**1 Extracting the phylogenetic dimension of coevolution reveals hidden**

**2 functional signal**

3

4 Alexandre Colavin[1,*], Esha Atolia[2,*], Anne-Florence Bitbol[3,4], Kerwyn Casey Huang[5,6,7,†]

5

6 [1]Biophysics Program, Stanford University School of Medicine, Stanford, CA 94305, USA

7 [2]Department of Chemical and Systems Biology, Stanford University School of Medicine,

8 Stanford, CA 94305, USA

9 [3]Sorbonne Université, CNRS, Institut de Biologie Paris-Seine, Laboratoire Jean Perrin

10 (UMR 8237), F-75005, Paris, France

11 [4]Institute of Bioengineering, School of Life Sciences, Ecole Polytechnique Fédérale de

12 Lausanne (EPFL), CH-1015 Lausanne, Switzerland

13 [5]Department of Bioengineering, Stanford University, Stanford, CA 94305, USA

14 [6]Department of Microbiology and Immunology, Stanford University School of Medicine,

15 Stanford, CA 94305, USA

16 [7]Chan Zuckerberg Biohub, San Francisco, CA 94158

17

18 [*]: Co-first authors.

19    *Keywords:* MirrorTree, protein sectors, normalized joint entropy, Direct Information,

20    Average Product Correction, entropy, MreB, H-Ras, ArgS, G6PD, Enolase, MAPK1,

21    Nested Coevolution

22

23    †To whom correspondence should be addressed: kchuang@stanford.edu

24 **Abstract**

25

26 Despite the structural and functional information contained in the statistical coupling

27 between pairs of residues in a protein, coevolution associated with function is often

28 obscured by artifactual signals such as genetic drift, which shapes a protein's

29 phylogenetic history and gives rise to concurrent variation between protein sequences

30 that is not driven by selection for function. Here, we introduce a method for explicitly

31 defining a phylogenetic dimension of coevolution signal, and demonstrate that

32 coevolution can occur on multiple phylogenetic timescales within a single protein. Our

33 method, Nested Coevolution (NC), can be applied as an extension to any coevolution

34 metric. We use NC to demonstrate that poorly conserved residues can nonetheless have

35 important roles in protein function. Moreover, NC improved structural-contact

36 prediction over gold-standard coevolution-based methods, particularly in subsampled

37 alignments with fewer sequences. NC also lowered the noise in detecting functional

38 sectors of collectively coevolving residues. Sectors of coevolving residues identified after

39 NC correction were more spatially compact and phylogenetically distinct from the rest

40 of the protein, and strongly enriched for mutations that disrupt protein activity. Our

41 conceptualization of the phylogenetic separation of coevolution represents an advance

42 from previous pragmatic attempts to reduce phylogenetic artifacts in measurements of

43 coevolution. Application of NC broadens the application of protein coevolution

44    measurements, particularly to eukaryotic proteins with fewer naturally available

45    sequences, and further elucidates relationships among protein evolution and genetic

46    diseases.

## Introduction

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 interfaces (12), intuiting novel protein conformations (5), understanding protein allostery (13), interpreting variants (14), identifying functional domains (15-18), and reprograming protein specificity (19). However, despite the increasing prevalence of sequencing data, sampling of the phylogenetic tree is necessarily 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.

Of the numerous existing methods for measuring protein coevolution, many implement methods for reducing the effects of phylogenetic noise. Although mutual information is extremely sensitive to the phylogenetic distribution of sequences and the conservation (measured via entropy) of individual positions, normalization by the joint entropy

67 reduces the influence of phylogeny and entropy and improves structural-contact

68 prediction (20). Statistical coupling analysis , which normalizes the covariance matrix by

69 a function of the entropy, provides sufficient information to specify a protein fold (21)

70 and to detect functional domains (6, 18). Direct coupling analysis (DCA) usually involves

71 down-weighting the coevolutionary signal contributions from over-represented

72 sequences, and attempts to deconvolve higher-order correlations to identify directly

73 interacting residue pairs (4, 22). Motivated by the observed strong relationship between

74 a position's average mutual information and the mutual information it exhibits with

75 specific positions, modifications such as the average product correction (APC) subtract

76 this average signal; this correction can be applied to any existing coevolution metric other

77 than mutual information. However, none of these strategies attempt to resolve the

78 evolutionary timescale of coevolution.

79

80 Even with affordable sequencing and widespread environmental sampling, coevolution

81 methods are often limited by the number of naturally occurring protein sequences

82 available. Successful predictions of structural contacts often require several thousand

83 sequences to align (3, 23), which is generally prohibitive for many mammalian proteins.

84 For other proteins, the phylogenetic distribution of available sequences is skewed by

85 sampling and is well recognized as a source of spurious signal in coevolution (20, 24).

86 Thus, methods that enable the separation of functional coupling from phylogenetic and

87    sampling noise would greatly expand the utility of coevolution, particularly for

88    applications to diseases involving human proteins with limited numbers of available

89    sequences.

90

91    Here, we introduce the concept of Nested Coevolution (NC), a correction that leverages

92    a well-defined null hypothesis to accurately measure the coevolutionary signal above

93    what is expected from phylogenetic distribution alone. We determined that NC results in

94    higher fidelity of the coevolutionary signal across gold-standard coevolution-based

95    metrics for structural prediction for many proteins, especially with fewer sequences. In

96    addition, we found that NC improves the detection of spatially contiguous groups of

97    collectively coevolving residues ("sectors") that are phylogenetically distinct from each

98    other and the protein itself, beyond differences in entropy alone. Finally, sectors

99    identified using NC were enriched for positions at which mutations are maximally

100   deleterious, highlighting the functional significance of signal from our method. Since our

101   method is agnostic to the underlying method of measuring coevolution, we anticipate

102   wide utility for the ability to resolve the temporal dimension of protein coevolution.

103 **Results**

104

105 *Background model of coevolution reveals temporal dimension of coevolution*

106 To interrogate the contribution of phylogenetic sampling to protein coevolution

107 measurements, we sought to separate the coevolution signal due to inter-clade and intra-

108 clade sequence comparisons (Fig. 1A,B). Given a multiple sequence alignment (MSA) for

109 a protein of interest (Fig. 1Ai), we first measure the total covariation ($C_T$) between every

110 pair of positions (Fig. 1Aii) using an established metric of residue-residue coupling such

111 as the normalized mutual information (NMI; Fig. 1A) (20):

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

113 where $H_i$ is the Shannon entropy (a measure of conservation) of position $i$, and $H_{ij}$ is the

114 joint Shannon entropy of positions $i$ and $j$. The quantity $H_i + H_j - H_{ij}$ is the mutual

115 information between positions $i$ and $j$, which measures the coupling between residues

116 (Fig. S1). The NMI residue pair covariation in Eq. 1 is a natural metric choice because

117 normalizing by $H_{ij}$ makes the mutual information independent of conservation (20). Note

118 that our algorithm can be applied to any covariation metric, and as we will show, our

119 main results are robust to metric choice.

120

121 The most straightforward null hypothesis for protein coevolution is that coevolutionary

122 coupling between pairs of proteins is completely absent—that is, that the probability of a

123    position having any particular amino acid identity is independent of any other position's

124    identity. Although this null hypothesis can be evaluated analytically for some methods

125    (SI), other methods have no known closed-form solution for the expected value of the

126    coevolution matrix under these conditions. Hence, we computationally compute the

127    average coevolution signal from many globally resampled MSAs in which each position

128    in each protein in the original MSA is replaced by the equivalent position from another

129    randomly chosen protein (resampled with replacement; Fig. 1Av,vi). We expect any

130    measured coevolution from these resampled matrices to represent signal due simply to

131    the distribution of amino acid identities at each position; any significant difference

132    between the coevolution signal measured in the original MSA and this null hypothesis

133    can potentially be attributed to coevolution.

134

135    However, this initial null hypothesis does not test for the phylogenetic structure of

136    sequences; in the globally resampled MSAs, every sequence is effectively evolutionarily

137    equidistant from one another. Previous attempts to remove the influence of phylogeny

138    such as APC (Fig. 1Aiii), which corrects the covariation matrix by subtracting the product

139    of its mean value across columns and rows for each pair of positions (Fig. 1Aiv), have

140    substantially improved contact prediction (20). However, the APC is a postulated

141    correction that does not directly take into account the phylogenetic structure of an MSA.

142    We sought to construct a null hypothesis-driven background model of the expected

143    coevolution in an MSA in which intra-clade coevolution is explicitly removed. We

144    achieve this goal by generating MSAs by resampling each position from sequences that

145    are closely related (Fig. 1Av,vi), thus removing correlations arising from recent

146    evolutionary history within each clade. We define a clade as the subset of sequences $S$

147    with a Jukes-Cantor distance below $d$, which we refer to as the phylogenetic cutoff. For

148    each value of $d$, we calculate the inter-clade covariation $(C_{S>d}^{i,j})$ from a resampled MSA

149    either analytically or via bootstrapping (Fig. 1Avii, S2A, Methods), where $C$ denotes the

150    chosen covariation measure (e.g. NMI). This inter-clade covariation thus measures the

151    expected value of covariation due solely to the comparison of sequences between clades

152    (Fig. 1B). We then average over many such null hypotheses (over many within-clade

153    resampled MSAs at fixed $d$), yielding the mean inter-clade covariation matrix $(C_{S>d})$ (Fig.

154    1Aviii), which represents the expected coevolution due to both the distribution of amino

155    acid identities at each position and the phylogenetic structure of the protein MSA (Fig.

156    1B). Significant differences between this background model and the baseline signal

157    measured from the original MSA represent signal that was contained in the intra-clade

158    comparison of closely related sequences. Since the difference between the background

159    model and the baseline signal qualitatively captures the significance of the baseline

160    measurement (Fig S2B, Methods), we subtract $C_{S>d}$ from the total covariation $C_T$ to obtain

161    the phylogenetic cutoff-dependent covariation signal $C_{S\leq d}$ (Fig. 1Aix); positive values

162    indicate that the total covariation is larger than expected by comparison of sequences

163    between clades, thus revealing covariation arising from recent evolutionary history in all

164    clades:

165
$$C_{S \leq d}^{i,j} \equiv C_T^{i,j} - C_{S>d}^{i,j}. \qquad (2)$$

166    We refer to the signal $C_{S \leq d}^{i,j}$ above the null hypothesis $C_{S>d}^{i,j}$ in Eq. 2 as a protein's "nested

167    coevolution" (NC), in that it separates coevolution signal into signal attributed to

168    comparison of sequences either within ($C_{S \leq d}^{i,j}$) or between ($C_{S>d}^{i,j}$) nested clades of a

169    phylogenetic tree. The only free parameter in the NC is the phylogenetic cutoff ($d$). As

170    we vary the cutoff value, many patterns of NC typically emerge, revealing distinct

171    windows of coevolution for a single protein MSA (Fig. 1C). The changes in NC observed

172    between two cutoffs represent the signal due to pairs of sequences whose distance is

173    between the cutoffs used to calculate each window. Hence, distinct evolutionary

174    timescales of protein coevolution are revealed as the phylogenetic cutoff is varied.

175

176    To test the relevance of NC windows to protein structure prediction, we measured the

177    enrichment of structural contacts from the pairs of residues with the highest 50 values in

178    the NC matrix $C_{S \leq d}^{i,j}$ for each value of $d$. Here, we applied NC as a correction to DCA, the

179    current gold standard for coevolution-based prediction of structural contacts (4, 22). We

180    employed the direct information (DI) metric to quantify coevolution (4, 22). In this and

181    subsequent analyses, we considered structural contacts to be within 5 Å at closest

182    approach, excluding pairs of residues within 5 amino acids on the sequence (Methods).

183    The NC phylogenetic cutoffs revealed a variety of improvements for the KH domain (Fig.

184    1D), which is present in a wide variety of nucleic acid-binding proteins (25). Some

185    windows generally outperformed DCA, without (Fig. 1E) or with (Fig. 1F) the APC.

