



# **From principal component to direct coupling analysis of coevolution in proteins: Low-eigenvalue modes are needed for structure prediction**

Simona Cocco, Remi Monasson, Martin Weigt

## **► To cite this version:**

Simona Cocco, Remi Monasson, Martin Weigt. From principal component to direct coupling analysis of coevolution in proteins: Low-eigenvalue modes are needed for structure prediction. PLoS Computational Biology, 2013, 9 (8), pp.e1003176. <10.1371/journal.pcbi.1003176>. <hal-00764377v2>

**HAL Id: hal-00764377**

**<https://hal.science/hal-00764377v2>**

Submitted on 27 Aug 2013

**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.



HAL Authorization

# From principal component to direct coupling analysis of coevolution in proteins: Low-eigenvalue modes are needed for structure prediction

Simona Cocco<sup>1</sup>, Remi Monasson<sup>2</sup>, Martin Weigt<sup>3,4,\*</sup>

**1** Laboratoire de Physique Statistique de l'Ecole Normale Supérieure - UMR 8550, associé au CNRS et à l'Université Pierre et Marie Curie, 24 rue Lhomond, 75005 Paris, France

**2** Laboratoire de Physique Théorique de l'Ecole Normale Supérieure - UMR 8549, associé au CNRS et à l'Université Pierre et Marie Curie, 24 rue Lhomond, 75005 Paris, France

**3** Université Pierre et Marie Curie, UMR 7238 - Laboratoire de Génomique des Microorganismes, 15 rue de l'Ecole de Médecine, 75006 Paris, France

**4** Human Genetics Foundation, Via Nizza 52, 10126 Torino, Italy

\* E-mail: martin.weigt@upmc.fr

## Abstract

Various approaches have explored the covariation of residues in multiple-sequence alignments of homologous proteins to extract functional and structural information. Among those are principal component analysis (PCA), which identifies the most correlated groups of residues, and direct coupling analysis (DCA), a global inference method based on the maximum entropy principle, which aims at predicting residue-residue contacts. In this paper, inspired by the statistical physics of disordered systems, we introduce the Hopfield-Potts model to naturally interpolate between these two approaches. The Hopfield-Potts model allows us to identify relevant 'patterns' of residues from the knowledge of the eigenmodes and eigenvalues of the residue-residue correlation matrix. We show how the computation of such statistical patterns makes it possible to accurately predict residue-residue contacts with a much smaller number of parameters than DCA. This dimensional reduction allows us to avoid overfitting and to extract contact information from multiple-sequence alignments of reduced size. In addition, we show that low-eigenvalue correlation modes, discarded by PCA, are important to recover structural information: the corresponding patterns are highly localized, that is, they are concentrated in few sites, which we find to be in close contact in the three-dimensional protein fold.

## Author Summary

Extracting functional and structural information about protein families from the covariation of residues in multiple sequence alignments is an important challenge in computational biology. Here we propose a statistical-physics inspired framework to analyze those covariations, which naturally unifies existing methods in the literature. Our approach allows us to identify statistically relevant 'patterns' of residues, specific to a protein family. We show that many patterns correspond to a small number of sites on the protein sequence, in close contact on the 3D fold. Hence, those patterns allow us to make accurate predictions about the contact map from sequence data only. Further more, we show that the dimensional reduction, which is achieved by considering only the statistically most significant patterns, avoids overfitting in small sequence alignments, and improves our capacity of extracting residue contacts in this case.

## Introduction

Thanks to the constant progresses in DNA sequencing techniques, by now more than 4,400 full genomes are sequenced [1], resulting in more than  $3.6 \cdot 10^7$  known protein sequences [2], which are classified into more than 14,000 protein domain families [3], many of them containing in the range of  $10^3 - 10^5$  homologous

(*i.e.* evolutionarily related) amino-acid sequences. These huge numbers are contrasted by only about 92,000 experimentally resolved X-ray or NMR structures [4], many of them describing the same proteins. It is therefore tempting to use sequence data alone to extract information about the functional and the structural constraints acting on the evolution of those proteins. Analysis of single-residue conservation offers a first hint about those constraints: Highly conserved positions (easily detectable in multiple sequence alignments corresponding to one protein family) identify residues whose mutations are likely to disrupt the protein function, *e.g.* by the loss of its enzymatic properties. However, not all constraints result in strong single-site conservation. As is well-known, compensatory mutations can happen and preserve the integrity of a protein even if single site mutations have deleterious effects [5, 6]. A natural idea is therefore to analyze covariations between residues, that is, whether their variations across sequences are correlated or not [7]. In this context, one introduces a matrix  $\Gamma_{ij}(a, b)$  of residue-residue correlations expressing how much the presence of amino-acid 'a' in position 'i' on the protein is correlated across the sequence data with the presence of another amino-acid 'b' in another position 'j'. Extracting information from this matrix has been the subject of numerous studies over the past two decades, see *e.g.* [5, 6, 8–21] and [7] for a recent up-to-date review of the field. In difference to these correlation-based approaches, Yeang *et al.* [22], proposed a simple evolutionary model which measures coevolution in terms of deviation from independent-site evolution. However, a full dynamical model for residue coevolution is still outstanding.

The direct use of correlations for discovering structural constraints such as residue-residue contacts in a protein fold has, unfortunately, remained of limited accuracy [5, 6, 9, 11, 13, 16]. More sophisticated approaches to exploit the information included in  $\Gamma$  are based on a *Maximum Entropy* (MaxEnt) [23, 24] modeling. The underlying idea is to look for the simplest statistical model of protein sequences capable of reproducing empirically observed correlations. MaxEnt has been used to analyze many types of biological data, ranging from multi-electrode recording of neural activities [25, 26], gene concentrations in genetic networks [27], bird flocking [28] etc. MaxEnt to model covariation in protein sequences was first proposed in a purely theoretical setting by Lapedes *et al.* [29], and applied to protein sequences in an unpublished preprint by Lapedes *et al.* [12]. It was used – even if not explicitly stated – by Ranganathan and coworkers to generate random protein sequences through Monte Carlo simulations, as a part of an approach called Statistical Coupling Analysis (SCA) [15]. Remarkably, many of those artificial proteins folded into a native-like state, demonstrating that MaxEnt modeling was able to statistically capture essential features of the protein family. Recently, one of us proposed, in a series of collaborations, two analytical approaches based on mean-field type approximations of statistical physics, called *Direct Coupling Analysis* (DCA), to efficiently compute and exploit this MaxEnt distribution ([17] uses message passing, [19] a computationally more efficient naive mean-field approximation), related approaches developed partially in parallel are [18, 20, 21]. Informally speaking, DCA allows for disentangling direct contributions to correlations (resulting from native contacts) from indirect contributions (mediated through chains of native contacts). Hence, DCA offers a much more accurate image of the contact map than  $\Gamma$  itself. The full potential of maximum-entropy modeling for accurate structural prediction was first recognized in [30] (quaternary structure prediction) and in [31] (tertiary structure prediction), and further applied by [32–38]. It became obvious that the extracted information is sufficient to predict folds of relatively long proteins and transmembrane domains. In [36] it was used to rationally design mutagenesis experiments to repair a non-functional hybrid protein, and thus to confirm the predicted structure.

Despite its success, MaxEnt modeling raises several concerns. The number of 'direct coupling' parameters necessary to define the MaxEnt model over the set of protein sequences, is of the order of  $L^2(q-1)^2$ . Here,  $L$  is the protein length, and  $q = 21$  is the number of amino acids (including the gap). So, for realistic protein lengths of  $L = 50 - 500$ , we end up with  $10^6 - 10^8$  parameters, which have to be inferred from alignments of  $10^3 - 10^5$  proteins. Overfitting the sequence data is therefore a major risk.

Another mathematically simpler way to extract information from the correlation matrix  $\Gamma$  is Principal Component Analysis (PCA) [39]. PCA looks for the eigenmodes of  $\Gamma$  associated to the largest eigenvalues. These modes are the ones contributing most to the covariation in the protein family. Combined with

clustering approaches, PCA was applied to identify functional residues in [8]. More recently PCA was applied to the SCA correlation matrix, a variant of the matrix  $\Gamma$  expressing correlations between sites only (and not explicitly the amino-acids they carry) and allowed for identifying groups of correlated (coevolving) residues – termed sectors – each controlling a specific function [40]. A fundamental issue with PCA is the determination of the number of relevant eigenmodes. This is usually done by comparing the spectrum of  $\Gamma$  with a null model, the Marcenko-Pastur (MP) distribution, describing the spectral properties of the sample covariance matrix of a set of independent variables [41]. Eigenvalues larger than the top edge of the MP distribution cannot be explained from sampling noise and are selected, while lower eigenvalues – inside the bulk of the MP spectrum, or even lower – are rejected.

In this article we show that there exists a deep connection between DCA and PCA. To do so we consider the Hopfield-Potts model, an extension of the Hopfield model introduced three decades ago in computational neuroscience [42], to the case of variables taking  $q > 2$  values. The Hopfield-Potts model is based on the concept of patterns, that is, of special directions in sequence space. These patterns show some similarities with sequence motifs or position-specific scoring matrices, but instead of encoding independent-site amino-acid preferences, they include statistical couplings between sequence positions. Some of these patterns are 'attractive', defining 'ideal' sequence motifs which real sequences in the protein family try to mimic. In distinction to the original Hopfield model [42], we also find 'repulsive' patterns, which define regions in the sequence space deprived of real sequences. The statistical mechanics of the inverse Hopfield model, studied in [43] for the  $q = 2$  case and extended here to the generic  $q > 2$  Potts case, shows that it naturally interpolates between PCA and DCA, and allows us to study the statistical issues raised by those approaches exposed above. We show that, in contradistinction with PCA, low eigenvalues and eigenmodes are important to recover structural information about the proteins, and should not be discarded. In addition, we propose a maximum likelihood criterion for pattern selection, not based on the comparison with the MP spectrum. We study the nature of the statistically most significant eigenmodes, and show that they exhibit remarkable features in term of localization: most repulsive patterns are strongly localized on a few sites, generally found to be in close contact on the three-dimensional structure of the proteins. As for DCA, we show that the dimensionality of the MaxEnt model can be very efficiently reduced with essentially no loss of predictive power for the contact map in the case of large multiple-sequence alignments, and with an improved accuracy in the case of small alignments containing too few sequences for standard mean-field DCA to work. These conclusions are established both from theoretical arguments, and from the direct application of the Hopfield-Potts model to a number of sample protein families.

## A short reminder of covariation analysis

Data are given in form of a *multiple sequence alignment* (MSA), in which each row contains the amino-acid sequence of one protein, and each column one residue position in these proteins, which is aligned based on amino-acid similarity. We denote the MSA by  $A = \{a_i^m | i = 1, \dots, L, m = 1, \dots, M\}$  with index  $i$  running over the  $L$  columns of the alignment (residue positions / sites), and  $m$  over the  $M$  sequences, which constitute the rows of the MSA. The amino-acids  $a_i^m$  are assumed to be represented by natural numbers  $1, \dots, q$  with  $q = 21$ , where we include the 20 standard amino acids and the alignment gap '-'.

In our approach, we do not use the data directly, but we summarize them by the amino-acid occupancies in single columns and pairs of columns of the MSA (cf. Methods for data preprocessing),

$$f_i(a) = \frac{1}{M} \sum_{m=1}^M \delta_{a, a_i^m} \quad (1)$$

$$f_{ij}(a, b) = \frac{1}{M} \sum_{m=1}^M \delta_{a, a_i^m} \delta_{b, a_j^m}, \quad (2)$$

with  $i, j = 1, \dots, L$  and  $a, b = 1, \dots, q$ . The Kronecker symbol  $\delta_{a,b}$  equals one for  $a = b$ , and zero else. Since frequencies sum up to one, we can discard one amino-acid value (*e.g.*  $a = q$ ) for each position without losing any information about the sequence statistics. We define the empirical covariance matrix through

$$C_{ij}(a, b) = f_{ij}(a, b) - f_i(a)f_j(b) , \quad (3)$$

with the position index  $i$  running from 1 to  $L$ , and the amino-acid index from 1 to  $q - 1$ . The covariance matrix  $C$  can therefore be interpreted as a square matrix with  $(q - 1)L$  rows and columns. We will adopt this interpretation throughout the paper, since the methods proposed become easier in terms of the linear algebra of this matrix.

### Maximum entropy modeling and direct couplings

Non-zero covariance between two sites does not necessarily imply the sites to directly interact for functional or structural purposes [13]. The reason is the following [17]: When  $i$  interacts with  $j$ , and  $j$  interacts with  $k$ , also  $i$  and  $k$  will show correlations even if they do not interact. It is thus important to distinguish between *direct* and *indirect* correlations, and to infer *networks of direct couplings*, which generate the empirically observed covariances. This can be done by constructing a (protein-family specific) statistical model  $P(a_1, \dots, a_L)$ , which describes the probability of observing a particular amino-acid sequence  $a_1, \dots, a_L$ . Due to the limited amount of available data, we require this model to reproduce empirical frequency counts for single MSA columns and column pairs,

$$f_i(a_i) = \sum_{\{a_k | k \neq i\}} P(a_1, \dots, a_L) \quad (4)$$

$$f_{ij}(a_i, a_j) = \sum_{\{a_k | k \neq i, j\}} P(a_1, \dots, a_L) , \quad (5)$$

*i.e.* marginal distributions of  $P(a_1, \dots, a_L)$  are required to coincide with the empirical counts up to the level of position pairs. Beyond this coherence, we aim at the *least constrained* statistical description. The *maximum-entropy principle* [23, 24] stipulates that  $P$  is found by maximizing the entropy

$$H[P] = - \sum_{a_1, \dots, a_L} P(a_1, \dots, a_L) \log P(a_1, \dots, a_L) , \quad (6)$$

subject to the constraints Eqs. (4) and (5). We readily find the analytical form

$$P(a_1, \dots, a_L) = \frac{1}{\mathcal{Z}(\{e_{ij}(a, b), h_i(a)\})} \exp \left\{ \frac{1}{2} \sum_{i,j} e_{ij}(a_i, a_j) + \sum_i h_i(a_i) \right\} , \quad (7)$$

where  $\mathcal{Z}$  is a normalization constant. The MaxEnt model thus takes the form of a (generalized)  $q$ -states Potts model, a celebrated model in statistical physics [44], or a Markov random field in a more mathematical language. The parameters  $e_{ij}(a, b)$  are the direct couplings between MSA columns, and the  $h_i(a)$  represent the local fields (position-weight matrices) acting on single sites. Their values have to be determined such that Eqs. (4) and (5) are satisfied. Note that, without the coupling terms  $e_{ij}(a, b)$ , the model would reduce to a standard position-specific scoring matrix. It would describe independent sites, and thus it would be intrinsically unable to capture residue covariation.

