

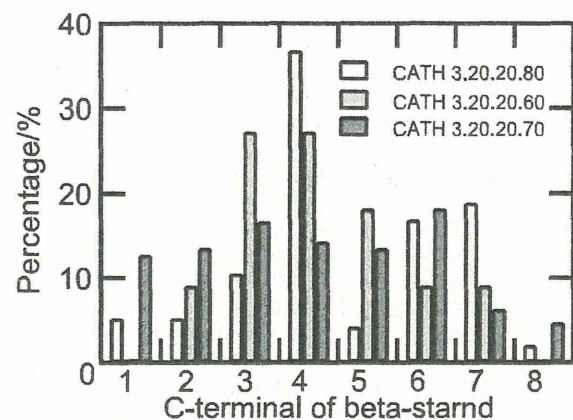
but almost unchanged between low and high functional diversity, suggesting that LBRs can discriminate functions in superfamilies with all ranges of functional diversity. The same tendency was observed with functional diversity defined by numbers of the fourth-digit EC number level functions (Figure S3 and Table S10). The similar tendencies between the two classification schemes, observed in prediction performance and the proportions of ASRs and LBRs, may be accounted for by the observation that superfamilies with high functional diversity at the third-digit level generally have many distinct fourth digits in each third-digit EC number function.

### Examples of superfamilies and enzymes

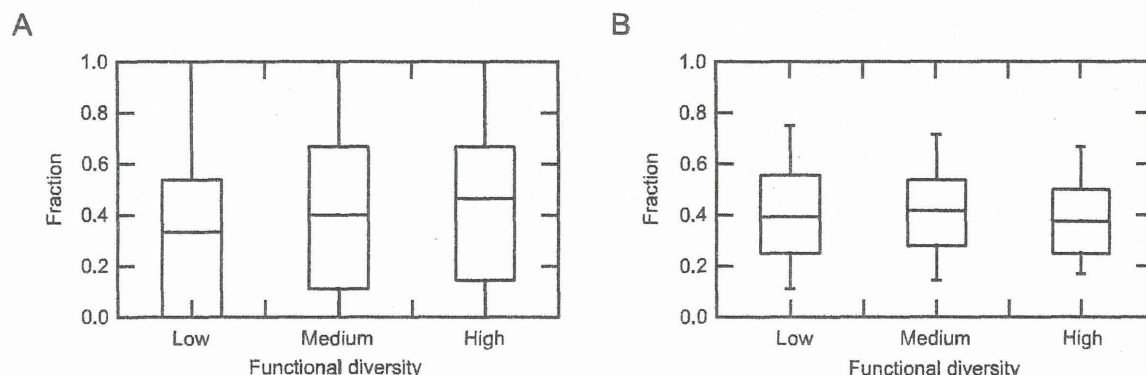
In this section, we describe a detailed investigation of the properties of the rf-SDRs in selected enzymes from superfamilies with different degrees of functional diversity. To remove potential biases associated with protein folds, we first show three superfamilies from a single fold, and next we show an additional example from a different fold. Only three folds, TIM barrel (CATH 3.20.20),  $\alpha$ - $\beta$ -plaits (CATH 3.30.70) and Rossmann fold (CATH 3.40.50), satisfied the condition of having superfamilies in each of all three classes of functional diversity and in each class, containing at least one enzyme, for which the ASR information was available. From these three, we selected the TIM barrel fold (CATH 3.20.20). The TIM barrel, ( $\alpha/\beta$ ) $_8$ -barrel fold, is one of the largest and oldest fold and in the enzymes belonging to this fold, all the active sites are located at the C-terminal ends of the  $\beta$ -strands. As typical examples of superfamilies with low and high functional diversity, we chose glycosidases (CATH 3.20.20.80) and aldolase class I (CATH 3.20.20.70), respectively. We then chose phosphoenolpyruvate-binding domains (CATH 3.20.20.60) as an example of the superfamilies with medium functional diversity, although the number of enzymes with available ASR information was limited and the proportion of ASRs to be selected as rf-SDRs was somewhat atypical. Therefore, we additionally examined the  $\alpha/\beta$ -hydrolase superfamily (CATH 3.40.50.1820) as a second example of the superfamilies with medium diversity, because this superfamily highlighted deviations from the average properties of this class of superfamilies explained by the well conserved catalytic triad.

