



HAL
open science

Extracting phylogenetic dimensions of coevolution reveals hidden functional signals

Alexandre Colavin, Esha Atolia, Anne-Florence Bitbol, Kerwyn Casey Huang

► **To cite this version:**

Alexandre Colavin, Esha Atolia, Anne-Florence Bitbol, Kerwyn Casey Huang. Extracting phylogenetic dimensions of coevolution reveals hidden functional signals. *Scientific Reports*, 2022, 12 (1), pp.820. 10.1038/s41598-021-04260-1 . hal-03534497

HAL Id: hal-03534497

<https://hal.sorbonne-universite.fr/hal-03534497>

Submitted on 19 Jan 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



OPEN

Extracting phylogenetic dimensions of coevolution reveals hidden functional signals

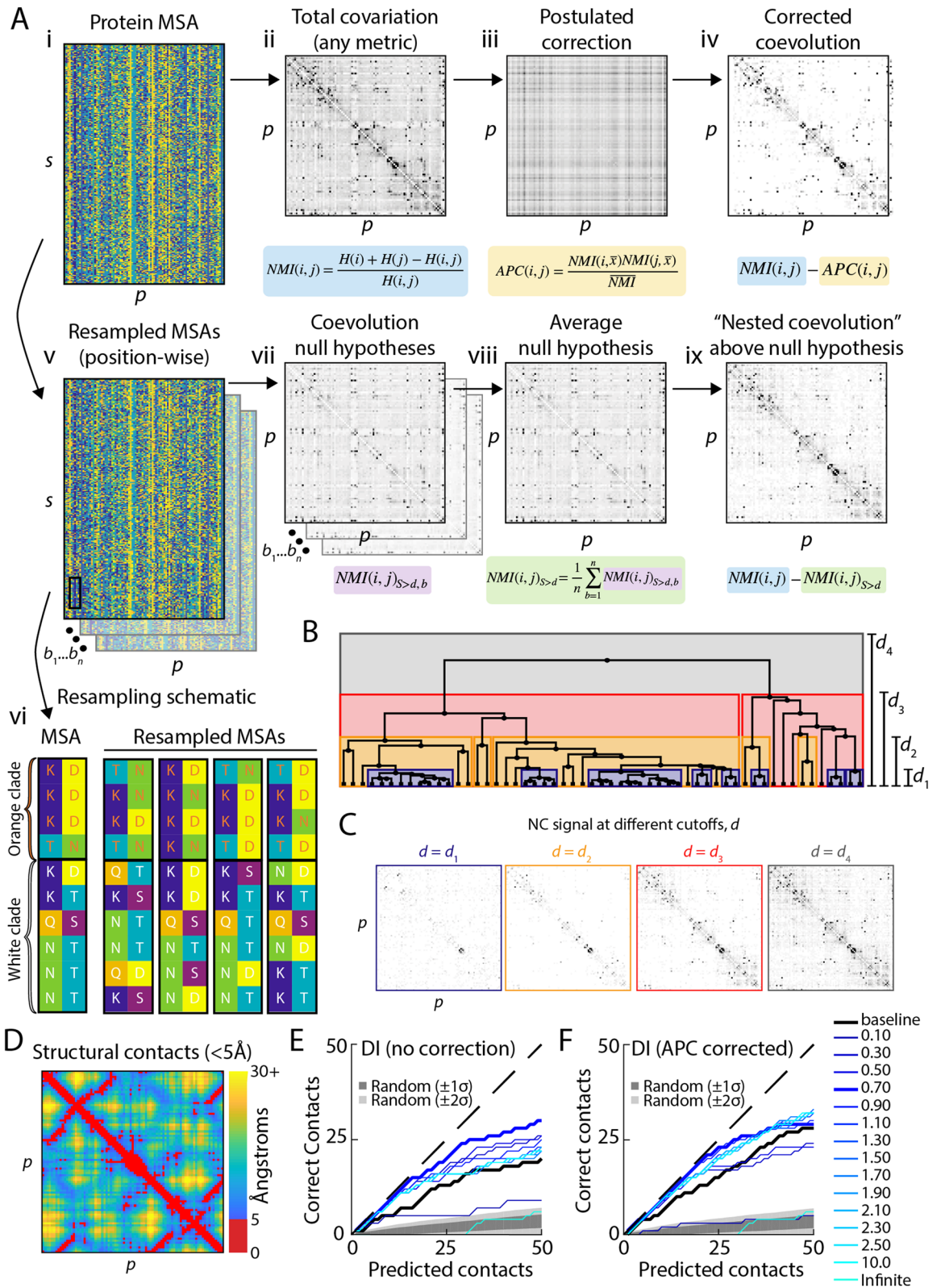
Alexandre Colavin^{1,9}, Esha Atolia^{2,9}, Anne-Florence Bitbol^{3,4,5} & Kerwyn Casey Huang^{6,7,8}✉

Despite the structural and functional information contained in the statistical coupling between pairs of residues in a protein, coevolution associated with function is often obscured by artifactual signals such as genetic drift, which shapes a protein's phylogenetic history and gives rise to concurrent variation between protein sequences that is not driven by selection for function. Here, we introduce a background model for phylogenetic contributions of statistical coupling that separates the coevolution signal due to inter-clade and intra-clade sequence comparisons and demonstrate that coevolution can be measured on multiple phylogenetic timescales within a single protein. Our method, nested coevolution (NC), can be applied as an extension to any coevolution metric. We use NC to demonstrate that poorly conserved residues can nonetheless have important roles in protein function. Moreover, NC improved the structural-contact predictions of several coevolution-based methods, particularly in subsampled alignments with fewer sequences. NC also lowered the noise in detecting functional sectors of collectively coevolving residues. Sectors of coevolving residues identified after application of NC were more spatially compact and phylogenetically distinct from the rest of the protein, and strongly enriched for mutations that disrupt protein activity. Thus, our conceptualization of the phylogenetic separation of coevolution provides the potential to further elucidate relationships among protein evolution, function, and genetic diseases.

It has long been appreciated that comparisons among homologous sequences of a protein of interest can provide key information about its function and structure. Just as evolutionarily conserved individual residues are generally crucial to a protein's proper function, the statistical covariation (arising from correlated evolution, i.e. coevolution) between pairs of residues^{1,2} carries information that is useful for predicting structural contacts^{3–7} and protein–protein interactions^{8–11} and their interfaces¹², intuiting novel protein conformations⁵, understanding protein allostery¹³, interpreting variants^{14,15}, identifying functional domains^{16–19}, and reprogramming protein specificity²⁰. However, despite the increasing prevalence of sequencing data, sampling of the phylogenetic tree is fundamentally limited and biased. Evolutionary events such as speciation can drive simultaneous changes that are statistically linked but may not reflect relevant functional coupling, for example when they arise from genetic drift. Hence, spurious covariation is more likely to arise in comparisons between distantly related sequences, hindering the ability of such studies to deliver functional insights.

Numerous methods exist for measuring protein coevolution. Statistical coupling analysis (SCA), which normalizes the covariance matrix by a function of the entropy, provides sufficient information to specify a protein fold²¹ and to detect functional domains^{6,19}. Mutual information (MI) with various corrections enables identification of some directly interacting residue pairs in the three dimensional protein structure²², and direct coupling analysis (DCA) has improved over MI by attempting to deconvolve higher-order correlations^{4,23}. All of these methods implement corrections for reducing the effects of phylogenetic noise. Although MI is extremely sensitive to the phylogenetic distribution of sequences and the conservation (measured via entropy) of individual positions, normalization by the joint entropy reduces the influence of phylogeny and entropy and improves structural-contact prediction²². DCA usually involves downweighting the coevolutionary signal contributions

¹Biophysics Program, Stanford University School of Medicine, Stanford, CA 94305, USA. ²Department of Chemical and Systems Biology, Stanford University School of Medicine, Stanford, CA 94305, USA. ³Laboratoire Jean Perrin (UMR 8237), Institut de Biologie Paris-Seine, CNRS, Sorbonne Université, 75005 Paris, France. ⁴Institute of Bioengineering, School of Life Sciences, Ecole Polytechnique Fédérale de Lausanne (EPFL), 1015 Lausanne, Switzerland. ⁵SIB Swiss Institute of Bioinformatics, 1015 Lausanne, Switzerland. ⁶Department of Bioengineering, Stanford University, Stanford, CA 94305, USA. ⁷Department of Microbiology and Immunology, Stanford University School of Medicine, Stanford, CA 94305, USA. ⁸Chan Zuckerberg Biohub, San Francisco, CA 94158, USA. ⁹These authors contributed equally: Alexandre Colavin and Esha Atolia. ✉email: kchuang@stanford.edu



◀ **Figure 1.** Nested coevolution (NC) introduces a phylogenetic dimension to existing coevolution metrics that removes noise and improves structural prediction. **(A)** Schematic illustrating the NC correction to traditional coevolution algorithms. The MSA **(i)** is used to generate a covariation matrix **(ii)** with a particular metric such as normalized mutual information (NMI) or direct information (DI). Previous studies have attempted to remove phylogenetic noise using the average product correction (APC) **(iii)**, which results in a corrected coevolution matrix **(iv)** that has lower levels of off-diagonal signal. For the NC correction, the MSA is resampled multiple times **(v)** within clades defined by a phylogenetic cutoff d **(vi)**, providing null hypotheses **(vii)** that are averaged **(viii)** to correct the covariation matrix **(i)**. The resulting difference **(ix)** is the NC matrix for a particular cutoff d . **(B)** The Jukes–Cantor phylogenetic distance between homologs defines clades (visualized as a tree) within the NC cutoff d . **(C)** NC signal at different cutoffs d as illustrated in **(B)** for the MSA of the KH domain from **(B)**. For small values of d , the NC matrix exhibits very little off-diagonal signal. **(D)** The structural contact map for KH, highlighting residue–residue contacts that are in close 3D proximity ($< 5 \text{ \AA}$, red), respectively. **(E,F)** NC with particular cutoffs d improves the prediction of structural contacts relative to DCA, applied to DI without **(E)** or after correction with APC **(F)** (black lines). All residues within five positions on the polypeptide sequence were excluded from the analysis. Black represents the predictions of the baseline NMI metric, and the shaded area represents the number of correct predictions expected by chance.

from over-represented sequences to reduce their bias^{4,23}. Motivated by the observed strong relationship between a position's average MI and the MI it exhibits with specific positions, the widely employed average product correction (APC) subtracts this average signal to address global phylogenetic patterns; this correction can be applied to any existing coevolution metric other than MI. While empirically enhancing the resolution of functional significance, none of these pragmatic strategies attempt to resolve how coevolution signal changes across the multiple potential phylogenetic timescales.

The fact that phylogenetic correlations can obscure functional coevolution signal has been frequently discussed and has motivated several recent studies. The Evolutionary Trace algorithm identifies residues whose variations during evolution correlate with major phylogenetic divergences^{24,25}, and has recently been applied to covariation²⁶ and functional prediction of oncogenes²⁷. A traditional means to address this issue in DCA contact prediction is to reweight sequences closer than a certain Hamming distance threshold^{4,23}. However, this correction and variants thereof were recently found to only yield small improvements, by contrast with the widely used APC (with which reweighting is generally combined) and with a similar entropy-based correction²⁸. The noise from phylogeny has been found to be smaller than that from entropy (i.e. conservation), but still of the same order of magnitude, hence reducing its influence on functional coevolution measurements remains an important goal²⁹. Mathematically, phylogenetic correlations strongly impact the modes of the covariance matrix with the largest eigenvalues, and suppressing these modes improves contact prediction relative to simple covariance³⁰. However, corrections (other than APC) designed to remove phylogenetic noise have not been found to substantially improve contact prediction by DCA thus far^{31,32}. Nonetheless, reweighting protein subfamilies can improve DCA results in cases where subfamilies feature structural differences, for example in homodimerization patterns^{33,34}.

Even with affordable sequencing and widespread environmental sampling, coevolution methods are often limited by the number of naturally occurring protein sequences available. Successful predictions of structural contacts using coevolution-based methods often require several thousand sequences to align^{3,35}, which is generally prohibitive for many mammalian proteins. For other proteins, the phylogenetic distribution of available sequences is skewed by sampling and is well recognized as a source of spurious signal in coevolution^{22,36}. Thus, methods that enable the separation of functional coupling from phylogenetic and sampling noise would greatly expand the utility of coevolution, particularly for applications to diseases involving human proteins with limited numbers of available sequences.

