

Protein Functional Site Prediction Using a Conservative Grade and a Proximate Grade

Yosuke Kondo and Satoru Miyazaki*

Department of Medicinal and Life Science, Faculty of Pharmaceutical Sciences, Tokyo University of Science, Japan

Abstract

So far, in order to predict important sites of a protein, many computational methods have been developed. In the era of big-data, it is required for improvements and sophistication of existing methods by integrating sequence data in the structural data. In this paper, we aim at two things: improving sequence-based methods and developing a new method using both sequence and structural data. Therefore, we developed an originally modified evolutionary trace method, in which we defined conservative grades calculated from a given multiple sequence alignment and a proximate grade in order to evaluate predicted active sites from a viewpoint of protein-ion, protein-ligand, protein-nucleic acid, protein-protein interaction by use of three-dimensional structures. In other words, the proximate grade also can evaluate an amino acid residue. When we applied our method to translation elongation factor Tu/1A proteins, it showed that the conservative grades are evaluated accurately by the proximate grade. Consequently, our idea indicated two advantages. One is that we can take into account various cocrystal structures for evaluation. Another one is that, by calculating the fitness between the given conservative grade and the proximate grade, we can select the best conservative grade.

Keywords: Evolutionary trace; Three-dimensional structure; Elongation factor

Introduction

When a protein works, a specific site to bind an ion or a molecule may exist. Identification of binding sites is important to investigate how the protein works and binds ions or molecules. In order to identify such an important site, it is necessary to prepare a mutant type of the protein, whose amino acid residue is mutated into another one, and then a difference of binding affinity between the mutant type and the wild type is investigated. However, mutating amino acid residues one by one takes an amount of time and costs. Therefore, it is effective for developing a method to narrow down the amino acid residues.

For electing the candidate sites, there are many computational methods, which are based on (i) sequence, (ii) structure and (iii) sequence and structure [1-6]. Sequence-based methods usually assume that such an important site is conservative against mutation and therefore important sites and others should have been mutated in different patterns. In order to detect such patterns, various methods have been developed [7]. One of the sequence-based methods is a method based on Shannon entropy (SE) [8,9]. However, the SE-based method may have three problems. The first one is that the SE-based method, in which twenty standard amino acids are regarded as characters, does not consider properties of amino acids. Therefore, a method based on SE of residue properties [10] or a sum of pairs [11] was proposed. The second one is that the SE-based method does not consider a background distribution of amino acids. Therefore, other information-theoretical method such as relative entropy [12] or Jensen Shannon divergence [13] was proposed. The third one is that the SE-based method, in which a rate of an amino acid is calculated, cannot take into account which amino acid is included in a sequence. Therefore, some methods based on windowing [13], weighting [14] or phylogenetic analysis was proposed. One of the methods based on a phylogenetic tree is an evolutionary trace (ET) method [15], which has been extended as weighted ET (WET) [16], integer-valued ET (iv-ET) and real-valued ET (rv-ET) methods [17]. Additionally, other methods based on phylogenetic trees are ConSurf [18] and Rate4Site [19,20] algorithms.

Although a variety of sequence-based methods have been already compared each other [13,21], what difference makes a difference is

difficult to understand because such methods do not be explained by an idea. Therefore, we consider a map, a mathematical formula, on a multiple sequence alignment (MSA) and aim at constructing an exhaustive method. As part of this effort, we propose a method currently including some existing methods such as the method based on SE or SE of residue properties, the method based on a sum of pairs with/without weighting and the iv-ET or the rv-ET method.

Even if a variety of methods are executable, how are the methods evaluable? There may exist two approaches: confirmation by site-directed mutagenesis and visualization onto a three-dimensional structure. The former is more consistent with identification of binding sites because the latter is verifiable that a site is proximate from ions or molecules. In spite of that, the latter has been still used because of indefinability of protein functional sites. Therefore, on the basis of benchmark sets such as catalytic sites, ligand-binding sites or protein-protein interfaces [13], the predictive ability has been evaluated. However, the latter is immature because of usually conducting only a structure [15,22]. This mainly causes two problems. The first one is that the latter neglects a protein which binds various ions or molecules because an entry in the Protein Data Bank (PDB) [23] does not always include all states of the protein structure. The second one is that the latter cannot take account of proteins which are derived from an ancestor. Therefore, protein structures derived from different organisms are incomparable with each other. To solve these problems, we consider another map, which measures proximity of amino acid residues and ions or molecules, and then two maps are integrated.

***Corresponding author:** Satoru Miyazaki, Department of Medicinal and Life Science, Faculty of Pharmaceutical Sciences, Tokyo University of Science, Yamazaki 2641, Noda-shi, Chiba 278-8510, Japan, Tel: +81-4-7121-3630; E-mail: smiyazak@rs.noda.tus.ac.jp

Received June 21, 2015; Accepted July 09, 2015; Published July 16, 2015

Citation: Kondo Y, Miyazaki S (2015) Protein Functional Site Prediction Using a Conservative Grade and a Proximate Grade. J Data Mining Genomics Proteomics 6: 175. doi:10.4172/2153-0602.1000175

Copyright: © 2015 Kondo Y, et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Mathematical Formulation of Mappings of an MSA

Notation of fundamental elements

Let $M=(mij)$ denote a given MSA and here mij denote an amino acid symbol of site j on sequence i in the MSA. Let ${}_kM = [m_{1k}, m_{2k}, \dots, m_{nk}]^t$ be column k on the MSA and we consider a mapping

$$f_x : M \rightarrow [0, \infty].$$

Mapping by a character type

In this section, we define mathematical formulation of a mapping by similarity of the amino acid symbols on ${}_iM$. Let ${}_iM \in {}_iM$ denote ${}_iM$ at time point $t = 1, 2, \dots, N + 1$, where N is a number of internal nodes on a phylogenetic tree reconstructed from the given MSA (Figure 1A), and be represented by a field of sets. For example, ${}_4M$;

$$\begin{aligned} {}_1M &= \{\{R, R, L, R, R, R\}\} \\ {}_2M &= \{\{R, R, L, R, R\}, \{R\}\} \\ {}_3M &= \{\{R, R, L, R\}, \{R\}, \{R\}\} \\ {}_4M &= \{\{R, R, L\}, \{R\}, \{R\}, \{R\}\} \\ {}_5M &= \{\{R, R\}, \{L\}, \{R\}, \{R\}, \{R\}\} \end{aligned}$$

Let there be $g_x : {}_iM \rightarrow [0, 1]$, which here maps to 1 if there exists ${}_iM_u \in {}_iM$ which comprises two or more types of characters and, in other cases, maps to 0. By $g_x({}_iM)$ only, ${}_1M - {}_5M$ are indistinguishable. If $g_x({}_1M)$ and $g_x({}_2M)$ are summed, the others are distinguishable. If $g_x({}_1M)$, $g_x({}_2M)$ and $g_x({}_3M)$ are summed, ${}_1M$ and ${}_2M$ and the others are distinguishable. Therefore, let

$$f_x({}_iM) := \sum_{t=1}^T g_x({}_iM) \quad (1)$$

Where $T=1, 2, \dots, N$.

