

Predicting Evolutionary Site Variability from Structure in Viral Proteins: Buriedness, Packing, Flexibility, and Design

Amir Shahmoradi · Dariya K. Sydykova · Stephanie J. Spielman ·
Eleisha L. Jackson · Eric T. Dawson · Austin G. Meyer · Claus O. Wilke

Received: 23 April 2014 / Accepted: 31 August 2014 / Published online: 13 September 2014
© Springer Science+Business Media New York 2014

Abstract Several recent works have shown that protein structure can predict site-specific evolutionary sequence variation. In particular, sites that are buried and/or have many contacts with other sites in a structure have been shown to evolve more slowly, on average, than surface sites with few contacts. Here, we present a comprehensive study of the extent to which numerous structural properties can predict sequence variation. The quantities we considered include buriedness (as measured by relative solvent accessibility), packing density (as measured by contact number), structural flexibility (as measured by B factors, root-mean-square fluctuations, and variation in dihedral angles), and variability in designed structures. We obtained structural flexibility measures both from molecular dynamics simulations performed on nine non-homologous viral protein structures and from variation in homologous variants of those proteins, where they were available. We obtained measures of variability in designed structures from flexible-backbone design in the Rosetta software. We found that most of the structural properties correlate with site variation in the majority of structures, though the

correlations are generally weak (correlation coefficients of 0.1–0.4). Moreover, we found that buriedness and packing density were better predictors of evolutionary variation than structural flexibility. Finally, variability in designed structures was a weaker predictor of evolutionary variability than buriedness or packing density, but it was comparable in its predictive power to the best structural flexibility measures. We conclude that simple measures of buriedness and packing density are better predictors of evolutionary variation than the more complicated predictors obtained from dynamic simulations, ensembles of homologous structures, or computational protein design.

Introduction

Patterns of amino-acid sequence variation in protein-coding genes are shaped by the structure and function of the expressed proteins (Wilke and Drummond 2010; Liberles et al. 2012; Marsh and Teichmann 2014). As the most basic reflection of this relationship, buried residues in proteins tend to be more evolutionarily conserved than exposed residues (Overington et al. 1992; Goldman et al. 1998; Mirny and Shakhnovich 1999; Dean et al. 2002). More specifically, when evolutionary variation is plotted as a function of Relative Solvent Accessibility (RSA, a measure of residue buriedness), the relationship falls, on average, onto a straight line with a positive slope (Franzosa and Xia 2009; Ramsey et al. 2011; Franzosa and Xia 2012; Scherrer et al. 2012). Importantly, however, this relationship represents on an average many sites and many proteins. At the level of individual sites in individual proteins, RSA is often only weakly correlated with evolutionary variation (Meyer and Wilke 2013; Meyer et al. 2013; Yeh et al. 2014b).

Electronic supplementary material The online version of this article (doi:10.1007/s00239-014-9644-x) contains supplementary material, which is available to authorized users.

A. Shahmoradi
Department of Physics, The University of Texas at Austin,
Austin, TX 78712, USA

A. Shahmoradi · D. K. Sydykova · S. J. Spielman ·
E. L. Jackson · E. T. Dawson · A. G. Meyer · C. O. Wilke (✉)
Department of Integrative Biology, Center for Computational
Biology and Bioinformatics, and Institute for Cellular
and Molecular Biology, The University of Texas at Austin,
Austin, TX 78712, USA
e-mail: wilke@austin.utexas.edu

Other structural measures, such as residue contact number (CN), have also been shown to correlate with sequence variability (Liao et al. 2005; Franzosa and Xia 2009; Yeh et al. 2014b), and some have argued that CN predicts evolutionary variation better than RSA (Yeh et al. 2014b, a). Because CN may be a proxy for residue and site-specific backbone flexibility (Halle 2002), a positive trend between local structural variability and sequence variability may also exist (Yeh et al. 2014b). Indeed, several authors have suggested that such protein dynamics may play a role in sequence variability (Liu and Bahar 2012; Nevin Gerek et al. 2013; Marsh and Teichmann 2014). However, a recent paper argued against the flexibility model, on the grounds that evolutionary rate is not linearly related to flexibility (Huang et al. 2014).

While RSA and CN can be calculated in a straightforward manner from individual crystal structures, measures of structural flexibility, either at the side-chain or the backbone level, are more difficult to obtain. Two viable approaches to measuring structural flexibility are (i) examining existing structural data or (ii) simulating protein dynamics. NMR ensembles may approximate physiologically relevant structural fluctuations. Similar fluctuations are observed in ensembles of homologous crystal structures (Maguida et al. 2008; Echave and Fernández 2010). The thermal motion of atoms in a crystal is recorded in B factors, which is available for every atom in every crystal structure. To measure protein fluctuations using a simulation approach, one can either use coarse-grained modeling, e.g., via Elastic Network Models (Sanejouand 2013), or atom-level modeling, e.g., via molecular dynamics (MD) (Karplus and McCammon 2002). However, it is not well understood which, if any, of these measures of structural flexibility provide insight into the evolutionary process, particularly into residue-specific evolutionary variation.

Here, we provide a comprehensive analysis of the extent to which numerous different structural quantities predict evolutionary sequence (amino-acid) variation. We considered two measures of evolutionary sequence variation: site entropy, as calculated from homologous protein alignments, and evolutionary rate. As structural predictors, we included buriedness (RSA), packing density (CN), and measures of structural flexibility, including B factors, several measures of backbone and side-chain variability obtained from MD simulations, and backbone variability obtained from alignments of homologous crystal structures. We additionally considered site variability, as predicted from computational protein design with Rosetta.

On a set of nine viral proteins, RSA and CN generally performed better at predicting evolutionary site variation than either measures of structural flexibility or computational protein design. Among the measures of structural flexibility, measures of side-chain variability performed

better than measures of backbone variability, possibly because the former are more tightly correlated with residue packing. Finally, site variability predicted from computational protein design performed worse than the best-performing measures of structural fluctuations.

Materials and Methods

Sequence Data, Alignments, and Evolutionary Rates

All viral sequences except influenza sequences were retrieved from <http://hfv.lanl.gov/components/sequence/HCV/search/searchi.html>. The sequences were truncated to the desired genomic region but did not restrict in any other way. Influenza sequences were downloaded from <http://www.fludb.org/brc/home.spg?decorator=influenza>. We only considered human influenza A, H1N1, excluding H1N1 sequences derived from the 2009 Swine Flu outbreak or any sequence before 1998, but we did not place any geographic restrictions.

