

The simplicity of protein sequence-function relationships

Yeonwoo Park^{a,d}, Brian P.H. Metzger^{b,e}, and Joseph W. Thornton^{b,c,1}

This manuscript was compiled on September 3, 2023

How complicated is the relationship between a protein's sequence and its function? High-order epistatic interactions among residues are thought to be pervasive, making a protein's function difficult to predict or understand from its sequence. Most prior studies, however, used methods that misinterpret measurement errors, small local idiosyncracies around a designated wild-type sequence, and global nonlinearity in the sequence-function relationship as rampant high-order interactions. Here we present a simple new method to jointly estimate global nonlinearity and specific epistatic interactions across a protein's genotype-phenotype map. Our reference-free approach calculates the effect of each amino acid state or combination by averaging over all genotypes that contain it relative to the global average. We show that this method is more accurate than any alternative approach and is robust to measurement error and partial sampling. We reanalyze 20 combinatorial mutagenesis experiments and find that main and pairwise effects, together with a simple form of global nonlinearity, account for a median of 96% of total variance in the measured phenotype (and > 92% in every case), and only a tiny fraction of genotypes are strongly affected by epistasis at third or higher orders. The genetic architecture is also sparse: the number of model terms required to explain the vast majority of phenotypic variance is smaller than the number of genotypes by many orders of magnitude. The sequence-function relationship in most proteins is therefore far simpler than previously thought, and new, more tractable experimental approaches, combined with reference-free analysis, may be sufficient to explain it in most cases.

Sequence-function relationship | genetic architecture | epistasis | reference-free analysis

If we had a comprehensive understanding of a protein's sequence-function relationship, we could predict the functional and evolutionary consequences of any mutation or novel amino acid sequence. Whether such knowledge is possible in practice depends on the extent of epistatic interactions. If all residues in a protein act independently, then knowing the effects of point mutations on any genetic background would suffice to predict the function of any possible sequence, and any mutation's evolutionary fate would be independent of the genetic context in which it occurs. A simple genetic architecture like this could be easily inferred using moderate-throughput experiments. At the opposite extreme, extensive high-order epistasis would cause each mutation to have idiosyncratic effects that depend absolutely on the particular sequence background into which it is introduced. In that case, assessing the protein's genetic architecture would require exhaustive assessment of every possible genotype, and the evolutionary accessibility of all mutations would change with every sequence substitution that occurs.

Deep mutational scanning (DMS) methods to characterize large libraries of protein variants have recently made it possible to assess the complexity of the sequence-function relationship, but studies to date disagree on the complexity of the sequence-function relationship. Some report extensive high-order interactions (1–7), while others find that they account for less than 10% of functional variance among sequences (8–16). Even pairwise interactions are pervasive and strong in some studies (7, 12, 17–21) but sparse and weak in others (9, 16, 22). In terms of overall complexity, some report a sparse genetic architecture in which only a small fraction of possible terms are important (13, 13, 16, 16) but others point to a much more complex mapping in which many different states and combinations shape the sequence-function relationship (7, 18, 20, 21).

These discrepancies may reflect the use of different methods to characterize epistasis. Two aspects of widely used approaches can yield overestimates of amino acid interactions. First, most studies to date have analyzed mutational data using reference-based models, which designate a single sequence as the reference against which all effects are measured: the main effects of mutations are estimated by introducing each one into a single reference genotype, and epistatic

Significance Statement

It is widely thought that a protein's function depends on complex interactions among amino acids. If so, it would be virtually impossible to predict the function of new variants, and understanding how proteins work genetically and biochemically would require huge combinatorial experiments. We show that prior studies overestimated complexity because they analyzed sequence-function relationships from the perspective of a single reference genotype and/or misinterpreted global phenotypic nonlinearities as complex amino acid interactions. By developing a new reference-free approach and using it to reanalyze 20 experimental datasets, we show that additive effects and pairwise interactions alone, along with a simple global nonlinearity, explain the vast majority of functional variation. Higher-order interactions are weak or rare, and a minuscule fraction of possible interactions shape each protein's function. Our work reveals that protein sequence-function relationships are surprisingly simple and suggests new strategies that are far more tractable than the massive experiments currently used.

Author affiliations: ^aCommittee on Genetics, Genomics, and Systems Biology, University of Chicago, Chicago, IL 60637; ^bDepartment of Ecology and Evolution, University of Chicago, Chicago, IL 60637; ^cDepartment of Human Genetics, University of Chicago, Chicago, IL 60637; ^dCurrent affiliation: Center for RNA Research, Seoul National University, Seoul, Republic of Korea 08826; ^eCurrent affiliation: Department of Biological Sciences, Purdue University, West Lafayette, IN 47907

Y.P., B.P.H.M., and J.W.T. designed research; Y.P. developed methods and analyzed data; Y.P. and J.W.T. wrote the paper.

The authors declare no competing interest.

¹To whom correspondence should be addressed. E-mail: joet1@uchicago.edu

125 interactions are calculated as the deviation of a protein variant
126 containing several mutations from the sum of the lower order
127 effects. A concern is that technical noise or small epistatic
128 idiosyncrasies in measurement of the reference genotype or
129 low-order variants can propagate into estimates of higher-
130 order effect terms, causing spurious higher-order interactions
131 to be inferred (23). Second, many studies have not fully
132 accounted for nonspecific epistasis, which arises from a global
133 nonlinear relationship between sequence and phenotype that
134 affects all mutations identically, such as diminishing fitness
135 returns or the relationship between protein stability and
136 protein function (24–27). If this nonlinearity is not adequately
137 addressed, spurious specific interactions must be invoked to
138 explain why every mutation’s effects differ among genetic
139 backgrounds.

140 We therefore developed a method that does not suffer
141 from these sources of error and used it to systematically
142 reexamine existing datasets. Advances have been made
143 in both potential areas of concern, but currently available
144 methods still have critical limitations. Fourier analysis (28,
145 29)—also known as simplex encoding (30) or graph Fourier
146 transform (31)—is reference-free: it averages the effects of
147 sequence states across many genetic backgrounds and defines
148 them relative to the global average over all genotypes, and is
149 therefore likely to improve robustness to measurement error
150 and local idiosyncrasies. This approach can be implemented
151 as simple linear regression when sampling is limited to just
152 two amino acid states per site (32). For datasets with more
153 than two states, however, current implementations require
154 complex matrix algebra, such as building and manipulating
155 large Hadamard matrices or constructing graph Fourier bases,
156 and the resulting model terms are intelligible only with respect
157 to these matrices. Because of this complexity, only one multi-
158 amino acid dataset has been analyzed using this approach
159 (31). A third formalism-background-averaging (BA) (23),
160 also known as the Walsh-Hadamard transform (2, 33)—has
161 also been developed. This approach, which has been applied
162 only to two-amino acid datasets (but see ref. (34) for an
163 application to tRNA), occupies a middle ground between
164 reference-based and Fourier analyses: it averages mutational
165 effects over backgrounds, but it defines them relative to a
166 particular reference state at each site rather than to a single
167 reference genotype.

168 Existing methods to address nonspecific epistasis also
169 have limitations. Sometimes molecular phenotypes can be
170 measured or transformed onto a scale that is not strongly
171 affected by nonspecific epistasis, such as free energy of binding
172 within the dynamic range of assay measurement (16, 35, 36).
173 But many phenotypes scale nonadditively because of multiple
174 and complex causes, and the appropriate transformation
175 to account for nonspecific epistasis can therefore seldom
176 be known in advance (37). Several studies have addressed
177 this problem by estimating from the data a transformation
178 that minimizes nonadditivity in the relationship between
179 the measured phenotype and the estimated main effects of
180 mutations (9, 11, 13, 22, 38–40). Many of these studies
181 use rigid convex or concave transformations that cannot
182 incorporate the most important kinds of nonlinearity, such
183 as the bounding of phenotypic measurements within upper
184 and lower limits, a pattern that has been observed in many
185 DMS studies (9, 22, 38); bounding can occur if measurement
186

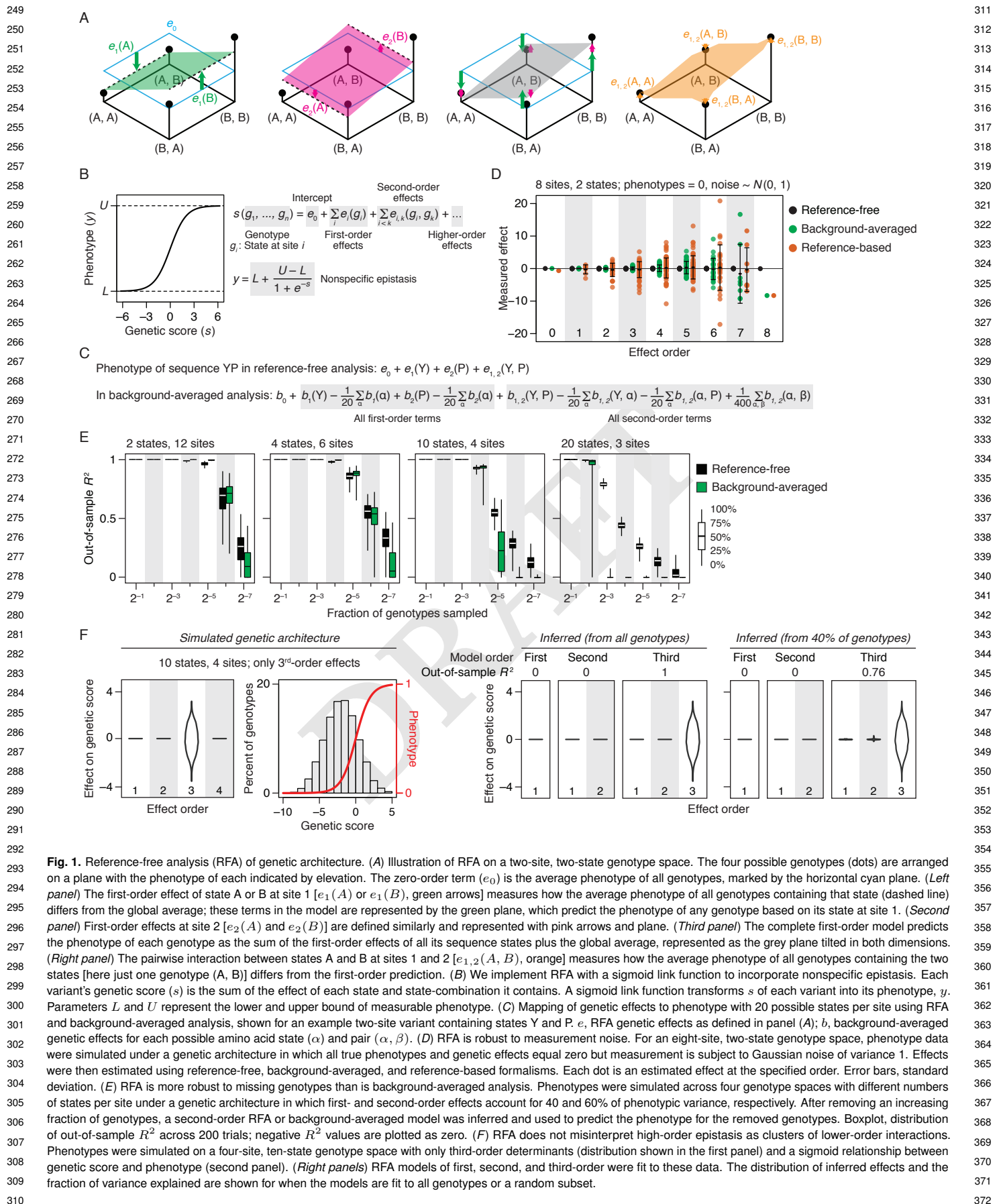
187 assays have limited dynamic range and/or the biochemical
188 processes that produce molecular phenotypes have an intrinsic
189 floor and/or ceiling, such as that produced by the relationship
190 between the free energy of folding/binding and the probability
191 that a protein occupies a functional state. Some studies have
192 used a flexible spline model or neural network (9, 22, 38) to
193 model nonspecific epistasis, but these methods have not been
194 widely adopted because they are cumbersome to implement
195 and difficult to interpret.

