotypic effects. Another possibility is that these mutations do have pleiotropic side effects
96 but these matter little to fitness in the evolution condition (Fig 1B). Note that here we do not
97 need to claim that these phenotypic effects *never* matter to fitness but rather that they do not
98 matter to fitness in the condition where they evolved. In fact the key prediction of this model is
99 that one should be able to detect such latent pleiotropy by showing that additional phenotypic
100 changes matter to fitness in other environments (Fig 1C).

101 If the model depicted in Fig 1B and C is true, then it is possible that adaptive mutations are
102 *locally modular* — that they affect very few phenotypes that matter to fitness in the evolution
103 condition — and *globally pleiotropic*. Under this model, the large number of distinct mutations
104 available to adaptation becomes important. Indeed while these mutations tend to influence
105 similar genes and pathways, their phenotypic effects do not simply collapse to a low
106 dimensional functional space. Instead this genetic diversity becomes a source of consequential
107 phenotypic diversity, but only once these genetic variants leave the local environment in which
108 they originated.

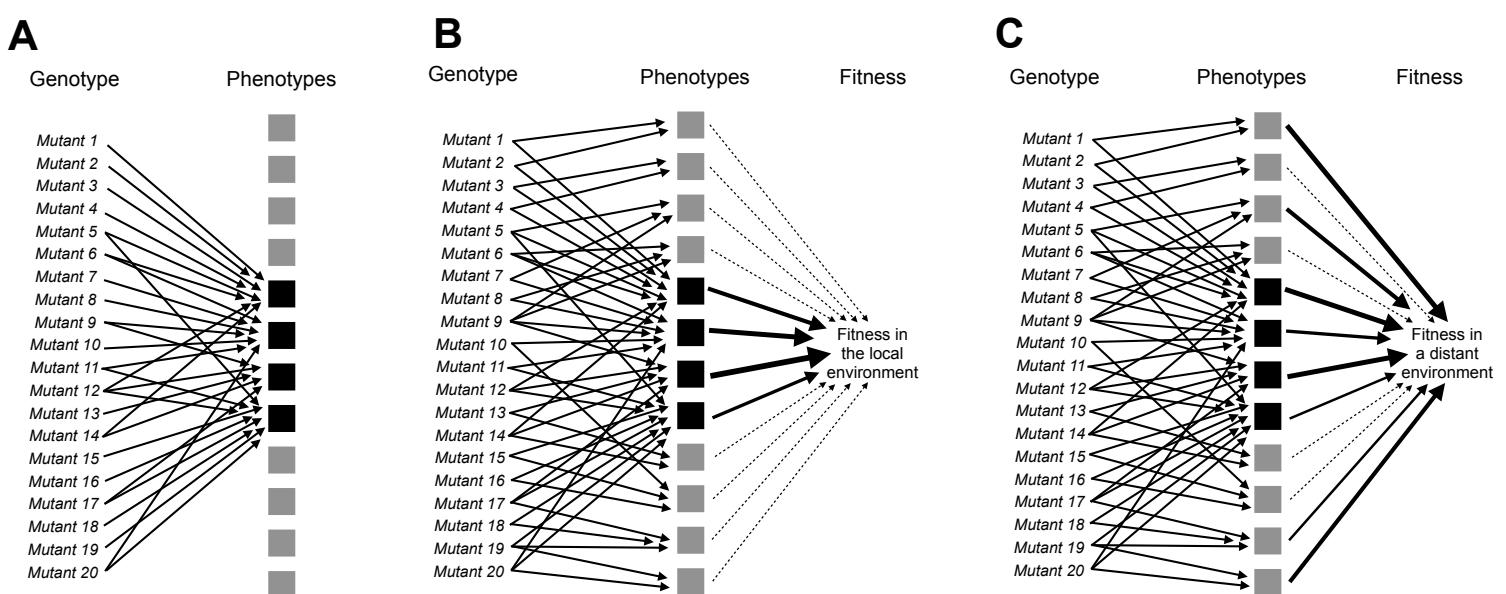
109 In order to test this model and better understand the genotype-phenotype-fitness map, we face
110 the difficult task of identifying which phenotypes adaptive mutations affect and then determining
111 how these phenotypes contribute to fitness. This is a challenging problem as the possible
112 number of phenotypes one can measure is virtually infinite, e.g. the expression level of every
113 gene or the quantity of every metabolite (Coombes et al., 2019; Mehlhoff et al., 2020). Further,
114 many measurable phenotypes are related in complex ways (Geiler-Samerotte et al., 2019).
115 Mapping their contribution to fitness requires a complete understanding of how genetic changes
116 lead to molecular changes and how these percolate to higher functional levels and ultimately
117 influence fitness (Kemble et al., 2020). This might be possible to do in some cases where the
118 phenotype to fitness mapping is simple (e.g. antibiotic resistance driven by a specific enzyme or
119 tRNA or protein folding mediating specific RNA or protein function) (Baeza-Centurion et al.,
120 2019; Cowperthwaite et al., 2005; Diss and Lehner, 2018; Domingo et al., 2019; Harmand et al.,
121 2017; Karageorgi et al., 2019; Li and Zhang, 2018; Otwinowski et al., 2018; Pressman et al.,
122 2019; Sarkisyan et al., 2016; Starr et al., 2018; Weinreich, 2006) but is exceptionally difficult for
123 complex phenotypes.

124 Moreover, to distinguish the model in Fig. 1A from Fig 1B, we need to understand these
125 genotype-phenotype-fitness maps not only in the environment in which adaptive mutants
126 evolved, but also in other environments, as depicted in Fig. 1C. And we need to do this for many
127 adaptive mutants so that we can assess the extent to which different mutants affect different
128 phenotypes. Considering the scope of this challenge, it is not surprising that despite much
129 theoretical discussion of modularity and pleiotropy as it relates to adaptation, experimental
130 approaches to address these questions have lagged behind.

131 Here we suggest a way to model the genotype-phenotype-fitness relationship that avoids the
132 problem of measuring each phenotype and its effect on fitness explicitly. We argue that it is
133 possible to investigate the genotype-phenotype-fitness map by comparing how the fitness
134 effects of many mutations change across a large number of environments. The way each
135 mutant's fitness varies across environments must be related to its phenotype, and thus the way
136 mutants co-vary in fitness across environments tells us whether they affect similar fitness-
137 relevant phenotypes. We can use the profiles of fitness across a set of environments to identify
138 the total number of fitness-relevant phenotypes affected across a collection of adaptive mutants,
139 the extent to which different mutants affect different phenotypes, and whether the contribution of
140 each phenotype to fitness changes across environments.

141

142 Here we build a genotype-phenotype-fitness model for hundreds of adaptive mutations that
 143 originally evolved in a glucose-limited environment. We use this model to accurately predict the
 144 fitness of these mutants across a set of 45 environments that vary in their similarity to the
 145 evolution condition. We find that the behavior of adaptive mutations can be described by a low-
 146 dimensional phenotypic model. In other words, these mutants affect a small number of
 147 phenotypes that matter to fitness in the glucose-limited condition in which they evolved. We find
 148 that this low-dimensional phenotypic model makes accurate predictions of mutant fitness in
 149 novel environments even when they are distant from the evolution condition. Moreover, we find
 150 that some phenotypes that contribute very little to fitness in the evolution condition become
 151 surprisingly important in some novel environments. This suggests that adaptive mutations are
 152 globally pleiotropic in that they affect many phenotypes overall, but that they are locally modular
 153 in that only a small number of these phenotypes have substantial effects on fitness in the
 154 environment they evolved in. Overall, we suggest that this set of adaptive mutations contains
 155 substantial and consequential latent phenotypic diversity, meaning that despite targeting similar
 156 genes and pathways, different adaptive mutants may respond differently to future evolutionary
 157 challenges. This finding has important consequences for understanding how directional
 158 selection can generate consequential phenotypic heterogeneity both in natural populations and
 159 also in the context of diseases such as cancer and viral or bacterial infections. In addition, our
 160 results show that our abstract, top-down approach is a promising route of analysis for
 161 investigating the phenotypic and fitness consequences of mutation.
 162



163 **Figure 1. Adaptive mutations can be locally modular and globally pleiotropic.** (A) A collection of
 164 adaptive mutations may affect a small number of phenotypes (four black squares). (B) Alternatively, these
 165 mutations may collectively (and individually) affect many phenotypes, but only a small number of
 166 phenotypes may matter to fitness (those indicated by black squares with thick arrows pointing to fitness),
 167 whereas the other phenotypes may make very small contributions to fitness (those indicated by the gray
 168 squares and thin, dashed lines leading to fitness). (C) Under the model in B, the contribution of each
 169 phenotype to fitness can change depending on the environment. Thus fitness differences between
 170 seemingly similar mutants can be revealed by measuring fitness in more environments. Such fitness
 171 differences suggest the presence of phenotypic differences between mutants.

172 RESULTS

173 Mutants that improve fitness under glucose limitation vary in their genotype-by- 174 environment interactions

176 A previous evolution experiment generated a collection of hundreds of independent mutations
177 that each provide a benefit to yeast cells growing in a glucose-limited environment (Levy et al.,
178 2015). Many of these mutants, which began the evolution experiment as haploids, underwent
179 whole-genome duplication to become diploid, which improved their relative fitness (Venkataram
180 et al 2016). Some of these diploids acquired additional mutations, including amplifications of
181 either chromosome 11 or 12 as well as point mutations, which generated additional fitness
182 benefits. The adaptive mutants that remained haploid acquired both gain- and loss-of-function
183 mutations in nutrient-response pathways (Ras/PKA and TOR/Sch9). Some other mutations
184 were also observed, including a mutation in the HOG pathway gene SSK2 (Venkataram et al.,
185 2016). Although these mutants have been well-characterized at the level of genotype and
186 fitness, it is unclear what phenotypes they affect. The first question we address is whether these
187 diverse mutations collectively affect a large number of phenotypes that matter to fitness, or
188 whether these mutants are functionally similar in that they collectively alter a small set of fitness-
189 relevant phenotypes.

190
191 Understanding the map from genotype to phenotype to fitness is extremely challenging because
192 each genetic change can influence multiple traits, not all of which are independent or contribute
193 to fitness in a meaningful way. We contend with this challenge by measuring how the relative
194 fitness of each adaptive mutant changes across a large collection of similar and dissimilar
195 environments, which we term the “fitness profile”. When a group of mutants demonstrate similar
196 responses to environmental change, we conclude that these mutants affect similar phenotypes.
197 By clustering mutants with similar fitness profiles across a collection of environments, we can
198 learn about which mutants influence similar phenotypes, as well as estimate the total number of
199 fitness-relevant phenotypes represented across all mutants and all investigated environments.
200

201 Because our mutant strains are barcoded, we can use previously-established methods to
202 measure their relative fitness in bulk and with high-precision (Venkataram et al., 2016).
203 Specifically, we compete a pool of the barcoded mutants against an ancestral reference strain
204 over the course of several serial dilution cycles. During each 48 hour cycle, the yeast are given
205 fresh glucose-limited media which supports 8 generations of exponential growth after which
206 glucose is depleted and cells transition to non-fermentable carbon sources. After every 48 hour
207 cycle, we transfer $\sim 5 \times 10^7$ cells to fresh media to continue the growth competition. We also
208 extract DNA from the remaining cells to PCR amplify and sequence their barcodes. We repeat
209 this process four times, giving us an estimate of the frequency of each barcode at five time-
210 points. By quantifying the log-linear changes in each barcode’s frequency over time and
211 correcting for the mean-fitness change of population, we can calculate the fitness of each
212 barcoded mutant relative to the reference strain (Fig 2A; Methods).

213
214 Using this method, we quantify the fitness of a large number of adaptive mutants in 45
215 environments. We focus on a set of 292 adaptive mutants that have been sequenced, show
216 clear adaptive effects in the glucose-limited condition in which these mutants evolved (hereafter
217 “evolution condition”; EC) (Fig 2B; Table S1), and for which we obtained high-precision fitness
218 measurements in all 45 environments. These environments include some experiments from
219 previously published work (Li et al., 2018; Venkataram et al., 2016), as well as 32 new
220 environments including replicates of the evolution condition, subtle shifts to the amount of
221 glucose, changes to the shape of the culturing flask, changes to the carbon source, and addition
222 of stressors such as drugs or high salt (Table S2).