From a computational point of view, however, it is not possible to solve Eqs. (4) and (5) exactly. The reason is that the calculations of  $\mathcal{Z}$  and of the marginals require summations over all  $q^L$  possible amino-acid sequences of length  $L$ . With  $q = 21$  and typical protein lengths of  $L = 50 - 500$ , the numbers of configurations are enormous, of the order of  $10^{65} - 10^{650}$ . The way out is an approximate determination

of the model parameters. The computationally most efficient way found so far is an approximation, called mean field in statistical physics, leading to the approach known as *direct coupling analysis* [19]. Within this mean-field approximation, the values for the direct couplings are simply equal to

$$e_{ij}(a, b) = (C^{-1})_{ij}(a, b) \quad \forall i, j \quad \forall a, b = 1, \dots, q-1, \quad (8)$$

and  $e_{ij}(a, q) = e_{ij}(q, a) = 0$  for all  $a = 1, \dots, q$ . Note that the couplings can be approximated with this formula in a time of the order of  $L^3(q-1)^3$ , instead of the exponential time complexity,  $q^L$ , of the exact calculation. On a single desktop PC, this can be achieved in a few seconds to minutes, depending on the length  $L$  of the protein sequences.

The problem can be formulated equivalently in terms of maximum-likelihood (ML) inference. Assuming  $P(a_1, \dots, a_L)$  to be a pairwise model of the form of Eq. (7), we aim at maximizing the log-likelihood

$$\mathcal{L}[\{e_{ij}(a, b), h_i(a)\}|A] = \frac{1}{M} \sum_{m=1}^M \log P(a_1^m, \dots, a_L^m) \quad (9)$$

of the model parameters  $\{e_{ij}(a, b), h_i(a)\}$  given the MSA  $A$ . This maximization implies that Eqs. (4) and (5) hold. In the rest of the paper, we will adopt the point of view of ML inference, cf. the details given in Methods. Note that, without restrictions on the couplings  $e_{ij}(a, b)$  ML and MaxEnt inference are equivalent, but under the specific form for  $e_{ij}(a, b)$  assumed in the Hopfield-Potts model, this equivalence will break down. More precisely, the ML model will fit Eqs. (4) and (5) only approximately

Once the direct couplings  $e_{ij}(a, b)$  have been calculated, they can be used to make predictions about the contacts between residues, details can be found in the Methods Section. In [19], it was shown that the predictions for the residue-residue contacts in proteins are very accurate. In other words, DCA allows to find a very good estimate of a partial contact map from sequence data only. Subsequent works have shown that this contact map can be completed by embedding it into three dimensions [31, 33].

### Pearson correlation matrix and principal component analysis

Another way to extract information about groups of correlated residues is the following. From the covariance matrix  $C$  given in Eq. (3), we construct the Pearson correlation matrix  $\Gamma$  through the relationship

$$\Gamma_{ij}(a, b) = \sum_{c, d=1}^{q-1} (D_i)^{-1}(a, c) C_{ij}(c, d) (D_j)^{-1}(d, b), \quad (10)$$

where the matrices  $D_i$  are the square roots of the single-site correlation matrices, *i.e.*

$$C_{ii}(a, b) = \sum_{c=1}^{q-1} D_i(a, c) D_i(c, b). \quad (11)$$

This particular form of the Pearson correlation matrix  $\Gamma$  in Eq. (10) results from the fact that we have projected the  $q$ -dimensional space defined by the amino-acids  $a = 1, \dots, q$  onto the subspace spanned by the first  $q-1$  dimensions. Alternative projections lead to modified but equivalent expressions of the Pearson matrix, cf. Text S1 (Sec. S1.3). Informally speaking, the correlation  $\Gamma_{ij}(a, b)$  is a measure of comparison of the empirical covariance  $C_{ij}(a, b)$  with the single-site fluctuations taken independently. Hence,  $\Gamma$  is normalized and coincides with the  $(q-1) \times (q-1)$  identity matrix on each site:  $\Gamma_{ii}(a, b) = \delta_{a, b}$ .

We further introduce the eigenvalues and eigenvectors ( $\mu = 1, \dots, L(q-1)$ )

$$\sum_{j=1}^L \sum_{b=1}^{q-1} \Gamma_{ij}(a, b) v_{jb}^\mu = \lambda_\mu v_{ia}^\mu, \quad (12)$$

where the eigenvalues are ordered in decreasing order  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{L(q-1)}$ . The eigenvectors are chosen to form an ortho-normal basis,

$$\sum_{ia} v_{ia}^\mu v_{ia}^\nu = L \delta_{\mu,\nu} , \quad (13)$$

for all  $\mu, \nu = 1, \dots, L(q-1)$ . Principal component analysis consists in a partial eigendecomposition of  $\Gamma$ , keeping only the eigenmodes contributing most to the correlations, *i.e.* with the largest eigenvalues. All the other eigenvectors are discarded. In this way, the directions of maximum covariation of the residues are identified.

PCA was used by Casari *et al.* [8] in the context of residue covariation to identify functional sites specific to subfamilies of a protein family given by a large MSA. To do so, the authors diagonalized the comparison matrix, whose elements  $\mathcal{C}(m, m')$  count the number of identical residues for each pair of sequences ( $m, m' = 1, \dots, M$ ). Projection of sequences onto the top eigenvectors of the matrix  $\mathcal{C}$  allows to identify groups of subfamily-specific co-conserved residues responsible for subfamily-specific functional properties, called specificity-determining positions (SDP). Up to date, PCA (or the closely related multiple correspondence analysis) is used in one of the most efficient tools, called S3det, to detect SDPs [45]. PCA was also used in an approach introduced by Ranganathan and coworkers [6, 40], called statistical coupling analysis (SCA). In this approach a modified residue covariance matrix,  $\tilde{C}^{SCA}$ , is introduced :

$$\tilde{C}_{ij}^{SCA}(a, b) = w_i^a C_{ij}(a, b) w_j^b \quad (14)$$

where the weights  $w_i^a$  favor positions  $i$  and residues  $a$  of high conservation. Amino-acid indices are contracted to define the effective covariance matrix,

$$\tilde{C}_{ij}^{SCA} = \sqrt{\sum_{a,b} \tilde{C}_{ij}^{SCA}(a, b)^2} . \quad (15)$$

The entries of  $\tilde{C}^{SCA}$  depend on the residue positions  $i, j$  only. In a variant of SCA the amino-acid information is directly contracted at the level of the sequence data. A binary variable is associated to each site: it is equal to one in sequences carrying the consensus amino-acid, to zero otherwise [40]. Principal component analysis can then be applied to the  $L$ -dimensional  $\tilde{C}^{SCA}$  matrix, and used to define so-called sectors, *i.e.* clusters of evolutionarily correlated sites.

## Results

To bridge these two approaches – DCA and PCA – we introduce the Hopfield-Potts model for the maximum likelihood modeling of the sequence distribution, given the residue frequencies  $f_i(a)$  and their pairwise correlations  $f_{ij}(a, b)$ . From a mathematical point of view, the model corresponds to a specific class of Potts models, in which the coupling matrix  $e_{ij}(a, b)$  is of low rank  $p$  compared to  $L(q-1)$ . It therefore offers a natural way to reduce the number of parameters far below what is required in the mean-field approximation of [19]. In addition, the solution of the Hopfield-Potts inverse problem, *i.e.* the determination of the low rank coupling matrix  $e$ , allows us to establish a direct connection with the spectral properties of the Pearson correlation matrix  $\Gamma$  and thus with PCA.

Here, we first give an overview over the most important theoretical results for Hopfield-Potts model inference, increasing levels of detail about the algorithm and its derivation are provided in Methods and Text S1. Subsequently we discuss in detail the features of the Hopfield-Potts patterns found in three different protein families, and finally assess our capacity to detect residue contacts using sequence information alone in a larger test set of protein families.



## Inference with the Hopfield-Potts model

The main idea of this work is that, though the space of sequences is  $L(q-1)$ -dimensional, the number of spatial directions being relevant for covariation is much smaller. Such a relevant direction is called *pattern* in the following, and given by a  $L \times q$  matrix  $\xi = \{\xi_i(a)\}$ , with  $i = 1, \dots, L$  being the site indices, and  $a = 1, \dots, q$  the amino acids. The *log-score* of a sequence  $(a_1, \dots, a_L)$  for one pattern  $\xi$  is defined as

$$S(a_1, \dots, a_L | \xi) = \left[ \sum_{i=1}^L \xi_i(a_i) \right]^2. \quad (16)$$

This expression bears a strong similarity with, but also a crucial difference to a position-specific scoring matrix (PSSM): As in a PSSM, the log-score depends on a sum over position and amino-acid specific contributions, but its non-linearity (the square in Eq. (16)) introduces residue-residue couplings, and thus is essential to take covariation into account.

In the Hopfield-Potts model, the probability of an amino-acid sequence  $(a_1, \dots, a_L)$  depends on the combined log-scores along a number  $p$  of patterns through

$$P(a_1, \dots, a_L) = \frac{1}{Z} \exp \left\{ \frac{1}{2L} \sum_{\mu=1}^{p_+} S(a_1, \dots, a_L | \xi^{+, \mu}) - \frac{1}{2L} \sum_{\nu=1}^{p_-} S(a_1, \dots, a_L | \xi^{-, \nu}) + \sum_{i=1}^L h_i(a_i) \right\}. \quad (17)$$

Patterns denoted with a  $+$ -superscript,  $\xi^{+, \mu}$  with  $\mu = 1, \dots, p_+$ , are said to be *attractive*, while the patterns labeled with a  $-$ -superscript,  $\xi^{-, \nu}$  for  $\nu = 1, \dots, p_-$ , are called *repulsive*. For the probability  $P(a_1, \dots, a_L)$  to be large, the log-scores  $S(a_1, \dots, a_L | \xi)$  for attractive patterns must be large, whereas the log-scores for repulsive patterns must be small (close to zero). As we will see in the following, the inclusion of such repulsive patterns is important: Compared to the mixed model (17), a model with only attractive patterns achieves a much smaller likelihood (at each given total number of parameters) and a strongly reduced predictivity of residue-residue contacts.

It is easy to see that Eq. (17) corresponds to a specific choice of the couplings  $e_{ij}(a, b)$  in Eq. (7), namely

$$e_{ij}(a, b) = \frac{1}{L} \sum_{\mu=1}^{p_+} \xi_i^{+, \mu}(a) \xi_j^{+, \mu}(b) - \frac{1}{L} \sum_{\nu=1}^{p_-} \xi_i^{-, \nu}(a) \xi_j^{-, \nu}(b), \quad (18)$$

where, without loss of generality, the  $q^{th}$  component of the patterns is set to zero,  $\xi_i^{+, \mu}(q) = \xi_i^{-, \nu}(q) = 0$ , for compatibility with the mean-field approach exposed above. Note that the coupling matrix, for linearly independent patterns, has rank  $p = p_+ + p_-$ , and is defined from  $p L(q-1)$  pattern components only, instead of  $\mathcal{O}(L^2(q-1)^2)$  parameters for the most general case of coupling matrices  $e_{ij}(a, b)$ . When  $p = L(q-1)$ , *i.e.* when all the patterns are taken into account, the coupling matrix  $e$  has full rank, and the Hopfield-Potts model is identical to the Potts model used to infer the couplings in DCA in [19]. All results of mean-field DCA are thus recovered in this limiting case.

The patterns are to be determined by ML inference, cf. Methods and Text S1 for details. In mean-field approximation, they can be expressed in terms of the eigenvalues and eigenvectors of the Pearson correlation matrix  $\Gamma$ , which were defined in Eq. (12). We find that attractive patterns correspond to the  $p_+$  largest eigenvalues ( $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{p_+} \geq 1$ ),

$$\xi_i^{+, \mu}(a) = \left( 1 - \frac{1}{\lambda_\mu} \right)^{1/2} \tilde{v}_{ia}^\mu, \quad \mu = 1, \dots, p_+, \quad (19)$$

and repulsive patterns to the  $p_-$  smallest eigenvalues ( $\lambda_{L(q-1)} \leq \lambda_{L(q-1)-1} \leq \dots \leq \lambda_{L(q-1)+1-p_-} \leq 1$ ),

$$\xi_i^{-, \nu}(a) = \left( \frac{1}{\lambda_{L(q-1)+1-\nu}} - 1 \right)^{1/2} \tilde{v}_{ia}^{L(q-1)+1-\nu}, \quad \nu = 1, \dots, p_-, \quad (20)$$



where, for all  $\mu = 1, \dots, L(q-1)$ ,

$$\tilde{v}_{ia}^\mu = \sum_{b=1}^{q-1} (D_i)^{-1}(a, b) v_{ib}^\mu . \quad (21)$$

The prefactor  $|1 - 1/\lambda|^{1/2}$  vanishes for  $\lambda = 1$ . It is not surprising that  $\lambda = 1$  plays a special role, as it coincides with the mean of the eigenvalues:

$$\frac{1}{L(q-1)} \sum_{\mu=1}^{L(q-1)} \lambda_\mu = \frac{1}{L(q-1)} \sum_{i=1}^L \sum_{a=1}^{q-1} \Gamma_{ii}(a, a) = 1 . \quad (22)$$

In the absence of any covariation between the residues  $\Gamma$  becomes the identity matrix, and all eigenvalues are unity. Hence all patterns vanish, and so does the coupling matrix (18). The Potts model (17) depends only on the local bias parameters  $h_i(a)$ , and it reduces to a PSSM describing independent sites.