186

187    To determine the added value of NC for other proteins and for another frequently utilized

188    coevolution metric, the Frobenius norm (26), which is frequently utilized in DCA as an

189    alternative to DI (27, 28), we carried out a DCA structural-contact analysis for 10 protein

190    family domains with DI or Frobenius norm (Methods). Across both metrics and all

191    proteins, NC improved the predictions of structural contacts (Fig. 2A), even relative to

192    the inclusion of APC (20). Hence, NC is a correction that enhances the predictive power

193    of state-of-the-art coevolution measurements.

194

195    *NC improves predictions of structural contacts using fewer sequences*

196    One common limitation for computing coevolution is the number of homologous

197    sequences available for constructing an MSA. To interrogate whether NC could still

198    accurately predict structural contacts with fewer sequences, we subsampled the MSAs of

199    10 proteins with different breadth (randomly selecting 10% or 1% of the sequences) or

200    depth (selecting the 10% or 1% of sequences most related to the protein used to construct

201    the MSA, Table S1) (Fig. 2B). NC improved structural contact prediction for a majority of

202    the subsampled MSAs when correcting DI without (Fig. 2C, S4A) or with (Fig. S3B)

203    application of APC. For the KH domain, more than twice as many true positives were

204    predicted after applying NC compared with DI+APC alone (Fig. S3A). Perhaps

205    unsurprisingly, breadth sampling generally performed better than depth sampling (Fig.

206    2C), indicating that accurate prediction is reliant on the sequences being sufficiently

207    distantly related. Nonetheless, for many proteins, the value of the NC correction was

208    enhanced when the number of homologous sequences was low, both for depth and

209    breadth samplings.

210

211    ***NC generates eigenvectors with increased fidelity, improving detection of spatially***

212    ***contiguous sets of coevolving residues***

213    Previous studies have utilized coevolution measurements to identify groups of residues

214    within a protein that are spatially contiguous on the tertiary structure and thus are

215    postulated to have a joint function (6, 18, 29-32). These "sectors" can by defined by a

216    variety of methods, such as the extreme-value residues of the eigenvectors of the

217    coevolution matrix with the largest eigenvalues (18), and have been proposed to reflect

218    independent biological properties such as catalytic efficiency and thermal stability (18).

219    Motivated by these successes, we sought to measure the effect of incorporating the

220    phylogenetic dimension revealed by NC when defining sectors of residues. Specifically,

221    we measured the NC- and APC-corrected coevolution using NMI across a range of

222    phylogenetic cutoffs, concatenating the results and performing eigendecomposition to

223    identify the most significant eigenvectors (Methods). The residues most strongly

224    associated with the positive or negative components of each resulting eigenvector are

225    considered a sector.

226

227    We first focused on MreB, an essential protein involved in cell-shape determination in

228    many rod-shaped bacteria (33). MreB belongs to a protein family that includes ParM,

229    FtsA, and MamK in bacteria, crenactin in archea, and actin in eukaryotes (34, 35). These

230    proteins are structural homologs characterized by a four-subdomain fold around an ATP-

231    binding pocket (35, 36), with very low sequence identity and disparate cellular functions.

232    Thus, we anticipated that the set of MreB homologs would have sufficient diversity to

233    support robust coevolution measurements, particularly functional sectors.

234

235    We compared NC-derived sectors with baseline sectors derived from eigenvectors of the

236    baseline coevolution matrix for MreB homologs. We identified the most closely related

237    baseline sectors for three NC eigenvectors with some of the highest eigenvalues, which

238    we refer to as eigenvectors A, B, and C (Methods). Each pair of NC and baseline

239    eigenvectors appeared similar, especially for the residues with the largest absolute

240    coefficients (Fig. 3A-C). However, the baseline eigenvectors exhibited much higher

241    variation of coefficients for residues across the protein (Fig. 3A-C). For eigenvectors A

242    and B, the NC-derived eigenvectors exhibited 32.8-fold and 38.3-fold lower standard

243    deviation (after removing the 50 highest and lowest coefficients) than the baseline-

244    derived eigenvectors, respectively (Fig. 3A,B). For eigenvector C, the baseline eigenvector

245    contained residues with both highly positive and highly negative coefficients, while the

246    high-magnitude coefficients of the NC eigenvector were solely positive (Fig. 3C); the

247    positive portion of the NC eigenvector again had substantially lower noise than the

248    baseline eigenvector (2.1-fold lower standard deviation, Fig. 3C).

249

250    Motivated by the distinct behaviors of the positive and negative components of

251    eigenvector C, we defined distinct positive and negative sectors (Methods) for each NC

252    and baseline eigenvector using a variable cutoff on the site contributions to adjust sector

253    size (as sectors are defined as the sets of amino acids with highest site contributions in a

254    given eigenvector). For different sector sizes, we quantified the spatial contiguity as the

255    mean pairwise distance between each residue within a sector. For sectors A-C (derived

256    from eigenvectors A-C), the first 5-9 residues exhibited approximately the same spatial

257    contiguity in the NC as in the baseline eigenvectors (Fig. 3D-F). However, as the cutoff

258    was increased, the NC sector remained more spatially compact than the baseline sector

259    (Fig. 3D-F). All three NC sectors were also more spatially contiguous than expected based

260    on random sampling for cutoffs yielding at least 50 residues (Fig. 3D-F), while the

261    baseline sector A was distributed across the protein structure (Fig. 3D,J). NC sectors A

262    and C were largely situated in subdomains IIA (Fig. 3G) and IA (Fig. 3I), respectively,

263   while sector B was localized to the ATP-binding pocket (Fig. 3H). Notably, sector C was

264   spatially contiguous (Fig. 3I) despite being spread across the protein sequence (Fig. 3C).

265   Baseline sectors B and C with 15 residues were qualitatively similar to the corresponding

266   NC sectors (Fig. 3K,L); the large background fluctuations of the baseline eigenvector

267   likely led to the inclusion of additional, erroneous residues into the sector prediction.

268   Thus, the phylogenetic correction of NC improves the fidelity of sector detection as

269   measured by the spatial contiguity of its constituent residues.

270

271   *Sectors display distinct phylogenetic signatures from the rest of the protein*

272   Since sectors have been postulated to reflect distinct evolutionary histories driven by

273   selection for particular biological functions (18), we sought to compare the phylogeny of

274   the residues within a sector with other sectors and the rest of the protein. The MirrorTree

275   algorithm (Methods) was originally developed to compare phylogenies of two proteins,

276   motivated by the assumption that similar histories signifies a common function, e.g.

277   through protein-protein interactions and/or acting in the same pathway (37, 38). After

278   computing a pairwise distance matrix of all sequences within an MSA for each of the two

279   proteins based on homologs in the same set of organisms, the MirrorTree score is defined

280   as the Pearson correlation coefficient between the entries in the two pairwise distance

281   matrices (37). We straightforwardly modified the MirrorTree method to compare the

282    complete protein MSA to the MSA filtered to include only the residues within the sector

283    of interest (Fig. 4A).

284

285    To broadly investigate sector identification, we identified 40 15-residue sectors for MreB

286    based on the positive and negative coefficients of the 20 eigenvectors with the highest

287    eigenvalues. As negative controls, we randomly sampled sets of residues of the same size

288    as each sector from across the protein. Sector-protein MirrorTree scores for sectors A-C

289    (Fig. 3) were substantially lower for sectors than for the random groups (Fig. 4B), which

290    all had MirrorTree scores close to 1, as expected (Fig. 4B). Baseline sectors A-C had

291    MirrorTree scores intermediate between those of the corresponding NC sector and

292    random groups (Fig. 4B), likely reflecting the noisy selection from baseline eigenvectors

293    of residues that functionally follow the phylogenetic history of the protein overall. To

294    evaluate the significance of the MirrorTree score and of the spatial contiguity of each

295    sector, we computed $z$-scores based on the mean and standard deviation of the two

296    metrics applied to the random groups of the same size as each sector. Sectors A-C had

297    MirrorTree scores <0.5 (Fig. 4B), indicating distinct phylogenetic histories from the

298    protein, and MirrorTree and spatial contiguity $z$-scores<-2 (Fig. 4C). There were four

299    other sectors (D-G) that had spatial contiguity $z$-scores<-2. These sectors largely

300    overlapped with A-C; we will return to this overlap in a later section. All other sectors

301    had spatial contiguity $z$-score>2, and all but five (H-L) had MirrorTree $z$-score>-2. Thus,

302 MirrorTree reveals that certain NC sectors have distinct evolutionary trajectories from

303 the protein itself, motivating us to focus on certain sectors (such as A, B, and C for MreB).

304

305 *Phylogenetic similarity and the role of entropy*

306 Conservation itself is a major determinant of protein function (39-41), and spatially

307 contiguous sets of residues can be identified solely on the basis of conservation (42). To

308 account for variation in entropy across a protein, previous studies have excluded

309 positions with high conservation (Shannon entropy<0.1) or composed of >25% gaps in

310 the MSA (43). For MreB, NC sectors A-C had lower entropy than baseline sectors or

311 random groups of the same size (Fig. 5A-C), albeit higher entropy than residues typically

312 considered highly conserved (entropy<0.1).

313

314 MirrorTree scores of NC sectors were also generally lower than those of baseline sectors

315 (Fig. 5D-F). To investigate the dependence of sector-protein MirrorTree scores on

316 entropy, we computed MirrorTree scores for thousands of random groups of the same

317 size as the sector (15 residues), biasing sampling using a Monte Carlo algorithm to obtain

318 a wide range of mean entropies; each random group was selected from residues that did

319 not overlap with the sector. For mean entropy $\lesssim 1$, MirrorTree scores were strongly

320 dependent on entropy (Fig. 5G-I). Thus, the low MirrorTree scores of the NC sectors were

321 due in part to their low entropy. Nonetheless, the MirrorTree score of NC sector A was

322     significantly lower than those of random groups with the same mean entropy ($z$-score -

323     3.5); the entropy of sector B was so low, presumably due to the high conservation of the

324     ATP-binding pocket (Fig. 3H, S4), that it was challenging to obtain random groups that

325     were not largely overlapping.

326

327     Since NC sector A displayed the greatest reduction in MirrorTree score relative to

328     random groups of the same mean entropy, we focused on this sector to investigate the

329     dependence of the sector-protein MirrorTree score on sector size. As the cutoff was

330     increased to include more residues, the MirrorTree score increased (Fig. 5D). To

331     disentangle whether this increase was due directly to the increase in size or to the

332     inclusion of residues that are more phylogenetically similar to the protein, we compared

333     the 10-residue version of sector A (Fig. 5G) with randomly selected groups of 10 residues

334     from 15- and 20-residue versions of sector A, as well as the entire protein. The mean

335     MirrorTree score increased as the size of the sampling group increased (Fig. 5J), even for

336     groups with similar entropy as the 10-residue sector (Fig. 5K). Moreover, 15-residue

337     versions of sectors B and C had similar entropy (Fig. 5B,C); hence, an approach driven by

338     entropy alone would not have divided these spatially separated clusters. Thus, the

339     strength of a residue's association in a sector of highly coevolving residues is associated

340     with more phylogenetic distinction from the rest of the protein than can be explained by

341     entropy alone.

342

### *Phylogenetic similarity highlights overlapping sectors*

344 The core residues of some MreB NC eigenvectors sometimes had high coefficients in

345 multiple eigenvectors (Fig. S5), suggesting that we should consider the union of the

346 sectors as a functional unit. To rationally identify sectors that should be merged, we again

347 exploited phylogenetic similarity by calculating MirrorTree correlation coefficients from

348 comparisons between pairs of sectors (Fig. 6A). MreB NC sectors A-C (Fig. 3) exhibited

349 low sector-sector MirrorTree scores with each other and with random groups (Fig. 6B),

350 as expected since they have low sector-protein MirrorTree scores (Fig. 6B). By contrast,

351 the random groups had MirrorTree scores close to 1 (Fig. 6B). NC sectors were also more

352 phylogenetically distinct from each other than baseline sectors (Fig. 6C). These data

353 suggest that the NC sectors were selected by evolutionary pressures that led to distinct

354 functions.