Here, we introduce the concept of nested coevolution (NC), a correction that leverages a well-defined null hypothesis to quantify the coevolutionary signal above what is expected from phylogenetic distribution alone. We determined that NC results in improved structural-contact prediction for several coevolution metrics across many proteins, especially those with fewer sequences. In addition, we found that NC improves the detection of spatially compact groups of collectively coevolving residues (“sectors”) that are phylogenetically distinct from each other and the protein itself, beyond differences in entropy alone. Finally, sectors identified using NC were enriched for positions at which mutations are maximally deleterious, suggesting that our method enhances the functional significance of coevolution signal. Since our method is agnostic to the underlying method of measuring coevolution, we anticipate wide utility for the ability to resolve the temporal dimension of protein coevolution.

Results

Background model of coevolution reveals temporal dimension of coevolution. To interrogate the contribution of phylogenetic sampling to protein coevolution measurements, we sought to separate the coevolution signal due to inter-clade and intra-clade sequence comparisons (Fig. 1A,B). Given a multiple sequence alignment (MSA) for a protein of interest (Fig. 1Ai), we first measure the total covariation (C_T) between every pair of positions (Fig. 1Aii) using an established metric of residue–residue coupling such as the normalized mutual information (NMI; Fig. 1A)²²:

$$C_T^{ij} = (H_i + H_j - H_{ij})/H_{ij}, \quad (1)$$

where H_i is the Shannon entropy (a measure of conservation) of position i , and H_{ij} is the joint Shannon entropy of positions i and j . The quantity $H_i + H_j - H_{ij}$ is the mutual information between positions i and j , which measures the coupling between residues (Fig. S1A). The NMI residue pair covariation in Eq. (1) is an attractive

choice of metric because normalizing by H_{ij} makes the MI independent of conservation²². Nonetheless, we note that our algorithm can be applied to any covariation metric, and as we will show, our main results are robust to metric choice.

The most straightforward null hypothesis for protein coevolution is that coevolutionary coupling between pairs of positions in a protein is completely absent—that is, that the probability of a position having any amino acid identity is independent of any other position's identity. Although this null hypothesis can be evaluated analytically for some methods, other methods have no known closed-form solution for the expected value of the coevolution matrix under these conditions. Hence, we computationally compute the average coevolution signal from many globally resampled MSAs in which each position in each sequence in the original MSA is replaced by the equivalent position from another randomly chosen sequence (resampled with replacement; Fig. 1Av,vi). We expect any measured coevolution from these resampled matrices to represent signal due simply to the distribution of amino acid identities at each position; any significant difference between the coevolution signal measured in the original MSA and this null hypothesis can potentially be attributed to coevolution.

This initial null hypothesis does not test for the phylogenetic structure of sequences; in the globally resampled MSAs, every sequence is effectively evolutionarily equidistant from one another. Previous attempts to remove the influence of phylogeny such as APC (Fig. 1Aiii), which corrects the covariation matrix by subtracting the product of its mean value across columns and rows for each pair of positions (Fig. 1Aiv), have substantially improved contact prediction²². However, the APC is a postulated correction that does not directly account for the phylogenetic structure of an MSA. We sought to construct a null hypothesis-driven background model of the expected coevolution in an MSA in which intra-clade coevolution is explicitly removed. We achieve this goal by incorporating the phylogenetic structure of the MSA into the null hypothesis, via generating MSAs by resampling each position from sequences that are closely related (Fig. 1Av,vi), thus removing correlations arising from recent evolutionary history within each clade. We define a clade as the subset of sequences S with a Jukes–Cantor distance below d , which we refer to as the phylogenetic cutoff. For each value of d , we calculate the inter-clade covariation $(C_{S>d}^{ij})$ from a resampled MSA either analytically or via bootstrapping (Fig. 1Avii, Fig. S2A, Methods), where C denotes the chosen covariation measure (e.g. NMI). This inter-clade covariation thus measures the expected value of covariation due solely to the comparison of sequences between clades (Fig. 1B). We then average over many such null hypotheses (over many within-clade resampled MSAs at fixed d), yielding the mean inter-clade covariation matrix $(C_{S>d})$ (Fig. 1Aviii), which represents the expected coevolution due to both the distribution of amino acid identities at each position and the phylogenetic structure of the protein MSA (Fig. 1B). Significant differences between this background model and the baseline signal measured from the original MSA represent signal that was contained in the intra-clade comparison of closely related sequences. Since the difference between the background null model $C_{S>d}$ and the baseline signal C_T qualitatively captures the significance of the baseline measurement (Fig. S2B, Methods), we subtract $C_{S>d}$ from the total covariation C_T to obtain the phylogenetic cutoff-dependent covariation signal $C_{S\leq d}$ (Fig. 1Aix):

$$C_{S\leq d}^{ij} \equiv C_T^{ij} - C_{S>d}^{ij}, \quad (2)$$

where positive values indicate that the total covariation is larger than expected by comparison of sequences between clades, thus revealing covariation arising from less divergent sequences in all clades. We refer to the signal $C_{S\leq d}^{ij}$ above the null hypothesis $C_{S>d}^{ij}$ in Eq. 2 as a protein's “nested coevolution” (NC), in that it separates coevolution into signal attributed to comparison of sequences within $(C_{S\leq d}^{ij})$ as compared with between $(C_{S>d}^{ij})$ nested clades of a phylogenetic tree. The only free parameter in the NC is the phylogenetic cutoff (d). As we vary the cutoff, typically many patterns of NC emerge, revealing distinct windows of coevolution for a single protein MSA (Fig. 1C). The changes in NC observed between two cutoffs represent the signal attributable to the comparison of pairs of sequences whose distance is between the cutoffs used to calculate each window; it is not the case that the signal for d_i is necessarily a subset of the signal for d_{i+1} . Hence, distinct signals of protein coevolution are revealed as the phylogenetic cutoff is varied.

To test the relevance of NC windows to protein structure prediction, we measured the enrichment of structural contacts from the pairs of residues with the highest 50 values in the NC matrix $C_{S\leq d}^{ij}$ for each value of d . Here, we applied NC as a correction to DCA, a current standard for coevolution-based prediction of structural contacts^{4,23}. We employed the direct information (DI) metric^{4,23} to quantify coevolution for the KH domain, which is present in a wide variety of nucleic acid-binding proteins³⁷. In this and subsequent analyses, we considered structural contacts to be within 5 Å at closest approach, excluding pairs of residues within 5 amino acids on the sequence (Methods); qualitatively similar results were obtained when defining contacts based on an 8 Å cutoff (Fig. S1D). The NC phylogenetic cutoffs revealed a variety of improvements (Fig. 1D). Some windows generally outperformed DCA, without (Fig. 1E) or with (Fig. 1F) the APC.

To determine the added value of NC for other proteins and for another coevolution metric (the Frobenius norm³⁸, which is frequently utilized in DCA as an alternative to DI^{39,40}), we carried out a DCA structural-contact analysis for 10 protein family domains with DI or Frobenius norm (Methods). Across both metrics and all proteins, NC improved the predictions of structural contacts (Fig. 2A), even relative to the inclusion of APC²². Hence, NC is a correction that generally enhances the predictive power of widespread coevolution measurements.

NC can improve predictions of structural contacts using fewer sequences. One common limitation for computing coevolution is the number of homologous sequences available for constructing an MSA. To interrogate whether NC could still accurately predict structural contacts with fewer sequences, we subsampled the MSAs of 10 proteins with different breadth (randomly selecting 10% or 1% of the sequences) or depth (selecting the 10% or 1% of sequences most related to the protein used to construct the MSA to mimic the lim-

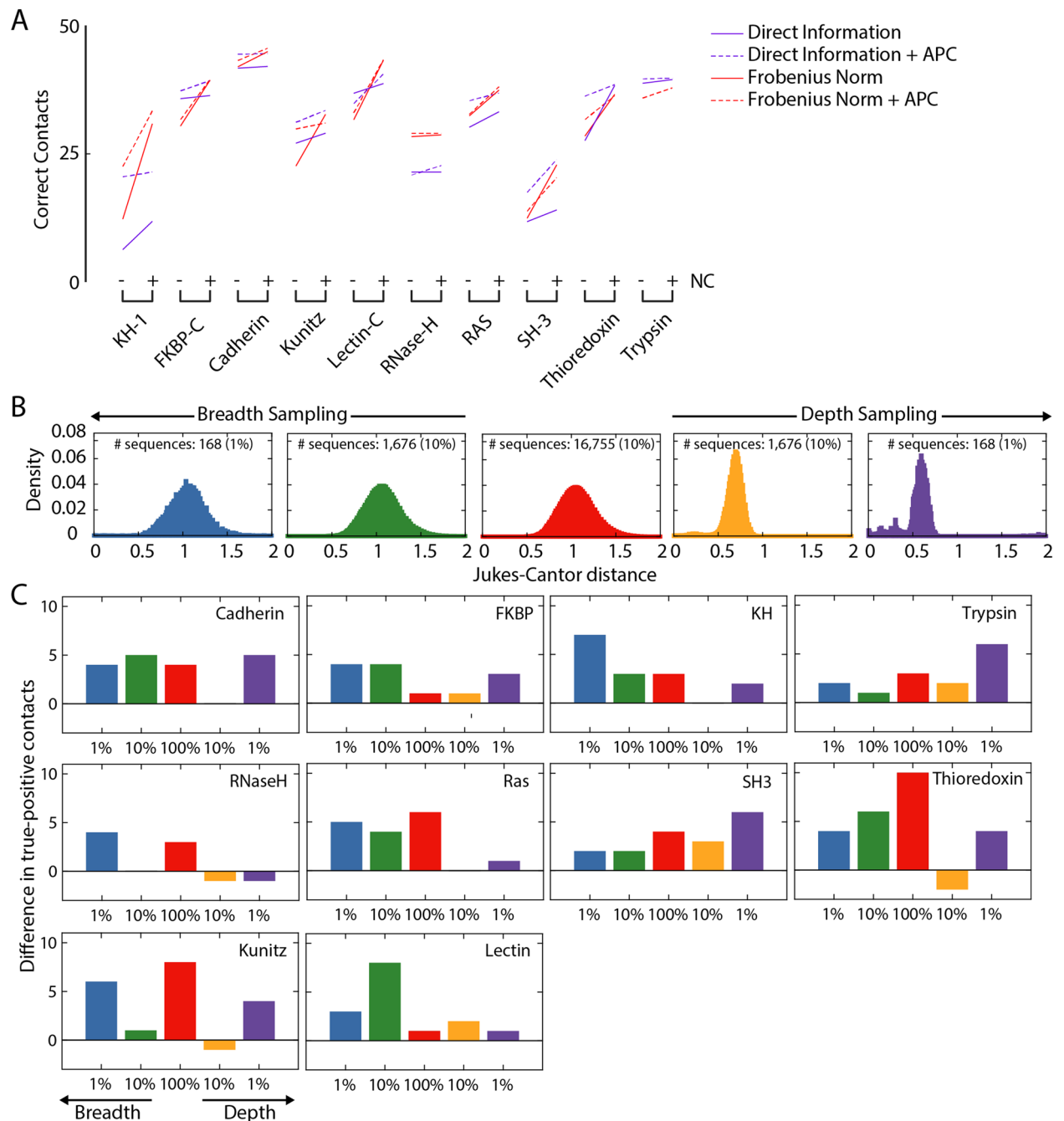


Figure 2. NC can improve predictions of structural contacts across proteins and coevolution methods, and resolve information loss due to subsampling of sequences. **(A)** NC increased the number of true-positive structural contacts among the first 50 predictions for 10 highly conserved proteins predicted by DCA using DI or Frobenius norm, without or with APC. We chose the maximum number of recovered true-positive structural contacts from NC applied across four phylogenetic cutoffs; the set of cutoffs was consistent across all proteins (0.2, 0.5, 0.8 and 1.1). **(B)** MSAs were subsampled across the breadth (random sampling) and depth (sorted sampling) of the MSA. Typically, the distribution of Jukes–Cantor distances in the MSA (red) remained essentially unchanged for breadth sampling (green and blue), while it shifted to lower values (as expected) for depth sampling (gold and purple); shown is the KH domain. **(C)** NC generally increased the number of true-positive structural contacts among the first 50 predictions relative to DCA employing DI (without APC) across proteins and both breadth and depth sampling (for DI with APC, see Fig. S3). Small decreases occurred for depth sampling of RNase H, thioredoxin, and Kunitz.

ited phylogenetic distribution of small protein families, Table S1) (Fig. 2B). NC improved structural contact prediction for a majority of the subsampled MSAs when correcting DI without (Fig. 2C, Fig. S3A) or with (Fig. S3B) application of APC. For the KH domain, more than twice as many true positives were predicted after applying NC compared with DI + APC alone (Fig. S3A). Perhaps unsurprisingly, breadth sampling generally performed better than depth sampling (Fig. 2C), indicating that accurate prediction is reliant on the sequences being sufficiently distantly related. Nonetheless, for many proteins, the value of the NC correction was enhanced when the number of homologous sequences was low, both for depth and breadth samplings.