The eigenvectors of the correlation matrix with large eigenvalues  $\lambda_\mu \gg 1$  contribute most to the covariation observed in the MSA (i.e. to the matrix  $\Gamma$ ), but they do not contribute most to the coupling matrix  $e$ . In the expression (18) for this matrix, each pattern carries a prefactor  $|1 - 1/\lambda_\mu|$ : Whereas this prefactor remains smaller than one for attractive patterns ( $\lambda_\mu > 1$ ), it can become very large for repulsive patterns ( $\lambda_\mu < 1$ ), see Fig. 1 (right panel). Thus, the contribution of a repulsive pattern to the  $e$  matrix may be much larger than the contribution of any attractive pattern.

Eqs. (19) and (20) *a priori* define  $L(q-1)$  different patterns, therefore we need a rule for selecting the  $p$  'best', i.e. most likely patterns. We show in Methods that the contribution of a pattern to the model's log-likelihood  $\mathcal{L}$  defined in Eq. (9) is a function of the associated eigenvalue  $\lambda$  only,

$$\Delta\mathcal{L}(\lambda) = \frac{1}{2} \left( \lambda - 1 - \log \lambda \right) . \quad (23)$$

As is shown in Fig. 1 (left panel), large contributions arrive from both the largest and the smallest eigenvalues, whereas eigenvalues close to unity contribute little. According to ML inference, we have to select the  $p$  eigenvalues with largest contributions. To this end, we define a threshold value  $\theta$  such that there are exactly  $p$  patterns with larger contributions  $\Delta\mathcal{L} > \theta$  to the log-likelihood; the  $L(q-1) - p$  patterns with smaller  $\Delta\mathcal{L}$  are omitted in the expression for the couplings Eq. (18). In accordance with Fig. 1, we determine thus the two positive real roots  $\ell_\pm$  ( $\ell_- < 1 < \ell_+$ ) of the equation

$$\Delta\mathcal{L}(\ell_\pm) = \theta , \quad (24)$$

and include all repulsive patterns with  $\lambda_{L(q-1)+1-\nu} < \ell_-$ , calling their number  $p_-$ , and all attractive patterns with  $\lambda_\mu > \ell_+$ , denoting their number by  $p_+$ . The total number of selected patterns is thus  $p = p_- + p_+$ .

## Features of the Hopfield-Potts patterns

We have tested the above inference framework in great detail using three protein families, with variable values of protein length  $L$  and sequence number  $M$ :

- The *Kunitz/Bovine pancreatic trypsin inhibitor* domain (PFAM ID PF00014) is a relatively short ( $L = 53$ ) and not very frequent ( $M = 2,143$ ) domain, after reweighting the effective number of diverged sequences is  $M_{eff} = 1,024$  (cf. Eq. (28) in Methods for the definition). Results are compared to the exemplary X-ray crystal structure with PDB ID 5pti [46].
- The bacterial *Response regulator* domain (PF00072) is of medium length ( $L = 112$ ) and very frequent ( $M = 62,074$ ). The effective sequence number is  $M_{eff} = 29,408$ . The PDB structure used for verification has ID 1nxw [47].

- The eukaryotic signaling domain *Ras* (PF00071) is the longest ( $L = 161$ ) and has an intermediate size MSA ( $M = 9,474$ ), leading to  $M_{eff} = 2,717$ . Results are compared to PDB entry 5p21 [48].

In a second step, we have used the 15 protein families studied in [33] to verify that our findings are not specific to the three above families, but generalize to other families. A list of the 15 proteins together with the considered PDB structures is provided in Text S1, Section 4.

To interpret the Hopfield patterns in terms of amino-acid sequences, we first report some empirical observations made for the patterns corresponding to the largest and smallest eigenvalues, *i.e.* to the most likely attractive and repulsive patterns. We concentrate our discussion in the main text on one protein family, the Trypsin inhibitor (PF00014). Analogous properties in the other two protein families are reported in Text S1.

The upper panel of Fig. 2 shows the spectral density. It is characterized by a pronounced peak around eigenvalue 1. The smallest eigenvalue is  $\lambda_m^{PF00014} \sim 0.1$ , the largest is  $\lambda_M^{PF00014} \sim 23$ . Large eigenvalues are isolated from the bulk of the spectrum, small eigenvalues are not.

To characterize the statistical properties of the patterns we define, inspired by localization theory in condensed matter physics, the inverse participation ratio (IPR) of a pattern  $\xi$  as

$$\text{IPR}(\xi) = \frac{\sum_{i,a} \xi_i(a)^4}{\left(\sum_{i,a} \xi_i(a)^2\right)^2}. \quad (25)$$

Possible IPR values range from one for perfectly localized patterns (only one single non-zero component) to  $1/(L(q-1))$  for a completely distributed pattern with uniform entries. IPR is therefore used as a localization measure for the patterns: Its inverse  $1/\text{IPR}(\xi)$  is an estimate of the effective number  $N_{eff}(\xi)$  of pairs  $(i, a)$ , on which the pattern has sizable entries  $\xi_i(a)$ . The middle panel of Fig. 2 shows the presence of strong localization for repulsive patterns (small eigenvalues) and for irrelevant patterns (around eigenvalue 1). A much smaller increase in the IPR is also observed for part of the large eigenvalues.

What is the typical contribution  $\delta e(\xi)$  of a pattern  $\xi$  to the couplings? Pattern  $\xi$  contributes  $\delta e_{ij}(a, b) = \frac{1}{L} \xi_i(a) \xi_j(b)$  to each coupling. Many contributions can be small, and others may be larger. An estimate of the magnitude of those relevant contributions can be obtained from the sum of the squared contributions normalized by the effective number  $N_{eff}(\xi)^2$  of pairs  $(i, a), (j, b)$  on which the patterns has large entries:

$$\delta e(\xi) = \sqrt{\frac{1}{N_{eff}(\xi)^2} \sum_{i,j,a,b} (\delta e_{ij}(a, b))^2} = \text{IPR}(\xi) \times \frac{1}{L} \sum_{i,a} \xi_i(a)^2. \quad (26)$$

The lower panel of Fig. 2 shows the typical contribution  $\delta e$  of a pattern as a function of its corresponding eigenvalue. Patterns with eigenvalues close to 1 have very small norms; they essentially do not contribute to the couplings. Highly localized patterns of large norm result in few and large contributions to the couplings ( $\lambda \ll 1$ ). Patterns associated to large eigenvalues  $\lambda \gg 1$  produce many weak contributions to the couplings.

### Repulsive patterns

In the upper row of Fig. 3 we display the three most localized repulsive patterns (smallest, 3rd and 4th smallest eigenvalues) for the trypsin inhibitor protein (PF00014). All three patterns have two very pronounced peaks, corresponding to, say, amino-acid  $a$  in position  $i$  and amino-acid  $b$  in position  $j$ , and some smaller minor peaks, resulting in IPR values above 0.3. For each pattern, the two major peaks are of opposite sign:  $\xi_i(a) \simeq -\xi_j(b)$ . As a consequence, amino-acid sequences carrying amino-acid  $a$  in position  $i$ , but not  $b$  in position  $j$  (as well as sequences carrying  $b$  in  $j$  but not  $a$  in  $i$ ) show large log-scores  $S \simeq [\xi_i(a)]^2$ , cf. Eq. (16). Their probability in the Hopfield-Potts model, given by (17), will be

strongly reduced as compared to the probability of sequences carrying either both amino-acids  $a$  and  $b$  in, respectively, positions  $i$  and  $j$ , or none of the two (scores  $S$  close to zero). Hence, we see that repulsive patterns do define repulsive directions in the sequence space, which tend to be avoided by sequences. A more thorough discussion of the meaning of repulsive patterns will be given in the Discussion Section.

In all three panels of Fig. 3, the two large peaks have highest value for the amino acid cysteine. Actually, for all of them, the pairs of peaks identify disulfide bonds, *i.e.* covariant bonds between two cysteines. They are very important for a protein’s stability and therefore highly conserved. The corresponding repulsive patterns forbid amino-acid configurations with a single cysteine in only one out of the two positions. Both residues are co-conserved. Note also that the trypsin inhibitor has only three disulfide bonds, *i.e.* all of them are seen by the most localized repulsive patterns. The second eigenvalues, which has a slightly smaller IPR, is actually found to be a mixture of two of these bonds, *i.e.* it is localized over four positions.

The observation of disulfide bonds is specific to the trypsin inhibitor. In other proteins, also the ones studied in this paper, we find similarly strong localization of the most repulsive patterns, but in different amino acid combinations. As an example, the most localized pattern in the response regulator domain connects a position with an Asp residue (negatively charged), with another position carrying either Lys or Arg (both positively charged), their interaction is thus coherent with electrostatics. In all observed cases, the consequence is a strong statistical coupling of these positions, which are typically found in direct contact.

### Attractive patterns

The strongest attractive pattern, *i.e.* the one corresponding to the largest eigenvalue  $\lambda_1$ , is shown in the leftmost panel of the lower row of Fig. 3. Its IPR is small ( $\sim 0.003$ ), implying that it is extended over most of the protein (a pattern of constant entries would have IPR  $1/(L(q-1)) \simeq 0.001$ ). As is shown in Text S1, strongest entries in  $\xi_i^1(a)$  correspond to conserved residues and these, even if they are distributed along the primary sequence, tend to form spatially connected and functionally important regions in the folded protein (*e.g.* a binding pocket), cf. left panel of Fig. 4. Clearly this observation is reminiscent of the protein sectors observed in [40], which are found by PCA applied to the before-mentioned modified covariance matrix. Note, however, that sectors are extracted from more than one principal component.

More characteristic patterns are found for the second and third eigenvalues. As is shown in Fig. 3, they show strong peaks at the extremities of the sequence, which become higher when approaching the first resp. last sequence position. The peaks are, for all relevant positions, concentrated on the gap symbol. These patterns are actually artifacts of the multiple-sequence alignment: Many sequences start or end with a stretch of gaps, which may have one out of at least three reasons: (1) The protein under consideration does not match the full domain definition of PFAM; (2) the local nature of PFAM alignments has initial and final gaps as algorithmic artifacts, a correction would however render the search tools less efficient; (3) in sequence alignment algorithms, the extension of an existing gap is less expensive than opening a new gap. The attractive nature of these two patterns, and the equal sign of the peaks, imply that gaps in equilibrium configurations of the Hopfield-Potts model frequently come in stretches, and not as isolated symbols. The finding that there are two patterns with this characteristic can be traced back to the fact that each sequence has two ends, and these behave independently with respect to alignment gaps.

### Theoretical results for localization in the limit case of strong conservation

The main features of the empirically observed spectral and localization properties of Fig. 2 can be found back in the limiting case of completely conserved sequences, which is amenable to an exact mathematical treatment. To this end, we consider  $L$  perfectly conserved sites, *i.e.* a MSA made from the repetition of

a unique sequence. As is shown in Text S1, Section 2, the corresponding Pearson correlation matrix  $\Gamma$  has only three different eigenvalues:

- a large and non-degenerate eigenvalue,  $\lambda_+$ , which is a function of  $q$  and  $L$  (and of the pseudocount used to treat the data, see Methods), whose corresponding eigenvector is extended;
- a small and  $(L-1)$ -fold degenerate eigenvalue,  $\lambda_- = (L-\lambda_+)/ (L-1)$ . The corresponding eigenspace is spanned by vectors which are perfectly localized in pairs of sites, with components of opposite signs;
- the eigenvalue  $\lambda = 1$ , which is  $L(q-2)$ -fold degenerate. The eigenspace is spanned by vectors, which are localized over single sites.

For a realistic MSA, *i.e.* without perfect conservation, degeneracies will disappear, but the features found above remain qualitatively correct. In particular, we find in real data a pronounced peak of eigenvalues around 1, corresponding to localized eigenmodes (Fig. 2). In addition, low-eigenvalue modes are found to be strongly localized, and the order of magnitude of  $\lambda_- \simeq 0.09$  is in good agreement with the smallest eigenvalues,  $\simeq 0.1$ , reported for the three analyzed domain families. Finally, the largest eigenmodes are largely extended, as found in the limit case above. Note that the eigenvalues found in the protein spectra, *e.g.*  $\lambda_1 \simeq 23$  for PF00014, are however smaller than in the limit case,  $\lambda_+ \simeq 48$ , due to only partial conservation in the real MSA.

## Residue-residue contact prediction with the Hopfield-Potts model

The most important feature of DCA is its ability to predict pairs of residues, which are distantly positioned in the sequence, but which form native contacts in the protein's tertiary structure, cf. the right panel of Fig. 4. Here, our contact prediction is based on the sampling-corrected Frobenius norm of the  $(q-1)$ -dimensional statistical coupling matrices  $e_{ij}$ , cf. Methods, which in [49] has been shown to outperform the direct-information measure used in [17]. This measure assigns a single scalar value for the strength of the direct coupling between two residue positions.

The contact map predicted from the 50 strongest direct couplings for the PF00014 family is compared to the native contact map in Fig. 5. In accordance with [19], a residue pairs is considered to be a true positive prediction if its minimal atom distance is below  $8\text{\AA}$  in the before mentioned exemplary protein crystal structures. This relatively large cutoff was chosen since DCA was found to extract a bimodal signal with pairs in the range below  $5\text{\AA}$  (turquoise in Fig. 5) and others with  $7-8\text{\AA}$  (grey in Fig. 5); both peaks contain valuable information if compared to typical distances above  $20\text{\AA}$  for randomly chosen residue pairs. To include only non-trivial contacts, we require also a minimum separation  $|i-j| > 4$  of at least 5 residues along the protein sequence. Remarkably the quality of the predicted contact map with the Hopfield-Potts model with  $p = 128$  patterns is essentially the same as with DCA, corresponding to  $p = L(q-1) = 1060$  patterns. In both cases predicted contacts spread rather uniformly over the native contact map, and 96% of the predicted contacts are true positives. This result is corroborated by the lower panels of Fig. 6, which show, for various values of the number  $p$  of patterns, the performance in terms of contact predictions for the three families studied here. The plots show the fraction of true-positives (TP), *i.e.* of native distances below  $8\text{\AA}$ , in between the  $x$  pairs of highest couplings, as a function of  $x$  [19].