As shown in Figure 1B, let $h_x : {}_iM \rightarrow [0, 1]$ be included in g_x and $g_x({}_iM)$ be represented as following three definitions:

$$g_1({}_iM) := \begin{cases} 0 & (\forall {}_iM_u \in {}_iM, h_x({}_iM_u) \leq \tau) \\ 1 & (\exists {}_iM_u \in {}_iM, h_x({}_iM_u) > \tau) \end{cases} \quad (2)$$

Where τ is a threshold of $h_x({}_iM)$,

$$g_2({}_iM) := \frac{1}{|{}_iM|} \sum_{{}_iM_u \in {}_iM} h_x({}_iM_u) \quad (3)$$

Where $|{}_iM|$ is a number of multisets in ${}_iM$ and

$$g_3({}_iM) := h_x({}_iM_{*u}) \quad (4)$$

Where ${}_iM_{*u}$ is a multiset which is separated at time point $t + 1$. For example, ${}_4M_{*u}$;

$$\begin{aligned} {}_4M_{*u} &= \{R, R, L, R, R, R\} \\ {}_4M_{*u} &= \{R, R, L, R, R\} \\ {}_4M_{*u} &= \{R, R, L, R\} \\ {}_4M_{*u} &= \{R, R, L\} \\ {}_4M_{*u} &= \{R, R\} \end{aligned}$$

Let A denote a field of sets of amino acid symbols and ${}_iG \subset {}_iM$ denote a field of sets of gaps in ${}_iM$. For example, A is definable as

$${}^{20}A = \{\{A\}, \{C\}, \{D\}, \{E\}, \{F\}, \{G\}, \{H\}, \{I\}, \{K\}, \{L\}, \{M\}, \{N\}, \{P\}, \{Q\}, \{R\}, \{S\}, \{T\}, \{V\}, \{W\}\}$$

$$\text{Or } {}^9A = \{\{M, L, V, I\}, \{H, R, K\}, \{S, T\}, \{A, G\}, \{D, E\}, \{Q, N\}, \{F, W, Y\}, \{P\}, \{C\}\}$$

and ${}_iG$ is definable as ${}_iG = \{\{^1\gamma, ^2\gamma, \dots, ^G\gamma\}\}$ or ${}_iG = \{\{^1\gamma\}, \{^2\gamma\}, \dots, \{^G\gamma\}\}$ where G is a number of gaps.

Let $h_x({}_iM_u)$ be represented as following four definitions:

$$h_x({}_iM_u) := \begin{cases} 0 & (\forall l \in {}_iM_u, \exists X \in A \cup {}_iG; l \in X) \\ 1 & (\text{otherwise}) \end{cases} \quad (5)$$

where $|A \cup {}_iG|$ is a number of sets in $A \cup {}_iG$ and

$$p({}_iM_u, X) = \frac{1}{|{}_iM_u|} \sum_{l \in {}_iM_u} \begin{cases} 0 & (l \notin X) \\ 1 & (l \in X) \end{cases} \quad (7)$$

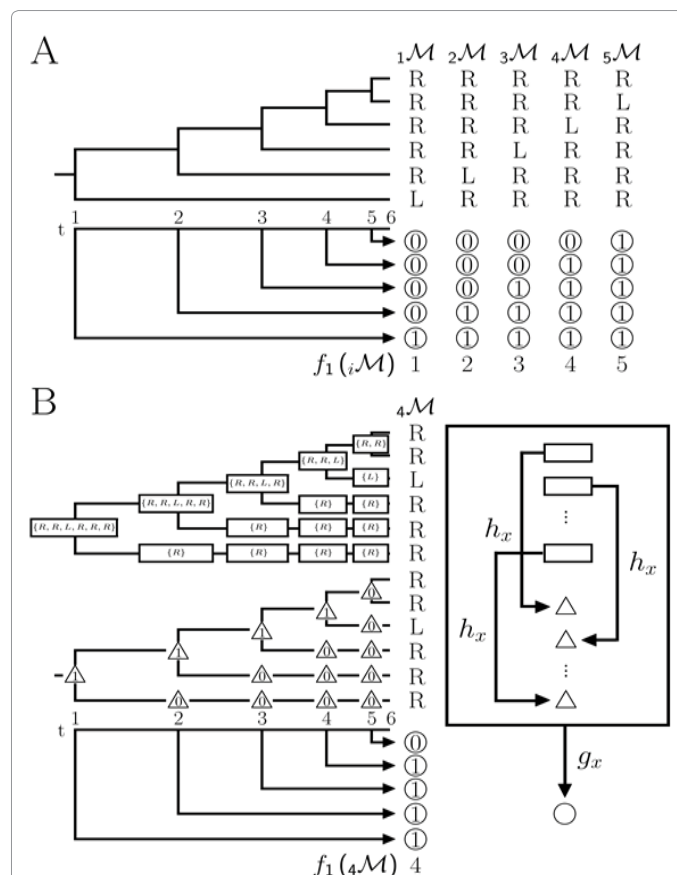


Figure 1: A Concept of a mapping by a character type. (A) A concept of $f_1({}_1M - {}_5M)$ are comprised of 5 R and 1 L and each character attaches a leaf node of a rooted phylogenetic tree under a hypothesis that the evolutionary rate is constant. Numbers in ascending order are assigned from the root to leaf nodes as time point t . In f_1 , after a value in a circle is assigned to ${}_iM$, values in circles are summed. (B) Concepts of g_x and h_x . In g_x , after h_x maps characters in a square to a value in a triangle, values in triangles are mapped to a value in a circle.

where $|^iM_u|$ is a number of characters in iM_u and if $p(^iM_u, X) = 0$, $p(^iM_u, X) \log_{|A \cup G|} p(^iM_u, X)$ is regarded as 0,

$$h_3(^iM_u) := \frac{1}{|^iM_u|^2} \sum_{l \in ^iM_u} \sum_{m \in ^iM_u} s(l, m) \quad (8)$$

Where

$$s(l, m) = \begin{cases} 0 (l \in X \in {}_iG \wedge m \in Y \in {}_iG \wedge X = Y) \\ \frac{S_{\max} - S_{\min}}{S_{\max}} (l \in X \in {}_iG \wedge m \in Y \in {}_iG \wedge X \neq Y) \\ \frac{S_{\max} - S_{\min}}{S_{\max}} (l \in X \in {}_iG \wedge m \in Y \in A) \\ \frac{S(l, l) - S_{\min}}{S(l, l)} (l \in X \in A \wedge m \in Y \in {}_iG) \\ \frac{S(l, l) - S(l, m)}{S(l, l)} (l \in X \in A \wedge m \in Y \in A) \end{cases} \quad (9)$$

Where S_{\max} , S_{\min} , $S(l, l)$ and $S(l, m)$ are the maximum, the minimum, a diagonal element and an off-diagonal element in an amino acid substitution matrix, respectively, and

$$h_4(^iM_u) := \frac{1}{|^iM_u|^2} \sum_{l \in ^iM_u} \sum_{m \in ^iM_u} s(l, m) w(l) \quad (10)$$

where w is a weight of sequence l .