196 Here we develop and implement a straightforward formula-
197 tion of reference-free analysis that is applicable to any number
198 of states, and we couple it in a joint estimation procedure
199 with an effective model of nonspecific epistasis. We first
200 explain our approach and compare its desirable properties
201 to existing approaches. We then use it to reanalyze 20
202 previously published combinatorial mutagenesis experiments
203 on proteins with diverse functions, and we use the results to
204 assess the complexity of the sequence-function relationship.
205 Finally, we explore strategies to infer sequence-function
206 relationships when only a fraction of possible genotypes can
207 be experimentally sampled.

208 Results

209 **Reference-free analysis of genetic architecture.** Our method
210 of reference-free analysis defines the causal factors in a
211 protein’s genetic architecture as sequence states rather than
212 mutations. This structure allows it to describe the genetic
213 causes of phenotypic variation across the ensemble of all
214 genotypes. In reference-based and background-averaged
215 analyses, the determinants of genetic architecture are mu-
216 tations—changes from the reference state to a different
217 state—rather than the states themselves. Proteins containing
218 a reference state therefore have no genetic determinant for
219 that state at any site or for any combination across sites that
220 includes even one reference state. For example, the “wild-
221 type” sequence contains the reference state at all sites: it has
222 no mutations, so it manifests no main effects or epistatic
223 interactions at all. All the single-step neighbors of the
224 reference are each subject to one main effect, but they contain
225 no combinations of mutations, so they cannot be affected
226 by epistasis at any order. Two-step mutants are subject to
227 one pairwise epistatic effect each but cannot be affected by
228 higher-order epistasis, and so on. In fact, all these “low-order”
229 genotypes are proteins too, and their genetic architecture is
230 just as interesting and complex as protein sequences distant
231 from the wild-type.

232 Reference-free analysis (RFA) allows all genotypes to
233 provide equally important evidence about the global genetic
234 architecture. RFA takes an ANOVA-like approach in which
235 every sequence state at every site is a causal factor that
236 can potentially affect the functional phenotype, and all such
237 factors can interact with each other. A combinatorial DMS
238 study represents a full factorial experiment from which all
239 possible causal factors and all possible interactions can be
240 quantified (Fig. 1A). In the absence of nonspecific epistasis,
241 the model is structured so that each protein’s phenotype is
242 the simple sum of the functional effects of all its states and
243 combinations. The model’s zero-order term, which affects all
244 sequences, is the average phenotype across all genotypes. The
245 first-order terms are the main effects of each amino acid state
246 at every variable site in the sequence, which are defined as
247
248



373 the difference between the average phenotype of all variants
374 containing a state of interest and the global average. The
375 interaction terms at each increasing order are the epistatic
376 effects of every pair, triplet, or higher-order combination,
377 defined as the difference between the average phenotype of
378 all variants containing that set of states and the expected
379 deviation from the global average given the relevant lower-
380 order effects.

381 To incorporate nonspecific epistasis, we use a generalized
382 linear model in which each protein's phenotype is a nonlinear
383 function of its genetic score—the sum of the specific effects
384 of the states and their combinations in the protein's sequence
385 (Fig. 1B). To incorporate phenotypic bounding, we use a sig-
386 moid link function, which contains only two parameters—the
387 maximum and minimum observable phenotype—to transform
388 genetic score into phenotype.

389 RFA has several desirable features. Setting aside the
390 link function for simplicity of explanation, the RFA model
391 at each order explains the maximum amount of phenotypic
392 variance across all measured genotypes that could possibly be
393 explained by any linear model of the same order (*SI Appendix*).
394 Consider the zero-order RFA model, in which the only term
395 is the mean phenotype across all genotypes; this estimator
396 minimizes the mean squared error between measurement and
397 prediction across all variants and therefore is the best possible
398 single-parameter predictor (Fig. 1A). In the first-order RFA
399 model, the predicted phenotype of a variant is the sum of
400 all the main effects of its constituent amino acids plus the
401 global average; because each main effect is calculated as the
402 deviation of the average phenotype of all variants containing
403 some amino acid state from the global average, this set of
404 predictors again minimizes the mean squared error across all
405 variants and maximizes the phenotypic variance explained
406 compared with any other first-order model (*SI Appendix*).
407 This model structure and its desirable properties extend to
408 each increasing order.

409 Reference-free analysis contrasts with reference-based
410 analysis (RBA), which defines each effect in the model using
411 single measurements rather than averages. The RBA zero-
412 order term is the phenotype of the designated reference
413 sequence; this estimator is a good predictor in the local
414 neighborhood of the reference but is less accurate across
415 sequence space than the global average. The first-order RBA
416 term for each state is the difference between the one mutant
417 that contains that state and the reference sequence, and
418 each higher-order term is the difference between the one
419 mutant containing a combination of states and the sum of the
420 estimated lower-order effects. These are good predictors
421 of the effects of introducing each state or combination
422 into the reference background, but they are suboptimal
423 estimators across the set of all genotypes. RFA also differs
424 from background-averaged analysis (BA), which designates a
425 particular state as the reference at each site; the main effect
426 of each amino acid is defined as the average difference in
427 phenotype of the set of variants containing that state and the
428 set of variants containing the reference state at the same site.

429 The structure of RFA has several additional advantages.
430 First, the mapping from reference-free effects to phenotype is
431 intuitive. Each variant's genetic score is a simple sum of the
432 effects of its sequence states and combinations. This contrasts
433 with BA and prior implementations of Fourier analysis,
434

435 where the genetic score of each variant is a complicated
436 weighted sum of every term in the entire model, including
437 the terms for states and combinations that the variant
438 does not contain (Fig. 1C). Second, RFA facilitates direct
439 quantification of the portion of all phenotypic variation
440 that is caused by any term or set of terms in the model
441 using a simple ANOVA-like framework. Because RFA terms
442 are defined as mean deviations from the global average,
443 they have a straightforward relationship to variance: The
444 variance attributable to any RFA term is the square of its
445 magnitude normalized by the fraction of all variants that
446 contain the state or combination. The contribution of any
447 set of terms—such as all terms at some particular order or
448 some set of sites—is the sum of the individual contributions
449 (*SI Appendix*).

450 **Robustness to measurement noise and partial sampling.** If

451 we had precise phenotypic measurement for every possible
452 variant, we could exactly compute the effects of genetic
453 states and combinations as they are encoded in any of the
454 formalisms. In reality, experimental data are always affected
455 by measurement noise, and in large libraries some variants
456 typically go unmeasured. RFA is designed to perform well in
457 the face of both these challenges.

458 To assess the performance of RFA versus RBA and BA
459 when measurements are noisy, we simulated phenotypic
460 measurements using a known genetic architecture and nor-
461 mally distributed measurement error. We then estimated
462 the genetic architecture from these data and compared the
463 estimated model terms to the true values under each approach
464 (Fig. 1D). We found that RFA yields estimated effect terms
465 that are precise and unbiased. By contrast, the average error
466 in RBA's model terms is 50 times greater than in RFA, and
467 the error increases systematically with epistatic order. For
468 background averaging, the error in first-order terms is about
469 twice that of RFA, but errors grow quickly as the order of
470 epistasis increases, reaching a maximum at high orders that
471 is 100-fold worse than RFA.

472 When data are incomplete, the model terms of RFA and
473 BA can still be estimated using regression because each
474 term is averaged over many particular genotypes, and the
475 phenotypes of unmeasured variants can then be predicted
476 from the estimated model. In both cases, terms estimated
477 by regression should converge to the true effects as sample
478 size increases, and the estimates are unbiased when variants
479 are sampled without bias (*SI Appendix*). (Regression cannot
480 be used with RBA, because any missing variant makes it
481 impossible to estimate the model term signified by that
482 variant and all terms above that order that depend on it.)

483 To characterize the relative power and accuracy of
484 regression-based RFA and BA with incomplete data, we
485 simulated data using a simple genetic architecture, removed
486 a variable fraction of variants from the dataset, fit the models
487 to the remaining data by regression, and used the best-fit
488 model to predict the phenotypes of the excluded variants (Fig.
489 1E). We found that when there are only two or four states per
490 site, both RFA and BA have high predictive accuracy, with a
491 decline only after the fraction of sampled genotypes drops to
492 0.1%, at which point RFA is slightly more accurate. When
493 there are 10 or 20 possible states, however, RFA predictions
494 were much more accurate and robust than BA, the accuracy
495 of which degraded rapidly as the sample size shrank. With 20
496

497 states per site, BA became completely uninformative when
 498 sample density dropped below 25%, whereas RFA maintained
 499 some predictive value even at much lower sampling densities.

500 The structure of the formalisms explains RFA's superior
 501 performance in the face of measurement noise and partial
 502 sampling. In RFA, every measurement in the dataset is
 503 used to calculate each model term. Averaging over so many
 504 measurements dramatically reduces the influence of individual
 505 errors: the expected error in RFA terms is always smaller than
 506 that of individual phenotypic measurements, is negligible for
 507 low- and medium-order terms, and increases slowly with
 508 epistatic order. By contrast, RBA calculates each term
 509 as the difference between individual variants, without any
 510 averaging; epistasis must be invoked whenever the phenotype
 511 of a variant deviates from the sum of its lower order effects,
 512 which themselves were calculated from the deviation of
 513 single genotypes from the reference. Because each RBA
 514 term is a chain of sums and differences of many individual
 515 measurements, error variance propagates: the expected error
 516 in any RBA term is always greater than that of individual
 517 measurements and it snowballs with order, so in practice
 518 high- and even medium-order terms cannot be estimated
 519 with reasonable accuracy. For the same reason, if there are
 520 small local idiosyncrasies in the phenotype of the wild-type
 521 or low-order mutants caused by higher-order epistasis, these
 522 deviations will propagate into increasingly large estimates of
 523 high-order interactions as distance from the reference grows.

524 By estimating each effect as an average across numerous
 525 genetic backgrounds, BA reduces error propagation compared
 526 to reference-based analysis. But differences are still defined
 527 relative to a particular reference state rather than the global
 528 average, so the number of genetic backgrounds for averaging is
 529 smaller than in RFA and the sensitivity to measurement noise
 530 in each term is therefore greater. The number of relevant
 531 genetic backgrounds for estimating each BA term declines
 532 exponentially with the epistatic order, so the expected error
 533 in those terms also increases exponentially, becoming as large

559 as the error of RBA at the highest orders. Moreover, BA
 560 predicts the phenotype of an unsampled variant as a weighted
 561 sum of every single term in the model, whereas RFA uses only
 562 the terms for the states and combinations in the variant's
 563 sequence (Fig. 1C). Errors in estimated model terms caused
 564 by noise therefore propagate in BA's phenotype predictions,
 565 and this effect is exacerbated as more states per site are
 566 considered, because the total number of terms in the model
 567 increases exponentially with the number of states. RFA is
 568 insensitive to the number of states, because it predicts a
 569 variant's phenotype using only the terms for the states that
 570 are contained in its sequence. Alternative implementations
 571 of Fourier analysis are structured similarly to BA in mapping
 572 the terms to phenotype (*SI Appendix*), so they are expected
 573 to be more sensitive to noise and partial sampling than RFA.