The three upper panels in Fig. 6 show the ratio between the selected pattern contributions to the log-likelihood,  $\sum_{\{\mu|\lambda_\mu \notin [\ell_-, \ell_+]\}} \Delta\mathcal{L}(\lambda_\mu)$ , and its maximal value obtained by including all  $L(q-1)$  patterns,  $\sum_{\mu=1}^{L(q-1)} \Delta\mathcal{L}(\lambda_\mu)$ . A large fraction of patterns can be omitted without any substantial loss in log-likelihood, but with a substantially smaller number of parameters. It is worth noting that, in Fig. 6, we do not find any systematic benefit of excluding patterns for the contact prediction, but the predictive power decreases initially only very slowly with decreasing pattern numbers  $p$ . For all three proteins, even with  $\sim 128$  patterns, very good contact predictions can be achieved, which are comparable to the ones

with  $L(q - 1) = 1060 - 3220$  patterns using the full DCA inference scheme of [19]. Almost perfect performance is reached, when the contribution of selected patterns to the log-likelihood is only at 60 – 80% of its maximal value. This could be expected from the fact that patterns corresponding to eigenvalues close to unity hardly contribute to the couplings, cf. lower panel in Fig. 2.

These findings are not restricted to the three test proteins, as is confirmed by the left panel of Fig. 7. In this figure, we average the TP rates for  $p = 8, 32, 128, 512$  and  $L(q - 1)$  (*i.e.* full mean-field DCA) for the 15 proteins studied in [33], which had been selected for their diversity in protein length and fold type. Further more, the discussion of the localization properties of repulsive patterns is corroborated by the results reported in Fig. 7, right panel. It compares the performance of the Hopfield-Potts model to predict residue-residue contacts, for the three cases where  $p = 100$  patterns are selected either according to the maximum likelihood criterion (patterns for eigenvalues  $\lambda < \ell_-$  and for  $\lambda > \ell_+$ ), or where only the strongest attractive ( $\lambda > \ell_+$ ) or only the strongest repulsive ( $\lambda < \ell_-$ ) patterns are taken into account. It becomes evident that repulsive patterns provide more accurate contact information, TP rates are almost unchanged between the curve of the  $p = 100$  most likely patterns, and the smaller subset of repulsive patterns. On the contrary, TP rates for contact prediction are strongly reduced when considering only attractive patterns, *i.e.* in the case corresponding most closely to PCA. This finding illustrates one of the most significant differences between DCA and PCA: Contact information is provided by the eigenvectors of the Pearson correlation matrix  $\Gamma$  in the lower tail of the spectrum.

As is discussed in the previous section, patterns with the largest contribution to the log-likelihood are dominated by (and localized in) conserved sites. Attractive patterns favor these sites to jointly assume their conserved values, whereas repulsive patterns avoid configurations where, in pairs of co-conserved sites, only one variable assumes its conserved value, but not the other one. However, we have also seen that an accurate contact prediction requires at least  $\sim 100$  patterns, *i.e.* it goes well beyond the patterns given by strongly conserved sites. In Fig. 4 we show, for the exemplary case of the Trypsin inhibitor, both the 10 sites of highest entry in the most attractive pattern  $\xi^{+,1}$  (corresponding to conserved sites), and the first 50 predicted intra-protein contacts using the full mean-field DCA scheme (results for  $p = 512$  are almost identical). It appears that many of the correctly predicted contacts are not included in the set of the most conserved sites. From a mathematical point of view, this is understandable - only variable sites may show covariation. From a biological point of view, this is very interesting, since it shows that highly variable residue in proteins are not necessarily functionally unimportant in a protein family, but they may undergo strong coevolution with other sites, and thus be very important for the structural stability of the protein, cf. also the Supporting Fig. S5 where the degree of conservation [50] is depicted for the highest-ranking DCA predicted contacts. In this figure we show that residues included in predicted contacts are found for all levels of conservation. It has, however, to be mentioned that in the considered MSA, there are no 100% conserved residues, the latter would not show any covariation. A small level of variability is therefore crucial for our approach.

A remark is necessary concerning the right panel of Fig. 4: Whereas conserved sites (which carry also the largest entries of the pattern with maximum eigenvalue) are collected in one or two spatially connected regions in the studied proteins, this is not necessarily true for all proteins. In particular complex domains with multiple functions and/or multiple conformations may show much more involved patterns. It is, however, beyond the scope of this paper to shed light onto the details of the biological interpretation of the principal components of  $\Gamma$ .

In which cases does the dimensional reduction achieved by selecting only a relatively small number of patterns provide an actual advantage over the standard mean-field DCA approach with  $p = L(q - 1)$  patterns? We have seen that for relatively large MSAs, where DCA gives very accurate results, the approach presented here achieves a strong dimensional reduction almost without loss in predictive power, but it did not improve the contact map prediction, cf. Figs. 5 and 6. However, when we reduce the number of sequences in the MSA, DCA undergoes a strong reduction in accuracy of prediction, see the full lines in Fig. 8 where DCA is applied to sub-alignments of the PF00014 domain family. Repeating

the same experiment with a finite number of patterns ( $p = 16$  in Fig. 8), the MSA-size dependence is strongly reduced. For very small alignments of only 10-30 sequences, the Hopfield-Potts model is still able to extract contacts with an astonishing TP rate of 70-80%, whereas DCA produces almost random results (TP rate ca. 30%). The success of the Hopfield-Potts approach for small MSA is not specific to the PF00014 domain, and holds for other protein families, see Fig. S15 in Text S1. Hopfield-Potts patterns are therefore an efficient means to reduce overfitting effects found in DCA, and to improve the signal-to-noise ratio.

## Discussion

In this paper we have proposed a method to analyze the correlation matrix of residue occurrences across multiple-sequence alignments of homologous proteins, based on the inverse Hopfield-Potts model. Our approach offers a natural interpolation between the spectral analysis of the correlation matrix, carried out in principal component analysis (PCA), and maximum entropy approaches which aim at reproducing those correlations within a global statistical model (*e.g.* DCA). The inverse Hopfield-Potts model requires to infer “directions” of particular importance in the sequence space, called patterns: The distribution of sequences belonging to a protein family tends to accumulate along attractive patterns (related to eigenmodes of the correlation matrix with large eigenvalues) and to get depleted around repulsive patterns (related to the low-eigenvalue modes). These patterns have some similarity with position-specific scoring matrices frequently used in the statistical modeling of sequences, but in contrast to the independence of different positions in PSSM, Hopfield-Potts patterns account for inter-position couplings, as needed for coevolutionary analysis.

Contrary to principal component analysis, which discards low-eigenvalue modes, we have shown that repulsive patterns are essential to characterize the sequence distribution, and in particular to detect structural properties (residue-residue contacts) of proteins from sequence data. In addition, we have shown how to infer not only the values of the patterns but also their statistical relevance from the sequence data. To do so, we have calculated the contribution of each pattern to the total likelihood of the Hopfield-Potts model given the data, establishing thus a clear criterion for pattern selection. The results of the application of the inverse Hopfield-Potts model to real sequence data confirm that most eigenmodes (with eigenvalues close to unity) can be discarded without affecting considerably the contact prediction (see Fig. 5 and Fig. 6). This makes our approach much less parameter-intensive than the full direct coupling analysis DCA. We have found empirically that it is sufficient to take into account the patterns contributing to  $\sim 60 - 80\%$  of the log-likelihood to achieve a very good contact map prediction in the case of large multiple-sequence alignments. In the case of reduced MSA size, we found that the dimensional reduction due selecting only the most likely patterns improves the signal-to-noise ratio of the inferred model, and therefore reaches a better contact prediction than mean-field DCA, down to very small numbers of sequences, see Fig. 8 and Fig. S15 in Text S1. Moreover the Hopfield-Potts approach can be very advantageous in terms of computational time. While DCA requires the inversion of the correlation matrix, which takes  $O(L^3(q-1)^3)$  time, computing the  $p$  patterns (corresponding to the largest and smallest eigenvalues) can be done in  $O(pL^2(q-1)^2)$  time only. The reduction in computational time can thus be very important for large proteins.

We have also studied the position-specific nature of patterns, taking inspiration from localization theory in condensed matter physics and random matrix theory (Fig. 3 and Figs. S8 and S12 in Text S1). Briefly speaking, a pattern is said to be localized if it is concentrated on a few sites of the sequence, and extended (over the sequence) otherwise. We have found that the principal attractive pattern (corresponding to the largest eigenvalue) is extended, with entries of largest absolute value in the most conserved sites (Figs. S3, S4, S9 & S13 in Text S1). Other strongly attractive patterns can be explained from the presence of extended gaps in the alignment, mostly found at the beginning or at the end of sequences. The other patterns of large likelihood contributions are repulsive, *i.e.* they correspond to small eigenvalues,



usually discarded by principal component analysis. Interestingly, these patterns appear to be strongly localized, that is, strongly concentrated in very few positions, which despite their separation along the sequence are found in close contact in the 3D protein structure. To give an example, in the Trypsin inhibitor protein, they are localized in position pairs carrying Cysteine, and being linked by disulfide bonds. Other amino-acid combinations were also found in the other protein families studied here, e.g. patterns connecting residues of opposite electrical charge. Taking into account only a number  $p$  of such repulsive patterns results in a predicted contact map of comparable quality to the one using maximum-likelihood selection, whereas the same number  $p$  of attractive patterns performs substantially worse (Fig. 7 and Fig. S7 & S11 in Text S1). The dimensional reduction of the Hopfield-Potts model compared to the Potts model (used in standard DCA) is thus even more increased as many relevant patterns are localized and contain only a few (substantially) non-zero components. As a consequence the couplings found with the Hopfield-Potts model are sparser than their DCA counterparts (Fig. S6 in Text S1).

It is important to stress that also distinct patterns, whether attractive or repulsive, can have large components on the same sites and residues. A general finding, supported by a theoretical analysis in the Results section, is that the more repulsive patterns are, the stronger they are localized, and the more conserved are the residues supporting them. Highly conserved sites therefore appear both in the most attractive pattern and, when covarying with other residues, in a few localized and repulsive patterns reflecting those covariations. As the number of patterns to be included to reach an accurate contact prediction is a few hundreds for the protein families considered here, the largest components of the weakly repulsive patterns, *i.e.* with the eigenvalues smaller than, but close to the threshold  $\theta$ , correspond to weakly conserved residues. In consequence many predicted contacts connect low-conservation residues. This statement is apparent from Fig. 4 and Figs. S10 and S14 in Text S1, which compare the sets of conserved sites and the pairs of residues predicted to be in contact by our analysis.

Why are repulsive patterns so successful in identifying contacts, in difference to attractive patterns? To answer this question, consider the simple case of a pattern  $\xi$  localized in two residues only, say it should prefer the co-occurrence of amino-acid  $a$  in position  $i$ , and of amino acid  $b$  in position  $j$ . We further assume that the two non-zero components  $\xi_i(a)$  and  $\xi_j(b)$  have the same amplitude and differ only by sign, *i.e.*  $\xi_i(a) = -\xi_j(b)$ . Now we consider a sequence of amino-acids  $(a_1, \dots, a_L)$  and ask whether it will have a large log-score  $S$  for pattern  $\xi$ , see Eq. (16). The outcome is given in the third column of Table 1. The log-score therefore corresponds to a XOR (exclusive or) between the presence of the two amino-acids  $a$  and  $b$  on their respective positions  $i$  and  $j$  in the sequence. If the pattern were attractive (cf. fourth column), it would favor sequences where exactly one of the two specified amino-acids is present. For a repulsive pattern (cf. fifth column), low log-score sequences are favored, *i.e.* either both  $a$  and  $b$  are present in positions  $i$  and  $j$ , or none of the two.

In case we assumed equal sign components, *i.e.*  $\xi_i(a) = \xi_j(b)$ , we would have found Table 2. This choice is poor in terms of enforcing covariation in the sequence: An attractive (resp. repulsive) pattern strongly favors (resp. disfavors) the simultaneous presence of amino acids  $a$  and  $b$  in positions  $i$  and  $j$ , but the likelihood is monotonous in the number of correctly present amino acids.

As a conclusion, we find that strong covariation can be efficiently enforced only by a repulsive pattern with opposite components (fifth column in Table 1). The acceptance of the (NO,NO) configuration is desirable, too: It signals the possibility of compensatory mutations, *i.e.* favorable double mutations changing both  $a$  and  $b$  in positions  $i$  and  $j$  to alternative amino acids. It is easy to generalize the above patterns to patterns having more than one favored amino-acid combination, *e.g.* favored pairings  $(a, b)$  and  $(c, d)$  can be enforced by a repulsive pattern with  $\xi_i(a) = -\xi_i(c) = -\xi_j(b) = \xi_j(d)$ . This theoretical argument explains why localized repulsive patterns critically encode for covariation. Remarkably the condition that the few, large components of repulsive patterns should sum up to zero agrees well with our findings in real MSAs, cf Fig. 3 and Figs. S8 and S12 in Text S1. Furthermore, it would be interesting to better understand the relationship between such localized patterns and specificity-determining positions [8, 45]: SDP are co-conserved in subfamilies of the full MSA, but vary from one family to another.



The most repulsive patterns are localized in residues, which are strongly conserved throughout the full alignment. We have also used S3det [45] to predict SDPs and to compare them to our 30 highest-scoring contact predictions, and we have not observed any particular signal. It would be interesting to extend the Hopfield-Potts approach to subfamilies and to investigate, if SDPs correspond to repulsive patterns in these subfamilies.

Last but not least, let us emphasize the importance of the prefactor  $|1 - \frac{1}{\lambda}|^{1/2}$  of the pattern, cf. Eqs. (19) and (20), where  $\lambda$  is the eigenvalue corresponding to the pattern. While this factor is at most equal to 1 for attractive patterns, it can take arbitrarily large values for repulsive patterns (Fig. 1, right panel). Moreover, repulsive patterns can be highly localized: they strongly contribute to a few couplings  $e_{ij}(a, b)$ , *e.g.* to one coupling between a single pair of positions  $i$  and  $j$  for patterns perfectly localized in two sites only (cf. Fig. 2, lower panel, and Fig. 3). Consequently those contributions are of particular importance in the ranking of couplings, which our contact prediction is based on. On the contrary, attractive patterns, even with sizeable norms, produce many weaker contributions to the couplings (cf. Fig. 2, lower panel), and do not alter their relative rankings as much as repulsive patterns do. This explains why contact prediction based on repulsive patterns only is much more efficient than when based on attractive patterns only (cf. Fig. 7).