223
224 In order to determine the total number of phenotypes that are relevant to fitness in the EC, we
225 focus on environments that are very similar to the EC but still induce small yet detectable
226 perturbations in fitness. We do so because the phenotypes that are the most relevant to fitness
227 may change with the environment (Fig 1B and Fig 1C). Thus, we partition the 45 environments
228 into a set of “subtle” perturbations, from which we will detect the phenotypes relevant to fitness

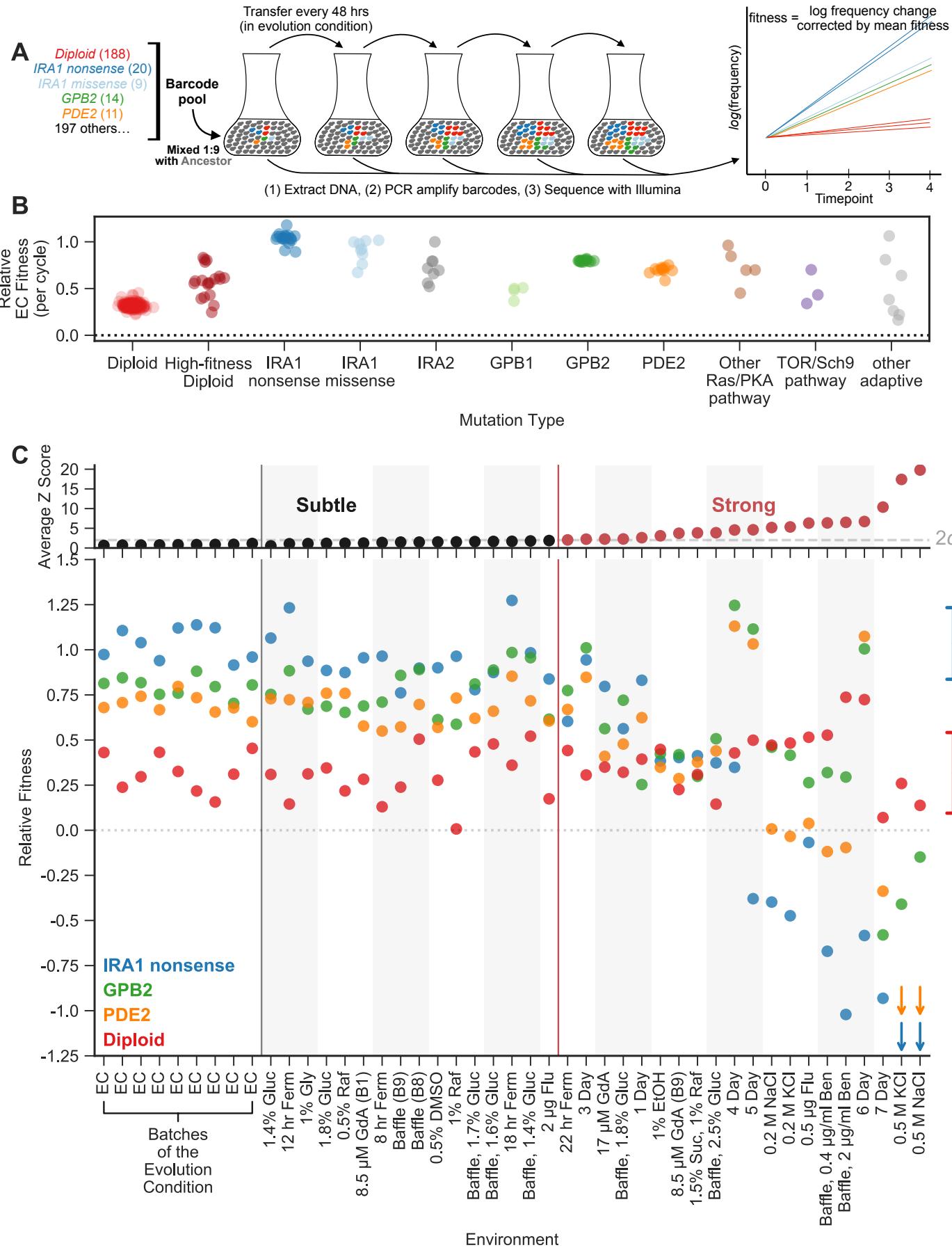
229 near the EC, and “strong” perturbations which we will use to study whether these mutants
230 influence additional phenotypes that matter in other environments (Fig 1C).

231
232 To partition environments into subtle and strong perturbations of the EC, we rely on the nested
233 structure of replicate experiments performed in the EC. We performed nine such replicates,
234 each at different times, which each included multiple replicates performed at the same time. We
235 observe much less variation across replicates performed simultaneously than across replicates
236 performed at different times ($p < 1e-5$ from permutation test). Variation across experiments
237 performed at different times is often referred to as “batch effects” and likely reflects
238 environmental variability that we were unable to control (e.g. slight fluctuations in incubation
239 temperature due to limits on the precision of the instrument). These environmental differences
240 between batches are as subtle as possible, as they represent the limit of our ability to minimize
241 environmental variation. Thus, variation in fitness across the EC batches serves as a natural
242 benchmark for the strength of environmental perturbations. If the deviations in fitness caused by
243 an environmental perturbation are substantially stronger than those observed across the EC
244 batches, we call that perturbation “strong”.

245
246 More explicitly, to determine whether a given environmental perturbation is subtle or strong, we
247 subtract the fitness of adaptive mutants in this environment from their average across the EC
248 batches. We then compare this difference to the variation in fitness observed across the EC
249 batches. Sixteen environmental perturbations provoked fitness differences that were similar to
250 those observed across EC batches (Z-score < 2). These environments, together with the nine
251 EC batches, make up a set of subtle environmental perturbations. The remaining 20
252 environments, where the average deviation in fitness is substantially larger than that observed
253 across batches (Z-score > 2), were classified as strong environmental perturbations (Fig 2C,
254 top; Methods).

255
256 The rank order of the fitnesses of many mutations is largely preserved across the 25
257 environments that represent subtle perturbations (Fig 2C, bottom). For example, *IRA1* nonsense
258 mutants, which are the most adaptive in the EC, generally remain the most adaptive across the
259 subtle perturbations. Additionally, the *GPB2* and *PDE2* mutants have similar fitness effects
260 across EC batches and only occasionally switch order across the subtle environmental
261 perturbations. In contrast, the 20 environments that represent strong perturbations reveal clear
262 genotype-by-environment interactions (Fig 2C, bottom). For example, altering the transfer time
263 from 48 to 24 hours (the “1 Day” environment in Fig 2C) affects *GPB2* mutants more strongly
264 compared to the other mutants in the Ras/PKA pathway, including *IRA1* and *PDE2*. The
265 strongest environmental perturbations reveal clear tradeoffs for some of these adaptive
266 mutants. For example, *PDE2* and *IRA1-nonsense* but not *GPB2* mutants are particularly
267 sensitive to osmotic stress as indicated by the NaCl and KCl environments. Additionally, *IRA1-*
268 *nonsense* mutants become strongly deleterious in the long transfer conditions that experience
269 stationary phase (5-, 6-, 7-Day environments) (Li et al., 2018). In contrast to complex behavior
270 exhibited by the adaptive haploids, the diploids appear to be relatively robust to strong tradeoffs,
271 appearing similarly adaptive across all perturbations, subtle and strong.

272
273 The observation that different mutants have different and fairly complex fitness profiles suggests
274 that they have different phenotypic effects. Even *PDE2* and *GPB2*, which have similar fitnesses
275 in the EC and are negative regulators of the same signalling pathway, have different fitness
276 profiles. Do these diverse phenotypic effects contribute to fitness in the EC? To examine how
277 many phenotypes matter to fitness in the EC, we test whether it is possible to create low
278 dimensional models that capture the complexity of the fitness profiles of all adaptive mutants
279 across all subtle perturbations.



389 average 85% of weighted variance for the test mutants in the left-out conditions. A model with
390 only the top five components explains 95.1%, and all eight components explain 96.2% of
391 variation. This suggests that the last few components have very small contributions to fitness in
392 the environments near the EC.

393

394 **A model including 8 fitness-relevant phenotypes recapitulates known features of 395 adaptive mutations**

396

397 We next ask whether the 8-dimensional phenotypic space clusters adaptive mutants found in
398 similar genes or pathways (e.g. Ras/PKA or TOR/Sch9), or that represent similar mutation types
399 (haploid v. diploid). We use Uniform Manifold Approximation and Projection (UMAP) to visualize
400 the distance between all the mutants in this phenotypic space. Since the first phenotypic
401 dimension captures the average fitness of each mutant in the EC, and since we already know
402 that mutations to the same gene have similar fitness in the EC (Fig. 2B), we exclude the first
403 phenotypic dimension from this analysis, though the inclusion of the first component does not
404 change the identity of the clusters (Fig S3). By focusing on the other 7 components, we are
405 asking whether genotype-by-environment interactions also cluster the mutants by gene,
406 mutation type, and pathway.

407

408 These 7 genotype-by-environment interactions indeed cluster the adaptive mutants by type and
409 by gene (Fig 3B). Specifically, the diploids, *IRA1-nonsense*, *GPB2*, and *PDE2* mutants each
410 form distinct clusters ($p = 0.0001$, $p = 0.006$, $p = 0.0001$, and $p = 0.0001$, respectively).
411 Interestingly, six high-fitness diploids (diploids with higher than average diploid fitness in the EC)
412 also form a distinct cluster ($p = 0.0001$) despite whole genome sequencing having revealed no
413 mutations in their coding sequences (Fig. 3B). To generate p-values, we calculated the median
414 pairwise distance, finding that multiple mutations in the same cluster are indeed more closely
415 clustered than randomly chosen groups of mutants.

416

417 Interestingly, the three smallest components, which capture very little variation in fitness across
418 the environments that reflect subtle perturbations of the EC, cluster some mutants by genotype
419 (Fig S3). Specifically, *PDE2*, *GPB2*, and *IRA1-nonsense* mutants are each closer to mutants of
420 their own type than to other adaptive haploids ($p = 0.0001$, $p = 0.0001$, and $p = 0.03$,
421 respectively). Note that the space defined by the three smallest components does not cluster
422 *IRA1-nonsense* mutants away from diploids ($p = 0.718$). This suggests that some mutants, e.g.
423 *IRA1-nonsense* and diploids, have smaller effects on these three phenotypic components.
424 Overall, our abstract phenotypic model, which reflects the way that each mutant's fitness
425 changes across environments, reveals that mutations to the same gene tend to interact similarly
426 with the environment. This indicates that our approach is a useful and unbiased way to identify
427 mutations that share functional effects (Li et al., 2018).

444 environmental perturbations reveals 8 fitness-relevant phenotypic components. The variance explained
445 by each component is indicated as a percentage of the total variance. The percentages in parentheses
446 indicate the relative amount of variation explained by each component when excluding the first
447 component. Each of these components explain more variation in fitness than do components that capture
448 variation across a simulated dataset in which fitness varies due to measurement noise. These simulations
449 were repeated 1000 times (grey lines) and used to define the limit of detection (dotted line). **(C)** An
450 abstract space containing 8 fitness-relevant phenotypic components reflects known biological features.
451 This plot shows the relationships of the mutants in a 7-dimensional phenotypic space that excludes the
452 first component, visualized using Uniform Manifold Approximation and Projection (UMAP). Mutants that
453 are close together have similar fitness profiles and are inferred to have similar effects on fitness-relevant
454 phenotypes. Mutants with mutations in the same gene tend to be closer together than random, in
455 particular *IRA1* nonsense mutants in dark blue, *GPB2* mutants in dark green, *PDE2* mutants in dark
456 orange, and diploid mutants in red. Six high-fitness diploid mutants also form a cluster despite no known
457 genetic similarities.
458

459 **Fitness variation across subtly different environments predicts fitness in substantially
460 different environments**

461 Now that we have identified the phenotypic components that contribute to fitness in
462 environments that represent subtle perturbations of the EC, we can test the ability of these
463 phenotypic components to predict fitness in more distant environments. Specifically, we can
464 measure how the contribution of each of these components to fitness changes in new
465 environments. We can also determine whether the phenotypic components that contribute very
466 little to explaining fitness variation near the EC might at times have large explanatory power in
467 distant environments (as depicted in Fig 1B and 1C).
468

469 To test this we performed bi-cross-validation, using the eight component model constructed
470 from fitness variation of 60 training mutants across 25 subtly different environments to predict
471 the fitness of 232 test mutants in the environments that represent strong perturbations of the
472 EC. To evaluate the predictive power of the model, we compare our model's fitness predictions
473 in each environment to predictions made using the average fitness in that environment. Thus,
474 negative prediction power indicates cases where the model predicts fitness worse than
475 predictions using this average.
476

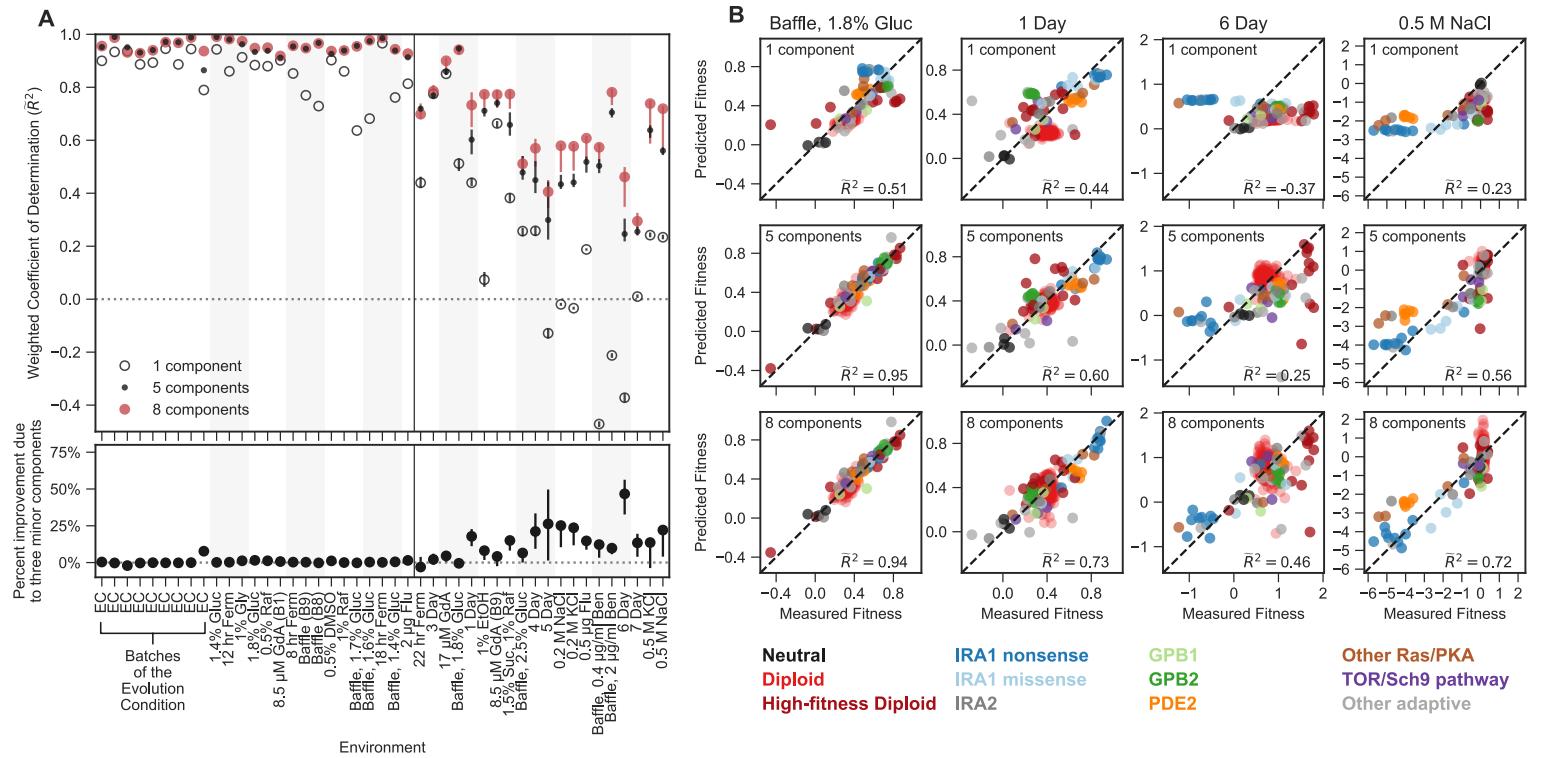
477 The 8-dimensional phenotypic model, which was generated exclusively with the data from
478 subtle environmental perturbations, has substantial predictive power in distant environments
479 (Figure 4). Predictions explain 29% to 95% of the variation in fitness of the 232 test mutants
480 across strong environmental perturbations. For instance, in an environment where glucose
481 concentration was increased from 1.5% to 1.8% and the flask was changed to one that
482 increases the oxygenation of the media (the “Baffle, 1.8% Glucose” environment), we predict
483 95% of weighted variance with the full 8-component phenotypic model, in contrast to 51% with
484 the 1-component model (Fig 4B). This ability to predict fitness is retained even when the first
485 component (effectively the fitness in EC) is a poor predictor of mutant fitness. For example, in
486 the environment where salt (0.5 M NaCl) was added to the media, the 1-component model
487 predicts fitness worse than predictions based on the average fitness for this environment,
488 resulting in negative variance explained (Fig 4A and 4B). This is due to the fact that mutant
489 fitness in this environment reflects extensive genotype-by-environment interactions, such that
490 the fitness of mutants in this environment is uncorrelated with EC fitness. However, our
491 predictions of mutant fitness in the 0.5 M NaCl environment improve when made using the 8-
492 component phenotypic model, which predicts 72% of weighted variance. Astoundingly, the 8-
493 component model captures strong tradeoffs between mutants with high fitness in the EC and
494 very low fitness in this high salt environment, specifically for *IRA1-nonsense* and, to a lesser
495 extent, *PDE2* mutants (Fig 4B). This was surprising because there appears to be very little
496 variation in fitness of these mutants across the subtle compared to the strong perturbations (Fig
497 2C).
498