Mapping by a coordinate type

Let \mathbb{R} denote a set of real numbers and there be $e: \mathbb{R}^3 \times \mathbb{R}^3 \rightarrow [0, \infty)$. Let $R \subset \mathbb{R}^3$, $Q \subset \mathbb{R}^3$ and

$$e(R, Q) := \min_{(r, q) \in R \times Q} (\|r - q\|_2) \quad (11)$$

where $\|\cdot\|_2$ is an Euclidean norm.

Let us consider structure k , which contains a protein and ions or molecules. Let ${}^kR \subset \mathbb{R}^3$ denote atomic coordinates of amino acid residue i in structure k and ${}^kQ \subset \mathbb{R}^3$ denote atomic coordinates of ions or molecules in structure k . Let K denote a number of structures and the sequences are aligned. Let $\{^iR, {}^2R, \dots, {}^K R\} \subseteq {}_iM$ denote a set of residues in iM denote a set of residues in iM and $\{^1\gamma, {}^2\gamma, \dots, {}^G\gamma\} = {}_iG \subset {}^iM$ denote a set of gaps in iM . Let

$$f_2({}^iM) := \min_{{}^kR \in {}_iM, {}^kQ \in G} [e({}^kR, {}^kQ)] \quad (12)$$

Materials and Methods

Data collection

In UniProtKB/Swiss-Prot release 2015_01 [24], entries which are annotated as 'Classic translation factor GTPase family. EF-Tu/EF-1A subfamily', do not include 'X' in the sequence and are not a fragment were 984 entries. In the PDB, entries which are referenced from above 984 entries and are determined by X-ray crystallography were 68 entries. 14 entries were excluded because of binding an immunoprotein [25] and forming a chimeric protein [26-29]. Consequently, as shown in Table 1, 54 entries including 103 chains were retained.

Computations of f_1 and f_2

As $N=984$ and $K=103$ in Figure 2, the sequences were aligned by the

MAFFT 7 program [30]. 477 iM were extracted because of including residues which have coordinate data.

A difference between two sequences was computed by the maximum likelihood method [31] using the Jones-Taylor-Thornton model [32] as a substitution matrix and the Dayhoff method [33] for computing equilibrium frequencies. From all combinations of the differences, a phylogenetic tree was written by the unweighted pair group method with arithmetic mean [34]. $f_1({}^iM)$ was computed by changing, T, \mathcal{G}_x , h_x , τ , A and iG . For h_3 or h_4 , the Gonnet matrix [35] was used. For h_2 , a weight was computed by the Sibbald and Algos algorithm [36] and the iteration number was 100,000.

By separating each asymmetric unit, $f_1({}^iM)$ was computed and, in each entry, representative ions or molecules were shown in Table 1. However, because of uncertain functions, we excluded the following ions or molecules; sodium ion, acetate ion, sulfate ion, ammonium ion, sugar (sucrose), di(hydroxyethyl)ether, glyoxylic acid, 5-bromofuran-2-carboxylic acid, β -mercaptoethanol and water [37-43].

Correlations between f_1 and f_2

Let $[0, \infty) \supset F \ni f_1({}^iM)$ denote a subset of non-negative real numbers and a set of $f_1({}^iM)$ and be represented as $F \ni v_1 < v_2 < \dots < v_J$. Let t_j denote a threshold and satisfy

$$t_j = \begin{cases} < v_1 (j=0) \\ \frac{v_j + v_{j+1}}{2} (j=1, 2, \dots, J-1) \\ > v_J (j=J) \end{cases} \quad (13)$$

Let c_2 denote a cutoff of $f_2({}^iM)$ and, in this study, $c_2=3 \text{ \AA}$. Let I_j denote a number of iM which satisfies $f_2({}^iM) > c_2$ and I_i denote a number of iM which satisfies $f_2({}^iM) \leq c_2$. Let $I_{fp}(t_j)$ denote a number of iM which satisfies $f_2({}^iM) > c_2$ and $f_1({}^iM) \leq t_j$ and $I_{ip}(t_j)$ denote a number of iM which satisfies $f_2({}^iM) \leq c_2$ and $f_1({}^iM) \leq t_j$. Let a false positive rate

$$p(t_j) = \frac{I_{fp}(t_j)}{I_j} \quad (14)$$

a true positive rate

$$q(t_j) = \frac{I_{ip}(t_j)}{I_i} \quad (15)$$

and an area under the curve

$$AUC = \frac{1}{2} \sum_{j=0}^{J-1} [p(t_{j+1}) - p(t_j)] \cdot [q(t_{j+1}) + q(t_j)] \quad (16)$$

Let $F_x \ni f_x({}^iM)$ denote a multiset of $f_x({}^iM)$ and represented as $F_x \ni {}^1V_x \leq {}^2V_x \leq \dots \leq {}^lV_x$, where l is a number of iM . Let r denote a rank function and

$$r({}^{j+k-1}V_x) = j - 1 + \frac{t_n + 1}{2} \quad (17)$$

where $j=1, 2, \dots, l$, $m=1, 2, \dots, l$, $k=1, 2, \dots, t_n$ and t_n is a size of the tied rank. Here, a Spearman's ρ [44] is defined as

$$\rho = \frac{1}{2\sqrt{T_1 T_2}} \left\{ T_1 + T_2 - \sum_{i=1}^l [r({}^iV_1) - r({}^iV_2)]^2 \right\} \quad (18)$$

