# Article

# Direct Calculation of Protein Fitness Landscapes through Computational Protein Design

Loretta Au[1,*] and David F. Green[2]
[1]Department of Statistics, The University of Chicago, Chicago, Illinois; and [2]Department of Applied Mathematics and Statistics, Stony Brook University, Stony Brook, New York

ABSTRACT   Naturally selected amino-acid sequences or experimentally derived ones are often the basis for understanding how protein three-dimensional conformation and function are determined by primary structure. Such sequences for a protein family comprise only a small fraction of all possible variants, however, representing the fitness landscape with limited scope. Explicitly sampling and characterizing alternative, unexplored protein sequences would directly identify fundamental reasons for sequence robustness (or variability), and we demonstrate that computational methods offer an efficient mechanism toward this end, on a large scale. The dead-end elimination and A* search algorithms were used here to find all low-energy single mutant variants, and corresponding structures of a G-protein heterotrimer, to measure changes in structural stability and binding interactions to define a protein fitness landscape. We established consistency between these algorithms with known biophysical and evolutionary trends for amino-acid substitutions, and could thus recapitulate known protein side-chain interactions and predict novel ones.

## INTRODUCTION

Protein mutagenesis studies can disentangle how native interactions in wild-type are functionally important, but incrementing the number of mutations for a variant results in a combinatorial expansion of the possible protein sequence space. Single mutant variants of a 350-amino-acid protein, for instance, would yield 6650 sequences, while changes as pairs or triplets would allow $>2.4 \times 10^7$ and $>5.7 \times 10^{10}$ unique sequences, respectively. The sheer magnitude of protein sequences raises many challenges for interpreting the role of primary structure in dictating protein structure and function, and although progress continues to be made toward this understanding, it remains incomplete. Existing methods offer a range of analytical results, varying in the type and number of sequences that are evaluated (Fig. 1 a). Comparative sequence analysis methods can measure sequence conservation, identify motifs, and evaluate evolutionary relationships of known, sequenced proteins (1–7), while primary structures that deviate away from biases of natural evolution can be created via mutagenesis protocols. As examples, alanine scanning replaces original amino-acid side chains with alanine (8–10), and even larger protein libraries are possible via directed evolution experiments (11–13), which can scale up to $10^{12}$ or more sequences for sampling; both approaches require additional resources for functional characterization. High sequence similarity by itself cannot guarantee that structural motifs or protein folds are shared (14–16), and this can affect how results derived

solely from sequence analysis should be interpreted. In contrast, mutagenesis studies may be more costly than comparative sequence analysis, but the protein expression and functional assays that accompany these methods provide a more comprehensive understanding of the biophysical requirements that are essential to sequence-function relationships. High-throughput and deep-sequencing methods for directed evolution have been improving (17–19), and continue to elucidate the functional requirements of protein fitness. However, financial and temporal costs may still impose some constraints, depending on problem size, which motivated our development of a resourceful computational approach that can still provide high-resolution data for analysis. In particular, our protocol methodically simulates mutant variants for computing the protein fitness requirements of a chosen protein system without selection bias, offering additional perspective to how the energetic landscape of sequence space is shaped.

The dead-end elimination and A* search algorithms (DEE/A*) were adapted here for large-scale in silico mutagenesis, and thus enabled us to explore protein sequences that would be inaccessible otherwise (Fig. 1 b; see Fig. S1 in the Supporting Material). By assessing all low-energy sequences and their corresponding structures, we could deconvolve the multiple contributions of wild-type amino acids to protein fitness, defined here as structure stabilizing and binding interactions. Our computational approach, demonstrated here for a G-protein heterotrimer, is applicable to any system. However, it requires a reliable structural template to define the wild-type sequence. Enhanced sampling of backbone conformations is also needed, to account for
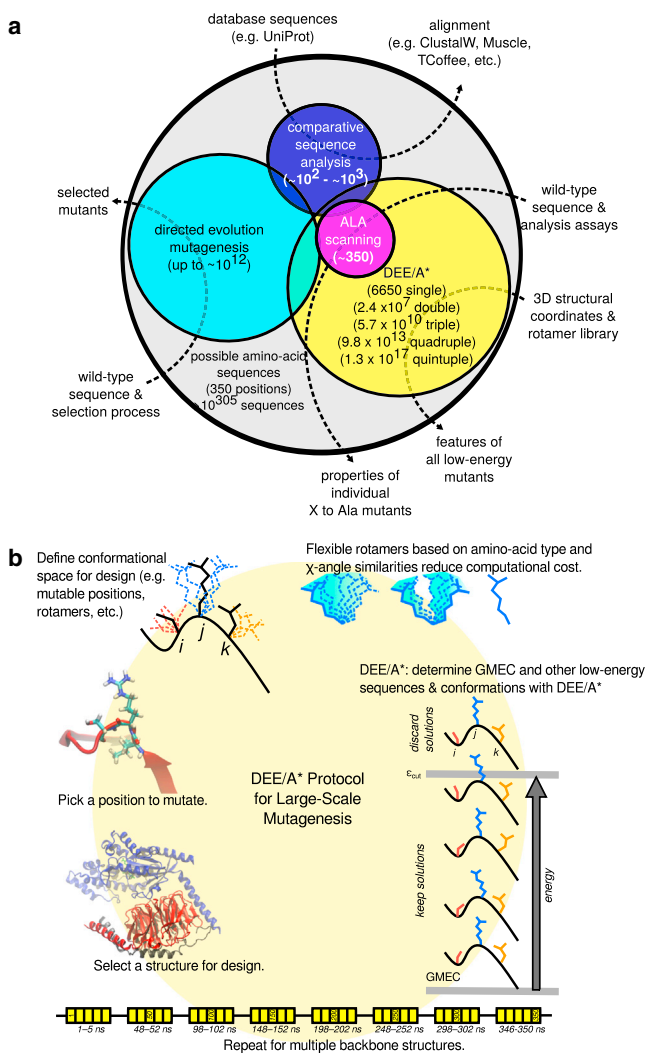
Editor: Amedeo Caflisch.

FIGURE 1 Systematic mutagenesis using computational protein design algorithms. (*a*) The mutant sequence space for a protein with 350 amino acids can become combinatorially large, and methods for exploring this space are shown here. Data for input and analytical results are indicated (*arched arrows* passing through the corresponding technique; note: *circles* are not drawn to scale). (*b*) From a molecular dynamics simulation of a given protein system, multiple snapshots are taken to create an ensemble of representative backbone structures. Each position in the protein may be substituted to any amino acid, and nearby side chains are flexible while others remain fixed. Dead-end elimination will discard rotamers that are incompatible with a low-energy structure for a given sequence, and the A* search will evaluate the combination of rotamers at all flexible positions that will yield the global minimum energy conformation (GMEC) and additional solutions below a designated energy cutoff ($\varepsilon_{\text{cut}}$), within the given constraints. To see this figure in color, go online.

slight variations in the protein microenvironment and measure the consistency of a mutational effect. This sampling was established using multiple conformations from a molecular dynamics simulation, but as an alternative, backbone flexibility could be accounted for using methods that introduce $\phi$- and $\psi$-angle perturbations to the backbone or by including multiple crystal structures of the protein (20–24).