499
500
501
502
503
504
505
506
507
508

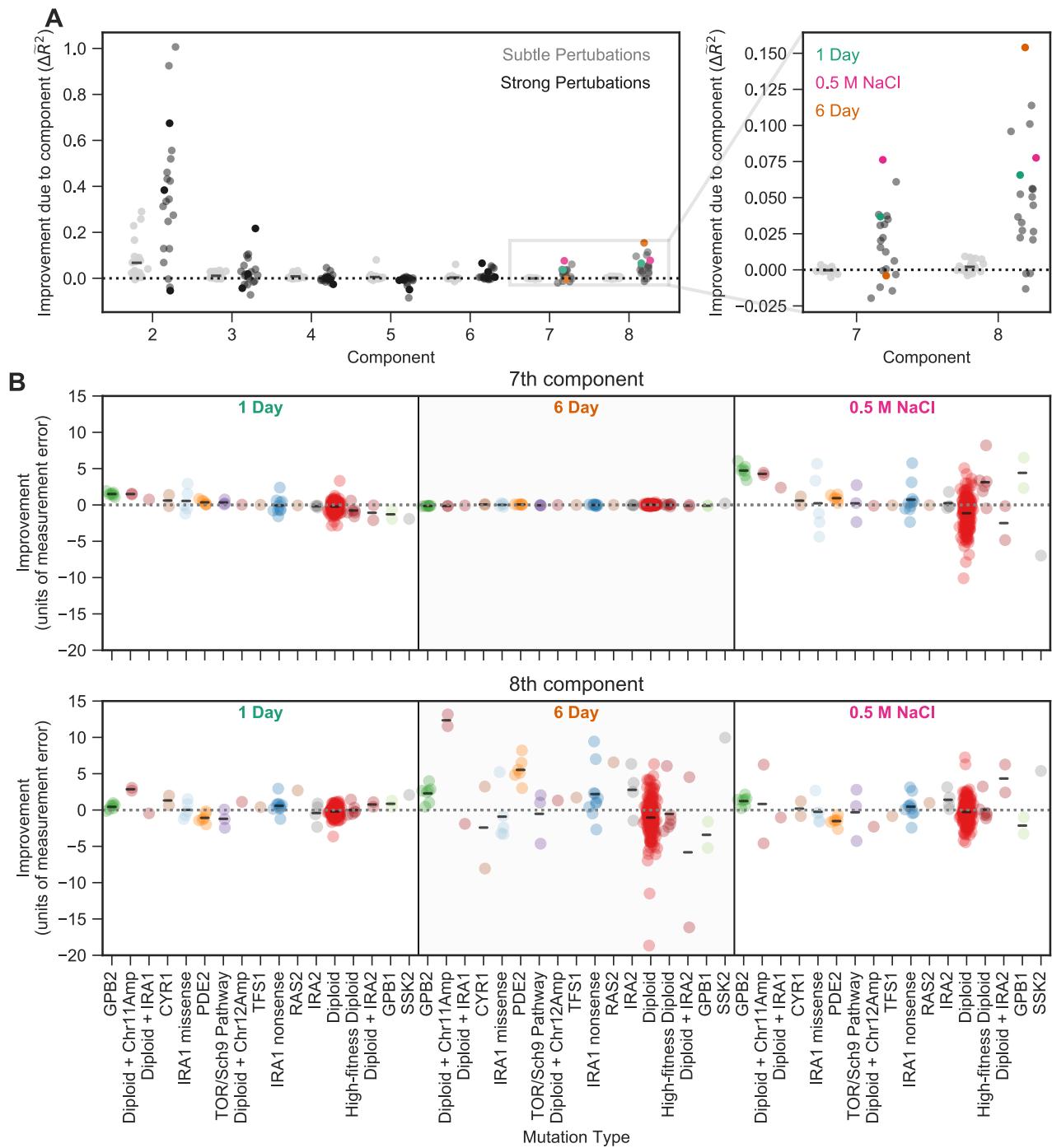
This ability to predict fitness is also observed for mutations in genes and pathways that are not represented in the 60 that comprise the training set (e.g. those with mutations in TOR/Sch9 and HOG pathway genes). For example, the 8-component model explains 93% of variation in the "Baffle, 1.8% Glucose" environment and 71% of variation in the 0.5M NaCl environment for these mutations, compared to 76% and 31% variance explained for the 1-component model, respectively. This indicates that our model is able to capture shared phenotypic effects that extend beyond gene identity. Altogether, our ability to accurately predict the fitness of new mutants in new environments suggests that the phenotypes our model identifies reflect causal effects on fitness.

509
510
511
512
513
514
515
516
517
518
519
520
521
522
523

Most strikingly, phenotypic models that include the three smallest phenotypic components, which together contribute only 1.1% to variance explained across the subtle environmental perturbations (Fig 4A), often explain a substantial amount of variance in the distant environments (Fig 4A; lower panel). For example, the three minor components contribute 17% of the overall weighted variance explained in the 1-Day condition ($\tilde{R}^2 = 0.6$ - 5-component model, $\tilde{R}^2 = 0.73$ - 8-component model; $(0.73-0.6)/0.73 = 0.17$) and 45% in the 6-Day environment, ($\tilde{R}^2 = 0.25$ - 5-component model, $\tilde{R}^2 = 0.46$ - 8-component model) (Fig 4A and 4B). In contrast, for other strong environments (e.g. Baffle - 1.8% Glucose, 8.5uM GdA (B9) and Baffle - 2.5% Glucose), the three smallest components do not add much explanatory power (Fig 4A). These observations demonstrate that phenotypic components that make very small contributions to fitness in the EC can contribute substantially to fitness in other environments. Overall, these observations suggest an answer to questions about how adaptation is possible when mutations have collateral effects on multiple phenotypes: not all of those phenotypes contribute substantially to fitness in the EC (Fig 1C).



524 **Figure 4. Mutant fitness variation across subtly different environments predicts mutant fitness in**
525 **novel and substantially different environments. (A) Top panel** vertical axis shows the accuracy of
526 fitness predictions in each of 45 environments on the horizontal axis. The accuracy is calculated as the
527 coefficient of determination, weighted such that each mutation type contributes equally. The left side of
528 this plot represents predictions of mutant fitness in subtle environmental perturbations. These predictions
529 are generated by holding out data from that environment when building the phenotypic model. The right
530 side of the plot displays predictions of mutant fitness in strong environmental perturbations. These
531 predictions are generated using a phenotypic model inferred from fitness variation across all 25 subtle
532 different environments (denoted by each of the points or open circles) and for each of the 25 leave-one-out
533 models (range of predictions is depicted with the error bars surrounding each point or open circle).
534 Predictions from the 8-component model (red point) are typically better than the 1-component mode
535 (open circle) and sometimes better than the 5-component model (black point). **Bottom panel** vertical axis
536 shows the percent of the 8-component model's improvement due to the three minor components
537 (calculated by the percent difference between the 5- and 8- component models). The left side shows the
538 improvement of the prediction in subtle environmental perturbations when that subtle perturbation was
539 held out. The right side shows the improvement of the prediction in strong environmental perturbations
540 when using the full model (dots) or the 25 leave-one-out models (the error bars represent the range of
541 improvement). **(B)** For each subplot, the horizontal axis shows the measured fitness value. The vertical
542 axis shows the predicted fitness value when predictions are made using the 1-component (top row), 5-
543 component (middle row), or 8-component (bottom row) models. Columns represent different
544 environments. Points are colored by the mutation type. Note that \bar{R}^2 less than zero indicates that the
545 prediction is worse than predictions using the mean fitness in that condition (see Methods).



597 **Figure 5. The contribution of a phenotypic component to fitness changes across environments**
598 **and differs for different types of mutants.** (A) Some phenotypic components improve fitness
599 predictions in some environments substantially more than they do in others. The vertical axis shows the
600 improvement in the predictive power of our 8-component phenotypic model due to the inclusion of each
601 component. For example, the improvement due to component 7 is calculated by the difference between
602 the 7-component model and the 6-component model. The improvement of predictive power for each of
603 the subtle environmental perturbations is shown as a gray point and for each of the strong perturbations
604 in black. Magnification shows improvement upon including each of the two smallest components, with
605 three strong perturbations highlighted. (B) Some phenotypic components improve fitness predictions for
606 some mutants substantially more than they do for others. For example, the 7th component explains little
607 variation in the 6-Day environment, but the 8th component explains a lot of variation in fitness in the 6-
608 Day environment and is particularly helpful in predicting the fitness of Diploid + Chromosome 11
609 Amplification mutations in this environment. Vertical axis shows the improvement in predictive power (in
610 units of standard deviation of measurement error) for each type of mutant (denoted on the horizontal axis)
611 in one of three environments (1 Day, 6 Day, and 0.5 M NaCl) when adding either the 7th (top panel) or
612 the 8th (bottom panel) component. Mutants are ordered by the improvement due to the 7th component in
613 the 1 Day environment. Since some types of mutants are more common, e.g. diploids, there are more
614 data points in that category.

615 **DISCUSSION**

616
617 Here we succeeded in building a low-dimensional statistical model that captures the relationship
618 from genotype to phenotype to fitness for hundreds of adaptive mutants. Mapping the complete
619 phenotypic and fitness impacts of genetic change is a key goal of biology. Such a map is
620 important in order to make meaningful predictions from genetic data (e.g. personalized
621 medicine) and to investigate the structure of biological systems (e.g. their degree of modularity
622 and pleiotropy) (500 Genomes Field Experiment Team et al., 2019; Collet et al., 2018; Eguchi et
623 al., 2019; Zan and Carlborg, 2020). Our model allows us to do both of these things. We made
624 accurate predictions about the fitness of unstudied mutants across multiple environments, and
625 we gained novel insights about the nature of pleiotropy of the adaptive process. Specifically, we
626 learned that adaptation is modular in the sense that hundreds of diverse adaptive mutants
627 collectively influence a small number of phenotypes that matter to fitness in the evolution
628 condition. We also learned that different mutants have distinct pleiotropic side effects that matter
629 to fitness in other conditions.

630
631 Building genotype-phenotype-fitness maps of adaptation has long been an elusive goal due to
632 both conceptual and technical difficulties. Indeed, the very first part of this task, namely the
633 identification of causal adaptive mutations, presents a substantial technical challenge (500
634 Genomes Field Experiment Team et al., 2019; Barrett et al., 2019, 2008). Fortunately, in some
635 systems such as in microbial experimental evolution and studies of cancer and resistance in
636 microbes and viruses, genomic methodologies combined with availability of repeated
637 evolutionary trials allow us now to detect with high confidence specific genetic changes
638 responsible for adaptation. In the context of microbial evolution experiments, lineage tracing and
639 genomics have opened up the possibility of not only detecting hundreds of specific adaptive
640 events but also measuring their fitness precisely and in bulk (Good et al., 2017; Levy et al.,
641 2015; Li et al., 2019, 2018; Nguyen Ba et al., 2019; Venkataram et al., 2016). Thus in these
642 cases we are coming close to solving the technical challenge of building the genotype to fitness
643 map of adaptation.

644
645 However, adding phenotype into this map remains a huge challenge even despite substantial
646 progress in mapping genotype to phenotype (Burga et al., 2019; Camp et al., 2019; Exposito-
647 Alonso et al., 2018; Geiler-Samerotte et al., 2016; Jakobson and Jarosz, 2019; Lee et al., 2019;
648 Paaby et al., 2015; Yengo et al., 2018; Ziv et al., 2017). In principle, we now have advanced
649 tools to measure a large number of phenotypic impacts of a genetic change, for instance
650 through high throughput microscopy, proteomics, or RNAseq (Manzoni et al., 2018; Ritchie et
651 al., 2015; Zhang and Kuster, 2019). The conceptual problem is how to define phenotypes given
652 the interconnectedness of biological systems (Geiler-Samerotte et al., 2019; Paaby and
653 Rockman, 2013). If a mutation leads to complex changes in cell size and shape, should each
654 change be considered a distinct phenotype? Or if a single mutation changes the expression of
655 hundreds or thousands of genes, should we consider each change as a separate phenotype?
656 Intuitively, it seems that we should seek higher order, more meaningful descriptions. For
657 example, perhaps these expression changes are coordinated and reflect the up-regulation of a
658 stress-response pathway. Unfortunately, defining the functional units in which a gene product
659 participates remains difficult, especially because these units re-wire across genetic
660 backgrounds, environments, and species (Geiler-Samerotte et al., 2019; Pavličev et al., 2017;
661 Sun et al., 2020; Zan and Carlborg, 2020).