For all viral sequences, we removed any sequence that was not in reading frame, any sequence which was shorter than 80 % of the longest sequence for a given viral protein (so as to remove all partial sequences), and any sequence containing any ambiguous characters. Alignments were constructed using amino-acid sequences with MAFFT (Katoh et al. 2002, 2005), specifying the—auto flag to select the optimal algorithm for the given data set, and then back-translated to a codon alignment using the original nucleotide sequence data.

To assess site-specific sequence variability in amino-acid alignments, we calculated the Shannon entropy (H_i) at each alignment column i :

$$H_i = - \sum_j P_{ij} \ln P_{ij}, \quad (1)$$

where P_{ij} is relative frequency of amino acid j at position i in the alignment.

For each alignment, we also calculated evolutionary rates, as described (Spielman and Wilke 2013). In brief, we generated a phylogeny for each codon alignment in RAxML (Stamatakis 2006) using the GTRGAMMA model. Using the codon alignment and phylogeny, we inferred evolutionary rates with a Random Effects Likelihood (REL) model, using the HyPhy software (Kosakovsky Pond et al. 2005). The REL model was a variant of the GY94 evolutionary model (Goldman and Yang 1994) with five ω rate categories as free parameters. We employed an Empirical Bayes approach (Yang 2000) to infer ω values for each position in the alignment. These ω values represent the evolutionary-rate ratio dN/dS at each site.

Table 1 PDB structures considered in this study

Viral protein	PDB ID	Chain	Sequence length	Number of sequences
Hemagglutinin precursor	1RD8	AB	503	1039
Dengue protease helicase	2JLY	A	451	2362
West Nile protease	2FP7	B	147	237
Japanese encephalitis helicase	2Z83	A	426	145
Hepatitis C protease	3GOL	A	557	1021
Rift valley fever nucleoprotein	3LYF	A	244	95
Crimean Congo nucleocapsid	4AQF	B	474	69
Marburg RNA binding domain	4GHA	A	122	42
Influenza nucleoprotein	4IRY	A	404	943

Protein Crystal Structures

A total of nine viral protein structures were selected for analysis, as tabulated in Table 1. Sites in the PDB structures were mapped to sites in the viral sequence alignments via a custom-built python script that creates a consensus map between a PDB sequence and all sequences in an alignment.

For each of the viral proteins, homologous structures were identified using the `blast.pdb` function of the R package Bio3D (Grant et al. 2006). BLAST hits were retained if they had $\geq 35\%$ sequence identity and $\geq 90\%$ alignment length. Among the retained hits, we subsequently identified sets of homologous structures with unique sequences and with mutual pairwise sequence divergences of $\geq 2\%$, $\geq 5\%$, and $\geq 10\%$.

Molecular Dynamics Simulations

Molecular dynamics (MD) simulations were carried out using the GPU implementation of the *Amber12* simulation package (Salomon-Ferrer et al. 2013) with the most recent release of the Amber fixed-charge force field (ff12SB; c.f., AmberTools13 Manual). Prior to MD production runs, all PDB structures were first solvated in a box of TIP3P water molecules (Jorgensen et al. 1983) such that the structures were at least 10\AA away from the box walls. Each individual system was then energy minimized using the steepest descent method for 1,000 steps, followed by conjugate gradient for another 1000 steps. Then, the structures were constantly heated from 0 to 300 K for 0.1ns, followed by 0.1ns constant pressure simulations with positional

harmonic restraints on all atoms to avoid instabilities during the equilibration process. The systems were then equilibrated for another 5ns without positional restraints, each followed by 15ns of production simulations for subsequent post-processing and analyses. All equilibration and production simulations were run using the SHAKE algorithm (Ryckaert et al. 1977). Langevin dynamics were used for temperature control.

Measures of Buriedness, Packing Density, and Structural Flexibility

As a measure of residue buriedness, we calculated Relative Solvent Accessibility (RSA). To calculate RSA, we first calculated the Accessible Surface Area (ASA) for each residue in each protein, using the DSSP software (Kabsch and Sander 1983). We then normalized ASA values by the theoretical maximum ASA of each residue (Tien et al. 2013) to obtain RSA. We considered two measures of local packing density, contact number (CN), and weighted contact number (WCN). We calculated CN for each residue as the total number of $C\alpha$ atoms surrounding the $C\alpha$ atom of the focal residue within a spherical neighborhood of a predefined radius r_0 . Following Yeh et al. (2014b), we used $r_0 = 13\text{\AA}$. We calculated WCN as the total number of surrounding $C\alpha$ atoms for each focal residue, weighted by the inverse square separation between the $C\alpha$ atoms of the focal residue and the contacting residue, respectively (Shih et al. 2012).

In most analyses, we actually used the inverse of CN and/or WCN, $iCN = 1/CN$ and $iWCN = 1/WCN$. Note that for Spearman correlations, which we use throughout here, replacing a variable by its inverse changes the sign of the correlation coefficient but not the magnitude.

As measures of structural flexibility, we considered RMSF, variability in backbone and side-chain dihedral angles, and B factors. We calculated RMSF for $C\alpha$ atoms based on both MD trajectories and homologous crystal structures. For MD trajectories, we calculated RMSF as

$$\text{RMSF}_j = \left[\sum_i \left(\mathbf{r}_i^{(j)} - \mathbf{r}_0^{(j)} \right)^2 \right]^{1/2} \quad (2)$$

where RMSF_j is the root-mean-square fluctuation at site j , $\mathbf{r}_i^{(j)}$ is the position of the $C\alpha$ atom of residue j at MD frame i , and $\mathbf{r}_0^{(j)}$ is the position of the $C\alpha$ atom of residue j in the original crystal structure.

To calculate RMSF from homologous structures, we first aligned the structures using the Bio3D package (Grant et al. 2006), and then we calculated

$$\text{RMSF}_j = \left[\sum_i w_i \left(\mathbf{r}_i^{(j)} - \langle \mathbf{r}^{(j)} \rangle \right)^2 \right]^{1/2}, \quad (3)$$

where $\mathbf{r}_i^{(j)}$ now stands for the position of the C α atom of residue j in structure i , $\langle \mathbf{r}^{(j)} \rangle$ is the mean position of that C α atom over all aligned structures, and w_i is a weight to correct for potential phylogenetic relationship among the aligned structures. The weights w_i were calculated using BranchManager (Stone and Sidow 2007), based on phylogenies built with RAxML as before.