**Glycosidase superfamily (CATH 3.20.20.80).** The glycosidase superfamily, where most enzymes belong to glycosidases (EC. 3.2.1), is a superfamily with low functional diversity. In our

dataset, this superfamily contained 16 different glycosidases (EC 3.2.1) and three different hexosyltransferases (EC 2.4.1) (Table S3). The white bars in Figure 6 shows the distribution of the positions of the active site residues at eight C-terminal ends of the  $\beta$ -strands in this superfamily, highlighting three main catalytic residues at the  $\beta$ -strands 4, 7 and 6. This observation is consistent with the fact that 12 of the 16 glycosidases in this superfamily have been characterized as members of a group known as “the 4/7 group” [47–49]. (In the literature, this group is normally referred to as “the 4/7 superfamily” but to avoid confusion, we use the term group here.) The enzymes in the 4/7 group utilize two conserved catalytic acidic residues located at the C-terminal ends of  $\beta$ -strands 4 (acid/base) and 7 (nucleophile), as well as residues at the end of  $\beta$ -strand 6, which modulate the nucleophile. This biased distribution is reflected in the proportion of ASRs to be selected as rf-SDRs (32.7%), which was lower than the average for the



**Figure 6. The distribution of active site residues at the end of eight  $\beta$ -strands of enzymes in the superfamily adopting the TIM barrel fold.** White bars represent the glycosidase superfamily (CATH 3.20.20.80), light gray bars represent the phosphoenolpyruvate-binding domain superfamily (CATH 3.20.20.60), and gray bars represent the aldolase class I superfamily (CATH 3.20.20.70). The percentages were calculated by using 18, three and 29 enzymes for glycosidases, phosphoenolpyruvate-binding domains and aldolase class I, respectively, for which active site information was available. doi:10.1371/journal.pone.0084623.g006



**Figure 5. Distributions of fractions of the rf-SDRs in active site residues (ASRs, A) and ligand binding residues (LBRs, B), observed in the superfamilies with low, medium and high degrees of functional diversity classified at the third-digit level of EC numbers.** The top and bottom of a box indicate 75th and 25th percentiles and the horizontal line in a box represents the median value. The top and bottom whiskers represent 90th and 10th percentiles. doi:10.1371/journal.pone.0084623.g005

group of superfamilies with low functional diversity (35.0%), (Tables S9 and S11).

Figure 7 shows two example enzymes of the 4/7 group, endo-1,4- $\beta$ -xylanase (EC 3.2.1.8, Figure 7A) and cellulase (EC 3.2.1.4, Figure 7B). In both enzymes, none of the two 4/7 catalytic residues (Glu 159, Glu 265 in Figure 7A and Glu 170, Glu 307 in Figure 7B, respectively) was selected as the rf-SDRs. The rf-SDRs included some residues on  $\beta$ -strand 6, His 236 in endo-1,4- $\beta$ -xylanase and His 254 and Tyr 256 in cellulase, which contact the nucleophiles and are invariant in each enzyme but different between the two enzymes [50–52]. The proportion of ASRs to be selected as rf-SDRs in endo-1,4- $\beta$ -xylanase is lower (0.25) than that in cellulase (0.5), possibly because the former enzyme share the active site residues (other than the 4/7 catalytic residues) with a larger number of other enzymes such as glucan 1,4- $\alpha$ -maltohydrolase (EC 3.2.1.133) and cyclomalto-dextrin glucanotransferase (EC 2.4.1.19) than the latter enzyme.

The rf-SDRs also included some LBRs, which are located in similar spatial positions but not equivalent in the sequence alignment, His 95 (endo-1,4- $\beta$ -xylanase) and His 122 (cellulase) [50] shown to be essential for ligand binding by mutagenesis experiments [53–55], and the residues critical for determining the substrate positions, Trp 241 at the +3 subsite [56], Asn 59 and Lys 62 at the -2 subsite [57], in endo-1,4- $\beta$ -xylanase.