NC eigenvectors exhibit reduced background noise, improving detection of spatially compact sets of coevolving residues.

Previous studies have utilized coevolution measurements to identify groups of residues within a protein that are spatially compact on the tertiary structure and thus are postulated to have a joint function^{6,19,41–44}. These “sectors” can be defined by a variety of methods, such as the extreme-value residues of the eigenvectors of the coevolution matrix with the largest eigenvalues¹⁹, and in the context of SCA have been proposed to reflect independent biological properties such as catalytic efficiency and thermal stability¹⁹. Motivated by these successes, we sought to measure the effect of incorporating the phylogenetic dimension revealed by NC when defining sectors of residues. Specifically, we measured the NC- and APC-corrected coevolution across a range of phylogenetic cutoffs, concatenated the resulting matrices, and performed eigendecomposition to identify the most significant eigenvectors (Methods). The residues most strongly associated with the positive or negative components of each resulting eigenvector are considered a sector. Here, we used NMI as a baseline coevolution metric, although NC could also be applied to SCA results.

We first focused on MreB, an essential protein involved in cell-shape determination in many rod-shaped bacteria⁴⁵. MreB belongs to a protein family that includes ParM, FtsA, and MamK in bacteria, crenactin in archaea, and actin in eukaryotes^{46,47}. These proteins are structural homologs characterized by a four-subdomain fold around an ATP-binding pocket^{47,48}, with a wide range of sequence identities and disparate cellular functions. Thus, we anticipated that the set of MreB homologs would have sufficient diversity to support robust coevolution measurements, particularly functional sectors.

We compared NC-derived sectors with sectors derived from eigenvectors of the NMI coevolution matrix for MreB homologs; hereafter, we define these NMI-derived sectors as the “baseline.” We identified the most closely related baseline sectors for three NC eigenvectors with some of the highest eigenvalues, which we refer to as eigenvectors A, B, and C (Methods). Each pair of NC and baseline eigenvectors appeared similar, especially for the residues with the largest absolute coefficients (Fig. 3A–C). However, the baseline eigenvectors exhibited much higher variation of coefficients across the protein (Fig. 3A–C). For eigenvectors A and B, the NC-derived eigenvectors exhibited 32.8-fold and 38.3-fold lower standard deviation (after removing the 50 highest and lowest coefficients) than the corresponding baseline-derived eigenvectors, respectively (Fig. 3A,B). For eigenvector C, the baseline eigenvector contained residues with both highly positive and highly negative coefficients, while the high-magnitude coefficients of the NC eigenvector were solely positive (Fig. 3C); the positive portion of the NC eigenvector again had substantially lower (2.1-fold) noise than the baseline eigenvector (Fig. 3C).

Motivated by eigenvector C, we defined distinct positive and negative sectors (Methods) for each NC and baseline eigenvector using a variable cutoff on the site contributions to adjust sector size (as sectors are defined as the sets of amino acids with highest site contributions in each eigenvector). For each sector size, we quantified the spatial compactness as the mean pairwise distance between alpha carbons of residues within a sector. For sectors A–C (derived from eigenvectors A–C), the first 5–9 residues exhibited approximately the same spatial compactness in the NC sectors as in the baseline sectors (Fig. 3D–F). However, as the cutoff was increased, the NC sector remained more spatially compact than the baseline sector (Fig. 3D–F). All three NC sectors were also more spatially compact than expected based on random sampling for cutoffs yielding up to at least 50 residues (Fig. 3D–F), while baseline sector A was distributed across the protein structure (Fig. 3D,J). NC sectors A and C were largely situated in subdomains IIA (Fig. 3G) and IA (Fig. 3I), respectively, while sector B was localized to the ATP-binding pocket (Fig. 3H). Notably, sector C was spatially compact (Fig. 3I) despite being spread across the protein sequence (Fig. 3C). Baseline sectors B and C with 15 residues were qualitatively similar to the corresponding NC sectors (Fig. 3K,L); the large background fluctuations of the baseline eigenvector likely led to the inclusion of additional, erroneous residues into the sector prediction. Thus, the phylogenetic correction of NC improves the spatial compactness of sectors.

Sectors display distinct phylogenetic signatures from the rest of the protein.

Since sectors have been postulated to reflect distinct evolutionary histories driven by selection for particular biological functions¹⁹, we sought to compare the phylogeny of the residues within a sector with other sectors and the rest of the protein. The MirrorTree algorithm (Methods) was originally developed to compare phylogenies of two proteins, motivated by the assumption that similar histories signify a common function, e.g. through protein–protein interactions and/or acting in the same pathway^{49,50}. After computing a pairwise distance matrix of all sequences within an MSA for each of the two proteins based on homologs in the same set of organisms, the MirrorTree score is defined as the Pearson correlation coefficient between the entries in the two pairwise distance matrices⁴⁹. We straightforwardly modified the MirrorTree method to compare the complete protein MSA to the MSA filtered to include only the residues within the sector of interest (Fig. 4A).

To broadly investigate sector identification, we identified 40 15-residue sectors for MreB based on the positive and negative coefficients of the 20 eigenvectors with the highest eigenvalues. As negative controls, we randomly sampled sets of residues of the same size as each sector from across the protein. Sector-protein MirrorTree scores for sectors A–C (Fig. 3) were substantially lower for sectors than for the random groups (Fig. 4B), which as expected all had MirrorTree scores close to 1 (Fig. 4B), meaning that the random sectors were phylogenetically

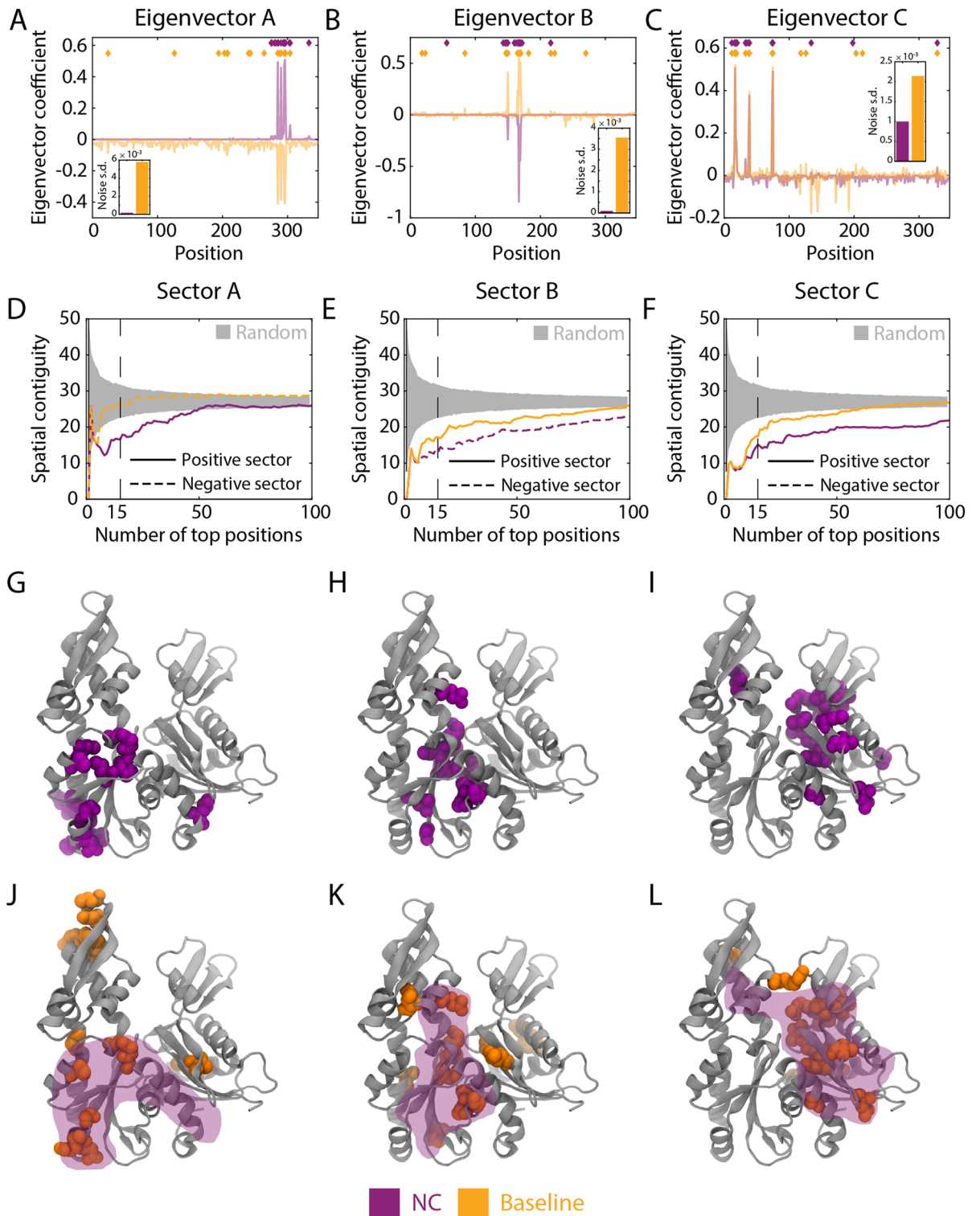


Figure 3. NC eigenvectors for the actin homolog MreB have lower noise and are more spatially contiguous than baseline eigenvectors. (A–C) Three eigenvectors with large eigenvalues were identified and paired between baseline coevolution (NMI with APC) and the NC correction for an MSA containing 9998 sequences of MreB. Aside from the residues with large coefficients, the NC eigenvectors exhibited lower signal variation than the baseline eigenvectors. Insets: standard deviations of the eigenvector coefficients after excluding the highest and lowest 50 values. (D–F) NC sectors are more spatially contiguous than the corresponding baseline sectors. Sectors were defined based on a sliding cutoff of the most positive or most negative coefficients of each eigenvector in (A–C). Spatial compactness was defined as the mean pairwise distance between each residue within a sector. Gray regions represent 95% confidence intervals of a randomly selected group of residues of the same size. (G–L) For the 15-residue versions of the NC and baseline sectors [vertical lines in (D–F)], the NC sectors (G–I) are more compact on the three-dimensional structure than the corresponding baseline sectors (J–L). The shaded purple regions in (J–L) represent the NC sector.

