

## MIT Open Access Articles

*On the Entropy of Protein Families*

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

**Citation:** Barton, John P., Arup K. Chakraborty, Simona Cocco, Hugo Jacquin, and Rémi Monasson. "On the Entropy of Protein Families." *Journal of Statistical Physics*, vol. 162, no. 5, January 2016, pp.1267–1293.

**As Published:** <http://dx.doi.org/10.1007/s10955-015-1441-4>

**Publisher:** Springer-Verlag

**Persistent URL:** <http://hdl.handle.net/1721.1/104853>

**Version:** Author's final manuscript: final author's manuscript post peer review, without publisher's formatting or copy editing

**Terms of use:** Creative Commons Attribution-Noncommercial-Share Alike



## On the Entropy of Protein Families

John P. Barton<sup>1,2,3,6</sup> · Arup K. Chakraborty<sup>1,2,3,4,5,6</sup> ·  
Simona Cocco<sup>7</sup> · Hugo Jacquin<sup>7</sup> · Rémi Monasson<sup>8</sup>

Received: 7 July 2015 / Accepted: 22 December 2015 / Published online: 13 January 2016  
© Springer Science+Business Media New York 2016

**Abstract** Proteins are essential components of living systems, capable of performing a huge variety of tasks at the molecular level, such as recognition, signalling, copy, transport, ... The protein sequences realizing a given function may largely vary across organisms, giving rise to a protein family. Here, we estimate the entropy of those families based on different approaches, including Hidden Markov Models used for protein databases and inferred statistical models reproducing the low-order (1- and 2-point) statistics of multi-sequence alignments. We also compute the entropic cost, that is, the loss in entropy resulting from a constraint acting on the protein, such as the mutation of one particular amino-acid on a specific site, and relate this notion to the escape probability of the HIV virus. The case of lattice proteins, for which the entropy can be computed exactly, allows us to provide another illustration of the concept of cost, due to the competition of different folds. The relevance of the entropy in relation to directed evolution experiments is stressed.

---

✉ Simona Cocco  
cocco@lps.ens.fr

Rémi Monasson  
monasson@lpt.ens.fr

<sup>1</sup> Ragon Institute of MGH, MIT & Harvard, Cambridge, MA, USA

<sup>2</sup> Department of Chemical Engineering, Massachusetts Institute of Technology, Cambridge, MA, USA

<sup>3</sup> Department of Physics, Massachusetts Institute of Technology, Cambridge, MA, USA

<sup>4</sup> Department of Chemistry, Massachusetts Institute of Technology, Cambridge, MA, USA

<sup>5</sup> Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge, MA, USA

<sup>6</sup> Institute for Medical Engineering & Science, Massachusetts Institute of Technology, Cambridge, MA, USA

<sup>7</sup> Laboratoire de Physique Statistique de l'ENS, UMR 8550, associé au CNRS et à l'Université P&M. Curie, 24 rue Lhomond, 75005 Paris, France

<sup>8</sup> Laboratoire de Physique Théorique de l'ENS, UMR 8549, associé au CNRS et à l'Université P&M. Curie, 24 rue Lhomond, 75005 Paris, France

**Keywords** Statistical inference · Entropy · Fitness landscape · Genomics · Hidden Markov models · Covariation · HIV virus

## 1 Introduction

Characterizing the statistical properties of a family of homologous protein sequences is a problem of fundamental importance in genomics. It is well known for instance that the frequencies of amino acids vary substantially along the sequence from site to site, as residues are generally strongly conserved in the protein cores and in binding pockets [1,2]. As the number of available sequences has hugely increased over the last years, higher-order statistical properties may now be accurately estimated. Correlations between pairs of residues in the sequence are known to reflect structural, functional, or phylogenetic constraints acting on the protein sequences [3–6]. Conservation, pairwise correlations, and possibly higher-order statistical constraints limit the number of putative proteins in a given family. It is of fundamental interest from an evolutionary point of view to be able to quantitatively estimate the diversity of proteins corresponding to a given family, and by extension, sharing the same biological function. The present paper, based on a variety of modeling approaches and of sequence data, is a modest attempt in this direction.

A natural way to quantify the diversity of proteins with the same function is through the Gibbs-Shannon entropy of the distribution of sequences in the corresponding protein family. Qualitatively, this entropy can be thought of as the logarithm of the number of sequences in the family, though there need not be a sharp divide between functional sequences (those belonging to the family) and dysfunctional ones. In the course of evolution, Nature has sampled many protein sequences across largely diverse organisms. Natural selection weeds out dysfunctional sequences, while amplifying those that perform their function efficiently. Current databases such as UniProt or PFAM [7–9] give us a sample of the diversity of those good sequences, i.e., ones that ensure large fitnesses to the organisms compared to other protein sequences. However, despite massive sequencing efforts the number of available sequences is likely to be incredibly small compared to all possible sequences with high fitnesses. That is, we only observe a subset of the true distribution of functional sequences. We are thus faced with the difficult task of estimating the entropy of a probability distribution over the sequence space in the presence of dramatic undersampling. This is only possible under strong assumptions on the smoothness of the sequence distribution. Here we explore several different approaches for estimating the entropy for protein families, given a limited sampling of sequences.

One popular approach in this context is to consider Maximum Entropy distributions [10–14] reproducing low-order statistics of the amino acids in the sequence databases, generally the single-site and pairwise frequencies. The corresponding distributions are smooth in the sense that they correspond to the Gibbs distributions associated to Potts Hamiltonians with local fields and pairwise couplings only. The difficulty in this approach is to compute those interaction parameters from the sequence statistics, and the corresponding entropies. In the present paper, we will resort to an approximate method allowing us to access those quantities, the Adaptive Cluster Expansion developed in [15,16], and will apply it to real protein data (homologous protein families in Sect. 3 and human immunodeficiency virus (HIV) sequence data in Sect. 4) and to synthetic, lattice-based protein models (Sect. 5) [17]. This method estimates the cross-entropy between the inferred Potts model and the data, which is equal to the entropy of the Potts model that reproduces the desired statistics from the data. In addition to giving us access to absolute estimates of the entropies of the protein families, our approach allows us to compute changes of entropies related to additional constraints acting

on the proteins. To illustrate this concept in the case of HIV, we will compute the variation in entropy as one amino acid is fixed to its consensus value. The loss in entropy, or entropy cost associated to this local constraint, is naturally related to the escape probability of a pathogen (virus or bacterium) from a drug or an immune attack. The latter can force mutations on one or multiple sites, and largely decreases the availability of putative escaping sequences. Another illustration will be provided by lattice-based proteins, where we will measure the decrease in entropy of a protein family, defined as the set of sequences folding properly into a given structure, resulting from the introduction of competing structures, i.e., alternative folding conformations.

We also estimate the entropy of the protein families in PFAM using their associated Hidden Markov Models (HMM) [1]. HMM define protein profiles, which are used to classify the families and answer sequence queries. HMM are, to some extent, similar to Maximum Entropy models reproducing 1-point statistics only, that is, to non-interacting Potts models with local fields only. However, HMM are capable of handling sequences of any length through the insertion of extra amino acids (for longer sequences than the length of the profile) or of gaps (for shorter sequences). As calculating exactly the value of the entropy of HMM models is generally a hard task, we will establish some bounds and approximations to this value in Sect. 2.

Last of all, in Sect. 6, we summarize our findings, and compare the values of the entropies found with our different approaches to previous estimates in the literature [18]. We comment in particular the possible relevance of our results for directed evolution experiments, where protein sequences are evolved and selected *in vitro*, starting from a pool of random sequences.

## 2 Wide-Scale Analysis of the Entropy of HMM Profiles Across the PFAM Database

### 2.1 Formalism for Hidden Markov Models

Hidden Markov Models are routinely used to define protein family profiles in databases, such as PFAM [9]. The underlying principle for HMM is the existence of hidden states, which condition the set of symbols (amino acids or gaps) composing the sequence. Briefly speaking, an HMM jumps from one hidden state  $\sigma$  to another state  $\tau$  in a sequential and stochastic way, depending on a set of transition rates. After each transition to a new hidden state, a symbol  $A$  may be produced, with an emission probability depending on the hidden state, and added to the sequence. When the last hidden state is reached the sequence is complete. A detailed description of HMM profiles can be found in [1], Chapter 5. Hereafter we briefly expose their salient aspects and introduce some notations.

In an HMM, for a profile of length  $N$ , the number of hidden states relevant to our purpose is  $N_s = 3N + 1$ . The initial and final states are denoted by, respectively,  $B$  and  $E$ . In between  $B$  and  $E$  the model includes  $N$  match states, denoted by  $M_j$ , with  $j = 1, 2, \dots, N$ ;  $N - 1$  insertion states  $I_j$ , with  $j = 1, 2, \dots, N - 1$ ;  $N$  deletion states  $D_j$ , with  $j = 1, 2, \dots, N$ . Amino-acid symbols are emitted by the  $I$  and  $M$  states. Insertion states  $I_j$  allow for the emission of excess symbols and to produce sequences longer than the profile length. Match states  $M_j$  emit symbols with probabilities dependent on the position ( $j$ ), and reflect the pattern of amino acid conservation along the profile. Deletion states  $D_j$  represent a gap in the sequence, i.e., the lack of correspondence to the site  $j$  of the profile. Note that the number of insertion states,  $N - 1$ , is different from the one ( $= N + 1$ ) in the description of HMMs in [1], Chapter 5, Fig. 5.2; the reason is that our definition of HMMs corresponds to the one of

the Matlab Bioinformatics toolbox we use to download the profiles from the PFAM database, and does not consider insertion states associated to the  $B$  and  $E$  states.

An HMM is fully defined by the transition rate matrix  $\mathcal{T}$ , of size  $N_s \times N_s$ , which gives the probability of jumping from any hidden state  $\sigma$  to another  $\tau$ , and by the emission matrix  $\mathcal{E}$ , which gives the probability of emitting symbol  $A$  given the hidden state  $\sigma$ . If the hidden state is an insertion or match state,  $\sigma = I_j$  or  $M_j$ , any of the 20 amino acids  $A$  may be emitted with probability  $\mathcal{E}(A|\sigma)$ ; if the hidden state is a deletion state,  $\sigma = D_j$ , the gap symbol is emitted with probability unity. The weight of the sequence  $\mathbf{A} = (A_1, A_2, \dots, A_L)$  with  $L$  emitted symbols is given by

$$P(\mathbf{A}; L) = \sum_{\sigma=(\sigma_1, \sigma_2, \dots, \sigma_L)} \mathcal{T}(B \rightarrow \sigma_1) \prod_{\ell=1}^L \left[ \mathcal{T}(\sigma_\ell \rightarrow \sigma_{\ell+1}) \mathcal{E}(A_\ell | \sigma_\ell) \right], \tag{1}$$

where we have defined  $\sigma_{L+1} \equiv E$ . The sum over all sequences  $\mathbf{A}$  of  $P(\mathbf{A}; L)$  is the probability to reach  $E$  from  $B$  through a path of length  $L$  across the hidden states; this sum, denoted by  $P(L)$ , is a priori smaller than unity, e.g., if the path length  $L$  is smaller than the model size  $N$  and  $E$  can not be reached from  $B$ . In practice, however,  $P(L)$  converges to unity as soon as  $L$  exceeds  $N$ , see below. Our goal is then to compute the entropy of the HMM,

$$S_1(L) = - \sum_{\mathbf{A}} P(\mathbf{A}; L) \log P(\mathbf{A}; L), \tag{2}$$

which is a function of the matrices  $\mathcal{T}, \mathcal{E}$  only. An exact calculation of  $S_1(L)$  is very difficult (see below), and we will instead compute the lower bound to the entropy,

$$S_1(L) > S_2(L), \tag{3}$$

and the approximation

$$S_1(L) \simeq 2 S_2(L) - S_3(L), \tag{4}$$

based on the Renyi entropies  $S_q(L)$ :

$$S_q(L) = \frac{1}{1-q} \log \left[ \sum_{\mathbf{A}} P(\mathbf{A}; L)^q \right]. \tag{5}$$

Note that  $2 S_2(L) - S_3(L)$  is not guaranteed to be a lower bound to  $S_1(L)$ , but is generally a closer estimate of  $S_1(L)$  than the guaranteed lower bound  $S_2(L)$ .