To assess variability in backbone and side-chain dihedral angles, we calculated $\text{Var}(\phi)$, $\text{Var}(\psi)$, and $\text{Var}(\chi_1)$. The variance of a dihedral angle was defined according to the most common definition in directional statistics: First, a unit vector \mathbf{x}_i is assigned to each dihedral angle α_i in the sample. The unit vector is defined as $\mathbf{x}_i = (\cos(\alpha_i), \sin(\alpha_i))$. The variance of the dihedral angle is then defined as

$$\text{Var}(\alpha) = 1 - \|\langle \mathbf{x} \rangle\|, \quad (4)$$

where $\|\langle \mathbf{x} \rangle\|$ represents the length of the mean $\langle \mathbf{x} \rangle$, calculated as $\langle \mathbf{x} \rangle = \sum_i \mathbf{x}_i / n$. Here, n is the sample size. The variance of a dihedral angle is, by definition, a real number in the range $[0, 1]$, with $\text{Var}(\alpha) = 0$ corresponding to the minimum variability of the dihedral angle and $\text{Var}(\alpha) = 1$ to the maximum, respectively (Berens 2009). Since the χ_1 angle is undefined for Ala and Gly, we excluded all sites with these residues in analyses involving χ_1 .

B factors were extracted from the crystal structures. We only considered the B factors of the C α atom of each residue.

Sequence Entropy from Designed Proteins

Designed entropy was calculated as described (Jackson et al. 2013). In brief, proteins were designed using RosettaDesign (Version 39284) (Leaver-Fay et al. 2011) using a flexible-backbone approach. This was done for all PDB structures in Table 1 as initial template structures. For each template, we created a backbone ensemble using the Backrub method (Smith and Kortemme 2008). The temperature parameter in Backrub was set to 0.6, allowing for an intermediate amount of flexibility. We had previously found in a different data set that intermediate flexibility gave the highest congruence between designed and observed site variability (Jackson et al. 2013).

For each of the nine template structures we designed 500 proteins.

Availability of Data and Methods

All details of simulations, input/output files, and scripts for subsequent analyses are available to view or

download at https://github.com/clauswilke/structural_prediction_of_ER.

Results

Data Set and Structural Variables Considered

Our goal in this work was to determine which structural properties best predict amino-acid variability at individual sites in viral proteins. To this end, we selected nine viral proteins for which we had both high-quality crystal structures and abundant sequences to assess evolutionary sequence variation (Table 1). We quantified evolutionary variability in two ways: by calculating sequence entropies for each alignment column, and by calculating site-specific evolutionary-rate ratios $\omega = dN/dS$ (see Methods for details). Throughout this paper, we primarily report results obtained for sequence entropy. Results for ω were largely comparable, with some specific caveats detailed below.

As predictors of evolutionary variability, we considered buriedness, packing density, and residue flexibility. We additionally considered the variation seen in computationally designed protein variants. Buriedness quantifies the extent to which a residue is protected from solvent. We determined residue buriedness by calculating the relative solvent accessibility (RSA), which represents the relative proportion of a residue's surface in contact with solvent.

Packing density quantifies how many other residues a given residue interacts with. We determined packing density by calculating contact number (CN) and weighted contact number (WCN). CN counts the number of contacts within a sphere of a given radius around the α -carbon of the focal residue, while WCN weights contacts by the distance between the two residues. Residue buriedness and packing density tend to be (anti-)correlated but measure qualitatively different properties of a residue. In particular, in the core of a protein, buriedness is always zero but packing density can vary. Because contact numbers decline as relative solvent accessibility increases, we replaced CN and WCN with their inverses, $iCN = 1/CN$ and $iWCN = 1/WCN$, in most analyses. Importantly, as Spearman rank correlations were used, this substitution only changed the sign of correlations but not the magnitude.

Measures of structural flexibility assess the extent to which a residue fluctuates in space as a protein undergoes thermodynamic fluctuations in solution. We quantified these fluctuations using several different measures. We considered B factors, which measure the spatial localization of individual atoms in a protein crystal, RMSF, the root-mean-square fluctuation of the C α atom over time, and variability in side-chain and backbone dihedral angles,

Table 2 Availability of homologous crystal structures

Viral Protein	BLAST hits ^a	Unique sequences			
		all	≥ 2% ^b	≥ 5% ^b	≥ 10% ^b
Hemagglutinin precursor	63	17	10	9	7
Dengue protease helicase	31	13	7	7	7
West Nile protease	21	16	10	7	6
Japanese encephalitis helicase	31	12	7	7	7
Hepatitis C protease	302	33	10	5	4
Rift valley fever nucleoprotein	95	9	5	5	5
Crimean Congo nucleocapsid	7	4	3	2	2
Marburg RNA binding domain	63	9	5	3	3
Influenza nucleoprotein	69	15	4	4	2

Although most viral proteins have many PDB structures available, the sequence divergence among these structures is low. Therefore, when calculating RMSF from crystal structures, we considered only those proteins with at least five homologous structures at 5 % pairwise sequence divergence (highlighted in bold).

^a BLAST hits against all sequences in the PDB, excluding hits with <35 % sequence identity and <90 % alignment length

^b Unique sequences at indicated minimum pairwise sequence divergence

including $\text{Var}(\chi_1)$, $\text{Var}(\phi)$, and $\text{Var}(\psi)$. We employed two broad approaches, one using PDB crystal structures and one using molecular dynamics (MD) simulations, to obtain these measurements. Crystal structures yielded measures for B factors and RMSF; we obtained B factors from individual protein crystal structures, given in Table 1, and we calculated RMSF from aligned homologous crystal structures for those proteins which had sufficient sequence variation among crystal structures (see Methods and Table 2 for details). MD simulations yielded measures for RMSF and variability in residue dihedral and side-chain angles. More specifically, we simulated MD trajectories for all crystal structures in Table 1. For each protein, we equilibrated the structure, simulated 15ns of chemical time, and recorded snapshots of the simulated structure every 10ps (see Methods for details). We obtained RMSF and angle variabilities from these snapshots. Additionally, we calculated time-averaged values of RSA, CN, and WCN. We also refer to these time-averaged measures as MD RSA, MD CN, and MD WCN, respectively. Unless specified otherwise, all results reported below were obtained using MD RSA, MD CN, and MD WCN.

