

# Title: The genotype-phenotype landscape of an allosteric protein

Authors: Drew S. Tack<sup>1</sup>, Peter D. Tonner<sup>1</sup>, Abe Pressman<sup>1</sup>, Nathanael D. Olson<sup>1</sup>, Sasha F. Levy<sup>2,3</sup>,  
Eugenia F. Romantseva<sup>1</sup>, Nina Alperovich<sup>1</sup>, Olga Vasilyeva<sup>1</sup>, David Ross<sup>1\*</sup>.

## Affiliations:

<sup>1</sup>National Institute of Standards and Technology, Gaithersburg, MD, 20899, USA.

<sup>2</sup>SLAC National Accelerator Laboratory, Menlo Park, CA, 94025, USA.

<sup>3</sup>Joint Initiative for Metrology in Biology, Stanford, CA, 94305, USA.

\*Correspondence to: david.ross@nist.gov.

## Abstract

Allostery is a fundamental biophysical mechanism that underlies cellular sensing, signaling, and metabolism. Yet a quantitative understanding of allosteric genotype-phenotype relationships remains elusive. Here we report the large-scale measurement of the genotype-phenotype landscape for an allosteric protein: the *lac* repressor from *Escherichia coli*, LacI. Using a method that combines long-read and short-read DNA sequencing, we quantitatively measure the dose-response curves for nearly 10<sup>5</sup> variants of the LacI genetic sensor. The resulting data provide a quantitative map of the effect of amino acid substitutions on LacI allostery and reveal systematic sequence-structure-function relationships. We find that in many cases, allosteric phenotypes can be quantitatively predicted with additive or neural-network models, but unpredictable changes also occur. For example, we were surprised to discover a new band-stop phenotype that challenges conventional models of allostery and that emerges from combinations of nearly silent amino acid substitutions.

## Introduction

Allostery is a fundamental biophysical mechanism that underlies cellular regulatory processes including sensing, signaling, and metabolism<sup>1-3</sup>. With allosteric regulation, ligand binding at one site on a biomolecule causes a conformational change that affects the activity of another, often distal, site. This conformational switching provides a sense-and-response function that defines the allosteric phenotype. Quantitative descriptions relating that phenotype to its causal genotype would improve our understanding of cellular function and evolution, and advance protein design and engineering<sup>4-6</sup>. However, the intramolecular interactions that mediate allosteric regulation are complex and distributed widely across the biomolecular structure, making the development of general quantitative descriptions challenging.

Recently described fitness landscape approaches have enabled the phenotypic characterization of 10<sup>4</sup> to 10<sup>5</sup> genotypes simultaneously<sup>7-12</sup>. Measurements at this scale facilitate the exploration of genotypes with widely distributed mutations, making them ideal for probing complex biological mechanisms like allostery. However, to quantitatively characterize the sense-and-response phenotypes inherent to allostery, a measurement must encompass the full dose-response curve that describes biomolecular activity as a function of ligand concentration.

1 Genetic sensors have served as a model of allosteric regulation for decades, and today are central to  
2 engineering biology. Genetic sensors are allosteric proteins that regulate gene expression in response to  
3 stimuli, giving cells the ability to regulate their metabolism and respond to environmental changes. Like  
4 many allosteric genetic sensors, the *lac* repressor, LacI, binds to DNA upstream of regulated genes,  
5 preventing transcription. Ligand binding to LacI causes a switch from the DNA-binding conformation to a  
6 non-binding conformation that allows transcription to proceed. This conformational switching results in  
7 the allosteric phenotype that is quantitatively defined by a dose-response curve relating the  
8 concentration of input ligand ( $L$ ) to the output response (the expression level of regulated genes,  $G$ ).  
9 Genetic sensors typically have sigmoidal dose-response curves following the Hill equation:

10

$$G(L) = G_0 + \frac{G_\infty - G_0}{1 + \left(\frac{EC_{50}}{L}\right)^n}$$

11 where  $G_0$  is basal gene expression in the absence of ligand,  $G_\infty$  is gene expression at saturating ligand  
12 concentrations,  $EC_{50}$  is the effective concentration of ligand that results in gene expression midway  
13 between  $G_0$  and  $G_\infty$ , and  $n$  quantifies the steepness of the dose-response curve (Fig. 1e). The dose-  
14 response curve is affected by several biophysical constants including ligand-binding affinity, DNA-binding  
15 affinity, and the allosteric constant (the equilibrium between the two conformations)<sup>2,13-15</sup>. These  
16 constants depend on amino acid residues and interactions spread widely across the protein, making it  
17 difficult to predict the effects of changes in protein sequence.

## 18 Results

### 19 Measuring the genotype-phenotype landscape

20 To measure the genotype-phenotype landscape for the allosteric LacI sensor, we first created a library of  
21 LacI variants using mutagenic PCR and attached a DNA barcode to the coding DNA sequence of each  
22 variant (Fig. 1a). We inserted the barcoded library into a plasmid where LacI regulates the expression of  
23 a tetracycline resistance gene (Supplementary Fig. 1a). Consequently, in the presence of tetracycline,  
24 the LacI dose-response modulates cellular fitness based on the concentration of the input ligand  
25 isopropyl- $\beta$ -D-thiogalactoside (IPTG). We then transformed the library into *E. coli* for the landscape  
26 measurement (Fig. 1b). To ensure that most variants in the library could regulate gene expression, we  
27 used fluorescence-activated cell sorting (FACS) to enrich the library for variants with low  $G_0$ . Then, using  
28 high-accuracy, long-read sequencing<sup>16</sup>, we determined the genotype for every variant in the library and  
29 indexed each variant to its attached DNA barcode (Fig. 1a).

30 The library contained 62,472 different LacI genotypes, with an average of 7.0 single nucleotide  
31 polymorphisms (SNPs) per genotype. Many SNPs were synonymous, i.e. coded for the same amino acid,  
32 so the library encoded 60,398 different amino acid sequences with an average of 4.4 amino acid  
33 substitutions per variant (Supplementary Fig. 2b).

34 To quantitatively determine the allosteric phenotype for every LacI variant in the library, we developed  
35 a new method to characterize the dose-response curves for large genetic sensor libraries. Briefly, we  
36 grew *E. coli* containing the library in 24 chemical environments (12 ligand concentrations, each with and  
37 without tetracycline). We used short-read sequencing of the DNA barcodes to measure the relative  
38 abundance of each variant at four timepoints during growth (Fig. 1c). We then used the changes in  
39 relative abundance to determine the fitness associated with each variant in each environment (Fig. 1d).

1 Finally, for each variant in the library, we used the fitness difference (with vs. without tetracycline) from  
2 all 12 ligand concentrations to quantitatively determine the dose-response curve using Bayesian  
3 inference (Fig. 1e). Most variants had sigmoidal dose-response curves (e.g. Supplementary Figs. 3-4),  
4 which we quantitatively described using the parameters of the Hill equation. We compared the results of  
5 variants with synonymous coding sequences and found that synonymous SNPs did not measurably  
6 impact the dose-response. So, for subsequent analysis we considered only amino acid substitutions.

7 To test the accuracy of the new method for library-scale dose-response curve measurements, we  
8 independently verified the results for over 100 LacI variants from the library. For each verification  
9 measurement, we chemically synthesized the coding DNA sequence for a single variant and inserted it  
10 into a plasmid where LacI regulates the expression of a fluorescent protein. We transformed the plasmid  
11 into *E. coli* and measured the resulting dose-response curve with flow cytometry (e.g. Fig. 1e). The flow  
12 cytometry results confirmed both the qualitative and quantitative accuracy of the new method  
13 (Supplementary Figs. 3-7).

#### 14 Effects of amino acid substitutions on LacI phenotype

15 During library construction, we chose the mutation rate to simultaneously achieve two objectives:  
16 exploration of a broad genotype-phenotype space, and acquisition of the single-amino-acid-substitution  
17 data most useful for building quantitative biophysical models of allosteric function<sup>2,13,14</sup>. Starting from  
18 the wild-type DNA sequence for LacI, there were 2110 possible SNP-accessible amino acid substitutions.  
19 Most of those substitutions were present in one or more variants within the library. However, nearly  
20 half were found only in combination with other substitutions. So, to comprehensively determine the  
21 impact of single amino acid substitutions, we constructed a deep neural network model (DNN) capable  
22 of accurately predicting the Hill equation parameters for LacI variants that were not directly measured.  
23 We adapted a recurrent architecture that captures the context dependence of mutational effects  
24 (Supplementary Fig. 8) and used an approximate Bayesian variational method to estimate uncertainties  
25 for the model predictions<sup>17</sup>.

26 We trained the model to predict the Hill equation parameters  $G_0$ ,  $G_\infty$ , and  $EC_{50}$  (Supplementary Fig. 9).  
27 For all three parameters, the root-mean-square error (RMSE) for the model predictions increases with  
28 the number of amino acid substitutions relative to the wild type (Supplementary Fig. 10). Importantly,  
29 for single-substitution variants, the model RMSE is comparable to the experimental measurement  
30 uncertainty (Supplementary Fig. 11). So, we could confidently integrate the experimental and DNN  
31 results to provide a nearly complete map of the effects of SNP-accessible amino acid substitutions.  
32 Furthermore, by integrating information about the causal substitutions from multiple genetic  
33 backgrounds, the model provided improved estimates of  $EC_{50}$  and  $G_\infty$  for variants with  $EC_{50}$  near or  
34 above the maximum ligand concentration measured (Supplementary Fig. 12).

35 The resulting map of single-substitution effects (Supplementary Data 1) includes quantitative point  
36 estimates and uncertainties for the Hill equation parameters for 94% of the possible SNP-accessible  
37 amino acid substitutions (1991 of 2110; 964 directly from measured data, and 1027 from DNN  
38 predictions). Most of the 119 substitutions missing from the dataset were probably excluded by FACS  
39 during library preparation because they caused a substantial increase in  $G_0$ . These include 83  
40 substitutions that have been shown to result in constitutively high  $G(L)$ <sup>18,19</sup>. Of the 1991 substitutions  
41 included in the dataset, 38% measurably affect the dose-response curve (beyond a 95% confidence  
42 bound).

1 The effect of any substitution depends strongly on its location within the protein structure, indicating  
2 systematic structure-function relationships underlying allostery (Fig. 2, using structural features as  
3 defined in references<sup>20-22</sup>). For example, substitutions that increase the basal expression,  $G_0$ , by more  
4 than 5-fold are located either in helix 4 of the DNA-binding domain, along the dimer interface, in the  
5 tetramerization helix, or at the protein start codon (Fig. 2a,d).  $G_0$  quantifies gene expression in the  
6 absence of ligand. So, apart from substitutions at the start codon that reduce the number of LacI  
7 proteins per cell<sup>23</sup>, these substitutions probably affect either the DNA-binding affinity, the allosteric  
8 constant, or both<sup>14</sup>. Interestingly, substitutions in helix 4 (R51C, Q54K, and L56M) and near the dimer  
9 interface (T68N, L71Q) that increase  $G_0$  also decrease  $EC_{50}$  approximately 10-fold, consistent with a  
10 change in the allosteric constant<sup>14</sup> (Supplementary Fig. 13a).

11 Amino acid substitutions that decrease ligand-saturated expression,  $G_\infty$ , by more than 5-fold are all  
12 located near the ligand-binding pocket or along the dimer interface (Fig. 2b,e). Six of these substitutions  
13 also increase  $EC_{50}$  more than 5-fold (A75T, D88N, S193L, Q248R, D275Y, and F293Y; Supplementary  
14 Fig. 13b). Except for D88N, which is at the dimer interface, these substitutions are in the ligand-binding  
15 pocket. Substitutions near the ligand-binding pocket probably change ligand-binding affinity, though  
16 studies with targeted substitutions have shown that they can also change the allosteric constant<sup>14</sup>.

17 Amino acid substitutions that change the effective concentration,  $EC_{50}$ , are the most numerous and are  
18 spread throughout the protein structure, with approximately 9% and 20% of all substitutions causing a  
19 greater than 5-fold or 2.5-fold shift in  $EC_{50}$ , respectively (Fig. 2c,f). The strongest effects are from  
20 substitutions in the DNA-binding domain, ligand-binding pocket, core-pivot domain, or dimer interface.  
21 Substitutions to the DNA-binding domain or dimer interface generally decrease  $EC_{50}$ . Substitutions to the  
22 ligand-binding pocket or core-pivot domain generally increase  $EC_{50}$ .

23 In addition to specific substitutions that affect both  $G_\infty$  and  $EC_{50}$ , we identified nine positions (N125,  
24 P127, D149, V192, A194, A245, N246, T276, Q291), where different substitutions either reduce  $G_\infty$  by  
25 more than 5-fold or increase  $EC_{50}$  by more than 5-fold, but not both. These positions are all in the ligand-  
26 binding pocket. We also identified five positions (H74, V80, K84, S97, M98) where different substitutions  
27 reduce either  $G_\infty$  or  $EC_{50}$  by more than 5-fold but not both. These positions are all located at the dimer  
28 interface.