We now turn to the computation of  $S_q(L)$ , where  $q$  is integer valued. According to (1), we have

$$\begin{aligned} \sum_{\mathbf{A}} P(\mathbf{A}; L)^q &= \sum_{\sigma^{(1)}, \sigma^{(2)}, \dots, \sigma^{(q)}} \prod_{m=1}^q \mathcal{T}(B \rightarrow \sigma_1^{(m)}) \prod_{\ell=1}^L \left[ \prod_{m=1}^q \mathcal{T}(\sigma_\ell^{(m)} \rightarrow \sigma_{\ell+1}^{(m)}) \right. \\ &\quad \left. \times R(\sigma_\ell^{(1)}, \sigma_\ell^{(2)}, \dots, \sigma_\ell^{(q)}) \right], \end{aligned} \tag{6}$$

where we have defined  $\sigma_{L+1}^{(1)} = \sigma_{L+1}^{(2)} = \dots = \sigma_{L+1}^{(q)} = E$ , and

$$R(\sigma^{(1)}, \sigma^{(2)}, \dots, \sigma^{(q)}) = \sum_A \mathcal{E}(A|\sigma^{(1)}) \mathcal{E}(A|\sigma^{(2)}) \dots \mathcal{E}(A|\sigma^{(q)}), \tag{7}$$

for any set of  $q$  states  $\sigma^{(1)}, \sigma^{(2)}, \dots, \sigma^{(q)}$ . We introduce the indices  $\hat{\sigma} = (\sigma^{(1)}, \sigma^{(2)}, \dots, \sigma^{(q)})$  to label the elements of the  $(N_s)^q \times (N_s)^q$ -dimensional effective transition rate matrix  $\mathcal{M}_q$ :

$$\mathcal{M}_q(\hat{\sigma} \rightarrow \hat{\tau}) = \prod_{m=1}^q \mathcal{T}(\sigma^{(m)} \rightarrow \tau^{(m)}) \times \begin{cases} 1 & \text{if } \hat{\sigma} = \hat{B} \text{ or } \hat{\tau} = \hat{E}, \\ R(\sigma^{(1)}, \sigma^{(2)}, \dots, \sigma^{(q)}) & \text{otherwise,} \end{cases} \tag{8}$$

where  $\hat{B} = (B, B, \dots, B)$ ,  $\hat{E} = (E, E, \dots, E)$ . Using those notations, we write

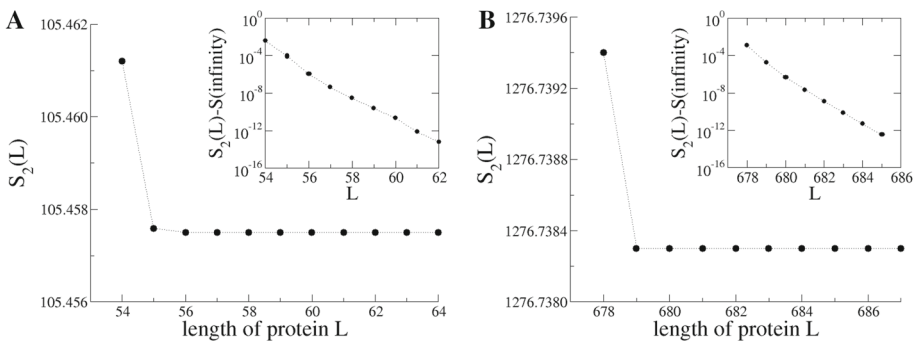
$$\sum_{\mathbf{A}} P(\mathbf{A}; L)^q = \mathcal{M}_q^L(\hat{B} \rightarrow \hat{E}) \quad \text{and} \quad S_q(L) = \frac{1}{1-q} \log \mathcal{M}_q^L(\hat{B} \rightarrow \hat{E}). \tag{9}$$

where  $\mathcal{M}_q^L$  denotes the  $L$ th-matrix power of  $\mathcal{M}_q$ . This formula explains why the calculation of  $S_q(L)$ , with  $q \geq 2$  is easier than the calculation of  $S_1(L)$ . Indeed, the hard computational step in the calculation of the entropy is obviously the summation over the enormous set of sequences  $\mathbf{A}$ , whose size grows exponentially with the length  $L$ . For integer values of  $q \geq 2$ , the summation can be split in  $L$  independent sums over the symbols  $A_\ell$ , which define the effective matrix  $R$  in (7). The Renyi entropies  $S_q$  (with  $q \geq 2$ ) can then be computed in a time growing linearly with  $L$  (and not exponentially) from the knowledge of the  $L$ th power of the transition matrix  $\mathcal{M}_q$  in (9). Unfortunately the size of  $\mathcal{M}_q$  grows exponentially with  $q$ , and this trick is limited to small values of  $q$ .

The formulas above were implemented in a Matlab routine. The rate matrix  $\mathcal{T}$  and the emission matrix  $\mathcal{E}$  were downloaded from PFAM profiles, and used to compute the  $\mathcal{M}_2$  and  $\mathcal{M}_3$  matrices. As the size of the latter grows as the cubic power of the number  $N_s$  of hidden states the computation of the Renyi entropy  $S_3(L)$  was done for moderate profile length only, i.e., in practice  $N \leq 100$ .

### 2.2 Convergence with $L$

We first study how our low bound for the HMM entropy depends on the length  $L$  of ‘emitted’ proteins. We plot in Fig. 1 the value of the Renyi entropy  $S_2(L)$  as function of  $L$  for two families: PF00014, a trypsin inhibitor, also studied in Sect. 3.2 and PF00063, a myosin-head protein. Those two families were chosen for the very different values of the lengths of their profiles:  $N = 53$  for PF00014 and  $N = 677$  for PF00063. The entropy  $S_2(L)$  is equal to minus infinity as long as  $L \leq N$ . The reason is that, in PFAM HMMs, the probabilities of transitions from any Match state  $M_j$  (with  $j < N$ ) to the end state  $E$  are zero.  $E$  can therefore not be reached in less than  $L = N + 1$  steps. The value of the entropy  $S_2(L = N + 1)$  corresponds



**Fig. 1** Lower bounds  $S_2(L)$  to the entropy of the Hidden Markov Models for families PF00014 (a) and PF00063 (b) as functions of the length of the proteins emitted by the models. For both families, the entropy has a finite value when  $L$  exceeds the profile length (53 for PF00014 and 677 for PF00063). *Insets* difference between  $S_2(L)$  and its asymptotic value versus  $L$

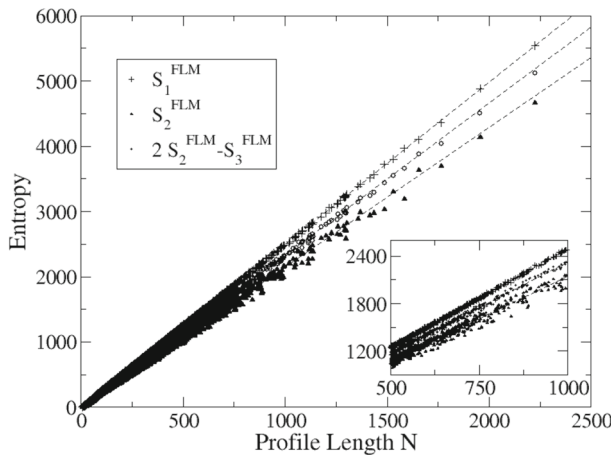
to the shortest transition path (through the Match or Deletion states) connecting  $B$  to  $E$ . As  $L$  increases, the probabilities of more and more processes (including self-transitions of Insertion states onto themselves, which have low but non-zero probabilities [1]) are collected, and the sum of the squared probabilities increases, which makes  $S_2$  decrease. Note that, once the state  $E$  is reached the system enters an infinite loop (through the transition  $E \rightarrow E$ ) and  $S_2(L)$  does include all the contributions coming from paths connecting  $B$  to  $E$  with length shorter or equal to  $L$  (by convention, in the calculation of  $R$  in (7), we consider that  $E$  emits empty symbols). We see that the entropy reaches an asymptotic plateau very quickly as  $L$  increases above  $N + 1$  (Inset of Fig. 1). In practice we choose  $L = 1.2 \times N$  to be sure that the convergence has been reached, and all paths from  $B$  to  $E$  have been taken into account ( $P(L) = 1$ ). A similar behaviour is observed for the Renyi entropy of order 3 as a function of the protein length  $L$  (not shown). To lighten notations we write in the following  $S_q$  for our estimate of the asymptotic entropy  $S_q(L \rightarrow \infty)$ .

### 2.3 Fixed-Length Model Built from HMM

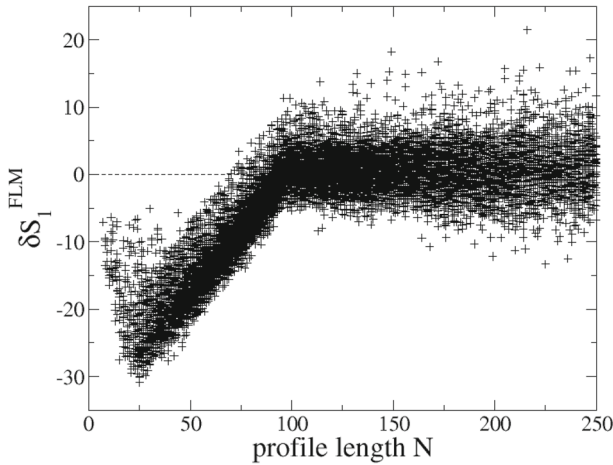
As a first step we ignore the possibility of insertion and deletion. The resulting simplified HMM model consists of the  $N$  Match states, visited one after the other in a sequential way. On each Match state  $M_j$ , an amino acid  $A_j$  is emitted according to the local probability of emission  $\mathcal{E}$ . In this simple Fixed-Length Model (FLM), symbols are emitted independently of each other. The entropy of the distribution of sequences produced with the FLM is therefore

$$S_1^{FLM} = - \sum_{j=1}^N \sum_{A_j} \mathcal{E}(A_j|M_j) \log \mathcal{E}(A_j|M_j). \tag{10}$$

In Fig. 2 we show the entropy  $S_1^{FLM}$  of the FLM for the 16,229 families PFnnnnn in the PFAM 28.0 database, released in May 2015, with numbers nnnnn smaller or equal to 17,126. A linear fit of the entropy as a function of the length  $N$  of the profile is excellent and gives,



**Fig. 2** Entropies of the Fixed-Length Model across the 16,229 families in PFAM 28.0 versus the length  $N$  of the profiles of the families. *Pluses*  $S_1^{FLM}$ ; *Circles* approximation  $2S_2^{FLM} - S_3^{FLM}$ ; *Triangles* lower bound  $S_2^{FLM}$ . The *continuous lines* show the linear fits (11), (13) and (14). *Inset* magnification of the plot in the region  $500 \leq N \leq 1000$



**Fig. 3** Difference between the entropy  $S_1^{FLM}$  of the Fixed-Length Model and its linear approximation (11), across the 16,229 families in PFAM 28.0, versus the length  $N$  of the profiles of the families

$$S_1^{FLM} \simeq \sigma_1^{FLM} \times N, \quad \text{where } \sigma_1^{FLM} = 2.4880 \pm 0.0006, \quad (11)$$

with 95 % confidence.

To investigate the finer statistics of the entropy as a function of  $N$ , we plot in Fig. 3 the residuals of the linear fit,

$$\delta S_1^{FLM} = S_1^{FLM} - \sigma_1^{FLM} \times N. \quad (12)$$

We observe a systematic and negative deviation for lengths  $N < 100$ . The reason for this variation is not entirely clear; one possible explanation could be the extra penalty introduced for shorter profiles [19]. This penalty term is intended to remove the sequences whose scores barely exceed the level of noise expected for random sequences. As a result, the shorter the profile length, the more sequences are removed, with a net effect of reducing the entropy of the family. For larger sizes, the deviation is on average zero, with a standard deviation comprised in a strip of width growing as the square root of  $L$ . This result is expected for independent-site models, due to the central limit theorem.