As an alternative to predicting evolutionary variation from simple structural measures such as contact density or

backbone flexibility, one can also predict evolutionary variation via a protein-design approach (Dokholyan and Shakhnovich 2001; Ollikainen and Kortemme 2013; Jackson et al. 2013). In this case, one takes the protein structure of interest, replaces all residue side chains with randomly chosen alternatives, and uses a coarse-grained or atom-level energy function to assess which side-chain choices are consistent with the backbone conformation of the focal structure. We have recently used this approach to compare natural and designed sequence variability in cellular proteins (Jackson et al. 2013), and we have found that (i) flexible-backbone design, where small backbone movements are allowed during the design phase, outperformed fixed-backbone design, and (ii) intermediate backbone flexibility, obtained via an intermediate design temperature, produced the highest congruence between designed and natural sequences. Similarly, Dokholyan and Shakhnovich (2001) had previously found that an intermediate temperature parameter gave the best agreement between designed and natural sequences in their model. Inspired by these prior results, we investigated here how protein design performed relative to simpler structural quantities. For all proteins in our study (Table 1), we used the Rosetta protein-design platform (Leaver-Fay et al. 2011) to generate 500 designed variants. We then calculated the sequence entropy at each alignment position of the designed variants. We refer to the resulting quantity as the *designed entropy*. We chose a design temperature of $T = 0.6$, which was near the optimal range in our previous work (Jackson et al. 2013).

Evaluating Structural Predictors of Evolutionary Sequence Variation

We began by comparing the Spearman correlations of sequence entropy with six different measures of local structural flexibility: B factors, RMSF obtained from MD simulations (MD RMSF), and RMSF obtained from crystal structures (CS RMSF), and variability in backbone and side-chain dihedral angles (ϕ , ψ , and χ_1). The correlation strengths of these quantities with entropy are shown in Fig. 1. Significant correlations ($P < 0.05$) are shown with filled symbols, and non-significant correlations are shown with empty symbols ($P \geq 0.05$). We found that the variability in backbone dihedral angles, $\text{Var}(\phi)$ and $\text{Var}(\psi)$, explained the least variation in sequence entropy, while the variability in the side-chain dihedral angle, $\text{Var}(\chi_1)$, explained, on average, more variation in sequence entropy than did any other measure of structural flexibility. B factors and the two measures of RMSF explained, on average, approximately the same amount of variation in entropy, even though the results for individual proteins were somewhat discordant (see also next sub-section).

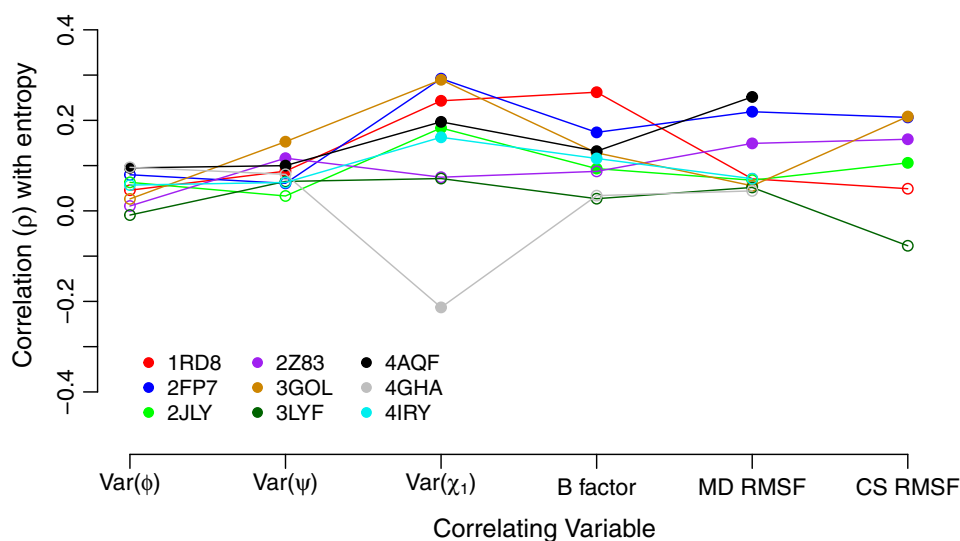


Fig. 1 Spearman correlation of sequence entropy with measures of structural variability. Each *symbol* represents one correlation coefficient for one protein structure. Significant correlations ($P < 0.05$) are shown as *filled symbols*, and insignificant correlations ($P \geq 0.05$) are shown as *open symbols*. The quantities $\text{Var}(\psi)$, $\text{Var}(\phi)$, $\text{Var}(\chi_1)$, and MD RMSF were obtained as time-averages over 15ns of MD

simulations. B factors were obtained from individual crystal structures. CS RMSF values were obtained from alignments of homologous crystal structures when available. Almost all structural measures of variability correlate weakly, but significantly, with sequence entropy

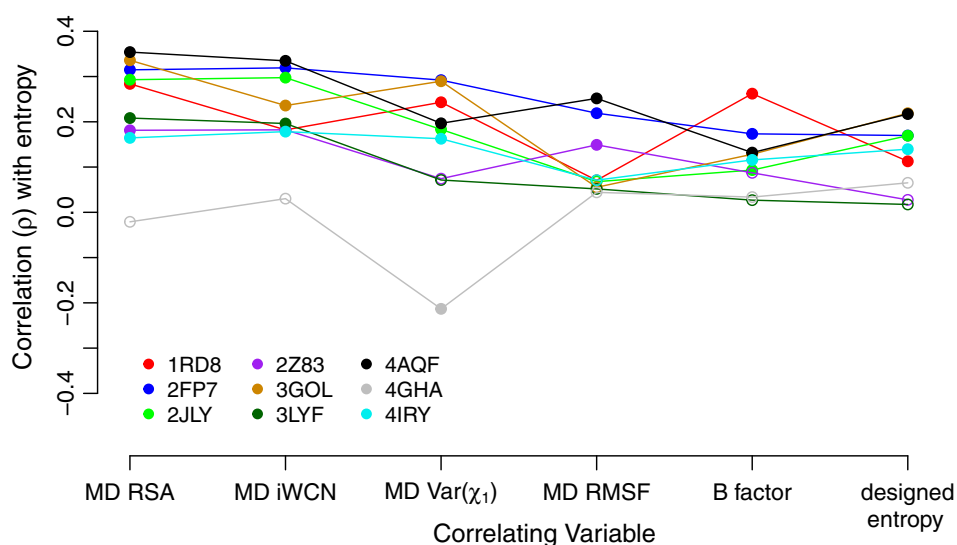


Fig. 2 Spearman correlation of sequence entropy with measures of buriedness, packing density, and structural flexibility, as well as with designed entropy. Each *symbol* represents one correlation coefficient for one protein structure. Significant correlations ($P < 0.05$) are shown as *filled symbols*, and insignificant correlations ($P \geq 0.05$) are shown as *open symbols*. The quantities MD RSA, MD iWCN, MD $\text{Var}(\chi_1)$, and MD RMSF were calculated as time-averages over 15ns of MD

simulations. B factors were obtained from crystal structures, and designed entropy was obtained from protein design in Rosetta. Compared to the measures of structural variability and to designed entropy, MD RSA and MD iWCN consistently show stronger correlations with sequence entropy. Note that results for MD iWCN are largely identical to those for MD iCN, so only MD iWCN was included here