Where $l=1, 2, \dots, l$, $m=1, 2, \dots, l$,

Subfamily	Organism	PDB ID	Resolution	Ions or molecules
EF-Tu	<i>Bos taurus</i> , mitochondrial	1D2E	1.94	GDP, Mg ²⁺
		1XB2	2.20	Elongation factor Ts mitochondrial
	<i>Escherichia coli</i>	1EFC	2.05	GDP, Mg ²⁺
		2HCJ	2.12	GDP, TAC, Mg ²⁺
		3U6B	2.12	GDP, Mg ²⁺
		2BVN	2.30	ENX, GNP, Mg ²⁺
		4G5G	2.30	Thiomuracin A derivative, GDP, Mg ²⁺
		1D8T	2.35	Thiocillin GE2270, GDP, Mg ²⁺
		3U6K	2.45	Thiocillin GE2270 analogue NVP-LDK733, GDP, Mg ²⁺
		1DG1	2.50	GDP, Mg ²⁺
1EFU		2.50	Elongation factor Ts	
1EFM		2.70	GD P	
3U2Q		2.70	Thiocillin GE2270 analogue NVP-LFF571, GDP, Mg ²⁺	
2HDN		2.80	GDP, TAC, Mg ²⁺	
1ETU		2.90	GDP, Mg ²⁺	
4Q7J	2.90	Elongation factor Ts, Q β replicase		
1OB2	3.35	Phe-tRNA, GNP, KIR, Mg ²⁺		
2FX3	3.40	GDP, Mg ²⁺		
<i>Pseudomonas putida</i> KT2440	4J0Q	2.29	GDP, MES, MPD, Mg ²⁺	
	4IW3	2.70	Putative uncharacterized protein, GDP, Mg ²⁺	
<i>Thermus aquaticus</i>	1EFT	2.50	GNP, Mg ²⁺	
	1B23	2.60	Cys-tRNA, GNP, Mg ²⁺	
	1TTT	2.70	Phe-tRNA, GNP, Mg ²⁺	
	1TUI	2.70	GDP, Mg ²⁺	
	1OB5	3.10	Phe-tRNA, ENX, GNP, Mg ²⁺	
<i>Thermus thermophilus</i>	2C78	1.40	GNP, PUL, Mg ²⁺	
	2C77	1.60	Thiocillin GE2270, GNP, Mg ²⁺	
	1EXM	1.70	GNP, Mg ²⁺	
	4LBW	1.74	GNP, Mg ²⁺	
	4H9G	1.93	GNP, Mg ²⁺	
	1HA3	2.00	GNP, Mg ²⁺	
	4LBV	2.03	GDP, MAU, Mg ²⁺	
	4LBZ	2.22	GNP, Mg ²⁺	
	4LC0	2.22	GNP, Mg ²⁺	
	4LBV	2.69	GNP, Mg ²⁺ GNP, Mg ²⁺	
	1AIP	3.00	Elongation factor Ts	
	4V5L	3.10	16S rRNA, 23S rRNA, Trp-tRNA, GCP, Mg ²⁺	
	4V5P	3.10	16S rRNA, 23S rRNA, Trp-tRNA	
	4V5Q	3.10	16S rRNA, 30S rpS12, Trp-tRNA, GDP, KIR	
	4V5R	3.10	16S rRNA, Trp-tRNA, GDP, KIR 16S rRNA, Trp-tRNA, GDP, KIR	
	4V5S	3.10	16S rRNA, 23S rRNA, Small protein B SMPB, tmRNA δ, GDP, KIR, Mg ²⁺	
	4V8Q	3.10	16S rRNA, 23S rRNA, 30S rpS12, Thr-tRNA, GDP, KIR, Mg ²⁺	
4V5G	3.60			
aEF1A	<i>Aeropyrum pernix</i>	3VMF	2.30	Peptide chain release factor subunit 1, GTP, Mg ²⁺
		3WXM	2.30	Protein pelota homologue, GTP, Mg ²⁺
	<i>Sulfolobus solfataricus</i>	1JNY		GD P
		1SKQ		GDP, Mg ²⁺
eEF1A	<i>Oryctolagus cuniculus</i>	4C0S		GDP, Mg ²⁺
	<i>Saccharomyces cerevisiae</i>	1F60		Elongation factor 1Bα
		2B7C		Elongation factor-1 β
		1G7C		Elongation factor 1- β, 5GP
		1IJE		Elongation factor 1- β, GD P
		2B7B		Elongation factor-1 β, GD P
		1IJF		Elongation factor 1- β, GD P

TAC; Tetracycline, ENX; Enacyloxin IIa, GNP; Phosphoaminophosphonic acid-guanylate ester, KIR; Kirromycin, MES; 2-(N-morpholino)-ethanesulfonic acid, MPD; (4S)-2-methyl-2,4-pentanediol, PUL; Pulvomycin, MAU; N-methyl kirromycin, GCP; Phosphomethylphosphonic acid guanylate ester, 5GP; Guanosine-5' -monophosphate.

Table 1: 54 PDB entries of EF-Tu/EF-1A proteins.

$$T_1 = \frac{1}{12} \left[I^3 - I - \sum_{n=1}^{N_1} (t_n^3 - t_n) \right] \quad (19)$$

And

$$T_2 = \frac{1}{12} \left[I^3 - I - \sum_{n=1}^{N_2} (t_n^3 - t_n) \right] \quad (20)$$

where N_1 and N_2 are numbers of tied ranks in F_1 and F_2 , respectively.

Visualization

$f_1(iM)$, $f_2(iM)$, AUC and Spearman's ρ were visualized by the matplotlib Python package [45]. A three-dimensional structure was visualized by the VMD program [46].

Results

Fitness between f_1 and f_2

If g_x , h_x , τ and A are same but iG is different, Table 2 shows that when $iG = {}^iG$, the AUC or the Spearman's ρ is smaller than $iG = {}^G G$. In the latter case, Figure 3 shows that when the time point increases, the AUC or the Spearman's ρ tends to increase.

Evaluation of predicted functional amino acid residues by f_2

Figure 4A shows that $iM \in M$ is classifiable in 4 by $f_1(iM)$ and $f_2(iM)$ using a receiver operating characteristic (ROC) curve [47] in Figure 4B. Figures 4C and 4D show that the left sides tend to have small $f_1(iM)$ and small $f_2(iM)$ but the right sides tend to have large $f_1(iM)$ and large $f_2(iM)$.

Discussion

Meanings of $f_1(iM)$, $f_2(iM)$, AUC and Spearman's ρ are as follows. $f_1(iM)$ becomes small when characters are only diverged in near to the root of the phylogenetic tree. $f_1(iM)$ becomes large when characters are diverged in far from the root. $f_1(iM)$ becomes small when at least one amino acid residue in iM is proximate from an ion or a molecule. $f_1(iM)$ becomes large when amino acid residues in iM are not proximate from ions or molecules in all cocrystal structures. If the AUC is 0.5, a correlation between $f_1(iM)$ and being proximate and being non-proximate under a cutoff of $f_2(iM)$ may not exist. If the AUC is close to 1, small $f_1(iM)$ and large $f_1(iM)$ correlate with being proximate and being non-proximate, respectively. If the AUC is close to 0, large $f_1(iM)$ and small $f_1(iM)$ correlate with being proximate and being non-proximate, respectively. If the Spearman's ρ is 0, a linear correlation between $f_1(iM)$ and $f_2(iM)$ may not exist. If the Spearman's ρ is close to 1 or -1, $f_1(iM)$ and $f_2(iM)$ have a positive or a negative linear correlation, respectively.

If $T=1$, $g_x = g_3$, $h_x = h_2$, $A = {}^{20}A$ and $iG = {}^iG$, the method is the method based on SE [8]. If $T=1$, $g_x = g_3$, $h_x = h_2$, $A = {}^9A$ and $iG = {}^iG$, the method is the method based on SE of residue properties [10]. If $T=1$ is changed to $T=N$ in the former and the latter, Figure 3 shows that the AUC is from 0.5779 to 0.6147 and the Spearman's ρ is from 0.0757 to 0.1241 and the AUC is from 0.5709 to 0.5992 and the Spearman's ρ is from 0.1152 to 0.1405, respectively. Therefore, in the former and the latter, distinguishing characters utilizing the phylogenetic tree is effective for improving the AUC and the Spearman's ρ .