Some aspects of the approach presented in this paper deserve further studies, and may actually lead to substantial improvements of our ability to detect residue contacts from statistical sequence analysis. The probably most important question is the capability of our approach to suppress noise in small MSAs, and to extract contact information in cases where mean-field DCA fails. This question is closely related to the determination of an optimal value for the pattern number  $p$  using sequence information alone. Second, the non-independence of sequences in the alignment, *e.g.* due to phylogenetic correlations, should be taken into account in a more accurate way than done currently by sequence reweighting. Third, the precise role of the – heuristically determined – large pseudo-count used to calculate the Pearson correlation matrix should also be elucidated. Fourth, while the use of the Frobenius norm for the coupling  $e_{ij}(a, b)$  (with the average-product correction, see Methods) has proven to be an efficient criterion for contact prediction, it remains unclear if there exist other contact estimators with better performance. In this context it would also be interesting to find a threshold for these contact scores, which separates a signal-rich from a noise-dominated region. And last but not least, it would be interesting to integrate prior knowledge about proteins, like *e.g.* amino-acid properties or predicted secondary structure, into the purely statistical inference approach presented here.

The MATLAB program necessary for the analysis of the data, the computation of the patterns, and the contact prediction is available as part of the Supporting Information. Users of the program are kindly requested to cite the present work.

## Methods

### Data preprocessing

Following the discussion of [19], we introduce two modifications into the definition Eq. (2) of the frequency counts  $f_i(a)$  and  $f_{ij}(a, b)$ :

- *Pseudocount regularization*: Some amino-acid combinations  $(a, b)$  do not exist in column pairs  $(i, j)$ , even if  $a$  is found in  $i$ , and  $b$  in  $j$ . This would formally lead to infinitely large coupling constants, and the covariance matrix  $C$  becomes non invertible. This divergence can be avoided by introducing a pseudocount  $\tilde{\nu}$ , which adds to the occurrence counts of each amino acid in each column of the MSA.
- *Reweighting*: The sampling of biological sequences is far from being identically and independently distributed (i.i.d.) , it is biased by the phylogenetic history of the proteins and by the human

selection of sequenced species. This bias will introduce global correlations. To reduce this effect, we decrease the statistical weight of sequences having many similar ones in the MSA. More precisely, the weight of each sequence is defined as the inverse number of sequences within Hamming distance  $d_H < xL$ , with an arbitrary but fixed  $x \in (0, 1)$ :

$$w_m = \frac{1}{||\{n | 1 \leq n \leq M; d_H[(a_1^n, \dots, a_L^n), (a_1^m, \dots, a_L^m)] \leq xL\}||} \quad (27)$$

for all  $m = 1, \dots, M$ . The weight equals one for isolated sequences, and becomes smaller the denser the sampling around a sequence is. Note that  $x = 0$  would account to removing double counts from the MSA. The total weight

$$M_{eff} = \sum_{m=1}^M w_m \quad (28)$$

can be interpreted as the effective number of independent sequences.

With these two modifications, frequency counts become

$$f_i(a) = \frac{1}{M_{eff} + \tilde{\nu}} \left[ \frac{\tilde{\nu}}{q} + \sum_{m=1}^M w_m \delta_{a, a_i^m} \right] \quad (29)$$

$$f_{ij}(a, b) = \frac{1}{M_{eff} + \tilde{\nu}} \left[ \frac{\tilde{\nu}}{q^2} + \sum_{m=1}^M w_m \delta_{a, a_i^m} \delta_{b, a_j^m} \right]. \quad (30)$$

Values  $\tilde{\nu} \simeq M_{eff}$  and  $x \simeq 0.2$  were found to work optimally across many protein families [19], we use these values. Besides these modifications, the Hopfield-Potts-model learning is performed as explained before.

## The number of independent model parameters

Amino-acid frequencies are not independent numbers. For instance, on each site  $i$ , the  $q$  amino-acid frequencies add up to one,

$$\sum_{a=1}^q f_i(a) = 1, \quad (31)$$

and two-site distributions have single-site distributions as marginals,

$$\sum_{a=1}^q f_{ij}(a, b) = f_j(b). \quad (32)$$

As a consequence, not all of the constraints (4) and (5) are independent, and the Potts model as given in Eq. (7) has more free parameters than needed to fulfill the constraints. Families of distinct parameter values result in the same model  $P(a_1, \dots, a_L)$  (in physics language, this corresponds to a gauge invariance: any function  $g_i(a)$  can be added to  $e_{ij}(a, b)$  and, simultaneously, be subtracted from  $h_i(a)$ , without changing the values of  $P$ ). As in [19], we remove this freedom by setting

$$e_{ij}(a, q) = e_{ij}(q, a) = h_i(q) = 0 \quad (33)$$

for all positions  $i, j$  and all amino acids  $a$ . Within this setting, each choice for the parameter values corresponds to a different outcome for  $P(a_1, \dots, a_L)$ . The parameters to be computed are therefore the couplings  $e_{ij}(a, b)$  and the fields  $h_i(a)$  with  $1 \leq a, b \leq q - 1$  only.

An different choice for the gauge is proposed in Text S1, Section 3, and leads to quantitatively equivalent predictions for the pattern structures and the contact map.

## Mean-field theory for determining the Hopfield-Potts patterns

The MaxEnt approach underlying DCA can be rephrased in a Bayesian framework. Assume the model to be given by Eq. (7), and assume the sequences in the MSA to be independently and identically sampled from  $P$ . The probability of the alignment for given model parameters (couplings and fields) is then given by

$$P[A|\{e_{ij}(a, b), h_i(a)\}] = \prod_{m=1}^M P(a_1^m, \dots, a_L^m) . \quad (34)$$

Plugging in Eq. (7) and defining the log-likelihood of the model parameters given the MSA  $A$ , we find

$$\begin{aligned} \mathcal{L}[\{e_{ij}(a, b), h_i(a)\}|A] &= \frac{1}{M} \log P[A|\{e_{ij}(a, b), h_i(a)\}] \\ &= \frac{1}{2} \sum_{i,j} \sum_{a,b} e_{ij}(a, b) f_{ij}(a, b) + \sum_{i,a} h_i(a) f_i(a) - \log \mathcal{Z}(\{e_{ij}(a, b), h_i(a)\}) \end{aligned} \quad (35)$$

One can readily see that the parameters  $\{e_{ij}(a, b), h_i(a)\}$  maximizing  $\mathcal{L}$  are solutions of Eqs. (4) and (5). The corresponding value for the maximum of  $\mathcal{L}$  coincides with the opposite of the entropy,  $-H[P]$ , for the MaxEnt distribution given by Eq. (7).

Following the study of the Ising model case ( $q = 2$ ) in [43], mean-field theory can be used to derive an approximate expression for the log-likelihood  $\mathcal{L}$  (35) when the couplings are chosen to obey Hopfield's prescription, Eq. (18). Calculations are presented in Text S1, Section 1. After optimization over the fields, we are left with the log-likelihood for the patterns only,

$$\begin{aligned} \mathcal{L}[\{\xi\}|A] &= \mathcal{L}_0 + \frac{1}{2L} \sum_{ij,ab} C_{ij}(a, b) \left( \sum_{\mu \leq p_+} \xi_i^{+, \mu}(a) \xi_j^{+, \mu}(b) - \sum_{\nu \leq p_-} \xi_i^{-, \nu}(a) \xi_j^{-, \nu}(b) \right) \\ &+ \frac{1}{2} \sum_{\mu \leq p_+} \log \left[ 1 - \frac{1}{L} \sum_{i,ab} \xi_i^{+, \mu}(a) C_{ii}(a, b) \xi_i^{+, \mu}(b) \right] + \frac{1}{2} \sum_{\nu \leq p_-} \log \left[ 1 + \frac{1}{L} \sum_{i,ab} \xi_i^{-, \nu}(a) C_{ii}(a, b) \xi_i^{-, \nu}(b) \right] \end{aligned} \quad (36)$$

where  $\mathcal{L}_0 = \sum_i \sum_{a=1}^q f_i(a) \log f_i(a)$ . So we find the trivial result that, for  $p = 0$  (no couplings), the log-likelihood is the negative of the sum of all single-column entropies,  $\mathcal{L}_0$ . The optimal patterns, *i.e.* those optimizing the log-likelihood  $\mathcal{L}$  are given by Eqs. (19) and (20). The total log-likelihood corresponding to this selection reads:

$$\mathcal{L}(p) = \mathcal{L}_0 + \sum_{\{\mu | \lambda_\mu \notin [\ell_-, \ell_+]\}} \Delta \mathcal{L}(\lambda_\mu) , \quad (37)$$

where function  $\Delta \mathcal{L}$  is defined in Eq. (23), and the bounds  $\ell_-$ ,  $\ell_+$  are defined in the Results Section.

The solution given in Eqs. (19) and (20) is defined up to a rotation in the pattern space, *i.e.* up to multiplication of all patterns with an indefinite orthogonal  $(p \times p)$ -matrix,  $\mathcal{O}$ , in  $O(p_+, p_-)$ . Indeed, the patterns  $\xi_i(a)$  and their rotated counterparts  $\hat{\xi}_i(a) = (\mathcal{O} \cdot \xi)_i(a)$  define the same set of couplings  $e_{ij}(a, b)$  through Eq. (18). Note that this invariance is specific to the Hopfield model, and should not be mistaken for the gauge invariance of the Potts model discussed in the Results Sections. We eliminate this arbitrariness according to the following procedure, detailed in Text S1: Our selection corresponds to the case where patterns are added one after the other, starting with the best possible single pattern, followed by the second best (orthogonal to the first one when single-site correlations  $C_{ii}(a, b)$  are factored out) etc.

## Contact prediction from couplings

Intuitively, residue position pairs with strong direct couplings are our best predictions for native contacts in the protein structure. To measure 'coupling strength', we need, however, to map the inferred  $q \times q$

coupling matrices  $e_{ij}$  onto a scalar parameter, for each  $1 \leq i < j \leq L$ . Whereas previous works on DCA have mainly used the so-called direct information [17, 19], it was recently observed that a different score actually improves the contact prediction starting from the same model parameters  $\{e_{ij}(a, b)\}$  [49]. To this end, we introduce the Frobenius norm

$$F_{ij} = \|e'_{ij}\|_2 = \sqrt{\sum_{a,b=1}^q \tilde{e}_{ij}(a, b)^2} \quad (38)$$

of the linearly transformed coupling matrices

$$\tilde{e}_{ij}(a, b) = e_{ij}(a, b) - e_{ij}(\cdot, b) - e_{ij}(a, \cdot) + e_{ij}(\cdot, \cdot) , \quad (39)$$

where ‘ $\cdot$ ’ denotes average over all amino acids and the gap in the concerned position. According to the above discussion, this corresponds to another gauge of the Hopfield-Potts model, more precisely to the gauge minimizing the Frobenius norm of each coupling matrix [17]. Further more, the norm is adjusted by an *average product correction* (APC) term, introduced in [16] to suppress effects from phylogenetic bias and insufficient sampling. Incorporating also this correction, we get our final scalar score:

$$F_{ij}^{APC} = F_{ij} - \frac{F_{\cdot j} F_{i \cdot}}{F_{\cdot \cdot}} , \quad (40)$$

where the ‘ $\cdot$ ’ now indicates a position average.

Sorting column pairs  $(i, j)$  by decreasing values of  $F_{ij}^{APC}$  calculated using standard mean-field DCA was shown to give accurate predictions for residue contacts in various proteins, *i.e.* in the case where all possible patterns are included ( $p = L(q-1)$ ) in Eq. (18). The Results Section shows how the performance in contact prediction varies when the number of patterns is  $p \ll L(q-1)$ .

Note that this criterion gives a coupling score to each pair of residue positions. The method itself does not provide a cutoff value for this score, below which predictions should not be considered any more. Results are therefore typically provided as parametric plots depending on the number of predicted contacts as a free parameter.

## Acknowledgments

We are grateful to R. Ranganathan and O. Rivoire for discussions.

## References

1. Pagani I, Liolios K, Jansson J, Chen I, Smirnova T, et al. (2012) The Genomes OnLine Database (GOLD) v.4: status of genomic and metagenomic projects and their associated metadata. *Nucleic Acids Res* 40: D571.
2. The Uniprot Consortium (2012) Reorganizing the protein space at the Universal Protein Resource (UniProt). *Nucleic Acids Res* 40: D71.
3. Punta M, Coghill PC, Eberhardt RY, Mistry J, Tate JG, et al. (2012) The Pfam protein families database. *Nucleic Acids Res* 40: D290.
4. Berman HM, Kleywegt GJ, Nakamura H, Markley JL (2012) The protein data bank at 40: Reflecting on the past to prepare for the future. *Structure* 20: 391 - 396.
5. Göbel U, Sander C, Schneider R, Valencia A (1994) Correlated mutations and residue contacts in proteins. *Proteins: Struct, Funct, Genet* 18: 309.
6. Lockless SW, Ranganathan R (1999) Evolutionarily conserved pathways of energetic connectivity in protein families. *Science* 286: 295-299.
7. de Juan D, Pazos F, Valencia A (2013) Emerging methods in protein co-evolution. *Nature Reviews Genetics* 14: 249-261.
8. Casari G, Sander C, Valencia A (1995) A method to predict functional residues in proteins. *Nature Structural Biology* 2: 171 -178.
9. Ortiz A, Hu W, Kolinski A, Skolnick J (1997) Method for low resolution prediction of small protein tertiary structure. *Pac Symp Biocomput* : 316 - 327.
10. Pazos F, Helmer-Citterich E, Ausiello G, Valencia A (1997) Correlated mutations contain information about protein- protein interaction. *J Mol Biol* 271: 511- 523.
11. Ortiz A, Kolinski A, Rotkiewicz P, Ilkowski B, Skolnick J (1999) Ab initio folding of proteins using restraints derived from evolutionary information. *PROTEINS: Structure, Function, and Genetics* 3: 177 - 185.
12. Lapedes AS, Giraud BG, Jarzynski C (2002) Using sequence alignments to predict protein structure and stability with high accuracy. LANL preprint .
13. Fodor AA, Aldrich RW (2004) Influence of conservation on calculations of amino acid covariance in multiple sequence alignments. *Proteins: Structure, Function, and Bioinformatics* 56: 211-221.
14. Socolich M, Lockless SW, Lee HL, Gardner K, Ranganathan R (2005) Evolutionary Information for Specifying a Protein Fold. *Nature* 437: 512-518.
15. Russ W, Lowery D, Mishra P, Yaffe M, Ranganathan R (2005) Natural-like Function in Artificial WW Domains. *Nature* 437: 579-583.
16. Dunn SD, Wahl LM, Gloor GB (2008) Mutual information without the influence of phylogeny or entropy dramatically improves residue contact prediction. *Bioinformatics* 24: 333.
17. Weigt M, White RA, Szurmant H, Hoch JA, Hwa T (2009) Identification of direct residue contacts in protein-protein interaction by message passing. *Proc Natl Acad Sci U S A* 106: 67.