Based on results from the above analysis, we proceeded to compare the relative explanatory power among the best-performing measures of structural flexibility ($\text{Var}(\chi_1)$, MD RMSF, and B factors) with buriedness (RSA), packing

density (iWCN), and designed entropy. Figure 2 shows the Spearman correlation coefficients between sequence entropy and each of the aforementioned quantities, for all proteins in our analysis. In this figure, several patterns

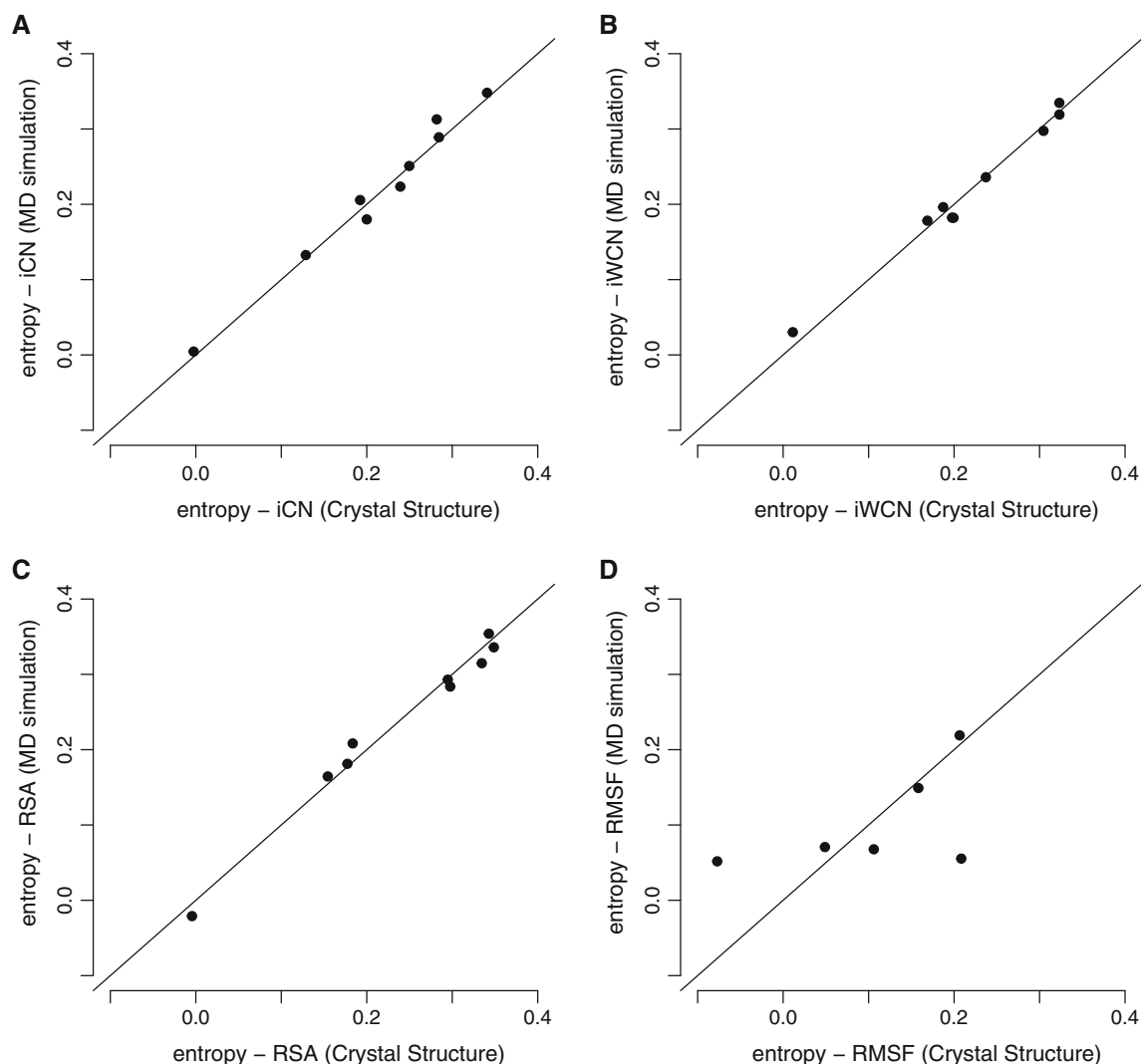


Fig. 3 Spearman correlations of sequence entropy with MD-derived and crystal-structure derived structural measures. The *vertical* axes in all *plots* represent the Spearman correlation of sequence entropy with one structural variable obtained from 15ns of molecular dynamics (MD) simulations. The *horizontal* axes represent the Spearman's rank correlation coefficient of sequence entropy with the same structural variable as in the *vertical* axes but measured from protein crystal

structures. Each *dot* represents one correlation coefficient for one protein structure. The quantities iCN, iWCN, and RSA have nearly identical predictive power for sequence entropy regardless of whether they are derived from MD simulations or from crystal structures. By contrast, MD RMSF yielded very different correlations than did CS RMSF

emerged. First, nearly all correlations were positive and most were statistically significant, with the main exception of the Marburg virus RNA binding domain (PDB ID 4GHA). This protein only showed a single significant negative correlation between sequence entropy and $\text{Var}(\chi_1)$. Second, correlations were generally weak, such that no correlation coefficient exceeded 0.4. Third, on average, correlations were strongest for RSA and iWCN, yielding average correlations of $\rho = 0.23$ and $\rho = 0.22$, respectively. Fourth, designed entropy performed worse than RSA or iWCN as a predictor of evolutionary sequence variability, but it performed roughly the same as the three flexibility measures in this figure; the values of designed

entropy, $\text{Var}(\chi_1)$, MD RMSF, and B factors showed average correlations of $\rho = 0.13$, $\rho = 0.14$, $\rho = 0.11$, and $\rho = 0.12$, respectively.