355

356 Of all sectors that had a MirrorTree $z$-score or a pairwise distance $z$-score<-2 (sectors A-

357 L, Fig. 4C), several pairs had a high sector-sector MirrorTree score. Hierarchical clustering

358 of the sectors based on their sector-sector MirrorTree profiles led to the identification of

359 five obvious "mega-sectors" from the sum of the clustered eigenvectors (Methods), which

360 we denote α, β, γ, δ, and ε (α, β, and γ contain sectors A, B, and C, respectively) (Fig. 6D).

361 The mega-sectors exhibited low sector-sector MirrorTree scores (Fig. 6E), and α, β, and γ

362    had both low sector-protein MirrorTree scores (Fig. 6F, G) and low spatial contiguity $z$-

363    scores (Fig. 6G). The 15-residue version of mega-sector α more compact than the 15-

364    residue version of A (Fig. 6H), and it  contained residues that interact with RodZ (Fig. 6I),

365    an MreB binding partner that modulates MreB filament nucleation (44) and curvature

366    (45). Notably, the regions of the 25-residue version of mega-sector α at the barbed and

367    pointed ends of the MreB subunit interact with each other in a polymerized MreB

368    filament (Fig. 6J), reinforcing the spatial contiguity of the mega-sector. Mega-sector *β* was

369    identical to sector B, surrounding the ATP-binding pocket (Fig. 6K). As with α and A, the

370    15-residue version of mega-sector γ was more compact than the 15-residue version of

371    sector C (Fig. 6L), indicating that clustering based on MirrorTree scores increases the

372    spatial contiguity of sectors.

373

374    ***NC identifies sectors that are not apparent from the full coevolution matrix***

375    To determine whether our findings about the properties of NC sectors applied to other

376    proteins, we performed similar sector calculations for enolase (the metalloenzyme

377    responsible for conversion of 2-phosphoglycerate to phosphoenolpyruvate during

378    glycolysis (46); Fig. 7A-C), the carbohydrate-processing enzyme glucose-6-phosphate

379    dehydrogenase (G6PD (47); Fig. 7D-F), and mitogen-activated protein kinase 1 (MAPK1)

380    (48, 49) (Fig. 7G-K). In each case, NC produced sectors with lower background noise and

381    higher spatial contiguity than baseline sectors.

382

383    Most of the MreB NC eigenvectors had strong signal for either positive or negative

384    coefficients, but not both (Fig. 3A-C). By contrast, one of the large-eigenvalue NC

385    eigenvectors for MAPK1 had groups of residues with both very positive and very

386    negative coefficients (Fig. 7G); these residues were located in distinct regions of the

387    protein (Fig. 7J,K). As validation for splitting the NC eigenvector into two sectors, the

388    sector-sector MirrorTree score (0.44) indicated that they are phylogenetically distinct;

389    moreover, the sector-sector MirrorTree score of the corresponding baseline sectors was

390    higher (0.71). Thus, NC eigenvectors can be interpreted as two phylogenetically distinct

391    sectors based on coefficient signs.

392

393    In addition to improving sector predictions by reducing background variation, we were

394    interested in determining whether NC is able to identify sectors that the full coevolution

395    matrix misses altogether. For the arginine tRNA ligase ArgS (50) and G6PD, the sector

396    with the most negative MirrorTree $z$-score had nearly the lowest spatial contiguity $z$-

397    scores (Fig. 7L,M) and no clear counterpart in any of the baseline eigenvectors (Methods).

398    For ArgS, the NC sector was spatially localized around the arginine binding site (Fig. 7N).

399    For G6PD, the NC sector was adjacent to one of the two NADPs that bind to the protein

400    (Fig. 7O). Thus, the NC correction reveals some sectors that are missed by the baseline

401    method.

402

### *NC sectors are enriched in damaging mutations*

404      To test the functional significance of NC sectors, we sought experimental datasets with

405      quantitative measurements of the consequences of mutations across a protein of interest.

406      Recent studies have pioneered the use of deep mutational scanning to systematically

407      generate and quantify the phenotypic or fitness effects of a large number of individual

408      mutations spanning entire domains or protein (29, 51-53), providing new insights into

409      structure-function relationships. Thus, we asked whether NC sectors were enriched in

410      residues for which mutation altered protein function and/or fitness.

411

412      The Ras superfamily of membrane-associated small G-proteins is highly conserved and

413      controls a broad range of cellular processes (54), has inactive and active states that are

414      regulated by a GTPase-activated protein (55), and has been implicated in cancer (56). A

415      recent deep mutational scanning study engineered plasmids to express mutant versions

416      of human H-Ras as well as the Ras-binding domain of human C-Raf (Raf-RBD) in

417      *Escherichia coli* (57), such that the binding of Ras·GTP to Raf-RBD led to transcription of a

418      chloramphenicol-resistance cassette. Thus, the binding efficacy of the Ras variant was

419      directly correlated with cellular growth rate. The effect of Ras mutations on fitness was

420      quantified by the logarithm of the enrichment of variants in the chloramphenicol-selected

421      versus the starting population, relative to wild-type. The distribution of fitness effects

422    was centered around zero, although there were some positions with mutations that

423    displayed significant functional effects (57).

424

425    To determine whether fitness-altering mutants in H-Ras are enriched at positions

426    identified by coevolution, we identified two high-eigenvalue sectors with obvious

427    corresponding baseline sectors. As in our previous analyses (Fig. 3A-C, 7A,D,G), aside

428    from the highly coevolving residues, the NC sectors had much lower noise than the

429    baseline sectors (Fig. 8A,B). The residues in the two NC sectors were non-overlapping,

430    and in both cases appeared to be concentrated in regions with low minimum relative

431    enrichment (Fig. 8C,D). Across cutoffs that defined sectors of various sizes, we computed

432    the minimum and maximum relative enrichment (representing deactivation and

433    activation, respectively) over all amino acid mutations for each position in the

434    NC/baseline sectors as well as for the residues with the lowest entropy, and compared to

435    the distribution over all residues. As expected, the residues with lowest entropy

436    consistently predicted significantly more negative minimum relative enrichment than

437    random sets of residues (Fig. 8E,F). The mean minimum relative enrichment in NC and

438    baseline versions of sector A was also significantly more negative than random residues,

439    with the NC sector outperforming the baseline sector and achieving similar enrichment

440    values to the lowest-entropy residues (Fig. 8E). NC sector B also exhibited mean

441    minimum relative enrichment significantly lower than random, by contrast to the

442    baseline sector (Fig. 8F). Thus, sectors A and B are more enriched for residues whose

443    mutation has the most potential for reducing fitness using NC versus baseline. The

444    maximum relative enrichment was highly similar for sectors and the protein overall (Fig.

445    S6A,B), suggesting that NC and baseline sectors are enriched for residues with the

446    potential for deactivating rather than activating mutations in the case of H-Ras. Thus, NC

447    sectors separate residues based on the maximum impact of mutations at these positions.

**Discussion**

449

450     Many existing coevolution methods build on correlation or mutual information,

451     sometimes employing ad-hoc corrections to partially remove the effects of entropy and

452     phylogeny. Our NC method harnesses phylogenetic distance between sequences as a

453     novel dimension in the measurement of protein coevolution, in order to increase

454     understanding of the functional relationships between amino acids in a protein. In

455     particular, here we demonstrated that coevolution can occur on multiple phylogenetic

456     timescales within a single protein. While the factors that determine whether pairs of

457     positions coevolve on short or long timescales are unknown, future studies using NC to

458     interrogate the specific biochemical functions of protein sectors may reveal general

459     patterns across diverse proteins. One interpretation of the variable contribution of

460     coevolution across phylogenetic distance within a single protein (Fig. 1C) is that the

461     frequency of mutation for coevolving residues within an NC sector is linked to the

462     timescale of change for the corresponding selective pressure on that sector. For example,

463     a sector that determines protein thermostability would be predicted to coevolve on a

464     timescale commensurate with the frequency of changes in environmental temperature,

465     whether these changes occur over long (e.g. glaciation and interglacial cycles of 100,000

466     years) or shorter (e.g. Atlantic multidecadal oscillations) timescales.

467

468 Importantly, NC and our repurposed MirrorTree methods are complementary to most

469 covariation metrics, and hence can enhance existing bioinformatics tools by defining a

470 phylogenetic dimension of coevolution and allowing focus on functional signal. We

471 anticipate that our approach will enable application of coevolution-based methods across

472 a much broader class of proteins, including those for which the set of sequences is limited

473 in number (Fig. 2) and/or for which the available homologous sequences are biased to a

474 particular segment of the phylogenetic tree (Fig. 1B). In particular, application to the

475 growing database of human exome sequences (58) may improve identification of rare

476 disease-causing mutations. NC may also enhance protein engineering tools by

477 highlighting targets for directed evolution. As we have demonstrated, NC expands our

478 ability to detect functional relationships between residues within proteins and to

479 determine the links between protein evolution and adaptation. In concert with deep

480 mutational scanning and other comprehensive functional screens (59), NC and

481 MirrorTree should provide deeper insight into the specific selective pressures under

482 which proteins have evolved.

483

484 The predominant application of coevolution so far has been structure prediction, from

485 using top DCA-predicted contacts as constraints (4) to employing DCA model

486 parameters as input training features for deep neural networks that seek to predict spatial

487 distances between amino acids (60). Here, we have shown that NC can improve contact

488    prediction by DCA. Moreover, the detection and interpretation of sectors as functional

489    units within proteins has been a growing research focus, particularly with respect to the

490    evolutionary origins of sectors. A recent theoretical study demonstrated that selection

491    acting on a functional property can give rise to a sector (28). Here, we showed that NC

492    better resolves sectors than baseline by reducing background noise (Fig. 3A-C), leading

493    to sectors with higher spatial contiguity (Fig. 3D-L) and lower MirrorTree scores (Fig. 4B).

494    Low MirrorTree scores reveal that residues within sectors have a different evolutionary

495    history from the rest of the protein, due to both entropy-dependent and entropy-

496    independent differences (Fig. 5). MirrorTree scores can further be used to evaluate NC

497    predictions in the absence of a known structure. Motivated by the original design

498    purpose of MirrorTree, we note that scores between sectors of two proteins could be used

499    to identify protein-protein interactions—potentially between hosts and microbes—due to

500    the improved performance of NC when the sampling of sequences is shallow (Fig. 2C).

501

502    Our observation that NC sectors, and moreover their cores, have high spatial contiguity

503    and low MirrorTree scores (Fig. 5J) supports the inferred link between coevolution and

504    spatial contiguity, and suggests that NC can help to guide experiments toward the

505    residues of highest importance for a sector's function (Fig. 8). Beyond the improvements

506    from lowering background signal, NC also predicts sectors that are otherwise difficult to

507    detect (Fig. 7L-O), thus highlighting its value. In addition, some studies have

508    demonstrated other applications such as protein engineering (19) and variant

509    interpretation (14). Improved detection of functional coevolution could even help to

510    refine MSA algorithms, which are ultimately a limiting factor in the detection of

511    coevolution. Our results suggest that the utility of coevolution as a signal for protein

512    science can be substantially improved by NC, opening new windows for broadly

513    understanding (and perhaps ultimately engineering) protein structure-function

514    relationships.

515 **Methods and Materials**

516

517 **MSA construction**

518 MSAs were constructed with BLAST (61) to identify up to 10,000 closest sequences to a

519 reference sequence, using the RefSeq database (62). Sequences were aligned with Clustal

520 Omega (63). Sequences with a Jukes-Cantor distance >1 from the reference sequence were

521 pruned. Redundant sequences and positions with >25% gaps were removed. Any

522 remaining gaps were filled with the closest amino acid from the closest sequence in terms

523 of Jukes-Cantor distance.