We chose these algorithms because of their proven success in redesigning proteins to improve existing function or introduce novel ones (25–30). DEE evaluates amino acids and side-chain configurations that are incompatible with a low-energy protein conformation, based on the target protein structure, a rotamer library, and the energetic model used (31). Consequently, the number of possible structures is reduced; the lowest-energy protein conformation and additional ones within a designated energetic cutoff can then be identified from the remaining rotamers using the A* search algorithm, a heuristic, best first search that estimates the energetic cost of different rotamer combinations (Fig. 1 *b*) (32,33). Unlike a stochastic algorithm, the deterministic nature of DEE/A* guarantees the same solution every time, although it could be prone to completing an exhaustive search before doing so (34,35). A hierarchy of energetic models with increasing accuracy can be used to refine the solutions from DEE/A*: beginning with coarser pairwise decomposable approximations, high-energy sequences can be discarded early so that more intensive implicit or explicit solvent computations are performed on fewer molecules, reducing computational expense (36).

## MATERIALS AND METHODS

### Molecular dynamics setup

An all-atom molecular dynamics simulation was performed on the wild-type G-protein heterotrimer $G_i\alpha_1\beta_1\gamma_2$ (PDB: 1GP2) (37,38). CHARMM and NAMD were used for the simulation (39,40), with periodic boundary and NPT ensemble conditions ($P = 1$ atm, $T = 300$ K) using PARAM22/27 parameters (41,42), the TIP3P model for solvation (43), and a 2-fs time step. The structure was prepared using the REDUCE program to define initial protonation states (44), and hydrogen atom coordinates were determined using the HBUILD module in CHARMM (45). Randomly selected water molecules were replaced with sodium and chloride ions to establish a relevant physiological ionic strength (145 mM), with a minimum 10 Å distance between solute and the box edge. The structure was minimized after 240 steps and 200 ps of equilibration using Langevin dynamics in NAMD. A 12 Å cutoff was used for short-range interactions, while long-range electrostatic interactions were accounted for using the particle-mesh Ewald method. Snapshots were saved at every time step to ensure correlation with the same Boltzmann distribution throughout the simulation.

### DEE/A* parameters

Forty wild-type conformations from the molecular dynamics trajectory were selected for analysis: the first 5 ns, the last 5 ns, and six additional intervals between them, each spanning 5 ns with a midpoint that was a multiple of 50 ns (Fig. S1). Side-chain orientations were defined using the original Dunbrack-Karplus rotamer library, augmented before use by adding $\pm 10°$ to each $\chi_1$- and $\chi_2$-angle for enhanced sampling (46).

We applied the generalized-Born implicit solvent model with switching from CHARMM (47), after preliminary pruning using a distance-dependent dielectric ($\varepsilon = 4r$) (48). A flexible rotamer model was also used to discard unfavorable orientations quickly, by averaging together rotamers with similar $\chi_1$- and $\chi_2$-angles, reducing the size of the conformational space searched by DEE/A* (49). Each position in the wild-type sequence was mutated, and all sequences within 30 kcal/mol from the global minimum energy conformation were kept and referenced to the corresponding

wild-type energy. Wild-type side chains within 5 Å of any GDP atom were included in the analysis of Gα-GDP interactions; not all side chains will interact with this ligand throughout the simulation due to backbone fluctuations, and free energy data were normalized accordingly.

## Computing protein fitness of each mutant sequence

Protein fitness was measured as a combination of structural stability and binding interactions, either between Gα with GDP or Gα with the $\beta\gamma$-heterodimer. Structural stability is defined here as the energetic difference between an amino-acid side chain in the context of a folded structure and the reference state, in which the amino acid is isolated, then N-acetylated and N-methylamidated at the N- and C-termini, respectively; binding is defined energetically as the difference between the folded protein bound to its interaction partner and the same folded protein, unbound. A 500-Å rigid-body translation was used to separate binding partners to compute unbound-state energies. A Boltzmann-weighted average at an effective temperature of 4500 K was computed to represent the overall effect of each mutation (see the Supporting Material). The use of a discrete rotamer library and discretely sampled protein backbones leads to an exaggeration of unfavorable energies and the neglect of conformational entropy terms, which often partially compensate enthalpic terms, can further overstate the energetics. The use of an elevated effective temperature accounts for some of this exaggeration, albeit in an ad hoc manner.

Energy minimization of each DEE/A* result was performed using a Newton-Rhapson algorithm in CHARMM for 4000 steps each to identify shortcomings in the rotamer library. Mutational free energy was computed by decomposing amino acids into the amino-, carboxyl-, and variable side chain (starting from C$\beta$) groups, and the energetic difference was computed against a hydrophobic isostere of the wild-type amino acid. Similarity matrices were constructed using theoretical amino-acid probabilities found in the wild-type $G_i\alpha_1\beta_1\gamma_2$, and counting the number of sequences that survive an energetically defined evolutionary pressure. Frequency of substitution from amino acid $i$ to $j$, $e_{ij}$, was computed using these counts, and compared to the expected frequency found in PAM120 and BLOSUM62 (see the Supporting Material).

## Measuring amino-acid substitution rates

Entries in any PAM or BLOSUM matrix is a score, on a half-bit scale, that indicates the probability of observing substitutions to wild-type amino acid $i$ with amino-acid $j$, $S_{ij}$. In a given set of protein sequences or within an aligned region of sequences (depending on the type of substitution matrix computed), the observed frequency of finding $i$ substituted by $j$, $e_{ij}$, is compared to a corresponding theoretical probability that the amino-acid exchange may happen, $p_i$, $p_j$ (where $p_i$ and $p_j$ are the natural, independent frequencies of occurrence for amino acids $i$ and $j$, respectively) and thus $S_{ij} = 2 \log_2(e_{ij}/p_ip_j)$. For comparison against these evolution-based observations, we defined $e_{ij}$ as the number of DEE/A* sequences that simultaneously satisfied the 1.5 kcal/mol cutoff for structural stability and binding interactions after mutation of amino acid $i$ into $j$ (sequences that survived DEE/A* fitness pressures). Algebraically, scores from PAM and BLOSUM matrices can be converted to $e_{ij}$ for comparison, because $e_{ij} = p_ip_j2^{(S_{ij}/2)}$; the values for $S_{ij}$ were provided by PAM120, BLOSUM62, or a randomly generated matrix, and wild-type amino-acid distributions of the entire heterotrimer were used to define $p_i$ and $p_j$ accordingly (see the Supporting Material).

## Statistical analysis for predictions

The Mann-Whitney-Wilcoxon statistical test was implemented using the *exactRankTests* library from the R statistical package. Neutral mutations were defined as changes from wild-type within a $-1.5$ and 1.5 kcal/mol range, and thus these values were set to zero before this analysis. An exact test was chosen to account for ties, and the null hypothesis (a zero vector, indicating no changes due to substitution) was compared to the empirical data collected for each position, a 20-dimensional vector representing the 20 possible amino acids underlying the cumulative distribution function (see the Supporting Material). A low $p$-value in these calculations suggests strong evidence that the position is mutationally sensitive for the aspect of fitness evaluated. Side-chain positions with known binding interactions were separated from all others to measure the true-positive rate of DEE/A*, based on a cutoff value, $p = 0.05$; this cutoff was also used as the premise for identifying additional side chains involved in binding interactions.