We also plot in Fig. 2 the lower bound  $S_2^{FLM}$  and the approximation  $2S_2^{FLM} - S_3^{FLM}$  to the entropy  $S_1^{FLM}$ . As the value of  $S_1^{FLM}$  is exactly known, we can assess the accuracy of the bound and of the approximation. This will be useful below in the case of HMM, where an exact computation of  $S_1$  is out of reach. We observe that both quantities increase on average linearly with the profile length  $N$ , with the slopes

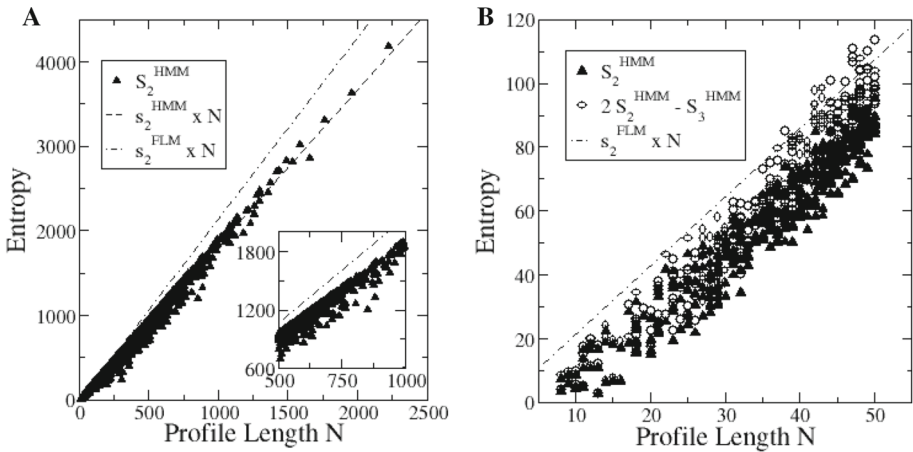
$$S_2^{FLM} \simeq \sigma_2^{FLM} \times N, \quad \text{where } \sigma_2^{FLM} = 2.1438 \pm 0.0012, \quad (13)$$

and

$$2 S_2^{FLM} - S_3^{FLM} \simeq \sigma_{2-3}^{FLM} \times N, \quad \text{where } \sigma_{2-3}^{FLM} = 2.3265 \pm 0.0006, \quad (14)$$

both with 95 % confidence. The deviations of the Renyi entropies  $S_2^{FLM}$  and  $S_3^{FLM}$  with respect to those linear fits show roughly the same behaviour as in the  $S_1^{FLM}$  case (not shown). However, for the latter entropies, deviations are larger and not Gaussianly distributed, as the central limit theorem does not apply to Renyi entropies of order different from unity.





**Fig. 4** Entropies of the Hidden Markov Model across the 16,229 families in PFAM database versus the length  $N$  of the profiles of the families. **a** Lower bound  $S_2^{HMM}$  (triangles). The lines show the linear fit (15) and the one for the Fixed-Length Model, see (13). *Inset* magnification of the plot in the region  $500 \leq N \leq 1000$ . **b** Profiles with lengths  $N \leq 50$  only across all families of index  $< 5000$  in PFAM 24.0. Circles approximation  $2S_2^{HMM} - S_3^{HMM}$ ; Triangles lower bound  $S_2^{HMM}$

### 2.4 Bounds and Approximation for the Full HMM Model

We plot in Fig. 4a the lower bound  $S_2^{HMM}$  to the true entropy  $S_1^{HMM}$  of the HMM model, which we are not able to compute exactly. We observe that  $S_2^{HMM}$  increases on average linearly with the profile length  $N$ , with the slopes

$$S_2^{HMM} \simeq \sigma_2^{HMM} \times N, \quad \text{where } \sigma_2^{HMM} = 1.8367 \pm 0.0015, \quad (15)$$

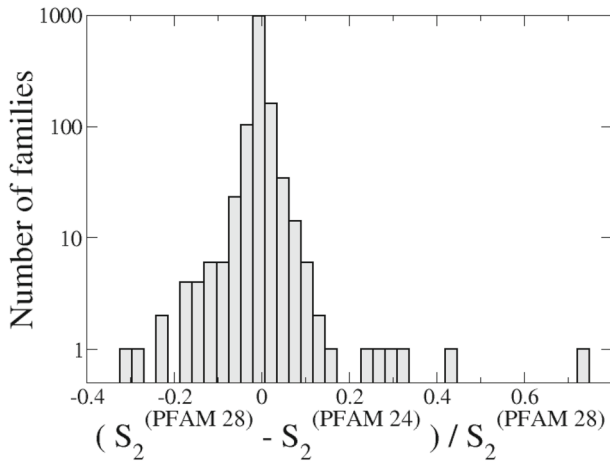
within 95 % accuracy. The slope is 14.3 % lower than its counterpart in the FLM. However, a refined approximation of the entropy based on the calculation of the Renyi entropy of order 3 gives

$$2 S_2^{HMM} - S_3^{HMM} \simeq \sigma_{2-3}^{FLM} \times N, \quad \text{where } \sigma_{2-3}^{FLM} = 2.236 \pm 0.008, \quad (16)$$

which is only 4 % less than the slope found for the FLM, see plot in Fig. 4b. Those results suggest that the entropy of the HMM is only weakly affected by the non-independence between the symbols due to the presence of gaps, and is very close to its FLM counterpart.

### 2.5 Comparison of HMM Entropies for Two Distributions of PFAM

Hereafter we study how the changes in the HMM from one PFAM release to another affect the value of the entropy. To do so we consider the current PFAM 28.0 release (May 2015) and release 24.0 (October 2009). To be as conservative as possible in our estimate of the change in entropy, we first identify 1343 families (among the families PF $n$ nnnn, with  $n$ nnnn  $< 5000$  in release 24.0), whose profile length have not changed from one release to another. The histogram of the relative variations of Renyi entropy  $S_2$  is shown in Fig. 5. The histogram is centered in zero, and is roughly symmetric around the origin. About 2 % of the families, that is, 28 families, show a relative change in entropy larger than 10 %. Those large changes show that some families are still affected by strong undersampling. Once rescaled by the length



**Fig. 5** Histogram of relative changes in the Renyi entropies  $S_2$  between PFAM releases 24.0 and 28.0 for more than 1300 protein families with unchanged profile lengths

$N$  of the profiles, the variations in entropy range between  $-0.18$  and  $0.08$ , which represent variations of about 5–10 % of the average slope  $\sigma_2^{HMM}$  computed above.

### 3 Potts Models for Protein Families and Inference with the Adaptive Cluster Expansion Algorithm

#### 3.1 Cross-Entropy and General Formalism

In this Section we use a Potts model to fit the probability distribution of the sequences  $\mathbf{A} = (a_1, a_2, \dots, a_N)$  associated to a specific protein family. The Potts distribution naturally arises in the Maximum Entropy framework as the least constrained (maximum entropy) distribution capable of reproducing the set  $\mathbf{p}$  of single-site and pairwise frequencies of amino-acids in the natural multi-sequence alignment (MSA) of a given protein family. The parameters of the Potts model are the local fields  $\mathbf{h} = \{h_i(a)\}$ , which may be interpreted as position weight matrices, and the couplings  $\mathbf{J} = \{J_{ij}(a, b)\}$  between the amino acids  $a$  and  $b$  at the sites  $i$  and  $j$ . Those parameters have to be fitted to reproduce the pairwise frequencies  $p_{ij}(a, b)$  and the single-site frequencies  $p_i(a)$  computed from the MSA.

The Potts parameters can be inferred through the minimization of the cross-entropy between the model (defined by its parameters  $\mathbf{J}, \mathbf{h}$ ) and the data ( $\mathbf{p}$ ), equal to minus the log-likelihood of the sequence data:

$$S^{cross}(\mathbf{J}, \mathbf{h}|\mathbf{p}) = \log \left[ \sum_{\mathbf{A}} \exp \left( \sum_i h_i(a_i) + \sum_{i < j} J_{ij}(a_i, a_j) \right) \right] - \sum_i \sum_{a=1}^{q_i} h_i(a) p_i(a) - \sum_{i < j} \sum_{a=1}^{q_i} \sum_{b=1}^{q_j} J_{ij}(a, b) p_{ij}(a, b). \quad (17)$$

In the expression above,  $q_i$  is the number of Potts symbols on site  $i$ . The maximum value of  $q_i$  is 21 (20 amino acids plus the gap symbol), but  $q_i$  can be sizeably smaller on sites where only

a few amino acids appear in the MSA. The cross-entropy  $S^{cross}$  is equivalent to the entropy  $S^{Potts}$  of the Potts model reproducing the 1- and 2-point statistics of the MSA (modulo two contributions introduced below). Hence, the cross-entropy can be used to quantify the diversity of sequences in the protein family in the same way as the entropy of the HMM in Sect. 2. To make sure that the minimum  $\mathbf{J}, \mathbf{h}$  is finite and unique we add to  $S^{cross}$  above the  $L_2$ -norm regularization term

$$\Delta S^{L_2}(\mathbf{J}, \mathbf{h}) = \frac{1}{200\gamma} \sum_{i,a} h_i(a)^2 + \frac{1}{2\gamma} \sum_{i < j} \sum_{a,b} J_{ij}(a,b)^2, \quad (18)$$

where  $10\sqrt{\gamma * M}$  is the magnitude of the largest couplings tolerated in the inference; the fields are allowed to take ten times larger values. The Adaptive Cluster Expansion (ACE) introduced in [15, 16, 20] is a way to calculate the minimal value of  $S^{cross} + \Delta S^{L_2}$  over  $\mathbf{J}$  and  $\mathbf{h}$  through the construction and the summation of many clusters of strongly interacting sites.

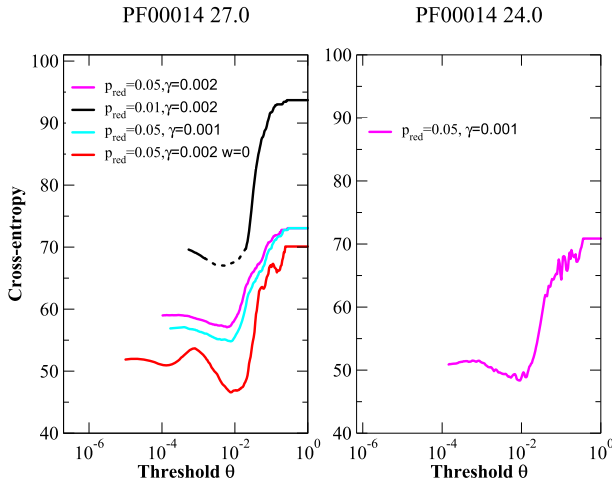
We will investigate how the approximate value of the entropy calculated with the ACE procedure depends on the detailed procedure followed to format the data. In particular, this formatting includes

- the reweighting [21] of similar sequences in the MSA to reduce the phylogenetic correlations of the sequences. The weight of each sequence is taken to be the inverse of the number of sequences with a Hamming distance smaller than  $wN$  (this number is always larger or equal to one, as the sequence itself is included). Hereafter, we will compare results obtained for  $w = 0.2$  and  $w = 0$  (no reweighting);
- the regularization term (18) with  $\gamma \simeq \frac{1}{M}$ , where  $M$  is the number of sequences in the MSA, which can be tuned to improve the convergence of the ACE;
- the reduction of the Potts alphabet. To speed up the algorithm and avoid overfitting we reduce the number of Potts symbols  $q_i$  on each site. To do so we consider only the observed amino acids as possible Potts symbols, and we may also lump together those which are not frequently observed in a unique, abstract Potts symbol. More precisely, we use a criterion based on the frequency, and group together all the amino acids observed with probability  $p < p_{red}$  or whose contributions to the site entropy is smaller than a fraction  $S_{red}$ . We will compare the efficiencies of both criteria for different values of the reduction parameters  $p_{red}$  and  $S_{red}$ .
- the effect of the substitution of the gaps in the MSA with amino acids, randomly drawn according to their frequencies in the sequences without gaps at each position.

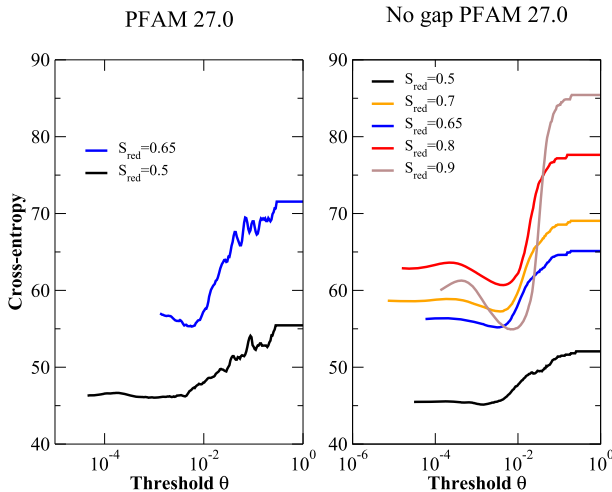
### 3.2 Application to Families PF00014 and PF00397

The inference procedure above will be applied to two small proteins: the trypsin inhibitor (PFAM family PF00014) and WW (PF00397). Trypsin inhibitor is a protein domain reducing the activity of trypsin, an enzyme involved in the breakdown of proteins during digestion; its PFAM profile includes 53 sites. It has been used as a benchmark for structural prediction based on amino-acid covariation [21–23]. WW is a protein domain able to bind peptides, and composed of 31 amino acids. WW was used as a benchmark to test the success of covariation-based procedures to design new folding and biologically functional sequences [24, 25]. We will compare the results of the Potts inference obtained from the MSA in PFAM releases 24.0 and 27.0.