If $T=1$, $g_x = g_3$, $h_x = h_3$ and $iG = {}^iG$, the method is the method based on a sum of pairs [11]. If $T=1$, $g_x = g_3$, $h_x = h_4$ and $iG = {}^iG$, the

method is the method based on a sum of pairs with weighting [11]. If $T=1$ is changed to $T=N$ in the former and the latter, Figure 3 shows that the AUC is from 0.6083 to 0.6276 and the Spearman's ρ is from 0.1982 to 0.1653 and the AUC is from 0.6093 to 0.6211 and the Spearman's ρ is from 0.2263 to 0.1502, respectively. Therefore, in the former and the latter, distinguishing characters utilizing the phylogenetic tree is effective for improving the AUC but not for the Spearman's ρ . However, in the above case, if $iG = {}^iG$ is changed to $iG = {}^G G$ in the former and the latter, Figure 3 shows that the AUC is from 0.6941 to 0.7349 and the Spearman's ρ is from 0.4981 to 0.5650 and the AUC is from 0.6846 to 0.7335 and the Spearman's ρ is from 0.4749 to 0.5637, respectively. Therefore, in the former and the latter, distinguishing characters utilizing the phylogenetic tree and considering that each gap is different are effective for improving the AUC and the Spearman's ρ .

If $T=N$, $g_x = g_3$, $h_x = h_1$, $A = {}^{20}A$ and $iG = {}^iG$, the method is the iv-ET method [17]. If $T=N$, $g_x = g_3$, $h_x = h_2$, $A = {}^{20}A$ and $iG = {}^iG$, the method is equivalent to the rv-ET method [17]. If $iG = {}^iG$ is changed to $iG = {}^G G$ in the former and the latter, Table 2 shows that the AUC is from 0.5896 to 0.6242 and the Spearman's ρ is from 0.1221 to 0.3650 and the AUC is from 0.6180 to 0.7417 and the Spearman's ρ is from 0.1308 to 0.5722, respectively. Therefore, in the former and the

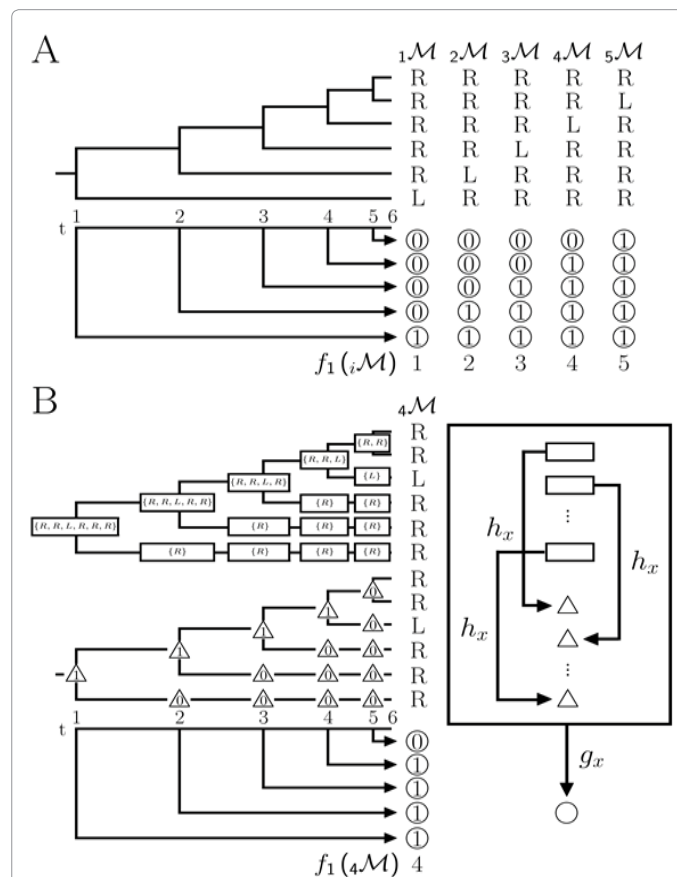


Figure 2: Computations of f_1 and f_2 (A) N sequences and K structures are extracted from the Swiss-Prot and the PDB, respectively. After all the sequences are aligned, $f_1(iM)$ and $f_2(iM)$ are computed by (B) and (C), respectively. (B) After a phylogenetic tree is written from sequences, $f_1(iM)$ is computed. (C) In structure k , kR and kQ denote coordinates of an amino acid residue and coordinates of ions or molecules, respectively. After proximity of kR and kQ is measured as $f_2(iM)$ and computed on K structures, $f_2(iM)$ is computed.

g_x	h_x	A	${}_i^1G$	r	AUC	Spearman's ρ
g_1	h_1	${}^{20}A$	${}_i^1G$	$0 < r < 1$ $0 < r < 1$	0.5896 0.6242	0.1221 0.3650
		9A	${}_i^1G$ ${}_i^G G$	$0 < r < 1$ $0 < r < 1$	0.5700 0.6184	0.1015 0.3509
	h_2	${}^{20}A$	${}_i^1G$	0.1	0.6036	0.1436
				0.2	0.6376	0.1977
				0.3	0.6420	0.2015
				0.4	0.6120	0.1585
		${}_i^G G$	0.1	0.7207	0.5566	
			0.2	0.7412	0.5773	
			0.3	0.7374	0.5613	
			0.4	0.7037	0.5028	
	9A	${}_i^1G$	0.2	0.5757	0.1171	
			0.3	0.6023	0.1659	
0.4			0.6151	0.1715		
0.1			0.5885	0.1160		
${}_i^G G$	0.1	0.6997	0.5185			
	0.2	0.6984	0.5097			
	0.3	0.7084	0.5063			
	0.4	0.6854	0.4797			
h_3	A	${}_i^1G$	0.1	0.5901	0.1288	
			0.2	0.6052	0.1467	
			0.3	0.6366	0.2315	
			0.4	0.6391	0.2763	
	${}_i^G G$	0.1	0.6758	0.4474		
		0.2	0.6975	0.4734		
		0.3	0.6905	0.4823		
		0.4	0.6865	0.4922		
h_4	A	${}_i^1G$	0.1	0.5872	0.1200	
			0.2	0.5939	0.1387	
			0.3	0.5916	0.1460	
			0.4	0.5995	0.1733	
	${}_i^G G$	0.1	0.6602	0.4162		
		0.2	0.6805	0.4470		
		0.3	0.6888	0.4599		
		0.4	0.6782	0.4675		
g_2	h_1	${}^{20}A$	${}_i^1G$	-	0.5916	0.0718
			${}_i^G G$	-	0.7399	0.5780
	h_2	9A	${}_i^1G$	-	0.5652	0.0587
			${}_i^G G$	-	0.7020	0.5145
	h_3	A	${}_i^1G$	-	0.6180	0.1308
			${}_i^G G$	-	0.7417	0.5722
	h_4	A	${}_i^1G$	-	0.5890	0.1257
			${}_i^G G$	-	0.7012	0.5091
	h_5	A	${}_i^1G$	-	0.6225	0.1579
			${}_i^G G$	-	0.7287	0.5517
	h_6	A	${}_i^1G$	-	0.6138	0.1412
			${}_i^G G$	-	0.7265	0.5501