## Computational resources

All DEE/A* mutations for an individual mutation were performed on a single 3.4 GHz Intel Pentium IV Xeon processor; most positions required ~4–5 h of computing time. There are 685 mutable positions in 1GP2, and using a cluster with 235 processing nodes, an average of 48 h was required to perform mutagenesis at all positions. Mutation free energy calculations required ~30 min of computing time on the same cluster.

## RESULTS

Heterotrimeric G-proteins are ubiquitous in eukaryotic cell-signaling pathways, and we have chosen this as a model system for our approach, with $G_i\alpha_1\beta_1\gamma_2$ (PDB: 1GP2) as our wild-type reference (37,50,51). This family of proteins has been shown to have unique patterns in interaction specificity between subunits that enable complex formation and biological function (52–55). As determinants of broad-spectrum biological function, we have focused on 1) the structural stability of the complete protein complex, 2) the binding interactions between the $\beta\gamma$-heterodimer and Gα, and 3) the binding of Gα to GDP (see Materials and Methods).

## Complete mutagenesis profiles calculated from using DEE/A*

Many mutations have a neutral effect on the protein (Figs. S2–S4; Tables S1–S3), but there is a tendency for mutations to be less favorable than wild-type. Approximately two-thirds of the sequences explored by DEE/A* are destabilizing to the wild-type structure, and greater energetic variance is seen in these sequences than those measured for changes in binding interactions (Fig. 2). This is due to both having fewer amino acids involved in binding (compared to stabilization), and having a broad range of microenvironments, from hydrophobic to highly solvent-exposed, available in the folded protein. A complete sequence profile for every position was established for our model system, identifying specific regions of unfavorable amino-acid substitution and highlighting those that are less sensitive to mutation (Figs. S5–S8, S10, and S11). Positions with several allowable and favorable substitutions usually have fewer geometric or electrostatic constraints; when very diverse functional groups cannot be accommodated at a position, it suggests

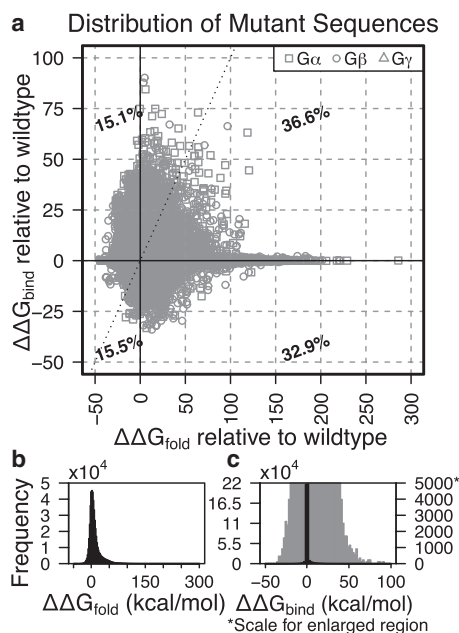## a  Distribution of Mutant Sequences



**b**

**c**

FIGURE 2 Protein fitness landscape for mutant sequences. Sequences are mapped according to energy relative to the wild-type sequence for structural stability ($\Delta\Delta G_{fold}$) and binding interactions ($\Delta\Delta G_{bind}$) in (*a*), and the relative proportions in each quadrant of this landscape are shown as percentages. The distribution of mutant sequences according to each aspect of protein fitness is shown in the bottom row. (*b*) For stability, there is a heavy tail in the distribution of $\Delta\Delta G_{fold}$, which indicates that most mutants are less stable than wild-type. (*c*) As for binding interactions, the value of $\Delta\Delta G_{bind}$ is often ~0 kcal/mol (shown in *black*, *y* axis on *left*), but a closer look at this histogram reveals that the distribution is skewed (shown in *gray*, *y* axis on *right*), with more sequences having a positive $\Delta\Delta G_{bind}$ value.

that unique side-chain interactions exist in the region and are required to maintain protein fitness.

## Conformational space adequately modeled with rotamer library and protein structures

The movement of protein side chains is most accurately simulated using small χ-angle perturbations, but amino-acid orientations have been examined and statistically determined to favor specific combinations of χ-angles, establishing the basis for rotamer libraries (56,57). Without this discretization, DEE/A* simply cannot work—the algorithm evaluates unique combinations of side-chain placement on a given backbone structure, and rotamer libraries provide clear definitions of these possibilities in a limited number. As an alternative, energy minimization algorithms can be used to find favorable side-chain orientations that may be unlisted in such libraries, and can work well when the number of structures needed for analysis is not overwhelming. A comparison was made between all DEE/A* sequences sharing the same backbone orientation and their corresponding minimized structures (each starting as a DEE/A* solution) to assess the influence of a discretized ro-

tamer space; a positive linear correlation was found between them (Fig. S15), suggesting energetic similarity despite slight differences in side-chain positioning. Approximately 60% of the data is found to be energetically unfavorable using DEE/A* and remains unfavorable after applying energy minimization, while roughly 20% of the data is favorable in both calculations. As the reference structure is derived from the wild-type sequence, a bias is found (and expected) toward mutant structures becoming more energetically favorable from using the minimization algorithm. However, as both wild-type and mutant sequences were minimized, relative energies from DEE/A* could be better than those from minimization (Fig. S15). Large discrepancies between the two methods of calculation (>20 kcal/mol) tended to involve substitutions to charged side chains or involve geometric constraints: in one particular β-sheet (GβAsp[247], GβThr[249], and GβArg[251]), nonaliphatic amino acids were disfavored to preserve directionality, hydrogen-bond interactions, and the size of ($i, i + 2$) side chains (58), suggesting that finer sampling could be beneficial in specific side-chain packing contexts (Fig. S17). Even so, ~80% of the sequences were within ± 5 kcal/mol of the alternative energy calculation, indicating a satisfactory evaluation of most sequences using DEE/A* without additional energy minimization (Fig. S16). The wild-type protein structures used were representative of major changes or fluctuations that the complete heterotrimer may undergo during simulation. The energetic variance of each mutant sequence was measured as the number of states in the structural ensemble increased (Fig. S12), and consistency in free energy was established first for sequences from densely packed, hydrophobic regions of the protein. Structural constraints within these regions were further reflected in the number of unsuitable mutations at these positions (Figs. S13 and S14; Table S4). In contrast, a greater number of backbone conformations was necessary to capture structural features of loops and other flexible regions of the heterotrimer, due to greater degrees of freedom.