524

525 **Calculation of the expected value of inter-clade covariation**

526 For our analyses, we define a pair of sequences to be within the same clade if the

527 phylogenetic distance is below a Jukes-Cantor distance $d$. The phylogenetic distance is

528 measured with respect to the aligned protein sequence (Table S1). We sought to measure

529 the expected value of residue-residue covariation due solely to the comparison of

530 sequences between clades, which we refer to as the inter-clade covariation $C_{s>d}$. Below,

531 we describe and compare measurement of the expected value of inter-clade covariation

532 in Eq. 1 of the main text both by approximation via bootstrapping and analytically.

533

534 *Bootstrapping*

535    In this approximate method, we bootstrap the original MSA: for every position, we

536    replace the amino acid with the identity of the same position from a random sequence in

537    the same clade. For example, in Fig. 1Avi we show two positions in an MSA, colored by

538    their clade membership for a given phylogenetic distance $d$. Note that the first position is

539    never a glutamine in the orange clade and is never a threonine in the white clade.

540    Similarly, the second position is never a serine in the orange clade and is never an

541    arginine in the white clade. The bootstrapped MSAs resample within clades, so as to not

542    change the phylogenetic structure of the MSA at distances $>d$; thus, the first position in

543    the bootstrapped MSAs still does not contain a glutamine, etc. The covariation measured

544    from each of the bootstrapped MSAs is averaged to obtain the matrix expected under the

545    hypothesis that there is no coupling between positions within the same clade. The

546    bootstrapping method can be applied for any coevolution heuristic.

547

548    *Analytical method*

549    To derive an analytical solution in place of bootstrapping the NMI metric, we rephrased

550    our aim as calculating the expected value of covariation between two positions under the

551    assumption that the two positions are independent within a clade.

552

553    Consider the Shannon entropy for position $i$:

554
$$H_i = -\sum_{k=1}^{20} p_{i=k} \ln p_{i=k},$$

555    where $p_{i=k}$ is the probability of finding amino acid $k$ at position $i$. The marginal

556    probabilities of positions $i$ and $j$ taking on a particular value in a bootstrapped MSA do

557    not change on average. However, the joint entropy, which relies on the joint probability,

558    will change:

559
$$H_{ij} = -\sum_{k,l=1}^{20} p_{i=k,j=l} \ln p_{i=k,j=l}.$$

560    We seek an expression for the joint entropy that captures the assumption that positions $i$

561    and $j$ are independent within clades. Since the joint probability of independent variables

562    is the product of the individual probabilities, we are left with calculating the sum of

563    probabilities from each clade $c$, weighted by the number of sequences $n_c$ in each clade:

564
$$p_{i=k,j=l}^{\text{null}} = \left( \sum_c n_c \, p_{i=k}^c p_{j=l}^c \right) \bigg/ \left( \sum_c n_c \right)$$

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

566    $i$.

567

568    A comparison of the bootstrapped and analytical methods for calculating NC for the

569    yeast actin protein is shown in Fig. S2.

570

571    *Estimating the statistical significance of nested coevolution*

572  The expectation value of our nested coevolution background model is described above

573  analytically only for normalized mutual information; other coevolution metrics do not

574  have a known closed form analytical solution, so we rely on bootstrapping to estimate

575  the expected value. Bootstrapping offers the additional advantage of providing an

576  estimate of the statistical significance of the observed raw coevolution signal by

577  measuring what fraction of bootstrapped MSAs achieve equal or greater coevolution

578  values. The accuracy of the significance estimate is limited by the number of bootstrap

579  measurements, since the maximum resolution is the reciprocal of the number of

580  bootstraps performed. Using hundreds of bootstraps, we compared significance

581  estimates with the absolute difference between the total and inter-clade covariation.

582  These values were highly correlated (Spearman's $\rho = 0.95$, Fig. S2B), indicating that either

583  the bootstrapping or analytical method of computing NC provides a surrogate for the

584  significance of the observation.

585

586  **Structural contact prediction**

587  Real structural contacts were determined by calculating the distance between the alpha

588  carbons of every pair of residues in the protein based on a crystal structure (Table S1).

589  All other atoms, including hydrogen atoms, were disregarded. To predict structural

590  contacts, we used mean-field DCA, and the value of the pseudocount is 0.5, and

591  sequences closer than 0.3 Hamming distance are reweighted (4, 22).

592

**Generation of NC sectors**

594     The output of NC is $n_d$ $p$-by-$p$ matrices (Fig. 1C), where $n_d$ is the number of phylogenetic

595     windows and $p$ is the number of amino acids in the protein. These $n_d$ matrices are

596     concatenated to obtain a supermatrix of dimension $pn_d$-by-$p$ (Fig. S7). Principal

597     component analysis using eigenvalue decomposition or singular value decomposition is

598     performed on the super matrix (thus avoiding the need to choose one value of the cutoff

599     distance $d$), with $pn_d$ observations and $p$ features. The eigenvectors are ordered highest to

600     lowest according to their associated eigenvalues. Each eigenvector is of length $p$, where

601     the $i^{th}$ coefficient corresponds to the importance of the $i^{th}$ amino acid in explaining the

602     variation in the direction of the respective eigenvector.

603

604     To extract the specific amino acids that are most responsible for explaining the variation

605     in a particular eigenvector, we identify the positions with the most positive or most

606     negative coefficients and define these groups of residues as two sectors. Sectors that have

607     <4 amino acids are ignored for downstream analysis.

608

609     NC and baseline sectors were paired if the dot product of the corresponding eigenvector

610     was >0.6.

611

612 **Calculating the spatial contiguity of a sector**

613 To quantify spatial contiguity, we calculate the mean distance between the alpha carbon

614 atoms of each pair of residues in the sector in the crystal structure.

615

616 **Adaptation of the MirrorTree algorithm**

617 Mirrortree was originally developed to predict protein-protein interactions based on the

618 similarity of phylogenetic trees (37). In brief, MSAs are calculated using protein

619 sequences from the same list of organisms for two proteins. For each MSA, the matrix of

620 pairwise Jukes-Cantor distances is calculated. The MirrorTree score is the Pearson

621 correlation coefficient of these two distance matrices. A high correlation indicates that the

622 two proteins have similar phylogenies and thus are likely to have experienced similar

623 functional selection. We adapted this method to compare the phylogenetic similarity of

624 protein sectors with the entire protein (Fig. 4A) or other sectors (Fig. 6A). To compute

625 sector-protein and sector-sector MirrorTree scores, filtered MSAs were created focusing

626 on the positions of a given sector.

627

628 Biased sampling of random sectors was accomplished via weighting of residues

629 according to their entropy.

630

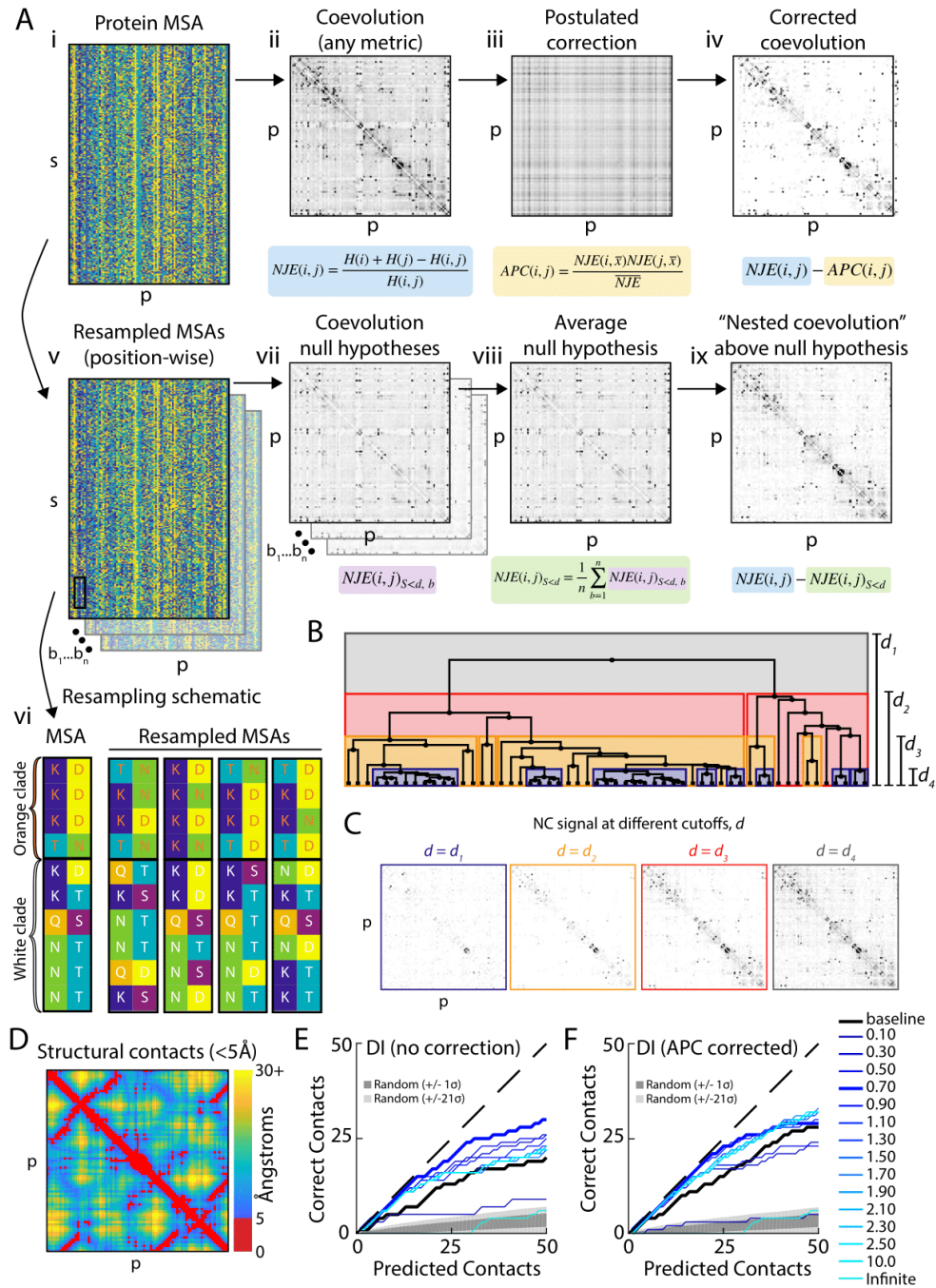631 **Calculation of megasectors**

632    Sets of sectors to be merged into megasectors were determined from hierarchical

633    clustering based on sector-sector MirrorTree scores. Merging was accomplished by

634    adding the corresponding eigenvectors after multiplying each sector by +1 or -1

635    corresponding to whether a positive or negative sector, respectively, was being merged.

636    The summed vector was then analyzed as if it were an eigenvector in order to define

637    megasectors at various size cutoffs.