18. Burger L, van Nimwegen E (2010) Disentangling Direct from Indirect Co-Evolution of Residues in Protein Alignments. *PLoS Comput Biol* 6: E1000633.
19. Morcos F, Pagnani A, Lunt B, Bertolino A, Marks DS, et al. (2011) Direct-coupling analysis of residue coevolution captures native contacts across many protein families. *Proc Natl Acad Sci U S A* 108: E1293.
20. Balakrishnan S, Kamisetty H, Carbonell JG, Lee SI, Langmead CJ (2011) Learning generative models for protein fold families. *Proteins: Struct, Funct, Bioinf* 79: 1061.
21. Jones DT, Buchan DWA, Cozzetto D, Pontil M (2012) PSICOV: precise structural contact prediction using sparse inverse covariance estimation on large multiple sequence alignments. *Bioinformatics* 28: 184.
22. Yeang C, Haussler (2007) Detecting coevolution in and among protein domains. *Plos Comput Biol* 31 (11): e211.
23. Jaynes ET (1957) Information Theory and Statistical Mechanics. *Physical Review Series II* 106: 620630.
24. Jaynes ET (1957) Information Theory and Statistical Mechanics II. *Physical Review Series II* 108: 171190.
25. Schneidman E, Berry M, Segev R, Bialek W (2006) Weak pairwise correlations imply strongly correlated network states in a neural population. *Nature* 440: 1007-1012.
26. Cocco S, Leibler S, Monasson R (2009) Neuronal couplings between retinal ganglion cells inferred by efficient inverse statistical physics methods. *Proc Natl Acad Sci U S A* 106: 14058-62.
27. Lezon TR, Banavar JR, Cieplak M, Maritan A, Fedoroff NV (2006) Using the principle of entropy maximization to infer genetic interaction networks from gene expression patterns. *Proc Nat Acad Sci* 103: 19033-19038.
28. Bialek W, Cavagna A, Giardina I, Mora T, Silvestri E, et al. (2012) Statistical mechanics for natural flocks of birds. *Proc Nat Acad Sci* .
29. Lapedes AS, Giraud BG, Liu L, Stormo GD (1999) Correlated mutations in models of protein sequences: Phylogenetic and structural effects. *Lecture Notes-Monograph Series: Statistics in Molecular Biology and Genetics* 33: 236-256.
30. Schug A, Weigt M, Onuchic JN, Hwa T, Szurmant H (2009) High-resolution protein complexes from integrating genomic information with molecular simulation. *Proc Natl Acad Sci USA* 106: 22124.
31. Marks DS, Colwell LJ, Sheridan RP, Hopf TA, Pagnani A, et al. (2011) Protein 3D Structure Computed from Evolutionary Sequence Variation. *PLoS ONE* 6: e28766.
32. Sadowski MI, Maksimiak K, Taylor WR (2011) Direct correlation analysis improves fold recognition. *Computational Biology and Chemistry* 35: 323 - 332.
33. Sulkowska JI, Morcos F, Weigt M, Hwa T, Onuchic JN (2012) Genomics-aided structure prediction. *Proc Natl Acad Sci* 109: 10340-10345.
34. Nugent T, Jones DT (2012) Accurate de novo structure prediction of large transmembrane protein domains using fragment-assembly and correlated mutation analysis. *Proceedings of the National Academy of Sciences* 109: E1540-E1547.

35. Hopf TA, Colwell LJ, Sheridan R, Rost B, Sander C, et al. (2012) Three-dimensional structures of membrane proteins from genomic sequencing. *Cell* 149: 1607-1621.
36. Dago AE, Schug A, Procaccini A, Hoch JA, Weigt M, et al. (2012) Structural basis of histidine kinase autophosphorylation deduced by integrating genomics, molecular dynamics, and mutagenesis. *Proc Natl Acad Sci USA* 109: 10148.
37. Marks DS, Hopf TA, Sander C (2012) Protein Structure Prediction from Sequence Variation. *Nature Biotechnology* 30: 1072-1080.
38. Taylor WR, Hamilton RS, Sadowski MI (2013) Prediction of contacts from correlated sequence substitutions. *Current Opinion in Structural Biology* : -.
39. Pearson K (1901) On lines and planes of closest fit to systems of points in space. *Philosophical Magazine* 2: 559-572.
40. Halabi N, Rivoire O, Leibler S, Ranganathan R (2009) Protein Sectors: Evolutionary Units of Three-Dimensional Structure. *Cell* 138: 774-786.
41. Bai Z, Silverstein JW (2009) Spectral analysis of large dimensional random matrices. London: Springer.
42. Hopfield JJ (1982) Neural networks and physical systems with emergent collective computational abilities. *Proc Natl Acad Sci* 79: 2554-2558.
43. Cocco S, Monasson R, Sessak V (2011) High-Dimensional Inference with the generalized Hopfield Model: Principal Component Analysis and Corrections. *Physical Review E* 83: 051123.
44. Wu FY (1982) The Potts Model. *Rev Mod Phys* 54: 235-268.
45. Rausell A, Juan D, Pazos F, Valencia A (2010) Protein interactions and ligand binding: From protein subfamilies to functional specificity. *Proceedings of the National Academy of Sciences* 107: 1995-2000.
46. Wlodawer A, Walter J, Huber R, Sjolín L (1984) Structure of bovine pancreatic trypsin inhibitor: Results of joint neutron and x-ray refinement of crystal form ii. *Journal of Molecular Biology* 180: 301 - 329.
47. Bent CJ, Isaacs NW, Mitchell TJ, Riboldi-Tunncliffe A (2004) Crystal Structure of the Response Regulator 02 Receiver Domain, the Essential YycF Two-Component System of *Streptococcus pneumoniae* in both Complexed and Native States. *J Bacteriol* 186: 2872-2879.
48. Pai EF, Krengel U, Petsko GA, Goody RS, Kabsch W, et al. (1990) Refined crystal structure of the triphosphate conformation of h-Ras p21 at 1.35 Å resolution: implications for the mechanism of gtp hydrolysis. *EMBO J* 9: 2351-2359.
49. Ekeberg M, Lövkvist C, Lan Y, Weigt M, E A (2013) Improved contact prediction in proteins: Using pseudolikelihoods to infer Potts models. *Phys Rev E* 87: 012707.
50. Ashkenazy H, Erez E, Martz E, Pupko T, Ben-Tal N (2010) ConSurf 2010: calculating evolutionary conservation in sequence and structure of proteins and nucleic acids. *Nucleic Acids Research* 38: W529-33.



## Tables

**Table 1.** Effect of a pattern with two non-zero and opposite components  $\xi_i(a) = -\xi_j(b)$ .

$a_i = a?$	$a_j = b?$	$\frac{S(a_1, \dots, a_L   \xi)}{\xi_i(a)^2}$	Favored by attractive pattern?	Favored by repulsive pattern?
NO	NO	0	NO	YES
YES	NO	1	YES	NO
NO	YES	1	YES	NO
YES	YES	0	NO	YES

**Table 2.** Effect of a pattern with two non-zero and equal components  $\xi_i(a) = \xi_j(b)$ .

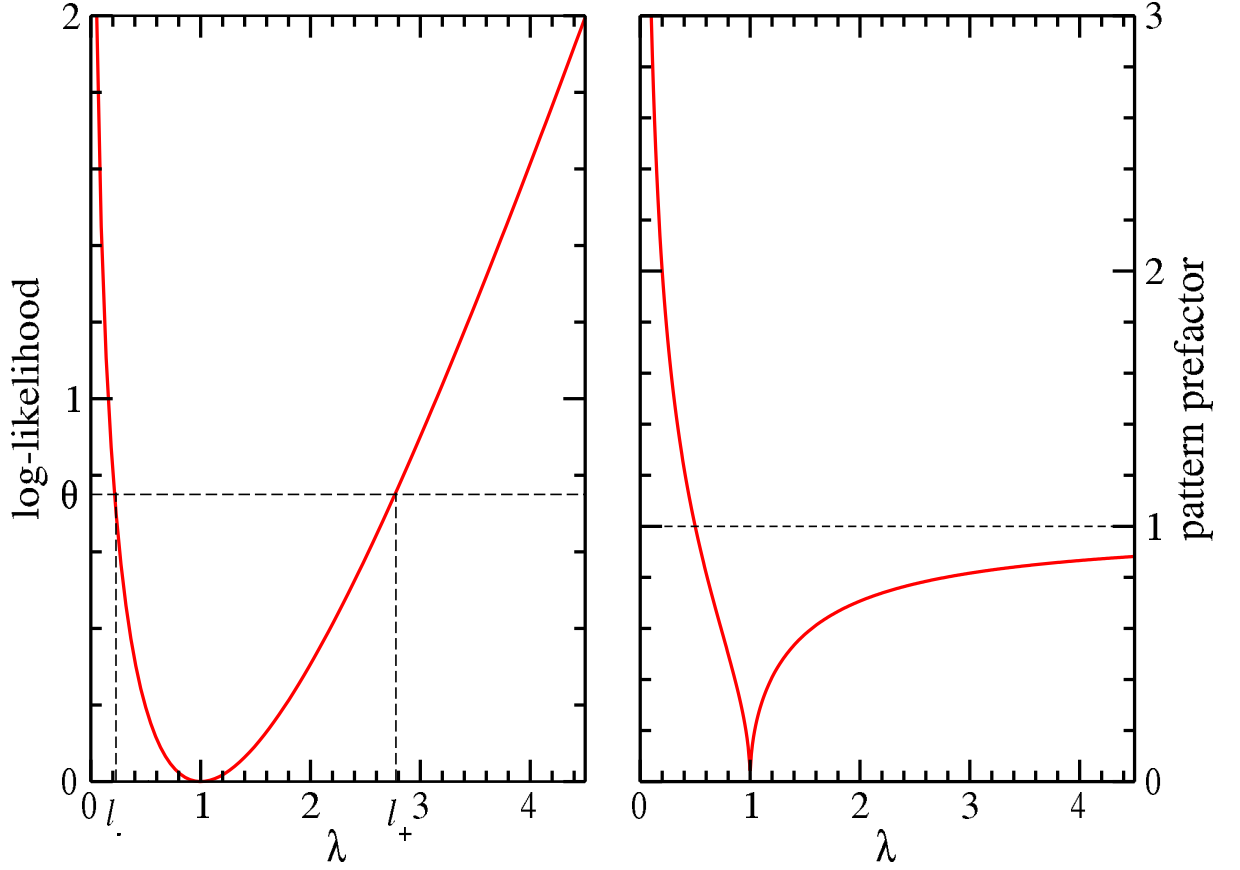
$a_i = a?$	$a_j = b?$	$\frac{S(a_1, \dots, a_L   \xi)}{\xi_i(a)^2}$	Favored by attractive pattern?	Favored by repulsive pattern?
NO	NO	0	NO	YES
YES	NO	1	NO	YES
NO	YES	1	NO	YES
YES	YES	4	YES	NO

## Supporting Information Files

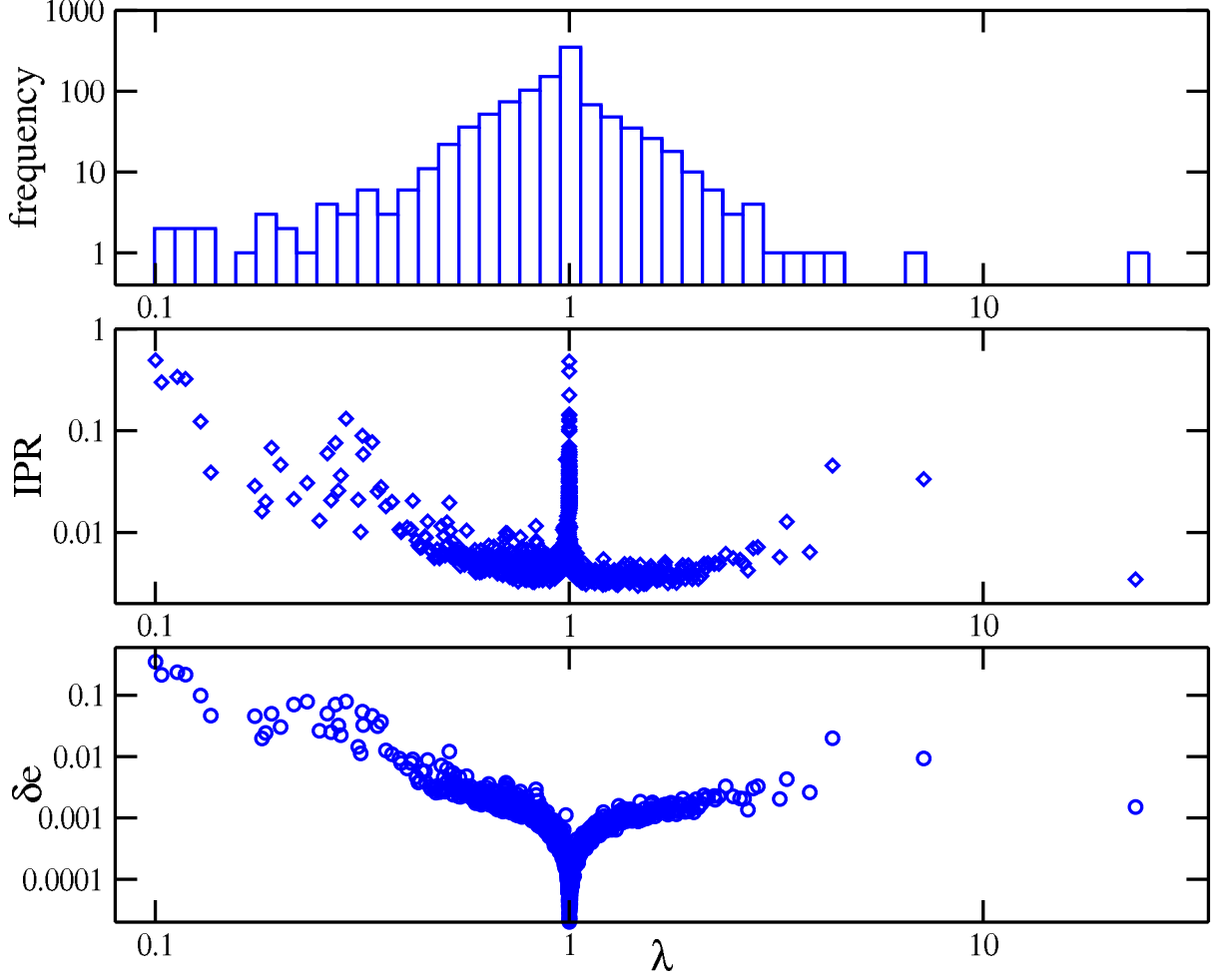
**Text S1.** Supporting Information for From principal component to direct coupling analysis of coevolution in proteins: Low-eigenvalue modes are needed for structure prediction.