In Figs. 6, 7 and 8 we show the entropies for the two proteins as functions of the threshold  $\theta$  used to select clusters in the ACE procedure, of the number of Potts symbols with the frequency or the entropy criteria, and of the regularization strength  $\gamma$ . At the starting threshold



**Fig. 6** PF00014: Cross-entropy obtained by the ACE algorithm as a function of the threshold  $\theta$  for selecting clusters, for the MSA of PFAM 27.0 and 24.0 (right). Amino acids with  $p < p_{sel}$  are regrouped in a single Potts state, and the regularization strength  $\gamma$  is varied. The reweighting factor is  $w = 0.2$  for all curves but one, where  $w = 0$ , see legend

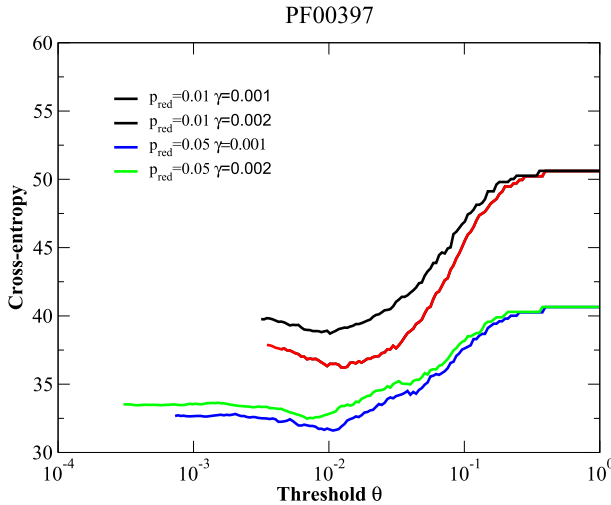


**Fig. 7** PF00014: Cross-entropy obtained by the ACE algorithm as a function of the threshold  $\theta$  for selecting clusters, for the MSA of PFAM 27.0 (left) and after removal of gaps through the randomization procedure (right). The reweighting factor is  $w = 0.2$ . Potts state reduction according to the entropy-based criterion, with cut-off  $S_{red}$

at  $\theta = 1$  only single-site clusters are selected, which corresponds to an independent-site model (IM), and the cross-entropy reads

$$S^{IM} = - \sum_i \sum_{a=1}^{q_i} p_i(a) \log p_i(a). \tag{19}$$

Upon lowering the selection threshold  $\theta$  more and larger clusters are summed in the expansion. The entropy decreases, possibly with some oscillations, and eventually converges at small



**Fig. 8** PF00397: Cross-entropy obtained by the ACE algorithm as a function of the threshold  $\theta$  for selecting clusters, for the MSA of PFAM 27.0, with reweighting  $w = 0.2$ . Amino acids with  $p < p_{sel}$  are regrouped in a single Potts state, and the regularization strength  $\gamma$  is varied. The reweighting factor is  $w = 0.2$

threshold. We note that such oscillations are possible because the algorithm produces an estimate for the entropy by summing up the contributions from many small clusters of sites, which can be either positive or negative (for details see [15, 16]). Once the entropy has reached convergence, we subtract from its value the contribution  $\Delta S^{L2}$  coming from the regularization, see (18), and add a contribution to partially correct for the clustering of amino acids in a unique Potts state,

$$\Delta S^{AGG} = - \sum_i \sum_{a=1}^{k_i} p_i(a) \log \left( \frac{p_i(a)}{p_i(r)} \right). \tag{20}$$

In the expression above,  $p_i(r)$  denotes the frequency of the abstract Potts symbol,  $r$ , which stands for the  $k_i$  Potts states lumped together. By definition,  $p_i(r) = \sum_{a=1}^{k_i} p_i(a)$ . The correction  $\Delta S^{AGG}$  vanishes if no amino acid have been grouped on site  $i$ , and  $p_i(r) = 0$ . If  $k_i \geq 1$ ,  $\Delta S^{AGG}$  is not equal to zero, and is equal to the entropy of  $k_i$  independent symbols with probabilities  $p_i(a)/p_i(r)$ , weighted by the probability  $p_i(r)$  of the abstract Potts symbol. It allows us to recover the full model from the reduced one in the IM case, see Table 2. The final expression for the entropy is therefore  $S^{Potts} = S^{cross} - \Delta S^{L2} + \Delta S^{AGG}$ .

As a general trend, we observe that the cross-entropy decreases when the reduction parameter  $p_{red}$  is made smaller or the cut-off fraction  $S_{red}$  is made larger, see Tables 1, 2, 3, 4 and 5. In other words, as the number of retained Potts symbols increases, the entropy decreases. This behaviour is easy to understand: keeping more Potts symbols amounts to reproducing more single-site and pairwise frequencies, and hence, to fulfil more constraints. Note also that due to the removal of the regularization contribution to the entropy,  $\Delta S^{L2}$ , the value of the cross-entropy depends only weakly on the regularization strength  $\gamma$ , see results for  $\gamma = 0.001$  and  $\gamma = 0.002$  in Figs. 6 and 8. Nevertheless, we observe that the cross-entropy increases slightly with  $\gamma$ , as the correlations are effectively less tightly constrained.

For PF00014, we obtain the entropy of the Independent Model (IM),  $S^{IM} = 96.75$ , and a corresponding entropy per site  $\sigma^{IM} = 1.82$  (Table 2). Naturally, as pairwise constraints

**Table 1** PF00014: Results of the Potts model inference with the ACE procedure

$p_{red}, \gamma, w$	Entropy $S^{Potts}$	Cross-entropy $S^{cross}$	$\Delta S^{L2}$	$\Delta S^{AGG}$
PFAM 27.0				
0.01	67.68	69.61	-5.00	3.07
0.05	81.02	59.00	-1.66	23.7
0.05, $\gamma = 10^{-3}$	79.11	56.89	-1.48	23.7
0.05, $\gamma = 10^{-3}, w = 0$	75.41	51.88	-2.74	26.28
PFAM 24.0				
0.05, $\gamma = 10^{-3}$	74.15	50.89	-3.79	27.05

Potts states were grouped together according to their frequencies, with cut-off  $p_{red}$ . Unless otherwise specified the reweighting factor is  $w = 0.2$  and the regularization strength  $\gamma = 0.002$

**Table 2** PF00014, release 27.0: Independent model with selection of the number of Potts states

$p_{red}$	Entropy $S^{IM}$	Cross-entropy $S^{IM-cross}$	$\Delta S^{L2}$	$\Delta S^{AGG}$
0.01	96.75	93.68	-0.014	3.07
0.05	96.75	73.06	-0.01	23.7

**Table 3** PF00014, release 27.0: Results of the Potts model inference with the ACE procedure

$S_{red}$	Entropy $S^{Potts}$	Cross-entropy $S^{cross}$	$\Delta S^{L2}$	$\Delta S^{AGG}$
0.9	62.97	60.01	-5.11	8.07
0.8	75.18	62.88	-3.56	15.86
0.7	81.07	58.65	-2.03	24.4
0.65	82.96	56.27	-1.68	28.37
0.6	83.81	51.43	-0.67	33.05
0.5	86.58	45.48	-0.31	41.43

Potts states were grouped together according to their contributions to the site entropy, with cut-off  $S_{red}$ . Gaps are replaced with randomly drawn amino acids, see text

**Table 4** PF00397, release 27.0: Entropies with the independent Potts model and selection of Potts states based on their frequencies

$p_{red}$	Entropy $S^{IM}$	Cross-entropy $S^{cross}$	$\Delta S^{L2}$	$\Delta S^{AGG}$
0.01	52.38	50.61	-0.01	1.78
0.05	52.37	40.65	-0.01	11.73

Reweighting factor:  $w = 0.2$ ; regularization strength is  $\gamma = 0.002$

are introduced the entropy of the corresponding model decreases. The entropy of the Potts model including couplings is  $S^{Potts} = 67.68$  ( $\sigma = 1.28$ ) for  $p_{red} = 0.01$ , and  $S^{Potts} = 75.4$  ( $\sigma = 1.42$ ) for  $p_{red} = 0.05$  (Table 1). The decrease in entropy is similar when regrouping the symbol according to the  $p_{red}$  or  $S_{red}$ -based criteria as long as the number of Potts symbols remaining are equivalent (Table 3). For instance, results for the reduction  $p_{red} = 0.05$  are similar to the ones with  $S_{red} = 0.7$ , and the ones with  $p_{red} = 0.01$  are similar to the ones

**Table 5** PF00397, release 27.0: Entropies with selection of Potts states based on their frequencies

$p_{red}$	Entropy $S^{Potts}$	Cross-entropy $S^{cross}$	$\Delta S^{L2}$	$\Delta S^{AGG}$
0.01	38.25	39.76	-3.25	1.78
0.01, $\gamma = 0.01$	37.64	37.56	-1.70	1.78
0.05	43.74	32.67	-0.66	11.73
0.05, $\gamma = 0.01$	43.821	32.67	-0.58	11.73

Unless otherwise specified the regularization strength is  $\gamma = 0.002$

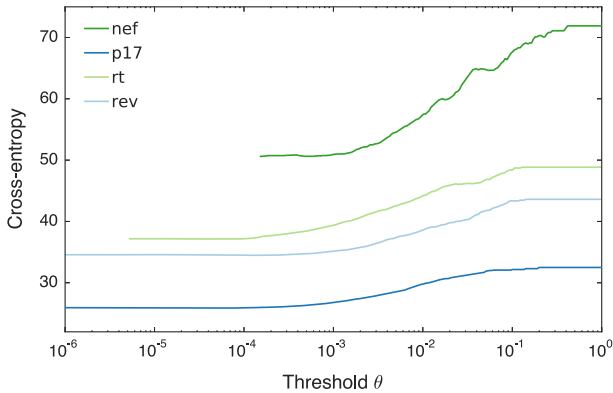
with  $S_{red} = 0.8 - 0.9$ . As noted above, enforcing more pairwise constraints (i.e., taking smaller  $p_{red}$  or larger  $S_{red}$ ) results in lower entropy.

We have also calculated the entropy in the family PF00014 from PFAM 24.0, containing  $M \simeq 2000$  sequences and an effective number of sequences  $M_{eff} = 1000$  after reweighting with  $w = 0.2$ , and compared it to the outcome of the calculation for PFAM 27.0, corresponding to  $M \simeq 4000$ ,  $M_{eff} = 2000$ . We find that the entropy has increased by about 5 % between the two releases, probably a result of the introduction of more diverse sequences in the database. In the absence of reweighting ( $w = 0$ ), the entropy decreases (by about 4 %), as similar sequences are given more weights and the distribution of sequences is more peaked, see Table 1. For PF00397, we obtain a similar behavior:  $S^{IM} = 52.38$  ( $\sigma = 1.69$ ) for the IM, and  $S^{Potts} = 38$  ( $\sigma = 1.22$ ) for  $p_{red} = 0.01$ ,  $S^{Potts} = 43.82$  ( $\sigma = 1.41$ ) for  $p_{red} = 0.05$ , see Tables 4 and 5.

## 4 Application to Phylogenetically Related HIV Proteins

Human immunodeficiency virus is distinguished by both a high mutation rate and a very short replication cycle, enabling the virus to rapidly generate mutations within specific parts of the viral sequence targeted by host immune responses, referred to as epitopes. Viruses bearing mutations within these epitopes are able to escape recognition by the immune system, thus preventing effective immune control of infection and allowing viral replication to continue unchecked. This process of rapid mutation also results in extraordinary HIV sequence diversity at the level of both individual hosts and of the virus population as a whole [26]. Together these factors contribute to the difficulty of designing an effective vaccine against HIV, which must be able to prime the immune system to combat diverse strains of the virus while also directing immune responses toward epitopes where escape is most difficult. Thus, quantifying the constraints on HIV proteins that limit their mutability represents an important step in the process of vaccine design.