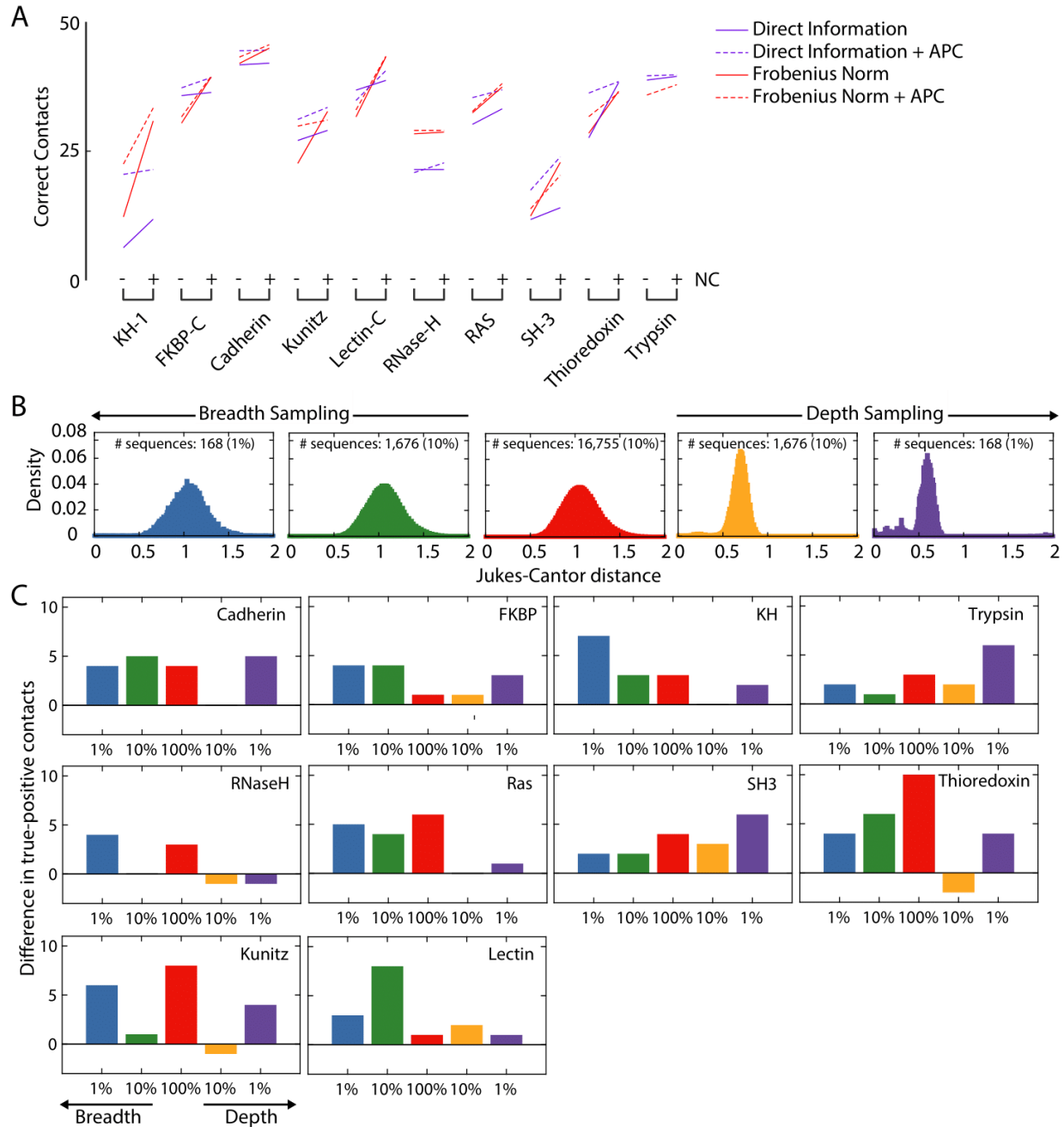## Acknowledgments

645 **Figure Legends**

646



647

648 **Figure 1: NC introduces a phylogenetic dimension to traditional coevolution metrics**

649 **that removes noise and improves structural prediction.**

650    A) Schematic illustrating the NC correction to traditional coevolution algorithms. The

651    MSA (i) is used to generate a covariation matrix (ii) with a particular metric such

652    as normalized joint entropy or direct information. Previous studies have

653    attempted to remove phylogenetic noise using the APC (iii), which results in a

654    corrected coevolution matrix (iv) that has lower levels of off-diagonal signal. For

655    the NC correction, the MSA is resampled multiple times (v) within clades defined

656    by a phylogenetic cutoff $d$ (vi), providing null hypotheses (vii) that are averaged

657    (viii) to correct the covariation matrix (i). The resulting difference (ix) is the NC

658    matrix for a particular cutoff $d$.

659    B) The Jukes-Cantor phylogenetic distance between homologs defines clades

660    (visualized as a tree) within the NC cutoff $d$.

661    C) NC signal at different cutoffs $d$ as illustrated in (B) for the MSA of the KH domain

662    from (B). For small values of $d$, the NC matrix exhibits very little off-diagonal

663    signal, signifying a reduction in noise.

664    D) The structural contact map for KH, highlighting contacts that are in close 3D

665    proximity (<5 Å, red), respectively.

666    E,F) NC with particular cutoffs $d$ improves the prediction of structural contacts

667    relative to DCA, applied to DI without (E) or after correction with APC (F) (black

668    lines). All residues within five positions on the polypeptide sequence were

669         excluded from the analysis. Gray represents the predictions of the baseline NMI
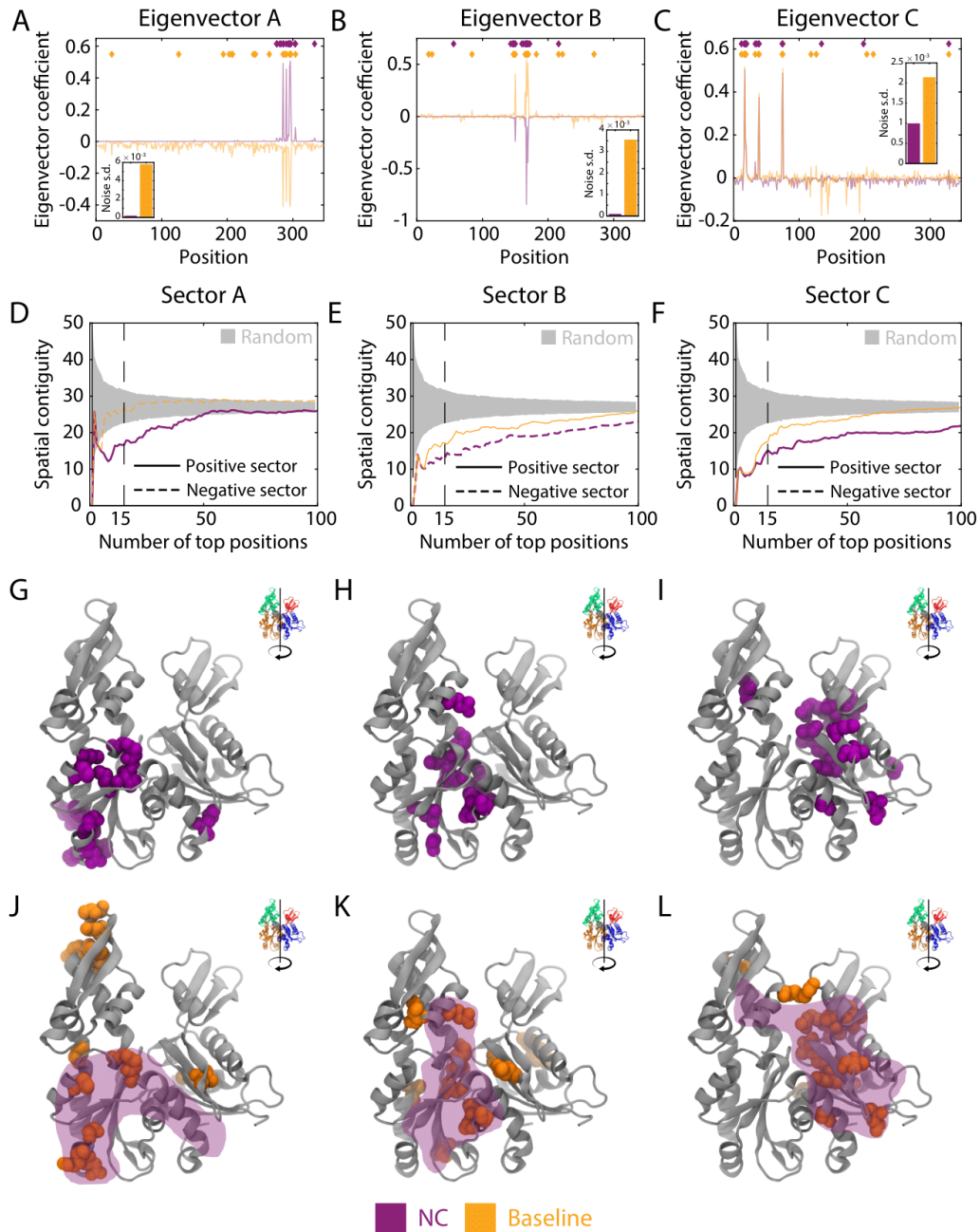
670         metric.

**Figure 2: NC improves predictions of structural contacts across proteins and coevolution methods, and resolves information loss due to subsampling of the set of sequences.**

675      A) NC increased the number of true-positive structural contacts among the first 50

676          predictions for 10 highly conserved proteins predicted by DCA using DI or

677          Frobenius norm, without or with APC.

678      B) MSAs were subsampled across breadth (random sampling) and depth (sorted

679          sampling) of the MSA. Typically, the distribution of Jukes-Cantor distances in the

680          MSA (red) remained essentially unchanged for breadth sampling (green and blue),

681          while it shifted to lower values (as expected) for depth sampling (gold and purple);

682          shown is the KH domain.

683      C) NC generally increased the number of true-positive structural contacts among the

684          first 50 predictions relative to DCA employing DI (without APC) across proteins

685          and both breadth and depth sampling (for DI with APC, see Fig. S3). Small

686          decreases occurred for depth sampling of RNase H, thioredoxin, and Kunitz.

687

**Figure 3: NC eigenvectors for the actin homolog MreB have lower noise and are more spatially contiguous than baseline eigenvectors.**
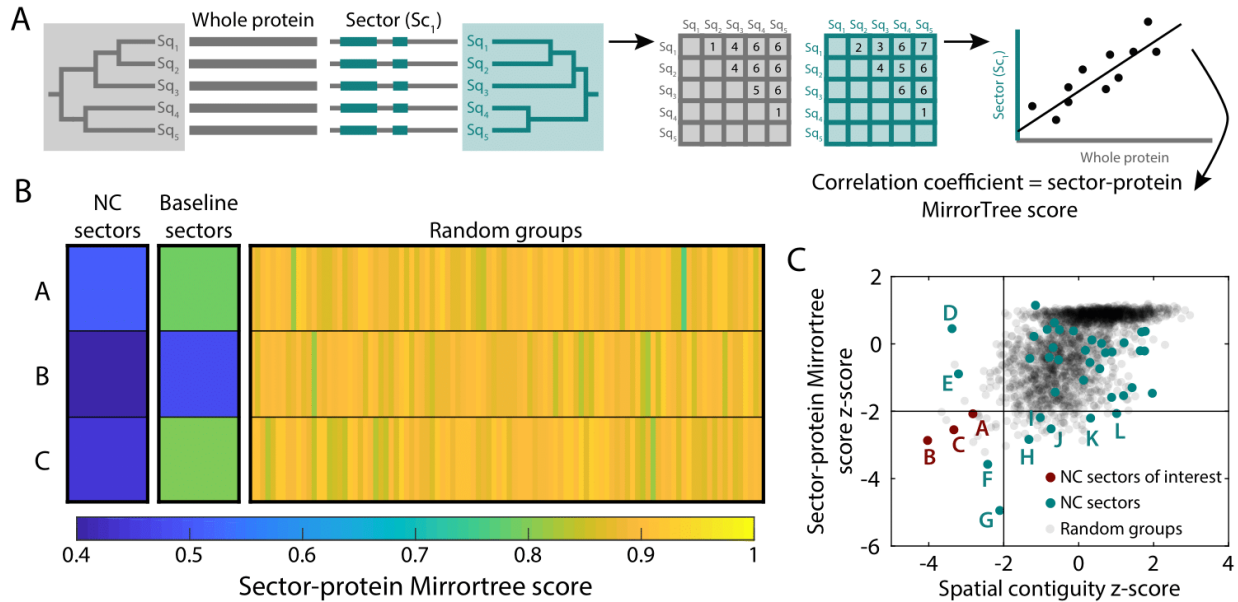
A-C) Three eigenvectors with large eigenvalues were identified and paired between baseline coevolution (NMI with APC) and the NC correction for an MSA

692     containing 9,998 sequences of MreB. Aside from the residues with large

693     coefficients, the NC eigenvectors exhibited lower background noise than the

694     baseline eigenvectors. Insets: standard deviations of the eigenvector coefficients

695     after excluding the highest and lowest 50 values.