## Amino-acid functional roles can be disentangled

Energetic profiles were created separately for stability and binding interactions using the free energy of all mutant sequences (see Materials and Methods). By measuring these two aspects of fitness independently, functional trade-offs in the protein could be identified, as demonstrated by the GDP-binding pocket of Gα (Fig. 3). If either requirement for stability or binding was not satisfied, the overall fitness of the protein was worse than wild-type—the maximum energy of either stability or binding, max($\Delta\Delta G_{fold}$, $\Delta\Delta G_{bind}$), could distinguish this for a given mutant. Asp[150], for example, could be easily replaced by most amino acids and remain functional, because the native orientation points the carboxyl group away from GDP, despite being near a guanine nitrogen (Fig. 3 *a*). Nearly all substitutions could
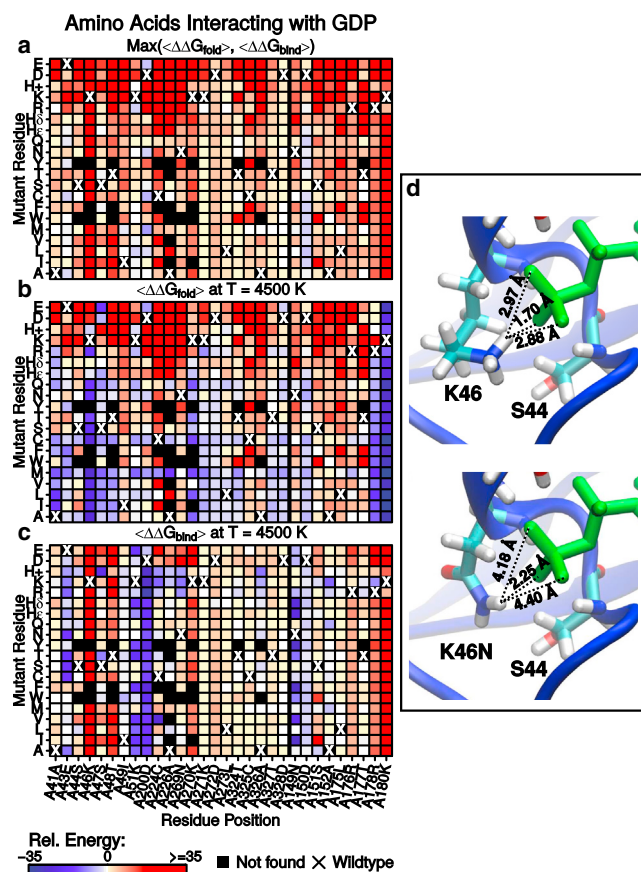
**FIGURE 3** Gα interactions with GDP. Side chains within 5 Å of GDP can have functional trade-offs as substitutions are made to wild-type. Energy referenced to wild-type is shown for each mutant. A black box for a sequence indicates that the corresponding structure had steric clashes, and was thus not found by DEE/A*. Black boxes with a large X highlight the wild-type amino acid. (a) Evaluating $\max(\langle\Delta\Delta G_{fold}\rangle, \langle\Delta\Delta G_{bind}\rangle)$ reveals difficulty in simultaneously satisfying both fitness criteria. (b) Average stability, $\langle\Delta\Delta G_{fold}\rangle$, and (c) average binding interactions, $\langle\Delta\Delta G_{bind}\rangle$, for Gα-GDP indicate varying degrees of mutational sensitivity at different positions. (d) Mutations often alter the proximity of important interactions, as seen in K46, in which hydrogen bonds are lost in K46N. To see this figure in color, go online.

improve protein stabilization for Ala[41], Lys[46], Lys[270], and Lys[180] in Gα, but the same mutations were poor candidates for binding GDP (Fig. 3 b). Similarly, side-chain substitutions in the same subunit could improve binding interactions relative to wild-type at positions Ile[49], Asp[200], and Asn[149], but doing so would generally destabilize the α-subunit (Fig. 3 c). Bulkier, aromatic amino acids did not fit well in this region, and charged side chains were also poor candidates because of their electrostatic requirements. These general trends were exhibited throughout the heterotrimer, and functional trade-offs were only a concern for the small number of positions present at the protein-binding interface or involved in protein-ligand interactions (Figs. S7 and S8). Most positions were sensitive to substitution, and this is expected for a highly evolved protein family. Gα positions at

the amino terminus or in switch II (residues 202–209) have the greatest energetic variation after mutation than other positions in the subunit, and these regions were known to interact with the β-subunit when inactive (50,59,60). Gβ, an example of a WD40 β-propeller protein, has positions at the binding interface that also show a similar trend, and where stability is lost after mutation is consistent with our expectations of the WD40 protein family (Fig. S9) (50,61).

## Mutant structures from DEE/A* provide practical computational models

Alternative methods for studying wild-type contributions to protein stability, some of which require significantly fewer computational resources than DEE/A*, do exist. Amino acids may be decomposed according to functional groups, for instance, so that the energy required to convert a side chain into its hydrophobic isostere can be measured, and this mutational free energy elucidates any underlying electrostatic interactions (62–66). Such calculations can be completed in ~30 min for a single heterotrimer, while DEE/A* would require ~48 h for the same system using the same computing cluster. The expense of using DEE/A* is well compensated for, however—all 20 amino-acid choices are evaluated when finding low-energy sequences and simultaneously modeling tertiary structures. Having up to 19 mutant variants thus provides multiple frames of reference for assessing how tolerant a wild-type side chain can be to different kinds of mutation. The outcomes also include visual examples of less intuitive substitutions and energetic data that can help rank mutational effects or quantify the mutational robustness of wild-type. Energetic comparisons made with hydrophobic isosteres has its advantages in efficiency, but relies on an artificial construction that is not found in biology; DEE/A* offers practical models in its representations of actual amino acids.

To illustrate the compatibility between these two kinds of calculations, and their differences, the mutational free energy of all positions involved in Gα-GDP interactions were compared to the sequences from DEE/A* (Fig. 4). Each aspect of fitness was treated independently for assessment; the number of stable states found (defined by energetic cutoff) and the mutant sequence energies distributed were compared to mutational free energies computed using hydrophobic isosteres; positions could be separated easily according to mutational robustness in this way. Lys[46] had negative mutational free energy, an indication that important interactions were made by this side chain to bind GDP, but from DEE/A*, we could understand that only wild-type would ever make these contributions—no other substitutions are allowed here. Conversely, we found that electrostatic contributions of Glu[43] were also important in the wild-type, but all mutations were allowed and tended to
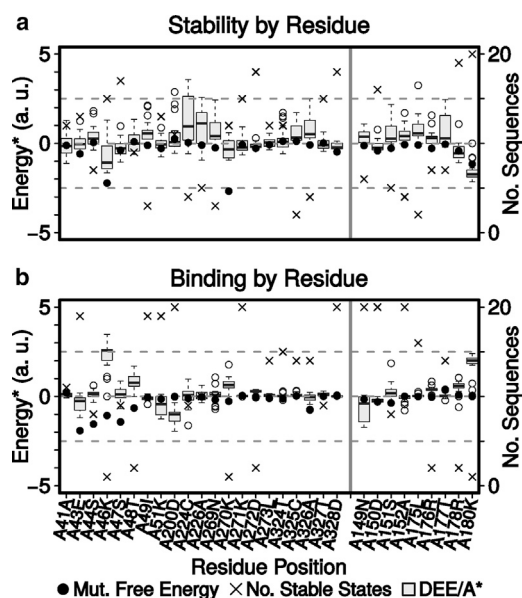
FIGURE 4 Energetic contributions of amino acids. Electrostatic calculations were performed for Gα amino acids that interact with GDP. Box-and-whisker plots represent the energetic distribution of all mutant sequences from DEE/A* at the specified position, and are overlaid onto mutational free energy data for (a) structural stability and (b) binding interactions. Arbitrary units (a.u.) were used for the y axes on the left at a scales: (i) 6 kcal/mol for stability mutation free energy, (ii) 1.6 kcal/mol for binding mutation free energy, and (iii) 16 kcal/mol for DEE/A* results in both contexts; the respective energetic ranges are thus (i) [−30,30] kcal/mol, (ii) [−8,8] kcal/mol, and (iii) [−80,80] kcal/mol. The number of sequences found at each position from DEE/A* are all those ≤1.5 kcal/mol from the wild-type energy; these quantities were marked with an X and follow the y axes on the right.

be more favorable than wild-type. Lys[180] could be substituted by anything to improve stability, and also have important native interactions; however, substitutions adversely affected binding interactions with GDP. All mutations were disallowed, even though the wild-type amino acid had little impact on binding, again demonstrating that geometry or interactions with neighboring residues play an important secondary role.