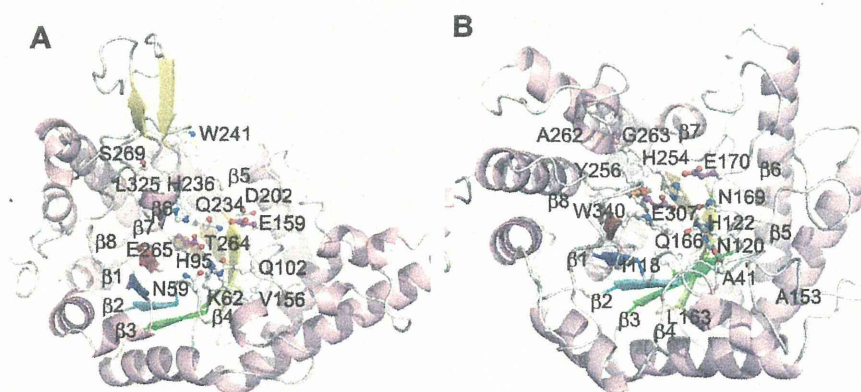
**Aldolase class I superfamily (CATH 3.20.20.70).** The Aldolase class I superfamily is known to be an old family including a variety of enzymes. In our dataset, predictors for 34 different enzymes were constructed in this superfamily (Table S3). These 34 enzymes included EC numbers with six different first-digits, showing the highest functional entropy in all the superfamilies. The ASR positions showed a broad distribution, indicating that the numerous functions are achieved by the active sites located at various ends of  $\beta$ -strands (Figure 6, dark gray bars). For instance, in 5-aminolevulinic acid dehydratase (ALADH, EC 4.2.1.24) [58], the catalytic Lys 195 and Lys 247 are positioned at the ends of  $\beta$ -7 and  $\beta$ -8, respectively and in phosphoribosylformimino-5-aminoimidazole carboxamide ribonucleotide (ProFAR) isomerase (HisA, EC 5.3.1.16) [59], the catalytic Asp 8 is positioned at the C-terminal end of  $\beta$ -1. Aldolase class I enzymes typically have substrates or cofactors with a phosphate-group, such as flavin mononucleotide (FMN), but enzymes in this superfamily also act

on a variety of other substrates. The proportion of ASRs to be selected as rf-SDRs (51.9%) was higher than the average for the group of superfamilies with high functional diversity (43.7%) (Tables S9 and S11). This observation suggests that the ASRs located differently among the enzymes can be used effectively for discriminating different functions in this superfamily.

Figures 8A and 8B show the rf-SDRs of quinolinate phosphoribosyltransferase (hQPRTase; EC 2.4.2.19) and  $\alpha$ -galactosidase ( $\alpha$ -Gal; EC 3.2.1.22) as examples of enzymes having dissimilar functions. The rf-SDRs of hQPRTase included one core residue of the phosphate binding motif [60] Ala 268 at the end of  $\beta$ -10, which corresponds to  $\beta$ -8 in a conventional ( $\alpha/\beta$ )<sub>8</sub> barrel (in Figure 8A, the numbering of the  $\beta$ -strands based on the conventional barrel), and one of the catalytic residues, Lys 140 on  $\beta$ -1. Leu 170 and Lys 172 on  $\beta$ -4, the conformational change of which was suggested to be important for the specificity and reaction mechanism [61], were also included (Figure 8A). On the other hand,  $\alpha$ -Gal recognizes the substrate having no phosphate moiety, mainly around the C-terminal ends of  $\beta$ -3 to  $\beta$ -6 [62]. In addition to the nucleophile Asp 130 at the end of  $\beta$ -4, many LBRs on these  $\beta$ -strands were selected as rf-SDRs (Figure 8B).

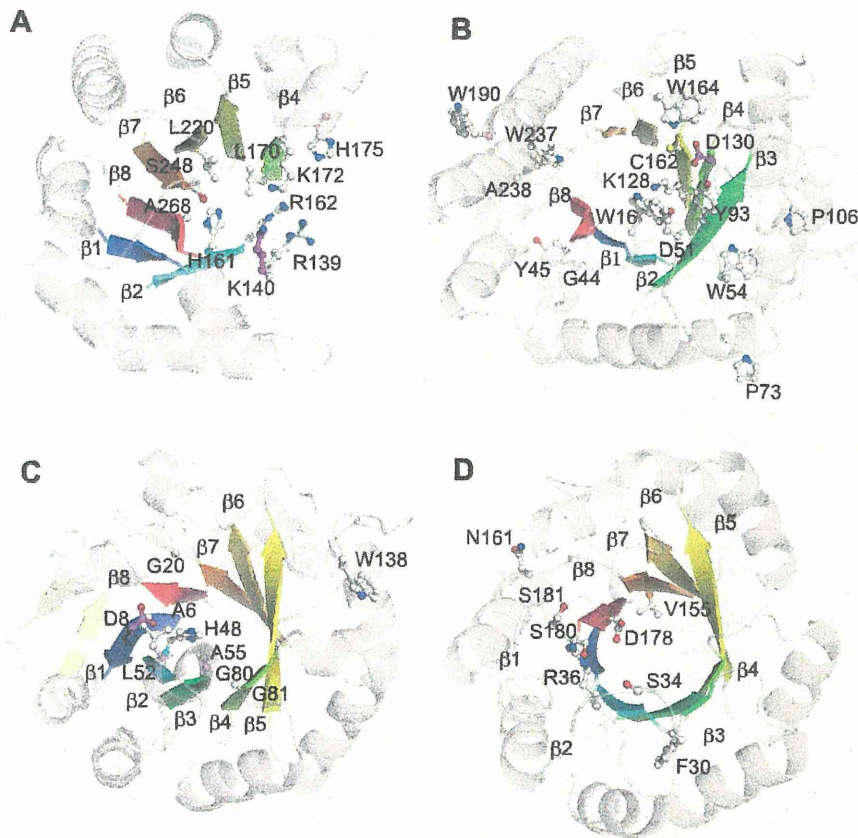
Figures 8C and 8D show ProFAR isomerase (HisA) (EC 5.3.1.16) and phosphoribosylanthranilate (PRA) isomerase (TrpF) (EC 5.3.1.24) as examples of enzymes having more similar functions. These enzymes catalyze the Amadori rearrangements of different substrates ProFAR and PRA by similar mechanisms [63,64]. These substrates share a ribose-5-phosphate moiety, and ProFAR has an additional ribose connected by imidazole and PRA has an anthranilate moiety. Also known are PriA, which can catalyze both reactions, and its close homologue subHisA, which lacks the TrpF activity [65].