29 Combining multiple substitutions in a single protein almost always has a log-additive effect on  $EC_{50}$ . Only  
30 0.57% (12 of 2101) of double amino acid substitutions have  $EC_{50}$  values that differ from the log-additive  
31 effects of the single substitutions by more than 2.5-fold (Fig. 3). This result, combined with the wide  
32 distribution of residues that affect  $EC_{50}$ , suggests that LacI allostery is controlled by a free energy  
33 balance with additive contributions from many residues and interactions.

### 34 Phenotypic innovation in an allosteric landscape

35 Beyond the comprehensive mapping of single-substitution effects, the LacI genotype-phenotype  
36 landscape measurement revealed a surprising number of variants with phenotypes that differ  
37 qualitatively from the wild type. For example, approximately 230 of the LacI variants have an inverted  
38 phenotype ( $G_0 > G_\infty$ , Fig. 1e), accounting for approximately 0.35% of the measured library  
39 (Supplementary Fig. 2a). We verified the dose-response curves for 10 inverted variants with flow  
40 cytometry (e.g. Supplementary Fig. 4). To understand the mutational basis for the inverted phenotype,  
41 we examined a set of 43 strongly inverted variants (with  $G_0/G_\infty > 2$ ,  $G_0 > G_{\infty,wt}/2$ , and  $EC_{50}$  between

1 3  $\mu\text{mol/L}$  and 1000  $\mu\text{mol/L}$ ). The results indicate that diverse substitutions can lead to the inverted  
2 phenotype. For example, we identified 10 amino acid substitutions associated with the inverted  
3 phenotype (S70I, K84N, D88Y, V96E, A135T, V192A, G200S, Q248H, Y273H, A343G; Fig. 4a,c). However,  
4 none of these substitutions are present in more than 12% of the strongly inverted variants, and 51% of  
5 the strongly inverted variants have none of these substitutions. Furthermore, the set of strongly  
6 inverted variants are more genetically distant from each other than randomly selected variants from the  
7 library (Fig. 4c, Supplementary Fig. 14).

8 The inverted LacI variants can provide specific insight into allosteric biophysics and structure-function  
9 relationships, since inversion of the dose-response curve requires inversion of both the allosteric  
10 constant<sup>13</sup> and the relative ligand-binding affinity between the two conformations<sup>2</sup>. Although the set of  
11 strongly inverted LacI variants are genetically diverse, many of them have substitutions in similar regions  
12 of the protein that may account for the requisite biophysical changes. First, 67% of the strongly inverted  
13 variants have substitutions near the ligand-binding pocket (within 7 Å), which likely contribute to the  
14 change in ligand-binding affinity. Surprisingly, 21% of the strongly inverted variants have no  
15 substitutions within 10 Å of the binding pocket, so binding affinity must be indirectly affected by distal  
16 substitutions in those variants. Second, nearly all strongly inverted variants have substitutions at the  
17 dimer interface (91%, compared to 54% for the full library), with most (70%) having substitutions in  
18 helix 5 (47%), helix 11 (28%), or both (5%, Fig. 4a,c). This suggests that residues in those structural  
19 features are important for modulating the allosteric constant.

## 20 Discovery of novel allosteric phenotypes

21 In addition to the inverted phenotypes, we were surprised to discover LacI variants with dose-response  
22 curves that did not match the sigmoidal form of the Hill equation. Specifically, we found variants with  
23 band-pass or band-stop dose-response curves, i.e. variants that repress or activate gene expression only  
24 over a narrow range of ligand concentrations (e.g. Fig. 1e). Approximately 200 of the LacI variants have  
25 band-stop or band-pass phenotypes, accounting for approximately 0.3% of the measured library  
26 (Supplementary Fig. 2a). We verified the dose-response curves of 13 band-stop variants and two band-  
27 pass variants using flow cytometry (e.g. Supplementary Fig. 5-6). To our knowledge, this is the first  
28 identification of single-protein genetic sensors with band-stop dose-response curves.

29 Phenotypic similarities between the band-stop and inverted LacI variants (i.e. high  $G_0$ , and initially  
30 decreasing gene expression as ligand concentration increases) imply similar biophysical requirements.  
31 However, amino acid substitutions associated with the band-stop phenotype are remarkably different  
32 from those for the inverted phenotype. While inverted variants often have substitutions near the ligand-  
33 binding pocket and dimer interface, a set of 31 strong band-stop variants are twice as likely as the full  
34 library to have substitutions in helix 9 (32% compared to 16%) and nearly four times as likely to have  
35 substitutions in strand J (13% compared to 3.4%). Helix 9 is on the periphery of the protein, and strand J  
36 is in the center of the C-terminal core domain. Furthermore, 100% of the strong band-stop variants have  
37 substitutions in the C-terminal core of the protein, compared with 79% of the full library (Fig. 4b,d).

38 To further investigate the band-stop phenotype, we chose a strong band-stop LacI variant with only  
39 three amino acid substitutions (R195H/G265D/A337D). We synthesized LacI variants with all possible  
40 combinations of those substitutions and measured their dose-response curves with flow cytometry.  
41 Although each single substitution resulted in a sigmoidal dose-response similar to wild-type LacI, the  
42 combination of two substitutions (R195H/G265D) gave rise to the band-stop phenotype (Fig. 5a,

1 Supplementary Fig. 15). To test whether this result applies to the band-stop phenotype generally, we  
2 used the single-substitution effects presented above to examine each of the substitutions associated  
3 with the strong band-stop phenotype. Individually, the substitutions associated with the band-stop  
4 phenotype are nearly silent, i.e. they have little or no effect on the dose-response curve; yet in  
5 combination with other substitutions, they result in the band-stop phenotype. In contrast, most of the  
6 individual substitutions associated with the inverted phenotype cause a large shift in either  $EC_{50}$ ,  $G_{\infty}$ , or  
7 both (Fig. 5b,c).

## 8 Discussion

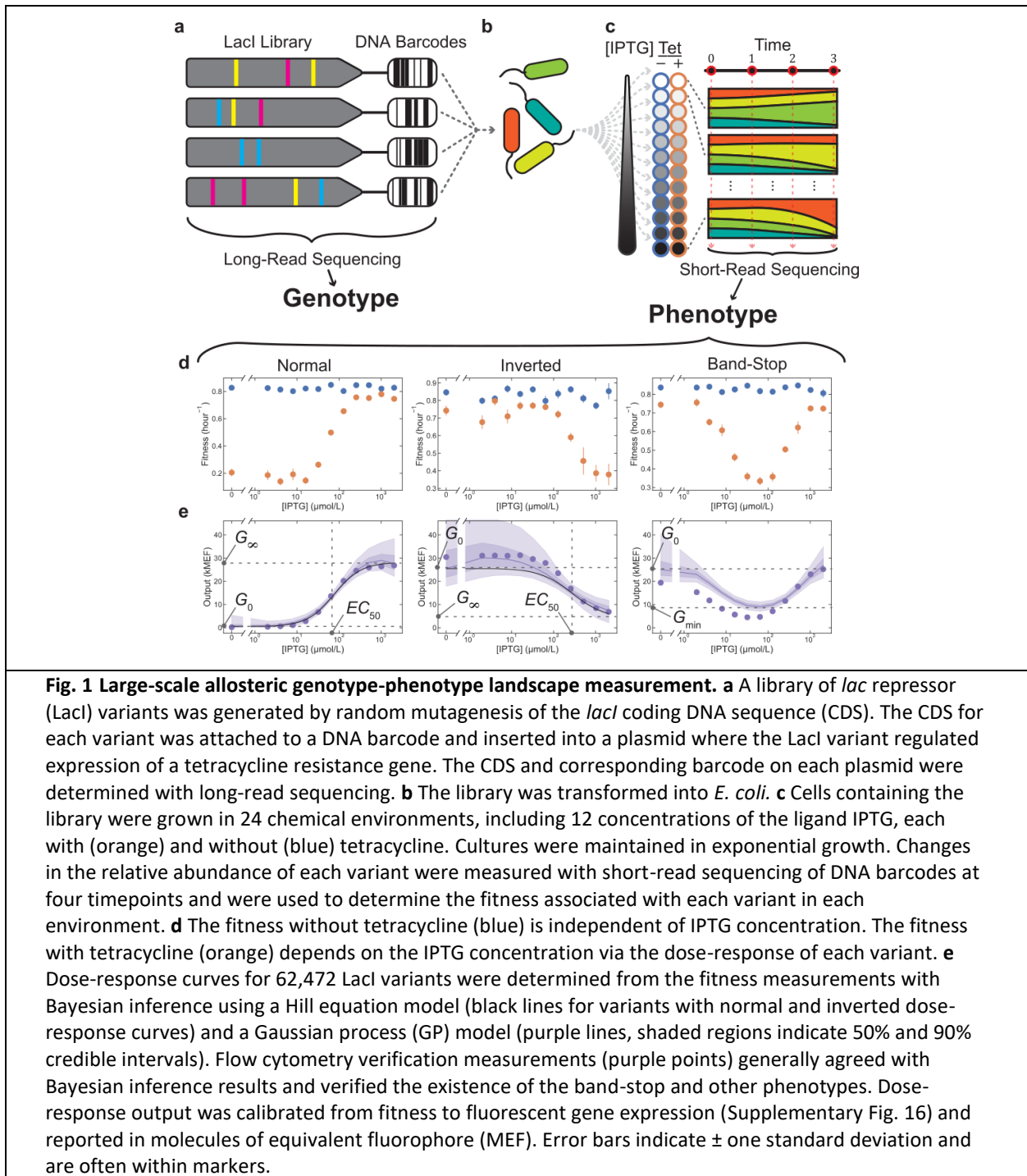
9 For the goal of an improved understanding of allostery, our results reveal the dual nature of the  
10 problem: First, the DNN model and the mapping of single-substitution effects demonstrate that large-  
11 scale measurements and analysis can overcome the challenges inherent to the structural complexity of  
12 allosteric function. They can provide accurate predictions for specific allosteric proteins and can also  
13 reveal systematic structure-function relationships that may be more generalizable (i.e. the importance  
14 of the dimer interface and the log-additivity of  $EC_{50}$ ). However, the band-stop phenotype highlights the  
15 limits of that predictability, as well as the constraints of conventional models of allostery. While the  
16 allosteric function of many LacI variants is well-described by the Monod-Wyman-Changeux (MWC)  
17 model of allostery<sup>2,13,14</sup>, the band-stop phenotype is inconsistent with that model. In particular, the  
18 biphasic dose-response of the band-stop variants suggests negative cooperativity and that the relevant  
19 free-energy changes may be more entropic than structural<sup>1</sup>. Our most surprising and unpredictable  
20 result is the emergence of the band-stop phenotype from combinations of nearly silent amino acid  
21 substitutions. However, with over one hundred genetically diverse band-stop variants, our dataset  
22 provides a basis for more systematic understanding even in this case. Furthermore, the relatively high  
23 abundance of inverted and band-stop variants (approximately 0.35% and 0.2% of the library,  
24 respectively, Supplementary Fig. 2a) with genotypes near the wild-type suggests that allosteric  
25 genotype-phenotype landscapes allow for rapid evolutionary innovation, a conclusion that is supported  
26 by the existence of natural transcription factors related to LacI with inverted phenotypes<sup>24,25</sup>.

27 Overall, our findings suggest that a surprising diversity of useful and potentially novel allosteric  
28 phenotypes exist with genotypes that are discoverable only via large-scale landscape measurements.