## DEE/A* substitutions are strongly correlated with known amino-acid exchanges

Finally, the overall DEE/A* substitution rates were compared with the PAM120 and BLOSUM62 similarity matrices to measure how well DEE/A* (and our choice of folding and binding as measures of fitness) can reflect protein evolutionary pressures. Each PAM and BLOSUM matrix accounts for a broad range of sequence evolution, and are standard matrices for use in sequence alignments (67–69). Energetically favorable DEE/A* sequences were used to derive the expected frequency of substituting amino acid $i$ with $j$, $e_{ij}$, for comparison to analogous values of $e_{ij}$ using PAM or BLOSUM (see the Supporting Material). We defined protein fitness to depend on a combination of
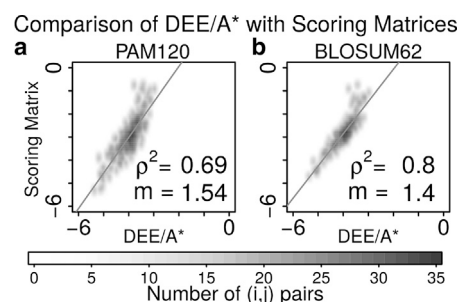


FIGURE 5 Comparison of DEE/A* and evolutionary substitution frequencies. Expected frequencies of substitution for any $(i, j)$ amino-acid pair were computed based on the number of mutant sequences satisfying the 1.5-kcal/mol cutoff. Amino-acid probabilities from wild-type provide a basis for deriving substitution rates for comparison to (a) PAM120 and (b) BLOSUM62.

protein stabilization and binding ability, at unknown proportions (Figs. S18–S20); a balanced weighting of both fitness criteria optimized the correlation between our DEE/A* data and a chosen similarity matrix, with $\rho^2 \approx 0.7$ and $\rho^2 \approx 0.8$ for PAM and BLOSUM, respectively (Fig. 5). To determine whether these correlations were meaningful, our DEE/A* data were also compared to random matrices constructed 1) from a uniform distribution bounded by the maximum and minimum scores of both PAM120 and BLOSUM62, and 2) by permuting the entries of each similarity matrix. Correlation was generally poor ($\rho^2 \approx 0.1$) between DEE/A* and a completely arbitrary matrix; correlation between PAM120 (or BLOSUM62) with a permuted version of itself was first established ($\rho^2 \approx 0.3$ for PAM120, and $\rho^2 \approx 0.5$ for BLOSUM62), then DEE/A* data were compared to the randomized version of each matrix and showed slight improvement ($\rho^2 \approx 0.4$ for PAM120 and $\rho^2 \approx 0.6$ for BLOSUM62) (Fig. S22). These data suggest that the general distribution of values in the DEE/A*-derived matrices is similar to that in the PAM and BLOSUM matrices (it is this that leads to nonzero correlation between randomized matrices, but there are deeper similarities in the detailed structure of the matrices; Figs. S21 and S22; Tables S5 and S6). Comparison to alternative PAM and BLOSUM did not yield any statistically significant differences, due to low variance between different versions of PAM and BLOSUM scores overall (not shown).

## Compatibility between alanine mutations from DEE/A* and thermal stability experiments

A full alanine scan was performed by Sun et al. (70) to understand the role of native interactions in stabilizing Gα, and we used these data for comparison with mutant alanine sequences from our DEE/A* calculations. In their experiment, each wild-type Gα residue was systematically mutated in the α-subunit to alanine (and wild-type alanine to glycine) for GDP- and GTP-bound states. The change in thermal stability

($\Delta T_m = T_{\text{mut}} - T_{wt}$) was measured for each single mutant, and a threshold of $\Delta T_m \leq 2°C$ was proposed by Sun et al. (70) for defining a destabilizing mutation. The analogous description from our DEE/A* calculations would be a mutation in which $\Delta\Delta G_{\text{fold}} > +1.5$ kcal/mol. We directly compared alanine substitutions from DEE/A* with the data provided by Sun et al. (70) using these two interpretations for destabilizing mutations (Fig. S23; see Table S7). We considered an unfavorable alanine mutation to be a positive outcome, and by these measurements, the sensitivity of DEE/A* was computed to be 0.52, while its specificity was 0.78. Energetic data for alanine mutants were then randomized and reassessed to establish a quantitative reference for these values, and we measured sensitivity and specificity to be $0.32 \pm 0.04$ and $0.67 \pm 0.02$, respectively, after 5000 independent trials (Table S8). Compared to all random trials performed, our DEE/A* calculations correctly classified important side-chain interactions (true-positives) and positions that were insensitive to mutation (true-negatives) at a consistently higher rate than any of the randomized cases (Fig. S24). (Consequently, the number of false-positives and false-negatives were both much lower than the randomized trials.) For all G$\alpha$ side chains, 75% (162 positions) were predicted to be mutationally insensitive by both our DEE/A* calculations and the data provided by Sun et al. (70) (Table S7). Based only on the structural data from computational simulations, the role of native interactions was correctly determined in ~68% of G$\alpha$ (the total number of true-positives and true-negatives) using only alanine substitutions. Although DEE/A* cannot perfectly replicate in cyto conditions and related assays, these statistics were strong indicators that DEE/A* can reasonably predict the importance of wild-type interactions.

## Consistency between DEE/A* and known point mutations

Oncogenic point mutations are available in public databases, such as COSMIC and cBioPortal (71–73), and several were found for the GNB1 gene, which encodes G$\beta_1$. In addition to compiling a complete list of these single mutants, Yoda et al. (74) discovered a few additional ones in their experiment that explored cancerous mutations affecting the $\beta$-subunit. These results from Yoda et al. (74) were used for comparison with corresponding DEE/A* mutants (Table S9). Either gain-of-function or loss-of-function mutations could be oncogenic, and both of these possibilities were considered in our comparison with DEE/A*. Furthermore, lethal mutations could affect function by altering heterotrimer stability or by disrupting proper association between G$\alpha$ and the $\beta\gamma$-heterodimer, but the distinction between these two mechanisms is not always known from the available data. Thus, an energetic definition that accounted for both gain-of-function and loss-of-function mutations, and the possible contexts for which mutations may

affect heterotrimer function, would be the maximum magnitude of either protein stability or binding interactions: $\max(|\Delta\Delta G_{\text{fold}}|, |\Delta\Delta G_{\text{bind}}|)$. A value $>1.5$ kcal/mol in our DEE/A* calculations would indicate either an activating or deactivating mutation. We found a positive correlation between our computational results and the set of known point mutations: of the 36 single mutants available for comparison, only three of them had a neutral change after mutation (both $|\Delta\Delta G_{\text{fold}}| \leq 1.5$ kcal/mol and $|\Delta\Delta G_{\text{bind}}| \leq 1.5$ kcal/mol (see the Supporting Material). The remaining 33 mutations (92%) were either activating or deactivating mutations in at least one of the fitness contexts. These results further demonstrated that our computational approach can capture important trends in mutational effects found in biological systems.