Here, we estimate the entropy of various HIV proteins through maximum entropy models capturing the frequency of mutations at each site and the pairwise correlation between mutations. Previously, such models have been successfully used in the context of HIV to predict the fitness (ability to replicate) of a library of mutant strains of the virus [27, 28], and to explore aspects of protein structure and stability [29, 30]. Unlike the protein families considered above, the HIV proteins we study here are more closely related phylogenetically, and thus the entropy that we compute may underestimate the true, potential variability of these proteins. However, this approach should still be successful in capturing nontrivial constraints on HIV. A variational mean-field theory calculation suggests that, while immune pressure due to genetically diverse individuals perturbs the inferred fields and phylogeny



**Fig. 9** Typical behavior of the cross-entropy obtained by the ACE algorithm as a function of the threshold  $\theta$ , shown for various example HIV proteins. All exhibit good convergence toward stable values of the cross-entropy as the threshold is lowered. Similar results are also seen for the other proteins not shown here

further modulates the fields in the maximum entropy model, the inferred couplings are not strongly affected [31].

The outline for this Section is as follows. In Sect. 4.1 we compare the entropy of all HIV proteins from different virus subtypes, except for the highly variable envelope subunit protein gp120. This analysis reveals a subset of strongly conserved HIV proteins that appear to be unusually constrained. In Sect. 4.2 we employ simulations to compute the loss in entropy when particular residues in these proteins are held fixed. There we show that the typical locations of escape mutations, which the virus uses to evade the host immune system, are associated with sites that strongly reduce the entropy when held fixed.

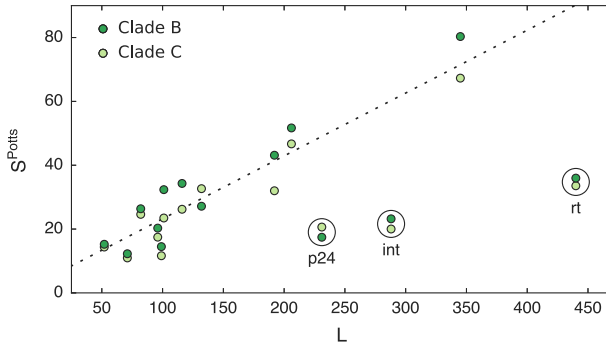
#### 4.1 Diversity Across the HIV Proteome

We inferred Potts models describing various HIV proteins from two prominent subtypes of the virus (clade B, dominant in Western Europe and the Americas, and clade C, dominant in Africa and parts of South Asia) using the ACE method, and through this method we obtained an estimate of the entropy for each protein. For more details on the inference method and computation of the entropy, see Sect. 3. The cross-entropy displayed good convergence for all proteins we tested, thus we do not expect large errors in our estimates of the entropy (Fig. 9). In all cases we used an entropy cutoff of  $S_{red} = 0.9$ , and regularization strength  $\gamma \simeq 1/M$ , where  $M$  is the number of unique patients from which the sequence data was collected (ranges from approximately 500 for some accessory proteins to 10,000 for protease).

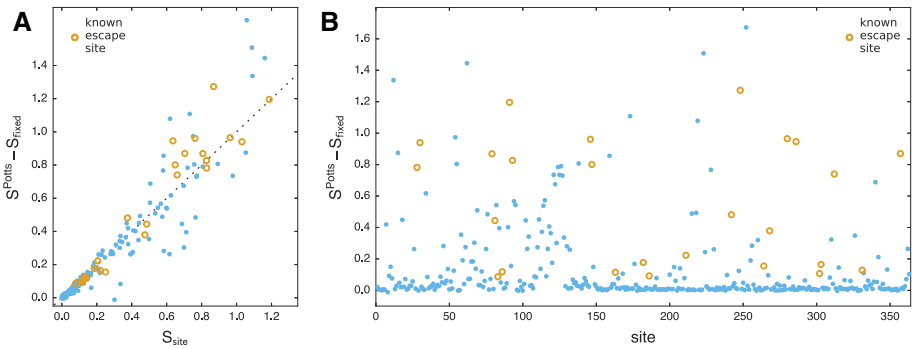
In Fig. 10 we show the entropy  $S^{Potts}$  of each protein versus its length  $L$  in amino acids. We find that most HIV proteins have a typical entropy per site of around 0.2, which holds for proteins obtained from both clade B and clade C viruses. Note that although the entropy scales roughly linearly with the protein length, this does not imply that variation at each site is independent; couplings between sites can contribute substantially to the entropy (see Fig. 11 below). In contrast with the typical scaling, there also exists a subset of HIV proteins that appear to be more highly constrained. Proteins p24, integrase, and reverse transcriptase have an entropy per site of roughly 0.08, substantially lower than for other proteins.

There are several factors that may contribute to the reduced entropy observed for these proteins. At first, the reduced variability of p24 may appear surprising because this protein is frequently targeted by host immune responses [32,33], which would encourage frequent





**Fig. 10** Entropy per site is comparable for most HIV proteins except for a subset that are more highly conserved. Here we show the entropy  $S^{Potts}$  versus length  $L$  for a set of HIV proteins. The typical entropy per site for these proteins, including the effects of coupling between sites, is around 0.2 (dotted line). In contrast, the proteins p24, integrase (int), and reverse transcriptase (rt) appear to be substantially constrained, with entropy per site of only 0.08 (circled). Note that the surface protein gp120 is not included in this analysis; this protein is highly variable, and may exhibit higher entropy per site than typical HIV proteins



**Fig. 11** Change in entropy  $S^{Potts} - S^{fixed}$  upon individually fixing each site in HIV proteins p17 and p24 equal to the consensus amino acid. Known escape sites are highlighted (open circles). **a** Reduction in entropy from fixing a site is typically similar to the single site entropy  $S_{site}$ , particularly for sites with low entropies, but sites with strong interactions can deviate significantly from this value. See main text for details. **b** Change in entropy as a function of position along the p17 (sites 1–132) and p24 (sites 133–363) proteins. Variation at known escape sites often contributes substantially more to the entropy than variation at other sites in the same epitope. Note that the  $CD8^+$  T cell epitopes considered here are usually 9–11 amino acids in length. Escape mutations can occur at sites within the epitope or at nearby flanking sites

mutation. This protein forms the viral capsid, however, and is therefore subject to strict conformational constraints. The mature capsid is composed of p24 hexamers and pentamers that bind together in a “fullerene cone” shape [34]. Previous work has shown that multiple mutations in residues along the p24 hexamer-hexamer interfaces, in particular, may be tightly constrained [35]. Epitopes in these regions are also frequently targeted by individuals who more effectively control HIV infection, possibly because mutations in these regions are more likely to damage viral fitness, thus decreasing the likelihood of escape [35].

In contrast to p24, reverse transcriptase and integrase are not frequently targeted by host immune responses [32, 33]. They are responsible for the reverse transcription of viral RNA to DNA and the integration of viral DNA into the host genome, respectively. These proteins do not appear to be under substantial pressure to widely explore the sequence space [36], which,

in addition to functional constraints, contributes to their reduced variability. Interestingly, we note that the conservation of reverse transcriptase observed here is also consistent with recent experimental studies that found extremely low tolerance for insertions in proteins involved in transcription for several different viruses [37–40], suggesting that such proteins may potentially operate under strong functional constraints in more general cases.

#### 4.2 Relationship Between the Entropy and Local Pressure

In addition to characterizing the entropy for HIV proteins, we can also explore how variation at individual sites within a protein contributes to its overall entropy. The simplest way to do this is just to compute the single site entropy  $S_{\text{site}}(i)$  of each site  $i$ , obtained from the empirical correlations

$$S_{\text{site}}(i) = - \sum_a p_i(a) \log p_i(a). \quad (21)$$

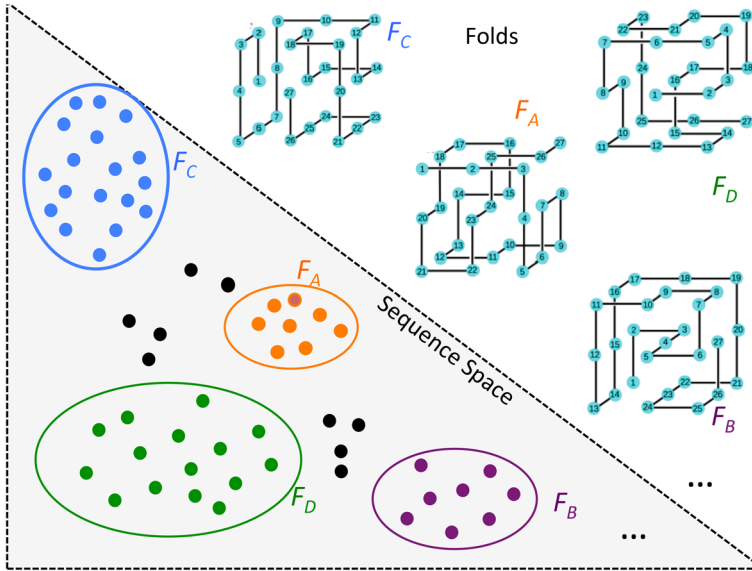
The drawback of this approach is that it neglects the effects of higher order constraints on protein sequences, such as those parameterized by the  $J_{ij}(a, b)$ , beyond just the frequency of amino acids observed at each site.

To capture some of these higher order constraints, we can use the Potts model inferred for each protein to generate an ensemble of sequences with the amino acid at certain sites held fixed. We can then compute the entropy  $S_{\text{fixed}}$  of this ensemble of sequences using the ACE method as before. The change in entropy  $\delta S = S^{\text{potts}} - S_{\text{fixed}}$  upon fixing a site to a given value then quantifies the contribution of variation at that site to the entropy, including the effects of the inferred pairwise interactions. In the following, we choose to fix sites to their consensus values (one at a time), but the approach could be extended to any specific wild-type sequence. In Fig. 11a, we see that the reduction in entropy from fixing most sites in the HIV proteins p17 and p24 is similar to the corresponding single site entropy  $S_{\text{site}}$ . The effect of interactions is difficult to discern at this level for sites with very low variability. However, as shown in Fig. 11a,  $\delta S$  deviates substantially from  $S_{\text{site}}$  for a number of more variable sites where the effects of mutations are strongly coupled to other sites in the protein (note the scale in the above figure). The reduction in entropy for sites that lie above the line in Fig. 11a is larger than expected from the single site entropy alone, indicating the presence of mostly positive (or, in the language of fitness, compensatory) couplings to other sites. For sites below the line  $\delta S$  is smaller than expected, indicating more negative (or deleterious) couplings that tend to suppress mutation. These latter sites may then be good targets for effective immune responses.

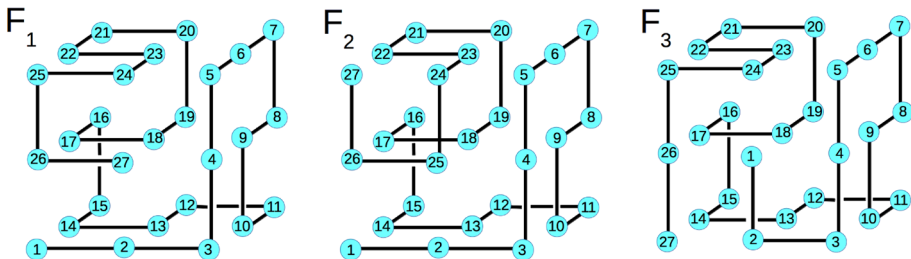
Entropy has previously been associated with immune escape in HIV. Generally, escape tends to occur more rapidly at epitopes where the average single site entropy is high [41,42]. We also observe a connection between the sites where escape mutations are typically observed in well-characterized epitopes (see [27]) and entropy. In Fig. 11b, we show the change in entropy upon fixing each site to the consensus amino acid in the p17 (sites 1–132) and p24 (sites 133–363) proteins, with known escape sites highlighted. Typically, these known escape sites contribute substantially more to the entropy than other sites within the same epitope. This result is intuitive: naturally we would expect that mutations at highly variable sites, or ones with many available compensatory interactions, should typically come at a low fitness cost to the virus, otherwise we would not frequently observe viruses with mutations at those sites. Mutations that both confer immune escape and which come at little fitness cost should then be selected more frequently.