MD Time-Averages Versus Crystal-Structure Snapshots

Except for analyses involving B factors and CS RMSF, we obtained structural measures by averaging quantities over MD trajectories. This approach, however, did not reflect conventional practice for measuring RSA, CN, or WCN, which are typically measured from individual crystal structures. Therefore, we examined whether MD time-

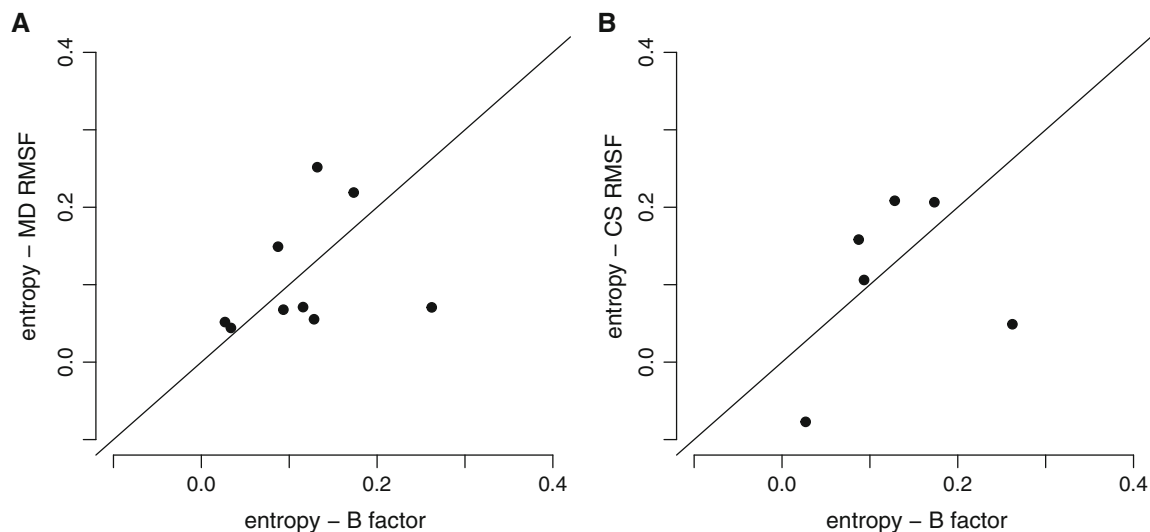


Fig. 4 Spearman correlations of sequence entropy with measures of structural variability. *Vertical* and *horizontal* axes represent Spearman correlations of the indicated quantities. Each *dot* represents one

correlation coefficient for one protein structure. MD RMSF, CS RMSF, and B factors all explain different amounts of variance in sequence entropy for different proteins

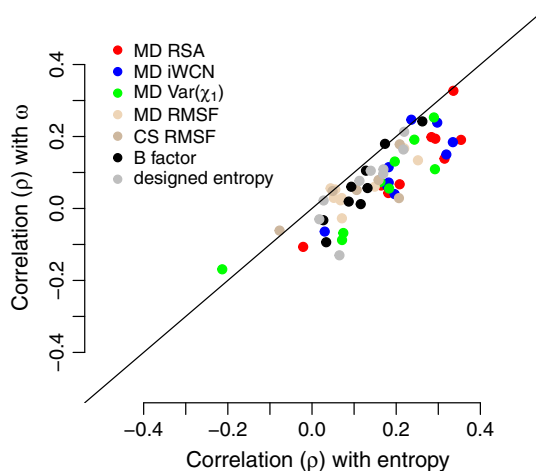


Fig. 5 Spearman correlations of structural quantities with sequence entropy and with the evolutionary-rate ratio ω . Nearly all points fall below the $x = y$ line, indicating that structural quantities generally predict as much as or more variation in sequence entropy than in ω

averages differed in any meaningful way from estimates obtained from crystal structures, and whether these estimates differed in their predictive power for evolutionary sequence variation.

As shown in Table 3, RSA, CN, and WCN from crystal structures were highly correlated with their corresponding MD trajectory time-averages, for all protein structures we examined (Spearman correlation coefficients of > 0.9 in all cases). Furthermore, the correlation coefficients we obtained when comparing the crystal-structure based measures to sequence entropy were virtually identical to coefficients obtained from the MD trajectory correlations

Table 3 Correlations between quantities obtained from MD trajectories and from crystal structures

Quantity	min ρ	max ρ	$\langle \rho \rangle$	SD (ρ)
RSA	0.937	0.981	0.948	0.012
CN	0.964	0.993	0.976	0.008
WCN	0.973	0.991	0.984	0.006
RMSF	0.218	0.723	0.502	0.181

For each quantity and each protein, we calculated the Spearman correlation ρ between the values obtained from MD time-averages and the values obtained from viral protein crystal structures. Note that crystal structures for all nine proteins were used for RSA, CN, and WCN calculations, but only the six proteins for which we had sufficient crystal structure variability were used for CS RMSF. We then calculated the minimum, maximum, mean, and standard deviation of these correlations

(Fig. 3a–c). Thus, in terms of predicting evolutionary variation, RSA, CN, and WCN values obtained from static structures performed as well as their MD equivalents averaged over short time scales.

By contrast, correlations between corresponding MD RMSF to CS RMSF measures were sometimes quite different, with correlation coefficients ranging from 0.218 to 0.723 (Table 3). Consequently, for the two proteins for which MD RMSF was the least correlated with CS RMSF (hepatitis C protease and Rift Valley fever nucleoprotein), the strength of correlation between site entropy and RMSF depended substantially on how RMSF was calculated (Figs. 1 and 3d).

Finally, we examined whether correlations between sequence entropy and B factors or the two RMSF measures were comparable (Fig. 4). Again, we found that