g_3	h_1	^{20}A	$\begin{matrix} {}^1_iG \\ G \\ {}^G_iG \end{matrix}$	- -	0.5792 0.7386	0.0378 0.5801
		9A	$\begin{matrix} {}^1_iG \\ G \\ {}^G_iG \end{matrix}$	- -	0.5655 0.7052	0.0501 0.5244
	h_2	^{20}A	$\begin{matrix} {}^1_iG \\ G \\ {}^G_iG \end{matrix}$	- -	0.6147 0.7393	0.1241 0.5723
		9A	$\begin{matrix} {}^1_iG \\ G \\ {}^G_iG \end{matrix}$	- -	0.5992 0.7059	0.1405 0.5170
	h_3	A	$\begin{matrix} {}^1_iG \\ G \\ {}^G_iG \end{matrix}$	- -	0.6276 0.7349	0.1653 0.5650
	h_4	A	$\begin{matrix} {}^1_iG \\ G \\ {}^G_iG \end{matrix}$	- -	0.6211 0.7335	0.1502 0.5637

Table 2: Correlations between f_1 and f_2 .

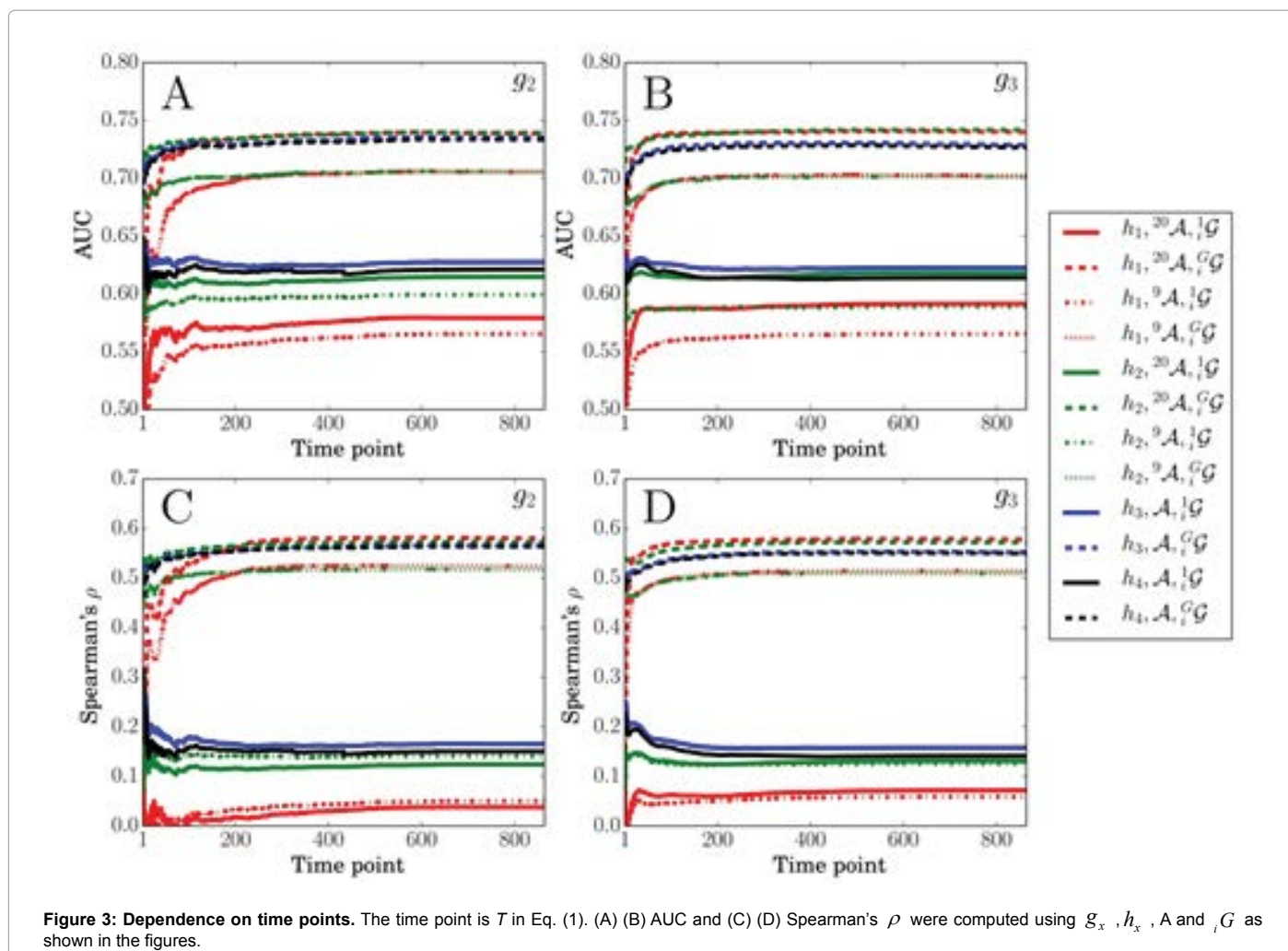


Figure 3: Dependence on time points. The time point is T in Eq. (1). (A) (B) AUC and (C) (D) Spearman's ρ were computed using g_x , h_x , A and ${}_iG$ as shown in the figures.