## High predictive ability of computational results

The Mann-Whitney-Wilcoxon statistical test was used for evaluating quantitative differences between mutational effects, and to assess mutational sensitivity of a wild-type side chain (see the Materials and Methods.) Binding interactions between subunits are documented for 49 positions (in total) of G$\alpha$ and G$\beta$ (50), and 38 (~78%) of these were detected by DEE/A* (Fig. 6 a), based on a threshold of $p = 0.05$. Five additional positions could be included in this count, if the threshold were adjusted to 0.10 instead, accounting for ~87% of known positions. Discrepancies for false-negatives from the remaining 11 positions (6 at threshold $p = 0.10$) were likely due to differences in conformational sampling. These positions tended to be in highly flexible protein regions: near the switch region of G$\alpha$ and near the N-terminal helix of G$\alpha$ (Fig. 6 b). The original observed interactions were established using x-ray crystallographic data, while our computational results were based on an ensemble of structures for the heterotrimer. Having this distinction for our data made it possible to computationally determine a broader regime of side chains involved in protein binding: 30 additional side chains were predicted by DEE/A* using the same statistical analysis (see Materials and Methods). The same metrics could be extrapolated for understanding structurally stabilizing interactions, and nearly all positions were found to have some significant contribution (Fig. S25; Tables S10 and S11). The molecular requirements for structural stability are not necessarily interchangeable with binding requirements, however, and further experimentation would be needed to verify the predictive ability of DEE/A* for this aspect of fitness and to separate side chains that are fundamental to stabilizing tertiary structure from less influential ones.

## DISCUSSION

Proteins must satisfy a number of conditions, including the ability to stably form an appropriate fold and associate with
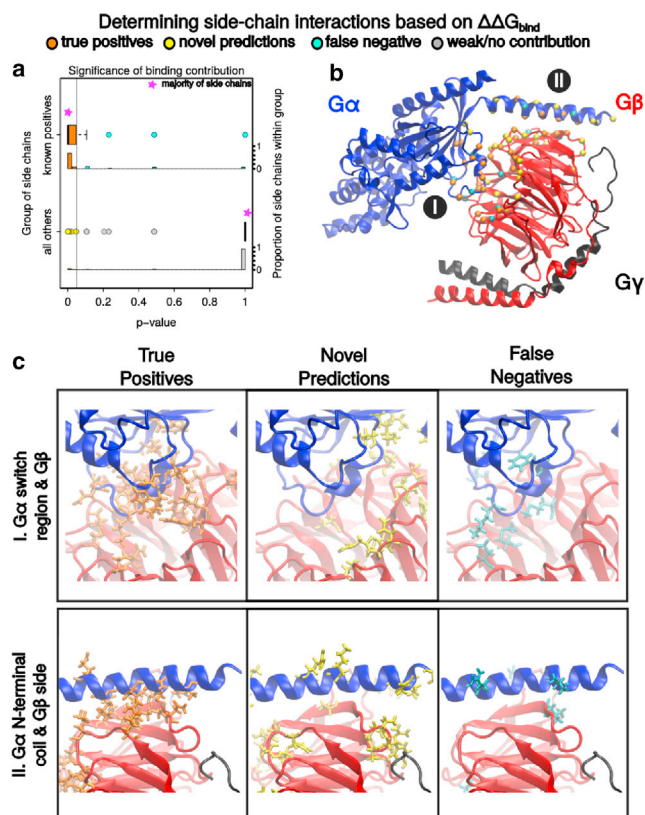
FIGURE 6  Recapitulation and prediction of side-chain interactions. Statistical differences based on all mutant $\Delta\Delta G_{bind}$ values were measured with the Mann-Whitney-Wilcoxon statistical test. (*a*) The distribution of *p* values for positions known to make binding interactions and all other positions are shown. (*Vertical line*) $p = 0.05$ for visual purposes. (*Magenta stars*) Nearly all positions of a subgroup that are found in the *p*-value distribution. (*b*) Positions for true-positives (*orange*, $p \leq 0.05$), novel predictions (*yellow*, $p \leq 0.05$), and false-negatives (*cyan*, $p \leq 0.05$) are mapped as spheres onto the heterotrimer for reference. (*Spheres*) The $\alpha$-carbon positions of each residue. (*c*) Structural examples of where true-positive, novel predictions, and false-negatives are typically found. To see this figure in color, go online.

various cognate binding targets, to be biologically functional and overcome different selection pressures. Many advances in high-throughput methods have considerably improved how protein sequence-function relationships can be studied, but the volume of possible sequence space remains inevitably greater. Our DEE/A* protocol provides a mechanism for studying mutations on a very large scale to help mend a part of this disparity, and provide a perspective that is different, although complementary, to these approaches. By analyzing novel sequence variants systematically, the energetic landscape of a protein was computed, and the functional role of each amino acid could be deconvolved. The performance of DEE/A* also relies on different resources than existing methods: when the number of sequences for alignment is inadequate or the evolutionary history of a given protein is not well understood, a thorough analysis of protein structural stability and binding interac-

tions is still possible. Whether alone or combined with existing methods, this level of detail can provide a better understanding of how protein design or engineering goals can be met.

Depending on the problem being considered, variations to our approach could be made to improve modeling details and computational accuracy. Longer molecular dynamics simulations and/or a greater number of representative structures would enhance sampling in highly flexible proteins or regions, for instance. Additional solvation models could also be used, e.g., by passing important sequences found using implicit-solvent models onto explicit-solvent simulations, to provide a more detailed explanation of electrostatic interactions. Furthermore, epistatic relationships from pairwise interactions and amino-acid covariation could be studied in depth using DEE/A* by introducing multiple mutations into the sequence at once. While these modifications may provide a better picture of how mutagenesis affects the wild-type protein, improvements in capturing amino-acid substitutions that completely reflect evolutionary biology might never be possible—evolutionary fitness pressures extend far beyond what can be measured by energetic change. Despite this, modeling side-chain substitutions with DEE/A* has shown consistency with structural studies and binding assays. DEE/A* thus complements comparative sequence analysis methods very well: sequence conservation can be directly linked to measurable aspects of fitness, and regions of allowable sequence variation can be explained.

## CONCLUSIONS

The dead-end elimination and A* search algorithms simultaneously search over protein sequence and conformational spaces, and we have leveraged these to elucidate many sequence-function relationships in a heterotrimeric G-protein. By adapting these two algorithms to find all low-energy single mutants, the multiple roles amino acids play in overall protein fitness could be deconvolved as a function of mutational robustness. Large-scale mutagenesis using this computational approach is able to capture many biophysical features of side-chain substitutions, and these changes in the initial wild-type structure satisfy expectations based on preexisting structural and experimental studies. DEE/A* reveals several relationships among primary structure, structural stability, and protein function, enhancing the utility of techniques in comparative sequence analysis and extending the boundaries of accessible protein sequence space.

## SUPPORTING MATERIAL

Supporting Materials and Methods, twenty-five figures, and eleven tables are available at http://www.biophysj.org/biophysj/supplemental/S0006-3495(15)01212-6.