**Code S1.** Matlab code for the Hopfield-Potts inference.

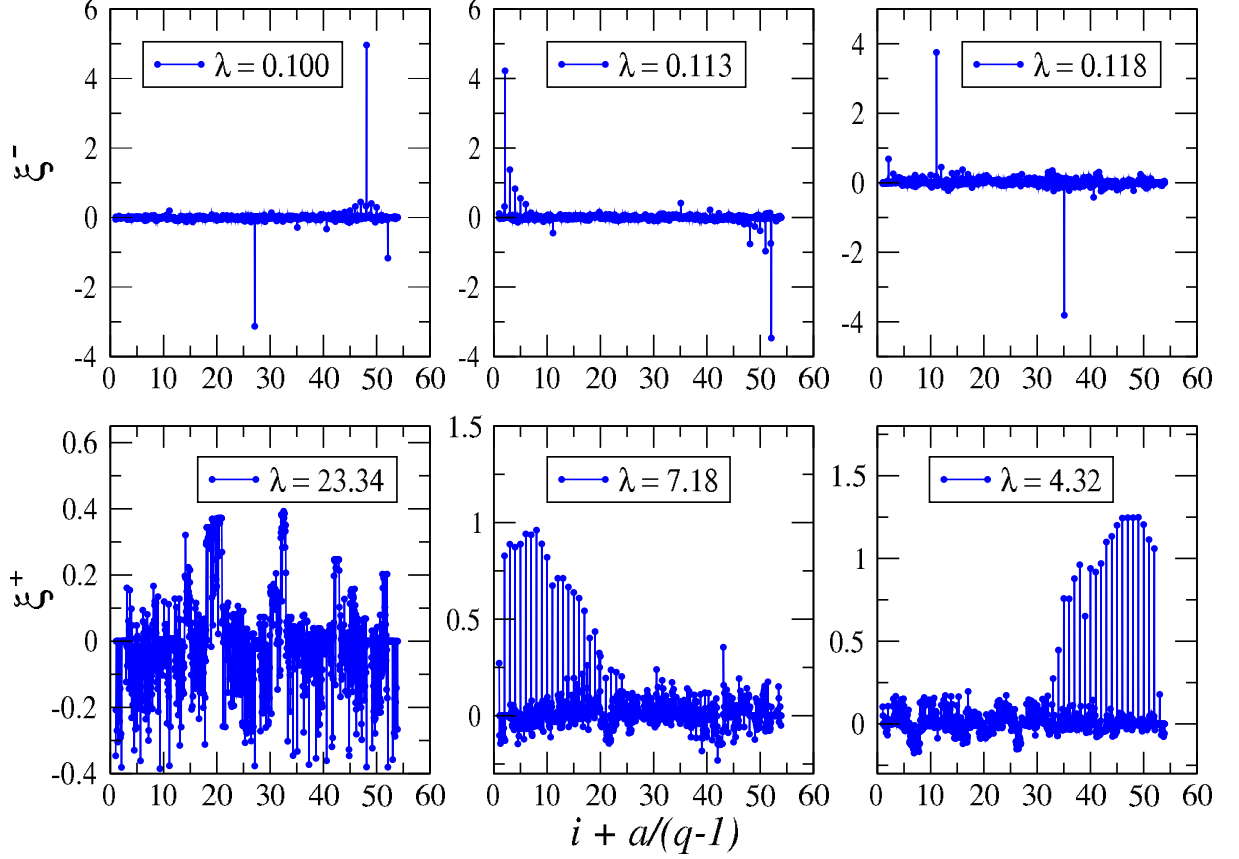
## Figures and Legends



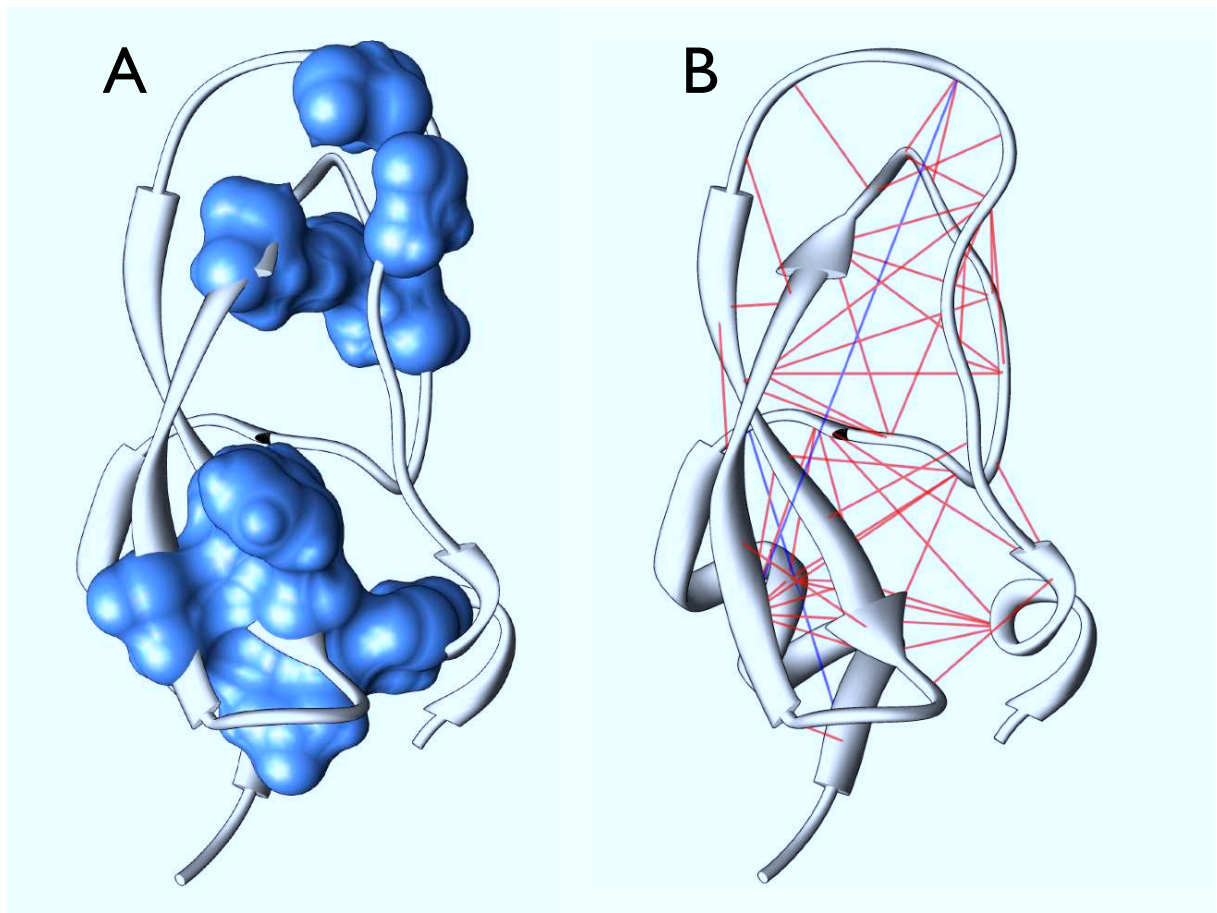
**Figure 1. Pattern selection by maximum likelihood and pattern prefactors:** (*left panel*) Contribution of patterns to the log-likelihood (full red line) as a function of the corresponding eigenvalues  $\lambda$  of the Pearson correlation matrix  $\Gamma$ . To select  $p$  patterns, a log-likelihood threshold  $\theta$  (dashed black line) has to be chosen such that there are exactly  $p$  patterns with  $\Delta\mathcal{L}(\lambda_\mu) > \theta$ . This corresponds to eigenvalues in the left and right tail of the spectrum of  $\Gamma$ . (*right panel*) Pattern prefactors  $|1 - \frac{1}{\lambda}|^{1/2}$  (full red line) as a function of the eigenvalue  $\lambda$ . Patterns corresponding to  $\lambda \simeq 1$  have essentially vanishing prefactors; patterns associated to large  $\lambda$  ( $\gg 1$ ) have prefactors smaller than 1 (dashed black line), while patterns corresponding to small  $\lambda$  ( $\ll 1$ ) have unbounded prefactors.



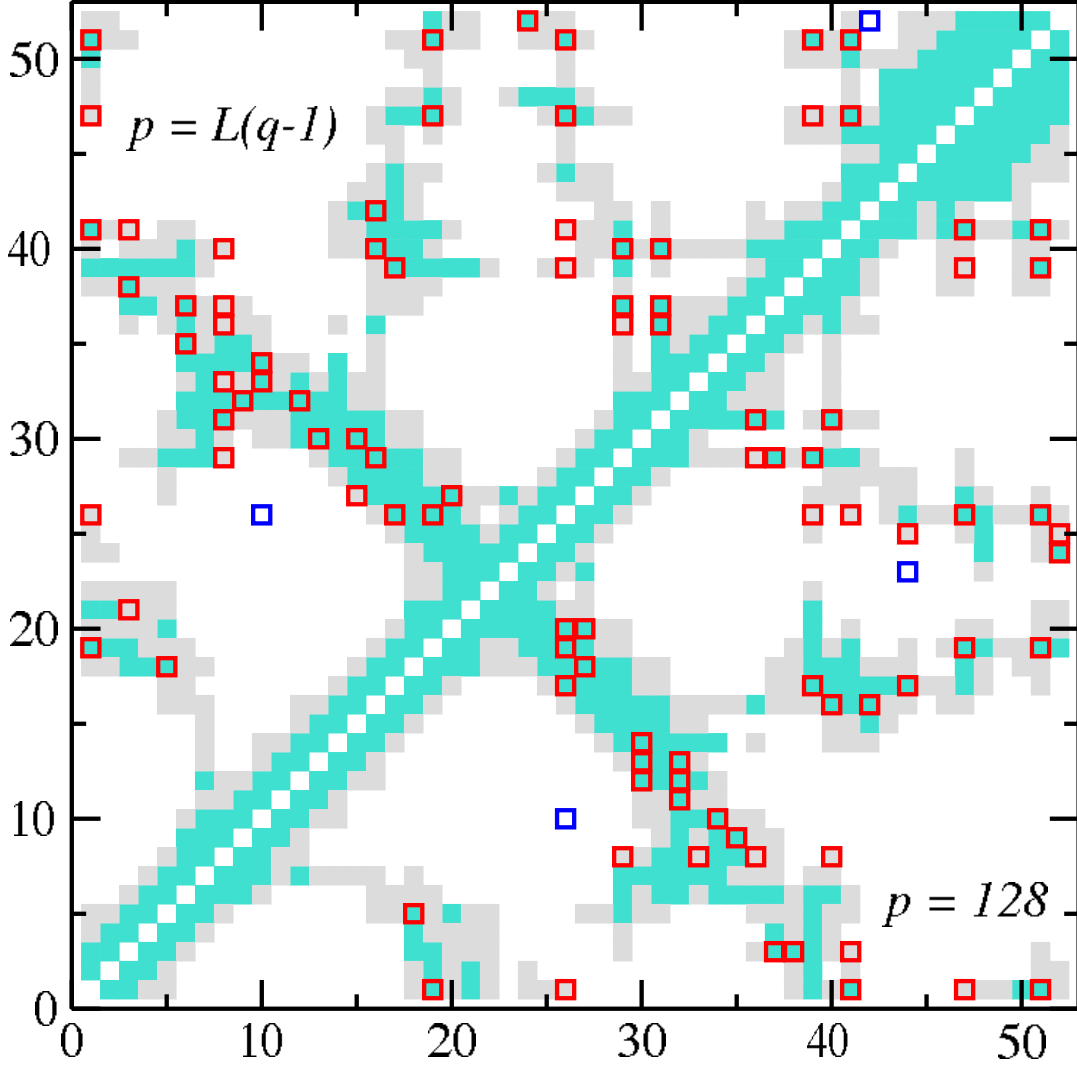
**Figure 2. Eigenvalues, localization and contributions to couplings for PF00014** (from top to bottom): (*top panel*) Spectral density as a function of the eigenvalues  $\lambda$ , note the existence of few very large eigenvalues, and a pronounced peak in  $\lambda = 1$ . (*middle panel*) Inverse participation ratio of the Hopfield patterns as a function of the corresponding eigenvalue  $\lambda$ . Large IPR characterizes the concentration of a pattern to few positions and amino acids. (*bottom panel*) Typical contribution  $\delta e$  to couplings due to each Hopfield pattern, defined in Eq. (26), as a function of the corresponding eigenvalue  $\lambda$ . Large contributions are mostly found for small eigenvalues, while patterns corresponding to  $\lambda \simeq 1$  do not contribute to couplings.