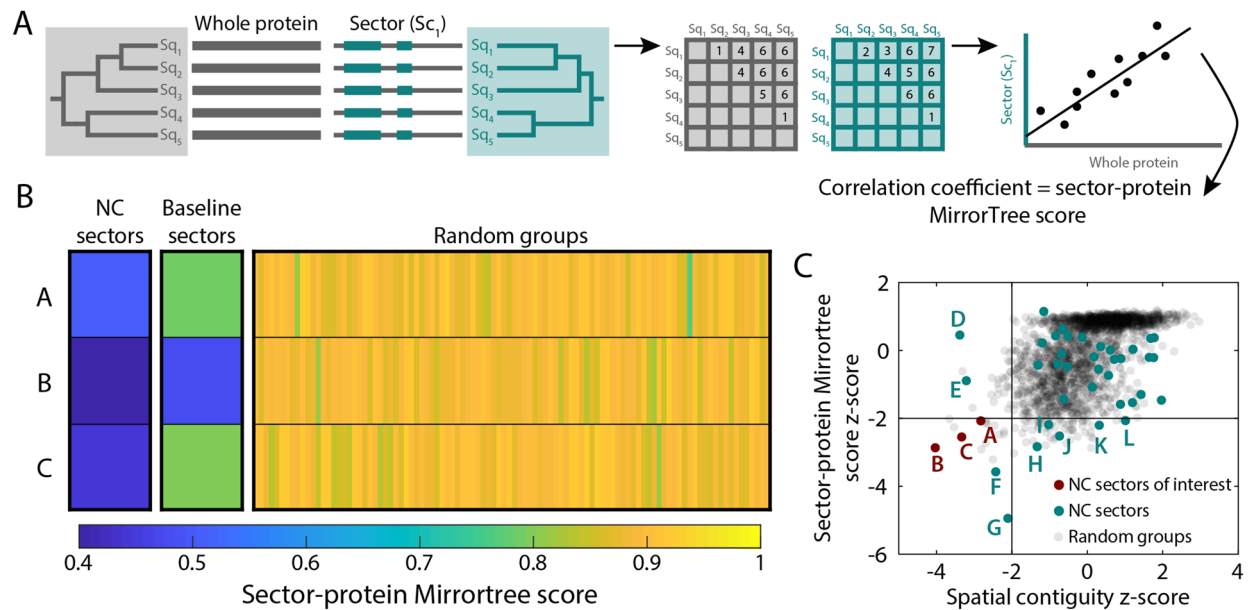


Figure 4. Sectors are phylogenetically distinct from the entire protein. **(A)** Schematic of how the MirrorTree algorithm⁴⁹ was repurposed to measure the phylogenetic similarity between sectors and the entire protein. The MirrorTree score is defined as the Pearson correlation coefficient between the entries in the two pairwise distance matrices of all sequences within an MSA for the protein versus only the residues in the sector. **(B)** MreB NC sectors A–C (Fig. 3) had lower sector-protein MirrorTree scores than the corresponding baseline sectors, while random groups of 15 residues had MirrorTree scores close to 1 (as expected). **(C)** MreB NC sectors were computed from the 15 most positive or negative coefficients of the 20 eigenvectors with the highest eigenvalues. Among these 40 sectors, the z-scores of the MirrorTree score and the spatial compactness were < -2 for sectors A–C. Sectors D–L substantially overlapped sectors A–C, and are considered in Fig. 6.

indistinguishable from the rest of the protein. Baseline sectors A–C had MirrorTree scores intermediate between those of the corresponding NC sector and random groups (Fig. 4B), likely reflecting dilution of signal due to noisy selection from baseline eigenvectors of residues that are more likely to functionally follow the phylogenetic history of the protein overall. To evaluate the significance of the MirrorTree score and of the spatial compactness of each sector, we computed z-scores based on the mean and standard deviation of the two metrics applied to the random groups of the same size as each sector. Sectors A–C had MirrorTree scores < 0.5 (Fig. 4B), indicating distinct phylogenetic histories from the protein, and MirrorTree and spatial compactness z-scores < -2 (Fig. 4C). There were four other sectors (D–G) that had spatial compactness z-scores < -2 . These sectors largely overlapped with A–C; we will return to this overlap in a later section. All other sectors had spatial compactness z-score > 2 , and all but five (H–L) had MirrorTree z-score > -2 . Thus, MirrorTree reveals that certain NC sectors have distinct evolutionary trajectories from the protein itself, motivating us to focus on certain sectors (such as A, B, and C for MreB).

Phylogenetic similarity and the role of entropy. Conservation itself is a major determinant of protein function^{51–53}, and spatially contiguous sets of residues can be identified solely on the basis of conservation⁵⁴. To account for variation in entropy across a protein, previous studies have excluded positions with high conservation (Shannon entropy < 0.1) or composed of $> 25\%$ gaps in the MSA⁵⁵. For MreB, NC sectors A–C had lower entropy than baseline sectors or random groups of the same size (Fig. 5A–C), albeit higher entropy than residues typically considered highly conserved (entropy < 0.1).

MirrorTree scores of NC sectors were also generally lower than those of baseline sectors (Fig. 5D–F). To investigate the dependence of sector-protein MirrorTree scores on entropy, we computed MirrorTree scores for thousands of random groups of the same size as the sector (15 residues), biasing sampling using a Monte Carlo algorithm to obtain a wide range of mean entropies; each random group was selected from residues that did not overlap with the sector. For mean entropy $\lesssim 1$, MirrorTree scores were strongly dependent on entropy (Fig. 5G–I). Thus, the low MirrorTree scores of the NC sectors were due in part to their low entropy. Nonetheless, the MirrorTree score of NC sector A was significantly lower than those of random groups with the same mean entropy (z-score -3.5); the entropy of sector B was so low, presumably due to the high conservation of the ATP-binding pocket (Fig. 3H, Fig. S4), that it was challenging to obtain random groups that were not largely overlapping.

Since NC sector A displayed the greatest reduction in MirrorTree score relative to random groups of the same mean entropy, we focused on this sector to investigate the dependence of sector-protein MirrorTree score on sector size. As the cutoff was increased to include more residues, the MirrorTree score increased (Fig. 5D). To disentangle whether this increase was due directly to the increase in size or to the inclusion of residues that are more phylogenetically similar to the protein, we compared the 10-residue version of sector A (Fig. 5G) with

randomly selected groups of 10 residues from 15- and 20-residue versions of sector A, as well as the entire protein. The mean MirrorTree score increased as the size of the sampling group increased (Fig. 5J), even for groups with similar entropy as the 10-residue sector (Fig. 5K). Moreover, 15-residue versions of sectors B and C had similar entropy (Fig. 5B,C); hence, an approach driven by entropy alone would not have divided these spatially separated clusters. Thus, the strength of a residue's association in a sector of highly coevolving residues is associated with more phylogenetic distinction from the rest of the protein than can be explained by entropy alone.

Phylogenetic similarity highlights overlapping sectors. The core residues of some MreB NC eigenvectors sometimes had high coefficients in multiple eigenvectors (Fig. S5), suggesting that we should consider the union of the sectors as a functional unit. To rationally identify sectors that should be merged, we again exploited phylogenetic similarity by calculating MirrorTree correlation coefficients from comparisons between pairs of sectors (Fig. 6A). MreB NC sectors A–C (Fig. 3) exhibited low sector–sector MirrorTree scores with each other and with random groups (Fig. 6B), as expected since they have low sector–protein MirrorTree scores (Fig. 6B). By contrast, the random groups had MirrorTree scores close to 1 (Fig. 6B). NC sectors were also more phylogenetically distinct from each other than baseline sectors (Fig. 6C). These data suggest that the NC sectors were selected by evolutionary pressures that led to distinct functions, which influenced their phylogeny in distinct manners.

Of all sectors that had a MirrorTree z -score or a pairwise distance z -score < -2 (sectors A–L, Fig. 4C), several pairs had a high sector–sector MirrorTree score. Hierarchical clustering of the sectors based on their sector–sector MirrorTree profiles led to the identification of five obvious “meta-sectors” from the sum of the clustered eigenvectors (Methods), which we denote α , β , γ , δ , and ϵ (α , β , and γ contain sectors A, B, and C, respectively) (Fig. 6D). The meta-sectors exhibited low sector–sector MirrorTree scores (Fig. 6E), and α , β , and γ had both low sector–protein MirrorTree scores (Fig. 6E,G) and low spatial compactness z -scores (Fig. 6G). The 15-residue version of meta-sector α was more compact than the 15-residue version of A (Fig. 6H), and it contained residues that interact with RodZ (Fig. 6I), an MreB binding partner that modulates MreB filament nucleation⁵⁶ and curvature⁵⁷. Notably, the regions of the 25-residue version of meta-sector α at the barbed and pointed ends of the MreB subunit interact with each other in a polymerized MreB filament (Fig. 6J), reinforcing the spatial compactness of the meta-sector. Meta-sector β was identical to sector B, surrounding the ATP-binding pocket (Fig. 6K). As with α and A, the 15-residue version of meta-sector γ was more compact than the 15-residue version of sector C (Fig. 6L), indicating that clustering based on MirrorTree scores increases the spatial compactness of sectors.

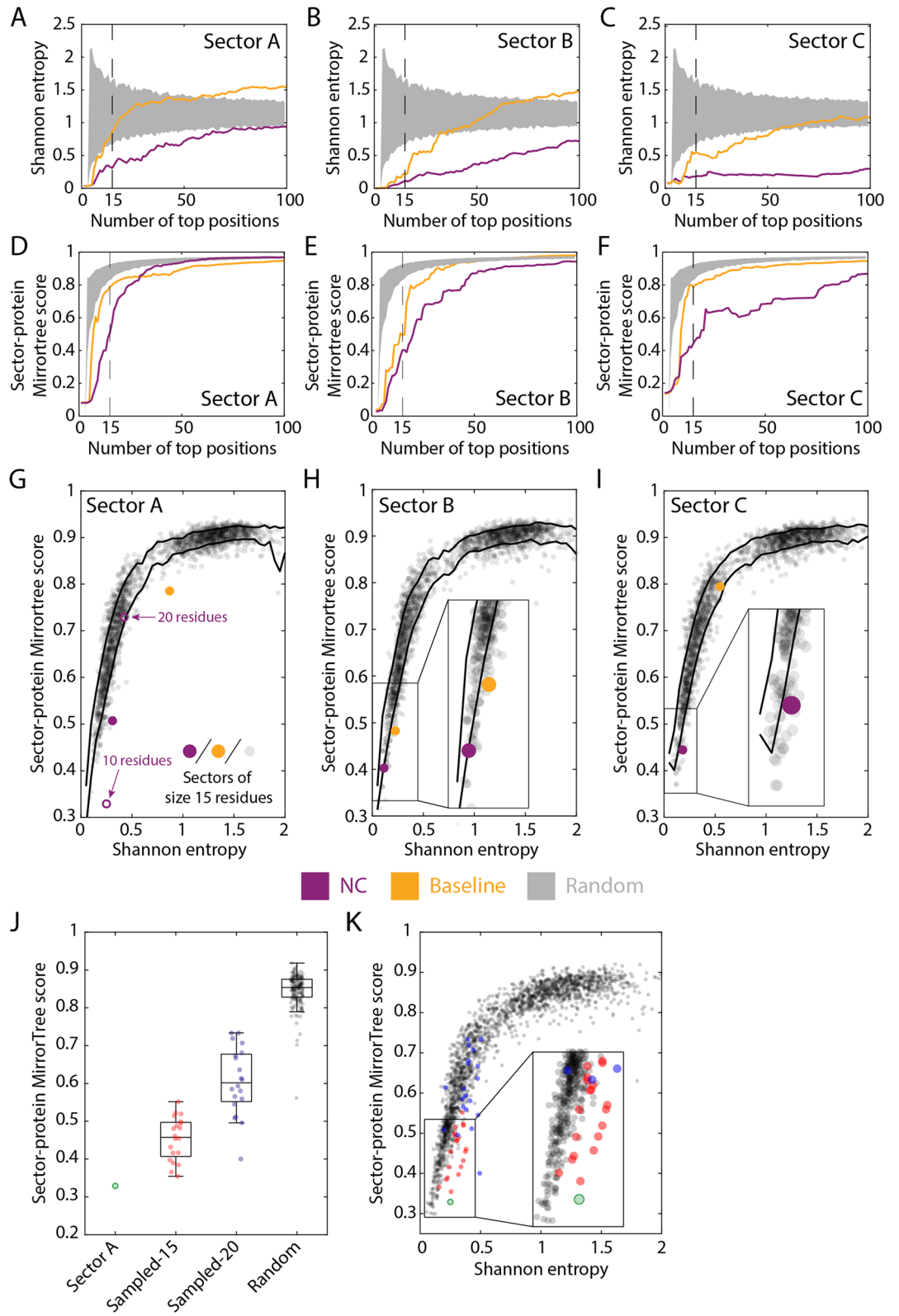
NC identifies sectors that are not apparent from the full coevolution matrix. To determine whether our findings about the properties of NC sectors applied to other proteins, we performed similar sector calculations for enolase (the metalloenzyme responsible for conversion of 2-phosphoglycerate to phosphoenolpyruvate during glycolysis⁵⁸; Fig. 7A–C), the carbohydrate-processing enzyme glucose-6-phosphate dehydrogenase (G6PD⁵⁹; Fig. 7D–F), and mitogen-activated protein kinase 1 (MAPK1)^{60,61} (Fig. 7G–K). In each case, NC produced sectors with lower background noise and higher spatial compactness than baseline sectors.

Most of the MreB NC eigenvectors had strong signal for either positive or negative coefficients, but not both (Fig. 3A–C). By contrast, one of the large-eigenvalue NC eigenvectors for MAPK1 had groups of residues with both very positive and very negative coefficients (Fig. 7G); these residues were located in distinct regions of the protein (Fig. 7J,K). As validation for splitting the NC eigenvector into two sectors, the sector–sector MirrorTree score (0.44) indicated that they are phylogenetically distinct; moreover, the sector–sector MirrorTree score of the corresponding baseline sectors was higher (0.71). Thus, NC eigenvectors can be interpreted as two phylogenetically distinct sectors based on coefficient signs.