662
663 If mutations influence more than one phenotype, then the mapping from phenotype to fitness
664 also becomes challenging. To investigate this map, we would need to find an artificial way to
665 perturb one phenotype without perturbing others such that we could isolate and measure effects
666 on fitness. Mapping phenotype to fitness is further complicated by the environmental
667 dependence of these relationships (Fragata et al., 2019). For example, a mutation that affects a

668 cell's ability to store carbohydrates for future use might matter far more in an environment where
669 glucose is re-supplied every 6 days instead of every 48 hours.

670
671 In our study, we turned the challenge of environment-dependence into the solution to the
672 seemingly intractable problem of interrogating the phenotype layer of the genotype-phenotype-
673 fitness map. We rely on the observation that the relative fitness of different mutations changes
674 across environments. We assume that differences in how mutant fitness varies across
675 environments must stem from differences in the phenotypes each mutation affects. Rather than
676 *a priori* defining the phenotypes that we think may matter, we use the similarities and
677 dissimilarities in the way fitness of multiple mutants vary across environments to define
678 phenotypes abstractly via their causal effects on fitness. This allows us to dispense with
679 measuring the phenotypes themselves and instead focus on measuring fitness with high
680 precision and throughput, since tools for doing so already exist (Venkataram et al., 2016). This
681 approach has the disadvantage of not identifying phenotypes in a traditional, more transparent
682 way. Still, it represents a major step forward in building genotype-phenotype-fitness maps
683 because it makes accurate predictions and provides novel insights about the phenotypic
684 structure of the adaptive response.

685
686 We successfully implemented this approach using a large collection of adaptive mutants
687 evolved in a glucose-limited condition. The first key result is that the map from adaptive mutant
688 to phenotype to fitness is modular, such that it is possible to create a genotype to phenotype to
689 fitness model that is low dimensional. Indeed, our model detects a small number (8) of fitness-
690 relevant phenotypes, the first two of which explain almost all of the variation in fitness (98.3%)
691 across 60 adaptive mutants in 25 environments representing subtle perturbations of the
692 glucose-limited evolution condition. This suggests that the hundreds of adaptive mutations we
693 study — including mutations in multiple genes in the Ras/PKA and TOR/Sch9 pathways,
694 genome duplication (diploidy), and various structural mutations — influence a small number of
695 phenotypes that matter to fitness in the evolution condition. This observation is consistent with
696 theoretical considerations suggesting that mutations that affect a large number of fitness-
697 relevant phenotypes are not likely to be adaptive. It also explains findings from other high-
698 replicate laboratory evolution experiments and studies of cancer that show hundreds of unique
699 adaptive mutations tend to hit the same genes and pathways repeatedly. Our work confirms the
700 intuition that these mutations all affect similar higher-order phenotypes (e.g. the level of activity
701 of a signaling pathway). This suggests that, despite the genetic diversity among adaptive
702 mutants, adaptation may be predictable and repeatable at the phenotypic level.

703
704 Note that although we detect only 8 fitness-relevant phenotypes, we expect the true number to
705 be much larger as the detectable number is limited by the precision of measurement. We expect
706 this partly because we know that if we had worse precision in this experiment we would have
707 detected fewer than 8 phenotypic components (Fig 3). Still, these additional undetected
708 components cannot be very consequential in terms of their contribution to fitness in the
709 evolution condition, given how well the first 8 components capture variation in environments that
710 are similar to the evolution condition.

711
712 Surprisingly, the model built only using subtle environmental perturbations was also predictive of
713 fitness in environments that perturbed fitness strongly. Indeed in some of these environments,
714 such as the environment where 0.5 M NaCl was added to the media or the time of transfer was
715 extended from two to six days, many of the mutants are no longer adaptive and some of them
716 become strongly deleterious. Here the fitness of the mutants in the evolving condition is a very
717 poor predictor of fitness but the full 8-dimensional phenotypic model explains from 29% to 95%
718 of the variance. What was particularly interesting is that the explanatory power of different
719 dimensions was very different for the strong compared to subtle perturbations. For instance, the
720 second dimension which explained 7% of weighted variation on average in the subtle

721 perturbations, explained 36% on average in the environments that represent strong
722 perturbations. The pattern was particularly striking for the smallest dimensions which at times
723 explained 15% in the strong environmental perturbations while again explaining at most 1% in
724 the subtle environments.

725
726 This discovery emphasizes that, although the smaller phenotypic dimensions contribute very
727 little to fitness in the evolution condition (Fig 1B), they can at times have a much larger
728 contribution in other environments (Fig 1C). This makes intuitive sense. For instance, we know
729 that some of the strongest adaptive mutations in our experiment, the nonsense mutations in
730 *IRA1*, appear to stop cells from shifting their metabolism towards carbohydrate storage when
731 glucose levels become low (Li et al., 2018). This gives these cells a head start once glucose
732 again becomes abundant and does not appear to come at a substantial cost, at least not until
733 these cells are exposed to stressful environments (e.g. high salt or long stationary phase) (Li et
734 al., 2018). This example, and more generally the observation that phenotypic effects that are
735 unimportant in the evolving condition can become much more important in other environments,
736 supports the idea that adaptation can happen through large effect mutations because many of
737 the pleiotropic phenotypic effects will be inconsequential in the local environment (Fig 1B – C).
738 We can thus argue that the low-dimensional phenotypic space near the evolution condition
739 hides *latent* and *consequential* phenotypic complexity across the collection of locally
740 phenotypically similar mutants. This complexity is hidden from natural selection in the evolution
741 condition but becomes important once the mutants leave the local environment and are
742 assessed globally for fitness effects. Thus, with respect to their effects on fitness-relevant
743 phenotypes, adaptive mutants may be locally modular, but globally pleiotropic.
744

745 The notion of latent phenotypic complexity is exciting as it generates a mechanism by which
746 directional selection generates rather than removes diversity. This suggests a solution to long-
747 standing questions in evolutionary biology about how diversity persists despite directional
748 selection (Walsh and Blows, 2009). Directional selection may promote multiple mutants that
749 affect similar fitness-relevant phenotypes in the evolution condition, but each mutant could have
750 disparate meaningful phenotypic effects that do not contribute immediately to fitness. When the
751 environment changes, these latent phenotypic effects may now matter, allowing for diverse
752 solutions to a variety of possible new environments. This latent phenotypic complexity also has
753 the potential to alter the future adaptive paths that a population takes even in a constant
754 environment. Indeed, these phenotypically diverse mutants are likely to affect the subsequent
755 direction of adaptation given that subsequent mutations can shift the context in which
756 phenotypes are important in the same way as do environmental perturbations. The end result is
757 that directional selection can enhance diversity both within a population and across populations
758 adapting to the same stressors.
759

760 The phenomenon of latent phenotypic complexity being driven by adaptation is dependent on
761 there being multiple mutational solutions to an environmental challenge, such that different
762 adaptive mutations might have different latent phenotypic effects. Latent phenotypic diversity
763 might be less apparent in cases where adaptation proceeds through mutations in a single gene
764 and certainly would not exist if adaptation relies on one unique mutation. Thus, in some ways
765 latent phenotypic diversity reflects redundancies in the mechanisms that allow cells to adapt to a
766 challenge. One such putative redundancy in the case investigated in this paper is that the
767 Ras/PKA pathway can be constitutively activated by loss of function mutations to a number of
768 negative regulators including *IRA1*, *PDE2*, and *GPB2*. Mutations in these genes might be
769 redundant in the sense that they influence the same fitness-relevant phenotype in the evolution
770 condition, which in this case is likely flux through the Ras/PKA pathway. This type of
771 redundancy is commonly observed in laboratory evolutions, as evidenced by studies that
772 combine adaptive mutants to find they are no more adaptive than the most fit single mutant
773 (Tenaillon et al., 2012) and the observation that subsequent adaptive mutations tend to be in

774 other pathways (Aggeli et al., 2020). The major insight from our paper is that we show that
775 mutations with redundant effects on fitness in the evolution condition are not necessarily
776 identical because they may influence different latent phenotypes. This observation adds to a
777 long list of examples demonstrating that redundancies, such as gene duplications and
778 dominance, allow evolution the flexibility to generate diversity.
779

780 One disadvantage of our approach is that the phenotypic components that we infer from our
781 fitness measurements are abstract. They represent causal effects on fitness, rather than
782 measurable features of cells. For this reason, perhaps we should not refer to them as
783 phenotypes but rather “fitotypes” (a mash of the terms “fitness” and “phenotype”) that act much
784 like the causal traits in Fisher’s geometric model (Fisher, 1930; Harmand et al., 2017; Lourenço
785 et al., 2011; Martin and Lenormand, 2006; Tenaillon, 2014; Tenaillon et al., 2007) or a
786 selectional pleiotropy model (Paaby and Rockman, 2013). Despite this limitation, these
787 fitotypes have proven useful in allowing us to understand the consequences of adaptive
788 mutation. In addition to insights explained above, we also learned that adaptive mutants in the
789 same gene do not always affect the same fitotypes. For example, we found that *IRA1-*
790 *missense* mutations have varied and distinct effects from *IRA1-nonsense* mutations. Further, we
791 believe that identifying fitotypes will ultimately prove useful in identifying the phenotypic effects
792 of mutation. The fitotypes can serve as a scaffold onto which a large number of phenotypic
793 measurements can be mapped. Even though fitotypes are independent with respect to their
794 contribution of fitness, and contribute to fitness linearly, the mapping of commonly measured
795 features of cells (e.g. growth rate, the expression levels of growth supporting proteins like
796 ribosomes) onto fitotypes may not be entirely straightforward. Nonetheless, methods such as
797 Sparse Canonical Correlation Analysis (Suo et al., 2017) hold promise in such a mapping and
798 might help us relate traditional phenotypes to fitotypes.
799

800 An important question for future research is how common is the pattern we detected in this
801 study and whether it applies to other cases of adaptation in other systems. The method we
802 described is generic and can be applied to any system as long as the fitness of a substantial set
803 of mutants can be profiled for fitness across a moderately large set of environments. This is
804 becoming possible to do in many systems. It is also already clear that the notion that diverse
805 genetic changes can have redundant effects in one environment but distinct and consequential
806 effects in other environments is important to our understanding of adaptation in other settings,
807 including in the context of antibiotic resistance and cancer. Indeed in cancer, tumors even within
808 a particular type of cancer, say lung adenocarcinoma, tend to be extremely genetically diverse
809 even if considering only driver mutations (The Cancer Genome Atlas Research Network, 2014) .
810 The driver mutations do fall into a smaller number of key driver genes and even fewer
811 pathways. While this apparent redundancy might suggest that all tumors are in fact functionally
812 similar in that they solve a small set of functional challenges (Hanahan and Weinberg, 2011,
813 2000), the notion of latent diversity we propose here suggests that the specific paths taken by
814 the tumors early might matter once the tumor encounters a new challenge such as when they
815 are treated by a cancer therapy. Substantial heterogeneity of tumor response to therapy is
816 consistent with this notion.
817

818 Despite the accumulation of large amounts of genomic and phenomic data, integrating this
819 information to identify the phenotypic consequences of mutation that are ultimately responsible
820 for fitness remains incredibly challenging. Our approach allows us to create an abstract
821 representation of the causal effects of genetic mutation and their changing contribution to fitness
822 across environments. This top-down view provides an opening to solving this problem, and
823 combining these approaches with phenotypic measurements will allow us to answer age-old
824 questions about the structure of biological systems and adaptation in a conceptually new and
825 technically powerful and high-throughput way.
826

827 **ACKNOWLEDGMENTS**

828 The authors thank Sandeep Venkataram for the BarcodeCounter2 script; Yuping Li, Monica
829 Sanchez, Tuya Yokoyama, Chris McFarland, and Dimitra Aggeli for technical assistance; Atish
830 Agarwala, Marc Salit, Sasha Levy, Gavin Sherlock, and all members of the Petrov Lab for
831 helpful comments and discussions. We are grateful to the twitter community that followed
832 #1BigBatch and provided us with very helpful feedback. We are grateful to Enrico Coen for the
833 suggestion of the term “fitnotype”. Some of the computing for this project was performed on the
834 Sherlock cluster. We would like to thank Stanford University and the Stanford Research
835 Computing Center for providing computational resources and support that contributed to these
836 research results. This work was supported by National Institutes of Health grant R35GM118165
837 (to DAP) and National Institutes of Health grant R35GM133674 (to KGS).

1077 **METHODS**1078
1079 **LEAD CONTACT AND MATERIALS AVAILABILITY**

1080 Further information and requests for resources and reagents should be directed to and will be
1081 fulfilled by the Lead Contact, Dmitri Petrov (dpetrov@stanford.edu).
1082

1083 **EXPERIMENTAL MODEL AND SUBJECT DETAILS**

1084 The yeast strains used in this study can be grown and maintained using standard methods (e.g.
1085 YPD media in test tubes, glycerol stocks for long term storage at -80°C), but should be
1086 propagated in the appropriate selection environment (a glucose-limited minimal media - M3
1087 medium for the evolution condition) for comparable fitness and phenotypic measurements. All of
1088 the strains we study are of genetic background MAT α , ura3Δ0, ybr209w::Gal-Cre-KanMX-
1089 1/2URA3-loxP-Barcode-1/2URA3-HygMX-lox66/71.

1090 Experiments were performed with barcoded mutants isolated from a previous evolution
1091 experiment (Levy et al., 2015). To measure their fitness, these mutants were competed against
1092 a constructed reference strain with a restriction site in the barcode region (Venkataram et al.,
1093 2016).

1094 Since we utilize some data from previous experiments (Li et al., 2018; Venkataram et al., 2016),
1095 the collection of adaptive barcoded mutants that we studied differed slightly across
1096 environments. These differences can be thought of as another parameter that varies across the
1097 environments (e.g. in addition to glucose or salt concentration). In some experiments, we used a
1098 collection containing 4,800 adaptive mutants that do not necessarily start at equal frequency
1099 (Venkataram et al., 2016). In other experiments, we used a collection containing a subset of 500
1100 of these mutants where each one starts at equal frequency (Li et al., 2018; Venkataram et al.,
1101 2016). Though the smallest collection of mutants we study comprises 500 strains, our work
1102 focuses on 292 of these (Table S1). We focus on strains for which we obtained fitness
1103 measurements in 45 environments and for which mutations conferring fitness advantages have
1104 been previously identified, either by whole genome sequencing or using a drug to test ploidy (Li
1105 et al., 2018; Venkataram et al., 2016).