696   D-F) NC sectors are more spatially contiguous than the corresponding baseline

697     sectors. Sectors were defined based on a sliding cutoff of the most positive or most

698     negative coefficients of each eigenvector in (A-C). Spatial contiguity was defined

699     as the mean pairwise distance between each residue within a sector.

700   G-L) For the 15-residue versions of the NC and baseline sectors (vertical lines in (D-

701     F)), the NC sectors (G-I) are more compact on the three-dimensional structure than

702     the corresponding baseline sectors (J-L). The shaded purple regions in (J-L)

703     represent the NC sector.

704

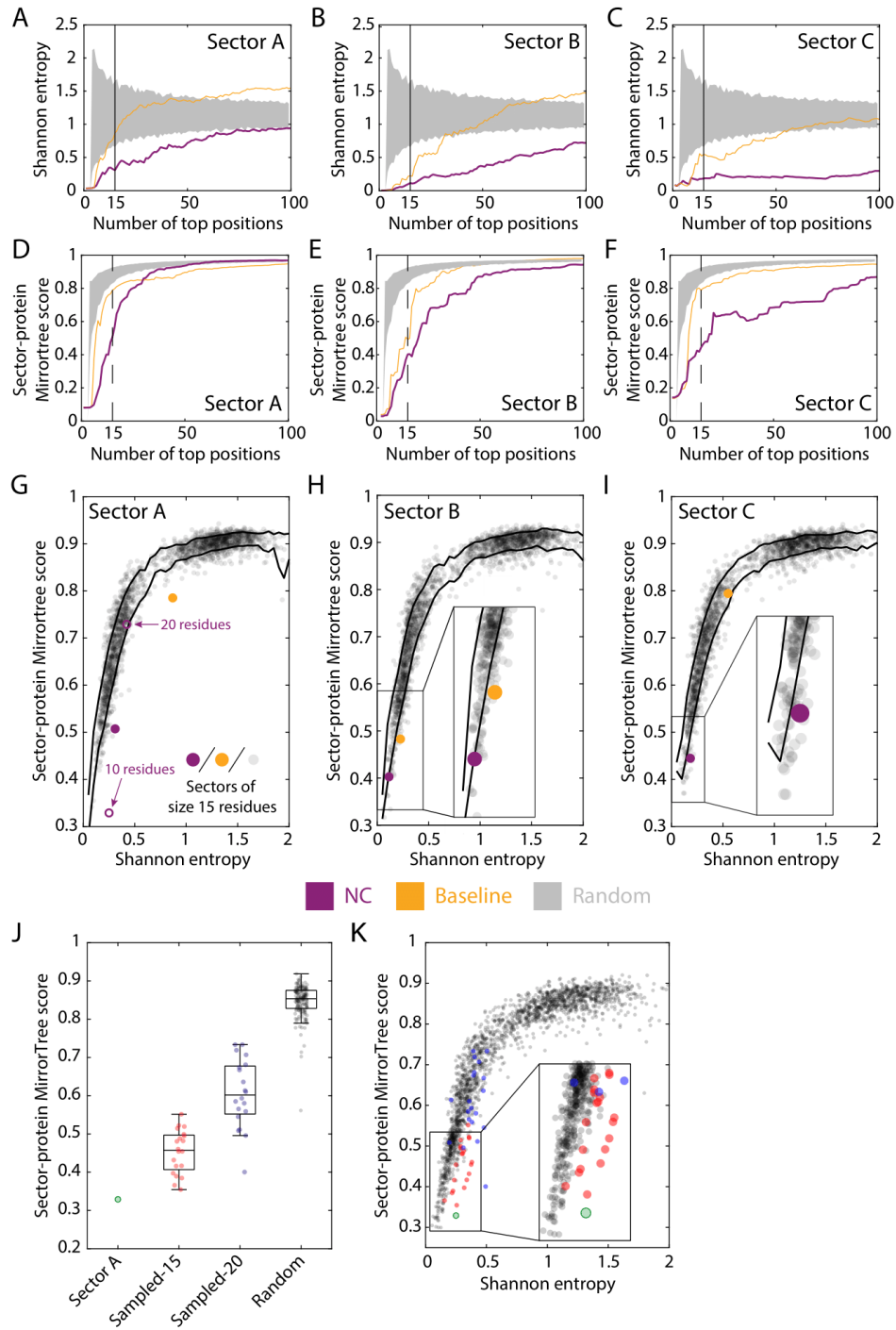**Figure 4: Sectors are phylogenetically distinct from the entire protein.**

A) Repurposing the MirrorTree algorithm (37) 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 $z$-

716     scores of the MirrorTree score and the spatial contiguity were <-2 for sectors A-C.

717     Sectors D-L substantially overlapped sectors A-C, and are considered in Fig. 6.

**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.**
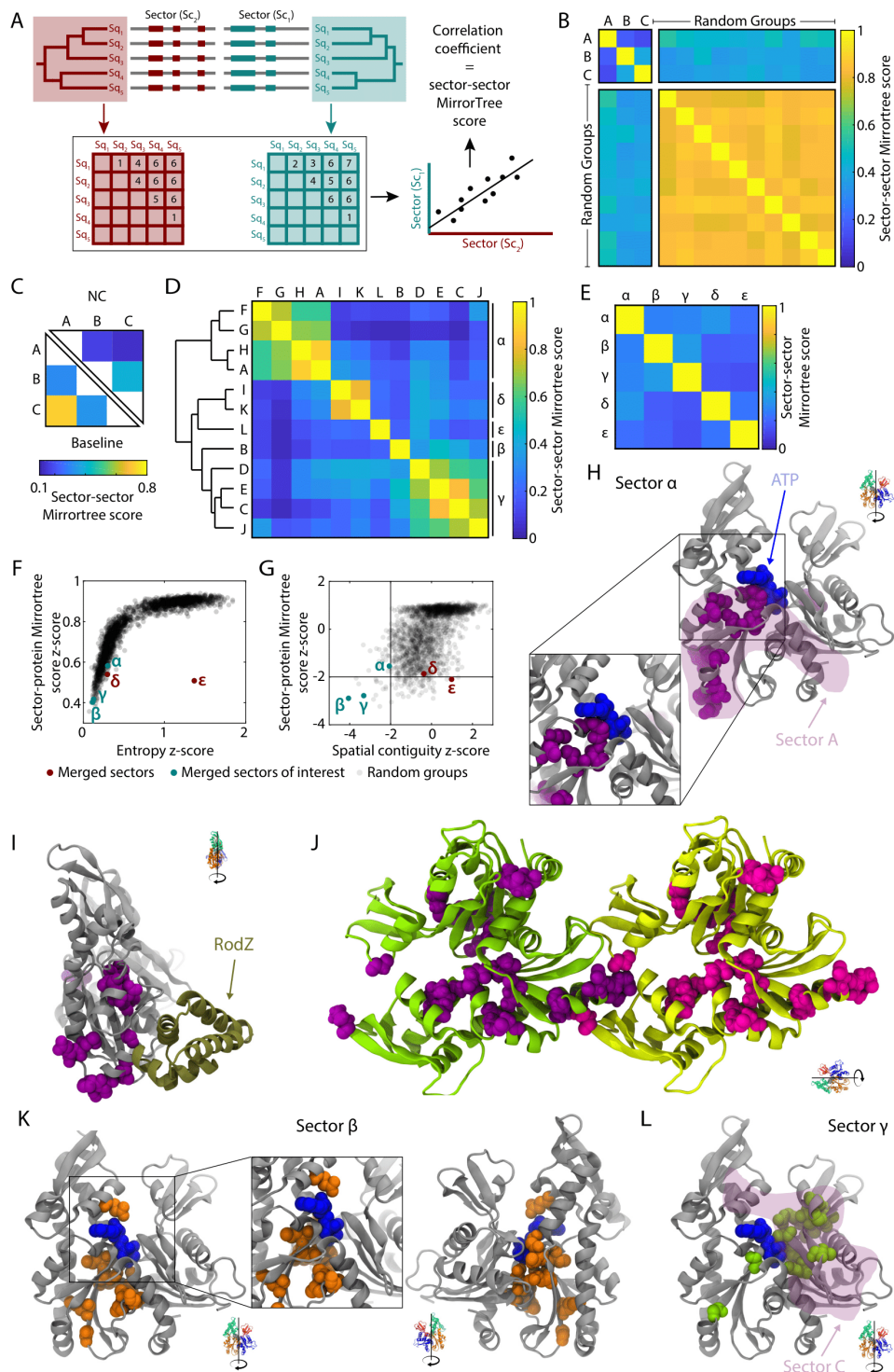
722    A-C) The Shannon entropy of MreB NC sectors A-C (Fig. 3) across size cutoffs is lower

723    than that of the corresponding baseline sectors, indicating that NC selects more

724    conserved residues (albeit entropy is still higher than the cutoff of <0.1 for typically

725    being considered highly conserved). Gray regions represent the entropy of a

726    randomly selected group of residues of the same size.

727    D-F) MirrorTree scores are lower for the NC sectors than for the corresponding

728    baseline sectors. Gray regions represent the MirrorTree scores of a randomly

729    selected group of residues of the same size.

730    G-I) The MirrorTree scores of sectors A-C (filled gold and purple circles) and of

731    random groups of 15 residues (gray). Although MirrorTree score is linked to

732    entropy, NC sectors A and C have MirrorTree scores significantly lower than

733    expected based on entropy alone. In (G), the open purple circles denote the

734    versions of sector A with 10 and 20 residues. Black curves indicate ±1 standard

735    deviation from the mean MirrorTree score for a given entropy.

736    J) The 10-residue version of NC sector A has lower MirrorTree score than sets of 10

737    residues selected from the 15- and 20-residue versions of the same sector, which

738    are lower than those of random groups of 10 residues. The central mark indicates

739    the median, and the bottom and top edges of the box indicate the 25th and 75th

740    percentiles, respectively. The whiskers extend to the most extreme data points not

741    considered outliers.

742    K) The 10-residue version of NC sector A has a lower MirrorTree score than 10-residue

743    subsets of the 15- and 20-residue versions of the same sector with similar entropy.