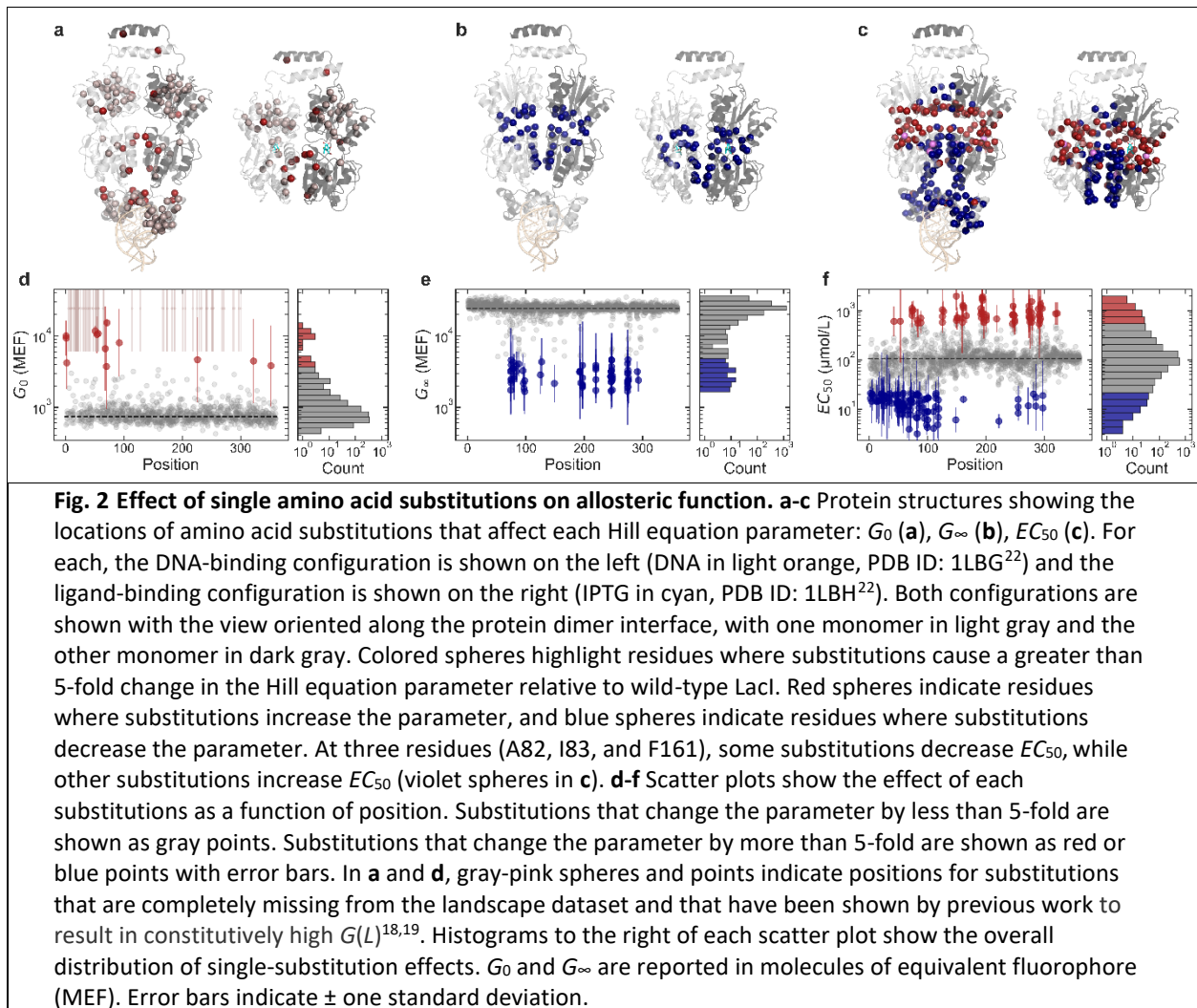
29

30

1 Figures

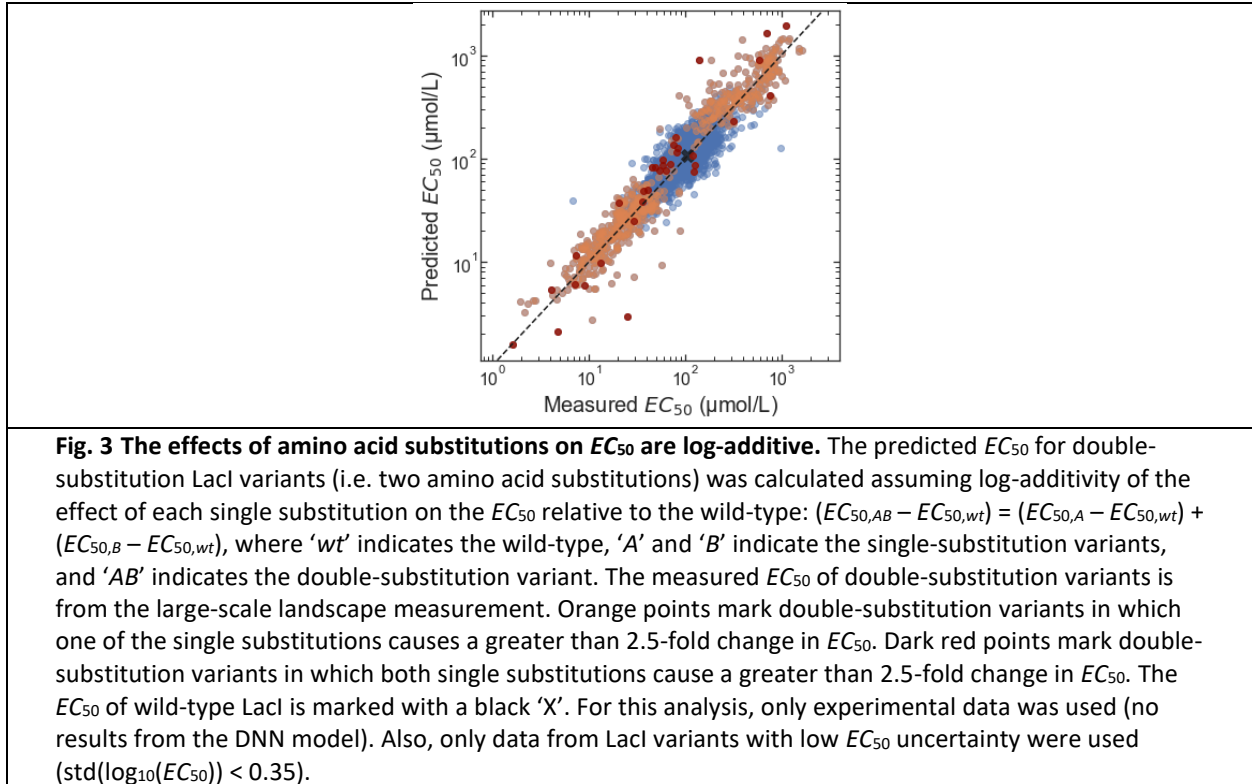


1  
2

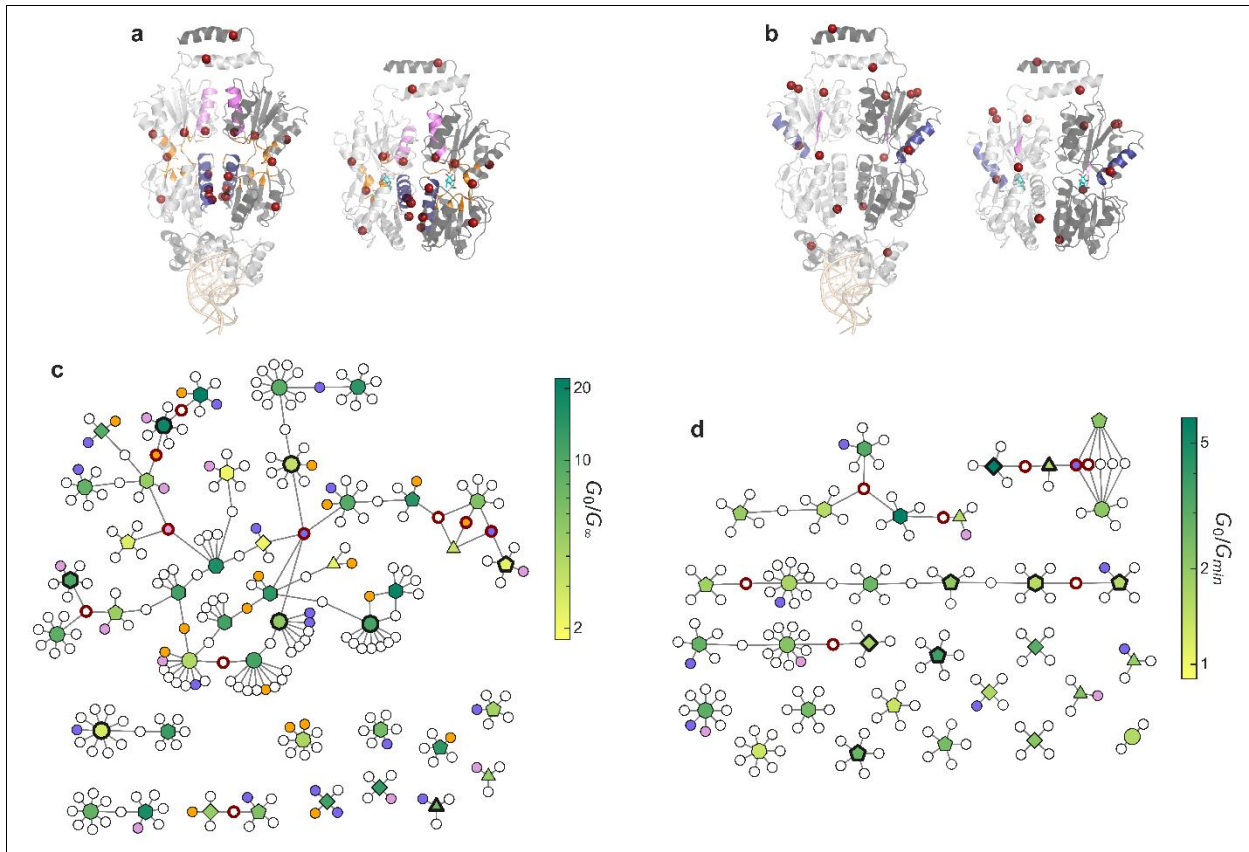


3  
4  
5



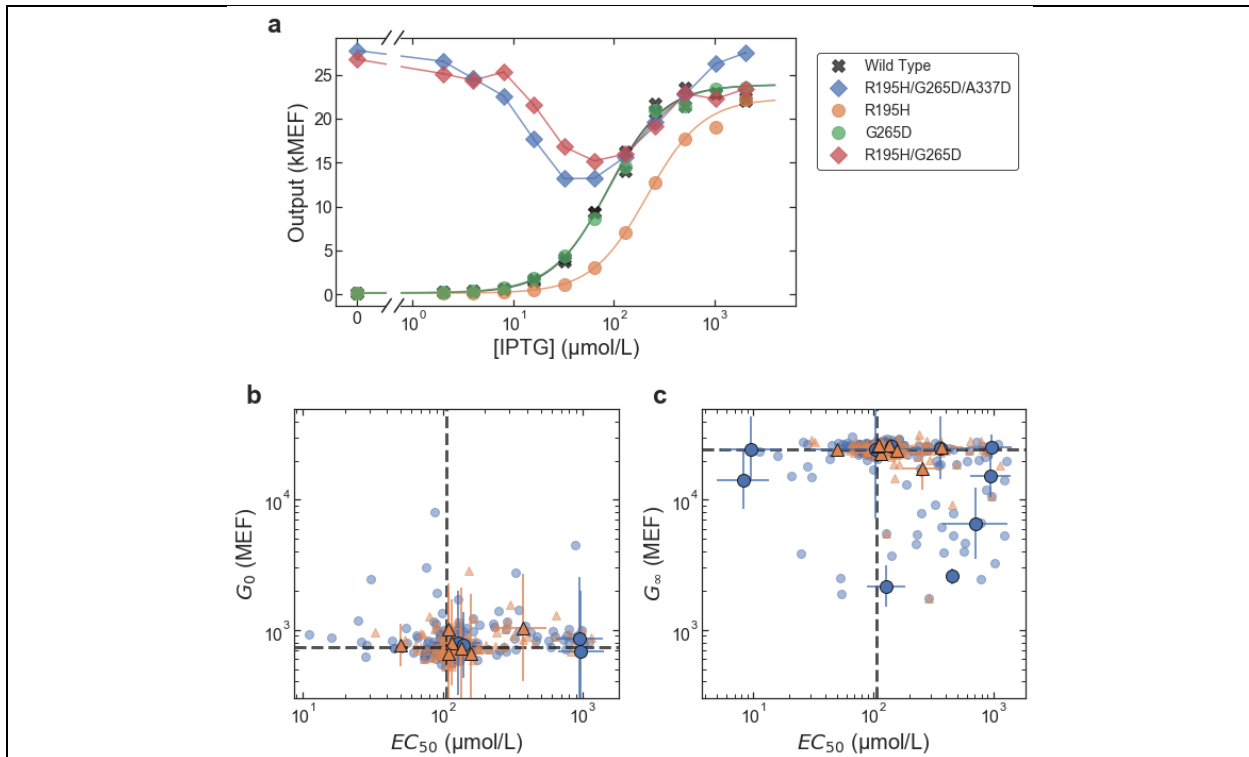


1



**Fig. 4 Analysis of inverted and band-stop genotypes.** **a-b** Location of amino acid substitutions associated with strongly inverted (**a**) and strong band-stop (**b**) phenotypes. For each plot, the DNA-binding configuration of Lacl is shown on the left (PDB ID: 1LGB), with the DNA operator at the bottom in light orange; the ligand-binding configuration is shown on the right (PDB ID: 1LBH), with IPTG in cyan. Both configurations are shown with the view oriented along the protein dimer interface, with one monomer in light gray and the other monomer in dark gray. The locations of associated (i.e. high-frequency) amino acid substitutions are highlighted as red spheres, and secondary structures where inverted or band-stop variants have amino acid substitutions at a significantly higher frequency than the full library are shaded with different colors. For strongly inverted variants (**a**), helix 5 is shaded blue, helix 11 is shaded violet, and the residues near the ligand-binding pocket are shaded orange. For strong band-stop variants (**b**), helix 9 is shaded blue, and strand J is shaded violet. **c-d** Network diagrams showing relatedness among genotypes for strongly inverted (**c**) and strong band-stop (**d**) variants. Within each network diagram, Lacl variants are represented by polygonal nodes, with a colormap indicating the  $G_0/G_{\infty}$  or  $G_0/G_{min}$  ratio (see Fig. 1e). The number of sides of the polygon indicates the number of amino acid substitutions relative to the wild-type, and bold outlines indicate variants that were verified with flow cytometry. Smaller circular nodes represent substitutions, with lines showing the substitutions for each variant. Bold red outlines on the substitution nodes indicate the associated substitutions shown as spheres in **a-b**, and the shading of substitution nodes matches the shading used to highlight secondary structures in **a-b**.