1106 In a few experiments, we spiked in re-barcoded mutants and additional neutral lineages as
1107 internal controls. Since re-barcoded mutants are identical, except for the barcode, these teach
1108 us about the precision with which we can measure a mutant's fitness. Specifically, we spiked in
1109 ten re-barcoded *IRA1-nonsense* mutants (each with a frameshift insertion AT to ATT mutation at
1110 bp 4090) and ten *IRA1-missense* mutants (each with a G to T mutation at bp 3776). Neutral
1111 lineages teach us about the behavior of the unmutated reference strain, which we must infer
1112 because it's barcode is eliminated from the experiment before sequencing. The spiked in
1113 neutrals include ten barcoded lineages from the original evolution experiment (Levy et al., 2015)
1114 for which whole genome sequencing did not reveal any mutations (Venkataram et al., 2016) and
1115 previous fitness measurements did not reveal any deviation from the reference (Li et al., 2018;
1116 Venkataram et al., 2016).
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127

1179 For each of these 45 conditions but three, we include between two and four replicates that were
1180 performed simultaneously (Table S2) such that overall we performed a total of 109 fitness
1181 measurements on our collection of adaptive mutants. Our replicate structure is nested in that
1182 some of our 45 conditions represent replicate experiments that we performed at different times.
1183 Variation across experiments performed at different times is often referred to as "batch effects"
1184 and likely reflects environmental variability that we were unable to control (e.g. slight fluctuations
1185 in incubation temperature due to limits on the precision of the instrument). In particular, we re-
1186 measured the fitness of the adaptive mutants in the EC on 9 different occasions, each time
1187 including 3 or more replicates. We refer to these 9 experiments as 'EC batches' in the main text.
1188 However, every set of experiments that was performed at the same time constitutes a separate
1189 "batch". There were slight differences across batches in the way we prepared barcodes for
1190 sequencing, which we detail in the relevant Methods sections. This variation across batches can
1191 be thought of as another parameter that varies across the 45 conditions (in addition to glucose
1192 or salt concentration). We report which experiments were performed in the same batch in Table
1193 S2.

1194

1195 Some conditions, including some Fluconazole conditions and Geldanamycin conditions, have
1196 unexpected orderings in the strength of perturbation (i.e. the smaller drug concentration shows
1197 a larger difference in fitness or similar concentrations seem to have different effects).
1198 Regardless of whether these observations reflect technical problems (e.g. degradation or poor
1199 solubility of the drug), we include these conditions because we use the effect of the realized
1200 perturbation on fitness to build low-dimensional phenotypic models. In other words, the identity
1201 of the perturbation does not matter in this study.

1202

1203

1204 **DNA Extraction of each sample**

1205 After a growth competition is complete, we extracted DNA from frozen samples following either
1206 a protocol described previously (for batches 1 – 6 and 10) (Venkataram et al., 2016) or a
1207 modified protocol that improves the ease and yield of extraction. Our modified protocol is as
1208 follows. For each sample, a single tube of the three that were frozen for each sample (see
1209 *Conducting the barcoded fitness measurements*) was removed from the freezer and thawed at
1210 room temperature. We extracted DNA from that sample using the following modification of the
1211 Lucigen MasterPure yeast DNA purification kit (#MPY80200). We transferred the thawed cells
1212 into a 15mL conical and centrifuge for 3 min at 4000 RPM. After discarding the supernatant, the
1213 pellet was then resuspended with 1.8 mL of the MasterPure lysis buffer, and 0.5 mm glass
1214 beads were added to help with disruption of the yeast cell wall. The mix of pellet, lysis buffer,
1215 and beads was then vortexed for 10 seconds and incubated for 45 minutes at 65°C, with
1216 periodic vortexing. The solution was then put on ice for 5 min and then 900 μ L of MPC Protein
1217 Reagent was mixed with the solution. We then separated protein and cell debris by
1218 centrifugation at 4000 RPM, transferring 1900 μ L of supernatant to a 2 mL centrifuge tube. We
1219 further separated remaining protein and cell debris by centrifuging at 13200 RPM for 5 min. The
1220 supernatant was then divided into two 2mL centrifuge tubes, with 925 μ L of the supernatant into
1221 each. Next, we added 1000 μ L of isopropanol to each tube, mixed by inversion, centrifuged at
1222 13200 RPM for 5 min, and discarded the supernatant. The pellet, containing the DNA was then
1223 resuspended in 250 μ L of Elution Buffer and 10 μ L of 5 ng/ μ L RNAase A was added. This was
1224 either left at room temperature overnight or incubated at 60°C for 15 min. Next the two tubes per
1225 sample were combined into a single tube and 1500 μ L of ethanol was added. This was then
1226 mixed by inversion, and strands of precipitating DNA appeared. This was centrifuged at 13200
1227 RPM for 2 min, and the supernatant was discarded. We again precipitated the DNA by
1228 resuspending with 750 μ L of ethanol, and collected the DNA by centrifuging 13200 RPM for 2
1229 min. The supernatant was discarded, and the tubes were left to air dry. Finally, we resuspended
1230 the pellet in Elution Buffer to a final concentration of 50 ng/ μ L for later use in PCR reactions
1231 (approximately 3600 ng of DNA were used for the PCR reactions).

1253 For the first step of PCR, we performed 8 reactions per sample to offset the effects of PCR
1254 jackpotting within each reaction. For each set of 8 reactions, we used the master mix:
1255 - 200 μ L OneTaq Hot Start 2X Master Mix with Standard Buffer (NEB M0484L)
1256 - 8 μ L 10uM Forward primer
1257 - 8 μ L 10uM Reverse primer
1258 - 72 μ L sample genomic DNA (diluted to 50ng/ μ L or all of sample if between 25-50ng/ μ L)
1259 - 16 μ L 50mM MgCl2
1260 - 96 μ L Nuclease Free Water (Fisher Scientific #AM9937)

1261
1262 We then aliquoted 50 μ L of the master mix into each of 8 PCR tubes, and ran on the
1263 thermocycler with the following cycle:

- 1264 1. 94°C for 10 min
- 1265 2. 94°C for 3 min
- 1266 3. 55°C for 1 min
- 1267 4. 68°C for 1 min
- 1268 5. Repeat steps 2-4 2x (for a total of 3 cycles)
- 1269 6. 68°C for 1 min
- 1270 7. Hold at 4°C

1271
1272 We then added 100 μ L of binding buffer from the ThermoScientific GeneJET Gel Extraction Kit
1273 (#K0692) to each PCR reaction, and performed a standard PCR purification protocol in one
1274 column per sample. In the final step, we eluted into 80 μ L of elution buffer.

1275
1276 For the second step of PCR, we use standard Nextera XT Index v2 primers (Illumina #FC-131-
1277 2004) to further label samples representing different conditions and timepoints with unique
1278 identifiers that allow for multiplexing on the same sequencing lane. We uniquely dual-indexed
1279 each sample using our nested scheme (see *Mitigating the effects of index hopping* section for
1280 details). We performed 3 reactions of the second step PCR per sample, using the master mix:

- 1281 - 1.5 μ L Q5 Polymerase (NEB #M0491L)
- 1282 - 30 μ L Q5 Buffer (NEB #M0491L)
- 1283 - 3 μ L 10mM dNTP (Fisher Scientific #PR-U1515)
- 1284 - 6.25 μ L i7 Nextera XT Primer ("N" primer)
- 1285 - 6.25 μ L i5 Nextera XT Primer ("S" primer)
- 1286 - 78 μ L purified step 1 PCR product
- 1287 - 25 μ L Nuclease Free Water (Fisher Scientific #AM9937)

1288
1289 This master mix was then divided into 3 PCR tubes per reaction, and run with the following
1290 protocol on a thermocycler:

- 1291 1. 98°C for 30 sec
- 1292 2. 98°C for 10 sec
- 1293 3. 62°C for 20 sec
- 1294 4. 72°C for 30 sec
- 1295 5. Repeat steps 2-4 at least 21 times and at most 27 times (for a total of 22 to 28 cycles)
- 1296 6. 72°C for 3 min
- 1297 7. Hold at 4°C

1298
1299
1300 We then added 100 μ L of binding buffer from the ThermoScientific GeneJET Gel Extraction Kit
1301 and purified the PCR product, eluting into 43 μ L. We found that increasing the number of cycles
1302 in the second step PCR beyond 21 did not seem to improve the amount of DNA recovered after

1303 gel extraction. For some samples, we experimented with a touch down procedure for the
1304 second step PCR where we started with a hotter annealing temperature and slowly decreased it
1305 over the course of 27 cycles. This also did not seem to increase the yield of DNA recovered
1306 from the PCR.

1307

1308 Removal of the Reference Strain via Digestion and Gel Purification

1309 To avoid the vast majority of our sequencing reads mapping only to the reference strain (and
1310 thus not being informative to relative fitness of the mutants), we use restriction digest to cut the
1311 ApaLI restriction site in the middle of the reference strain's barcode region. We mixed 43 μ L of
1312 the second step PCR product with 2 μ L of ApaLI (NEB #R0507L) and 5 μ L of 10X Cutsmart and
1313 incubated at 37°C for at least 2 hours (up to overnight). After digestion, we conducted size
1314 selection by running the digested sample on a gel, removing all product less than 300bp, and
1315 isolating the DNA using a standard ThermoScientific GeneJET Gel Extraction protocol. Our
1316 expected product is 350bp. We did not remove longer sequences via gel extraction because of
1317 the possibility that some barcode sequences may selectively form complexes with themselves
1318 or other barcodes.

1319
1320 Note that for some samples, we also digested the reference strain before PCR, in addition to
1321 after PCR, to decrease the amount of reference strain barcode. For these samples, we mixed
1322 80 μ L of genomic DNA (at concentration 50ng/ μ L) with 10 μ L of 10X Cutsmart and 2 μ L of ApaLI
1323 and incubated 37°C for at least 2 hours (up to overnight). This product was then used as the
1324 template for PCR step 1 (with appropriate water volume adjustments to ensure 50 μ L reactions).

1325

1326 Sample pooling and Amplicon Sequencing

1327 We used the Qubit High Sensitivity (ThermoFisher #Q32854) method to quantify the
1328 concentration of the final product for each sample, then pooled samples with different dual
1329 indices in equal frequency for sequencing. Our samples were then sent to either Novogene
1330 (<https://en.novogene.com/>) or Admera Health (<https://www.admerahealth.com/>) for quality
1331 control (qPCR and either Bioanalyzer or TapeStation) and sequencing. We used 2x150 paired-
1332 end sequencing along with index sequencing reads on Illumina HiSeq machines using patterned
1333 flow cells (either HiSeq 4000 or HiSeq X). We also used Illumina Nextseq machines with
1334 unpatterned flow cells. We found that the former was more subject to index hopping errors,
1335 please see *Mitigating the effects of index hopping* for a discussion of how our dual indexing
1336 reduces effects of index hopping. All amplicon samples were sequenced with at least 20%
1337 genomic DNA spiked in (either whole genomes from an unrelated project or phi-X) to ensure
1338 adequate diversity on the flow cell.

1339

1340 Mitigating the effects of index hopping

1341 To reduce the effects of index hopping observed on Illumina patterned flow cell technology
1342 (including HiSeq 4000, HiSeq X, and Novaseq machines) (Illumina, 2017; Sinha et al., 2017),
1343 we devise a nested unique-dual-indexing approach. This approach uses a combination of inline
1344 indices attached during the first step of PCR, as well as Nextera indices attached during the
1345 second step of PCR. The latter indices are not part of the sequencing read (they are read in a
1346 separate Index Read). This process uniquely labels both ends of all DNA strands such that DNA
1347 strands from multiple samples can be multiplexed on the same flow cell. Had we only labeled
1348 one end of each DNA strand, index hopping could have caused us to incorrectly identify some
1349 reads as coming from the wrong sample.

1350
1351 One approach to label samples with unique-dual-indices is to use 96 forward primers, each of
1352 which is paired to one of 96 reverse primers, instead our nested approach allows us to uniquely
1353 dual-index samples with only 40 total primers (12 forward inline, 8 reverse inline, 12 Nextera i7,
1354 8 Nextera i5). Specifically, we can use combinations of the Nextera and inline primers. One way

1355 to think of this is that there are 96 possible ways to combine the forward inline and Nextera i5
1356 primers that are on the same side of the read, effectively creating 96 unique labels for that end
1357 of the read.

1358
1359 To reduce the effect of index hopping contamination on our results, we included only samples
1360 that were sequenced on non-patterned flow cell technology (HiSeq 2000 and 2500 for samples
1361 in batches 1-6, 10, NextSeq for samples in batch 9) or were sequenced on patterned flow cell
1362 technology (patterned flow cell HiSeq) with nested unique-dual indexing.
1363

1364 **Processing of Amplicon Sequencing Data**

1365 We processed the amplicon sequencing data by first using the index tags to de-multiplex reads
1366 representing different conditions and timepoints. Then, using Bowtie2 (Langmead and Salzberg,
1367 2012), we mapped reads to a known list of barcodes generated by Venkataram et al. (2016),
1368 removed PCR duplicates using the UMIs from the first-step primers, and counted the number of
1369 reads for each barcode in each sample. The source code for this step can be found at
1370 <https://github.com/sandeepvenkataram/BarcodeCounter2>. We processed all raw data for this
1371 study using this pipeline, including re-processing the raw sequencing files for data from previous
1372 studies (Li et al., 2018; Venkataram et al., 2016) so that all data was processed together using
1373 the most recent version of the code.
1374

1375 Several samples included technical replicates where the sample was split at various times in the
1376 process, including before DNA extraction, before PCR, and prior to sequencing. Read counts
1377 across these technical replicates were merged in order to calculate the best estimate of barcode
1378 frequencies. Counts were merged after appropriately accounting for PCR duplicates as
1379 identified from Unique Molecular Identifiers.
1380

1381 **QUANTIFICATION AND STATISTICAL ANALYSIS**

1382 **Fitness Estimate Inference**

1383 The amplicon sequencing data shows the relative frequency of each barcode in each time-point
1384 of every one of our 109 fitness measurement experiments. To estimate the fitness of each
1385 barcoded mutant in each experiment, we calculate how barcode frequencies change over time.
1386 We do this using previously described methods (Venkataram et al., 2016).
1387

1388 Briefly, we first calculated the log-frequency change of each barded adaptive mutant for each
1389 subsequent pair of time-points. This log-frequency change must be corrected by the mean
1390 fitness of the population, such that it represents the relative fitness of each mutant relative to the
1391 reference strain, which makes up the bulk of the population. Since we destroyed barcodes
1392 pertaining to the reference strain by digesting them, we infer how the mean fitness of the
1393 population changes at each time-point using barcoded lineages that are known to be neutral
1394 (see *Identification of neutral lineages*). Once we calculated the change in the relative fitness of
1395 each barcoded mutant across each pair of consecutive time-points, we took a weighted average
1396 across all pairs as our final estimate of each adaptive mutant's relative fitness for a given
1397 experiment. We weighted each pairwise fitness estimate using an uncertainty measure
1398 generated from a noise model (see *Noise model* section below).
1399

1400
1401 This results in 109 fitness measurements per each barcoded mutant, with some of the 45
1402 conditions having more representation than others due to having more replicates. In cases
1403 where we have replicates, we averaged the fitness values across the replicates, weighted by
1404 the measurement uncertainty, resulting in our final 45 fitness estimates per each adaptive
1405 mutant lineage.
1406

1407 We included only timepoints with at least 1,000 reads for which at least 400 mutants were
1408 detected to have at least 1 read. Furthermore, we required that fitness measurement
1409 experiments must include at least three timepoints to be included in our analysis.