## AUTHOR CONTRIBUTIONS

L.A. and D.F.G. designed the experiments; L.A. performed the experiments; L.A. and D.F.G. contributed analytical tools; L.A. analyzed the data; and L.A. and D.F.G. wrote the article.

## ACKNOWLEDGMENTS

## REFERENCES

1. Thompson, J. D., T. J. Gibson, …, D. G. Higgins. 1997. The CLUSTAL_X windows interface: flexible strategies for multiple sequence alignment aided by quality analysis tools. *Nucleic Acids Res.* 25:4876–4882.

2. Notredame, C., D. G. Higgins, and J. Heringa. 2000. T-Coffee: a novel method for fast and accurate multiple sequence alignment. *J. Mol. Biol.* 302:205–217.

3. Edgar, R. C. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* 32:1792–1797.

4. Lockless, S. W., and R. Ranganathan. 1999. Evolutionarily conserved pathways of energetic connectivity in protein families. *Science.* 286:295–299.

5. Süel, G. M., S. W. Lockless, …, R. Ranganathan. 2003. Evolutionarily conserved networks of residues mediate allosteric communication in proteins. *Nat. Struct. Biol.* 10:59–69.

6. Magliery, T. J., and L. Regan. 2005. Sequence variation in ligand binding sites in proteins. *BMC Bioinformatics.* 6:240.

7. McLaughlin, R. N., Jr., F. J. Poelwijk, …, R. Ranganathan. 2012. The spatial architecture of protein function and adaptation. *Nature.* 491:138–142.

8. Cunningham, B. C., and J. A. Wells. 1989. High-resolution epitope mapping of hGH-receptor interactions by alanine-scanning mutagenesis. *Science.* 244:1081–1085.

9. Massova, I., and P. A. Kollman. 1999. Computational alanine scanning to probe protein-protein interactions: a novel approach to evaluate binding free energies. *J. Am. Chem. Soc.* 121:8133–8143.

10. Weiss, G. A., C. K. Watanabe, …, S. S. Sidhu. 2000. Rapid mapping of protein functional epitopes by combinatorial alanine scanning. *Proc. Natl. Acad. Sci. USA.* 97:8950–8954.

11. Sidhu, S. S., and S. Koide. 2007. Phage display for engineering and analyzing protein interaction interfaces. *Curr. Opin. Struct. Biol.* 17:481–487.

12. Ernst, A., D. Gfeller, …, S. S. Sidhu. 2010. Coevolution of PDZ domain-ligand interactions analyzed by high-throughput phage display and deep sequencing. *Mol. Biosyst.* 6:1782–1790.

13. Araya, C. L., D. M. Fowler, …, S. Fields. 2012. A fundamental protein property, thermodynamic stability, revealed solely from large-scale measurements of protein function. *Proc. Natl. Acad. Sci. USA.* 109:16858–16863.

14. Cordes, M. H. J., R. E. Burton, …, R. T. Sauer. 2000. An evolutionary bridge to a new protein fold. *Nat. Struct. Biol.* 7:1129–1132.

15. Newlove, T., J. H. Konieczka, and M. H. J. Cordes. 2004. Secondary structure switching in Cro protein evolution. *Structure.* 12:569–581.

16. van Dorn, L. O., T. Newlove, …, M. H. J. Cordes. 2006. Relationship between sequence determinants of stability for two natural homologous proteins with different folds. *Biochemistry.* 45:10542–10553.

17. Fowler, D. M., C. L. Araya, …, S. Fields. 2010. High-resolution mapping of protein sequence-function relationships. *Nat. Methods.* 7:741–746.

18. Fowler, D. M., J. J. Stephany, and S. Fields. 2014. Measuring the activity of protein variants on a large scale using deep mutational scanning. *Nat. Protoc.* 9:2267–2284.

19. Stiffler, M. A., D. R. Hekstra, and R. Ranganathan. 2015. Evolvability as a function of purifying selection in TEM-1 β-lactamase. *Cell.* 160:882–892.

20. Harbury, P. B., J. J. Plecs, …, P. S. Kim. 1998. High-resolution protein design with backbone freedom. *Science.* 282:1462–1467.

21. Davis, I. W., W. B. Arendall, 3rd, …, J. S. Richardson. 2006. The backrub motion: how protein backbone shrugs when a sidechain dances. *Structure.* 14:265–274.

22. Georgiev, I., and B. R. Donald. 2007. Dead-end elimination with backbone flexibility. *Bioinformatics.* 23:i185–i194.

23. Smith, C. A., and T. Kortemme. 2008. Backrub-like backbone simulation recapitulates natural protein conformational variability and improves mutant side-chain prediction. *J. Mol. Biol.* 380:742–756.

24. Smith, C. A., and T. Kortemme. 2011. Predicting the tolerated sequences for proteins and protein interfaces using RosettaBackrub flexible backbone design. *PLoS One.* 6:e20451.

25. Dahiyat, B. I., and S. L. Mayo. 1997. De novo protein design: fully automated sequence selection. *Science.* 278:82–87.

26. Shimaoka, M., J. M. Shifman, …, T. A. Springer. 2000. Computational design of an integrin I domain stabilized in the open high affinity conformation. *Nat. Struct. Biol.* 7:674–678.

27. Bolon, D. N., and S. L. Mayo. 2001. Enzyme-like proteins by computational design. *Proc. Natl. Acad. Sci. USA.* 98:14274–14279.

28. Sarkar, C. A., K. Lowenhaupt, …, D. A. Lauffenburger. 2002. Rational cytokine design for increased lifetime and enhanced potency using pH-activated "histidine switching". *Nat. Biotechnol.* 20:908–913.

29. Looger, L. L., M. A. Dwyer, …, H. W. Hellinga. 2003. Computational design of receptor and sensor proteins with novel functions. *Nature.* 423:185–190.

30. Bolon, D. N., R. A. Grant, …, R. T. Sauer. 2005. Specificity versus stability in computational protein design. *Proc. Natl. Acad. Sci. USA.* 102:12724–12729.

31. Desmet, J., M. De Maeyer, …, I. Lasters. 1992. The dead-end elimination theorem and its use in protein side-chain positioning. *Nature.* 356:539–542.

32. Leach, A. R., and A. P. Lemon. 1998. Exploring the conformational space of protein side chains using dead-end elimination and the A* algorithm. *Proteins.* 33:227–239.

33. Pierce, N. A., J. A. Spriet, and S. L. Mayo. 2000. Conformational splitting: a more powerful criterion for dead-end elimination. *J. Comput. Chem.* 21:999–1009.

34. Desjarlais, J. R., and N. D. Clarke. 1998. Computer search algorithms in protein modification and design. *Curr. Opin. Struct. Biol.* 8:471–475.

35. Voigt, C. A., D. B. Gordon, and S. L. Mayo. 2000. Trading accuracy for speed: a quantitative comparison of search algorithms in protein sequence design. *J. Mol. Biol.* 299:789–803.

36. Green, D. F. 2010. A statistical framework for hierarchical methods in molecular simulation and design. *J. Chem. Theory Comput.* 6:1682–1697.

37. Wall, M. A., D. E. Coleman, …, S. R. Sprang. 1995. The structure of the G protein heterotrimer G$_i$α1β1γ2. *Cell.* 83:1047–1058.