574 **Reference-free analysis does not oversimplify genetic architecture.** We explored the possibility that RFA might
 575 oversimplify genetic architecture by misinterpreting high-
 576 order interactions as clusters of lower-order effects. The
 577 model is structured so that each order of reference-free effects
 578 produces a distinct pattern of phenotypic variation, and
 579 the pattern produced by effects at one order cannot be
 580 explained by model terms at another (*SI Appendix*). High-
 581 order variation appears as noise around the mean at lower
 582 orders, so a truncated low-order RFA model cannot explain
 583 any phenotypic variation caused by unmodeled higher-order
 584 interactions. The complexity of genetic architecture can
 585 therefore be accurately gauged by fitting truncated models
 586 and determining how much phenotypic variance is explained
 587 (*SI Appendix*).

588 To verify that RFA in practice does not oversimplify
 589 genetic architecture—particularly when nonspecific epistasis
 590 is present and sampling is incomplete—we used simulations
 591 in which phenotypes are generated by a genetic architecture
 592 that contains only third-order effects plus nonspecific epistasis.
 593 We then fit RFA models truncated at various orders to these
 594 data (Fig. 1F). First- and second-order truncated models
 595
 596
 597
 598
 599
 600

601 **Table 1. Combinatorial mutagenesis datasets analyzed in this study.**

| 602 Protein | 603 Genotype space | 604 Phenotype | 605 Ref. |
|------------------------------------|---|--|----------|
| 606 Methyl-parathion hydrolase | 607 2^5 (32) | 608 Catalytic activity | 609 (46) |
| 610 β -lactamase | 611 2^5 (32) | 612 Antibiotics resistance (MIC) | 613 (48) |
| 614 Dihydrofolate reductase | 615 3×2^4 (48) | 616 Antibiotics resistance (IC ₇₅) | 617 (3) |
| 618 Influenza A H3N2 hemagglutinin | 619 $2^2 \times 3^2 \times 4^2$ (576) | 620 Viral replication fitness | 621 (39) |
| 622 Antibody CR6261 | 623 2^{11} (2,048) | 624 Affinity for influenza hemagglutinin strain H1 | 625 (40) |
| 626 Antibody CR6261 | 627 2^{11} (2,048) | 628 Affinity for influenza hemagglutinin strain H9 | 629 (40) |
| 630 Bacterial antitoxin ParD3 | 631 20^3 (8,000) | 632 Fitness conferred by binding to toxin ParE3 | 633 (41) |
| 634 Bacterial antitoxin ParD3 | 635 20^3 (8,000) | 636 Fitness conferred by binding to toxin ParE2 | 637 (41) |
| 638 Aequorea victoria GFP (avGFP) | 639 2^13 (8,192) | 640 Fluorescence | 641 (13) |
| 642 Bacterial antitoxin ParD3 | 643 $13 \times 12 \times 10 \times 6$ (9,360) | 644 Fitness conferred by binding to toxin ParE3 | 645 (49) |
| 646 Bacterial antitoxin ParD3 | 647 $13 \times 12 \times 10 \times 6$ (9,360) | 648 Fitness conferred by binding to toxin ParE2 | 649 (49) |
| 650 SARS-CoV-2 spike protein | 651 2^{15} (32,768) | 652 Affinity for human ACE2 | 653 (7) |
| 654 Antibody CH65 | 655 2^{16} (65,536) | 656 Affinity for influenza hemagglutinin strain MA90 | 657 (21) |
| 658 Antibody CH65 | 659 2^{16} (65,536) | 660 Affinity for influenza hemagglutinin strain MA90-G189E | 661 (21) |
| 662 Antibody CH65 | 663 2^{16} (65,536) | 664 Affinity for influenza hemagglutinin strain SI06 | 665 (21) |
| 666 Antibody CR9114 | 667 2^{16} (65,536) | 668 Affinity for influenza hemagglutinin strain B | 669 (40) |
| 670 Antibody CR9114 | 671 2^{16} (65,536) | 672 Affinity for influenza hemagglutinin strain H1 | 673 (40) |
| 674 Antibody CR9114 | 675 2^{16} (65,536) | 676 Affinity for influenza hemagglutinin strain H3 | 677 (40) |
| 678 Transcription factor ParB | 679 20^4 (160,000) | 680 Fitness conferred by transcription | 681 (47) |
| 682 Protein G B1 domain (GB1) | 683 20^4 (160,000) | 684 Binding enrichment for IgG-Fc | 685 (10) |

621 correctly explain zero phenotypic variance and detect no
 622 first- or second-order effects. When the third-order model
 623 is used, all variance is correctly attributed to third-order
 624 interactions. Similar results obtain when variants are only
 625 partially sampled.

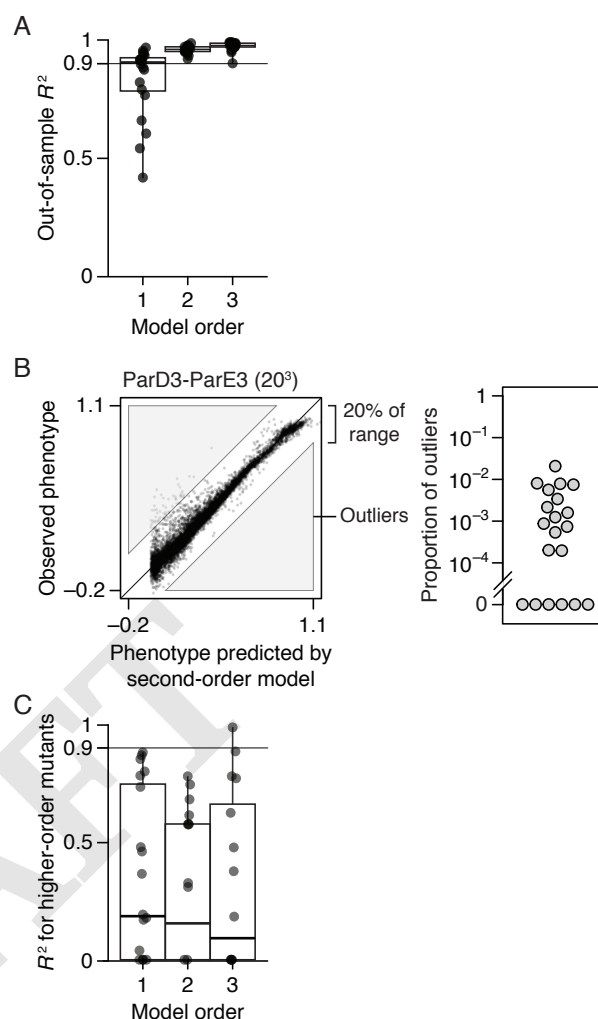
626 **Simplicity of protein sequence-function relationships.** To
 627 understand the genetic architecture of real proteins, we used
 628 RFA to analyze 20 published experiments that characterized
 629 mutant libraries in a variety of protein families with different
 630 types of functions: antibodies, enzymes, fluorescent proteins,
 631 transcription factors, viral surface proteins, and toxin-
 632 antitoxin systems. We considered only datasets in which
 633 combinatorial libraries were used and measurements had high
 634 reproducibility ($r^2 > 0.9$ among replicates; Table 1). We
 635 focused primarily on deep mutational scans of large libraries,
 636 but we included three small datasets in which high-order
 637 epistasis has been reported. The datasets range in size from
 638 32 to 160,000 possible genotypes, with the number of variable
 639 sites ranging from 3 to 16 and the number of states per site
 640 from 2 to 20.

641 We first assessed the extent to which main effects alone
 642 explain the genetic architecture by fitting a truncated first-
 643 order reference-free model, with the sigmoid link function to
 644 incorporate nonspecific epistasis. Using cross-validation to
 645 estimate the fraction of phenotypic variance explained, we
 646 found that the first-order model achieves a median out-of-
 647 sample R^2 of 0.91 across all 20 datasets, a maximum of 0.97,
 648 and > 0.75 in all but four cases (Fig. 2A). There is no clear
 649 relationship between the amount of variance explained by
 650 main effects and the number of sites or states assayed (*SI*
 651 *Appendix*, Fig. S1): the 11 datasets with $R^2 > 0.9$ include two-
 652 state, 16-site experiments in which up to 16th-order epistasis
 653 is theoretically possible (CR9114-B and H3) and a four-site,
 654 20-state experiment in which the 80 main effects account for
 655 92% of phenotypic variance (ParB). The additive effects of
 656 individual amino acids therefore account for the majority of
 657 genetic variation in protein function in most cases.

658 When second-order terms are included, virtually all genetic
 659 variance is explained, with a median cross-validation R^2
 660 of 0.96 and a minimum of 0.92 across all datasets (Fig.
 661 2A). Adding third-order terms offers only marginal or no
 662 improvement in fit (median change in out-of-sample R^2 of
 663 0.02, maximum 0.04). The small fraction of phenotypic
 664 variance unexplained by the third-order model represents
 665 some combination of fourth- and higher-order epistasis,
 666 measurement noise, and limitations in the sigmoid link
 667 function to accurately capture nonspecific epistasis.

668 Although high-order epistasis is negligible for the majority
 669 of genotypes, there could still be a subset of genotypes shaped
 670 by strong high-order epistasis. To address this possibility,
 671 we analyzed the residuals of the second-order model, which
 672 represent the sum of all higher-order epistatic interactions
 673 and measurement noise. Genotypes with a residual greater
 674 than 20% of the phenotype range were considered candidates
 675 for strong higher-order epistasis (Fig. 2B), although erratic
 676 measurement noise cannot be excluded. The proportion of
 677 such genotypes is zero in six datasets and between 0.02% and
 678 2% in the others. Strong high-order epistasis therefore affects
 679 a tiny fraction of genotypes.

680 These data establish that protein sequence-function rela-
 681 tionships are surprisingly simple: estimating just additive
 682



683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744

Fig. 2. Simplicity of protein sequence-function relationships. (A) RFA of 20 combinatorial mutagenesis datasets (Table 1). First-, second-, and third-order models with the sigmoid link function were evaluated by cross-validation—by inferring the model from a subset of data and predicting the rest of data. Each dot shows the mean out-of-sample R^2 for one dataset; boxplots show the median, interquartile range, and total range across datasets. *SI Appendix*, Fig. S1, shows the R^2 for individual datasets. (B) Variants possibly affected by strong high-order epistasis were identified as outliers in the second-order model (residual greater than 20% of the phenotype range). (Left) Outliers in the ParD3-ParE3 (20^3) dataset. Each point is a variant, plotted by its observed and predicted phenotype. (Right) Proportion of outliers in each dataset. (C) Reference-based analysis of the 20 datasets. Each model was evaluated by predicting the phenotypes of higher-order mutants. Nonspecific epistasis was accounted for as in (A), and the wild-type genotype was used as reference. Negative R^2 values are plotted as zero.

effects and pairwise interactions, coupled with a simple model of nonspecific epistasis, is sufficient for high-accuracy phenotypic prediction across the entire ensemble of protein variants. Third- and higher-order interactions are not completely absent, but these effects are typically weak, and each one affects a small number of genotypes.