In the rf-SDRs of HisA, the only known catalytic residue (Asp 8) was selected. In TrpF, the corresponding active site, Cys 7, was not selected and the reason is unclear. In LBRs, some residues interacting with different moieties of each substrate were selected to be rf-SDRs: Ser 34 and Arg 36 of TrpF, which interact with the anthranilate moiety of the substrate [66], Gly 20 and Leu 52 of HisA, which would interact with the imidazole and attached amide moieties (inferred from the homologous PriA structure). Additionally, the rf-SDRs included His 48 and Trp 138 of HisA, likely to be important for the catalytic activity for PRA (also



**Figure 7. The rf-SDRs for (A) endo-1,4-xylanase (EC 3.2.1.8, CATH domain: 1r87A00) and (B) cellulase (EC 3.2.1.4, CATH domain: 1edgA00) in the glycosidase superfamily (CATH 3.20.20.80).** The rf-SDRs are represented by balls and sticks, where nitrogen atoms are colored blue, oxygen atoms are red, sulfur atoms are yellow and carbon atoms are white. The carbon atoms of the active sites selected as rf-SDRs are colored magenta. Eight  $\beta$ -strands in a conventional barrel are colored blue, cyan, green, lemon, yellow, yelloworange, orange, and red, from the N-terminal to the C-terminal. In both enzymes, none of the two catalytic acid residues common in many enzymes in the superfamily, colored magenta, was selected.

doi:10.1371/journal.pone.0084623.g007



**Figure 8. The rf-SDRs for (A) quinolinate phosphoribosyltransferase (hQPRTase; EC 2.4.2.19, CATH domain: 1qprF02), (B)  $\alpha$ -galactosidase ( $\alpha$ -Gal; EC 3.2.1.22, CATH domain: 1uasA01), (C) phosphoribosylformimino-5-aminoimidazole carboxamide ribonucleotide isomerase (HisA) (EC 5.3.1.16, CATH domain: 1qo2A00) and (D) phosphoribosylanthranilate isomerase (TrpF) (EC 5.3.1.24, CATH domain: 1nsjA00) in aldolase class I superfamily (CATH 3.20.20.70). The rf-SDRs are represented by balls and sticks, where nitrogen atoms are colored blue, oxygen atoms are red, sulfur atoms are yellow and carbon atoms are white. The carbon atoms of the active sites selected as rf-SDRs are colored magenta. Eight  $\beta$ -strands in a conventional barrel are colored blue, cyan, green, lemon, yellow, yelloworange, orange, and red, from the N-terminal to the C-terminal. The rf-SDRs in the figures A and B clearly show that the rf-SDRs for hQPRTase include the phosphate binding motif located in  $\beta$ -7 and  $\beta$ -8 in the conventional barrel structure but those for  $\alpha$ -Gal are mainly located after  $\beta$ -1 to  $\beta$ -5. The figure D shows the residues interacting with different moieties in substrates between HisA and TrpF, Ser 34 and Arg 36.**  
doi:10.1371/journal.pone.0084623.g008

inferred from the PriA structure) [67]. In addition to these residues, different residues in different enzymes were selected, from those interacting with common parts of the substrates such as the phosphate moiety.