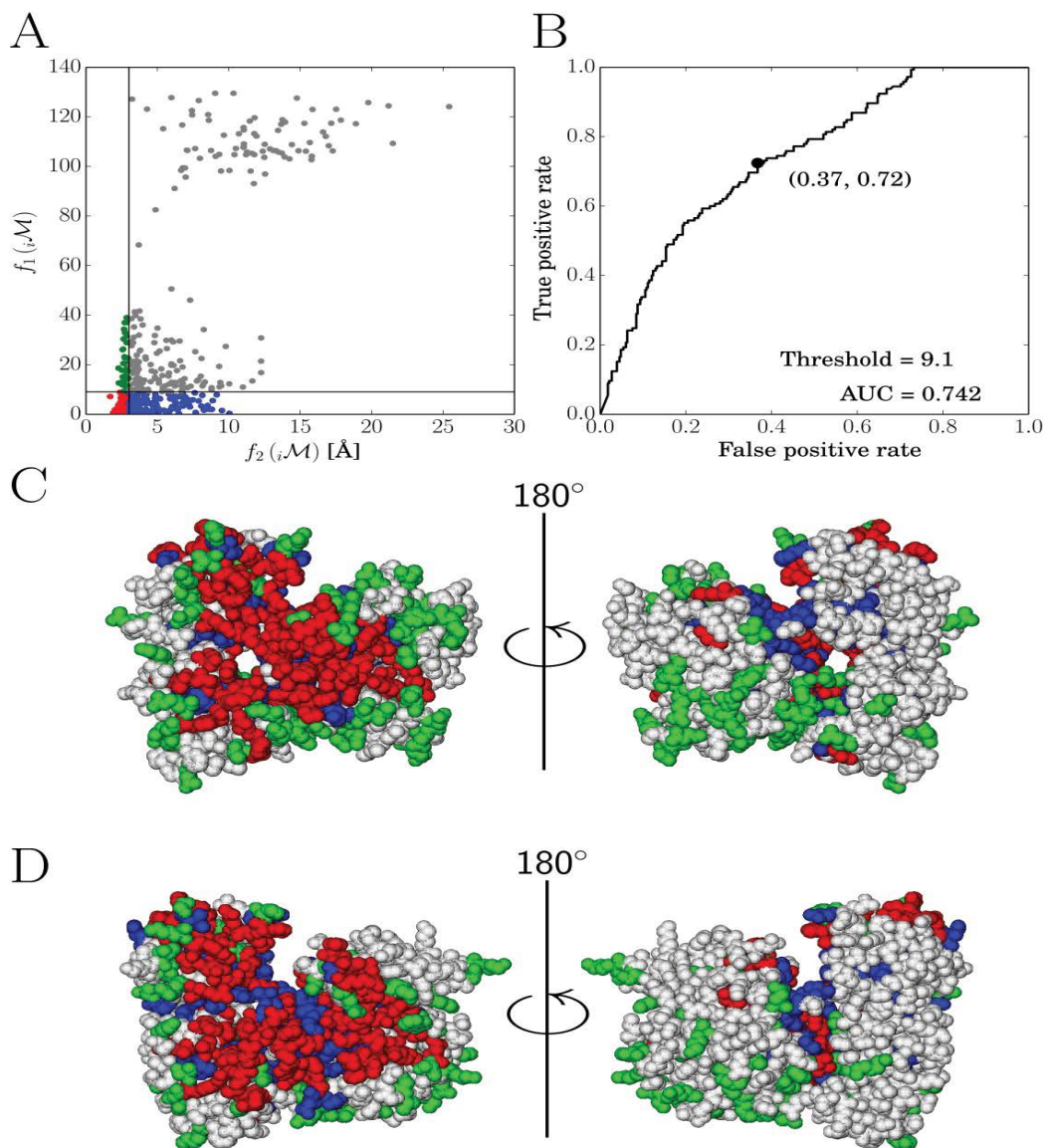


Figure 4: A scatter plot, an ROC curve and three-dimensional structures. $f_1(iM)$ was computed using g_2 , h_2 , ^{20}A and $^G G$ (A) A scatter plot of $f_1(iM)$ and $f_2(iM)$ However, one point whose $f_1(468M)$ 122.2 and $f_2(468M)$ 45.56 was not shown. $iM \in M$ was classified into 2 by whether $f_2(iM)$ is equal to or smaller than 3 Å or larger than 3 Å. (B) By regarding the former as true and the latter as false, the ROC curve was written using $f_1(iM)$ The threshold was determined so that (true positive rate + 1 - false positive rate) is maximum and, eventually, $M \in M$ was classified into 4, which were visualized onto three-dimensional structures of (C) *Thermus thermophilus* EF-Tu [42] and (D) *Saccharomyces cerevisiae* EF1A [48].

latter, considering that each gap is different is effective for improving the AUC and the Spearman's ρ Thus, $f_1(iM)$ is evaluable by $f_2(iM)$ and our methods improved some existing methods.

EF-Tu/EF-1A proteins are responsible for protein biosynthesis [42,48] and we selected cocrystal structures involving the function. Therefore, if $f_2(iM)$ is small, an amino acid residue in iM is proximate from a region involving protein biosynthesis. If $f_2(iM)$ is large, the amino acid residues in iM are not proximate from the region. Figures 4A, 4C and 4D show the proximate region and the non-proximate region and Figure 4B shows that, on the ROC curve of $f_1(iM)$, the AUC is 0.742, which indicates that the proximate region tends to become small $f_1(iM)$ but the non-proximate region tends to

become large $f_1(iM)$ In addition, Table 2 shows that the Spearman's ρ is 0.5722, which indicates that $f_1(iM)$ tends to be small if $f_2(iM)$ is small and $f_1(iM)$ tends to be large if $f_2(iM)$ is large. However, a complete linear correlation between

and $f_2(iM)$ was not obtainable and therefore not all of $f_1(iM)$ can explain $f_2(iM)$. This may indicate that $f_1(iM)$ and $f_2(iM)$ can measure a similar thing each other but cannot always measure a same thing and, by $f_1(iM)$ and $f_2(iM)$, measurable things such as importance for binding ions or molecules or importance for maintaining the structure may be different. Thus, from a different point of view, $f_1(iM)$ and $f_2(iM)$ can evaluate an amino acid residue.

Conclusions

Methods to map an MSA, which is represented as a character type and a coordinate type, were described and we propose two usages. The first one is to assess fitness between the first map and the second map. The second one is to evaluate predicted functional amino acid residues by use of the second map. Our methods show a better performance and reliability for functional site prediction of EF-Tu/EF-1A proteins.