Finally, we asked whether using RBA instead of RFA would produce spurious inference of epistasis from these datasets. We fit first-, second-, and third-order RBA models (including the sigmoid link function) to a designated wild-type and all single, double, and triple mutants; the phenotypes of other genotypes were then predicted using the best-fit

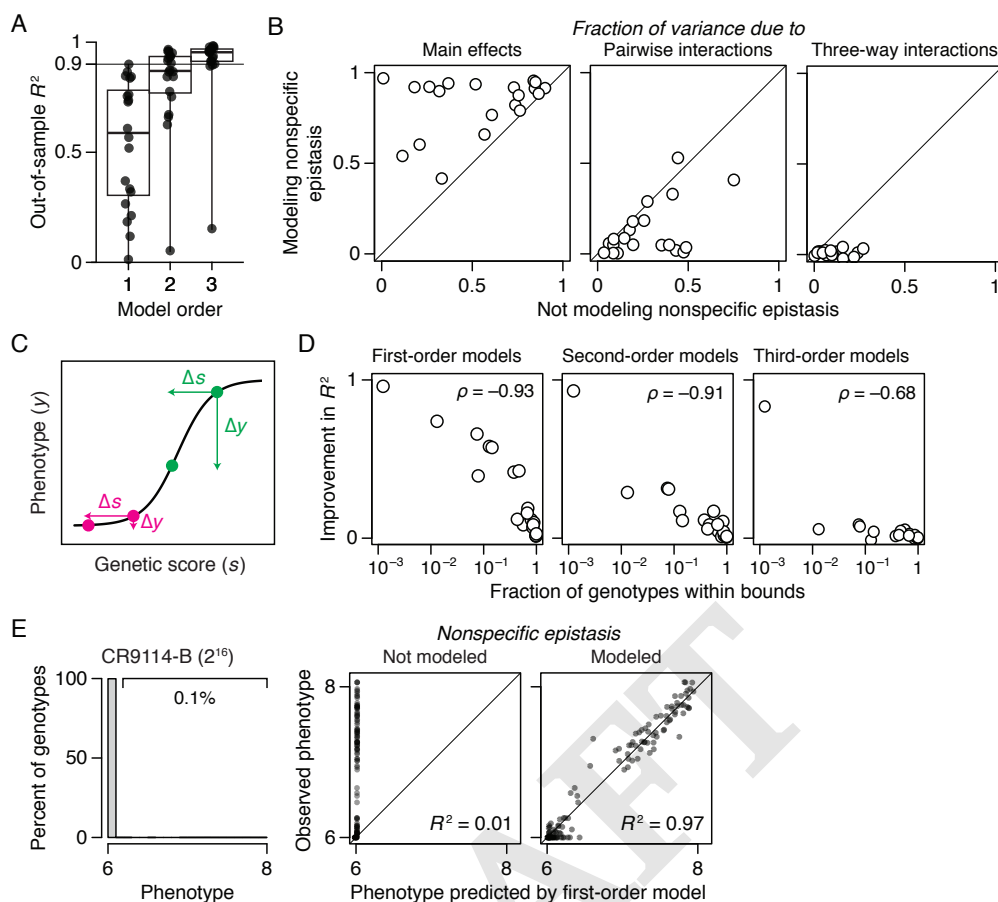


Fig. 3. The primary cause of nonspecific epistasis is phenotype bounding. (A) RFA of the 20 datasets without incorporating nonspecific epistasis, shown as in Fig. 2A. (B) Incorporating nonspecific epistasis reduces the amount of phenotypic variance attributable to pairwise and higher-order interactions. Each dot shows the variance component for one dataset computed with or without incorporating nonspecific epistasis. (C) Nonspecific epistasis causes the phenotypic effect of a mutation (Δy) to vary among genetic backgrounds (magenta versus green) even when the effect on genetic score (Δs) is constant. Phenotype bounding is a particularly strong form of nonspecific epistasis that causes mutations to appear neutral on backgrounds near the bounds but not on others. (D) The extent to which the sigmoid link function improves the model fit (comparing out-of-sample R^2 in Fig. 3A versus 2A) is proportional to the fraction of genotypes at or beyond the phenotype bounds. (E) In a dataset where only 0.1% of genotypes are within the bounds, incorporating nonspecific epistasis raises the fraction of phenotypic variance attributable to main effects from 0.01 to 0.97.

model parameters, and the R^2 was calculated. We found that RBA's accuracy is dramatically lower than RFA's: The median R^2 across datasets is less than 0.2 at all orders, leaving the vast majority of genetic variance to be explained by higher-order epistasis (Fig. 2C). The fraction of variance attributable to each epistatic order fluctuates dramatically with the protein chosen as the reference (*SI Appendix*, Fig. S2). Using the published “wild-type” sequence does not systematically attribute less or more variation to epistatic orders compared with using random reference sequences.

The primary cause of nonspecific epistasis is phenotype bounding. We next characterized the effect of incorporating nonspecific epistasis in the 20 datasets by comparing the results of RFA with and without the sigmoid link function. We found that incorporating nonspecific epistasis dramatically improves phenotype prediction, increases the variance attributable to main and low-order epistatic effects, and reduces that attributed to high-order specific epistasis (Fig. 3, A and B). For the first-order reference-free models, using the link function improves the median out-of-sample R^2 from 0.59 to 0.92. With second-order models, the sigmoid link

function improves the median R^2 from 0.87 to 0.96. With third-order models, median R^2 improves from 0.95 to 0.98.

The dramatic improvement in fit conferred by the simple sigmoid function suggests that phenotype bounds—biological or technical limits on the dynamic range over which genetic states have measurable effects on function—are the primary cause of nonspecific epistasis in most proteins (Fig. 3C). Corroborating this conclusion, the degree of improvement in R^2 when the sigmoid link function is used is tightly correlated with the fraction of genotypes at or beyond the phenotype bounds (Fig. 3D). For example, in the CR9114-B dataset, 99.9% of genotypes are at the lower bound and the out-of-sample R^2 of the first-order model rises from 0.01 to 0.97 by incorporating nonspecific epistasis (Fig. 3E). By contrast, modeling nonspecific epistasis has a modest impact when most genotypes lie within the dynamic range.

Taken together, these findings indicate that limited range of measurable phenotypic variation is the primary cause of nonspecific epistasis in deep mutational scanning datasets, and that incorporating it using a simple link function can yield a dramatic improvement in fit and reduce spurious inferences of specific epistasis, including at high orders. Although the

mechanisms underlying global nonlinearity in the genotype-phenotype relationship are likely to be complex and to vary among proteins, the simple sigmoid link function effectively captures its most salient features.

Sparsity of protein sequence-function relationships. Next, we asked whether protein function across the 20 datasets tends to be dictated by a few large-effect amino acid states/combinations or by many determinants of small effects. To quantify the sparsity of each protein's genetic architecture, we estimated the minimal number of model terms required to predict the function with 90% accuracy (T_{90}). We calculated each protein's T_{90} by ranking all the effects in the protein's genetic architecture by their contribution to phenotypic variance, constructing increasingly complex RFA models by sequentially including each effect term, and estimating the predictive accuracy of each model using cross-validation (Fig. 4A).

We found that genetic architecture is very sparse (Fig. 4B). T_{90} ranges from just 6 to 44 terms across all datasets except for the GB1 dataset (282 terms), in which the mutated sites were specifically chosen to be enriched for epistatic interactions (10). T_{90} increases very slowly with the size of genotype space, so the fraction of all possible terms that must be included to reach R^2 of 0.90 (F_{90}) declines approximately linearly as the number of possible genotypes rises (Fig. 4C). This relationship holds irrespective of the number of states per variable site. Taken together, our findings suggest that even very large genetic architectures should be describable with a compact set of important terms. For example, for a genotype space of two states at 100 variable sites ($\sim 10^{30}$ genotypes and the same number of possible model terms), the expected T_{90} is less than 10,000 terms.

Estimating genetic architecture by random sampling. Even though only a small fraction of terms is important in proteins' genetic architecture, finding them may be challenging. Experimentally analyzing exhaustive libraries is intractable for more than a small number of sites. A critical question is therefore whether genetic architecture can be estimated by a sparse learning approach that characterizes a relatively

small random sample of possible genotypes and uses penalized regression to estimate from these data the most important effects in the genetic model (13).

To characterize the fraction of genotypes that must be sampled to reconstruct the genetic architecture of each dataset, we simulated sparse learning by randomly sampling a variable number of genotypes and using penalized regression to estimate the RFA terms. We then predicted the phenotypes of the unsampled genotypes, calculated the out-of sample R^2 , and determined the minimum sample size required for R^2 of 0.9 (N_{90} ; Fig. 5A).

We found that genetic architecture cannot be reliably estimated by sparse random sampling. Excluding the three small datasets, N_{90} ranges from 0.2 to 25% of the total number of genotypes, with a median of 5% (Fig. 5B). Even the lowest end of this range does not bode well for estimating genetic architecture in large sequence spaces that contain billions or more genotypes.

We explored several factors that might determine the required sampling density: the total number of genotypes in the sequence space, the sparseness of the architecture, and the fraction of genotypes with phenotypes in the dynamic range of measurement. First, the genetic model for a larger sequence space entails more potential terms at every epistatic order, so estimating it might require sampling a larger library. We found that N_{90} does increase with the number of total possible genotypes, but there is considerable scatter in this relationship (Fig. 5B). Second, one might expect that estimating a simple genetic architecture requires a smaller sample than a more complex architecture. We found a weak relationship between the number of model terms required to explain 90% of the phenotypic variance (T_{90}) and the number of genotypes that must be sampled to achieve this level of explanation (N_{90}) (Fig. 5C). An extreme case is the CR9114-B dataset (total $2^{16} = 65,536$ genotypes), in which just ten main effects explain 90% of the variance but 16,000 genotypes—about 25% of the total—must be sampled to find them.

Finally, we considered whether the masking of phenotype by the upper or lower bound might be a factor in the effectiveness of sampling strategies. Genotypes with phenotypes at or near these limits contribute little quantitative

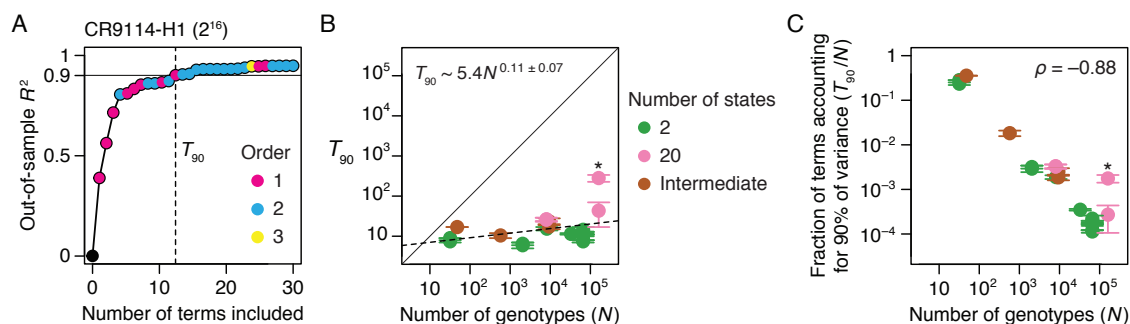


Fig. 4. Sparsity of protein sequence-function relationships. (A) Measuring the sparsity of genetic architecture illustrated on the CR9114-H1 dataset. Reference-free effects were estimated using a third-order model and then ranked by the fraction of variance they explain. Models of increasing complexity were then constructed by sequentially including each effect term, and each model was evaluated by cross-validation. Each dot represents a model, colored by the order of the last term added. Vertical line marks T_{90} , the minimal number of terms required for an out-of-sample R^2 of 0.9. (B) T_{90} as a function of the total number of genotypes. Dotted line, best-fit power function. Asterisk, GB1 dataset. Each T_{90} was estimated in two ways: as the number of terms required to reach R^2 of 0.9 (upper error bar)—an overestimate because measurement noise prevents any model from attaining out-of-sample R^2 of 1—and as the number of terms required for an R^2 equal to 90% of that of the complete third-order model (lower error bar). Circles show the average of the two estimates. (C) Fraction of all possible reference-free terms that account for 90% of phenotypic variance plotted versus the total number of genotypes. Asterisk, GB1 dataset.