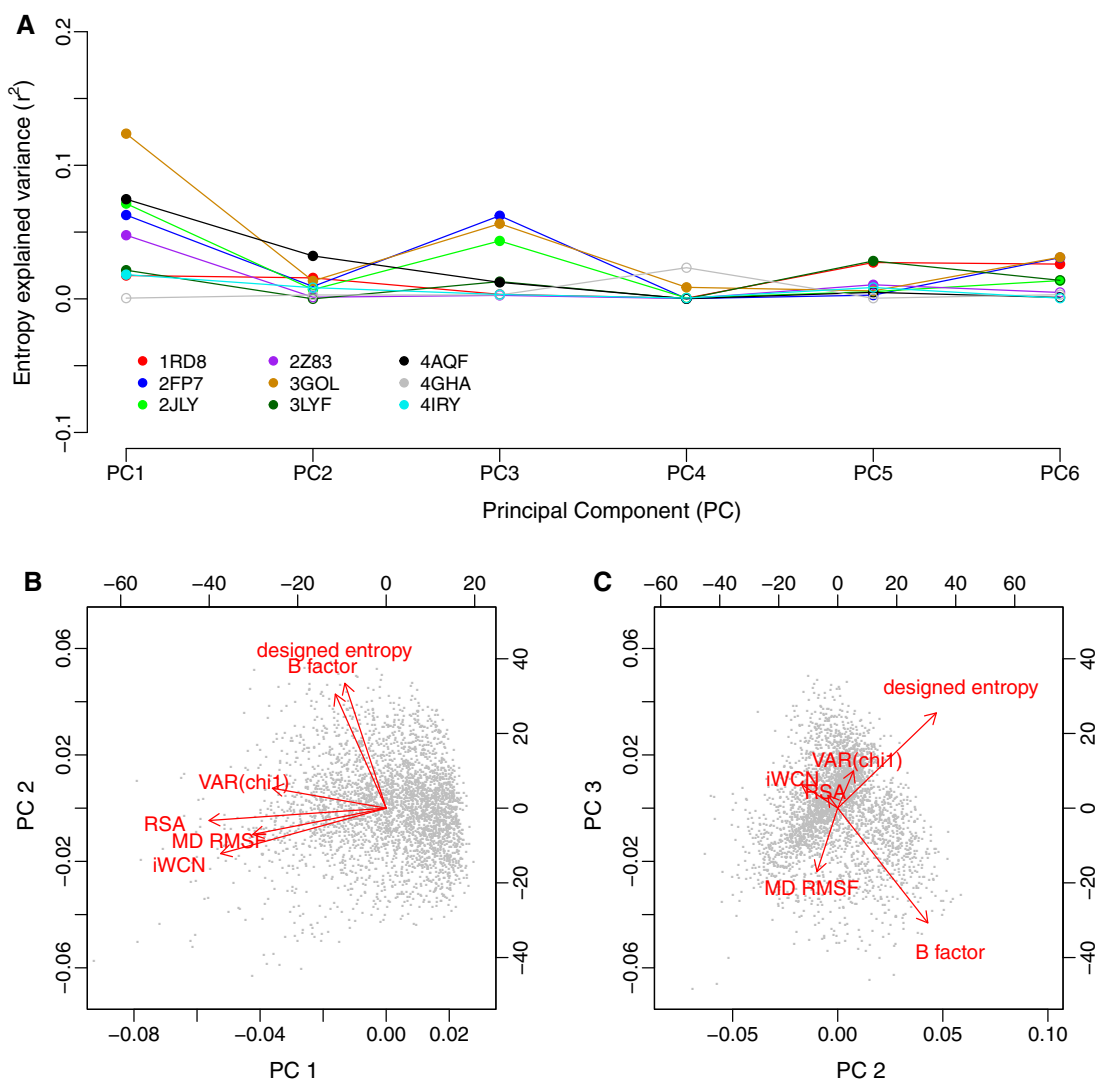


Fig. 6 Principal Component (PC) Regression of sequence entropy against the structural variables. **a** Variance in entropy explained by each principal component. For most proteins, PC1 and PC3 show the strongest correlations with sequence entropy. Significant correlations ($P < 0.05$) are shown as filled symbols, and insignificant correlations ($P \geq 0.05$) are shown as open symbols. **b** and **c** Composition of the three leading components. Red arrows represent the loadings of each

of the structural variables on the principal components; black dots represent the amino acid sites in the PC coordinate system. The variables RSA, iWCN, MD RMSF, and $\text{Var}(\chi_1)$ load strongly on PC1 and weakly on PC2, while B factor and designed entropy load strongly on PC2 and weakly on PC1. Interestingly, B factor and designed entropy also load strongly on PC3, but in opposite directions

correlations between sequence entropy and B factors were generally different from those obtained for both MD RMSF and CS RMSF. This result highlighted that, while B factors, MD RMSF, and CS RMSF all measure backbone flexibility, they each contain distinct information about evolutionary sequence variability in our data set.

Sequence Entropy Versus Evolutionary-Rate Ratio ω

In the previous subsections, we used sequence entropy as a measure of site-wise evolutionary variation. While sequence entropy is a simple and straightforward measure of site variability, it has two potential drawbacks. First,

while measured from homologous protein alignments, sequence entropy doesn't correct for the phylogenetic relationship of those alignment sequences. Hence, entropy can be biased if some parts of the phylogeny are more densely sampled than others. Second, entropy does not take the actual substitution process into account. As a result, a single substitution near the root of the tree can result in a comparable entropy to a sequence of substitutions toggling back and forth between two amino acids.

To consider an alternative quantity of evolutionary variation that doesn't suffer from either of these drawbacks, we calculated the evolutionary-rate ratio $\omega = dN/dS$ for all proteins at all sites, and repeated all analyses

with ω instead of entropy. We found that results generally carried over, but with somewhat weaker correlations. Figure 5 plots, for each protein, the Spearman correlations between ω and our various predictors versus the correlation between entropy and our predictors. Most data points fall below the $x = y$ line and are shifted downwards by approximately 0.1. Thus, correlations of structural quantities and designed entropy with ω are, on average, approximately 0.1 smaller than correlations of the same quantities with sequence entropy.

Multi-Variate Analysis of Structural Predictors

The various structural quantities we have considered are by no means independent of each other. Measures of buriedness and packing density co-vary with each other, as do measures of structural flexibility. Further, the latter co-vary with the former, as does designed entropy. Therefore, we conducted a joint multivariate analysis, which included most structural quantities considered in this work. We employed this strategy to determine the extent to which these quantities contained independent information about sequence variability while additionally assessing whether combining multiple structural quantities yielded improved predictive power. We employed a principal component (PC) regression approach, which has previously been used successfully to disentangle genomic predictors of whole-protein evolutionary rates (Drummond et al. 2006; Bloom et al. 2006). For each analysis described below, we first carried out a PC analysis of the predictor variables (i.e., the structural quantities such as RSA and RMSF), and we subsequently regressed the response (either sequence entropy or ω) against the individual components. Note that variables were not rank-transformed for this analysis.

For a first PC analysis, we pooled all structural quantities and then regressed entropy against each PC separately, for all proteins in our data set. This strategy allowed us to analyze all proteins in our data set individually but in such a way that results were comparable from one protein to the next. We excluded CS RMSF from this analysis, so that we could include results from all nine viral proteins. The results of this analysis are shown in Fig. 6. The first component (PC1) explained, on average, the largest amount of variation in sequence entropy (see Fig. 6a). PC3 yielded the second-highest r^2 value, on average, while all other components explained very little variation in sequence entropy. When looking at the composition of the components, we found that RSA, iWCN, RMSF, and $\text{Var}(\chi_1)$ all loaded strongly on PC1, while PC2 and PC3 were primarily represented by designed entropy and B factors (see Fig. 6b and c). RMSF also had moderate loadings on PC3. Interestingly, designed entropy and B

factors load with equal signs on PC2 but with opposite signs on PC3.