### 5 Exact and Approximate Values of the Entropy for Lattice-Based Proteins

In this section, we compute the entropy of the families of lattice-based proteins (LP) [17, 43–45]. Lattice proteins considered here are composed of 27 amino acids occupying the sites of a  $3 \times 3 \times 3$  cube, see Figs. 12 and 13. There are  $N_{\text{fold}} = 103,346$  possible folds  $F$  (conformations of the protein backbone) unrelated by symmetry. Given a fold  $F$ , each amino-acid sequence  $A = (a_1, \dots, a_{27})$  is assigned an energy



**Fig. 12** Pictorial representation of the sequence space (*bottom left corner*) and of four depicted folds (*top right corner*) among the  $N_{\text{fold}}$  possible structures  $F$ . Sequences  $A$  that fold in one of the four structures, say,  $F$ , e.g., such that  $P_{\text{nat}}(F|A) > 0.995$ , see (23), are shown by *coloured dots*, with the *same colors* as the corresponding structures. *Dark dots* correspond to unfolded sequences, i.e., having low values of  $P_{\text{nat}}$  with all structures. The logarithm of the volume in the sequence space associated to each fold defines its entropy, otherwise called *designability* [43]. The entropies  $S^{\text{Potts}}$  of the four folds shown in the figure have been calculated in [47], using the pairwise Potts models inferred from the families of sequences associated to the folds, with the ACE expansion and are recalled in Table 6



**Fig. 13** The three folds considered here:  $F_1$ ,  $F_2$  and  $F_3$ . It is easy to see that  $F_2$  is obtained from  $F_1$  by simply exchanging the sites 25 and 27, while  $F_3$  is obtained from  $F_1$  by exchanging 1 and 27. We also see that the only sites affected by these exchanges are the nine sites 1, 2, 4, 14, 18, 22, 24, 25 and 27

$$E_{LP}(\mathbf{A}|F) = \sum_{i < j} c_{ij}^{(F)} E(a_i, a_j) \quad (22)$$

where  $E(a, b)$  is the Miyazawa-Jernigan statistical energy matrix [46]. The matrix  $c_{ij}^{(F)}$  is the contact matrix associated with the fold  $F$ : the entry is equal to unity if  $i$  and  $j$  are in contact, i.e., are nearest neighbors on the cube, and zero otherwise.

The probability that a given sequence  $\mathbf{A}$  folds in conformation  $F$  is defined following [17] as:

$$P_{\text{nat}}(F|\mathbf{A}) = \frac{e^{-E_{LP}(\mathbf{A}|F)}}{\sum_{F'=1}^{\mathcal{N}_{\text{fold}}} e^{-E_{LP}(\mathbf{A}|F')}} = \frac{1}{1 + \sum_{F'(\neq F)} e^{-[E_{LP}(\mathbf{A}|F') - E_{LP}(\mathbf{A}|F)]}} \quad (23)$$

From Eq. (23) it is clear that the fold  $F^*$ , which maximize the probability that a given sequence  $\mathbf{A}$  folds in it, is the one, among all the other possible and competing structures  $F'$ , with minimal energy  $E_{LP}(\mathbf{A}|F^*)$ . However the sequence is said to be folded in this structure  $F^*$  only if  $P_{\text{nat}}(F^*|\mathbf{A})$  is very large, typically larger than 0.995. Therefore the requirement for a sequence to fold in a structure  $F^*$  is the existence of a large energy gap  $E_{LP}(\mathbf{A}|F') - E_{LP}(\mathbf{A}|F^*)$  (at least of the order of five, in units of the temperature, set to unity here) with the other competing structures. Given a fold, this gap condition is generally satisfied by many sequences, see sketch in Fig. 12, which define the protein family. The fold attached to this set of sequences is called native fold, while the structures that have the smallest energy gap with the sequences in the set are said to be its close competitors.

## 5.1 Designability and Entropy

An important characteristic of a structure is the volume (cardinality) of the set of attached sequences, see Fig. 12, called designability [43,45]. The logarithm of the numbers of sequences folding in a given structure informally corresponds to the entropy defined here, see introduction. In [43] it was shown by numerical simulations that the designability depends on the structure: as sketched in Fig. 12, some structures are associated to a large volume in the sequence space, while some correspond to smaller volumes. In [45], it was proposed that the largest eigenvalue of the contact map  $c_{ij}$  of a structure is indicative of its designability.

In a recent work [47], large alignments of size  $\mathcal{O}(10^4)$  for the four structures ( $F_A, F_B, F_C, F_D$ ) in Fig. 12, were generated, and used to infer the Maximum Entropy Potts models reproducing the 1- and 2-point statistics with the Adaptive Cluster Expansion described in Sect. 3. We summarize here the procedure we have followed to generate the alignments of sequences folding in the four structures of Fig. 12, and the results we have obtained for their entropies. To generate a MSA attached to a fold, say,  $F$ , we perform a search in the sequence space to find sequences  $\mathbf{A}$  with large folding probability  $P_{\text{nat}}(F|\mathbf{A}) > 0.995$  [48]. To this aim we have used a Monte Carlo procedure to sample the Gibbs distribution associated to the effective Hamiltonian

$$\mathcal{H}_W(\mathbf{A}|F) = -\ln P_{\text{nat}}(F|\mathbf{A}), \quad (24)$$

in the sequence space at large inverse temperature ( $\beta = 10^3$ ). Here  $W$  denotes the world of proteins, that is, the set of all possible structures; in [47] 10,000 folds among the  $\mathcal{N}_{\text{fold}}$  were randomly chosen. The sampled sequences form the MSA, which gives access to the 1- and 2-point statistics of the family. We have then inferred the pairwise Maximum-Entropy Potts model and computed its cross-entropy with the ACE procedure. Results are given in Table 6.

**Table 6** Estimates of how designable are the proteins families associated to structures  $F_A, F_B, F_C, F_D$  (ranked in increasing order of their entropies): largest eigenvalues of the contact map matrix  $\mathbf{c}$ , entropy of the inferred Potts model obtained by ACE (2nd), and mean percentage of identity between sequences (3rd)

Fold	Top eigenvalue of $\mathbf{c}$	Potts entropy (ACE)	Mean % identity btw seq.	Dist. to nearest struc.
$F_B$	2.6	50.2	24	14
$F_A$	2.5	50.9	23	11
$F_D$	2.7	55.4	21	9
$F_C$	2.9	58.4	19	4

We also give the distance to the nearest structure (4th column). For the identity calculation, we average the number of amino acids that take their consensus values, and divide by the number of amino acids in the protein ( $=27$ )

**Table 7** Entropies for the family associated to the fold  $F_1$  and for the protein worlds with one, two and three structures, as calculated exactly, by fitting an independent model or by fitting a Potts model, either with the ACE algorithm or with exact calculation

Protein world	Exact	Ind. model	Potts (ACE)	Potts (exact)
$F_1$	80.8848	80.8848	80.8848	80.8848
$[F_1; F_2]$	77.1560	77.5035	77.1060	77.2504
$[F_1; F_3]$	77.1560	77.5035	77.1060	77.2485
$[F_2; F_3]$	77.2054	77.8174	77.2294	
$[F_1; F_2 F_3]$	75.3762	75.7432	75.3331	

For the protein world  $[F_2; F_3]$ , the entropy is that of the family associated to the fold  $F_2$ . Empty cells signal entropies that would have been too costly to compute, see main text

The Potts entropy is bounded from above by  $27 \times \log 20 \simeq 80.9$ ; the difference between this upper bound (corresponding to a set of  $L = 27$  fully unconstrained amino acids) and the Potts entropy is a measure of the structural constraints acting on the sequences. As reported in Table 6 we have also compared the Potts entropy to different estimators, such as the maximal eigenvalue of the contact matrix  $\mathbf{c}^{(F)}$  of the target fold under consideration [45], and the mean sequence variability in the alignment (average Hamming distance to the consensus sequence across the alignment), see Supplementary Information in [47]. The general picture that arises from [47] is that the presence of competing folds that are close (either in terms of the contact matrix or in terms of energy gaps) to the native fold globally constrains the protein sequences and reduces the entropy of the family, hence defining an entropy cost associated to the competition in folding. Hereafter we show that this cost can be accurately computed in the case of a very small number of competing structures. This simple ‘protein world’ can be used, in turn, as a testbed for the inference algorithm and the approach developed in Sect. 3.

## 5.2 Exact Calculation of the Entropy for Pairs or Triplets of Proteins

We start from the simplest case, that of a unique possible fold,  $F_1$  in Fig. 13. In that case, any sequence  $A$  will necessarily fold into  $F_1$ , and the corresponding effective Hamiltonian  $\mathcal{H}_{F_1}(A|F)$  vanishes. The amino acids can be assigned randomly on each site, and the entropy is simply (Table 7, top line):

$$S(F_1) = \ln(20^{27}) = 80.8848. \tag{25}$$

In a more complex protein world made of two proteins,  $F_1$  and  $F_2$ , the probability that a sequence  $A$  folds into  $F_1$  now defines the effective Hamiltonian:

$$\mathcal{H}_{[F_1; F_2]}(A|F_1) = -\log\left(1 + e^{-E_{LP}(A|F_2)+E_{LP}(A|F_1)}\right) \tag{26}$$

where  $[F_1; F_2]$  denotes the two-protein world made of  $F_1$  and  $F_2$ , with  $F_1$  chosen as the reference fold. On our small cube, the contact matrices are uniquely defined by a set of 28 contacts (pairs of neighbours on the cube, excluding contiguous sites on the protein backbone). We have found a large number of pairs of protein folds that share 24 out of 28 of those contacts. Choosing  $F_1$  and  $F_2$  to have 24 common contacts (Fig. 13), we have only 4 pairs of sites that are relevant in the calculation of the energy difference in (26). The effective Hamiltonian will be constraining 8 sites (2 for each contact) at most, and will not depend on the amino acids on the other sites. It turns out that out of those 8 sites, the 4 differing contacts are carried by only 6 distinct sites. The calculation of the partition function associated to  $\mathcal{H}_{[F_1; F_2]}(A|F_1) = \mathcal{H}_{[F_1; F_2]}(a_1, a_2, \dots, a_6|F_1)$  is numerically tractable as it involves a summation over  $20^6$  configurations only,

$$Z_{[F_1; F_2]} = 20^{21} \sum_{a_1=1}^{20} \dots \sum_{a_6=1}^{20} e^{-\beta \mathcal{H}_{[F_1; F_2]}(a_1, \dots, a_6|F_1)}, \tag{27}$$

and the corresponding entropy for the fold  $F_1$  is

$$S([F_1; F_2]) = \ln Z_{[F_1; F_2]} - \frac{d}{d\beta} \ln Z_{[F_1; F_2]}. \tag{28}$$

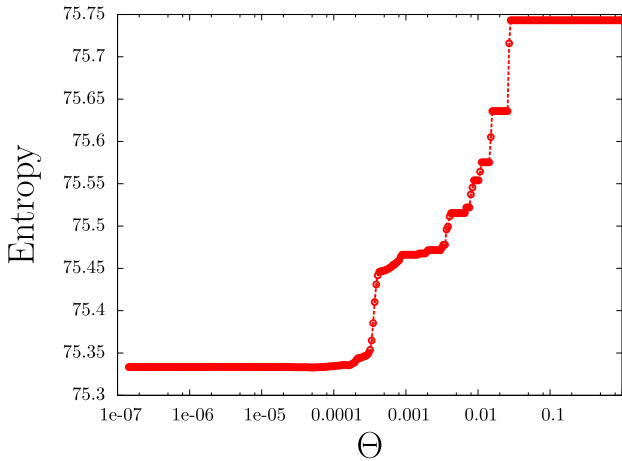
The value of this entropy is given in Table 7 (1st column), and is close to 77.16. The decrease with respect to  $S(F_1)$  in (25) measures the loss in entropy due to the introduction of the competing fold  $F_2$ . Using a conversion in log base 20 the entropic cost is of  $\approx 1.3$  site.

We then consider another fold  $F_3$ . This third structure is also close to  $F_1$ , see Fig. 13, and a bit further away from  $F_2$ . We have calculated the entropy of the twofold world comprised of  $F_1$  and  $F_3$ : we find that  $S([F_1; F_3])$  is identical to  $S([F_1; F_2])$  as  $F_2$  and  $F_3$  both share 24 contacts with  $F_1$ . The entropy  $S([F_2; F_3])$  is slightly larger than  $S([F_1; F_3])$  (Table 7), as can be expected from the fact that  $F_2$  and  $F_3$  are further apart, with only 22 common contacts. The energy gap between  $F_2$  and  $F_3$  is therefore larger than between  $F_1$  and  $F_3$ , and sequences folding in  $F_2$  are less constrained by the presence of the competing fold  $F_3$  than the sequences folding in  $F_1$  in the presence of the competing fold  $F_3$  too. This result agrees with the qualitative findings of [47].

When this third fold  $F_3$  is added to the protein world, the effective Hamiltonian (associated to the folding into  $F_1$ ) reads

$$\mathcal{H}_{[F_1; F_2 F_3]}(A|F_1) = -\log\left(1 + e^{-E_{LP}(A|F_2)+E_{LP}(A|F_1)} + e^{-E_{LP}(A|F_3)+E_{LP}(A|F_1)}\right) \tag{29}$$

and depends on the values of nine amino acids on the sequence only. The calculation of the partition function and of the entropy  $S([F_1; F_2 F_3])$  can be done along the lines above; it now requires to sum up over  $20^9$  configurations, and was done in one CPU day on a desktop computer. Addition of a thirdfold leads to a value of the entropy of 75.37, which shows that an additional 0.8 site has been constrained in the process, see Table 7.