1



**Fig. 5 The band-stop phenotype emerges from combinations of nearly silent amino acid substitutions.**

**a** Dose-response curves measured with flow cytometry for selected LacI variants: wild-type LacI (grey 'X's), a strong band-stop variant identified from the library with only three amino acid substitutions (R195H/G265D/A337D; blue diamonds), LacI variants containing the single substitutions R195H (orange circles) and G265D (green circles), LacI variant with the double substitutions R195H/G265D (red diamonds). The single substitution R195H (orange) or G265D (green) results in sigmoidal dose-response curves similar to wild-type LacI, but the combination of the two, R195H/G265D (red), results in a band-stop phenotype. The complete set of permutations of R195H, G265D, and A337D are shown in Supplementary Fig. 15. **b-c** Effects of individual amino acid substitutions associated with inverted and band-stop phenotypes. Each plot shows the joint effect of individual amino acid substitutions on two Hill equation parameters. The blue circles plotted with error bars show the effects of substitutions associated with the strongly inverted phenotype and the orange triangles plotted with error bars show the effects of substitutions associated with the strong band-stop phenotype. Most substitutions associated with the inverted phenotype cause a large shift in either  $EC_{50}$ ,  $G_\infty$ , or both, consistent with the biophysical requirements for inverting the dose-response curve. In contrast, most of the amino acid substitutions associated with the band-stop phenotype are nearly silent. Light blue circles and light orange triangles show the effects for all amino acid substitutions found in the sets of strongly inverted and strong band-stop variants, respectively. Dashed gray lines mark the wild-type parameter values. Plotted data includes a combination of direct experimental measurements and DNN model predictions and is included in Supplementary Data 1. Error bars indicate  $\pm$  one standard deviation.

1

2

## 1 Methods

### 2 Strain, plasmid, and library construction

3 All reported measurements were completed using *E. coli* strain MG1655 $\Delta$ lac<sup>26</sup>. Briefly, strain  
4 MG1655 $\Delta$ lac was constructed by replacing the lactose operon of *E. coli* strain MG1655 (ATCC #47076)  
5 with the bleomycin resistance gene from *Streptoalloteichus hindustanus* (*Shble*).

6 Two plasmids were used for this work: a library plasmid (pTY1, Supplementary Fig. 1a) used for the  
7 measurement of the genotype and phenotype of the entire LacI library, and a verification plasmid  
8 (pVER, Supplementary Fig. 1b) used to verify the function of over 100 LacI variants from the library  
9 chosen to test the accuracy of the library-scale dose-response curve measurement method. A step-by-  
10 step description of the plasmid assembly protocol is available<sup>27</sup>. The sequences are available in GenBank  
11 (MT702633, MT702634, for pTY1 and pVER, respectively).

12 Plasmid pTY1 contained the *lacI* coding DNA sequence (CDS) and the lactose operator (*lacO*) regulating  
13 the transcription of a tetracycline resistance gene, *tetA*, which, in the presence of tetracycline, confers a  
14 measurable change in fitness connected with the expression level of the regulated genes. Plasmid pTY1  
15 also encoded Enhanced Yellow Fluorescent Protein (YFP), which was used during library construction to  
16 select a library in which most of the LacI variants could function as allosteric repressors (see below).

17 Plasmid pVER contained a similar system in which LacI and *lacO* regulate the transcription of only YFP.  
18 Plasmid pVER was used to measure dose-response curves of clonal LacI variants using flow cytometry.  
19 Each variant chosen from the library for verification was chemically synthesized (Twist Biosciences),  
20 inserted into pVER, and transformed into *E. coli* strain MG1655 $\Delta$ lac for flow cytometry measurements to  
21 confirm the dose-response curve inferred from the library-scale measurements.

22 The LacI library was generated by error-prone PCR of the wild-type *lacI*. The library was inserted into  
23 pTY1 along with randomly synthesized DNA barcodes. Each barcode consisted of 54 random nucleotides  
24 introduced with PCR primers (Integrated DNA Technologies). Most of the variants in the initial library  
25 had high  $G(0)$ , i.e. the  $\Gamma$  phenotype<sup>18</sup>. To generate a library in which most of the LacI variants could  
26 function as allosteric repressors, we used fluorescence activated cell sorting (FACS) to select a portion of  
27 the library with low fluorescence in the absence of ligand (Sony SH800S Cell Sorter). To allow  
28 comprehensive long-read sequencing of the library (PacBio sequel II, see Long-read sequencing section,  
29 below), we further reduced the library size by dilution of the FACS-selected library to create a  
30 population bottleneck of the desired size. For the work reported here, we used a library of  
31 approximately  $10^5$  LacI variants (determined by serial plating and colony counting).

32 A spike-in control strain was used to normalize the DNA barcode read counts for the sequencing-based  
33 fitness measurement (see Library-scale fitness measurement section, below). The spike-in control strain  
34 contained the Library Plasmid with a LacI variant that had a constant, high *tetA* expression level. The  
35 fitness of the spike-in control was determined from OD<sub>600</sub> data acquired during growth of clonal cultures  
36 with the same automated growth protocol as used for the genotype-phenotype landscape  
37 measurement (see Growth protocol for landscape measurement section, below). The fitness of the  
38 spike-in control was measured in all 24 chemical environments and was independent of IPTG  
39 concentration but was slightly lower with tetracycline (0.75 hour<sup>-1</sup>) than without tetracycline  
40 (0.81 hour<sup>-1</sup>).

1

## 2 Culture conditions

3 Unless otherwise noted, *E. coli* cultures were grown in a rich M9 media (3 g/L KH<sub>2</sub>PO<sub>4</sub>, 6.78 g/L  
4 Na<sub>2</sub>HPO<sub>4</sub>, 0.5 g/L NaCl, 1 g/L NH<sub>4</sub>Cl, 0.1 mmol/L CaCl<sub>2</sub>, 2 mmol/L MgSO<sub>4</sub>, 4% glycerol, and  
5 20 g/L casamino acids) supplemented with 50 µg/mL kanamycin.

6 *E. coli* cultures were grown in a laboratory automation system that controlled preparation of 96-well  
7 culture plates with media and additives (i.e. IPTG and tetracycline). Cultures were grown in clear-bottom  
8 96-well plates with 1.1 mL square wells (4titude, 4ti-0255). The culture volume per well was 0.5 mL.  
9 Before incubation, an automated plate sealer (4titude, a4S) was used to seal each 96-well plate with a  
10 gas permeable membrane (4titude, 4ti-0598). Cultures were incubated in a multi-mode plate reader  
11 (BioTek, Neo2SM) at 37 °C with a 1 °C gradient applied from the bottom to the top of the incubation  
12 chamber to minimize condensation on the inside of the membrane. During incubation, the plate reader  
13 was set for double-orbital shaking at 807 cycles per minute. Optical density at 600 nm (OD<sub>600</sub>) was  
14 measured every 5 minutes during incubation, with continuous shaking applied between measurements.  
15 After incubation, an automated de-sealer (Brooks, XPeel) was used to remove the gas permeable  
16 membrane from each 96-well plate.

## 17 Growth protocol for landscape measurement

18 To measure the fitness and dose-response curve of every LacI variant in the library, a culture of *E. coli*  
19 containing the LacI library was mixed at a 99:1 ratio with a culture of the *E. coli* spike-in control. The  
20 culture was loaded into the automated microbial growth and measurement system where it was  
21 distributed across a 96-well plate and then grown to stationary phase (12 hours). Cultures were then  
22 diluted 50-fold into a new 96-well plate, Growth Plate 1, containing 11 rows with a 2-fold serial dilution  
23 gradient of IPTG with concentrations ranging from 2 µmol/L to 2048 µmol/L and one row without IPTG.  
24 Growth in IPTG allowed each variant to reach a steady state tetA expression level in each IPTG  
25 concentration. Growth Plate 1 was grown for 160 minutes, corresponding to approximately  
26 3.3 generations, and then diluted 10-fold into Growth Plate 2. Growth Plate 2 contained the same IPTG  
27 gradient as Growth Plate 1 with the addition of tetracycline (20 µg/mL) to alternating rows in the plate,  
28 resulting in 24 chemical environments, with 4 duplicate wells for each environment. Growth Plate 2 was  
29 grown for 160 minutes and then diluted 10-fold into Growth Plate 3, which contained the same  
30 24 chemical environments as Growth Plate 2. This process was repeated for Growth Plate 4, which also  
31 contained the same 24 chemical environments. The total growth time for the fitness measurements in  
32 the 24 chemical environments, 480 minutes across Growth Plates 2-4, corresponded to approximately  
33 10 generations for the fastest-growing cultures. The 50-fold dilution factor from stationary phase into  
34 Growth Plate 1 and the 160 minute growth time per plate were chosen to maintain the cultures in  
35 exponential growth for the entire 480 minutes. During each 160 minute incubation, the cultures without  
36 tetracycline increased approximately 10-fold in optical density, to a final OD<sub>600</sub> of approximately 0.5  
37 (corresponding to an estimated cell density of  $4 \times 10^8$  cells/mL).

38 After each growth plate was used to seed the subsequent plate (or at the end of 160 minutes for  
39 Growth Plate 4), the remaining culture volumes for each chemical environment (approximately  
40 450 µL/well, four duplicates per plate) were combined and pelleted by centrifugation (3878 g for  
41 10 minutes at 23 °C). Plasmid DNA was then extracted from the 24 combined samples with a custom  
42 method using reagents from the QIAprep Miniprep Kit (Qiagen cat. #27104) on an automated liquid

1 handler equipped with a positive-pressure filter press (step-by-step protocol available<sup>28</sup>). After  
2 extraction, DNA was eluted into a final volume of approximately 50  $\mu$ L and the concentration of DNA in  
3 each sample ranged from undetectable up to approximately 1.5 ng/ $\mu$ L. This corresponds to an estimated  
4 maximum of approximately  $10^{10}$  plasmids per sample.

### 5 Barcode sequencing

6 After plasmid extraction, each set of 24 plasmid DNA samples was prepared for barcode sequencing  
7 using a custom sequencing sample preparation method on a second automated liquid handler (step-by-  
8 step protocol is available<sup>29</sup>). Briefly, the plasmid DNA was linearized with Apal restriction enzyme. Then,  
9 a 3-cycle PCR was performed to attach sample multiplexing tags to the resulting amplicons so the  
10 different samples could be distinguished when pooled and run on the same sequencing flow cell. Eight  
11 forward index primers and 12 reverse index primers were used to label the amplicons from each sample  
12 across the 24 chemical environments and the four time points. After a magnetic-bead-based cleanup  
13 step, a second, 15-cycle PCR was run to attach the standard Illumina paired-end adapter sequences and  
14 to amplify the resulting amplicons for sequencing. After a second magnetic-bead-based cleanup, the  
15 24 samples from each time point were pooled and stored at 4 °C until sequencing. For sequencing, DNA  
16 was diluted to a final concentration of approximately 5 nmol/L and combined with 20% phiX control  
17 DNA. DNA from each of the 4 time points was sequenced in a separate lane on an Illumina HiSeqX using  
18 paired-end mode with 150 bp in each direction.

19 To count DNA barcodes and estimate the fitness associated with each LacI variant, the sequencing data  
20 was analyzed using custom software written in C# and Python, and the Bartender1.1 barcode clustering  
21 algorithm<sup>30</sup> ([https://github.com/djross22/nist\\_lacI\\_landscape\\_analysis](https://github.com/djross22/nist_lacI_landscape_analysis)).