38. Carrascal, N., and D. F. Green. 2010. Energetic decomposition with the generalized-Born and Poisson-Boltzmann solvent models: lessons from association of G-protein components. *J. Phys. Chem. B.* 114:5096–5116.

39. Brooks, B. R., C. L. I. Brooks, 3rd, …, M. Karplus. 2009. CHARMM: the biomolecular simulation program. *J. Comput. Chem.* 30:1545–1614.

40. Phillips, J. C., R. Braun, …, K. Schulten. 2005. Scalable molecular dynamics with NAMD. *J. Comput. Chem.* 26:1781–1802.

41. MacKerell, A. D. J., J. Wiórkiewicz-Kuczera, and M. Karplus. 1995. An all-atom empirical energy function for the simulation of nucleic acids. *J. Am. Chem. Soc.* 117:11946–11975.

42. MacKerell, A. D., D. Bashford, …, M. Karplus. 1998. All-atom empirical potential for molecular modeling and dynamics studies of proteins. *J. Phys. Chem. B.* 102:3586–3616.

43. Jorgensen, W. L., J. Chandrasekhar, …, M. L. Klein. 1983. Comparison of simple potential functions for simulating liquid water. *J. Chem. Phys.* 79:926.

44. Word, J. M., S. C. Lovell, …, D. C. Richardson. 1999. Asparagine and glutamine: using hydrogen atom contacts in the choice of side-chain amide orientation. *J. Mol. Biol.* 285:1735–1747.

45. Brünger, A. T., and M. Karplus. 1988. Polar hydrogen positions in proteins: empirical energy placement and neutron diffraction comparison. *Proteins.* 4:148–156.

46. Dunbrack, R. L., Jr., and M. Karplus. 1993. Backbone-dependent rotamer library for proteins. Application to side-chain prediction. *J. Mol. Biol.* 230:543–574.

47. Im, W., M. S. Lee, and C. L. Brooks, 3rd. 2003. Generalized Born model with a simple smoothing function. *J. Comput. Chem.* 24:1691–1702.

48. Lippow, S. M., K. D. Wittrup, and B. Tidor. 2007. Computational design of antibody-affinity improvement beyond in vivo maturation. *Nat. Biotechnol.* 25:1171–1176.

49. Mendes, J., A. M. Baptista, …, C. M. Soares. 1999. Improved modeling of side-chains in proteins with rotamer-based methods: a flexible rotamer model. *Proteins.* 37:530–543.

50. Wall, M. A., B. A. Posner, and S. R. Sprang. 1998. Structural basis of activity and subunit recognition in G protein heterotrimers. *Structure.* 6:1169–1183.

51. Neves, S. R., P. T. Ram, and R. Iyengar. 2002. G protein pathways. *Science.* 296:1636–1639.

52. Fawzi, A. B., D. S. Fay, …, J. K. Northup. 1991. Rhodopsin and the retinal G-protein distinguish among G-protein $\beta$ $\gamma$ subunit forms. *J. Biol. Chem.* 266:12194–12200.

53. Schmidt, C. J., T. C. Thomas, …, E. J. Neer. 1992. Specificity of G protein $\beta$ and $\gamma$ subunit interactions. *J. Biol. Chem.* 267:13807–13810.

54. Rens-Domiano, S., and H. E. Hamm. 1995. Structural and functional relationships of heterotrimeric G-proteins. *FASEB J.* 9:1059–1066.

55. Yan, K., V. Kalyanaraman, and N. Gautam. 1996. Differential ability to form the G protein $\beta$ $\gamma$ complex among members of the $\beta$ and $\gamma$ subunit families. *J. Biol. Chem.* 271:7141–7146.

56. Dunbrack, R. L., Jr. 2002. Rotamer libraries in the 21st century. *Curr. Opin. Struct. Biol.* 12:431–440.

57. Petrella, R. J., and M. Karplus. 2001. The energetics of off-rotamer protein side-chain conformations. *J. Mol. Biol.* 312:1161–1175.

58. Murzin, A. G. 1992. Structural principles for the propeller assembly of $\beta$-sheets: the preference for seven-fold symmetry. *Proteins.* 14:191–201.

59. Conklin, B. R., and H. R. Bourne. 1993. Structural elements of G $\alpha$ subunits that interact with G $\beta$ $\gamma$, receptors, and effectors. *Cell.* 73:631–641.

60. Neer, E. J. 1995. Heterotrimeric G proteins: organizers of transmembrane signals. *Cell.* 80:249–257.

61. Wu, X.-H., Y. Wang, …, Y.-D. Wu. 2012. Identifying the hotspots on the top faces of WD40-repeat proteins from their primary sequences by $\beta$-bulges and DHSW tetrads. *PLoS One.* 7:e43005.

62. Hendsch, Z. S., and B. Tidor. 1994. Do salt bridges stabilize proteins? A continuum electrostatic analysis. *Protein Sci.* 3:211–226.

63. Archontis, G., T. Simonson, and M. Karplus. 2001. Binding free energies and free energy components from molecular dynamics and Poisson-Boltzmann calculations. Application to amino acid recognition by aspartyl-tRNA synthetase. *J. Mol. Biol.* 306:307–327.

64. Elcock, A. H. 2001. Prediction of functionally important residues based solely on the computed energetics of protein structure. *J. Mol. Biol.* 312:885–896.

65. Hendsch, Z. S., M. J. Nohaile, …, B. Tidor. 2001. Preferential heterodimer formation via undercompensated electrostatic interactions. *J. Am. Chem. Soc.* 123:1264–1265.

66. Green, D. F., and B. Tidor. 2005. Design of improved protein inhibitors of HIV-1 cell entry: optimization of electrostatic interactions at the binding interface. *Proteins.* 60:644–657.

67. Dayhoff, M. O., R. M. Schwartz, and B. C. Orcutt. 1978. Chapter 22: A model of evolutionary change in proteins. *In* Atlas of Protein Sequence and Structure. National Biomedical Research Foundation, Washington, DC, pp. 345–352.

68. Henikoff, S., and J. G. Henikoff. 1992. Amino acid substitution matrices from protein blocks. *Proc. Natl. Acad. Sci. USA.* 89:10915–10919.

69. Altschul, S. F. 1991. Amino acid substitution matrices from an information theoretic perspective. *J. Mol. Biol.* 219:555–565.

70. Sun, D., T. Flock, …, D. B. Veprintsev. 2015. Probing G$\alpha_i$1 protein activation at single-amino acid resolution. *Nat. Struct. Mol. Biol.* 22:149–170.

71. Forbes, S. A., D. Beare, …, P. J. Campbell. 2015. COSMIC: exploring the world's knowledge of somatic mutations in human cancer. *Nucleic Acids Res.* 43:D805–D811.

72. Cerami, E., J. Gao, …, N. Schultz. 2012. The cBio cancer genomics portal: an open platform for exploring multidimensional cancer genomics data. *Cancer Discov.* 2:401–404.

73. Gao, J., B. A. Aksoy, …, N. Schultz. 2013. Integrative analysis of complex cancer genomics and clinical profiles using the cBioPortal complementary data sources and analysis options. *Sci. Signal.* 6:1–20.

74. Yoda, A., G. Adelmant, …, A. A. Lane. 2015. Mutations in G protein $\beta$ subunits promote transformation and kinase inhibitor resistance. *Nat. Med.* 21:71–75.