**Fig. 14** Entropy of the protein family associated to  $F_1$  as a function of the threshold  $\theta$  as computed by the ACE procedure in the case  $[F_1; F_2 F_3]$ . The entropy saturates to a value very close to the exact one, see Table 7. The plateau at the beginning of the calculation (large  $\theta$ ) corresponds to the entropy of the IM

### 5.3 Approximate Calculation of the Entropy Based on the Inferred Potts Models

As the models above are exactly solvable, they can be used as a simple testbed for our Maximum Entropy-ACE approach, which we have already used for real protein data. To do so, we have computed the one- and two-point marginals,  $p_i(a)$  and  $p_{ij}(a, b)$  for the protein worlds  $[F_1; F_2]$ ,  $[F_1; F_3]$ ,  $[F_2; F_3]$  and  $[F_1; F_2 F_3]$ , from very large MSA with  $5 \times 10^5$  sequences generated through Monte Carlo sampling. We first fit the 1-point statistics only with an IM. The corresponding entropies are given by (19), with  $q_i = 20$  for all 27 sites, with values listed in the second column of Table 7.

We then take into account the 2-point statistics, and infer the corresponding Maximum-Entropy Potts models with the ACE algorithm. We show in Fig. 14 the behaviour of the entropy predicted by the ACE algorithm for the case  $[F_1; F_2 F_3]$ , as a function of the threshold  $\theta$  in the algorithm (see Sect. 3). Similar curves are obtained for the protein worlds made of two structures. The entropy converges to a value very close to the exact value calculated through complete enumeration of the Potts states, see Sect. 5.2. Even though our sampled alignment is very large here, correlations are still not exactly measured, leading to seemingly significant correlations on pairs of sites outside the restricted subset involved in the exact partition function. Due to those spurious constraints, the entropy is slightly lower than its exact value (Table 7, third column).

Last of all, we can determine the coupling parameters of the Potts model by brute force optimization of the cross-entropy (17), as the number of sites effectively involved is small (see discussion above). This computation provides us with the exact entropy of the Potts model associated to the MSA we have generated, see results in Table 7, fourth column. The computation takes about one hour on a desktop computer when the number of relevant sites is 6 (for the worlds  $[F_1; F_2]$  and  $[F_1; F_3]$ ), but would require several days when the number of relevant sites is 9 (for the worlds  $[F_2; F_3]$  and  $[F_1; F_2 F_3]$ ). As expected, the entropy of the exact Potts model is now larger than the exact entropy (Table 7, first column): the Potts model is, indeed, less constrained than the many-body model defined by the Hamiltonians in (26,29).

## 6 Discussion

In this paper we have used different methods to calculate the entropy of protein families. One of our main findings, obtained from the wide-scale comparative analysis of a large number of PFAM families, is that the entropy primarily depends on the length  $N$  of the family profile. More precisely, we find a linear scaling  $S \simeq \sigma N$ . The value of the slope  $\sigma$  depends on the method we have used. For the HMM model we find  $\sigma \simeq 1.9 - 2.2$ . Maximum Entropy modelling of a few protein families with pairwise interaction Potts models give values of  $\sigma$  ranging between 1.2 (when all amino acids present in the multiple sequence alignment are kept in the modelling) to 1.7 (for large reduction in the number of a.a. used), while the independent-site model gives  $\sigma \simeq 1.7 - 1.8$ .

Those estimates for  $\sigma$  are compatible with previous results in the literature. The authors of [18] estimated  $\sigma \simeq 1.9$  based on the following modelling of a protein family. Given the contact map  $\mathbf{c} = \{c_{ij}\}$  of the native fold corresponding to the family (supposed to be perfectly conserved throughout the family), the energy of an amino-acid sequence  $\mathbf{A}$  is approximated as a sum of pairwise energetic parameters between amino-acids in contact on the structure (relative distance smaller than 6.5 Å),

$$E_{AP}(\mathbf{A}, \mathbf{c}) = \sum_{i < j} E(a_i, a_j) c_{ij} \quad (30)$$

The energetic parameters  $E(a, a')$  describe the magnitude of the interaction between amino acids  $a$  and  $a'$ , and are given by the Miyazawa-Jernigan energetic matrix; variants of this statistically derived energy matrix  $E$  were proposed without affecting much the value of  $\sigma$ . The Gibbs distribution associated to this energy is the sequence distribution for the family. By computing the average energy  $\langle E \rangle(T)$  at different temperatures  $T$  with Monte Carlo simulations, one can obtain the value of the entropy through thermodynamic integration:

$$S(T) - S(\infty) = \frac{\langle E \rangle(T)}{T} - \int_T^\infty dt \frac{\langle E \rangle(t)}{t^2}. \quad (31)$$

In the formula above,  $S(\infty)$  is the entropy of the system at infinite temperature, and is equal to  $N$  times the entropy of the background amino acid distribution,  $s_{BG} = -\sum_{a=1}^{20} p(a) \log p(a)$ . As a result the estimate of  $\sigma \approx 1.9$  was found, see Fig. 2 of [18].

The entropies we have found with the HMM models are larger than with the Maximum Entropy Potts approach. One possible explanation is that, while HMM are routinely used to identify families, they are not supposed to reproduce faithfully the statistics of the MSAs when used as generative models. More precisely, HMM generate sequences that are more variable than the ones found in natural MSA, even at the level of single-site frequencies. In the Maximum Entropy Potts approach, we find smaller values of the entropy, especially when increasing the number of Potts states on each site (up to the number of amino acids observed at least once in the MSA). The reason is that increasing the number of pairwise correlations to reproduce corresponds to increasing the number of constraints to satisfy, and therefore leads to a decrease in entropy. However, this may also lead to overfitting the data if the number of sequences in the MSA is too small.

In the case of phylogenetically related HIV sequences we find a tenfold decrease for the entropy per site,  $\sigma \simeq 0.2$ . This small value reflects the high phylogenetic correlations between sequences and the poor variability in the MSA. To better understand how this value compares to the ones we have found for protein families, we have considered the example of the RT (reverse-transcriptase, PF00078), a long protein with more than 500 amino-acids, which is unusually conserved in the HIV data (entropy per site = 0.08). We have looked at one



domain of this RT protein, known as PF06817, the so-called RT thumb domain, composed of a four-helix bundle. In HIV data, the first 10–15 sites of the domain, not counting gaps, tend to be quite conserved, while the latter part is more variable. The resulting entropy is very low. Conversely, the HMM profile shows much less conservation. The full alignment on the PFAM database contains sequences from many different viruses, so this might also contribute to the observed variability (especially if viewing the representative proteomes on PFAM). It might be the case that, while (at least part of) this protein is well conserved in HIV, it is not as conserved across many different viruses. RT is not thought to be often targeted by human immune responses, so that will contribute to the reduced variability in HIV, in addition to functional constraints. Intuitively we would expect that this protein as a whole should be functionally constrained, but perhaps either the constraints are virus-specific, or the frequency of targeting by the immune factor is the dominant reason why it appears more conserved.

While the entropy computed in the presence of high phylogenetic correlations is not representative of the diversity in the protein family which may be observed across distant organisms, it can be used to characterize the constraints acting on the different sites of a given protein, on the different proteins of HIV. In particular we have computed the cost in entropy corresponding to fixing the amino acid content on one site, e.g., to its consensus value. While this cost is close to the entropy of the single-site amino-acid frequencies for most sites, the two quantities differ on some sites, which signals the presence of strong coupling effects (epistasis). Computing the entropy cost offers another potential avenue to investigate the fitness landscape of the virus. Sites associated to high entropies are likely to be the sites of escape mutations for the virus, in response to host immune pressure. Note that, from a computational point of view, the limited variability in the MSA helps for the inference of the Potts model. The system is, in physical jargon, in a paramagnetic phase with large local fields, and the Independent Site Model already provides a good starting point for the inference.

In the artificial lattice-based protein models we have studied, the entropy is very large,  $\sigma \simeq 3$ , due to the extremely reduced protein worlds we have considered (only a few proteins coexist and compete), in order to be amenable to exact calculations. Calculations taking into all the possible competing structures on the cubic lattice show a drastic reduction in the entropy per site, and give  $\sigma \simeq 1.8 - 2.1$  [47], a value close to the one found for real protein families. It is important to underline that, while the lattice-protein model does not contain only 2-body interactions, the true entropy is very accurately recovered with the pairwise Potts model, see Table 7.

An important question is whether our values for the entropy can be confronted to experiments. In directed evolution experiments, starting from a pool of random sequences, sequences are selected according to their *in vitro* fitness, such as binding affinity against a target. The fittest sequences are mutated, amplified, and another round of selection can take place. One fundamental issue is the size of the initial pool of sequences allowing for the selection of (at least one) fit protein(s). In one experiment [49] Keefe and Szostak started from a pool of  $6 \times 10^{12}$  proteins with 80 amino acids each, and selected them according to their ATP binding affinity. After 4 cycles of selection and mutation (made possible by the RNA tags attached to the proteins) they found 4 different sequences of new ATP binding proteins. The authors estimate that 1 in  $10^{11}$  random sequences has ATP-binding activity comparable to the one isolated in the study. Assuming that this ratio corresponds to the ratio of the number of proteins in the ‘ATP-binding family’ over the number of sequences with 80 amino acids, we obtain that the entropy of this putative family is  $S = \ln(10^{-11} \times 20^{80}) \simeq 214.3$ . The entropy per site is therefore  $\sigma \simeq 2.67$ . This estimate is large compared to the values we have found in the analysis of the natural protein families in the present work. One possible

explanation is that the definition of ‘ATP-binding family’ is actually too loose compared to natural families, which would lead to high apparent values for the entropy. We believe that further work to connect estimates of the entropy and in vitro directed evolution experiments in a quantitative way would be very useful.

Last of all, while we have considered here the entropy of the distribution of amino acid sequences, we should not forget that those sequences are coded at the DNA level by nucleotides. The redundancy of the genetic code adds extra entropy to the value we have computed. This additive contribution depends on the amino acid content, as the degeneracy of amino acids varies from 1 to 6. In addition, it also depends on the organisms and on the tissue where the protein are expressed through the codon bias. More subtle effects, e.g., resulting from the pressure exerted by the innate immune system, also limit the diversity of the nucleotide sequences at fixed amino-acid content [50].

**Acknowledgments** S.C., H.J. and R.M. were partly funded by the Agence Nationale de la Recherche Coevstat project (ANR-13-BS04-0012-01).