22 The sequence of the nominal Illumina compatible amplicon was (with Illumina adapters and flow cell  
23 binding sequences in gray):

```
24 AATGATACGGCGACCACCGAGATCTACACTCTTTCCCTACACGACGCTCTTCCGATCTZZZZZZZZXXXXXX  
25 XXXXXCATCGGGT GAGCCCCGGGCTGTCCGCGTNNNTNNNANNTNNNANNTNNNANNTNNNANNTNNNANNTNNN  
26 TGCCAGCAGGCCGGCCACGCTNNNTNNNANNTNNNANNTNNNANNTNNNANNTNNNANNTNNNANNTNNNANNTNNN  
27 GGCCGCACGATGCGTCCGGCGTAA GAGGXXXXXXXXXXZZZZZZZAGATCGGAAGAGCGGTTCAGCAGGAA  
28 TGCCGAGACCGATCTCGTATGCCGTCTTCTGCTTG
```

29 The nominal forward and reverse reads from paired-end barcode sequencing were:

```
30 ZZZZZZZZXXXXXXXXXXCATCGGGT GAGCCCCGGGCTGTCCGCGTNNNTNNNANNTNNNANNTNNNANNTNNN  
31 ANNTNNNANNATATG
```

32 and

```
33 ZZZZZZZZXXXXXXXXXXCCTCTACGCCGACGCATCGTGC GGCCGGCCGGCCGGCC CACCGNNNTNNNANNTNNN  
34 ANNTNNNTNNNTNNNANNTNNNANNAGCGT
```

35 The Z's at the beginning of each read are random nucleotides used as unique molecular identifiers  
36 (UMIs) to correct for PCR jackpotting<sup>31</sup>, the X's are the sample multiplexing tag sequences, and the N's  
37 are the random nucleotides of the DNA barcodes. To minimize the chances of barcode crosstalk, we  
38 used dual barcodes, with independent random barcode sequences on the forward and reverse reads  
39 and 27 random nucleotides in each of the forward and reverse barcodes.

- 1 The raw sequences were parsed, and sequences were kept for further analysis only if they passed the  
2 following quality criteria for both the forward and reverse reads:
- 3 1. The four bases after the multiplexing tag (highlighted in yellow above) must match the nominal  
4 sequence with one allowed mismatch, and the multiplexing tag sequence (highlighted in pink  
5 above) must match the nominal sequence for one of the multiplexing tags used with up to three  
6 allowed mismatches.
  - 7 2. The five flanking bases before and after the barcodes (highlighted in cyan above) must match  
8 the nominal sequence with one allowed mismatch per set of five bases, and the number of  
9 bases in the barcode (highlighted in green above) must be between 35 and 41 (inclusive).
  - 10 3. The mean Illumina quality score for the barcode and the five flanking bases before and after the  
11 barcode must be greater than 30.

12 For the four lanes of HiSeq data, there were 2,024,537,456 raw reads, of which 1,576,168,836 reads  
13 passed the quality criteria (78%). Note that 20% of the DNA sample loaded onto the HiSeq instrument  
14 was phiX DNA.

15 True barcode sequences were identified using the Bartender1.1 clustering algorithm<sup>30</sup> with the following  
16 parameter settings: maximum cluster distance = 4, cluster merging threshold = 8, cluster seed  
17 length = 5, cluster seed step = 1, frequency cutoff = 500. Barcodes from the forward and reverse reads  
18 were clustered independently. The Bartender1.1 clustering algorithm identified 43,259 distinguishable  
19 forward barcode clusters and 31,055 distinguishable reverse barcode clusters.

20 To correct for insertion-deletion read errors, barcode clusters of different length were considered for  
21 merging. First, barcode clusters with sequences that were sub-strings of one another were automatically  
22 merged. Second, pairs of barcode clusters with a DNA sequence Levenshtein distance of 1 or 2 were  
23 merged if the ratio of the smaller cluster read count to the total read count of both clusters was less  
24 than 0.001 and 0.0001, respectively. Third, all barcode clusters with a Levenshtein distance less than 7  
25 from the barcode for the spike-in control were merged.

26 After merging barcode clusters of different lengths, there were 43,169 distinguishable forward barcode  
27 clusters and 30,931 distinguishable reverse barcode clusters. The random positions within the forward  
28 and reverse barcodes had approximately equal probabilities for each nucleotide, with a mean entropy  
29 per position of 1.9799 bits  $\pm$  0.0066 bits.

30 After barcode clustering and merging, the barcode sequencing reads were sorted based on the sample  
31 multiplexing tags and the barcode read counts were corrected for PCR jackpotting effects. Sets of  
32 multiple barcode reads were treated as PCR jackpot duplicates if they had the same UMI sequence, the  
33 same multiplexing tag, and the same barcode sequence for both forward and reverse barcode reads. In  
34 the corrected barcode count, each set of PCR jackpot duplicates was counted as a single read.  
35 Approximately 15% of the total barcode sequencing reads were found to be PCR jackpot duplicates.

36 The forward and reverse barcodes were then combined to give the DNA barcodes used to measure the  
37 relative abundance of each *LacI* variant in the library. An additional barcode count threshold was  
38 applied, keeping only DNA barcodes with a total read count (across all 24 environments and 4 time  
39 points) greater than 2000. A small number (139) of DNA barcodes were identified as likely chimeras with

1 forward and reverse barcodes combined from different plasmid templates<sup>32–34</sup>. The likely chimera  
2 barcodes were not used in further analysis.

3 Finally, 14 pairs of DNA barcodes were found with DNA sequence Hamming distance of one (across both  
4 forward and reverse barcodes). Only one DNA barcode from each pair was also found in the long-read  
5 sequencing data (see Long-read sequencing section, below). In addition, the fitness curves (vs. IPTG  
6 concentration) were very similar for both barcodes in each pair. Based on this, the read counts  
7 associated with each of those 14 pairs of dual barcodes were merged, and each pair was treated as a  
8 single DNA barcode.

9 The final set of 67,730 DNA barcodes was used for all subsequent analysis to extract estimates of the  
10 fitness and dose-response curve associated with each barcode.

### 11 Long-read sequencing

12 The full sequence of the Library Plasmid for every *LacI* variant in the library was measured using PacBio  
13 circular consensus HiFi sequencing. The HiFi sequencing data was used to determine the consensus *lacI*  
14 sequence for each variant and the corresponding DNA barcode. Of the 67,731 distinct DNA barcodes  
15 (see Barcode sequencing section, above), the HiFi sequencing data was used to determine the *lacI*  
16 sequences for 63,064 (93%), 3,878 with a single HiFi *lacI* read, and 59,186 with multiple HiFi *lacI* reads.

17 In addition, the full plasmid sequence was used to detect unintended mutations in the plasmid, i.e.  
18 mutations to plasmid regions other than the *lacI* CDS. For analysis of the HiFi read data, the full plasmid  
19 sequence was divided into 11 non-overlapping regions that roughly correspond to different functional  
20 elements of the plasmid (Supplementary Table 1), and the sequences for each region were extracted  
21 from the HiFi reads using a custom bioinformatic pipeline  
22 ([https://github.com/djross22/nist\\_lacI\\_landscape\\_analysis](https://github.com/djross22/nist_lacI_landscape_analysis)). The number of unintended mutations to  
23 plasmid regions other than the *lacI* CDS was relatively low (Supplementary Table 1), so it was not  
24 possible to examine mutational effects with base-pair- or residue-level resolution. However, by pooling  
25 the mutational information for each region, significant region-specific effects could be detected. To  
26 determine if mutations in a region of the plasmid had a significant effect, the estimated Hill equation  
27 parameters were compared for all variants with one or more mutations in a given plasmid region vs. all  
28 variants with zero mutations in that region. Significant differences in the geometric mean of one or  
29 more Hill equation parameters were found for variants with mutations in the following regions: tetA (p-  
30 value for  $\log_{10}(G_{\infty})$ :  $2 \times 10^{-56}$ ), KAN (p-value for  $\log_{10}(G_{\infty})$ :  $4 \times 10^{-11}$ ), origin of replication (p-value for  
31  $\log_{10}(G_{\infty})$ :  $6 \times 10^{-14}$ ), and YFP (p-value for  $\log_{10}(G_0)$ :  $4 \times 10^{-109}$ ; p-value for  $\log_{10}(G_{\infty})$ :  $5 \times 10^{-10}$ ; p-value for  
32  $\log_{10}(EC_{50})$ :  $2 \times 10^{-74}$ ), where the p-values given are for Welch's unequal-variances t-test.

33 In addition, 43 of the 535 variants with the wild-type *LacI* amino acid sequence had mutations in the  
34 regulatory region (containing the  $P_{lacI}$  and  $P_{tacI}$  promoters, the *lacO* operator, the *riboJ* insulator, and the  
35 RBS sites for both *lacI* and *tetA*). Of those 43 variants, 3 had  $EC_{50}$  values that differed by approximately  
36 2-fold or more from the geometric mean value for the wild-type  $EC_{50}$ . The Kolmogorov-Smirnov test was  
37 used to compare the distributions of  $EC_{50}$  values between the wild-type variants with and without  
38 mutations in the regulatory region; the results indicated a significant difference (p-value: 0.024).

39 To avoid biasing the results of the machine learning and other quantitative phenotypic analyses, variants  
40 were excluded from those analyses if they had one or more mutations in the non-*lacI* regions that show  
41 significant mutational effects: tetA, YFP, KAN, the origin of replication, and the regulatory region. After



1 applying this data quality filter in addition to those described above, there were 54,162 variants that we  
2 used for further quantitative analysis.

### 3 Library-scale fitness measurement

4 The experimental approach for this work was designed to maintain bacterial cultures in exponential  
5 growth phase for the full duration of the measurements. So, in all analysis, the Malthusian definition of  
6 fitness was used, i.e. fitness is the exponential growth rate<sup>35</sup>.

7 The fitness of cells containing each Lacl variant was calculated from the change in the relative  
8 abundance of DNA barcodes over time. The spike-in control was used to normalize the DNA barcode  
9 count data to enable the determination of the absolute fitness for each Lacl variant in the library.

10 Briefly, for each Lacl variant in each of the 24 chemical environments, the ratio of the barcode read  
11 count to the spike-in read count was fit to a function assuming exponential growth and a lag in the onset  
12 of the fitness impact of tetracycline. The fitness associated with each variant in each of the 24 chemical  
13 environments was determined as a parameter in the corresponding least-squares fit as detailed below.

14 The barcode sequencing data was analyzed with a model based on the assumption that the number of  
15 cells containing each Lacl variant grows with an exponential expansion rate that is independent of all  
16 other variants. So, for each sample, at the end of the incubation cycle for Growth Plate  $j$ , the number of  
17 cells with Lacl variant  $i$  is:

$$N_{i,j} = \frac{N_{i,j-1}}{d} \exp(\mu_{i,j} \Delta t) \quad (1)$$

18 where,  $d$  ( $= 10$ ) is the dilution factor used in transferring the cell culture from Growth Plate  $j - 1$  to  
19 Growth Plate  $j$ ,  $\Delta t$  ( $\approx 165$  minutes) is the total incubation time for each growth plate (including time  
20 required for automated cell passaging), and  $\mu_{i,j}$  is the fitness (ie mean exponential growth rate) of cells  
21 with Lacl variant  $i$  in Growth Plate  $j$ .

22 For samples without tetracycline, the chemical composition of the media was the same for all growth  
23 plates, so the fitness is assumed to be constant,  $\mu_{i,j} = \mu_i^0$ , where  $\mu_i^0$  is the fitness associated with Lacl  
24 variant  $i$  in the absence of tetracycline. Consequently, the number of cells in each Growth Plate for  
25 samples grown without tetracycline is:

$$\log(N_{i,j}^0) = \log(N_{i,0}^0) + j(\mu_i^0 \Delta t - \log(d)) \quad (2)$$

26 where  $N_{i,j}^0$  is the number of cells with Lacl variant  $i$  at the end of Growth Plate  $j$  for samples grown  
27 without tetracycline.

28 For samples grown with tetracycline, the tetracycline was only added to the culture media for Growth  
29 Plates 2-4. Because of the mode of action of tetracycline (inhibition of translation), there was a lag in its  
30 effect on cell fitness. Accordingly, the analysis assumes fitness varies as a function of time:

$$\mu_{i,j} = \mu_i^0 + (\mu_i^{tet} - \mu_i^0) e^{-\alpha j} \quad (3)$$

31 where  $\mu_i^{tet}$  is the steady-state fitness with tetracycline, and  $\alpha$  is a transition rate. Based on test  
32 measurements with a small-scale library, the transition rate was kept fixed at  $\alpha = \log(5)$ . From Eq. (3),  
33 the number of cells in each Growth Plate for samples grown with tetracycline is:

$$\log(N_{i,j}^{tet}) = \log(N_{i,0}^{tet}) + j(\mu_i^{tet}\Delta t - \log(d)) + \frac{\Delta t}{\alpha}(\mu_i^0 - \mu_i^{tet} + (\mu_i^{tet} - \mu_i^0)e^{-\alpha j}) \quad (4)$$

1 The barcode read count for variant  $i$  in Growth Plate  $j$  was assumed to be proportional to the cell  
2 number:

$$R_{i,j} = a_i b_j N_{i,j} \quad (5)$$

3 where  $a_i$  is a proportionality constant associated with variant  $i$ , and  $b_j$  is a proportionality constant  
4 associated with Growth Plate  $j$ . The proportionality constant  $a_i$  can be different for each variant  $i$  due to  
5 differences in PCR amplification efficiency resulting from variations in the barcode sequences on each  
6 amplicon. Similarly, the proportionality constant  $b_j$  can be different for each Growth Plate because of  
7 sample-to-sample variations in the DNA extraction efficiency or differences in PCR efficiency associated  
8 with different sample multiplexing tag sequences.