1410

1411 **Identification of Neutral Lineages**

1412 Previous work using this fitness measurement method focused on a larger collection of 4800
1413 barcoded yeast lineages, where the vast majority of these lineages were neutral (Li et al., 2018;
1414 Venkataram et al., 2016). In order to increase the number of reads per adaptive lineage, we
1415 used a smaller pool of 500 lineages for most experiments. However, this prevents us from
1416 identifying neutral lineages as was done in previous studies, by rejecting outlier lineages with
1417 higher than typical fitness values. Instead, we used a set of 35 high-confidence neutral lineages
1418 to infer mean fitness (see *Experimental model and subject details*). These lineages showed no
1419 fitness differences from the neutral expectation in previous studies and were shown to possess
1420 no mutations detectable via whole genome sequencing. These high-confidence neutral lineages
1421 were present in all experiments, and were spiked into experiments from batch 9 to increase their
1422 frequency. We used these neutrals to perform the fitness inference in two steps. First, we
1423 inferred fitness using this collection of high-confidence neutrals to make a first pass at inferring
1424 the fitness values. Next, we included lineages with similar behavior to the high-confidence
1425 neutrals to improve our estimate of mean fitness.

1426

1427 **Noise model**

1428 To quantify the uncertainty for each fitness measurement, we used the noise model as outlined
1429 in Venkataram et al., 2016.

1430

1431 Briefly, this noise model accounts for the uncertainty coming from several sources of noise. The
1432 first type of noise scales with the number of reads for a given lineage. This noise stems from
1433 stochasticity in population dynamics (coming from the inherent stochasticity in growth and noise
1434 associated with dilution), from counting noise associated with a finite coverage, and technical
1435 noise from DNA extraction and PCR. We fit this noise by quantifying the variation in the
1436 frequency of neutral lineages (see *Identification of neutral lineages*). There is additional variation
1437 in fitness observed for high-frequency lineages between replicate experiments (here we refer to
1438 variation across replicates that were performed simultaneously, not variation across batches).
1439 We also accounted for this uncertainty following previous studies. Specifically, we fit an
1440 additional frequency-independent source of noise using between-replicate variation.

1441

1442 **Checks on noise model**

1443 Because our ability to count the phenotypes that matter to fitness hinges upon measurement
1444 error, we further assessed the accuracy of our noise model. We did so by using barcoded
1445 lineages that should have the same fitness because they are genetically identical. Since our
1446 fitness estimates are imperfect (*i.e.* they contain some noise), we estimated each of these
1447 lineages as having slightly different fitness. We then asked if the variation in fitness across
1448 these lineages is explained by our noise model, or if there is more variation than our noise
1449 model can account for. We did this explicitly by calculating, for each lineage, how far our fitness
1450 estimate is from the best guess for the true underlying fitness value (the group's mean) in units
1451 of the estimate's measurement precision. We then calculated the percent of lineages that are a
1452 given distance from the group's average to understand the accuracy of the model. For instance,
1453 if the noise model perfectly captures the uncertainty of each measurement, then 10% of the
1454 diploid lineages should have a difference from the weighted diploid mean in the 10th percentile,
1455 20% in the 20th percentile, etc. Because 188 of our 292 barcode mutants are diploids without
1456 additional mutations, diploids are an ideal group to use to assess the accuracy of the noise
1457 model. This procedure shows that, for the vast majority of replicates, the noise model is
1458 conservative. That is, diploid lineages tend to have less variation in fitness than expected by the

1459 noise model (Fig S1).

1460

1461 Classifying mutants by mutation type

1462 Some types of mutants are present more than others. For example, 188 of our 292 mutants are
1463 diploids and 30 mutants are in the IRA1 gene. If not properly accounted for, this imbalance can
1464 lead to some unfairness in predictions for our model. For example, if we use mostly diploid
1465 lineages to train our model, we will be very good at predicting the fitness of diploids but poor at
1466 predicting other types of mutants. This means that we must classify our mutants by mutation
1467 type in order to properly balance them. We classified mutants following previous work
1468 (Venkataram et al., 2016) that classified mutants as either diploids, or if haploid, by the gene
1469 possessing the putative causal mutation. Because previous work finds differences in fitness
1470 between missense and nonsense/frameshift/indel mutations in IRA1, here we classified these
1471 mutants into “missense” and “nonsense” classes, where mutants with frameshift and indel
1472 mutations were classified as “nonsense”. We also classified diploid mutants with additional
1473 mutations in nutrient-response genes or chromosomal amplifications as separate groups.
1474 Additionally, we created a separate class for “high-fitness diploid” mutants that possess no
1475 additional detected mutations (other than being diploid) but have very high fitness in the EC. To
1476 be classified as a high-fitness diploid, a diploid mutant must have an average fitness across all 9
1477 EC batches that is greater than 2 standard deviations above the average of all diploids.

1478

1479 Calculation of Weighted Average Z Score

1480 To partition environments into subtle and strong perturbations of the EC, we relied on the 9
1481 experiments performed in the EC. Since each of these experiments was performed at a different
1482 time, variation in fitness across these experiments represents batch effects, and we therefore
1483 refer to these 9 experiments as “EC batches”. Environmental differences between batches are
1484 as subtle as possible, as they represent the limit of our ability to minimize environmental
1485 variation. Thus, variation in fitness across the EC batches serves as a natural benchmark for the
1486 strength of environmental perturbations. If the deviations in fitness caused by an environmental
1487 perturbation are substantially stronger than those observed across the EC batches, we call that
1488 perturbation “strong”.

1489
1490 More explicitly, to determine whether a given environmental perturbation is subtle or strong, we
1491 first quantified the typical variation in fitness for each mutant, across the EC batches:

$$1492 \sigma_i = \frac{1}{n_{batches}} \sum_j^{batches} |f_{ij} - \bar{f}_i|$$

1493 where σ_i^2 represents the variance in fitness across the EC batches for mutant i , and
1494 \bar{f}_i represents the average fitness of mutant i across the EC batches.

1495
1496 To ensure that each mutation type contributes equally to our classification of how different each
1497 environment is from the evolution condition, we weighed each mutant’s contribution to this
1498 difference. We did so based on the number of mutants with the same mutation type, such that
1499 the mutation-type-weighted average Z-score for a given environment j is given by:

$$1500 z_j = \sum_i^{mutants} \frac{|f_{ij} - \bar{f}_i|}{n_{type(i)} \sigma_i}$$

1501 where $n_{type(i)}$ represents the number of mutants that are the same mutation type as mutant i .

1502
1503 We then classified the environmental perturbations based on this Z-score. Sixteen environments
1504 provoked fitness differences resulting in a Z-score of less than two, and we classified these
1505 environmental perturbations as “subtle”. The remaining 20 environments had Z-scores greater

than 2, which we classified as “strong” environmental perturbations.

Model of phenotypes that contribute to fitness

In order to count the phenotypes that affect fitness in our collection of mutants, we explored a low-dimensional phenotypic model. We explicitly used a model of fitness-relevant phenotypes such that each mutant is represented as having a fixed effect on each phenotype, represented by a vector of k phenotypes, e.g. mutant i is represented by the vector $(p_{i1}, p_{i2}, p_{i3}, \dots, p_{ik})$. In addition, each environment is represented by a vector of phenotypic weights, representing the importance of each of the k phenotypes to fitness in that environment, e.g. environment j represented by the column vector $(e_{1j}, e_{2j}, e_{3j}, \dots, e_{kj})$. The fitness effect of mutant i in a given environment j is the linear combination of that mutant’s phenotypes, each weighted by its importance in environment j :

$$f_{ij} = p_{i1}e_{1j} + p_{i2}e_{2j} + p_{i3}e_{3j} + \dots + p_{ik}e_{kj}$$

Our fitness measurements reflect mutant fitness relative to a reference strain, therefore, our model places the reference strain (which has fitness 0 by definition) at the origin of this multi-dimensional space. Our model only includes phenotypes that differ between the reference strain and least one mutant. This is sensible given that our reference strain is a modified version of the ancestor of all of these mutant lineages. Thus, if there exists a phenotype that contributes to fitness, but none of the adaptive mutants altered that phenotype, our model will not detect it. More explicitly, a phenotype that contributes to fitness would have a non-zero value of e , but if no mutant alters that phenotype from the reference, all mutants would have a zero value of p for that phenotype. Thus, the non-zero value of e would always be multiplied by a zero value for p and this phenotypic dimension would not be represented in our model. This is not to say that if only a single mutant of the 292 alters a particular phenotype we would include it as a phenotypic dimension. Our power to add dimensions to our model is limited by measurement noise. We only include dimensions that capture more variation in fitness than do dimensions that capture measurement noise (see *Estimating the detection threshold using measurement error*).

Similarly, because we measure fitness, and not phenotype, our model is blind to any phenotypic effect that does not contribute to fitness in at least one of the 45 environments we studied. If a mutant has large phenotypic effects, but they do not cause that mutant’s fitness to differ from the reference strain in any of these 45 environments, this phenotypic effect will not be represented in our low-dimensional phenotypic model. More explicitly, mutants may have non-zero phenotypic effects p , but if these do not influence their fitness in any environment we study, e will be zero for all 45 environments. Thus, p times e will also be zero and we will not include this phenotypic dimension in our model.

Importantly, the phenotypic dimensions that we infer from our fitness measurements are abstract entities. They represent causal effects on fitness, rather than measurable features of cells. For this reason, they might be called “fitotypes” (a mash of the terms “fitness” and “phenotype”). Even though the fitotypes are independent with respect to their contribution of fitness, and contribute to fitness linearly, the mapping of commonly measured features of cells (e.g. growth rate, the expression levels of growth supporting proteins like ribosomes) onto fitotypes may be more complicated. For instance, a commonly measured cellular feature that has a complicated nonlinear mapping to fitness could be detected as many, linearly-contributing fitotypes. This is another reason that our phenotypic dimensions are not necessarily comparable to what people traditionally think of as a “phenotype”.

Using Singular Value Decomposition to decompose the fitness matrix

Our goal is to use fitness measurements to learn about the phenotypic effects of mutations as well as the contribution of these phenotypes to fitness in different environments. We conducted fitness measurements for 292 mutants in each of 45 environments and organized these data

1558 into a fitness matrix, F , where every row corresponds to a mutant, every column corresponds to
1559 an environment, and every entry is a fitness measurement. Because our model (see *Model of*
1560 *phenotypes that contribute to fitness*) represents fitness in a given environment as the sum of
1561 multiple phenotypes, each scaled by their contribution to fitness in that environment, we can use
1562 Singular Value Decomposition (SVD) to decompose the fitness matrix F as:

$$P\Sigma E^T = F$$

1564 The left hand side of this equation consists of three matrices: P , which represents the positions
1565 of the mutants in our low-dimensional model of phenotypic space, E^T , which represents the
1566 contribution of a phenotype to fitness in a given environment, and Σ , a diagonal matrix
1567 representing the singular values of the fitness matrix F . Though the singular values are
1568 informative in this separation of three matrices, particularly for the amount of variation captured
1569 by each of the inferred components, we can also think of this as a decomposition into two
1570 matrices, where we fold the singular values into either the mutant phenotypes or the
1571 environment weights, as described in the main text. Either way, this decomposition captures the
1572 data represented in the fitness matrix F , including measurement error as well as the underlying
1573 biological signals.

1574
1575 Importantly, the dimensions in the phenotypic model we built using SVD are detected in the
1576 order of their explanatory power. Moreover, the first dimension is the best, linear 1-component
1577 model that explains the data (if evaluated by mean squared error). This is true for any set of the
1578 first k components. This means, for example, that the model with the first eight components is
1579 the best possible 8-component linear model for the observed data (Eckart and Young, 1936).

1580
1581 One issue in this type of analysis is that adding more components always improves the
1582 explanatory power of the model, even when those components capture variation that is primarily
1583 due to measurement noise. This type of overfitting problem is common in statistics, and several
1584 methods have been devised to select the appropriate number of components to include. We use
1585 two such methods here.

1586
1587 **Estimating the detection threshold using measurement error**
1588 One method to select the appropriate number of components to include in the model and
1589 prevent overfitting (*i.e.* prevent fitting a component that primarily represents noise) is to use
1590 measurement error as a type of control. This is only possible if the amount of measurement
1591 error is known. We estimated the amount of noise in our fitness measurements using a
1592 previously described noise model (see *Noise Model*) (Venkataram et al., 2016). Since this noise
1593 model includes counting noise, every fitness measurement may have a different amount of
1594 noise. For example, mutants present at low frequency will be subject to more stochasticity
1595 resulting from counting noise. We used this noise model to simulate fitness tables (F) where
1596 mutant fitnesses vary exclusively due to measurement noise. We simulated 1000 noise-only
1597 matrices, where each entry is pulled from a normal distribution centered at zero and with
1598 variance equal to the estimated measurement noise of the corresponding entry in the true
1599 fitness matrix F . We then applied SVD to each noise-only matrix, which gave us a set of singular
1600 values generated only by noise. From many such simulations, we took the average size of the
1601 largest component, which reveals how much variation can be explained by a component that
1602 captures only noise. We found that the largest noise-components are of the size that they would
1603 capture 0.07% of variation in our true fitness matrix. Thus, we set this as our limit of detection.
1604 In other words, in order for us to include 8 components in our low-dimensional model, all of
1605 them must explain more than 0.07% of the variation in fitness. This approach is analogous to
1606 identifying a threshold when measurement noise is known but not identical for all entries in the
1607 matrix (Josse and Sardy, 2014).