In addition to improving sector predictions by reducing background variation, we were interested in determining whether NC can identify sectors that the full coevolution matrix misses altogether. For the arginine tRNA ligase ArgS⁶² and G6PD, the sector with the most negative MirrorTree z -score had nearly the lowest spatial compactness z -score (Fig. 7L,M) and no clear counterpart in any of the baseline eigenvectors (Methods). For ArgS, the NC sector was spatially localized around the arginine binding site (Fig. 7N). For G6PD, the NC sector was adjacent to one of the two NADPs that bind to the protein (Fig. 7O). Thus, the NC correction reveals some sectors that have apparent functional significance but are missed by the baseline method.

NC sectors are enriched in damaging mutations. To more rigorously test the functional significance of NC sectors, we sought experimental datasets with quantitative measurements of the consequences of mutations across a protein of interest. Recent studies have pioneered the use of deep mutational scanning to systematically generate and quantify the phenotypic or fitness effects of many individual mutations spanning entire domains or proteins^{41,63–65}, thereby providing new insights into structure–function relationships. Thus, we asked whether NC sectors were enriched in residues for which mutation altered protein function and/or fitness.

The Ras superfamily of membrane-associated small G-proteins is highly conserved and controls a broad range of cellular processes⁶⁶, has inactive and active states that are regulated by a GTPase-activated protein⁶⁷, and has been implicated in cancer⁶⁸. A recent deep mutational scanning study engineered plasmids to express mutant versions of human H-Ras as well as the Ras-binding domain of human C-Raf (Raf-RBD) in *Escherichia coli*⁶⁹, such that the binding of Ras-GTP to Raf-RBD led to transcription of a chloramphenicol-resistance cassette. Thus, the binding efficacy of the Ras variant was directly correlated with cellular growth rate in the presence of chloramphenicol. The effect of Ras mutations on fitness was quantified by the logarithm of the enrichment of variants in the chloramphenicol-selected versus the starting population, relative to wild-type. The distribution of fitness effects was centered around zero, although there were some positions with mutations that displayed significant functional effects⁶⁹.



◀**Figure 5.** Sector-protein MirrorTree scores of residue groups are correlated with entropy, but NC sectors have lower MirrorTree scores than expected from entropy alone. (A–C) The Shannon entropy of MreB NC sectors A–C (Fig. 3) across size cutoffs is lower than that of the corresponding baseline sectors, indicating that NC selects more conserved residues (albeit entropy is still higher than the cutoff of <0.1 for typically being considered highly conserved). Gray regions represent 95% confidence intervals of a randomly selected group of residues of the same size. (D–F) MirrorTree scores are lower for the NC sectors than for the corresponding baseline sectors. Gray regions represent the MirrorTree scores of a randomly selected group of residues of the same size. (G–I) The MirrorTree scores of sectors A–C (filled gold and purple circles) and of random groups of 15 residues (gray). Although MirrorTree score is linked to entropy, NC sectors A and C have MirrorTree scores significantly lower than expected based on entropy alone. In (G), the open purple circles denote the versions of sector A with 10 and 20 residues. Black curves indicate ± 1 standard deviation from the mean MirrorTree score for a given entropy. (J) The 10-residue version of NC sector A has lower MirrorTree score than sets of 10 residues selected from the 15- and 20-residue versions of the same sector, which are lower than those of random groups of 10 residues. The central mark indicates the median, and the bottom and top edges of the box indicate the 25th and 75th percentiles, respectively. The whiskers extend to the most extreme data points not considered outliers. (K) The 10-residue version of NC sector A has a lower MirrorTree score than 10-residue subsets of the 15- and 20-residue versions of the same sector with similar entropy. Same data as in (J). Thus, the 10-residue sector represents a “core” of the most highly coevolving residues.

To determine whether fitness-altering mutations in H-Ras are enriched at positions identified by coevolution, we identified two high-eigenvalue sectors with obvious corresponding baseline sectors. As in our previous analyses (Figs. 3A–C and 7A,D,G), aside from the highly coevolving residues, the NC sectors had much lower noise than the baseline sectors (Fig. 8A,B). The residues in the two NC sectors were non-overlapping, and in both cases appeared to be concentrated in regions with low minimum relative enrichment (Fig. 8C,D). Across cutoffs that defined sectors of various sizes, we computed the minimum and maximum relative enrichment (representing deactivation and activation, respectively) over all amino acid mutations for each position in the NC/baseline sectors as well as for the residues with the lowest entropy and compared to the distribution over all residues. As expected, the residues with lowest entropy consistently predicted significantly more negative minimum relative enrichment than random sets of residues (Fig. 8E,F). The mean minimum relative enrichment in NC and baseline versions of sector A was also significantly more negative than random residues, with the NC sector outperforming the baseline sector and achieving similar enrichment values to the lowest-entropy residues (Fig. 8E). NC sector B also exhibited mean minimum relative enrichment significantly lower than random, by contrast to the baseline sector (Fig. 8F). Thus, sectors A and B are more enriched for residues whose mutation has the most potential for reducing fitness using NC versus baseline. The maximum relative enrichment was highly similar for sectors and the protein overall (Fig. S6A,B), suggesting that NC and baseline sectors are enriched for residues with the potential for deactivating rather than activating mutations in the case of H-Ras. Thus, NC sectors can separate residues based on the maximum impact of mutations at these positions.

Discussion

Many existing coevolution methods build on correlation or mutual information, sometimes employing ad hoc corrections to partially remove the effects of entropy and phylogeny. Our NC method harnesses phylogenetic distance between sequences as a novel dimension in the measurement of protein coevolution, in order to increase understanding of the functional relationships between amino acids in a protein. While the factors that determine whether pairs of positions coevolve on short or long timescales are unknown, future studies using NC to interrogate the specific biochemical functions of protein sectors may reveal general patterns across diverse proteins. One interpretation of the variable contribution of coevolution across phylogenetic distance within a single protein (Fig. 1C) is that the frequency of mutation for coevolving residues within an NC sector is linked to the timescale of change for the corresponding selective pressure on that sector. For example, a sector that determines protein thermostability would be predicted to coevolve on a timescale commensurate with the frequency of changes in environmental temperature, whether these changes occur over long (e.g. glaciation and interglacial cycles of 100,000 years) or shorter (e.g. Atlantic multidecadal oscillations) timescales. NC can simultaneously incorporate effects of multiple phylogenetic timescales via the choice of cutoffs as well as investigate global phylogenetic patterns.

Importantly, NC and our repurposed MirrorTree methods are complementary to most covariation metrics, even phylogenomic methods such as Evolutionary Trace²⁶, and hence can enhance existing bioinformatics tools by defining a phylogenetic dimension of coevolution and improving resolution of functional signal. NC will also likely benefit from continued improvement in methods of alignment. In the future, it would be interesting to apply NC to other DCA implementations such as pseudo-likelihood optimization-based (plmDCA, GREMLIN, CCMpred) or generative algorithms (bmDCA, ACE) to examine how NC improves contact prediction with those methods^{70,71}. We anticipate that our approach will enable application of coevolution-based methods across a much broader class of proteins, including those for which the set of sequences is limited in number (Fig. 2) and/or for which the available homologous sequences are biased to a particular segment of the phylogenetic tree (Fig. 1B). In particular, application to the growing database of human exome sequences⁷² may improve identification of rare disease-causing mutations. The phylogenetic cutoff can be tuned based on the entropy and phylogenetic structure of the protein of interest to focus on different properties of the coevolution matrix, and a coevolution matrix with a specific cutoff could be used to query how sector identification changes as a function of phylogenetic depth. NC may also enhance protein engineering tools by highlighting targets for directed

Figure 6. MreB NC sectors are generally phylogenetically distinct, and those with phylogenetic overlap collectively overlap with functionally important regions. (A) Schematic of how the MirrorTree algorithm was repurposed to measure the phylogenetic similarity between sectors. (B) MreB NC sectors A, B, and C exhibited low sector–sector MirrorTree scores with each other, but high values with random groups of 15 residues (which also exhibited high MirrorTree scores with each other). (C) NC sectors A–C have lower sector–sector MirrorTree scores with each other than baseline sectors A–C with each other, indicating that they are more phylogenetically distinct. (D) Hierarchical clustering of MreB NC sectors A–L (Fig. 4C) based on sector–sector MirrorTree profiles suggests five distinct meta-sectors. (E–G) The MreB meta-sectors defined by the sum of the clustered eigenvectors exhibited low sector–sector MirrorTree scores with each other (E) as well as low sector–protein MirrorTree scores (F). Meta-sectors α , β , and γ (similar to sectors A–C) exhibited high spatial compactness (z -score < -2). (H,I) Meta-sector α was more spatially contiguous than sector A (shaded purple region) (H), and contained residues around the interface with MreB's binding partner RodZ (I). (J) The 25-residue version of meta-sector α connects the pointed and barbed ends of each subunit in a protofilament. (K) Meta-sector β (identical to sector B) surrounds the ATP binding pocket. (L) Meta-sector γ is more spatially contiguous than sector C (shaded purple region).

evolution. As we have demonstrated, NC expands our ability to detect functional relationships between residues within proteins, which could shed light in the future on the links between protein evolution and adaptation. In concert with deep mutational scanning and other comprehensive functional screens⁷³, NC and MirrorTree may be able to provide deeper insight into the specific selective pressures under which proteins have evolved.

The predominant application of coevolution so far has been structure prediction, from using top DCA-predicted contacts as constraints⁴ to employing DCA model parameters as input training features for deep neural networks that seek to predict spatial distances between amino acids⁷⁴. Here, we have shown that NC can improve contact prediction by DCA in an interpretable manner by removing correlations from pure phylogeny, complementing deep learning algorithms sometimes described as black boxes. Moreover, the detection and interpretation of sectors as functional units within proteins has been a growing research focus, particularly with respect to the evolutionary origins of sectors. A recent theoretical study demonstrated that selection acting on a functional property can give rise to a sector⁴⁰. MirrorTree scores reveal that residues within sectors have a different evolutionary history from the rest of the protein, due to both entropy-dependent and entropy-independent differences (Fig. 5). MirrorTree scores can further be used to evaluate NC predictions in the absence of a known structure. Motivated by the original design purpose of MirrorTree, we note that scores between sectors of two proteins could be used to identify protein–protein interactions—potentially between hosts and microbes—enhanced by the improved performance of NC when the sampling of sequences is shallow (Fig. 2C).

Our observation that residues most strongly associated with sectors exhibit higher spatial compactness and lower MirrorTree scores when NC is applied (Fig. 5J) supports the inferred link between coevolution and spatial compactness, and suggests that NC can help to guide experiments toward the residues of highest importance for a sector's function (Fig. 8). Beyond the improvements from lowering background signal, NC also predicts sectors that are otherwise difficult to detect (Fig. 7L–O). In addition, some studies have demonstrated other applications of coevolution such as protein engineering²⁰ and variant interpretation¹⁴. Our results suggest that the utility of coevolution as a signal for protein science can be substantially improved by NC, opening new windows for broadly understanding protein structure–function relationships.

Methods and materials

MSA construction. MSAs were constructed with BLAST⁷⁵ to identify up to 10,000 closest sequences to a reference sequence, using the RefSeq database⁷⁶. Sequences were aligned with Clustal Omega⁷⁷. Sequences with a Jukes–Cantor distance > 1 from the reference sequence were pruned. Redundant sequences and positions with $> 25\%$ gaps were removed. Any remaining gaps were filled with the amino acid from the closest sequence in terms of Jukes–Cantor distance.

Calculation of the expected value of inter-clade covariation. For our analyses, we define a pair of sequences to be within the same clade if the phylogenetic distance is below a Jukes–Cantor distance d . The phylogenetic distance is measured with respect to the aligned protein sequence (Table S1). We sought to measure the expected value of residue–residue covariation due solely to the comparison of sequences between clades, which we refer to as the inter-clade covariation $C_{S>d}$. Below, we describe and compare measurement of the expected value of inter-clade covariation in Eq. (1) of the main text both by approximation via bootstrapping and analytically.