9 The logarithm of the read count normalized by the spike-in read count was used to estimate the fitness  
10 of each variant from its associated barcode read count:

$$\log(r_{i,j}) \equiv \log\left(\frac{R_{i,j}}{R_{spike,j}}\right) \quad (6)$$

11 For samples without tetracycline,  $\mu_i^0$  was estimated for each variant using a weighted linear least-  
12 squares fit to the log-count ratio vs.  $j$ :

$$\log(r_{i,j}^0) = \log(r_{i,0}^0) + j\Delta\mu_i^0\Delta t \quad (7)$$

13 where  $r_{i,j}^0 \equiv \frac{a_i}{a_{spike}} \frac{N_{i,0}^0}{N_{spike,0}^0}$ , and  $\Delta\mu_i^0 \equiv \mu_i^0 - \mu_{spike}^0$  is the difference between the fitness of variant  $i$  and the  
14 spike-in fitness without tetracycline.

15 For samples grown with tetracycline,  $\mu_i^{tet}$  was estimated for each variant with a weighted least-squares  
16 fit to the non-linear form for the log-count ratio:

$$\log(r_{i,j}^{tet}) = \log(r_{i,0}^{tet}) + j\Delta\mu_i^{tet}\Delta t + \frac{\Delta t}{\alpha}(\Delta\mu_i^0 - \Delta\mu_i^{tet} + (\Delta\mu_i^{tet} - \Delta\mu_i^0)e^{-\alpha j}) \quad (8)$$

17 where  $r_{i,j}^{tet} \equiv \frac{a_i}{a_{spike}} \frac{N_{i,0}^{tet}}{N_{spike,0}^{tet}}$ , and  $\Delta\mu_i^{tet} \equiv \mu_i^{tet} - \mu_{spike}^{tet}$  is the difference between the fitness of variant  $i$  and the  
18 spike-in fitness with tetracycline.

19 For the least-squares fits to determine both  $\mu_i^0$  and  $\mu_i^{tet}$ , the fits were weighted based on the propagated  
20 uncertainties of  $r_{i,j}^0$  and  $r_{i,j}^{tet}$  calculated assuming that the uncertainty of each read count was dominated  
21 by Poisson sampling.

22 For the fitness landscape measurement, there were a large number of outliers for the read count  
23 measurements from three of the samples: Growth Plate 3, without tetracycline, [IPTG] = 8  $\mu\text{mol/L}$ ;  
24 Growth Plate 4, without tetracycline, [IPTG] = 64  $\mu\text{mol/L}$  and [IPTG] = 2048  $\mu\text{mol/L}$ . These three samples  
25 were excluded from the analysis.

## 1 Dose-response curve measurements

2 Plasmids pTY1 and pVER were engineered to provide two independent measurements of the dose-  
3 response curve for LacI variants. First, in pTY1, LacI regulates the expression of a tetracycline resistance  
4 gene (*tetA*) that enables determination of the dose-response from barcode sequencing data by  
5 comparing the fitness measured with tetracycline to the fitness measured without tetracycline. Second,  
6 in pVER, the LacI regulates the expression of a fluorescent protein (YFP) that enables direct  
7 measurement of the dose-response curve with flow cytometry.

8 A set of nine randomly selected LacI variants were used to calibrate the estimation of regulated gene  
9 expression output from the barcode-sequencing fitness measurements (Supplementary Fig. 16). The  
10 calibration data consisted of the fitness data for each calibration variant from the library barcode  
11 sequencing measurement (using the library plasmid, pTY1) and flow cytometry data for each calibration  
12 variant prepared as a clonal culture (using the verification plasmid, pVER). This data was fit to a Hill  
13 equation model for the fitness impact of tetracycline as a function of the regulated gene expression  
14 level,  $G$ :

$$\frac{\mu^{tet}}{\mu^0} - 1 = \Delta f \left( \frac{G^{n_f}}{G_{50}^{n_f} + G^{n_f}} - 1 \right) \quad (9)$$

15 where  $\mu^{tet}$  is the fitness with tetracycline,  $\mu^0$  is the fitness without tetracycline,  $\Delta f$  is the maximal fitness  
16 impact of tetracycline (when  $G = 0$ ),  $G_{50}$  is the gene expression level that produces a 50% recovery in  
17 fitness, and  $n_f$  characterizes the steepness of the fitness calibration curve. Because the fitness calibration  
18 curve, Eq. (9), is nonlinear, it cannot be directly inverted to give the regulated gene expression level for  
19 all possible fitness measurements. So, two Bayesian inference models were used to estimate the dose-  
20 response curves for every LacI variant in the library using the barcode sequencing fitness measurements.  
21 Source code for both models is included in the software archive at  
22 [https://github.com/djross22/nist\\_lacI\\_landscape\\_analysis](https://github.com/djross22/nist_lacI_landscape_analysis). Both inference models used Eq. (9) to  
23 represent the relationship between fitness and regulated gene expression. The parameters  $\Delta f$ ,  $G_{50}$ , and  
24  $n_f$  were included in both inference models as parameters with informative priors. Priors for  $G_{50}$  and  $n_f$   
25 were based on the results of the fit to the fitness calibration data (Supplementary Fig. 16:  
26  $G_{50} \sim \text{normal}(\text{mean}=13,330, \text{std}=500)$ ,  $n_f \sim \text{normal}(\text{mean}=3.24, \text{std}=0.29)$ ). We chose the prior for  $\Delta f$   
27 based on an examination of  $\mu^{tet}/\mu^0 - 1$  measured with zero IPTG:  $\Delta f \sim \text{exponentially-modified-}$   
28  $\text{normal}(\text{mean}=0.720, \text{std}=0.015, \text{rate}=14)$ . The use of a prior for  $\Delta f$  with a broad right-side tail was  
29 important to accommodate variants in the library for which  $\mu^{tet}/\mu^0 - 1$  was systematically less  
30 than -0.722.

31 The first Bayesian inference model assumed that the dose-response curve for each LacI variant was  
32 described by the Hill equation. The Hill equation parameters for each variant,  $G_\infty$ ,  $G_0$ ,  $EC_{50}$ , and  $n$  and  
33 their associated uncertainties were determined using Bayesian parameter estimation by Markov Chain  
34 Monte Carlo (MCMC) sampling with PyStan<sup>36</sup>. Broad, flat priors were used for  $\log_{10}(G_0)$ ,  $\log_{10}(G_\infty)$ , and  
35  $\log_{10}(EC_{50})$ , with error function boundaries to constrain those parameter estimates to within the  
36 measurable range ( $100 \text{ MEF} \leq G_0$ ,  $G_\infty \leq 50,000 \text{ MEF}$ ;  $0.1 \mu\text{mol/L} \leq EC_{50,i} \leq 40,000 \mu\text{mol/L}$ ). The prior for  $n_i$   
37 was a gamma distribution with shape parameter of 4.0 and inverse scale parameter of 3.33. The  
38 inference model was run individually for each LacI variant, with four independent chains, 1000 iterations  
39 per chain (500 warmup iterations), and the `adapt_delta` parameter set to 0.9. Testing with data from a  
40 set of randomly selected variants indicated that these settings for the Stan sampling algorithm typically

1 produced a Gelman-Rubin  $\hat{R}$  diagnostic less than 1.05 and number of effective iterations greater than  
2 100.

3 The second Bayesian inference model was a non-parametric Gaussian process (GP) model<sup>37</sup> that  
4 assumed only that the dose-response curve for each *Lacl* variant was a smooth function of IPTG  
5 concentration. The GP model was used to determine which variants had band-pass or band-stop  
6 phenotypes. The GP model was also implemented using MCMC sampling with PyStan<sup>36</sup>. The GP  
7 inference model was run individually for each variant, with four independent chains, 1000 iterations per  
8 chain (500 warmup iterations), and the `adapt_delta` parameter set to 0.9. Testing with data from a set of  
9 randomly selected variants indicated that these settings for the Stan sampling algorithm of the GP  
10 model typically produced a Gelman-Rubin  $\hat{R}$  diagnostic less than 1.02 and number of effective iterations  
11 greater than 200.

## 12 Flow cytometry measurements

13 Over 100 *Lacl* variants from the library were chosen for flow cytometry verification of the dose-response  
14 curves. The CDSs of these variants were chemically synthesized (Twist Bioscience), cloned into the  
15 verification plasmid, pVER, and then transformed into MG1655 $\Delta lac$ . Transformants were plated in LB  
16 supplemented with kanamycin and 0.2% glucose. *Lacl* variant sequences were verified with Sanger  
17 sequencing (Psomagen USA). For flow cytometry measurements of dose-response curves, a culture of  
18 *E. coli* containing pVER with a chosen variant sequence was distributed across 12 wells of a 96-well plate  
19 and grown to stationary phase using the automated microbial growth system. After growth to stationary  
20 phase, cultures were diluted 50-fold into a plate containing the same 12 IPTG concentrations used  
21 during the fitness landscape measurement (0  $\mu\text{mol/L}$  to 2048  $\mu\text{mol/L}$ ). In some cases, higher IPTG  
22 concentrations were used to capture the full dose-response curves of selected variants (e.g.  
23 Supplementary Figs. 4-5). Cultures were then grown for 160 minutes ( $\sim 3.3$  generations) before being  
24 diluted 10-fold into the same IPTG gradient and grown for another 160 minutes. Then, 5  $\mu\text{L}$  of each  
25 culture was diluted into 195  $\mu\text{L}$  of PBS supplemented with 170  $\mu\text{g/mL}$  chloramphenicol and incubated at  
26 room temperature for 30-60 minutes to halt the translation of YFP and allow extant YFP to mature in the  
27 cells.

28 Samples were measured on an Attune NxT flow cytometry with autosampler using a 488 nm excitation  
29 laser and a 530 nm  $\pm$  15 nm bandpass emission filter. Blank samples were measured with each batch of  
30 cell measurements, and an automated gating algorithm was used to discriminate cell events from non-  
31 cell events (Supplementary Fig. 17a-b). With the Attune cytometer, the area and height parameters for  
32 each detection channel are calibrated to give the same value for singlet events. So, to identify singlet  
33 cell events and exclude multiplet cell events, a second automated gating algorithm was applied to select  
34 only cells with side scatter area  $\cong$  side scatter height (Supplementary Fig. 17c-d). All subsequent analysis  
35 was performed using the singlet cell event data. Fluorescence data was calibrated to molecules of  
36 equivalent fluorophore (MEF) using fluorescent calibration beads (Spherotech, part no. RCP-30-20A).  
37 The cytometer was programmed to measure a 25  $\mu\text{L}$  portion of each cell sample, and the 40-fold  
38 dilution used in the cytometry sample preparation resulted in approximately 20,000 singlet cell  
39 measurements per sample. The geometric mean of the YFP fluorescence was used as a summary  
40 statistic to represent the regulated gene expression level as a function of the input ligand concentration,  
41 [IPTG] for each *Lacl* variant.

## 1 Calculation of abundance for *Lacl* phenotypes

2 The relative abundance of the various *Lacl* phenotypes (Supplementary Fig. 2) was estimated using the  
3 results of both Bayesian inference models (Hill equation and GP). Variants were labeled as “flat  
4 response” if the Hill equation model and the GP model agreed (i.e. if the median estimate for the Hill  
5 equation dose-response curve was within the central 90% credible interval from the GP model at all 12  
6 IPTG concentrations) and if the posterior probability for  $G_0 > G_\infty$  was between 0.05 and 0.95 (from the  
7 Hill equation model inference). Variants were labeled as having a negative response if the slope,  $\partial G/\partial L$ ,  
8 was negative at one or more IPTG concentrations with 0.95 or higher posterior probability (from the GP  
9 model inference). To avoid false positives from end effects, this negative slope criteria was only applied  
10 for IPTG concentrations between 2  $\mu\text{mol/L}$  IPTG and 1024  $\mu\text{mol/L}$ . Variants were labeled as “always on”  
11 (the  $I^-$  phenotype from reference<sup>18</sup>) if they were flat-response and if  $G(0)$  was greater than 0.25 times  
12 the wild-type  $G_\infty$  value with 0.95 or higher posterior probability (from the GP model inference). Variants  
13 were labeled as “always off” (the  $I^S$  phenotype from reference<sup>18</sup>) if they were flat-response but not  
14 always on. Variants were labeled as band-stop or band-pass if the slope,  $\partial G/\partial L$ , was negative at some  
15 IPTG concentrations and positive at other IPTG concentrations, both with 0.95 or higher posterior  
16 probability (from the GP model inference). Band-stop and band-pass variants were distinguished by the  
17 ordering of the negative-slope and positive-slope portions of the dose-response curves. Variants that  
18 had a negative response but that were not band-pass or band-stop, were labeled as inverted. False-  
19 positive rates were estimated for each phenotypic category by manually examining the fitness vs. IPTG  
20 data for *Lacl* variants with less than three substitutions. Typical causes of false-positive phenotypic  
21 labeling included unusually high noise in the fitness measurement and biased fit results due to outlier  
22 fitness data points. Estimated false-positive rates ranged between 0.001 and 0.005. The relative  
23 abundance values shown in Supplementary Fig. 2a were corrected for false positives using the estimated  
24 rates.