1608
1609

1610 **Estimating detection threshold using bi-cross-validation**

1611 Another method for identifying the appropriate number of components is to use their predictive
 1612 power. This method relies on the intuition that measurement error is uncorrelated across
 1613 different mutants and different environments. Therefore, a component that represents
 1614 measurement error should not contain information that can help predict the fitnesses of these
 1615 mutants in new environments. It should also not contain information that can help predict the
 1616 fitness of unstudied mutants. We used a bi-cross-validation scheme of the SVD devised by
 1617 Owen and Perry (2009) which divides the mutants and environments into distinct groups of
 1618 training and testing sets. This subsequently divided our matrix of fitness measurements into 4
 1619 submatrices: the fitness of the training mutants in the training environments (D), the fitness of
 1620 the training mutants in the testing environments (C), the fitness of the testing mutants in the
 1621 training environments (B), and the fitness of the testing mutants in the testing environments (A).
 1622

$$1623 F = \begin{pmatrix} A & \substack{\text{Test Mutants} \\ \text{Test Environments}} & B & \substack{\text{Test Mutants} \\ \text{Train Environments}} \\ C & \substack{\text{Train Mutants} \\ \text{Test Environments}} & D & \substack{\text{Train Mutants} \\ \text{Train Environments}} \end{pmatrix}$$

1624 We carried out SVD on the training data (submatrix D), which returned a set of singular values
 1625 and corresponding components that captured the fitness data in D . We then used these
 1626 components to predict the fitness of the testing mutants in the testing environments (submatrix
 1627 A). First, we tried to predict these fitness values by only using the first component. That is, we
 1628 fixed this first component and the first singular value for the training mutants. We then found the
 1629 best first component for the testing environments based on the fitness values of the training
 1630 mutants in these environments (i.e. using the information in submatrix D), given the constraint
 1631 that the training mutants can only be represented by the one component. We then conducted an
 1632 analogous procedure to find the first component of the testing mutants by fixing the first
 1633 component of the training environments by using the information in submatrix B . Then, we tried
 1634 to predict the fitness of the testing mutants in the testing environments using the first component
 1635 independently fit for each. We subsequently repeated this procedure, giving the testing mutants
 1636 access to more of the training components each time. If the components detected by the
 1637 training components represent biological signal, then this should improve the ability to predict
 1638 the fitness of the testing mutants in the testing environments. However, once the components
 1639 primarily represent measurement error, their inclusion should harm predictive power. Therefore,
 1640 we use the number of components with the best ability to predict the held-out data (submatrix A)
 1641 as the number of components that represent biological signal in our data.
 1642

1643 For computational efficiency, we explicitly used the formulation proposed by Owen and Perry
 1644 (2009) for the prediction of the held-out submatrix A :

$$1645 \hat{A} = B(\hat{D}^{(k)})^+ C$$

1646 where $(\hat{D}^{(k)})^+$ denotes the Moore-Penrose inverse of the rank k approximation of sub-matrix D .
 1647 This prediction is equivalent to the procedure outlined above, provided that least-squares
 1648 regression is used to identify the components of the testing mutants and testing conditions,
 1649 conditional upon the training components (Owen and Perry, 2009).
 1650

1651 We divided our mutants into fixed training and testing sets (see *Division of Mutants into Training*
 1652 and *Testing Sets*) and used these sets throughout our study. As for training versus testing
 1653 environments, these changed depending on our goal. For validating the number of components
 1654 to include in our phenotypic model, we held out each of the 25 subtle environmental
 1655 perturbations, using it as the testing environment and the other 24 for training. For making
 1656 predictions of the fitness of the testing mutants in the strong environmental perturbations, we
 1657 used all 25 subtle environmental perturbations as the training set, though we also show how
 1658

1659 these predictions vary when each of the 25 subtle environmental perturbations is held out from
1660 the training set.

1662 Division of Mutants into Training and Testing Sets

1663 In order to perform bi-cross-validation on our data, we need to divide our data into training and
1664 testing sets. Because some mutation types, in particular diploids and Ras/PKA mutants, are
1665 present more than others in our collection of mutants, we sampled the training set such that
1666 each mutation type is represented roughly equally (see *Classifying mutants by mutation type*).
1667 Specifically, we designated half of each mutation type, with a maximum of 20 representatives of
1668 each type, as belonging to the training set. The remaining mutants comprise the test set. For
1669 example, there are 188 diploids included in the 292 adaptive mutants. We included 20 in the
1670 training set and 168 in the test set. There are 20 *IRA1-nonsense* mutants included in the 292,
1671 and we included 10 in the training and 10 in the test set. Additionally, genes that are
1672 represented only once in the set of mutations are placed in the test set. This results in a training
1673 set of 60 mutants and a testing set of 232 mutants (see Table S1).

1675 Clustering mutants in phenotype space

1676 After inferring the low-dimensional model of phenotype space using SVD, we used Uniform
1677 Manifold Approximation and Projection (UMAP) to visualize how the mutants cluster in that
1678 space. For this analysis, we used the 8-component phenotypic model that we built from the 60
1679 training mutants and the 25 subtle perturbations. We did this to avoid the model being
1680 dominated by variation in very common mutations, specifically the diploids, which make up
1681 188/292 of our adaptive mutants. We added more mutants in the visualization by finding the
1682 location of each of the testing mutants (except diploids) by least sum of squares optimization.
1683 To do so we fixed the coordinates for the 25 environments and found the coordinates for each
1684 mutant that best estimated its fitness in all environments. To further avoid our visualization
1685 being dominated by the diploids, we included only the diploids present in the training set in our
1686 visualization. For UMAP, we specified that 20 neighbors are used.

1687 Though UMAP tends to preserve both local and global structure (McInnes et al., 2018) it is not
1688 necessarily representative of the distance between objects in high-dimensional space. Thus, to
1689 quantify more precisely the clustering by gene observed, we explicitly compared the median
1690 pairwise distance between these apparent clusters to 10000 randomly chosen sets of the same
1691 size and calculated empirical p-values. Because there are many diploids such that they will be
1692 the most prevalent type of mutant drawn in these randomly chosen sets, we only drew from
1693 strains that have other mutations besides or in addition to diploidy.

1696 Calculation of Weighted Coefficient of Determination

1697 Because mutants are present in unequal numbers in the test set, standard measures of
1698 variance explained are likely to be representative of our ability to predict mutants that have
1699 many barcoded lineages present in the data, for instance diploid and IRA1-nonsense mutations.
1700 These measures would be less representative of mutants with few lineages present, i.e.
1701 TOR/Sch9 pathway mutants. Thus, we use a measure of predictability (\tilde{R}^2) that weights the
1702 contribution of each mutant to overall variance explained based on the number of lineages that
1703 share its mutation type (diploids, IRA1 nonsense, IRA1 missense, GPB2, etc.). This effectively
1704 measures our ability to predict the fitness of each mutation type, rather than each mutant. For
1705 overall predictive power across all mutants and conditions, we used the measure:

$$1706 \tilde{R}^2 = 1 - \frac{\sum_i^{mutants} \sum_j^{conditions} \frac{1}{n_{type(i)}} (f_{ij} - \hat{f}_{ij})^2}{\sum_i^{mutants} \sum_j^{conditions} \frac{1}{n_{type(i)}} (f_{ij} - \bar{f})^2}$$

1707 where \bar{f} denotes the average fitness for all evaluated mutants and evaluated conditions.

1708 We used a similar measure to quantify the ability to predict fitness for each environment j . This
1709 is given by:

$$1710 \tilde{R}_j^2 = 1 - \frac{\sum_i^{\text{mutants}} \frac{1}{n_{\text{type}(i)}} (f_{ij} - \hat{f}_{ij})^2}{\sum_i^{\text{mutants}} \frac{1}{n_{\text{type}(i)}} (f_{ij} - \bar{f}_j)^2}$$

1711 where \bar{f}_j denotes the average fitness across all evaluated mutants in condition j .

1712
1713 Note that this measure explicitly compares a model's fitness prediction in each environment to
1714 predictions made using the average fitness in that environment, such that if the model's fitness
1715 prediction is the same as the average fitness, \tilde{R}^2 is zero. It is possible that a given model's
1716 fitness prediction is worse than that of the average fitness in that environment, resulting in
1717 negative values of \tilde{R}^2 . In our work, negative \tilde{R}^2 values occur for the 1-component model when
1718 predicting the fitness of mutants in some of the strong environmental perturbations. In particular,
1719 this occurs when fitness in that environment is uncorrelated with EC fitness, which is captured
1720 by the first component, such that the EC fitness is unable to make reasonable predictions of
1721 fitness in this environment.

1722
1723 Note that we observe qualitatively similar results to this measure when we use a standard
1724 variance explained measure and exclude diploids, which dominate the test set (see Fig S5).

1726 Calculating mutant-specific improvement

1727 It is possible that all 292 of our adaptive mutants each affect all 8 of the phenotypic components
1728 in our low-dimensional model, however, it is also possible that some mutants influence some
1729 phenotypes more strongly than others. In order to quantify how much a specific component
1730 lends to the ability to predict the fitness of each mutant in each environment, we need a metric
1731 to calculate the difference in predictive accuracy for the model with and without this component.
1732 Specifically, to assess the impact of the inclusion of the k th component, we compared the
1733 prediction accuracy of the k -component model to the model that includes the first $k-1$
1734 components.

1735
1736 Because fitness estimates vary in their reliability due to finite coverage and other sources (see
1737 *Noise model* section), we should factor this uncertainty in our measure of prediction
1738 improvement. For example, a small improvement in prediction accuracy for a very uncertain
1739 fitness estimate is less meaningful than the same improvement in prediction accuracy for a
1740 fitness estimate that we are quite confident in. Thus, we scale the difference in prediction
1741 accuracy by the amount of uncertainty in the underlying fitness estimate.

1742
1743 This gives us the measure of improvement in the estimate of the fitness of mutant i in condition j
1744 due to the inclusion of the n th component as:

$$1745 I_{ij}^k = \frac{(\hat{f}_{ij}^{k-1} - f_{ij}) - (\hat{f}_{ij}^k - f_{ij})}{\epsilon_{ij}}$$

1746 where \hat{f}_{ij}^k and \hat{f}_{ij}^{k-1} represent the estimate of the fitness of mutant i in condition j for the model
1747 with k and $k-1$ components, respectively. f_{ij} and ϵ_{ij} represent the measured fitness value and
1748 measurement uncertainty for the fitness of mutant i in condition j , respectively.

1749
1750
1751
1752
1753

1754 **DATA AND CODE AVAILABILITY**

1755

1756 **Data Resource**

1757 The raw Illumina sequencing data for the fitness measurement assays conducted in this study
1758 can be found under NIH BioProject: PRJNA641718. Sequencing data previously published in
1759 Venkataram et al., 2016 can be found under NIH BioProject: PRJNA310010. Sequencing data
1760 previously published in Li et al., 2018 can be found under NIH BioProject: PRJNA388215.

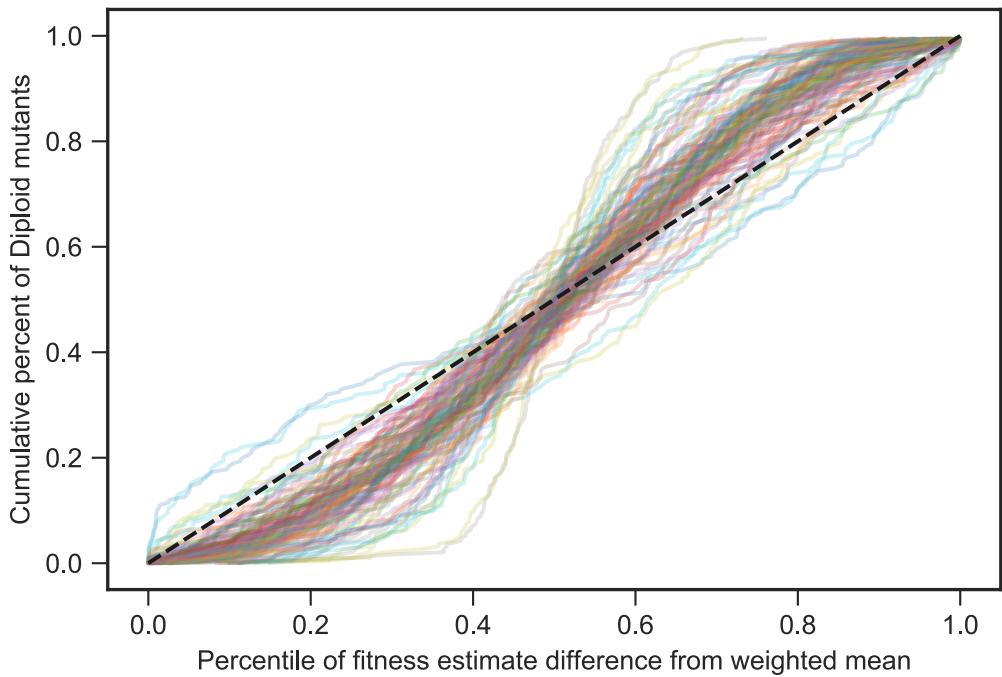
1761

1762 **Code**

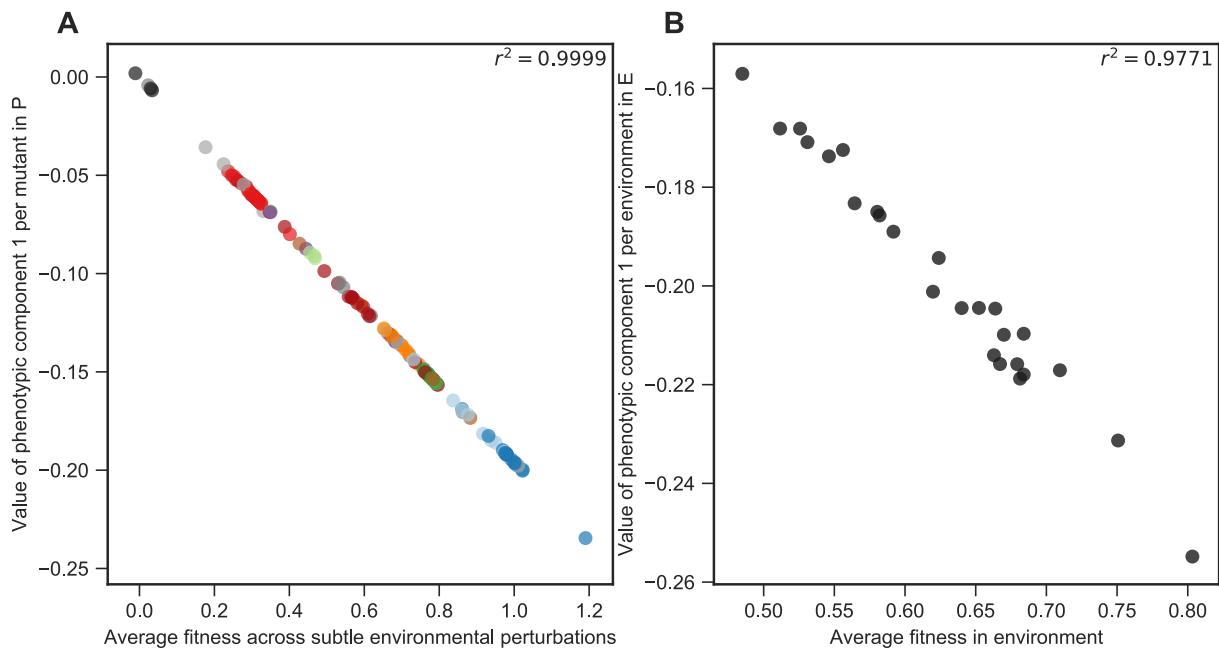
1763 The software repository for the barcode counting code can be found at
1764 <https://github.com/sandeepvenkataram/BarcodeCounter2>.