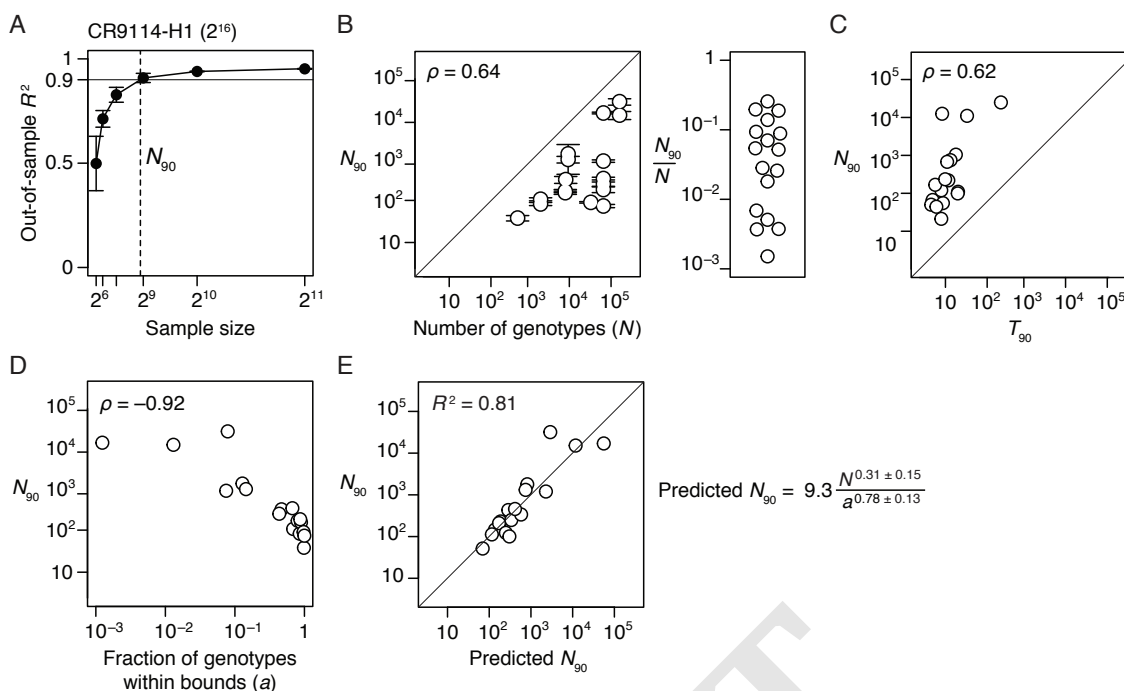


Fig. 5. Learning the genetic architecture by random sampling. (A) Learning by random sampling illustrated on the CR9114-H1 dataset. Up to third-order reference-free effects were inferred from a varying number of randomly sampled genotypes and were evaluated by predicting the phenotypes of all unsampled genotypes. For each sample size, mean and standard deviation of out-of-sample R^2 across 10 trials are shown. Dashed line marks N_{90} , the minimal sample size required for mean out-of-sample R^2 of 0.9. (B) (Left) N_{90} as a function of the total number of genotypes (N). Error bars were computed as in Fig. 4B. The three datasets with 48 or fewer genotypes are not shown. (Right) Distribution of the fraction of genotypes that must be sampled to account for 90% of phenotypic variance. (C) N_{90} as a function of T_{90} , the minimal number of reference-free effects required to explain 90% of phenotypic variance (Fig. 4). (D) N_{90} as a function of the fraction of genotypes within phenotype bounds. (E) Modeling N_{90} as a power function of the total number of genotypes (N) and the fraction of genotypes within phenotype bounds (a). The best-fit curve is shown along with standard errors.

information about the effects of the states they contain, so if most variants in a library are at the bounds, then very large samples might be required to obtain information about the genetic architecture. We found a strong negative relationship between N_{90} and the fraction of variants in the dynamic range (Fig. 5D). In the CR9114-B dataset discussed above, for example, 99.9% of all variants are at the lower bound, so the 16,000 variants required to reach N_{90} only contain about 16 genotypes in the dynamic range. Conversely, in the CH65-MA90 dataset, there are > 65,000 total genotypes, but the architecture can be estimated from a sample of just 99 variants because virtually all of the data are within the dynamic range.

The size of sequence space (N) and the fraction of variants in dynamic range (a) are therefore the key factors that determine how well a genetic architecture can be reconstructed by random sampling. To quantify the effects of these factors, we modeled N_{90} as a function of N and a across the datasets (Fig. 5E). The inferred relationship allows us to predict how large a sample should be required to estimate a genetic architecture given the size of the sequence space and the fraction of variants in dynamic range. If all genotypes in the CR9114-B dataset were in the dynamic range, a sample of only 300 variants would need to be measured, rather than the 16,000 actually required. But some sequence spaces are so large that estimating their genetic architecture by random sampling would not be practical, even if dynamic range were unlimited: for the two-state, 100-variable-site protein described above, it would still be necessary to measure

20 billion variants, even though only $\sim 10,000$ terms are expected to account for 90% of phenotypic variance.

We conclude that despite the simplicity of proteins' genetic architecture, its most important causal factors cannot be efficiently estimated by random sampling using experimental libraries, in which the majority of variants are typically nonfunctional. It is therefore important to develop an efficient non-random sampling strategy to identify the important main and pairwise effects in a protein's genetic architecture. Characterizing libraries of low-order combinations in diverse functional homologs, rather than attempting complete combinatorial scans in a single protein, may be effective. Improvements that expand the dynamic range of deep mutational scan experiments will also help.

Understanding genetic architecture. A benefit of combining the sigmoid link function with RFA is that specific genetic effects can then be expressed in simple terms that are comparable across datasets (Fig. 6A). The sigmoid model describes the observed phenotype of a protein variant as an equilibrium between "functional" and "nonfunctional" states that depends on s , the variant's genetic score; the upper and lower bounds represent ensembles in which the fraction of proteins occupying each state approaches the measurable limits. The relative occupancy of the functional state (the ratio of its occupancy to that of the nonfunctional state) is e^s , and its fractional occupancy is $(1 + e^{-s})^{-1}$. This relationship is analogous to the Boltzmann equation, with s taking on the role of $-\Delta G$, the Gibbs free energy difference between the states, expressed in units of kT . When s equals 0, the

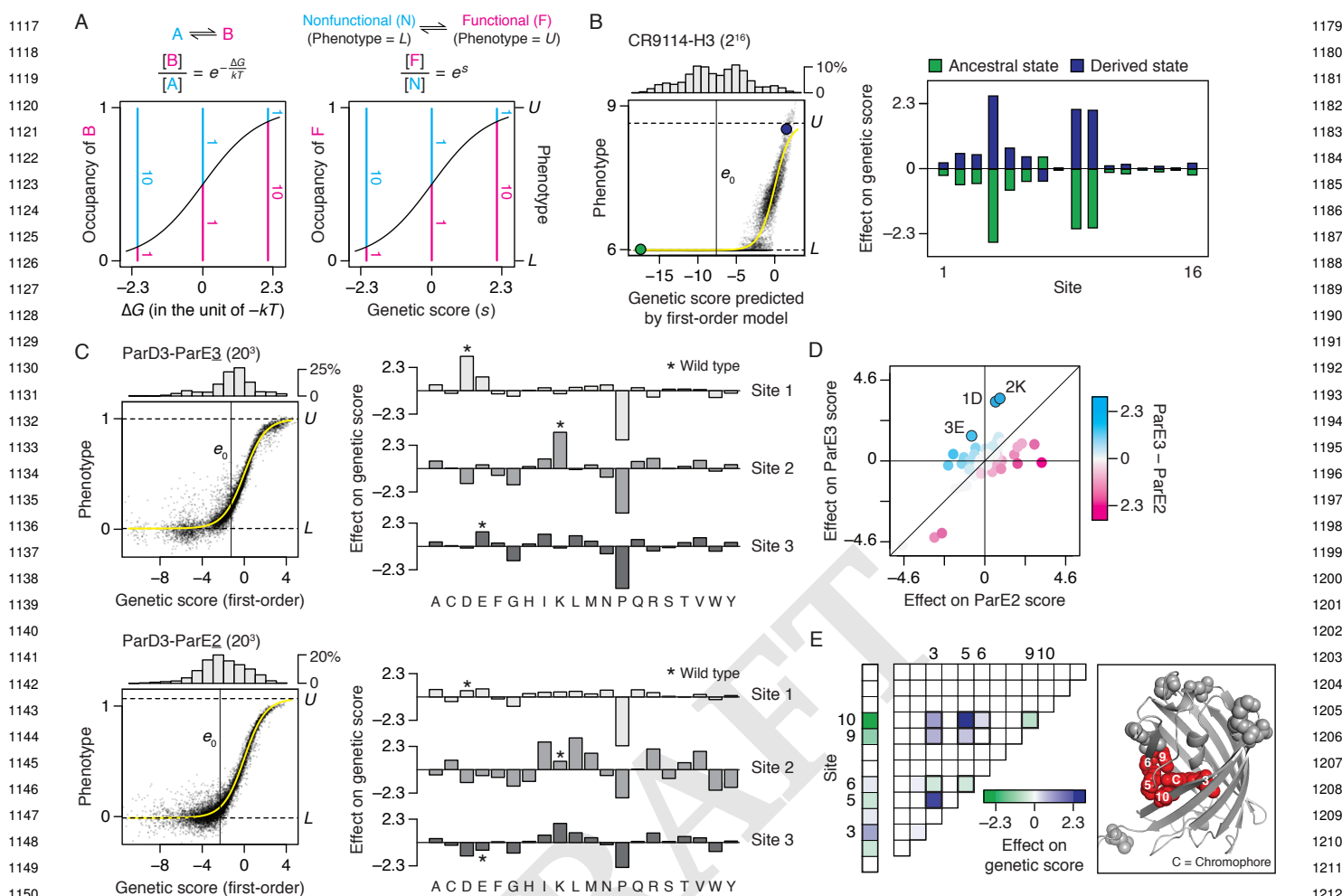


Fig. 6. Understanding genetic architecture. (A) Interpreting genetic score (s) as Gibbs free energy difference (ΔG). (Left) Relative occupancy of two thermodynamic states as a function of their ΔG . k , Boltzmann constant; T , absolute temperature. (Right) Our sigmoid model of nonspecific epistasis corresponds to an equilibrium between two states—the “functional” state, of phenotype of U , and the “nonfunctional” state, of phenotype of L . Their relative occupancy (pink and blue lines) equals e^s , allowing s to be interpreted as $-\Delta G$ in the unit of kT . (B) Analysis of the CR9114-H3 dataset, which measures the affinity of all possible combinations of ancestral and derived amino acids at 16 sites in an antibody towards an influenza hemagglutinin. (Left) First-order RFA. Each dot is a genotype, plotted by its measured phenotype and estimated genetic score. Histogram, distribution of genetic score; yellow curve, inferred nonspecific epistasis; horizontal lines, phenotype bounds; vertical line, global average; green and purple dots, ancestral and derived genotypes. (Right) First-order effects of amino acids at each site. (C) ParD3-ParE3 and ParD3-ParE2 (20^3) datasets, which measure how all possible variants of the protein ParD3 at three sites bind to ParE3, the cognate substrate, or ParE2, a noncognate substrate. (Left) First-order RFA shown as in (B). (Right) First-order effects of amino acids at each site. Asterisk, wild type. (D) Comparing the effect of each amino acid on ParE3 versus ParE2 binding. Wild type amino acids are marked. (E) avGFP dataset, which measures the fluorescence of all possible combinations of pairs of amino acids at 13 sites. (Left) Main effects and pairwise interactions, which account for 57 and 38% of phenotypic variance, respectively. Values are shown only for one of the two of amino acids in each site. The ten pairwise interactions possible among sites 3, 5, 6, 9, and 10 are outlined. (Right) Crystal structure of avGFP (PDB ID: 3e5w). Spheres, the 13 mutated residues; red, the chromophore and the five residues with the strongest phenotypic contribution.

functional and nonfunctional states are equally populated, and the phenotype is midway between the upper and lower bounds. An amino acid that increases the score by 2.3 always causes a ten-fold increase in the relative occupancy of the measurable functional state, corresponding to an apparent $\Delta\Delta G$ of -1.4 kcal/mol at 37°C . This relationship holds across proteins, functions, and assay systems, which all display the same scaling relationship between a variant’s genetic score and its phenotype, mediated via the probability of occupying the functional state.