## References

1. Durbin, R., Sean Eddy, R., Krogh, A., Mitchison, G.: *Biological Sequence Analysis Probabilistic Models of Proteins and Nucleic Acids*. Cambridge University Press, London (1998)
2. Ashkenazy, H., Erez, E., Martz, E., Pupko, T., Ben-Tal, N.: ConSurf 2010: calculating evolutionary conservation in sequence and structure of proteins and nucleic acids. *Nucl. Acids Res.* **38**, W529–W533 (2010)
3. Lapedes, A.S., Giraud, B.G., Liu, L., Stormo, G.D.: Correlated mutations in models of protein sequences: phylogenetic and structural effects. *Lect. Notes-Monogr. Ser.* **33**, 236–256 (1999)
4. Rausell, A., Juan, D., Pazos, F., Valencia, A.: Protein interactions and ligand binding: from protein subfamilies to functional specificity. *Proc. Natl. Acad. Sci.* **107**(5), 1995–2000 (2010)
5. Pazos, F., Helmer-Citterich, E., Ausiello, G., Valencia, A.: Correlated mutations contain information about protein-protein interaction. *J. Mol. Biol.* **271**, 511–523 (1997)
6. de Juan, D., Pazos, F., Valencia, A.: Emerging methods in protein co-evolution. *Nat. Rev. Genet.* **14**, 249–261 (2013)
7. Berman, H.M., Kleywegt, G.J., Nakamura, H., Markley, J.L.: The protein data bank at 40: reflecting on the past to prepare for the future. *Structure* **20**(3), 391–396 (2012)
8. The UniProt Consortium: Reorganizing the protein space at the Universal Protein Resource (UniProt). *Nucl. Acids Res.* **40**, D71 (2012)
9. Punta, M., Coggill, P.C., Eberhardt, R.Y., Mistry, J., Tate, J.G., Boursnell, C., Pang, N., Forslund, K., Ceric, G., Clements, J., Heger, A., Holm, L., Sonnhammer, E.L.L., Eddy, S.R., Bateman, A.I., Finn, R.D.: The Pfam protein families database. *Nucl. Acids Res.* **40**, D290 (2012)
10. Jaynes, E.T.: On the rationale of maximum-entropy methods. *Proc. IEEE* **70**(9), 939–952 (1982)
11. Bialek, William: *Biophysics: Searching for Principles*. Princeton University Press, Princeton (2012)
12. Weigt, Martin, White, Robert A., Szurmant, Hendrik, Hoch, James A., Hwa, Terence: Identification of direct residue contacts in protein-protein interaction by message passing. *Proc. Natl. Acad. Sci. USA* **106**(1), 67–72 (2009)
13. Burger, L., van Nimwegen, E.: Disentangling Direct from Indirect Co-Evolution of Residues in Protein Alignments. *PLoS Comput. Biol.* **6**, E1000633 (2010)
14. Balakrishnan, S., Kamisetty, H., Carbonell, J.G., Lee, S.I., Langmead, C.J.: Learning generative models for protein fold families. *Proteins: Struct. Funct. Bioinf.* **79**, 1061 (2011)
15. Cocco, Simona, Monasson, Rémi: Adaptive cluster expansion for inferring Boltzmann machines with noisy data. *Phys. Rev. Lett.* **106**, 090601 (2011)
16. Cocco, Simona, Monasson, Rémi: Adaptive cluster expansion for the inverse Ising problem: convergence, algorithm and tests. *J. Stat. Phys.* **147**(2), 252–314 (2012)
17. Shakhnovich, E., Gutin, A.: Enumeration of all compact conformations of copolymers with random sequence of links. *J. Chem. Phys.* **93**, 5967–5971 (1990)
18. Shakhnovich, E.: Protein design: a perspective from simple tractable models. *Fold. Des.* **3**, R45–R58 (1998)

19. Finn, Robert D., Mistry, Jaina, Tate, John, Coghill, Penny, Heger, Andreas, Pollington, Joanne E., Luke Gavin, O., Gunasekaran, Prasad, Ceric, Goran, Forslund, Kristoffer, Holm, Liisa, Sonnhammer, Erik L.L., Eddy, Sean R., Bateman, Alex: The pfam protein families database. *Nucl. Acids Res.* **38**(suppl 1), D211–D222 (2010)
20. Barton, J.P., Cocco, S., De Leonardi, E., Monasson, R.: Large pseudocounts and L2-norm penalties are necessary for the mean-field inference of Ising and Potts models. *Phys. Rev. E* **90**(1), 012132 (2014)
21. Morcos, F., Pagnani, A., Lunt, B., Bertolino, A., Marks, D.S., Sander, C., Zecchina, R., Onuchic, J.N., Hwa, Terence, Weigt, Martin: Direct-coupling analysis of residue coevolution captures native contacts across many protein families. *Proc. Natl. Acad. Sci. USA* **108**(49), E1293–E1301 (2011)
22. Ekeberg, M., Lovkvist, C., Lan, Y., Weigt, M., Aurell, E.: Improved contact prediction in proteins: using pseudolikelihoods to infer potts models. *Phys. Rev. E* **87**, 012707 (2013)
23. Cocco, S., Monasson, R., Weigt, M.: From principal component to direct coupling analysis of coevolution in proteins: Low-eigenvalue modes are needed for structure prediction. *PLoS Comput. Biol.* **9**, E1003176 (2013)
24. Russ, W., Lowery, D.M., Mishra, P., Yaffe, M.B., Ranganathan, R.: Natural-like function in artificial WW domains. *Nature* **437**, 579–583 (2005)
25. Socolich, Michael, Lockless, Steve W., Russ, William P., Lee, Heather, Gardner, Kevin H., Ranganathan, Rama: Evolutionary information for specifying a protein fold. *Nature* **437**(7058), 512–518 (2005)
26. Korber, Bette, Gaschen, Brian, Yusim, Karina, Thakallapally, Rama, Keşmir, Can, Detours, Vincent: Evolutionary and immunological implications of contemporary HIV-1 variation. *Br. Med. Bull.* **58**(1), 19–42 (2001)
27. Ferguson, Andrew L., Mann, Jaclyn K., Omarjee, Saleha, Ndung'u, Thumbi, Walker, Bruce D., Chakraborty, Arup K.: Translating HIV sequences into quantitative fitness landscapes predicts viral vulnerabilities for rational immunogen design. *Immunity* **38**(3), 606–617 (2013)
28. Mann, Jaclyn K., Barton, John P., Ferguson, Andrew L., Omarjee, Saleha, Walker, Bruce D., Chakraborty, Arup K., Ndung'u, Thumbi: The fitness landscape of HIV-1 Gag: advanced modeling approaches and validation of model predictions in vitro testing. *PLoS Comput. Biol.* **10**(8), e1003776 (2014)
29. Haq, Omar, Andrej, Michael, Morozov, Alexandre V., Levy, Ronald M.: Correlated electrostatic mutations provide a reservoir of stability in HIV protease. *PLoS Comput. Biol.* **8**(9), e1002675 (2012)
30. Flynn, William F., Chang, Max W., Tan, Zhiqiang, Oliveira, Glenn, Yuan, Jinyun, Okulicz, Jason F., Torbett, Bruce E., Levy, Ronald M.: Deep sequencing of protease inhibitor resistant HIV patient isolates reveals patterns of correlated mutations in gag and protease. *PLoS Comput. Biol.* **11**(4), e1004249 (2015)
31. Shekhar, K., Ruberman, C.F., Ferguson, A.L., Barton, J.P., Kardar, M., Chakraborty, A.K.: Spin models inferred from patient-derived viral sequence data faithfully describe HIV fitness landscapes. *Phys. Rev. E* **88**(6), 062705 (2013)
32. Addo, M.M., Yu, X.G., Rathod, A., Eldridge, R.L., Strick, D., Johnston, M.N., Corcoran, C., Fitzpatrick, C.A., Feeney, M.E., Rodriguez, W.R., Basgoz, N., Draenert, R., Brander, C., Goulder, P.J.R., Rosenberg, E.S., Altfeld, Marcus, Walker, Bruce D.: Comprehensive epitope analysis of human immunodeficiency virus type 1 (HIV-1)-specific T-cell responses directed against the entire expressed HIV-1 genome demonstrate broadly directed responses, but no correlation to viral load. *J. Virol.* **77**(3), 2081–2092 (2003)
33. Streeck, H., Jolin, J.S., Qi, Ying, Yassine-Diab, B., Johnson, R.C., Kwon, D.S., Addo, M.M., Brumme, C., Routy, J.P., Little, S., Jessen, H.K., Kelleher, A.D., Hecht, F.M., Sekaly, R.P., Rosenberg, E.S., Walker, Bruce D., Carrington, Mary, Altfeld, Marcus: Human immunodeficiency virus type 1-specific CD8+ T-cell responses during primary infection are major determinants of the viral set point and loss of CD4+ T cells. *J. Virol.* **83**(15), 7641–7648 (2009)
34. Zhao, Gongpu, Perilla, Juan R., Yufenyuy, Ernest L., Meng, Xin, Chen, Bo, Ning, Jiying, Ahn, Jinwoo, Gronenborn, Angela M., Schulten, Klaus, Aiken, Christopher, et al.: Mature hiv-1 capsid structure by cryo-electron microscopy and all-atom molecular dynamics. *Nature* **497**(7451), 643–646 (2013)
35. Dahirel, V., Shekhar, K., Florencia, P., Miura, T., Artyomov, M., Talsania, S., Allen, T.M., Altfeld, M., Carrington, M., Irvine, D.J., Walker, B.D., Chakraborty, A.K.: Coordinate linkage of HIV evolution reveals regions of immunological vulnerability. *Proc. Natl. Acad. Sci.* **108**(28), 11530–11535 (2011)
36. Barton, John P., Kardar, Mehran, Chakraborty, Arup K.: Scaling laws describe memories of host pathogen riposte in the HIV population. *Proc. Natl. Acad. Sci.* **112**(7), 1965–1970 (2015)
37. Beitzel, B.F., Bakken, R.R., Smith, J.M., Schmaljohn, C.S.: High-resolution functional mapping of the venezuelan equine encephalitis virus genome by insertional mutagenesis and massively parallel sequencing. *PLoS Pathog.* **6**(10), e1001146 (2010)
38. Heaton, Nicholas S., Sachs, David, Chen, Chi-Jene, Hai, Rong, Palese, Peter: Genome-wide mutagenesis of influenza virus reveals unique plasticity of the hemagglutinin and ns1 proteins. *Proc. Natl. Acad. Sci.* **110**(50), 20248–20253 (2013)

39. Remenyi, R., Qi, H., Su, S.Y., Chen, Z., Wu, N.C., Arumugaswami, V., Truong, S., Chu, V., Stokelman, T., Lo, H.H., Olson, A., Wu, T.T., Chen, S.H., Lin, C.Y., Sun, R.: A comprehensive functional map of the hepatitis c virus genome provides a resource for probing viral proteins. *mBio* **5**, e01469-14 (2014)
40. Fulton, B.O., Sachs, D., Beaty, S.M., Won, S.T., Lee, B., Palese, P., Heaton, N.S.: Mutational analysis of measles virus suggests constraints on antigenic variation of the glycoproteins. *Cell Rep.* **11**(9), 1331–1338 (2015)
41. Ferrari, Guido, Korber, Bette, Goonetilleke, Nilu, Liu, Michael K.P., Turnbull, Emma L., Salazar-Gonzalez, Jesus F., Hawkins, Natalie, Self, Steve, Watson, Sydeaka, Betts, Michael R., Gay, Cynthia, McGhee, Cynthia, Pellegrino, Pierre, Williams, Ian, Tomaras, Georgia D., Haynes, Barton F., Gray, Clive M., Borrow, Persephone, Roederer, Mario, McMichael, Andrew J., Weinhold, Kent J.: Relationship between functional profile of HIV-1 specific CD8 T cells and epitope variability with the selection of escape mutants in acute HIV-1 infection. *PLoS Pathog.* **7**(2), e1001273 (2011)
42. Liu, M.K.P., Hawkins, N., Ritchie, A.J., Ganusov, V.V., Whale, V., Brackenridge, S., Li, H., Pavlicek, J.W., Cai, F., Rose-Abrahams, M., Treurnicht, F., Hraber, P., Riou, C., Gray, C., Ferrari, G., Tanner, R., Ping, L.H., Anderson, J.A., Swanstrom, R., Cohen, M., Abdool Karim, S.S., Haynes, B., Borrow, P., Perelson, A.S., Shaw, G.M., Hahn, B.H., Williamson, C., Korber, B.T., Gao, F., Self, S., McMichael, A., Goonetilleke, N.: Vertical T cell immunodominance and epitope entropy determine HIV-1 escape. *J. Clin. Investig.* **123**(1), 380–393 (2013)
43. Li, H., Helling, R., Tang, C., Wingreen, N.: Emergence of preferred structures in a simple model of protein folding. *Science* **273**, 666–669 (1996)
44. Li, H., Tang, C., Wingreen, N.: Designability of protein structures: a lattice-model study using the miyazawa-jernigan matrix. *Proteins* **49**, 403–412 (2002)
45. England, Jeremy L., Shakhnovich, Eugene I.: Structural determinant of protein designability. *Phys. Rev. Lett.* **90**, 218101 (2003)
46. Miyazawa, A., Jernigan, R.: Estimation of effective interresidue contact energies from protein crystal structures: quasi-chemical approximation. *Macromolecules* **18**, 534 (1985)
47. Jacquin, H., Gilson, A., Shakhnovich, E., Cocco, S., Monasson, R.: Benchmarking inverse statistical approaches for protein structure and design with exactly solvable models. available on Biorxiv, 2015. doi:[10.1101/028936](https://doi.org/10.1101/028936)
48. Berezovsky, I.N., Zeldovich, K.B., Shakhnovich, E.: Positive and negative design in stability and thermal adaptation of natural proteins. *PLoS Comput. Biol.* **3**(32), e52 (2007)
49. Keefe, Anthony, Szostak, W.Jack: Functional proteins from a random-sequence library. *Nature* **410**(6829), 715–718 (2001)
50. Greenbaum, B., Cocco, S., Levine, A., Monasson, R.: A quantitative theory of entropic forces acting on constrained nucleotide sequences applied to viruses. *Proc. Natl. Acad. Sci. USA* **111**, 5054–5059 (2014)