Bootstrapping. In this approximate method, we bootstrap the original MSA: for every position, we replace the amino acid with the identity of the same position from a random sequence in the same clade. For example, in Fig. 1Avi we show two positions in an MSA, colored by their clade membership for a given phylogenetic distance d . Note that the first position is never a glutamine in the orange clade and is never a threonine in the white clade. Similarly, the second position is never a serine in the orange clade and is never an arginine in the white clade. The bootstrapped MSAs resample within clades, so as to not change the phylogenetic structure of the MSA at distances $> d$; thus, the first position in the bootstrapped MSAs still does not contain a glutamine, etc. The covariation measured from each of the bootstrapped MSAs is averaged to obtain the matrix expected under the

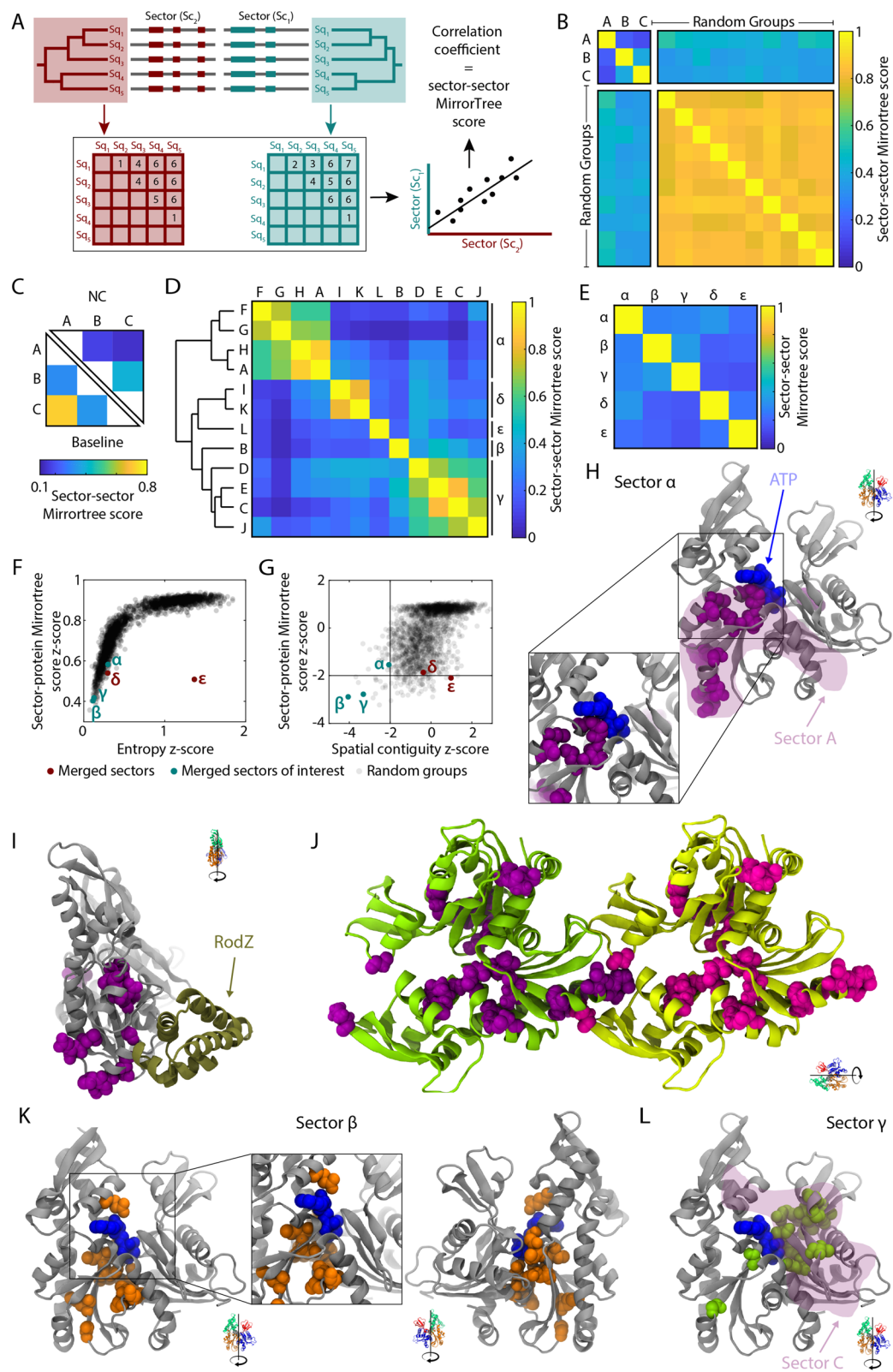


Figure 7. NC eigenvectors generally improve sector prediction across proteins, and enable identification of sectors that are not detectable using the baseline method. **(A,D,G)** NC eigenvectors for enolase **(A)**, G6PD **(D)**, and MAPK1 **(G)** exhibit lower background noise than the corresponding baseline (NMI with APC) eigenvectors. **(B,E,H,I)** The appropriate NC sectors (positive or negative values of the eigenvector) associated with the eigenvectors in **(A,D,G)** are more spatially contiguous across size cutoffs than the baseline sectors. Note that the MAPK1 eigenvector was split into a positive sector **(H)** and a negative sector **(I)**. Gray regions represent 95% confidence intervals of a randomly selected group of residues of the same size. **(C,F)** The 15-residue versions of the sectors in **(B,E)** on the crystal structures of enolase **(C)** and G6PD **(F)** illustrate the more compact nature of the NC sectors as compared with the baseline sectors. **(J,K)** The 50- and 20-residue versions of the NC sectors in **(H,I)** are more spatially compact on the structure than the corresponding baseline sectors, and occupy distinct parts of the protein. **(L–O)** For ArgS **(L)** and G6PD **(M)**, certain high-eigenvalue NC sectors had no obvious baseline counterpart. These NC sectors had low MirrorTree and spatial compactness z -scores **(L,M)**, and 15-residue versions occupied spatially compact regions around ligands [arginine in **(N)**, NADP in **(O)**] on the structure **(N,O)**. Thus, NC enables the detection of sectors that are otherwise obscured by phylogenetic bias.

hypothesis that there is no coupling between positions within the same clade. The bootstrapping method can be applied for any coevolution heuristic.

Analytical method. To derive an analytical solution in place of bootstrapping the NMI metric, we rephrased our aim as calculating the expected value of covariation between two positions under the assumption that the two positions are independent within a clade.

Consider the Shannon entropy for position i :

$$H_i = - \sum_{k=1}^{20} p_{i=k} \ln p_{i=k},$$

where $p_{i=k}$ is the probability of finding amino acid k at position i . The marginal probabilities of positions i and j taking on a particular value in a bootstrapped MSA do not change on average. However, the joint entropy, which relies on the joint probability, will change, as described below:

$$H_{ij} = - \sum_{k,l=1}^{20} p_{i=k,j=l} \ln p_{i=k,j=l}.$$

We seek an expression for the joint entropy that captures the assumption that positions i and j are independent within clades. Since the joint probability of independent variables is the product of the individual probabilities, we are left with calculating the sum of probabilities from each clade c , weighted by the number of sequences n_c in each clade:

$$p_{i=k,j=l}^{\text{null}} = \left(\sum_c n_c p_{i=k}^c p_{j=l}^c \right) / \left(\sum_c n_c \right)$$

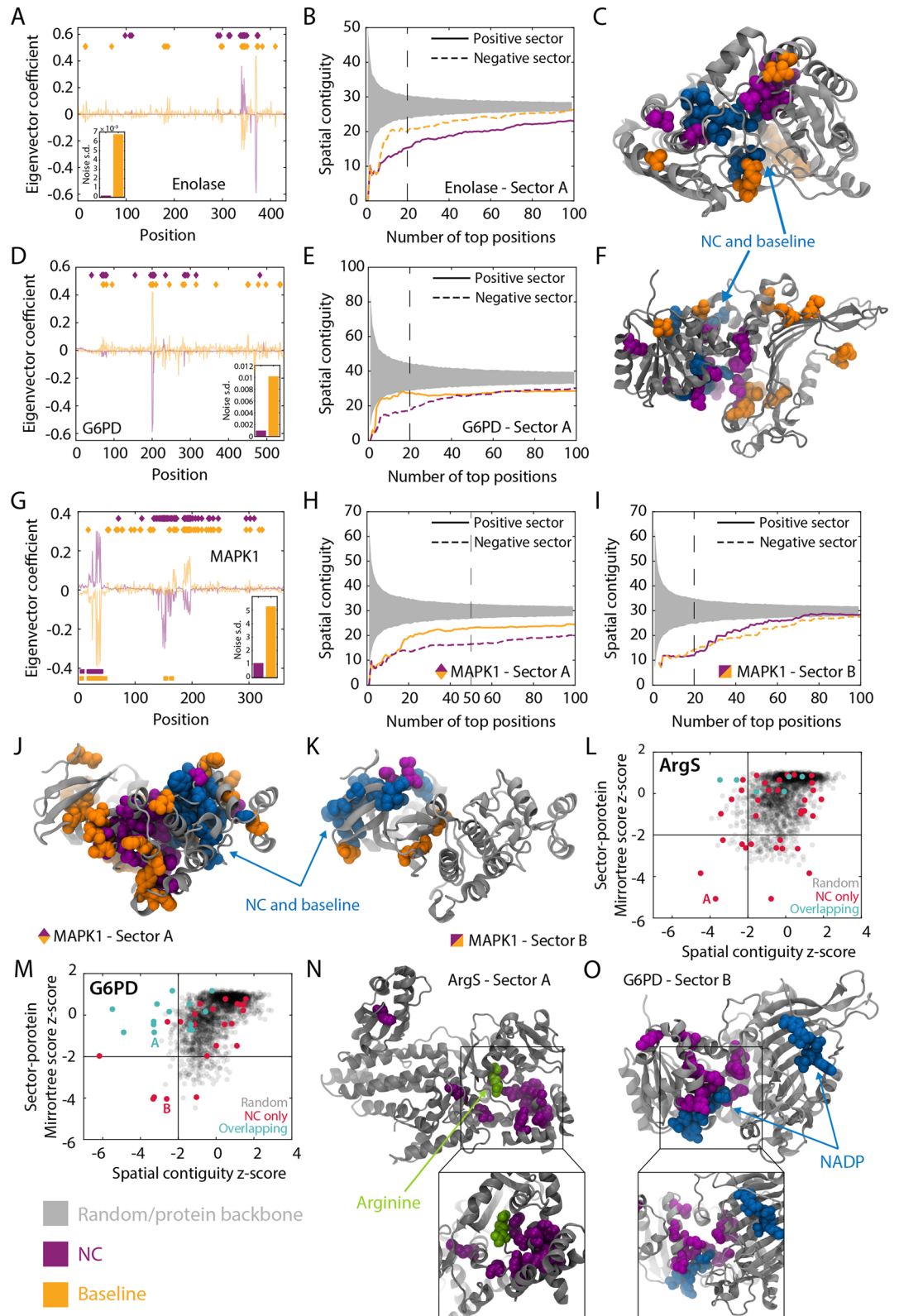
where $p_{i=k}^c$ is the marginal probability of finding amino acid k within clade c at position i .

A comparison of the bootstrapped and analytical methods for calculating NC for the yeast actin protein is shown in Fig. S2.

Estimating the statistical significance of nested coevolution. The expectation value of our NC background model is described above analytically only for NMI; other coevolution metrics do not have a known closed-form analytical solution, so we rely on bootstrapping to estimate the expected value. Bootstrapping offers the additional advantage of providing an estimate of the statistical significance of the observed raw coevolution signal by measuring the fraction of bootstrapped MSAs that achieves equal or greater coevolution values. The accuracy of the significance estimate is limited by the number of bootstrap measurements, since the maximum resolution is the reciprocal of the number of bootstraps performed. Using hundreds of bootstraps, we compared significance estimates with the absolute difference between the total and inter-clade covariation. These values were highly correlated (Spearman's $\rho = 0.95$, Fig. S2B), indicating that the difference between the baseline signal and either the bootstrapping or analytical method of computing NC provides a surrogate for the significance of the observation.

Structural contact prediction. Real structural contacts were determined by calculating the distance between the alpha carbons of every pair of residues in the protein based on a crystal structure (Table S1). All other atoms, including hydrogen atoms, were disregarded. To predict structural contacts, we used mean-field DCA with pseudocount value 0.5, and sequences closer than 0.3 Hamming distance were reweighted^{4,23}.