We used this framework to interpret the genetic architecture of several example proteins. First, the CR9114-H3 dataset (Fig. 6B) consists of affinity measurements for binding of hemagglutinin to each of 2^{16} antibodies (all

possible combinations of ancestral and derived amino acids at 16 sites that evolved during affinity maturation). The vast majority of variants in this library are at or near the lower bound of detectable binding; as a result, the average genetic score is -7.8 , corresponding to just 0.04% occupancy of the measurable functional state (ΔG of 4.7 kcal/mol). Even the highest genetic score in the entire library is only 2.6 — 93% occupancy of the functional state. Main effects at three key sites explain the most phenotypic variance: Substituting any of these from the ancestral to derived state increases the genetic score by between 4.2 and 5.2, corresponding to an increase in relative occupancy of the functional state by 70- to 180-fold and a $\Delta\Delta G$ of ~ 2 to 3 kcal/mol each. Other sites make modest contributions: The five next-largest effects

1241 each change the genetic score by about 1 (0.7 kcal/mol) when
1242 mutated back to the ancestral state, shifting the relative
1243 occupancy by 36% each, but reducing the absolute occupancy
1244 to just 8% when all five change together. There is virtually no
1245 specific epistasis in this genetic architecture (Fig. S1). This
1246 means that there are many different combinations of the five
1247 moderate-effect sites that provide a sufficient genetic score
1248 to confer measurable affinity, but only if the derived state at
1249 all three large-effect sites are present. The remaining eight
1250 sites have negligible effects on binding and are completely
1251 degenerate among functional antibodies.

1252 Second, the genetic architecture of specificity in a protein
1253 can be understood by analysis of genetic scores with different
1254 substrates (Fig. 6C). A deep mutational scan was performed
1255 on the ParD3 protein (20 states at 3 sites in the binding
1256 interface) for binding to its cognate ligand ParE3 and a
1257 noncognate ligand ParE2 (41). In both cases, first-order
1258 determinants account for the vast majority of genetic variance,
1259 with main effects on genetic score ranging from strongly
1260 positive (3.6) to strongly negative (-4.8); this corresponds
1261 to changes in ΔG on the order of -2 to 3 kcal/mol and
1262 changes in relative occupancy ranging from a 36-fold increase
1263 to 120-fold decrease. Effects on specificity can be quantified
1264 as the difference in a state's effect on genetic score for the two
1265 substrates. Eight different states distributed across the three
1266 variable sites change the genetic score in favor of one ligand
1267 or the other by more than 1.6, meaning that they change the
1268 relative occupancy of the two substrates by at least 5-fold
1269 each (Fig. 6D). For example, the states in the wild-type
1270 ParD3 (Asp [D], Lys [K], and Glu [E] in the three variable
1271 sites) increase specificity for the cognate ligand by scores
1272 of 2.8, 2.7, and 2.2, respectively, corresponding to a 10-fold
1273 change in relative occupancy by each; two of these states
1274 (1D and 2K) achieve this by increasing the affinity for both
1275 ligands with a stronger effect on cognate versus noncognate
1276 binding, whereas 3E increases cognate binding but reduces
1277 noncognate binding.

1278 Finally, RFA can be used to characterize the scale of
1279 epistatic effects on function. In the avGFP dataset (13),
1280 pairwise interactions account for 38% of phenotypic variance.
1281 Out of 13 sites analyzed, however, main and pairwise effects
1282 involving just five sites account for the vast majority of the
1283 variance explained (Fig. 6E). These sites, which tightly
1284 surround the chromophore in the avGFP crystal structure,
1285 engage in a dense network of epistatic interactions in which
1286 nine of the ten possible pairwise interactions are non-zero.
1287 Although only three of these effects are strong (changing the
1288 genetic score by > 1), the total impact is substantial: A total
1289 change in genetic score of 2.8 caused by main effects and 7.5
1290 by pairwise interactions, corresponding to 16- and 1,700-fold
1291 increases in the relative occupancy of functional state (1.7
1292 and -4.5 kcal/mol), respectively.

1293 Discussion

1294 Our finding that main and pairwise interactions account
1295 for virtually all genetic variation within proteins contrasts
1296 with many earlier reports (1-7). This difference is likely
1297 attributable to overestimation of epistasis in prior studies,
1298 the vast majority of which used reference-based analysis
1299 and/or have not fully decoupled specific epistasis from global
1300 nonlinearity in the genotype-phenotype relationship. It is
1301
1302

1303 possible that higher-order epistasis is more important in
1304 some other proteins not examined here, but this seems
1305 unlikely, given the consistency of the pattern we observed
1306 across 20 different deep mutational scans in a wide variety
1307 of proteins with different architectures and functions. Most
1308 of the studies we examined focused on a small or moderate
1309 number of sites selected a priori because they vary between
1310 two functional proteins of interest or they are in important
1311 structural positions (e.g., at binding interfaces or active sites).
1312 In some cases the sites are clustered, and in others they are
1313 dispersed across the protein structure. We therefore have no
1314 reason to expect that the sites examined in the studies we
1315 analyzed are depleted for higher-order epistasis.

1316 Our analyses assessed the genetic architecture of a single
1317 function per protein, rather than the determinants of
1318 functional specificity when multiple functions are measured.
1319 It is possible that higher-order interactions could be more
1320 important in determining functional specificity. Reference-
1321 free analysis could easily be expanded to identify the genetic
1322 architecture of specificity using DMS studies of multiple
1323 functions; a recent study used a similar approach and found
1324 that higher-order interactions within a transcription factor are
1325 unimportant for determining its specific preferences among
1326 DNA binding sites (42). Higher-order epistasis might be
1327 more important among loci than it is within proteins, but
1328 this is an open question. It is not obvious, for example,
1329 that contacts across interfaces between molecules should
1330 produce more higher-order genetic interactions than the
1331 physically similar contacts that occur within proteins, or
1332 that dependencies among molecules in signal transduction
1333 or metabolic pathways should involve more higher-order
1334 interactions than within the complex environment of a single
1335 protein, once the global nonlinearities imposed by these
1336 pathways are accounted for.

1337 The lack of higher-order epistasis within proteins may seem
1338 surprising, given the complexity of proteins' three-dimensional
1339 structure, in which clusters of three or more residues often
1340 contact each other directly. Our findings suggest that the
1341 effects of most such clusters can be largely explained by the
1342 sum of the pairwise interactions they comprise. But these
1343 couplings themselves depend on conformation, which itself is
1344 determined by the state at other sites; if a mutation alters
1345 the conformation, it will change some pairwise couplings
1346 and produce higher-order epistasis. In the datasets we
1347 examined, this kind of conformational epistasis appears to be
1348 relatively unimportant. A possible explanation is that in these
1349 experiments the majority of sites—and therefore presumably
1350 the protein's overall fold—were held constant. Ultimately,
1351 the folding of a protein into its native conformation and
1352 the couplings that result would seem to require higher-order
1353 interactions, and these might be revealed if a large scan of a
1354 protein that can adopt multiple conformations were possible.
1355 The importance of these interactions in the overall sequence-
1356 function relationship relative to lower-order effects, however,
1357 is an open question.

1358 The effectiveness of the Boltzmann-like sigmoid function
1359 to model nonspecific epistasis seems surprising, because
1360 nonlinearity in the genotype-phenotype relationship almost
1361 certainly arises from complex biological and technical causes
1362 that vary among proteins, functions, and measurement
1363 techniques. Our analyses indicate that upper and lower
1364

1365 bounds on the dynamic range over which a phenotype
 1366 can be produced and measured are the primary cause of
 1367 nonspecific epistasis within proteins. Whether or not the
 1368 sigmoid transformation is “true,” our findings indicate that
 1369 accounting for this form of nonlinearity—irrespective of the
 1370 factors that produce it—is sufficient to allow a low-order model
 1371 of specific epistasis to provide a parsimonious explanation
 1372 of genetic architecture that captures virtually all phenotypic
 1373 variation across all the proteins we examined.

1374 Our finding that RFA outperforms RBA in providing a
 1375 compact and accurate characterization of the global genotype-
 1376 phenotype map does not mean that RBA is never useful.
 1377 There are some settings in which the object of interest is not
 1378 a protein’s genetic architecture but particular interactions
 1379 among mutations in the sequence neighborhood immediately
 1380 around a designated wild-type or ancestral protein. In these
 1381 cases RBA is appropriate, but it should be used with caution
 1382 because of its propensity to infer spurious interactions as
 1383 distance from the reference sequence increases.

1384 Epistasis can make evolutionary trajectories contingent on
 1385 the chance occurrence of permissive and restrictive epistatic
 1386 modifiers (27, 43, 44). It was recently shown that the effects
 1387 of most mutations drift gradually as proteins accumulate
 1388 substitutions over long-term evolutionary time (45). Our
 1389 results imply that this drift is likely attributable to the
 1390 cumulative effect of many small pairwise interactions rather
 1391 than higher-order modulations. The relative unimportance
 1392 of high-order epistasis implies that the pairwise dependencies
 1393 that make evolution contingent on prior mutations are likely
 1394 to remain largely stable over evolutionary time, rather than
 1395 being idiosyncratically rewired with every substitution that
 1396 occurs at other sites.

1397 For scientists who would like to understand how proteins
 1398 work, our findings are reassuring, but they clarify a major
 1399 challenge ahead. Proteins’ genetic architecture is intelligible;
 1400 a small fraction of main and pairwise effects provides a
 1401 compact and efficient explanation of 90 to 95% of functional
 1402 genetic variation across the vast space of possible sequences.
 1403 Complete combinatorial experiments are intractable for many
 1404 states at more than a few sites or even two states at a
 1405 moderate number of sites, but the unimportance of high-
 1406 order epistasis means that it is unnecessary to assay the vast
 1407 array of triplets, quartets, and so on. The challenge is that
 1408 the small set of key first- and second-order determinants
 1409 cannot be efficiently identified from a random sample of
 1410 variants, because sequence space is huge and most random
 1411 polypeptides are virtually nonfunctional—particularly when
 1412 the dynamic range of measurement is limited—so they do not
 1413 provide useful quantitative information about the sequence
 1414 states and pairs that they contain. Assessing low-order
 1415 effects in a single sequence neighborhood is not sufficient,
 1416 because the resulting estimates would be subject to the
 1417 same kind of errors and idiosyncracies that plague reference-
 1418 based estimates. An effective strategy may therefore be to
 1419 perform comprehensive single- and double-mutant scans using
 1420 a diverse set of functional proteins as starting points, and then
 1421 analyze the results using RFA. A critical issue is to determine
 1422 just how diverse the proteins used as starting points must
 1423 be, while continuing to improve the efficiency and dynamic
 1424 range of experimental methods. The potential power of a
 1425 relatively practical strategy like this has been overlooked to
 1426

1427 date, presumably because protein architecture is not nearly
 1428 as complex as it was previously thought to be.