## 25 Comparison of synonymous mutations

26 The library contained a set of 39 variants with the wild-type *lacl* CDS (but different DNA barcodes), and a  
27 set of 310 variants with only synonymous nucleotide changes (i.e. no amino acid substitutions). Both  
28 sets had long-read sequencing coverage for the entire plasmid and were screened to retain only variants  
29 with zero unintended mutations in the plasmid (i.e. no mutations in regions of the plasmid other than  
30 the *lacl* CDS). The Hill equation fit results for those two sets were compared to determine whether  
31 synonymous nucleotide changes significantly affected the phenotype. The Kolmogorov-Smirnov test was  
32 used to compare the distributions of Hill equation parameters between these two sets. The resulting p-  
33 values (0.71, 0.40, 0.28, and 0.17 for  $G_0$ ,  $G_\infty$ ,  $EC_{50}$ , and  $n$  respectively) indicate that there were no  
34 significant differences between them. Additionally, the library contained 40 sets of variants, each with  
35 four or more synonymous CDSs (including the set of synonymous wild-type sequences and 39 non-wild-  
36 type sequences). A hierarchical model was used to compare the Hill equation parameters within each  
37 set of synonymous CDSs. Within each set, the uncertainty associated with individual variants was  
38 typically larger than the variant-to-variant variability estimated by the hierarchical model. Overall, these  
39 results indicate that synonymous SNPs did not measurably impact the *Lacl* phenotype, so only the amino  
40 acid sequences were considered for any subsequent quantitative genotype-to-phenotype analysis.

## 1 Analysis of single-substitution data

2 The single amino acid substitution results presented in Fig. 2, Fig. 5b-c, Supplementary Fig. 13, and  
3 included in Supplementary Data 1 are a combination of direct experimental observations, DNN model  
4 results, and estimates of  $G_0$  for missing substitutions.

5 For direct experimental observations, multiple LacI variants were often present in the library with the  
6 same single substitution. To ensure that the highest quality data was used for the single-substitution  
7 analysis, only data for variants with more than 5000 total barcode reads were used (see Barcode  
8 sequencing section, above). For each single substitution, if there was only one LacI variant with more  
9 than 5000 barcode reads, the median and standard deviation for each parameter were used directly  
10 from the Bayesian inference using the Hill equation model. If there was more than one LacI variant with  
11 a given single substitution and more than 5000 barcode reads, the consensus Hill equation parameter  
12 values and standard deviations for that substitution were calculated using a hierarchical model based on  
13 the eight schools model<sup>38,39</sup>. The hierarchical model was applied separately for each Hill equation  
14 parameter. The logarithm of the parameter values was used as input to the hierarchical model, and the  
15 input data were centered and normalized by  $1.15 \times$  the minimum measurement uncertainty. The  
16 standard normal distribution was used as a loosely informative prior for the consensus mean effect, and  
17 a half-normal prior (mean = 0.5, std = 1) was used for the normalized consensus standard deviation (i.e.  
18 hierarchical standard deviation). These priors and normalization were chosen so that the model gave  
19 intuitively reasonable results for the consensus of two LacI variants (i.e. close to the results for the LacI  
20 variant with the lowest measurement uncertainty). Results for the hierarchical model were determined  
21 using Bayesian parameter estimation by Markov Chain Monte Carlo (MCMC) sampling with PyStan<sup>36</sup>.  
22 MCMC sampling was run with 4 independent chains, 10,000 iterations per chain (5,000 warmup  
23 iterations), and the `adapt_delta` parameter set to 0.975.

24 For  $G_0$ , the direct experimental results were used for the 1047 substitutions plotted as gray points or red  
25 points and error bars in Fig. 2d and Supplementary Fig. 13. In addition, estimated values were used for  
26 the 83 missing substitutions that have been previously shown to result in an “always on” LacI phenotype  
27 (i.e., the  $I^-$  phenotype<sup>18,19</sup>). For these substitutions, plotted as pink-gray points and error bars in Fig. 2d,  
28 the median value was estimated to be equal to the wild-type value for  $G_\infty$  (24,000 MEF), and the  
29 geometric standard deviation was estimated to be 4-fold, both based on information from previous  
30 publications<sup>18,19</sup>. Note that these 83 substitutions are completely missing from the experimental  
31 landscape dataset, i.e. they are not found in any LacI variant, as single substitutions or in combination  
32 with other substitutions.

33 For  $G_\infty$  and  $EC_{50}$ , the direct experimental results were used for the 964 substitutions that are found as  
34 single substitutions in the library and that have a consensus standard deviation for  $\log_{10}(EC_{50})$  less than  
35 0.35. An additional 74 substitutions are found as single substitutions in the library, but with higher  $EC_{50}$   
36 uncertainty. For these substitutions, either  $EC_{50}$  is comparable to or higher than the maximum ligand  
37 concentration used for the measurement (2048  $\mu\text{mol/L}$  IPTG), or  $G_\infty$  is comparable to  $G_0$  (or both).  
38 Consequently, the dose-response curve is flat or nearly flat across the range of concentrations used, and  
39 the Bayesian inference used to estimate the Hill equation parameters results in  $EC_{50}$  and  $G_\infty$  estimates  
40 with large uncertainties. The DNN model can provide a better parameter estimate for these flat-  
41 response variants because it uses data and relationships from the full library (e.g. the log-additivity  
42 of  $EC_{50}$ ) to predict parameter values for each single substitution. So, the DNN model results were used

1 for these 74 substitutions. Finally, the DNN model results were used for an additional 953 substitutions  
2 that are found in the library, but only in combination with other substitutions (i.e. not as single  
3 substitutions).

#### 4 Identification of high-frequency substitutions and structural features associated with 5 inverted and band-stop phenotypes

6 The set of 43 strongly inverted Lacl variants discussed in the main text and used for the plots in Fig. 4a,c  
7 were identified by the following criteria:  $G_0/G_\infty \geq 2$ ,  $G_0 > G_{\infty,wt}/2$ ,  $G_\infty < G_{\infty,wt}/2$ , and  $EC_{50}$  between  
8  $3 \mu\text{mol/L}$  and  $1000 \mu\text{mol/L}$ . The set of 31 strong band-stop variants discussed in the main text and used  
9 for the plots in Fig. 4b,d were identified by the following criteria:  $G_0 > G_{\infty,wt}/2$ ,  $G_{min} < G_{\infty,wt}/2$ , and the  
10 slope,  $\partial \log(G)/\partial \log(L)$ , of less than  $-0.07$  at low IPTG concentrations and greater than zero at higher IPTG  
11 concentrations, both with 0.95 or higher posterior probability (from the GP model inference). In  
12 addition, the sets of strongly inverted and strong band-stop variants were manually screened for likely  
13 false positives due to outlier fitness data points.

14 A hypergeometric test was used to determine the amino acid substitutions that occur more frequently  
15 in the set of strongly inverted or strong band-stop variants than in the full library. For each possible  
16 substitution, the cumulative hypergeometric distribution was used to calculate the probability of the  
17 observed number of occurrences of that substitution in the set of inverted or band-stop variants under a  
18 null model of no association. This probability was used as a p-value for the null hypothesis that the  
19 observed number of inverted or band-stop variants with that substitution resulted from an unbiased  
20 random selection of variants from the full library. Substitutions were considered to occur at significantly  
21 higher frequency if they had a p-value less than 0.005 and if they occurred more than once in the set of  
22 inverted or band-stop variants. In the set of strongly inverted variants, ten associated (higher frequency)  
23 amino acid substitutions were identified: S70I, K84N, D88Y, V96E, A135T, V192A, G200S, Q248H, Y273H,  
24 and A343G. In the set of strong band-stop variants, eight associated substitutions were identified: V4A,  
25 A92V, H179Q, R195H, G178D, G265D, D292G, and R351G. To estimate the number of false-positives,  
26 random sets of Lacl variants were chosen with the same sample size as the strongly inverted (43) or the  
27 strong band-stop (31) variants and the same significance criteria was applied. From 300 independent  
28 iterations of the random selection, the estimated mean number of false-positive substitutions was 2.1  
29 and 2.3 for the inverted and band-stop phenotypes, respectively.

30 A similar procedure was used to determine which structural features within the protein are mutated  
31 with higher frequency in the inverted or band-stop Lacl variants. The structural features considered  
32 were the secondary structures from the complete crystal structure of Lacl<sup>22</sup>, as well as larger structural  
33 features (N-terminal core domain, C-terminal core domain, DNA-binding domain, dimer interface) and  
34 functional domains (ligand-binding, core-pivot). The p-value threshold used for significance was 0.025.  
35 For the strongly inverted variants, six domains were identified with a higher frequency of amino acid  
36 substitutions: the dimer interface, residues within 7 Å of the ligand-binding pocket, helix 5, helix 11,  
37 strand I, and the N-terminal core. For the strong band-stop variants, three features were identified: the  
38 C-terminal core, strand J, and helix 9. From 300 independent random selections of variants from the full  
39 library, the estimated mean number of false-positive features was 0.39 and 0.50 for the inverted and  
40 band-stop phenotypes, respectively.

## 1 Deep neural network (DNN) modeling

2 The dataset was pruned to a set of high-quality sequences for DNN modeling. Specifically, data for a *LacI*  
3 variant was only used for modeling if it satisfied the following criteria:

- 4 1. No mutations were found in the long-read sequencing results for the regions of the plasmid  
5 encoding kanamycin resistance, the origin of replication, the *tetA* and YFP genes, and the  
6 regulatory region containing the promoters and ribosomal binding sites for *lacI* and *tetA*  
7 (Supplementary Table 1).
- 8 2. The total number of barcode read counts for a *LacI* variant was greater than 3000.
- 9 3. The number of amino acid substitutions was less than 14.
- 10 4. The measurement uncertainty for  $\log_{10}(G_{\infty})$  was less than 0.7.
- 11 5. The results of the Hill equation model and the GP model agreed at all 12 IPTG concentrations.  
12 More specifically, data were only used if the median estimate for the dose-response curve from  
13 the Hill equation model was within the central 90% credible interval from the GP model at all  
14 12 IPTG concentrations.

15 After applying the quality criteria listed above, 47,462 *LacI* variants remained for DNN modeling. The  
16 data were used to train the DNN model to predict the Hill equation parameters  $G_0$ ,  $G_{\infty}$ , and  $EC_{50}$  as  
17 detailed below.

18 Amino acid sequences were represented as one-hot encoded vectors of length  $L = 2536$ , and with  
19 mutational paths represented as  $K \times L$  tensors for a sequence with  $K$  substitutions. The logarithm of the  
20 Hill equation parameter values were normalized to a standard deviation of 1, and then shifted by the  
21 corresponding value of the wild-type sequence in order to correctly represent the prediction goal of the  
22 change in each parameter relative to wild-type *LacI*. A long-term, short-term recurrent neural network  
23 was selected for the underlying model<sup>17</sup>, with 16 hidden units, a single hidden layer, and hyperbolic  
24 tangent (tanh) non-linearities. Inference was performed in pytorch<sup>40</sup> using the Adam optimizer<sup>41</sup>. For  
25  $EC_{50}$  and  $G_0$ , the contribution of individual data points to the regression loss were weighted inversely  
26 proportional to their experimental uncertainty. Model selection was performed with 10-fold cross-  
27 validation on the training set (80% of all available data). Approximate Bayesian inference was performed  
28 with the Bayes-by-backprop approach<sup>42</sup>. Briefly, this substitutes the point-estimate parameters of the  
29 neural network with variational approximations to a Bayesian model, represented as a mean and  
30 variance of a normal random variable. Effectively, this only doubles the number of parameters in the  
31 model. A mixture of two normal distributions was used as a prior for each parameter weight, with the  
32 two mixture components having high and low variance respectively. This prior emulates a sparsifying  
33 spike-slab prior while remaining tractable for inference based on back-propagation. Posterior means of  
34 each weight were used to calculate posterior predictive means, while Monte-Carlo draws from the  
35 variational posterior were used to calculate the model prediction uncertainty (Supplementary Fig. 10).

36 Variational approximations typically underestimate uncertainty. So, to correct the uncertainty  
37 estimates, the model prediction uncertainty obtained from the variational approximation was compared  
38 to the model root-mean-square error (RMSE) (i.e. the root-mean-square difference between the model  
39 prediction and the experimental measurement). For all three Hill equation parameters ( $G_0$ ,  $G_{\infty}$ , and  
40  $EC_{50}$ ), both the prediction uncertainty and the RMSE increase with the number of amino acid



1 substitutions relative to wild-type sequence (Supplementary Fig. 10a-b), and the RMSE at each  
2 substitutional distance is an approximately linear function of the median model uncertainty  
3 (Supplementary Fig. 10c). So, for the single-substitution analysis (Fig. 2, Fig. 5b-c, Supplementary Fig. 13,  
4 Supplementary Data 1), the uncertainties from the variational approximation were multiplied by a factor  
5 of 3.8. This rescaled the uncertainties so that the median uncertainty was approximately equal to the  
6 RMSE for each substitutional distance.