References

- Huang YF, Golding GB (2014) Phylogenetic Gaussian process model for the inference of functionally important regions in protein tertiary structures. *PLoS Comput Biol* 10.
- Fang C, Noguchi T, Yamana H (2014) Simplified sequence-based method for ATP-binding prediction using contextual local evolutionary conservation. *Algorithms Mol Biol* 9.
- Gultas M, Duzgun G, Herzog S, Juger SJ, Meckbach C, et al. (2014) Quantum coupled mutation finder: predicting functionally or structurally important sites in proteins using quantum Jensen-Shannon divergence and CUDA programming. *BMC Bioinformatics* 15:96.
- Janda JO, Popal A, Bauer J, Busch M, Klocke M, et al. (2014) H2rs: deducing evolutionary and functionally important residue positions by means of an entropy and similarity based analysis of multiple sequence alignments. *BMC Bioinformatics* 15.
- Lee TW, Yang ASP, Brittain T, Birch NP (2015) An analysis approach to identify specific functional sites in orthologous proteins using sequence and structural information: application to neuroserpin reveals regions that differentially regulate inhibitory activity. *Proteins* 83: 135-152.
- Huang YF, Golding GB (2015) FuncPatch: a web server for the fast Bayesian inference of conserved functional patches in protein 3D structures. *Bioinformatics* 31: 523-531.
- Valdar W (2002) Scoring residue conservation. *Proteins* 48: 227-241
- Sander C, Schneider R (1991) Database of homology-derived protein structures and the structural meaning of sequence alignment. *Proteins* 9: 56-68.
- Shenkin P, Erman B, Mastrandrea L (1991) Information-theoretical entropy as a measure of sequence variability. *Proteins* 11: 297-313.
- Williamson RM (1995) Information theory analysis of the relationship between primary sequence structure and ligand recognition among a class of facilitated transporters. *J Theor Biol* 174: 179-188.
- Valdar W, Thornton J (2001) Protein-protein interfaces: Analysis of amino acid conservation in homodimers. *Proteins* 42: 108-124.
- Wang K, Samudrala R (2006) Incorporating background frequency improves entropy-based residue conservation measures. *BMC Bioinformatics* 7.
- Capra JA, Singh M (2007) Predicting functionally important residues from sequence conservation. *Bioinformatics* 23: 1875-1882.
- Henikof S, Henikof J (1994) Position-based sequence weights. *J Mol Biol* 243: 574-578.
- Lichtarge O, Bourne H, Cohen F (1996) An evolutionary trace method defines binding surfaces common to protein families. *J Mol Biol* 257: 342-358.
- Landgraf R, Fischer D, Eisenberg D (1999) Analysis of heregulin symmetry by weighted evolutionary tracing. *Protein Eng* 12: 943-951.
- Mihalek I, Res I, Lichtarge O (2004) A family of evolution-entropy hybrid methods for ranking protein residues by importance. *J Mol Biol* 336: 1265-1282.
- Armon A, Graur D, Ben-Tal N (2001) ConSurf: An algorithmic tool for the identification of functional regions in proteins by surface mapping of phylogenetic information. *J Mol Biol* 307: 447-463.
- Pupko T, Bell RE, Mayrose I, Glaser F, Ben-Tal N (2002) Rate4site: an algorithmic tool for the identification of functional regions in proteins by surface mapping of evolutionary determinants within their homologues. *Bioinformatics* 18: S71-S77.
- Mayrose I, Graur D, Ben-Tal N, Pupko T (2004) Comparison of site-specific rate-inference methods for protein sequences: Empirical Bayesian methods are superior. *Mol Biol Evol* 21: 1781-1791.
- Johansson F, Toh H (2010) A comparative study of conservation and variation scores. *BMC Bioinformatics* 11.
- Ashkenazy H, Erez E, Martz E, Pupko T, Ben-Tal N (2010) ConSurf 2010: calculating evolutionary conservation in sequence and structure of proteins and nucleic acids. *Nucleic Acids Res* 38: W529-W533.
- Berman H, Westbrook J, Feng Z, Gilliland G, Bhat T, et al. (2000) The protein data bank. *Nucleic Acids Res* 28: 235-242.
- Consortium TU (2015) Uniprot: a hub for protein information. *Nucleic Acids Res* 43: D204-D212, 2015.
- Dai S, Crawford F, Marrack P, Kappler JW (2008) The structure of HLA-DR52c: comparison to other HLA-DRB3 alleles. *Proc Natl Acad Sci U S A* 105: 11893-11897.
- Kidmose RT, Vasiliev NN, Chetverin AB, Andersen GR, Knudsen CR (2010) Structure of the Q β replicase, an RNA-dependent RNA polymerase consisting of viral and host proteins. *Proc Natl Acad Sci U S A* 107: 10884-10889.
- Takehita D, Tomita K (2010) Assembly of Q β viral RNA polymerase with host translational elongation factors EF-Tu and -Ts. *Proc Natl Acad Sci U S A* 107: 15733-15738.
- Takehita D, Tomita K (2012) Molecular basis for RNA polymerization by Q β replicase. *Nat Struct Mol Biol* 19: 229-237.
- Takehita D, Yamashita S, Tomita K (2012) Mechanism for template-independent terminal adenylation activity of Q β replicase. *Structure* 20: 1661-1669.
- Katoh K, Standley DM (2013) MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol* 30: 772-780.
- Kishino H, Miyata T, Hasegawa M (1990) Maximum-likelihood inference of protein phylogeny and the origin of chloroplasts. *J Mol Evol* 31: 151-160.
- Jones D, Taylor W, Thornton J (1992) The rapid generation of mutation data matrices from protein sequences. *Comput Appl Biosci* 8: 275-282.
- Dayhof MO, Schwartz RM (1978) Chapter 22: A model of evolutionary change in proteins. In: *Atlas of Protein Sequence and Structure*.
- Sneath PHA, Sokal RR (1973) Numerical taxonomy: the principles and practice of numerical classification. W. H. Freeman and company, San Francisco.
- Benner S, Cohen M, Gonnet G (1994) Amino-acid substitution during functionally constrained divergent evolution of protein sequences. *Protein Eng* 7: 1323-1332.
- Sibald P, Argos P (1990) Weighting aligned protein or nucleic acid sequences to correct for unequal representation. *J Mol Biol* 216: 813-818.
- Hefron SE, Mui S, Aorora A, Abel K, Bergmann E, et al. (2006) Molecular complementarity between tetracycline and the GTPase active site of elongation factor Tu. *Acta Crystallogr D Biol Crystallogr* 62: 1392-1400.
- Nissen P, Thirup S, Kjeldgaard M, Nyborg J (1999) The crystal structure of Cys-tRNACys-EF-Tu-GDPNP reveals general and specific features in the ternary complex and in tRNA. *Structure* 7: 143-156.
- LaMarche MJ, Leeds JA, Amaral K, Brewer JT, Bushell SM, et al. (2011) Antibacterial optimization of 4-aminothiazolyl analogues of the natural product GE2270 A: identification of the cycloalkylcarboxylic acids. *J Med Chem* 54: 8099-8109.
- Kobayashi K, Saito K, Ishitani R, Ito K, Nureki O (2012) Structural basis for translation termination by archaeal RF1 and GTP-bound EF1 α complex. *Nucleic Acids Res* 40: 9319-9328.
- Groftehaug MK, Therkelsen MO, Taaning R, Skrydstrup T, Morth JP, et al. (2013) Identifying ligand-binding hot spots in proteins using brominated fragments. *Acta Crystallogr F Struct Biol Commun* 69: 1060-1065.
- Parmeggiani A, Krab I, Watanabe T, Nielsen R, Dahlberg C, et al. (2006) Enacyloxin IIa pinpoints a binding pocket of elongation factor Tu for development of novel antibiotics. *J Biol Chem* 281: 2893-2900.
- Vogele L, Palm G, Mesters J, Hilgenfeld R (2001) Conformational change of elongation factor Tu (EF-Tu) induced by antibiotic binding - crystal structure of the complex between EF-Tu•GDP and aureodox. *J Biol Chem* 276: 17149-17155.
- Spearmen C (2010) The proof and measurement of association between two things. *Int J Epidemiol* 39: 1159-1161.

45. Hunter JD (2007) Matplotlib: a 2D graphics environment. *Computing Sci Eng* 9: 90-95.
46. Humphrey W, Dalke A, Schulten K (1996) VMD: visual molecular dynamics. *J Mol Graph Model* 14: 33-38.
47. Lasko T, Bhagwat J, Zou K, Ohno-Machado L (2005) The use of receiver operating characteristic curves in biomedical informatics. *J Biomed Inform* 38: 404-415.
48. Andersen G, Pedersen L, Valente L, Chatterjee I, Kinzy T, et al. (2000) Structural basis for nucleotide exchange and competition with tRNA in the yeast elongation factor complex eEF1A : eEF1B α . *Mol Cell* 6: 1261-1266.