744    Same data as in (J). Thus, the 10-residue sector represents a "core" of the most
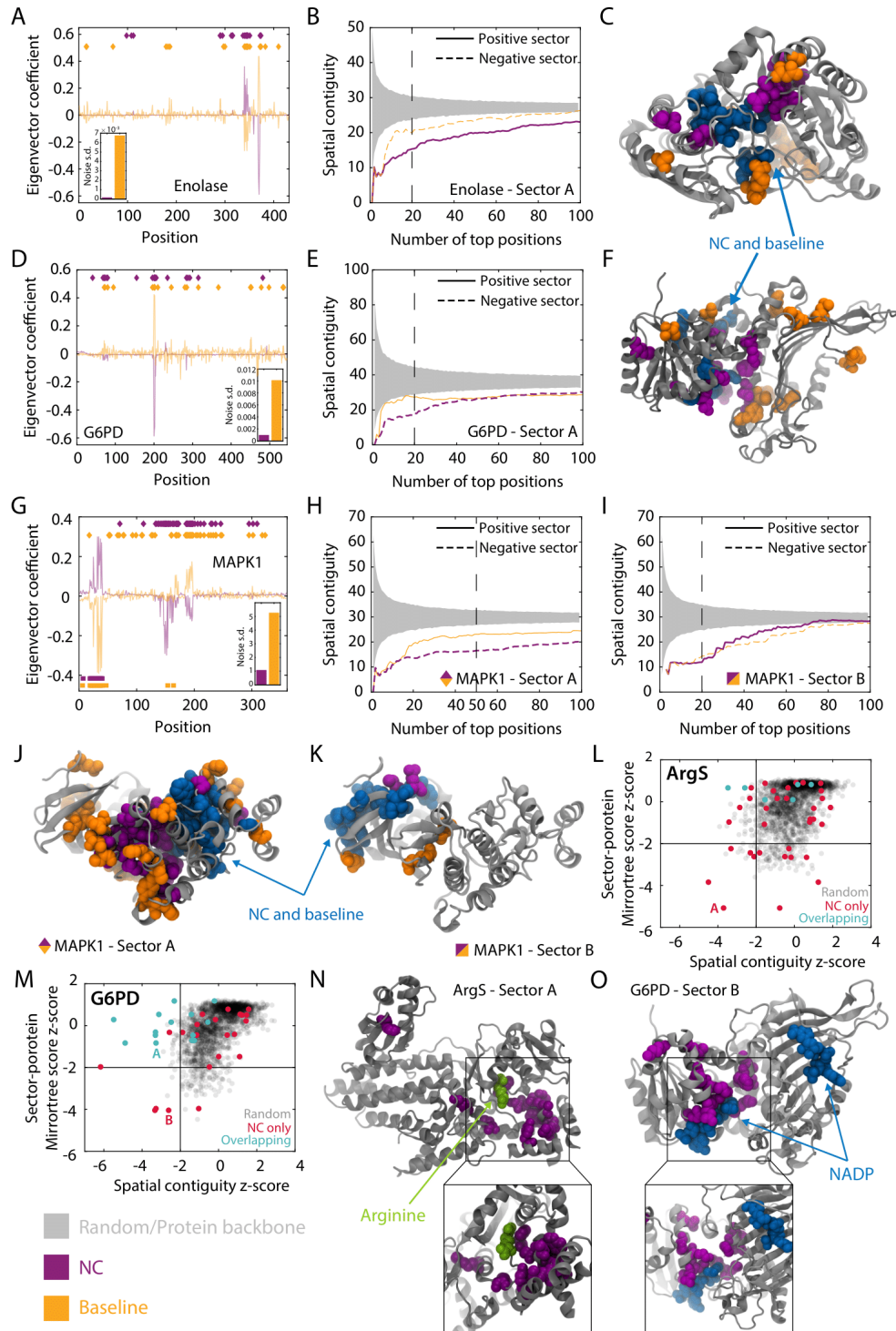
745    highly coevolving residues.

746

**Figure 6: MreB NC sectors are generally phylogenetically distinct, and those with phylogenetic overlap collectively overlap with functionally important regions.**

749    A) Repurposing the MirrorTree algorithm to measure the phylogenetic similarity

750        between sectors.

751    B) MreB NC sectors A, B, and C exhibited low sector-sector MirrorTree scores with

752        each other, but high values with random groups of 15 residues (which also

753        exhibited high MirrorTree scores with each other).

754    C) NC sectors A-C have lower sector-sector MirrorTree scores with each other than

755        baseline sectors A-C with each other, indicating that they are more

756        phylogenetically distinct.

757    D) Hierarchical clustering of MreB NC sectors A-L (Fig. 4C) based on sector-sector

758        MirrorTree profiles suggests five distinct mega-sectors.

759    E-G) The MreB mega-sectors defined by the sum of the clustered eigenvectors

760        exhibited low sector-sector MirrorTree scores with each other (E) as well as low

761        sector-protein MirrorTree scores (F). Mega-sectors $\alpha$, $\beta$, and $\gamma$ (similar to sectors A-

762        C) exhibited high spatial contiguity ($z$-score<-2).

763    H,I) Mega-sector $\alpha$ was more spatially contiguous than sector A (shaded purple

764        region) (H), and contained residues around the interface with MreB's binding

765        partner RodZ (I).

766    J) The 25-residue version of mega-sector $\alpha$ connects the pointed and barbed ends of

767        each subunit in a protofilament.

768    K) Mega-sector $\beta$ (identical to sector B) surrounds the ATP binding pocket.

769     L) Mega-sector $\gamma$ is more spatially contiguous than sector C (shaded purple region).

770

**Figure 7: NC eigenvectors generally improve sector prediction across proteins, and enable identification of sectors that are not detectable using the baseline method.**
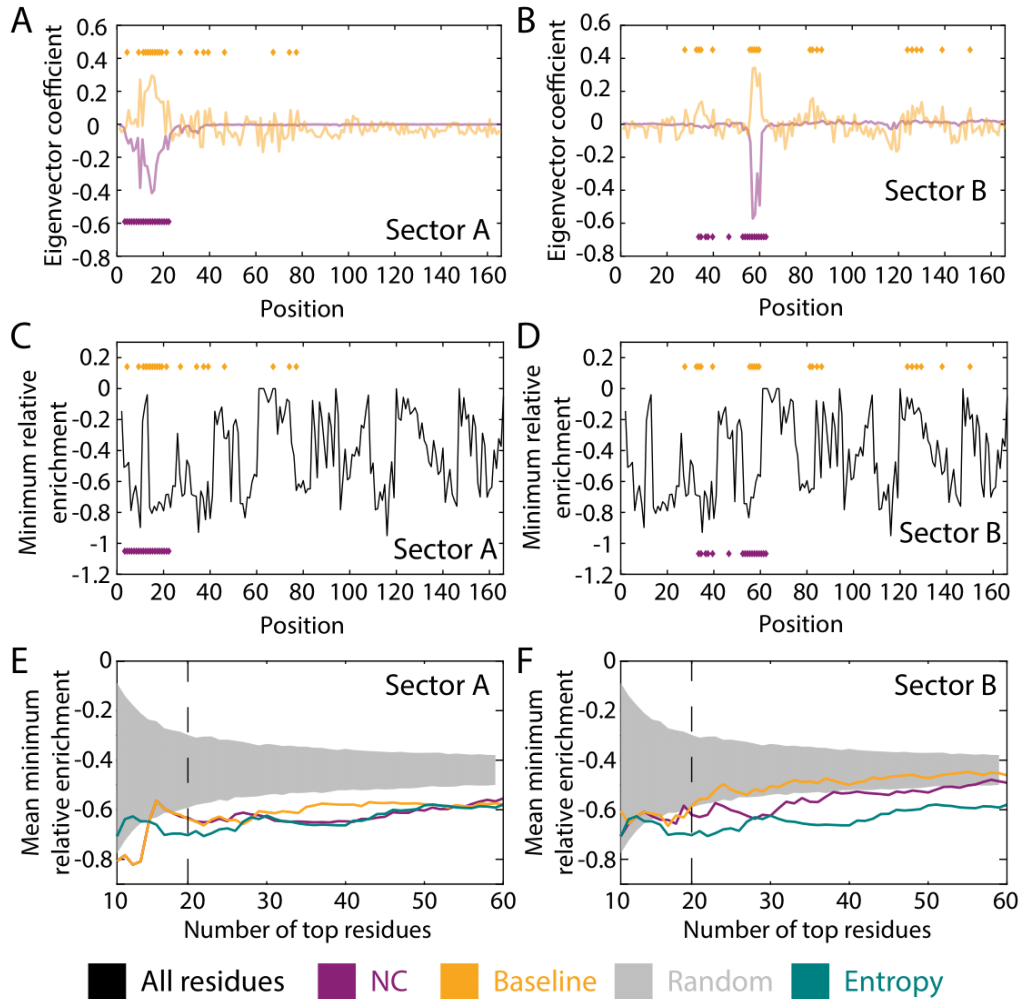
773    A,D,G) NC eigenvectors for enolase (A), G6PD (D), and MAPK1 (G) exhibit lower

774         background noise than the corresponding baseline (NMI with APC)

775         eigenvectors.

776    B,E,H,I) The appropriate NC sectors (positive or negative values of the eigenvector)

777         associated with the eigenvectors in (A,D,G) are more spatially contiguous across

778         size cutoffs than the baseline sectors. Note that the MAPK1 eigenvector was split

779         into a positive sector (H) and a negative sector (I).

780    C,F) The 15-residue versions of the sectors in (B,E) on the crystal structures of enolase

781         (C) and G6PD (F) illustrate the more compact nature of the NC sectors as

782         compared with the baseline sectors.

783    J,K) The 50- and 20-residue versions of the NC sectors in (H,I) are more spatially

784         compact on the structure than the corresponding baseline sectors, and occupy

785         distinct parts of the protein.

786    L-O) For ArgS (L) and G6PD (M), certain high-eigenvalue NC sectors had no

787         obvious baseline counterpart. These NC sectors had low MirrorTree and spatial

788         contiguity $z$-scores (L,M), and 15-residue versions occupied spatially compact

789         regions around ligands (arginine in (N), NADP in (O)) on the structure (N,O).

790         Thus, NC enables the detection of sectors that are otherwise hidden.

791

**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-RBD (57). 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.

800      E,F)Across most sector size cutoffs, the mean minimum relative enrichment was

801      significantly lower than random (gray) for NC sectors A and B and comparable

802      that of the residues with the lowest entropy (teal). NC sectors also outperformed

803      their baseline counterparts.

## References

1. M. Hollstein, B. Sidransky, B. Vogelstein, C. C. Harris, P53 Mutations in Human Cancers. *Science* **253**, 49-53 (1991).

2. M. J. Zvelebil, G. J. Barton, W. R. Taylor, M. J. E. Sternberg, Prediction of protein secondary structure and active sites using the alignment of homologous sequences. *Journal of Molecular Biology* **195**, 957-961 (1987).

3. R. S. Dwyer, D. P. Ricci, L. J. Colwell, T. J. Silhavy, N. S. Wingreen, Predicting functionally informative mutations in Escherichia coli BamA using evolutionary covariance analysis. *Genetics* **195**, 443-455 (2013).

4. D. S. Marks, L. J. Colwell, R. Sheridan, T. A. Hopf, A. Pagnani, R. Zecchina, C. Sander, Protein 3D structure computed from evolutionary sequence variation. *PLoS ONE* **6** (2011).

5. F. Morcos, B. Jana, T. Hwa, J. N. Onuchic, Coevolutionary signals across protein lineages help capture multiple protein conformations. *Proceedings of the National Academy of Sciences of the United States of America* **110**, 20533-20538 (2013).

6. K. a. Reynolds, R. N. McLaughlin, R. Ranganathan, Hot spots for allosteric regulation on protein surfaces. *Cell* **147**, 1564-1575 (2011).

868   7.    M. Weigt, R. A. White, H. Szurmant, J. A. Hoch, T. Hwa, Identification of direct

869         residue contacts in protein-protein interaction by message passing. *Proc Natl*

870         *Acad Sci U S A* **106**, 67-72 (2009).

871   8.    A. F. Bitbol, Inferring interaction partners from protein sequences using mutual

872         information. *PLoS Comput Biol* **14**, e1006401 (2018).

873   9.    A. F. Bitbol, R. S. Dwyer, L. J. Colwell, N. S. Wingreen, Inferring interaction

874         partners from protein sequences. *Proc Natl Acad Sci U S A* **113**, 12180-12185

875         (2016).

876   10.   L. Burger, E. van Nimwegen, Accurate prediction of protein-protein interactions

877         from sequence alignments using a Bayesian method. *Mol Syst Biol* **4**, 165 (2008).

878   11.   T. Gueudre, C. Baldassi, M. Zamparo, M. Weigt, A. Pagnani, Simultaneous

879         identification of specifically interacting paralogs and interprotein contacts by

880         direct coupling analysis. *Proc Natl Acad Sci U S A* **113**, 12186-12191 (2016).

881   12.   S. Ovchinnikov, H. Kamisetty, D. Baker, Robust and accurate prediction of

882         residue-residue interactions across protein interfaces using evolutionary

883         information. *eLife*  (2014).

884   13.   O. Rivoire, Parsimonious evolutionary scenario for the origin of allostery and

885         coevolution patterns in proteins. *Phys Rev E* **100**, 032411 (2019).

886    14.    T. A. Hopf, J. B. Ingraham, F. J. Poelwijk, C. P. Scharfe, M. Springer, C. Sander, D.

887            S. Marks, Mutation effects predicted from sequence co-variation. *Nat Biotechnol*

888            **35**, 128-135 (2017).