## 7 Data Availability

8 The raw sequence data for long-read and short-read DNA sequencing have been deposited in the NCBI  
9 Sequence Read Archive and are available under the project accession number PRJNA643436. Plasmid  
10 sequences have been deposited in the NCBI Genbank under accession codes MT702633, and MT702634,  
11 for pTY1 and pVER, respectively.

12 The processed data table containing information for each LacI variant in the library is publicly available  
13 via the NIST Science Data Portal, with the identifier ark:/88434/mds2-2259  
14 (<https://data.nist.gov/od/id/mds2-2259> or <https://doi.org/10.18434/M32259>).

## 15 Code Availability

16 All custom data analysis code is available at [https://github.com/diross22/nist\\_lacl\\_landscape\\_analysis](https://github.com/diross22/nist_lacl_landscape_analysis).

## 17 References

- 18 1. Motlagh, H. N., Wrabl, J. O., Li, J. & Hilser, V. J. The ensemble nature of allostery. *Nature* **508**, 331–  
19 339 (2014).
- 20 2. Razo-Mejia, M. *et al.* Tuning Transcriptional Regulation through Signaling: A Predictive Theory of  
21 Allosteric Induction. *Cell Systems* **6**, 456-469.e10 (2018).
- 22 3. Fenton, A. W. Allostery: an illustrated definition for the ‘second secret of life’. *Trends Biochem. Sci.*  
23 **33**, 420–425 (2008).
- 24 4. Raman, S., Taylor, N., Genuth, N., Fields, S. & Church, G. M. Engineering Allostery. *Trends Genet* **30**,  
25 521–528 (2014).
- 26 5. Huang, P.-S., Boyken, S. E. & Baker, D. The coming of age of de novo protein design. *Nature* **537**,  
27 320–327 (2016).
- 28 6. He, X. & Liu, L. Toward a prospective molecular evolution. *Science* **352**, 769–770 (2016).
- 29 7. Sarkisyan, K. S. *et al.* Local fitness landscape of the green fluorescent protein. *Nature* **533**, 397–401  
30 (2016).
- 31 8. Puchta, O. *et al.* Network of epistatic interactions within a yeast snoRNA. *Science* **352**, 840–844  
32 (2016).
- 33 9. Li, C., Qian, W., Maclean, C. J. & Zhang, J. The fitness landscape of a tRNA gene. *Science* **352**, 837–  
34 840 (2016).
- 35 10. Pressman, A. D. *et al.* Mapping a Systematic Ribozyme Fitness Landscape Reveals a Frustrated  
36 Evolutionary Network for Self-Aminoacylating RNA. *J. Am. Chem. Soc.* **141**, 6213–6223 (2019).

- 1 11. Domingo, J., Diss, G. & Lehner, B. Pairwise and higher-order genetic interactions during the  
2 evolution of a tRNA. *Nature* **558**, 117–121 (2018).
- 3 12. Li, C. & Zhang, J. Multi-environment fitness landscapes of a tRNA gene. *Nat Ecol Evol* **2**, 1025–1032  
4 (2018).
- 5 13. Monod, J., Wyman, J. & Changeux, J. P. ON THE NATURE OF ALLOSTERIC TRANSITIONS: A PLAUSIBLE  
6 MODEL. *J. Mol. Biol.* **12**, 88–118 (1965).
- 7 14. Chure, G. *et al.* Predictive shifts in free energy couple mutations to their phenotypic consequences.  
8 *PNAS* **116**, 18275–18284 (2019).
- 9 15. Daber, R., Sochor, M. A. & Lewis, M. Thermodynamic analysis of mutant lac repressors. *J. Mol. Biol.*  
10 **409**, 76–87 (2011).
- 11 16. Wenger, A. M. *et al.* Accurate circular consensus long-read sequencing improves variant detection  
12 and assembly of a human genome. *Nature Biotechnology* **37**, 1155–1162 (2019).
- 13 17. Hochreiter, S. & Schmidhuber, J. Long Short-Term Memory. *Neural Comput.* **9**, 1735–1780 (1997).
- 14 18. Markiewicz, P., Kleina, L. G., Cruz, C., Ehret, S. & Miller, J. H. Genetic Studies of the lac Repressor.  
15 XIV. Analysis of 4000 Altered Escherichia coli lac Repressors Reveals Essential and Non-essential  
16 Residues, as well as ‘Spacers’ which do not Require a Specific Sequence. *Journal of Molecular*  
17 *Biology* **240**, 421–433 (1994).
- 18 19. Pace, H. C. *et al.* Lac repressor genetic map in real space. *Trends in Biochemical Sciences* **22**, 334–  
19 339 (1997).
- 20 20. Flynn, T. C. *et al.* Allosteric transition pathways in the lactose repressor protein core domains:  
21 Asymmetric motions in a homodimer. *Protein Sci* **12**, 2523–2541 (2003).
- 22 21. Perturbation from a distance: mutations that alter LacI function through long-range effects -  
23 PubMed. <https://pubmed.ncbi.nlm.nih.gov/14636069/>.
- 24 22. Lewis, M. *et al.* Crystal structure of the lactose operon repressor and its complexes with DNA and  
25 inducer. *Science* **271**, 1247–1254 (1996).
- 26 23. Hecht, A. *et al.* Measurements of translation initiation from all 64 codons in E. coli. *Nucleic Acids Res*  
27 **45**, 3615–3626 (2017).
- 28 24. Myers, G. L. & Sadler, J. R. Mutational inversion of control of the lactose operon of Escherichia coli.  
29 *Journal of Molecular Biology* **58**, 1–28 (1971).
- 30 25. Rolfes, R. J. & Zalkin, H. Purification of the Escherichia coli purine regulon repressor and  
31 identification of corepressors. *Journal of Bacteriology* **172**, 5637–5642 (1990).
- 32 26. Sarkar, S., Tack, D. & Ross, D. Sparse estimation of mutual information landscapes quantifies  
33 information transmission through cellular biochemical reaction networks. *Communications Biology*  
34 **3**, 1–8 (2020).
- 35 27. Tack, D. S., Alperovich, N., Vasilyeva, O. & Ross, D. Assembly of plasmids for LacI genotype-  
36 phenotype landscape measurement. *protocols.io* [https://www.protocols.io/view/assembly-of-](https://www.protocols.io/view/assembly-of-plasmids-for-laci-genotype-phenotype-l-bjxxkkpn)  
37 [plasmids-for-laci-genotype-phenotype-l-bjxxkkpn](https://www.protocols.io/view/assembly-of-plasmids-for-laci-genotype-phenotype-l-bjxxkkpn) doi:10.17504/protocols.io.bjxxkkpn.
- 38 28. Alperovich, N., Romantseva, J., Vasilyeva, O. & Ross, D. Automation Protocol for Plasmid DNA  
39 Extraction from E. coli. *protocols.io* [https://www.protocols.io/view/automation-protocol-for-](https://www.protocols.io/view/automation-protocol-for-plasmid-dna-extraction-fro-bjvkkn6)  
40 [plasmid-dna-extraction-fro-bjvkkn6](https://www.protocols.io/view/automation-protocol-for-plasmid-dna-extraction-fro-bjvkkn6) doi:10.17504/protocols.io.bjvkkn6.

- 1 29. Alperovich, N., Tack, D. S., Vasilyeva, O., Levy, S. F. & Ross, D. Automation Protocol for DNA Barcode  
2 Sequencing Library Preparation. *protocols.io* [https://www.protocols.io/view/automation-protocol-](https://www.protocols.io/view/automation-protocol-for-dna-barcode-sequencing-lib-bjjzkkp6)  
3 [for-dna-barcode-sequencing-lib-bjjzkkp6](https://www.protocols.io/view/automation-protocol-for-dna-barcode-sequencing-lib-bjjzkkp6) doi:10.17504/protocols.io.bjjzkkp6.
- 4 30. Zhao, L., Liu, Z., Levy, S. F. & Wu, S. Bartender: a fast and accurate clustering algorithm to count  
5 barcode reads. *Bioinformatics* **34**, 739–747 (2018).
- 6 31. Kivioja, T. *et al.* Counting absolute numbers of molecules using unique molecular identifiers. *Nature*  
7 *Methods* **9**, 72–74 (2012).
- 8 32. Schlecht, U., Liu, Z., Blundell, J. R., St.Onge, R. P. & Levy, S. F. A scalable double-barcode sequencing  
9 platform for characterization of dynamic protein-protein interactions. *Nature Communications* **8**,  
10 15586 (2017).
- 11 33. Omelina, E. S., Ivankin, A. V., Letiagina, A. E. & Pindyurin, A. V. Optimized PCR conditions minimizing  
12 the formation of chimeric DNA molecules from MPRA plasmid libraries. *BMC Genomics* **20**, 536  
13 (2019).
- 14 34. Smyth, R. P. *et al.* Reducing chimera formation during PCR amplification to ensure accurate  
15 genotyping. *Gene* **469**, 45–51 (2010).
- 16 35. Wu, B., Gokhale, C. S., van Veelen, M., Wang, L. & Traulsen, A. Interpretations arising from  
17 Wrightian and Malthusian fitness under strong frequency dependent selection. *Ecol Evol* **3**, 1276–  
18 1280 (2013).
- 19 36. Carpenter, B. *et al.* Stan: A Probabilistic Programming Language. *Journal of Statistical Software* **76**,  
20 1–32 (2017).
- 21 37. Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning) | Guide  
22 books. <https://dl.acm.org/doi/book/10.5555/1162254>.
- 23 38. Rubin, D. B. Estimation in Parallel Randomized Experiments. *Journal of Educational Statistics* **6**, 377–  
24 401 (1981).
- 25 39. Diagnosing Biased Inference with Divergences. [https://mc-stan.org/users/documentation/case-](https://mc-stan.org/users/documentation/case-studies/divergences_and_bias.html)  
26 [studies/divergences\\_and\\_bias.html](https://mc-stan.org/users/documentation/case-studies/divergences_and_bias.html).
- 27 40. Paszke, A. *et al.* PyTorch: An Imperative Style, High-Performance Deep Learning Library.  
28 *arXiv:1912.01703 [cs, stat]* (2019).
- 29 41. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. *arXiv:1412.6980 [cs]* (2017).
- 30 42. Blundell, C., Cornebise, J., Kavukcuoglu, K. & Wierstra, D. Weight Uncertainty in Neural Networks.  
31 *arXiv:1505.05424 [cs, stat]* (2015).

32

## 33 Acknowledgements

34 We would like to thank Vanya Paralanov, Daniel Samarov, Ben Scott, Zvi Kelman, Gilad Kusne, and  
35 Swarnavo Sarkar for thoughtful discussions during planning and execution of this work. We would also  
36 like to thank Jayan Rammohan, William Brad O’Dell, and Elizabeth Strychalski for insights during the  
37 experimental work, as well as improving the manuscript.

## 1 Author Contributions

2 D.S.T., and D.R. conceived of the process.

3 D.S.T, S.L., and D.R. developed the experimental workflow.

4 D.S.T. designed, built, and tested genetic constructs.

5 E.F.R., and D.R. programmed automated protocols.

6 D.S.T., E.F.R., N.A., O.V., and D.R. performed landscape and verification experiments.

7 P.D.T. and D.R. performed Bayesian inference and model fitting.

8 P.D.T. designed and evaluated the recurrent architecture for machine learning.

9 P.D.T., N.D.O, and D.R. contributed to long-read sequencing analysis.

10 D.S.T., P.D.T, A.P., and D.R. wrote the manuscript.

11 All authors contributed to the manuscript.

## 12 Supplementary information

13 Supplementary Information: This file includes 17 supplementary figures and 1 supplementary table.

14 Supplementary Data 1: Single-substitution Hill equation parameters. The table contains the estimated  
15 Hill equation parameter for all of the single-substitution LaCl variants analyzed. The values listed in the  
16 table are the base-10 logarithm for each parameter. The column headings indicate the parameter and  
17 the source of the estimate as follows: "exp\_": experimental values, "dnn\_": values predicted by the DNN  
18 model, "est\_": values for missing substitutions estimated based on previously published results, "\_err":  
19 the uncertainty (1 standard deviation) of the  $\log_{10}(\text{parameter})$ . Column headings that start with "best\_"  
20 contain the values used for the analysis and plots contained in the manuscript. See Methods for more  
21 details.

## 22 Disclaimer

23 The authors declare no competing interests

24 Certain commercial equipment, instruments, or materials are identified to adequately specify  
25 experimental procedures. Such identification neither implies recommendation nor endorsement by the  
26 National Institute of Standards and Technology nor that the equipment, instruments, or materials  
27 identified are necessarily the best for the purpose.

28 Supplementary Information is available for this paper.

29 Correspondence and requests for materials should be addressed to David Ross ([david.ross@nist.gov](mailto:david.ross@nist.gov)).

30