1429 Methods

1430 **Reference-free analysis (RFA).** Here we define RFA and state
 1431 its key properties. A detailed exposition with proofs is pro-
 1432 vided in *SI Appendix*, and scripts and tutorials for performing
 1433 RFA are available on GitHub ([github.com/whatdoidohaha/](https://github.com/whatdoidohaha/RFA)
 1434 *RFA*).

1435 We consider a simple genotype space defined by q states at
 1436 each of n sites, but RFA can also be applied when the number
 1437 of states varies among sites. Let \mathbf{g} denote a genotype, $y(\mathbf{g})$
 1438 its phenotype, and G the set of all genotypes. The global
 1439 average phenotype is denoted

$$1440 e_0 = \langle y | G \rangle,$$

1441 where the brackets indicate averaging of y across G . RFA
 1442 decomposes the phenotype into the contribution of individual
 1443 states and their interactions. The first-order effect of state s
 1444 at site i is the difference between the average phenotype of
 1445 the subset of genotypes sharing that state (denoted G_i^s) and
 1446 the global average:

$$1447 e_i(s) = \langle y | G_i^s \rangle - e_0.$$

1448 The pairwise interaction between states s_1 and s_2 at sites i_1
 1449 and i_2 is the difference between the average phenotype of the
 1450 subset of genotypes sharing that state-pair ($G_{i_1, i_2}^{s_1, s_2}$) and the
 1451 global average after accounting for the main effects:

$$1452 e_{i_1, i_2}(s_1, s_2) = \langle y | G_{i_1, i_2}^{s_1, s_2} \rangle - [e_0 + e_{i_1}(s_1) + e_{i_2}(s_2)].$$

1453 Similarly, higher-order effects are the difference between
 1454 the average phenotype of the subset of genotypes sharing a
 1455 particular set of states and the global average after accounting
 1456 for the lower-order effects.

1457 RFA predicts the phenotype of a genotype of interest by
 1458 summing the effects of the states present in that genotype.
 1459 For a genotype with state g_i in each site i , the predicted
 1460 phenotype under RFA of order k is

$$1461 y_k(\mathbf{g}) = e_0 + \sum_i e_i(g_i) + \sum_{i_1 < i_2} e_{i_1, i_2}(s_1, s_2) + \dots +$$

$$1462 \sum_{i_1 < \dots < i_k} e_{i_1, \dots, i_k}(g_{i_1}, \dots, g_{i_k}).$$

1463 The overall accuracy of this prediction can be quantified by
 1464 the sum of squared errors:

$$1465 \epsilon_G = \sum_{\mathbf{g} \in G} [y(\mathbf{g}) - y_k(\mathbf{g})]^2.$$

1466 Among all possible ways of predicting the phenotype using
 1467 effects of order up to k —including reference-based analysis
 1468 under any choice of reference genotype and background-
 1469 averaged analysis under any choice of reference states—RFA
 1470 minimizes ϵ_G for any k for any genetic architecture. For
 1471 example, when k equals zero—that is, when all phenotypes
 1472 are predicted by a single number— ϵ_G is minimized by the
 1473 global average phenotype, which is the RFA zero-order term.
 1474

By explaining as much phenotypic variance as possible at any order of approximation, RFA provides the simplest description of genetic architecture.

A key task in the analysis of genetic architecture is to quantify the contribution of individual states and interactions to the phenotype. RFA facilitates this task by decomposing the total phenotypic variance into the contribution of each factor:

$$\text{Var}(y|G) \left(= \frac{1}{q^n} \sum_{g \in G} [y(g) - \langle y|G \rangle]^2 \right) = \sum_{e \neq e_g} \frac{e^2}{q^{O(e)}},$$

where e denotes an effect, $O(e)$ its order, and the summation involves all nonzero-order effects. An effect of order k affects the phenotype of one in every q^k genotypes. The expression above therefore states that the amount of phenotypic variance attributable to an effect is the square of its magnitude, normalized by the fraction of genotypes it affects.

A corollary of the definition of reference free effects is that the first-order effects of all states at a site sum to zero:

$$\sum_{1 \leq s \leq q} e_i(s) = 0.$$

We call this the zero-mean property. The second-order effects of all state-pairs in one site-pair also sum to zero, as do all higher-order effects at a combination of sites.

Inferring reference-free effects from noisy and incomplete data. When individual phenotypes are subject to measurement error of variance ω , reference-free effects of order k computed from these measurements have an error of variance

$$\frac{(q-1)^k}{q^n} \omega.$$

By definition $k \leq n$, so the variance of computed effects is always less than ω and is miniscule when k is small relative to n . Therefore, reference-free effects can be robustly determined from noisy phenotypic measurements, thanks to the averaging of effects over large numbers of genotypes. By contrast, the error associated with reference-based effects of order k is $2^k \omega$, which is always greater than ω and typically too large to distinguish effects from errors when $k > 2$. The error associated with background-averaged effects of order k is $(2q)^k / q^n \times \omega$, which is greater than the error for reference-free effects of the same order and exceeds ω as k increases.

When measurement is incomplete, reference-free effects can be inferred by regression. To infer the effects in a truncated model that contains terms of order up to k , we model

$$y(\mathbf{g}) = y_k(\mathbf{g}) + \epsilon(\mathbf{g}),$$

where the error $\epsilon(\mathbf{g})$ is the sum of all higher-order effects and measurement noise. Regression estimates are obtained by minimizing the sum of squared errors across the set of sampled genotypes (G^*):

$$\epsilon_{G^*} = \sum_{g \in G^*} [y(\mathbf{g}) - y_k(\mathbf{g})]^2.$$

Because RFA minimizes the sum of squared errors across all genotypes, the regression estimates converge to the true values

as more genotypes are sampled. Furthermore, the regression estimates are unbiased provided that genotypes are randomly missing. This is because $\epsilon(\mathbf{g})$ is unbiased—equals zero when averaged across all genotypes. This in turn derives from the zero-mean property, which implies that the net phenotypic contribution of any order of effects is zero when averaged across all genotypes; unmodeled higher-order interactions do not bias the regression because they appear as noise to lower-order models.

Nonspecific epistasis. We account for nonspecific epistasis by assuming that the effects of sequence states are transformed by a global link function into the observed phenotype (25). The net effect of the sequence states in a genotype is referred to as its genetic score (s). We model the link function by a sigmoid that is defined by two parameters, L and U , which represent the lower and upper bound of phenotype:

$$y(\mathbf{g}) = L + \frac{U - L}{1 + e^{-s(\mathbf{g})}}.$$

Implementation. Reference-free effects and nonspecific epistasis were jointly inferred by L1-regularized regression. The optimal L1 penalty was determined by maximizing the out-of-sample R^2 in cross-validation. Except for four datasets, cross-validation was performed by randomly partitioning genotypes into training and test sets. For the three datasets with 48 or fewer genotypes and the CR9114-B dataset where only 81 genotypes are above the lower phenotype bound, cross-validation was performed by leaving out each measurement replicate in turn. The R package *lbfgs* was used for numerical optimization. All scripts for inference and analysis are available on GitHub (github.com/whatdoidohaha/RFA).

To jointly infer reference-based effects and nonspecific epistasis, we devised a two-step approach. This was necessary because reference-based analysis is incompatible with regression. For example, regression infers a first-order model by assigning values to the effects of point mutations that best predict the phenotype for both point and higher-order mutants. However, the effect of a point mutation is defined solely by the phenotype of the one variant that contains only that mutation; the regression estimate can be far from true depending on the exact phenotypes of higher-order mutants. For each candidate set of nonspecific epistasis parameters, we computed the reference-based effects on genetic score that exactly recapitulate the phenotypes of mutants up to model order. The effects were then used to predict the phenotype for higher-order mutants. We only predicted higher-order mutants for which all relevant lower-order effects are measured; for example, when a point mutant is missing, any double or higher-order mutant involving that mutation was excluded from prediction. This procedure was repeated for different values of nonspecific epistasis parameters, resulting in values that maximize the R^2 .

Background-averaged analysis was originally developed only for binary state spaces. To implement it for spaces with more than two states per site, we extended the recursive matrix formalism of ref. (23) and implemented it in a custom R script. The same multi-state formalism has been independently derived and published recently (34).

Combinatorial mutagenesis datasets. We systematically mined the literature for mutagenesis experiments with a

1613 combinatorially complete design. Among the many datasets
1614 comprising fewer than 100 genotypes, we chose three datasets
1615 where high-order epistasis has been reported. Any larger
1616 dataset in which precise measurement ($r^2 > 0.9$ between
1617 replicates) is available for at least 40% of possible genotypes
1618 was included for analysis. Several datasets were edited as
1619 follows.

1620 The methyl-parathion hydrolase activity (46) was mea-
1621 sured in the presence of seven different metal cofactors. In
1622 every case, second-order reference-free analysis coupled with
1623 the sigmoid model of nonspecific epistasis explained more
1624 than 90% of phenotypic variance. Only the Ni²⁺ dataset, in
1625 which epistasis accounts for the greatest fraction of phenotypic
1626 variance, is presented here.

1627 The original dihydrofolate reductase dataset (3) includes
1628 a noncoding mutation for a total of 96 variants. We only
1629 analyzed the 48 coding site variants fixed for the mutant state
1630 in the noncoding site. IC₇₅—the antibiotics concentration
1631 that reduces the growth rate by 75%—was reported in
1632 logarithmic scale, set arbitrarily as -2 when the variant is
1633 unviable at any concentration. We reverted the logarithm,
1634 making IC₇₅ equal to 0 when the variant is unviable.

1635 The hemagglutinin study (39) characterized an identical
1636 set of genetic variants in six different genetic backgrounds.
1637 We only analyzed the genetic background for which the
1638 measurement is most precise (Bei89).

1639 In the avGFP dataset (13), fluorescence is systematically
1640 higher in the second measurement replicate by a factor of
1641 1.31. This difference was normalized when combining the two
1642 replicates.

1643 The ParB study (47) measures how the transcription factor
1644 ParB binds to two DNA motifs, parS and NBS. Because
1645 measurement r^2 is less than 0.9 for the NBS dataset, only
1646 the parS dataset was analyzed. The absolute fitness of each
1647 variant was inferred by comparing the read count before
1648 and after the bulk competition assay. Variants with the pre-
1649 competition read count fewer than 15 were excluded, resulting
1650 in 42.2% coverage of the 160,000 possible genotypes—down
1651 from 97.0% in the original study.

1652 The extent of measurement noise in the GB1 dataset (10)
1653 could not be directly determined because measurement was
1654 not replicated, but comparison to an independent dataset for
1655 a subset of variants showed that measurement r^2 is greater
1656 than 0.9. Variants with a pre-competition read count fewer
1657 than 100 were excluded, resulting in 68.6% coverage of the
1658 160,000 possible genotypes—down from 93.4% in the original
1659 study.