Generation of NC sectors. The output of NC is n_d p -by- p matrices (Fig. 1C), where n_d is the number of phylogenetic windows and p is the number of amino acids in the protein. These n_d matrices were concatenated to obtain a supermatrix of dimension $p n_d$ -by- p (Fig. S7). Eigenvalue decomposition or singular value decompo-



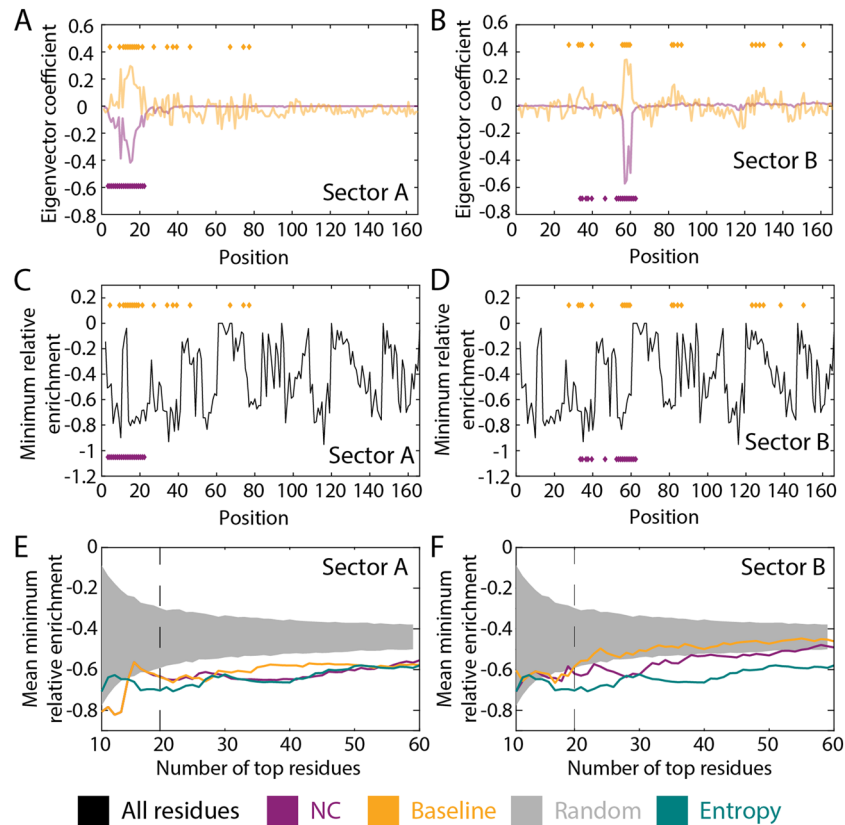


Figure 8. NC sectors predict deactivating mutations in H-Ras. **(A,B)** NC predicts two eigenvectors with much lower background noise than the baseline counterparts. The purple and gold diamonds represent the locations of residues in sectors of size 20. **(C,D)** Fitness data from a screen of binding efficacy of H-Ras to Raf-RBD⁶⁹. Shown is the minimum enrichment over all mutations at each position (thus representing maximum deactivation). The purple and gold diamonds represent the locations of residues in sectors of size 20. **(E,F)** Across most sector size cutoffs, the mean minimum relative enrichment was significantly lower than random (gray) for NC sectors A and B and comparable to that of the residues with the lowest entropy (teal). NC sectors also outperformed their baseline counterparts.

sition of the covariation in this super matrix was performed (thus avoiding the need to choose one value of the cutoff distance d), with pn_d observations and p features. The eigenvectors were ordered highest to lowest according to their associated eigenvalues. Each eigenvector is of length p , where the i th coefficient corresponds to the importance of the i th position in explaining the variation in the respective eigenvector.

To extract the specific positions that are most responsible for explaining the variation in a particular eigenvector, we identified the positions with the most positive or most negative coefficients and defined these groups of residues as two sectors. Sectors that had < 4 amino acids were ignored for downstream analysis.

NC and baseline sectors were paired if the dot product of the corresponding eigenvector was > 0.6 .

Calculating the spatial compactness of a sector. To quantify spatial compactness, we calculated the mean distance between the alpha carbon atoms of each pair of the most highly associated residues in the sector in the crystal structure. The calculation was repeated with increasing numbers of the most highly associated residues to avoid having to arbitrarily choose a single sector size.

Adaptation of the MirrorTree algorithm. Mirrortree was originally developed to predict protein–protein interactions based on the similarity of phylogenetic trees⁴⁹. In brief, MSAs are calculated using protein sequences from the same list of organisms for two proteins. For each MSA, the matrix of pairwise Jukes–Cantor distances is calculated. The MirrorTree score is the Pearson correlation coefficient of these two distance matrices. A high correlation indicates that the two proteins have similar phylogenies and thus are likely to have experienced similar functional selection. We adapted this method to compare the phylogenetic similarity of protein sectors with the entire protein (Fig. 4A) or other sectors (Fig. 6A). To compute sector–protein and sector–sector MirrorTree scores, filtered MSAs were created focusing on the positions of a given sector.

Biased sampling of random sectors was accomplished via weighting of residues according to their entropy.

Calculation of meta-sectors. Sets of sectors to be merged into meta-sectors were determined from hierarchical clustering based on sector–sector MirrorTree scores. Merging was accomplished by adding the corresponding eigenvectors after multiplying each sector by +1 or –1 corresponding to whether a positive or negative sector, respectively, was being merged. The summed vector was then analyzed as if it were an eigenvector in order to define meta-sectors at various size cutoffs.

Data availability

All multiple sequence alignments and associated data are available upon request from the corresponding author.

Code availability

All code is available at <https://github.com/acolavin/nested-coevolution>.

Received: 2 September 2021; Accepted: 17 December 2021

Published online: 17 January 2022

References

- Hollstein, M., Sidransky, B., Vogelstein, B. & Harris, C. C. P53 mutations in human cancers. *Science* **253**, 49–53 (1991).
- Zvelebil, M. J., Barton, G. J., Taylor, W. R. & Sternberg, M. J. E. Prediction of protein secondary structure and active sites using the alignment of homologous sequences. *J. Mol. Biol.* **195**, 957–961 (1987).
- Dwyer, R. S., Ricci, D. P., Colwell, L. J., Silhavy, T. J. & Wingreen, N. S. Predicting functionally informative mutations in *Escherichia coli* BamA using evolutionary covariance analysis. *Genetics* **195**, 443–455. <https://doi.org/10.1534/genetics.113.155861> (2013).
- Marks, D. S. *et al.* Protein 3D structure computed from evolutionary sequence variation. *PLoS ONE* <https://doi.org/10.1371/journal.pone.0028766> (2011).
- Morcos, F., Jana, B., Hwa, T. & Onuchic, J. N. Coevolutionary signals across protein lineages help capture multiple protein conformations. *Proc. Natl. Acad. Sci. U.S.A.* **110**, 20533–20538. <https://doi.org/10.1073/pnas.1315625110> (2013).
- Reynolds, K. A., McLaughlin, R. N. & Ranganathan, R. Hot spots for allosteric regulation on protein surfaces. *Cell* **147**, 1564–1575. <https://doi.org/10.1016/j.cell.2011.10.049> (2011).
- Weigt, M., White, R. A., Szurmant, H., Hoch, J. A. & Hwa, T. Identification of direct residue contacts in protein–protein interaction by message passing. *Proc. Natl. Acad. Sci. U.S.A.* **106**, 67–72. <https://doi.org/10.1073/pnas.0805923106> (2009).
- Bitbol, A. F. Inferring interaction partners from protein sequences using mutual information. *PLoS Comput. Biol.* **14**, e1006401. <https://doi.org/10.1371/journal.pcbi.1006401> (2018).
- Bitbol, A. F., Dwyer, R. S., Colwell, L. J. & Wingreen, N. S. Inferring interaction partners from protein sequences. *Proc. Natl. Acad. Sci. U.S.A.* **113**, 12180–12185. <https://doi.org/10.1073/pnas.1606762113> (2016).
- Burger, L. & van Nimwegen, E. Accurate prediction of protein–protein interactions from sequence alignments using a Bayesian method. *Mol. Syst. Biol.* **4**, 165. <https://doi.org/10.1038/msb4100203> (2008).
- Gueudre, T., Baldassi, C., Zamparo, M., Weigt, M. & Pagnani, A. Simultaneous identification of specifically interacting paralogs and interprotein contacts by direct coupling analysis. *Proc. Natl. Acad. Sci. U.S.A.* **113**, 12186–12191. <https://doi.org/10.1073/pnas.1607570113> (2016).
- Ovchinnikov, S., Kamisetty, H. & Baker, D. Robust and accurate prediction of residue–residue interactions across protein interfaces using evolutionary information. *Elife* **3**, e02030 (2014).
- Rivoire, O. Parsimonious evolutionary scenario for the origin of allostery and coevolution patterns in proteins. *Phys. Rev. E* **100**, 032411. <https://doi.org/10.1103/PhysRevE.100.032411> (2019).
- Hopf, T. A. *et al.* Mutation effects predicted from sequence co-variation. *Nat. Biotechnol.* **35**, 128–135. <https://doi.org/10.1038/nbt.3769> (2017).
- Frazer, J. *et al.* Disease variant prediction with deep generative models of evolutionary data. *Nature* **599**, 91–95. <https://doi.org/10.1038/s41586-021-04043-8> (2021).
- Altschuh, D., Vernet, T., Moras, D. & Nagai, K. Coordinated amino acid changes in homologous protein families. *Protein Eng.* **2**, 193–199 (1988).
- Atchley, W., Wollenberg, K., Fitch, W., Terhalle, W. & Dress, A. Correlations among amino acid sites in bHLH protein domains: An information theoretic analysis. *Mol. Biol. Evol.* **17**, 164–178 (2000).
- Göbel, U., Sander, C., Schneider, R. & Valencia, A. Correlated mutations and residue contacts in proteins. *Proteins* **18**, 309–317. <https://doi.org/10.1002/prot.340180402> (1994).
- Halabi, N., Rivoire, O., Leibler, S. & Ranganathan, R. Protein sectors: Evolutionary units of three-dimensional structure. *Cell* **138**, 774–786. <https://doi.org/10.1016/j.cell.2009.07.038> (2009).
- Skerker, J. M. *et al.* Rewiring the specificity of two-component signal transduction systems. *Cell* **133**, 1043–1054. <https://doi.org/10.1016/j.cell.2008.04.040> (2008).
- Socolich, M. *et al.* Evolutionary information for specifying a protein fold. *Nature* **437**, 512–518. <https://doi.org/10.1038/nature03991> (2005).
- Dunn, S. D., Wahl, L. M. & Gloor, G. B. Mutual information without the influence of phylogeny or entropy dramatically improves residue contact prediction. *Bioinformatics (Oxford, England)* **24**, 333–340. <https://doi.org/10.1093/bioinformatics/btm604> (2008).
- Morcos, F. *et al.* Direct-coupling analysis of residue coevolution captures native contacts across many protein families. *Proc. Natl. Acad. Sci. U.S.A.* **108**, E1293–1301. <https://doi.org/10.1073/pnas.1111471108> (2011).
- Lichtarge, O., Bourne, H. R. & Cohen, F. E. An evolutionary trace method defines binding surfaces common to protein families. *J. Mol. Biol.* **257**, 342–358. <https://doi.org/10.1006/jmbi.1996.0167> (1996).
- Wilkins, A., Erdin, S., Lua, R. & Lichtarge, O. Evolutionary trace for prediction and redesign of protein functional sites. *Methods Mol. Biol.* **819**, 29–42. https://doi.org/10.1007/978-1-61779-465-0_3 (2012).
- Sung, Y. M., Wilkins, A. D., Rodriguez, G. J., Wensel, T. G. & Lichtarge, O. Intramolecular allosteric communication in dopamine D2 receptor revealed by evolutionary amino acid covariation. *Proc. Natl. Acad. Sci. U.S.A.* **113**, 3539–3544. <https://doi.org/10.1073/pnas.1516579113> (2016).
- Katsonis, P. & Lichtarge, O. A formal perturbation equation between genotype and phenotype determines the evolutionary action of protein-coding variations on fitness. *Genome Res.* **24**, 2050–2058. <https://doi.org/10.1101/gr.176214.114> (2014).
- Hockenberry, A. J. & Wilke, C. O. Phylogenetic weighting does little to improve the accuracy of evolutionary coupling analyses. *Entropy (Basel)* <https://doi.org/10.3390/e21101000> (2019).
- Vorberg, S., Seemayer, S. & Soding, J. Synthetic protein alignments with CCMgen quantify noise in residue–residue contact prediction. *PLoS Comput. Biol.* **14**, e1006526. <https://doi.org/10.1371/journal.pcbi.1006526> (2018).
- Qin, C. & Colwell, L. J. Power law tails in phylogenetic systems. *Proc. Natl. Acad. Sci. U.S.A.* **115**, 690–695. <https://doi.org/10.1073/pnas.1711913115> (2018).