1765 The software repository for the fitness estimate inference can be found at
1766 <https://github.com/barcoding-bfa/fitness-assay-python>.

1767 The code for all downstream analysis, including figure generation can be found at
1768 <https://github.com/grantkinsler/1BigBatch>.

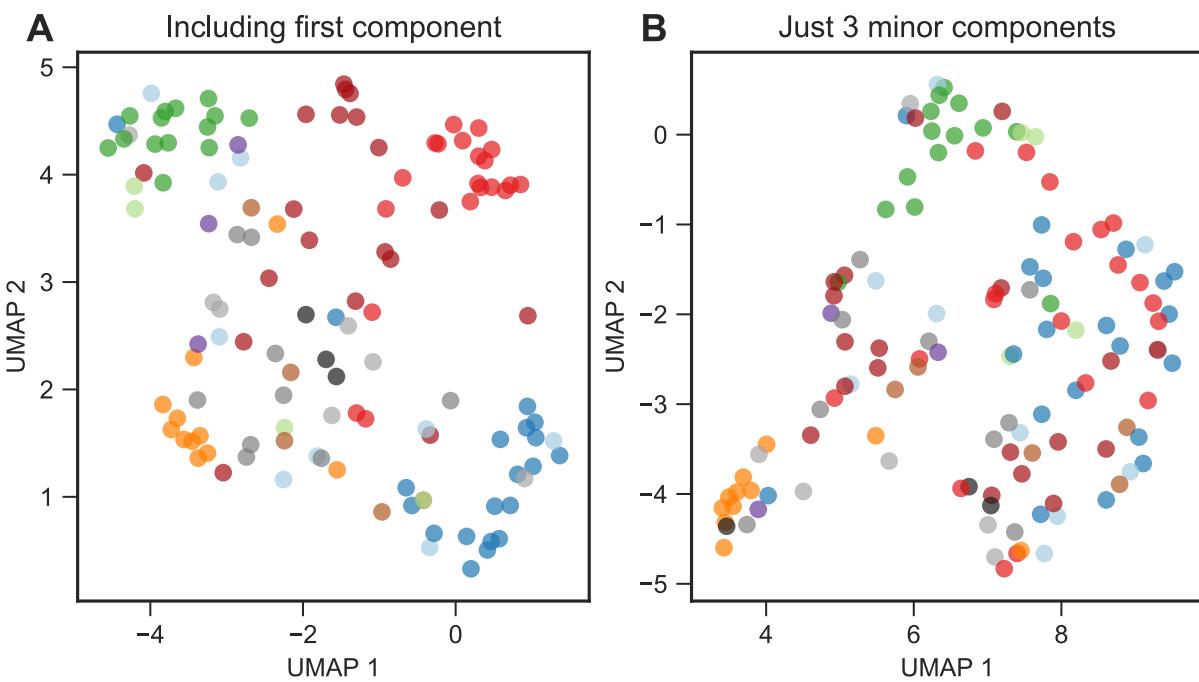


1770
 1771 **Fig S1. Noise model is a conservative measure of uncertainty.** Fitness differences among strains that
 1772 are genetically identical and have very similar fitness effects tell us about the amount of measurement
 1773 noise. Our strain collection includes 188 diploids that have similar fitnesses and possess no mutations
 1774 other than diploidy. For each diploid fitness estimate, we calculated the percentile of deviation from the
 1775 weighted average of all diploid fitness estimates in a particular environment. This is shown on the
 1776 horizontal axis. The vertical axis shows the cumulative percent of diploids with deviations listed on the
 1777 horizontal axis. If the noise model perfectly captures the uncertainty of each measurement, then it should
 1778 be represented by the black dashed line, as, for instance, 20% of the diploids should have a difference
 1779 from the mean in the 20th percentile. Each line represents a single experiment (we have 45 environments
 1780 each with several replicates for a total of 109 experiments, see Methods). For the vast majority of
 1781 experiments, the diploids are closer to the mean than predicted by our noise model, as indicated by each
 1782 line's sigmoidal shape. This indicates that the noise model is conservative.



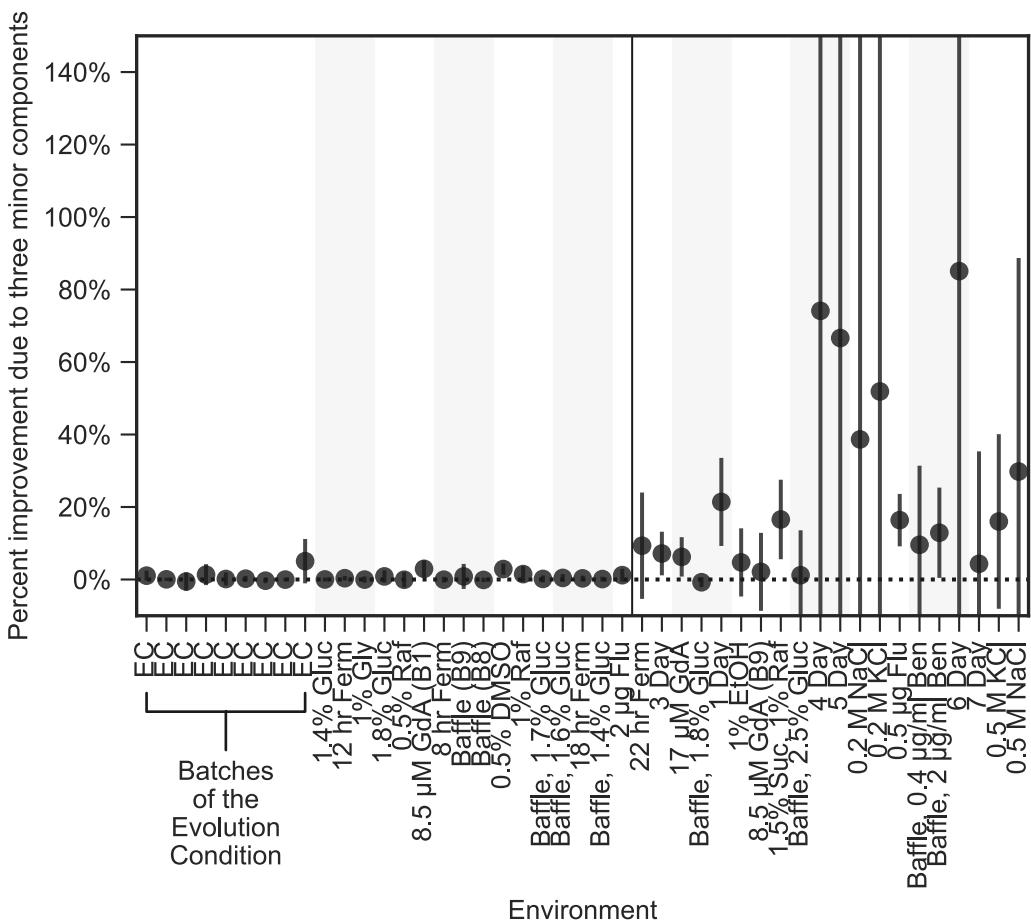
1783
1784
1785
1786
1787
1788
1789
1790

Figure S2. The first component represents the mean fitness of each mutant in the 25 subtle perturbations, as well as the mean impact of each perturbation on fitness. (A) The horizontal axis shows the average fitness of each mutant across all 25 environments that represent subtle perturbations. The vertical axis shows the value of the first phenotypic component for each mutant. Mutants are colored as in Figure 2. **(B)** The horizontal axis shows the average fitness of all 292 mutants in each environment. The vertical axis shows the value of the first phenotypic component in the environment weight space E.



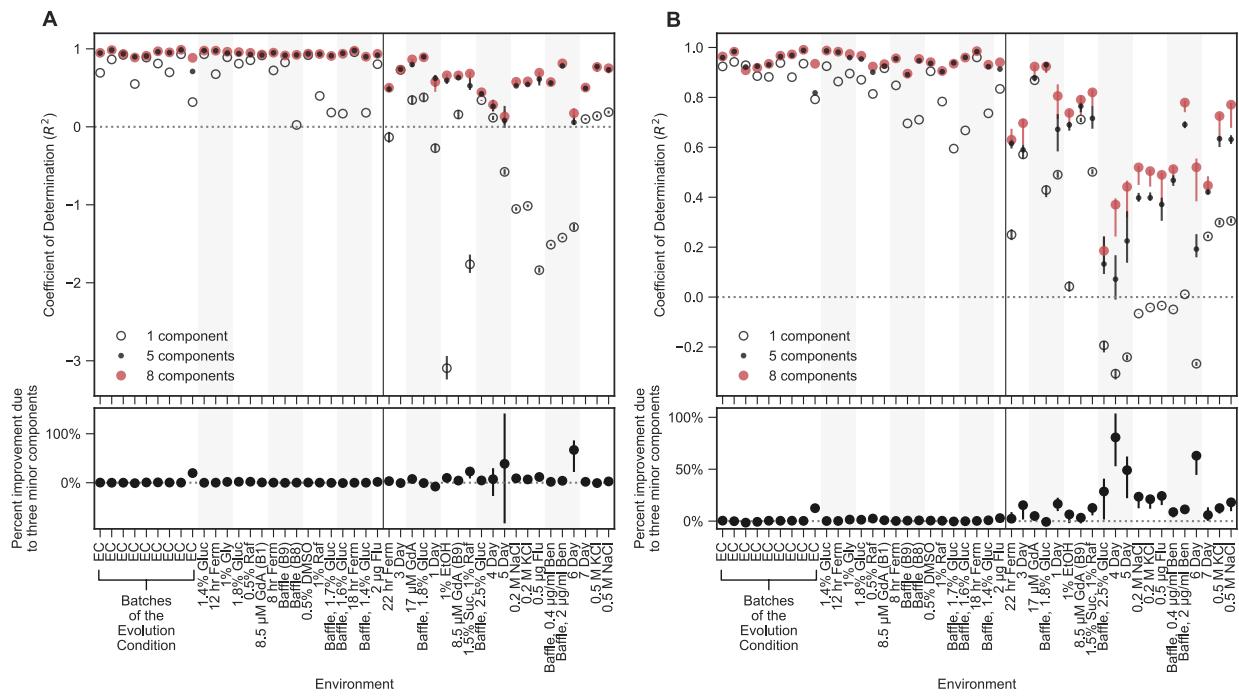
1791
1792
1793
1794
1795
1796
1797

Fig S3. Low-dimensional phenotypic models, and subsets of such models, cluster mutants by gene and mutation type. (A) UMAP clusters mutants visually by gene when using the full 8-component phenotype space. **(B)** UMAP also shows some clustering when using only the three components that explain the least variation in mutant fitness. Though the clustering is clear for *PDE2* and *GPB2*, it less clearly delineates *IRA1-nonsense* and diploid mutants. This suggests these mutants do not have substantial effects on these phenotypic components.



1798
1799
1800
1801
1802

Fig S4. Improved fitness predictions when including the three smallest phenotypic components is not specific to choice of training mutants. This plot is similar to the lower panel of figure 4A, except here, black dots indicate the average improvement across 100 choices of the training and test sets. Error bars indicate two standard deviations from the mean.



1803
1804
1805
1806
1807
1808

Fig S5. Prediction ability using unweighted coefficient of determination. These plots are similar to figure 4A except here the vertical axis displays prediction power using a standard, rather than a weighted, coefficient of determination measure. Because diploids dominate the number of mutants in the collection, there are large differences between panel A (which shows all mutants) and panel B (which omits diploids).

1809 **SUPPLEMENTARY TABLES**

Mutation Type	Total Number	Number in Training Set	Number in Test Set
Diploid	188	20	168
High-Fitness Diploid			
Diploid + Chr11Amp	3	1	2
Diploid + Chr12Amp	1	0	1
Diploid + IRA1	1	0	1
Diploid + IRA2	3	1	2
No known add'l mutations	11	5	6
IRA1			
IRA1 nonsense	20	10	10
IRA1 missense	9	4	5
IRA1 other	1	0	1
IRA2	8	4	4
GPB1	4	2	2
GPB2	14	7	7
PDE2	11	5	6
Other Ras/PKA pathway			
CYR1	3	1	2
TFS1	1	0	1
RAS2	1	0	1
TOR/Sch9 pathway			
KOG1	1	0	1
SCH9	1	0	1
TOR1	1	0	1
Other Adaptive	7	0	7
Neutral	3	0	3
TOTAL	292	60	232

1810

Table S1. List of all mutants included in this study.

1811