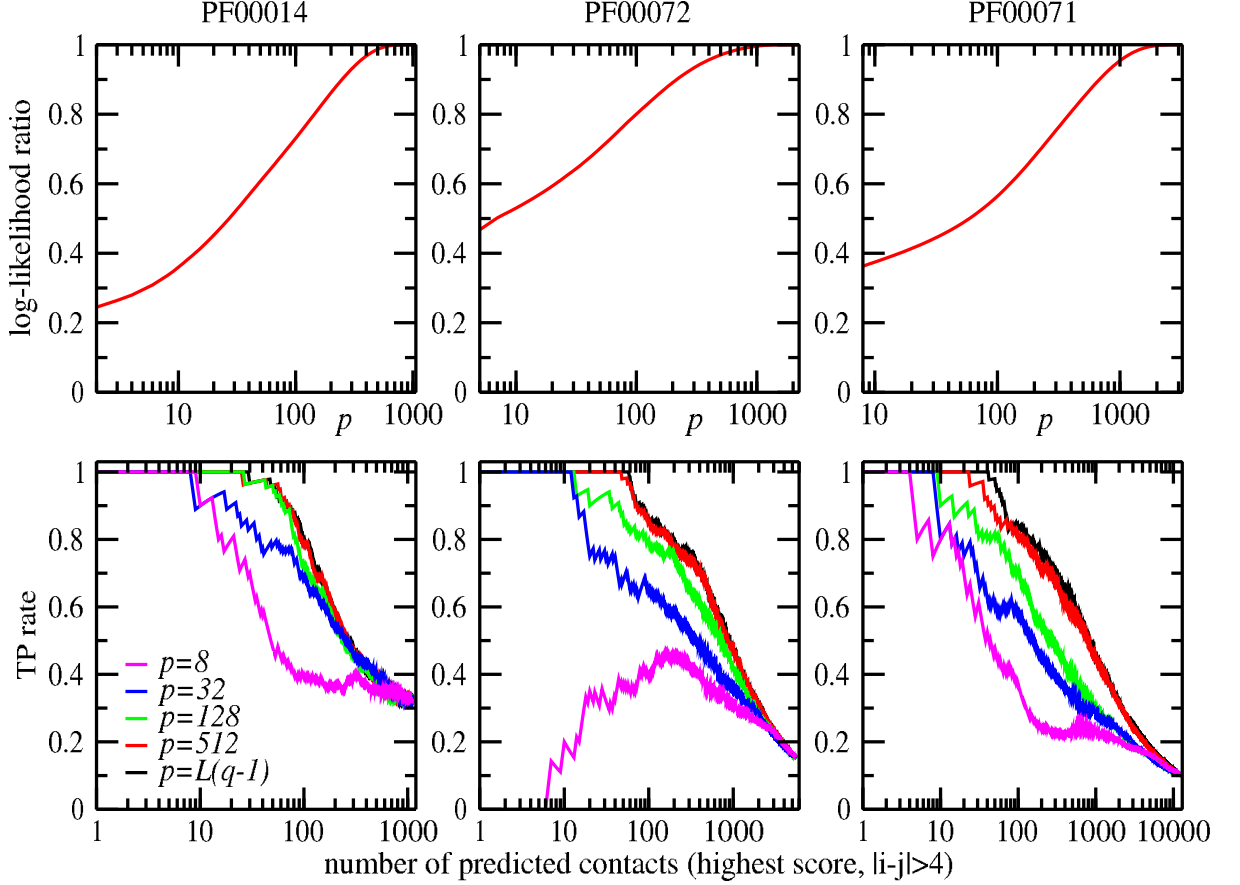
**Figure 3. Attractive and repulsive patterns for PF00014:** (*upper panels*) The most localized repulsive patterns (corresponding to the first, third and fourth smallest eigenvalues and inverse participation ratios 0.49, 0.34, 0.32 respectively) are strongly concentrated in pairs of positions. (*lower panels*) The most attractive patterns (corresponding to the three largest eigenvalues); the top pattern is extended, with inverse participation ratio 0.003, while the second and third patterns, with inverse participation ratios 0.033, 0.045 respectively, have essentially non-zero components over the gap symbols only which accumulate on the edges of the sequence. Note the  $x$ -coordinates  $i + a/(q - 1)$ ; its integer part is the site index,  $i$ , and the fractional part multiplied by  $q - 1$  is the residue value,  $a$ .



**Figure 4. The principal component and predicted contacts visualized on the 3D structure of the trypsin inhibitor protein domain PF00014.** (A) The 10 positions (residue ID 5,12,14,22,23,30,35,40,51,55) of largest entries in the most attractive Hopfield pattern (largest eigenvalue of  $\Gamma$ , corresponding to the principal component) are shown in blue, they correspond also to very conserved sites. Note that, while they are distant along the protein backbone, they cluster into spatially connected components in the folded protein. (B) The 50 residue pairs with strongest couplings (ranked according to the Frobenius norms Eq. (40), with at least 5 positions separation along the backbone, are connected by lines. Only two out of these pairs are not in contact (blue links), all other 48 are thus true-positive contact predictions (red links). Many contacts link pairs of not conserved positions. Note that links are drawn between C-alpha atoms, whereas contacts are defined via minimal all-atom distances, making some red lines to appear rather long even if corresponding to native contacts.

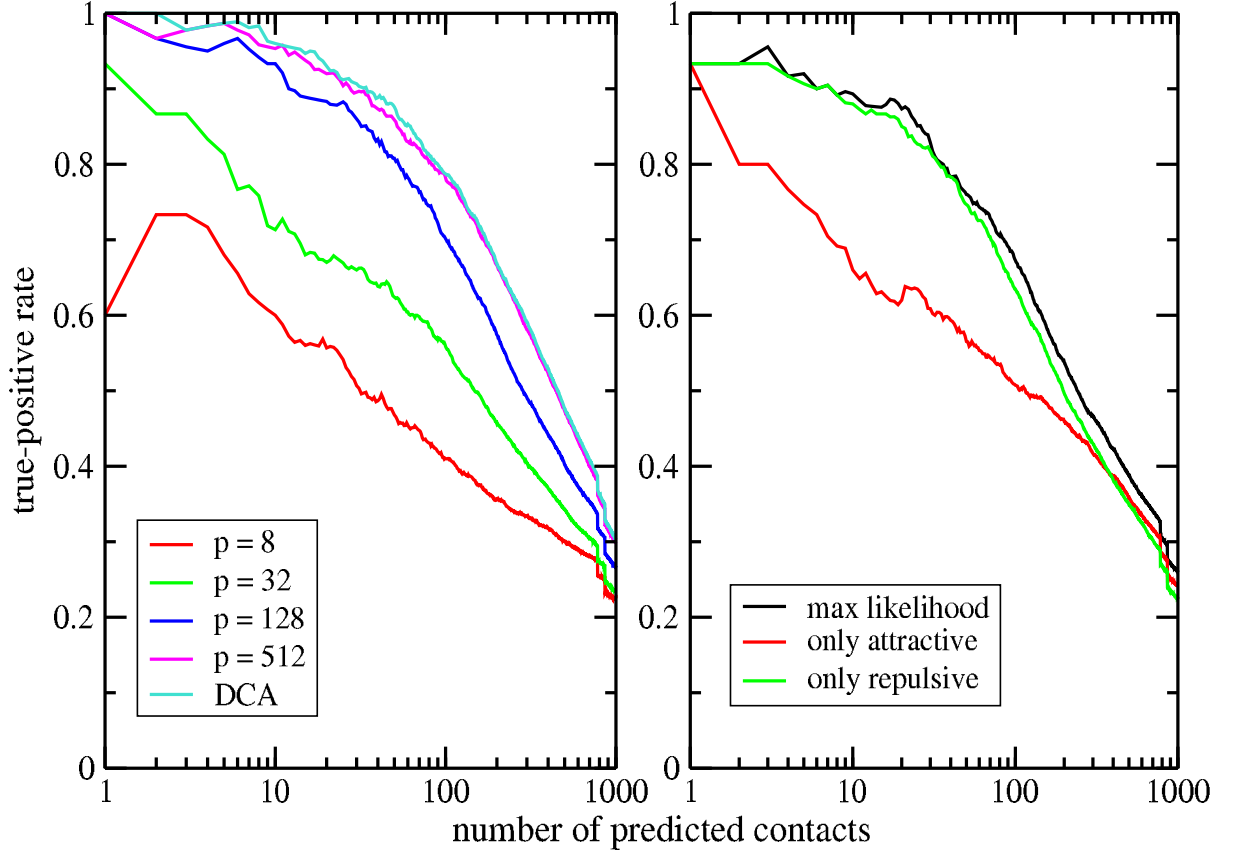


**Figure 5. Contact map for the PF00014 family.** Filled squares represent the native contact map on the 3D fold (PDB 5pti, with turquoise squares signaling all-atom distances below 5Å, and grey ones distances between 5Å and 8Å). The 50 top predicted contacts with minimal separation of 5 positions along the sequence ( $|i - j| \geq 5$ ) are shown with empty squares: true-positive predictions (distance < 8Å) are colored in red, and false-positive predictions in blue. Predictions are made with the Hopfield-Potts model with  $p = 128$  patterns (bottom right corner) and with  $p = L(q - 1) = 1060$  patterns (DCA, top left corner). For both values of  $p$  there are 48 true-positive and 2 false-positive predictions.

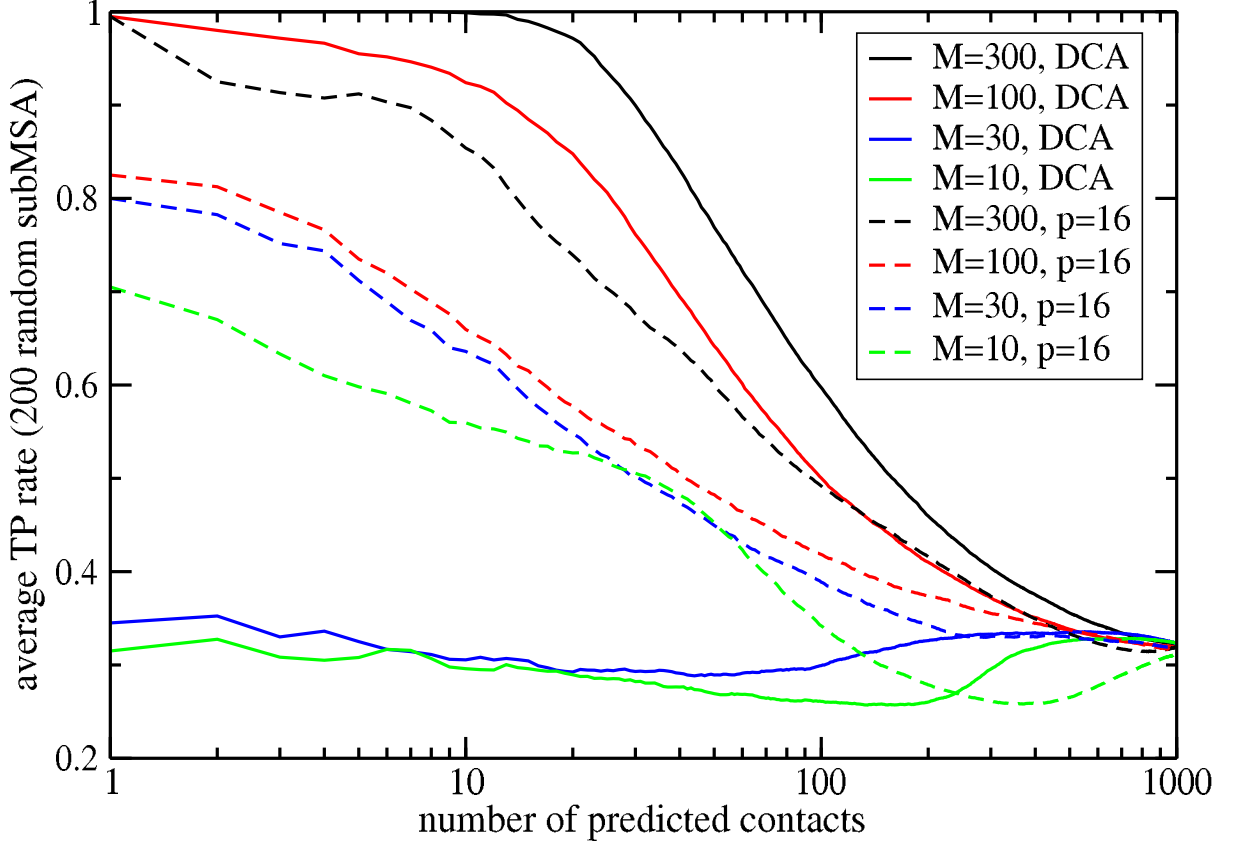


**Figure 6. Contact predictions for the three considered protein families.** The upper panels show the fraction of the interaction-based contribution to the log-likelihood of the model given the MSA, defined as the ratio of the log-likelihood with  $p$  selected patterns over the maximal log-likelihood obtained by including all  $L(q-1)$  patterns, as a function of the number  $p$  of selected patterns, it reaches one for  $p = (q-1)L$  corresponding to the Potts model used in DCA. The lower panels show the TP rates as a function of the predicted residue contacts, for various numbers  $p$  of selected patterns, where selection was done using the maximum-likelihood criterion.  $p = (q-1)L$  gives the contact predictions obtained by DCA approach. Only non-trivial contacts between sites  $i, j$  such that  $|i-j| > 4$  are considered in the calculation of the TP rate.





**Figure 7. Contact predictions across 15 protein families.** (*left panel*) TP rates for the contact prediction with variable numbers  $p$  of Hopfield-Potts patterns, averaged over 15 distinct protein families. (*right panel*) TP rates for the contact prediction using only the repulsive (green line) resp. attractive (red line) patterns, which are contained in the  $p = 100$  most likely patterns (black line), averaged over 15 protein families. It becomes obvious that the contact prediction remains almost unchanged when only the subset of repulsive patterns is used, whereas it drops substantially by keeping only attractive patterns.



**Figure 8. Noise reduction due to pattern selection in reduced data sets.** (*full lines*) TP rates of mean-field DCA for sub-MSAs of family PF00014 with  $M = 10, 30, 100, 300$  sequences; each curve is averaged over 200 randomly selected sub-alignments. Whereas for  $M = 100$  and  $M = 300$  the accuracy of the first predictions is close to one, mean-field DCA does not extract any reasonable signal for  $M = 10$  and  $M = 30$ . (*dashed lines*) The same sub-MSA are analyzed with the Hopfield-Potts model using  $p = 16$  patterns (maximum-likelihood selection). Whereas this selection reduces the accuracy for  $M \geq 100$ , it results in increased TP rates for  $M \leq 30$ . Dimensional reduction by pattern selection has lead to an efficient noise reduction.