We interpreted PC1 to represent a buriedness/packing-density component. By definition, PC1 measures the largest amount of variation among the structural quantities, and all structural quantities reflect to some extent the buriedness of residues and the number of residue-residue contacts. PC2 and PC3 were more difficult to interpret. Since designed entropy and B factors loaded strongly on both but with two different combinations of signs, we concluded that the most parsimonious interpretation was to consider PC2 as a component representing sites with high designed entropy and high spatial fluctuations (as measured by B factors) and PC3 representing sites with high designed entropy and low spatial fluctuations. Using these interpretations, our PC regression analysis suggested that of all the structural quantities considered here, residue buriedness/packing was the best predictor of evolutionary variation. Designed entropy was a useful predictor as well, but it tended to perform better at sites with low spatial fluctuations.

For a second PC analysis, we included the predictor CS RMSF, which therefore restricted the data set to include only six proteins (see Table 2). This analysis, which retained sequence entropy as the response variable, yielded comparable results to the first PC analysis. The main differences occurred in PC2 and PC3, where CS RMSF generally loaded in the opposite direction of B factor, and either in the same (PC2) or the opposite (PC3) direction of designed entropy (Fig. S1).

Finally, we redid the two PC analyses described above, but instead with ω as the response variable (Figs. S2 and S3). Again, these results were largely comparable to results from PC analyses with sequence entropy as the response.

Discussion

We have carried out a comprehensive analysis of the extent to which different structural quantities predict sequence evolutionary variation in nine viral proteins. We found that measures of buriedness and local packing generally performed better than measures of structural flexibility. Further, the former measures also performed better than a computational protein-design approach that employed a sophisticated all-atom force field to determine allowed amino-acid distributions at each site. Finally, there was no difference in predictive power between structural quantities obtained from averaging structural quantities over 15ns of MD simulations versus taking the same quantities from individual crystal structures.

Our results are broadly in agreement with recent work by Echave and collaborators (Yeh et al. 2014b; Huang et al.

2014). These authors found that RSA and CN showed comparable correlation strengths with evolutionary sequence variation (Yeh et al. 2014b). Further, they demonstrated that the observed relationship between evolutionary variation and residue–residue contacts was not consistent with a flexibility model that puts evolutionary variability in proportion to structural flexibility (Huang et al. 2014). Instead, a mechanistic stress model, in which amino-acid substitutions cause physical stress in proportion to the number of residue–residue contacts affected, could explain all the observed data (Huang et al. 2014).

The correlation strengths we observed were consistently lower than those observed previously (Jackson et al. 2013; Yeh et al. 2014b). We believe that this result was due to our choice of analyzing viral proteins instead of the cellular proteins or enzymes used in prior works. First, while viral sequences are abundant, their alignments may not be as diverged as alignments that can be obtained for sequences from cellular organisms. For example, our influenza sequences spanned only approximately one decade. Despite the high mutation rates observed in RNA viruses, the evolutionary variation that can accumulate over this time span is limited. This relatively lower evolutionary divergence makes resolving differences between more and less conserved sites much more difficult. Second, many viral proteins experience a substantial amount of selection pressure to evade host immune responses. The resulting positive selection on viral sequences may mask evolutionary constraints imposed by structure. For example, influenza hemagglutinin displays positive selection throughout the entire sequence, regardless of the extent of residue burial (Meyer and Wilke 2013; Meyer et al. 2013; Suzuki 2006; Bush et al. 1999). However, the results we obtained here for viral proteins are broadly consistent with the results obtained earlier for cellular proteins (Dokholyan and Shakhnovich 2001; Franzosa and Xia 2009; Jackson et al. 2013; Yeh et al. 2014b), indicating that viral proteins evolve under many of the same biophysical selection pressures that cellular proteins experience.

We have found here that correlations between sequence entropy and structural quantities were consistently higher than correlations between the evolutionary-rate ratios ω and structural quantities. Surprisingly, in a recent study on cellular proteins, Yeh et al. (2014a) found that entropy performed worse than quantities assessing substitution rates. One possible explanation for this discrepancy is again our choice of viral sequences. Our sequence alignments almost certainly contained some polymorphisms, whereas the sequences of Yeh et al. (2014a) likely did not. It is known that polymorphisms may diminish the reliability of ω estimates (Kryazhimskiy and Plotkin 2008). While the effect of polymorphisms on sequence entropy is not known, it seems plausible that entropy would be less

sensitive to them than ω is. Alternatively, since viral proteins frequently experience positive selection, rate estimates may be confounded by this selection pressure and thus less reflective of constraints imposed by protein structure. By contrast, even under positive selection amino-acid distributions at sites would have to be consistent with the constraints imposed by the protein structure, and entropy would remain sensitive to these constraints.

We found that simple measures of buriedness or packing density, such as RSA or CN, were better predictors of evolutionary variation than was sequence variability predicted from computational protein design. In other words, simple quantities that can be obtained trivially from PDB structures performed better than a sophisticated protein-design strategy that makes use of an all-atom energy function and requires thousands of CPU-hours to complete. This result highlights that, even though computational protein design has yielded impressive results in specific cases (Kuhlman et al. 2003; Röthlisberger et al. 2008; Fleishman et al. 2011), this approach remains limited in its ability to predict evolutionary variation. Similarly, we have previously found that flexible-backbone design with Rosetta produced designs whose surface and core were too similar (Jackson et al. 2013). We attributed this discrepancy to either the solvation model or the model of backbone flexibility we used (Backrub, see Smith and Kortemme 2008). The results we found here suggest that the model of backbone flexibility may indeed be the cause of at least some of the discrepancies between predicted and observed site variability. In particular, in our PC regression analysis, the component in which designed entropy loaded opposite to B factor and MD RMSF generally had the second-highest predictive power for evolutionary variability, after the component representing buriedness/packing density. In sum, designed entropy was a better predictor for evolutionary sequence variability for sites with less structural flexibility compared to sites with more flexibility.

Even though RSA and CN remain the best currently known predictors of evolutionary variation, neither quantity has particularly high predictive power. One reason why predictive power may be low is that neither quantity accounts for correlated substitutions at interacting sites. Yet such correlated substitutions happen regularly. For example, covariation among sites encodes information about residue–residue contacts and 3D structure (Halabi et al. 2009; Burger and Nimwegen 2010; Marks et al. 2011; Jones et al. 2014), and evolutionary models that incorporate residue–residue interactions tend to perform better than models that do not (Rodrigue et al. 2005; Bordner and Mittelman 2014). An improved predictor of evolutionary variation would have to correctly predict this covariation from structure. In principle, computational protein design, which takes into consideration the atom-level details of the