889    15.    D. Altschuh, T. Vernet, D. Moras, K. Nagai, Coordinated amino acid changes in

890            homologous protein families. *Protein Engineering* **2**, 193-199 (1988).

891    16.    W. Atchley, K. Wollenberg, W. Fitch, W. Terhalle, A. Dress, Correlations among

892            amino acid sites in bHLH protein domains: an information theoretic analysis.

893            164-178 (2000).

894    17.    U. Göbel, C. Sander, R. Schneider, a. Valencia, Correlated mutations and residue

895            contacts in proteins. *Proteins* **18**, 309-317 (1994).

896    18.    N. Halabi, O. Rivoire, S. Leibler, R. Ranganathan, Protein sectors: evolutionary

897            units of three-dimensional structure. *Cell* **138**, 774-786 (2009).

898    19.    J. M. Skerker, B. S. Perchuk, A. Siryaporn, E. a. Lubin, O. Ashenberg, M. Goulian,

899            M. T. Laub, Rewiring the Specificity of Two-Component Signal Transduction

900            Systems. *Cell* **133**, 1043-1054 (2008).

901    20.    S. D. Dunn, L. M. Wahl, G. B. Gloor, Mutual information without the influence of

902            phylogeny or entropy dramatically improves residue contact prediction.

903            *Bioinformatics (Oxford, England)* **24**, 333-340 (2008).

904    21.    M. Socolich, S. W. Lockless, W. P. Russ, H. Lee, K. H. Gardner, R. Ranganathan,

905         Evolutionary information for specifying a protein fold. *Nature* **437**, 512-518

906         (2005).

907    22.    F. Morcos, A. Pagnani, B. Lunt, A. Bertolino, D. S. Marks, C. Sander, R. Zecchina,

908         J. N. Onuchic, T. Hwa, M. Weigt, Direct-coupling analysis of residue coevolution

909         captures native contacts across many protein families. *Proc Natl Acad Sci U S A*

910         **108**, E1293-1301 (2011).

911    23.    L. C. Martin, G. B. Gloor, S. D. Dunn, L. M. Wahl, Using information theory to

912         search for co-evolving residues in proteins. *Bioinformatics* **21**, 4116-4124 (2005).

913    24.    K. R. Wollenberg, W. R. Atchley, Separation of phylogenetic and functional

914         associations in biological sequences by using the parametric bootstrap. *Proc Natl*

915         *Acad Sci U S A* **97**, 3288-3291 (2000).

916    25.    M. F. Garcia-Mayoral, D. Hollingworth, L. Masino, I. Diaz-Moreno, G. Kelly, R.

917         Gherzi, C. F. Chou, C. Y. Chen, A. Ramos, The structure of the C-terminal KH

918         domains of KSRP reveals a noncanonical motif important for mRNA

919         degradation. *Structure* **15**, 485-498 (2007).

920    26.    G. H. Golub, C. F. Van Loan, *Matrix computations*, Johns Hopkins studies in the

921         mathematical sciences (Johns Hopkins University Press, Baltimore, ed. 3rd, 1996),

922         pp. xxvii, 694 p.

27.  M. Ekeberg, C. Lovkvist, Y. Lan, M. Weigt, E. Aurell, Improved contact prediction in proteins: using pseudolikelihoods to infer Potts models. *Phys Rev E Stat Nonlin Soft Matter Phys* **87**, 012707 (2013).

28.  S. W. Wang, A. F. Bitbol, N. S. Wingreen, Revealing evolutionary constraints on proteins through sequence analysis. *PLoS Comput Biol* **15**, e1007010 (2019).

29.  R. N. McLaughlin, Jr., F. J. Poelwijk, A. Raman, W. S. Gosal, R. Ranganathan, The spatial architecture of protein function and adaptation. *Nature* **491**, 138-142 (2012).

30.  M. Novinec, M. Korenc, A. Caflisch, R. Ranganathan, B. Lenarcic, A. Baici, A novel allosteric mechanism in the cysteine peptidase cathepsin K discovered by computational methods. *Nat Commun* **5**, 3287 (2014).

31.  O. Rivoire, K. A. Reynolds, R. Ranganathan, Evolution-Based Functional Decomposition of Proteins. *PLoS Comput Biol* **12**, e1004817 (2016).

32.  R. G. Smock, O. Rivoire, W. P. Russ, J. F. Swain, S. Leibler, R. Ranganathan, L. M. Gierasch, An interdomain sector mediating allostery in Hsp70 molecular chaperones. *Mol Syst Biol* **6**, 414 (2010).

33.  H. Shi, B. P. Bratton, Z. Gitai, K. C. Huang, How to Build a Bacterial Cell: MreB as the Foreman of E. coli Construction. *Cell* **172**, 1294-1305 (2018).

941   34.   T. Izore, R. Duman, D. Kureisaite-Ciziene, J. Lowe, Crenactin from Pyrobaculum

942          calidifontis is closely related to actin in structure and forms steep helical

943          filaments. *FEBS Lett* **588**, 776-782 (2014).

944   35.   F. van den Ent, L. A. Amos, J. Lowe, Prokaryotic origin of the actin cytoskeleton.

945          *Nature* **413**, 39-44 (2001).

946   36.   F. van den Ent, L. Amos, J. Lowe, Bacterial ancestry of actin and tubulin. *Curr*

947          *Opin Microbiol* **4**, 634-638 (2001).

948   37.   R. A. Craig, L. Liao, Phylogenetic tree information aids supervised learning for

949          predicting protein-protein interaction based on distance matrices. *BMC*

950          *Bioinformatics* **8**, 6 (2007).

951   38.   F. Pazos, A. Valencia, Similarity of phylogenetic trees as indicator of protein-

952          protein interaction. *Protein Eng* **14**, 609-614 (2001).

953   39.   C. L. Araya, C. Cenik, J. A. Reuter, G. Kiss, V. S. Pande, M. P. Snyder, W. J.

954          Greenleaf, Identification of significantly mutated regions across cancer types

955          highlights a rich landscape of functional molecular alterations. *Nat Genet* **48**, 117-

956          125 (2016).

957   40.   Z. Hu, B. Ma, H. Wolfson, R. Nussinov, Conservation of polar residues as hot

958          spots at protein interfaces. *Proteins* **39**, 331-342 (2000).

959   41.   O. B. Ptitsyn, Protein folding and protein evolution: common folding nucleus in

960          different subfamilies of c-type cytochromes? *J Mol Biol* **278**, 655-666 (1998).

961   42.   T. Teşileanu, L. J. Colwell, S. Leibler, Protein Sectors: Statistical Coupling

962         Analysis versus Conservation. *PLOS Computational Biology* **11**, e1004091 (2015).

963   43.   I. Anishchenko, S. Ovchinnikov, H. Kamisetty, D. Baker, Origins of coevolution

964         between residues distant in protein 3D structures. *Proc Natl Acad Sci U S A* **114**,

965         9122-9127 (2017).

966   44.   B. P. Bratton, J. W. Shaevitz, Z. Gitai, R. M. Morgenstein, MreB polymers and

967         curvature localization are enhanced by RodZ and predict E. coli's cylindrical

968         uniformity. *Nat Commun* **9**, 2797 (2018).

969   45.   A. Colavin, H. Shi, K. C. Huang, RodZ modulates geometric localization of the

970         bacterial actin MreB to regulate cell shape. *Nat Commun* **9**, 1280 (2018).

971   46.   T. G. Spring, F. Wold, The purification and characterization of Escherichia coli

972         enolase. *J Biol Chem* **246**, 6797-6802 (1971).

973   47.   D. N. Wright, W. R. Lockhart, Effects of Growth Rate and Limiting Substrate on

974         Glucose Metabolism in Escherichia Coli. *J Bacteriol* **89**, 1082-1085 (1965).

975   48.   S. L. Pelech, J. S. Sanghera, M. Daya-Makin, Protein kinase cascades in meiotic

976         and mitotic cell cycle control. *Biochem Cell Biol* **68**, 1297-1330 (1990).

977   49.   T. W. Sturgill, J. Wu, Recent progress in characterization of protein kinase

978         cascades for phosphorylation of ribosomal protein S6. *Biochim Biophys Acta* **1092**,

979         350-357 (1991).

980   50.   I. N. Hirshfield, H. P. Bloemers, The biochemical characterization of two mutant

981         arginyl transfer ribonucleic acid synthetases from Escherichia coli K-12. *J Biol*

982         *Chem* **244**, 2911-2916 (1969).

983   51.   S. L. Dove, J. K. Joung, A. Hochschild, Activation of prokaryotic transcription

984         through arbitrary protein-protein contacts. *Nature* **386**, 627-630 (1997).

985   52.   J. K. Joung, E. I. Ramm, C. O. Pabo, A bacterial two-hybrid selection system for

986         studying protein-DNA and protein-protein interactions. *Proc Natl Acad Sci U S A*

987         **97**, 7382-7387 (2000).

988   53.   W. A. Lim, R. T. Sauer, Alternative packing arrangements in the hydrophobic

989         core of lambda repressor. *Nature* **339**, 31-36 (1989).

990   54.   C. W. Johnson, D. Reid, J. A. Parker, S. Salter, R. Knihtila, P. Kuzmic, C. Mattos,

991         The small GTPases K-Ras, N-Ras, and H-Ras have distinct biochemical properties

992         determined by allosteric effects. *J Biol Chem* **292**, 12981-12993 (2017).

993   55.   C. Wellbrock, M. Karasarides, R. Marais, The RAF proteins take centre stage. *Nat*

994         *Rev Mol Cell Biol* **5**, 875-885 (2004).

995   56.   I. A. Prior, P. D. Lewis, C. Mattos, A comprehensive survey of Ras mutations in

996         cancer. *Cancer Res* **72**, 2457-2467 (2012).

997   57.   P. Bandaru, N. H. Shah, M. Bhattacharyya, J. P. Barton, Y. Kondo, J. C. Cofsky, C.

998         L. Gee, A. K. Chakraborty, T. Kortemme, R. Ranganathan, J. Kuriyan,

Deconstruction of the Ras switching cycle through saturation mutagenesis. *Elife* **6** (2017).

58. R. Do *et al.*, Exome sequencing identifies rare LDLR and APOA5 alleles conferring risk for myocardial infarction. *Nature* **518**, 102-106 (2015).

59. H. Q. Nguyen, J. Roy, B. Harink, N. P. Damle, N. R. Latorraca, B. C. Baxter, K. Brower, S. A. Longwell, T. Kortemme, K. S. Thorn, M. S. Cyert, P. M. Fordyce, Quantitative mapping of protein-peptide affinity landscapes using spectrally encoded beads. *Elife* **8** (2019).

60. A. W. Senior, R. Evans, J. Jumper, J. Kirkpatrick, L. Sifre, T. Green, C. Qin, A. Zidek, A. W. R. Nelson, A. Bridgland, H. Penedones, S. Petersen, K. Simonyan, S. Crossan, P. Kohli, D. T. Jones, D. Silver, K. Kavukcuoglu, D. Hassabis, Improved protein structure prediction using potentials from deep learning. *Nature* **577**, 706-710 (2020).

61. T. Madden, The BLAST Sequence Analysis Tool. 2002 Oct 9 [Updated 2003 Aug 13]. *The NCBI Handbook [Internet]. Bethesda (MD): National Center for Biotechnology Information (US)* (2002).

62. T. Tatusova, S. Ciufo, B. Fedorov, K. O'Neill, I. Tolstoy, RefSeq microbial genomes database: new representation and annotation strategy. *Nucleic Acids Res* **42**, D553-559 (2014).

1018    63.    F. Sievers, A. Wilm, D. Dineen, T. J. Gibson, K. Karplus, W. Li, R. Lopez, H.

1019           McWilliam, M. Remmert, J. Soding, J. D. Thompson, D. G. Higgins, Fast, scalable

1020           generation of high-quality protein multiple sequence alignments using Clustal

1021           Omega. *Mol Syst Biol* **7**, 539 (2011).

1022