31. Horta, E. R. & Weigt, M. On the effect of phylogenetic correlations in coevolution-based contact prediction in proteins. *PLoS Comput. Biol.* **17**, 032601 (2021).
32. Rodriguez Horta, E., Barrat-Charlaix, P. & Weigt, M. Toward inferring Potts models for phylogenetically correlated sequence data. *Entropy (Basel)* **21**, 1090 (2019).
33. Malinverni, D. & Barducci, A. Coevolutionary analysis of protein subfamilies by sequence reweighting. *Entropy (Basel)* **21**, 1127. <https://doi.org/10.3390/e21111127> (2020).
34. Malinverni, D., Marsili, S., Barducci, A. & De Los Rios, P. Large-scale conformational transitions and dimerization are encoded in the amino-acid sequences of Hsp70 chaperones. *PLoS Comput. Biol.* **11**, e1004262. <https://doi.org/10.1371/journal.pcbi.1004262> (2015).
35. Martin, L. C., Gloor, G. B., Dunn, S. D. & Wahl, L. M. Using information theory to search for co-evolving residues in proteins. *Bioinformatics* **21**, 4116–4124. <https://doi.org/10.1093/bioinformatics/bti671> (2005).
36. Wollenberg, K. R. & Atchley, W. R. Separation of phylogenetic and functional associations in biological sequences by using the parametric bootstrap. *Proc. Natl. Acad. Sci. U.S.A.* **97**, 3288–3291. <https://doi.org/10.1073/pnas.070154797> (2000).
37. Garcia-Mayoral, M. F. *et al.* The structure of the C-terminal KH domains of KSRP reveals a noncanonical motif important for mRNA degradation. *Structure* **15**, 485–498. <https://doi.org/10.1016/j.str.2007.03.006> (2007).
38. Golub, G. H. & Van Loan, C. F. *Matrix Computations* 3rd edn. (Johns Hopkins University Press, 1996).
39. Ekeberg, M., Lovkvist, C., Lan, Y., Weigt, M. & Aurell, E. Improved contact prediction in proteins: Using pseudolikelihoods to infer Potts models. *Phys. Rev. E Stat. Nonlinear Soft Matter Phys.* **87**, 012707. <https://doi.org/10.1103/PhysRevE.87.012707> (2013).
40. Wang, S. W., Bitbol, A. F. & Wingreen, N. S. Revealing evolutionary constraints on proteins through sequence analysis. *PLoS Comput. Biol.* **15**, e1007010. <https://doi.org/10.1371/journal.pcbi.1007010> (2019).
41. McLaughlin, R. N. Jr., Poelwijk, F. J., Raman, A., Gosal, W. S. & Ranganathan, R. The spatial architecture of protein function and adaptation. *Nature* **491**, 138–142. <https://doi.org/10.1038/nature11500> (2012).
42. Novinec, M. *et al.* A novel allosteric mechanism in the cysteine peptidase cathepsin K discovered by computational methods. *Nat. Commun.* **5**, 3287. <https://doi.org/10.1038/ncomms4287> (2014).
43. Rivoire, O., Reynolds, K. A. & Ranganathan, R. Evolution-based functional decomposition of proteins. *PLoS Comput. Biol.* **12**, e1004817. <https://doi.org/10.1371/journal.pcbi.1004817> (2016).
44. Smock, R. G. *et al.* An interdomain sector mediating allostery in Hsp70 molecular chaperones. *Mol. Syst. Biol.* **6**, 414. <https://doi.org/10.1038/msb.2010.65> (2010).
45. Shi, H., Bratton, B. P., Gitai, Z. & Huang, K. C. How to build a bacterial cell: MreB as the foreman of *E. coli* construction. *Cell* **172**, 1294–1305. <https://doi.org/10.1016/j.cell.2018.02.050> (2018).
46. Izore, T., Duman, R., Kureisaite-Ciziene, D. & Lowe, J. Crenactin from *Pyrobaculum calidifontis* is closely related to actin in structure and forms steep helical filaments. *FEBS Lett.* **588**, 776–782. <https://doi.org/10.1016/j.febslet.2014.01.029> (2014).
47. van den Ent, F., Amos, L. A. & Lowe, J. Prokaryotic origin of the actin cytoskeleton. *Nature* **413**, 39–44. <https://doi.org/10.1038/35092500> (2001).
48. van den Ent, F., Amos, L. & Lowe, J. Bacterial ancestry of actin and tubulin. *Curr. Opin. Microbiol.* **4**, 634–638. [https://doi.org/10.1016/s1369-5274\(01\)00262-4](https://doi.org/10.1016/s1369-5274(01)00262-4) (2001).
49. Craig, R. A. & Liao, L. Phylogenetic tree information aids supervised learning for predicting protein–protein interaction based on distance matrices. *BMC Bioinform.* **8**, 6. <https://doi.org/10.1186/1471-2105-8-6> (2007).
50. Pazos, F. & Valencia, A. Similarity of phylogenetic trees as indicator of protein–protein interaction. *Protein Eng.* **14**, 609–614. <https://doi.org/10.1093/protein/14.9.609> (2001).
51. Araya, C. L. *et al.* Identification of significantly mutated regions across cancer types highlights a rich landscape of functional molecular alterations. *Nat. Genet.* **48**, 117–125. <https://doi.org/10.1038/ng.3471> (2016).
52. Hu, Z., Ma, B., Wolfson, H. & Nussinov, R. Conservation of polar residues as hot spots at protein interfaces. *Proteins* **39**, 331–342 (2000).
53. Ptitsyn, O. B. Protein folding and protein evolution: Common folding nucleus in different subfamilies of c-type cytochromes?. *J. Mol. Biol.* **278**, 655–666. <https://doi.org/10.1006/jmbi.1997.1620> (1998).
54. Teşileanu, T., Colwell, L. J. & Leibler, S. Protein sectors: Statistical coupling analysis versus conservation. *PLOS Comput. Biol.* **11**, e1004091. <https://doi.org/10.1371/journal.pcbi.1004091> (2015).
55. Anishchenko, I., Ovchinnikov, S., Kamisetty, H. & Baker, D. Origins of coevolution between residues distant in protein 3D structures. *Proc. Natl. Acad. Sci. U.S.A.* **114**, 9122–9127. <https://doi.org/10.1073/pnas.1702664114> (2017).
56. Bratton, B. P., Shaevitz, J. W., Gitai, Z. & Morgenstein, R. M. MreB polymers and curvature localization are enhanced by RodZ and predict *E. coli*'s cylindrical uniformity. *Nat. Commun.* **9**, 2797. <https://doi.org/10.1038/s41467-018-05186-5> (2018).
57. Colavin, A., Shi, H. & Huang, K. C. RodZ modulates geometric localization of the bacterial actin MreB to regulate cell shape. *Nat. Commun.* **9**, 1280. <https://doi.org/10.1038/s41467-018-03633-x> (2018).
58. Spring, T. G. & Wold, F. The purification and characterization of *Escherichia coli* enolase. *J. Biol. Chem.* **246**, 6797–6802 (1971).
59. Wright, D. N. & Lockhart, W. R. Effects of growth rate and limiting substrate on glucose metabolism in *Escherichia coli*. *J. Bacteriol.* **89**, 1082–1085 (1965).
60. Pelech, S. L., Sanghera, J. S. & Daya-Makin, M. Protein kinase cascades in meiotic and mitotic cell cycle control. *Biochem. Cell Biol.* **68**, 1297–1330. <https://doi.org/10.1139/o90-194> (1990).
61. Sturgill, T. W. & Wu, J. Recent progress in characterization of protein kinase cascades for phosphorylation of ribosomal protein S6. *Biochim. Biophys. Acta* **1092**, 350–357. [https://doi.org/10.1016/s0167-4889\(97\)00012-4](https://doi.org/10.1016/s0167-4889(97)00012-4) (1991).
62. Hirshfield, I. N. & Bloemers, H. P. The biochemical characterization of two mutant arginyl transfer ribonucleic acid synthetases from *Escherichia coli* K-12. *J. Biol. Chem.* **244**, 2911–2916 (1969).
63. Dove, S. L., Joung, J. K. & Hochschild, A. Activation of prokaryotic transcription through arbitrary protein–protein contacts. *Nature* **386**, 627–630. <https://doi.org/10.1038/386627a0> (1997).
64. Joung, J. K., Ramm, E. I. & Pabo, C. O. A bacterial two-hybrid selection system for studying protein–DNA and protein–protein interactions. *Proc. Natl. Acad. Sci. U.S.A.* **97**, 7382–7387. <https://doi.org/10.1073/pnas.110149297> (2000).
65. Lim, W. A. & Sauer, R. T. Alternative packing arrangements in the hydrophobic core of lambda repressor. *Nature* **339**, 31–36. <https://doi.org/10.1038/339031a0> (1989).
66. Johnson, C. W. *et al.* The small GTPases K-Ras, N-Ras, and H-Ras have distinct biochemical properties determined by allosteric effects. *J. Biol. Chem.* **292**, 12981–12993. <https://doi.org/10.1074/jbc.M117.778886> (2017).
67. Wellbrock, C., Karasarides, M. & Marais, R. The RAF proteins take centre stage. *Nat. Rev. Mol. Cell Biol.* **5**, 875–885. <https://doi.org/10.1038/nrml498> (2004).
68. Prior, I. A., Lewis, P. D. & Mattos, C. A comprehensive survey of Ras mutations in cancer. *Cancer Res.* **72**, 2457–2467. <https://doi.org/10.1158/0008-5472.CAN-11-2612> (2012).
69. Bandaru, P. *et al.* Deconstruction of the Ras switching cycle through saturation mutagenesis. *Elife* <https://doi.org/10.7554/eLife.27810> (2017).
70. Cocco, S., Monasson, R. & Weigt, M. From principal component to direct coupling analysis of coevolution in proteins: Low-eigenvalue modes are needed for structure prediction. *PLoS Comput. Biol.* **9**, e1003176. <https://doi.org/10.1371/journal.pcbi.1003176> (2013).

71. Rivoire, O. Elements of coevolution in biological sequences. *Phys. Rev. Lett.* **110**, 178102. <https://doi.org/10.1103/PhysRevLett.110.178102> (2013).
72. Do, R. *et al.* Exome sequencing identifies rare LDLR and APOA5 alleles conferring risk for myocardial infarction. *Nature* **518**, 102–106. <https://doi.org/10.1038/nature13917> (2015).
73. Nguyen, H. Q. *et al.* Quantitative mapping of protein–peptide affinity landscapes using spectrally encoded beads. *Elife* <https://doi.org/10.7554/eLife.40499> (2019).
74. Senior, A. W. *et al.* Improved protein structure prediction using potentials from deep learning. *Nature* **577**, 706–710. <https://doi.org/10.1038/s41586-019-1923-7> (2020).
75. Madden, T. The BLAST sequence analysis tool. 2002 Oct 9 [Updated 2003 Aug 13]. In *The NCBI Handbook [Internet]*. (National Center for Biotechnology Information (US), 2002).
76. Tatusova, T., Ciufu, S., Fedorov, B., O'Neill, K. & Tolstoy, I. RefSeq microbial genomes database: New representation and annotation strategy. *Nucleic Acids Res.* **42**, D553–559. <https://doi.org/10.1093/nar/gkt1274> (2014).
77. Sievers, F. *et al.* Fast, scalable generation of high-quality protein multiple sequence alignments using Clustal Omega. *Mol. Syst. Biol.* **7**, 539. <https://doi.org/10.1038/msb.2011.75> (2011).

Acknowledgements

The authors thank the Huang lab for useful discussions. This work was supported by Stanford Graduate Fellowships (to A.C. and E.A.), a Gerald J. Lieberman Fellowship (to A.C.), grant 851173 from the European Research Council under the European Union's Horizon 2020 research and innovation programme (to A.-F.B.), NSF CAREER Award MCB-1149328 (to K.C.H.), and the Allen Discovery Center at Stanford on Systems Modeling of Infection (to K.C.H.). K.C.H. is a Chan Zuckerberg Biohub Investigator.

Author contributions

A.C., E.A., and K.C.H. designed the research; A.C. and E.A. performed the research; A.C., E.A., A.-F.B., and K.C.H. analyzed the data; and A.C., E.A., A.-F.B., and K.C.H. wrote the paper. All authors reviewed it before submission.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-021-04260-1>.

Correspondence and requests for materials should be addressed to K.C.H.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022