1660 Acknowledgments

1661 We thank members of the Thornton Laboratory and R.
1662 Ranganathan and S. Kuehn at the University of Chicago
1663 for discussion, and the University of Chicago Research
1664 Computing Center for high-performance computing. This
1665 work was supported by the National Institutes of Health
1666 grants R01GM131128 (J.W.T.), R01GM121931 (J.W.T.),
1667 and F32GM122251 (B.P.H.M.) and Samsung Scholarship
1668 (Y.P.).
1669

1670 References

1671 1. Sadvovsky E, Yifrach O (2007) Principles underlying ener-
1672 getic coupling along an allosteric communication trajectory
1673
1674

of a voltage-activated K⁺ channel. *Proc Natl Acad Sci USA* 1675
104(50):19813–19818. 1676

2. Weinreich DM, Lan Y, Wylie CS, Heckendorn RB 1677
(2013) Should evolutionary geneticists worry about higher- 1678
order epistasis. *Curr Opin Genet Dev* 23(6):700–707. 1679

3. Palmer AC et al. (2015) Delayed commitment to 1680
evolutionary fate in antibiotic resistance fitness landscapes. 1681
Nat Commun 6:7385. 1682

4. Sailer ZR, Harms MJ (2017) Molecular ensembles 1683
make evolution unpredictable. *Proc Natl Acad Sci USA* 1684
114(45):11938–11943. 1685

5. Guerrero RF, Scarpino SV, Rodrigues JV, Hartl DL, 1686
Ogbunugafor CB (2019) Proteostasis Environment Shapes 1687
Higher-Order Epistasis Operating on Antibiotic Resistance. 1688
Genetics 212(2):565–575. 1689

6. Lozovsky ER, Daniels RF, Heffernan GD, Jacobus DP, 1690
Hartl DL (2021) Relevance of Higher-Order Epistasis in Drug 1691
Resistance. *Mol Biol Evol* 38(1):142–151. 1692

7. Moulana A et al. (2022) Compensatory epistasis 1693
maintains ACE2 affinity in SARS-CoV-2 Omicron BA.1. *Nat* 1694
Commun 13(1):7011. 1695

8. Chen J, Stites WE (2001) Higher-order packing 1696
interactions in triple and quadruple mutants of staphylococcal 1697
nuclease. *Biochemistry* 40(46):14012–14019. 1698

9. Sarkisyan KS et al. (2016) Local fitness landscape of 1699
the green fluorescent protein. *Nature* 533(7603):397–401. 1700

10. Wu NC, Dai L, Olson CA, Lloyd-Smith JO, Sun R 1701
(2016) Adaptation in protein fitness landscapes is facilitated 1702
by indirect paths. *eLife* 5:e16965. 1703

11. Sailer ZR, Harms MJ (2017) Detecting High-Order 1704
Epistasis in Nonlinear Genotype-Phenotype Maps. *Genetics* 1705
205(3):1079–1088. 1706

12. Adams RM, Kinney JB, Walczak AM, Mora T (2019) 1707
Epistasis in a fitness landscape defined by antibody-antigen 1708
binding free energy. *Cell Syst* 8(1):86–93. e3. 1709

13. Poelwijk FJ, Socolich M, Ranganathan R (2019) 1710
Learning the pattern of epistasis linking genotype and 1711
phenotype in a protein. *Nat Commun* 10(1):4213. 1712

14. Tamer YT et al. (2019) High-Order Epistasis in 1713
Catalytic Power of Dihydrofolate Reductase Gives Rise to a 1714
Rugged Fitness Landscape in the Presence of Trimethoprim 1715
Selection. *Mol Biol Evol* 36(7):1533–1550. 1716

15. Yang G et al. (2019) Higher-order epistasis shapes 1717
the fitness landscape of a xenobiotic-degrading enzyme. *Nat* 1718
Chem Biol 15(11):1120–1128. 1719

16. Ballal A et al. (2020) Sparse Epistatic Patterns 1720
in the Evolution of Terpene Synthases. *Mol Biol Evol* 1721
37(7):1907–1924. 1722

17. Hinkley T et al. (2011) A systems analysis of 1723
mutational effects in HIV-1 protease and reverse transcriptase. 1724
Nat Genet 43(5):487–489. 1725

18. Olson CA, Wu NC, Sun R (2014) A comprehensive 1726
biophysical description of pairwise epistasis throughout an 1727
entire protein domain. *Curr Biol* 24(22):2643–2651. 1728

19. Podgornaia AI, Laub MT (2015) Pervasive degen- 1729
eracy and epistasis in a protein-protein interface. *Science* 1730
347(6222):673–677. 1731

20. Diss G, Lehner B (2018) The genetic landscape of a 1732
physical interaction. *eLife* 7:e32472. 1733

- 1737 21. Phillips AM et al. (2023) Hierarchical sequence-affinity
1738 landscapes shape the evolution of breadth in an anti-influenza
1739 receptor binding site antibody. *eLife* 12:e83628. 1799
- 1740 22. Otwinowski J, McCandlish DM, Plotkin JB (2018)
1741 Inferring the shape of global epistasis. *Proc Natl Acad Sci*
1742 *USA* 115(32):E7550–E7558. 1800
- 1743 23. Poelwijk FJ, Krishna V, Ranganathan R (2016) The
1744 Context-Dependence of Mutations: A Linkage of Formalisms.
1745 *PLoS Comput Biol* 12(6):e1004771. 1801
- 1746 24. Kryazhimskiy S, Rice DP, Jerison ER, Desai MM
1747 (2014) Global epistasis makes adaptation predictable despite
1748 sequence-level stochasticity. *Science* 344(6191):1519–1522. 1802
- 1749 25. Otwinowski J, Plotkin JB (2014) Inferring fitness land-
1750 scapes by regression produces biased estimates of epistasis.
1751 *Proc Natl Acad Sci USA* 111(22):E2301–E2309. 1803
- 1752 26. Johnson MS, Reddy G, Desai MM (2023) Epistasis
1753 and evolution: recent advances and an outlook for prediction.
1754 *BMC Biol* 21(1):120. 1804
- 1755 27. Starr TN, Thornton JW (2016) Epistasis in protein
1756 evolution. *Protein Sci* 25(7):1204–1218. 1805
- 1757 28. Weinberger ED (1991) Fourier and Taylor series on
1758 fitness landscapes. *Biol Cybern* 65(5):321–330. 1806
- 1759 29. Stadler PF (1996) Landscapes and their correlation
1760 functions. *J Math Chem* 20(1):1–45. 1807
- 1761 30. Stormo GD (2011) Maximally efficient modeling of
1762 DNA sequence motifs at all levels of complexity. *Genetics*
1763 187(4):1219–1224. 1808
- 1764 31. Brookes DH, Aghazadeh A, Listgarten J (2022) On
1765 the sparsity of fitness functions and implications for learning.
1766 *Proc Natl Acad Sci USA* 119(1):e2109649118. 1809
- 1767 32. Miton CM, Chen JZ, Ost K, Anderson DW, Tokuriki
1768 N (2020) Statistical analysis of mutational epistasis to reveal
1769 intramolecular interaction networks in proteins. *Methods*
1770 *Enzymol* 643:243–280. 1810
- 1771 33. Weinreich DM, Lan Y, Jaffe J, Heckendorn RB (2018)
1772 The influence of higher-order epistasis on biological fitness
1773 landscape topography. *J Stat Phys* 172(1):208–225. 1811
- 1774 34. Faure AJ, Lehner B, Miró Pina V, Serrano Colome C,
1775 Weghorn D (2023) An extension of the Walsh-Hadamard
1776 transform to calculate and model epistasis in genetic
1777 landscapes of arbitrary shape and complexity. *bioRxiv*
1778 2023.03.06.531391. 1812
- 1779 35. Anderson DW, McKeown AN, Thornton JW (2015)
1780 Intermolecular epistasis shaped the function and evolution
1781 of an ancient transcription factor and its DNA binding sites.
1782 *eLife* 4:e07864. 1813
- 1783 36. Starr TN et al. (2020) Deep Mutational Scanning of
1784 SARS-CoV-2 Receptor Binding Domain Reveals Constraints
1785 on Folding and ACE2 Binding. *Cell* 182(5):1295–1310.e20. 1814
- 1786 37. Domingo J, Diss G, Lehner B (2018) Pairwise and
1787 higher-order genetic interactions during the evolution of a
1788 tRNA. *Nature* 558(7708):117–121. 1815
- 1789 38. Pokusaeva VO et al. (2019) An experimental
1790 assay of the interactions of amino acids from orthologous
1791 sequences shaping a complex fitness landscape. *PLoS Genet*
1792 15(4):e1008079. 1816
- 1793 39. Wu NC et al. (2020) Major antigenic site B of human
1794 influenza H3N2 viruses has an evolving local fitness landscape.
1795 *Nat Commun* 11(1):1–10. 1817
- 1796 40. Phillips AM et al. (2021) Binding affinity landscapes
1797 constrain the evolution of broadly neutralizing anti-influenza
1798 antibodies. *eLife* 10:e71393. 1818
- 1799 41. Lite TV et al. (2020) Uncovering the basis of protein-
1800 protein interaction specificity with a combinatorially complete
1801 library. *eLife* 9:e60924. 1819
- 1802 42. Metzger BPH, Park Y, Starr TN, Thornton JW
1803 (2023) Epistasis facilitates functional evolution in an ancient
1804 transcription factor. *eLife* 12:RP88737. 1820
- 1805 43. Wagner A (2008) Neutralism and selectionism: a
1806 network-based reconciliation. *Nat Rev Genet* 9(12):965–974. 1821
- 1807 44. Shah P, McCandlish DM, Plotkin JB (2015) Contingency
1808 and entrenchment in protein evolution under purifying
1809 selection. *Proc Natl Acad Sci USA* 112(25):E3226–E3235. 1822
- 1810 45. Park Y, Metzger BPH, Thornton JW (2022) Epistatic
1811 drift causes gradual decay of predictability in protein evolu-
1812 tion. *Science* 376(6595):823–830. 1823
- 1813 46. Anderson DW, Baier F, Yang G, Tokuriki N (2021)
1814 The adaptive landscape of a metallo-enzyme is shaped by
1815 environment-dependent epistasis. *Nat Commun* 12(1):3867. 1824
- 1816 47. Jalal ASB et al. (2020) Diversification of DNA-Binding
1817 Specificity by Permissive and Specificity-Switching Mutations
1818 in the ParB/Noc Protein Family. *Cell Rep* 32(3):107928. 1825
- 1819 48. Weinreich DM, Delaney NF, Depristo MA, Hartl DL
1820 (2006) Darwinian evolution can follow only very few muta-
1821 tional paths to fitter proteins. *Science* 312(5770):111–114. 1826
- 1822 49. Aakre CD et al. (2015) Evolving new protein-protein
1823 interaction specificity through promiscuous intermediates.
1824 *Cell* 163(3):594–606. 1827
- 1825 1828
- 1826 1829
- 1827 1830
- 1828 1831
- 1829 1832
- 1830 1833
- 1831 1834
- 1832 1835
- 1833 1836
- 1834 1837
- 1835 1838
- 1836 1839
- 1837 1840
- 1838 1841
- 1839 1842
- 1840 1843
- 1841 1844
- 1842 1845
- 1843 1846
- 1844 1847
- 1845 1848
- 1846 1849
- 1847 1850
- 1848 1851
- 1849 1852
- 1850 1853
- 1851 1854
- 1852 1855
- 1853 1856
- 1854 1857
- 1855 1858
- 1856 1859
- 1857 1860