**Phosphoenolpyruvate-binding domain superfamily (CATH 3.20.20.60).** The phosphoenolpyruvate-binding domain superfamily mainly consists of transferases (EC 2) and lyases (EC 4). Most of these enzymes have substrates or cofactors with a phosphate-moiety, while the phosphate binding sites are distributed over the C-terminal ends of  $\beta$ -strands 2 to 6. The predictors for six different enzymes consisting of two phosphotransferases with paired acceptors (EC 2.7.9), two oxo-acid-lyases (EC 4.1.3) and other transferases (EC 2) were constructed (Table S3). This superfamily was classified into the group of medium functional diversity.

Despite generally dissimilar active sites among these enzymes (Figure 6, light gray bars), the proportion of ASRs to be selected as rf-SDRs (23.5%) was lower than the average for the group of superfamilies with medium functional diversity (43.4%) (Tables S9 and S11). This result may be explained by the conservation of some of the active site residues. For example, pyruvate phosphate

dikinase (EC 2.7.9.1) has the only known active site, Cys 831 [68] and this position in the alignment was also occupied by cysteine in pyruvate water dikinase (EC 2.7.9.2) (although no active site information is available for the latter enzyme). This position was not selected to be an rf-SDR, decreasing the average proportion of ASRs to be selected.

**$\alpha/\beta$ -hydrolase superfamily (CATH 3.40.50.1820).**  $\alpha/\beta$ -hydrolase superfamily is one of the large superfamilies, containing a wide variety of enzymes such as carboxylic acid ester hydrolases, peptidases, lipid hydrolases and haloalkane dehalogenases. In our dataset, predictors for 13 enzymes were constructed (Table S3). All these enzymes shared the first digit of the EC number (EC3; hydrolases) and this superfamily belonged to the group of superfamilies with medium functional diversity. A variety of functions are achieved by the conserved catalytic triad: a nucleophile (Ser, Cys or Asp) positioned after  $\beta$ -5, an acidic residue after  $\beta$ -7 and histidine after the last  $\beta$ -8 strand, and the versatile substrate binding sites by insertions and deletions at the C-terminal ends of  $\beta$ -3, 4, 6, 7 or 8 [69,70]. Such a conserved catalytic triad and a similar chemical reaction mechanism are reflected in the proportion of ASRs to be selected as rf-SDRs

(26.2%), which was lower than the average value (43.4%) for the group of medium functional diversity (Tables S9 and S11).

For instance, acetylcholine esterase (AChE, EC 3.1.1.7) shown in Figure 9 has the conventional catalytic triad, Ser, Glu, and His, and a deep and narrow cavity around the catalytic site called “active site gorge” formed by large insertions, which is considered to determine the specificity for acetylcholine [71]. In 15 rf-SDRs, no residue of the catalytic triad was selected and about 40% of the rf-SDRs were located in the active site gorge. Trp 84 and Phe 330 are known as the anionic site to bind the choline moiety and Tyr 121, Trp 279 and Phe 290 are important for determining the gorge conformation [72–75]. Phe 290 causes steric hindrance with a large acyl group in the acyl pocket and plays a critical role in stabilizing the methyl moiety of acetylcholine [76].

These examples show whether each residue can be selected as an rf-SDR or not depends on whether it is conserved within a superfamily regardless of what roles the equivalent residues play in other enzymes. A residue may be conserved and used as a catalytic residue for the same chemical reaction in other enzymes and thus, it tends not to be selected as an rf-SDR, as observed in the glycosidase superfamily. A conserved residue may be used for catalyzing different chemical reaction but because of its conservation, it cannot be selected to be an rf-SDR, as observed in the  $\alpha/\beta$ -hydrolase superfamily. In some superfamilies, different amino acid residues are used for catalyzing different chemical reactions or binding different ligands, in which case, these functional residues can be selected for rf-SDRs, as observed in the aldolase class I superfamily.

## Conclusion

We have developed EFPrf, a novel method based on random forests for predicting enzyme functions at the fourth-digit level of

the EC number in each CATH homologous superfamily. As input attributes, we used amino acid residue similarities at ASRs, LBRs and CSRs, in addition to similarity in the full-length sequence. The prediction performance of EFPrf improved significantly over the decision trees constructed using BLAST scores alone (the simple model), especially in the low MTTSI regions, where it is known to be difficult to distinguish detailed functions by sequence similarity alone. This observation suggested that the information about functionally important sites would be useful for predicting detailed functions. During the construction of EFPrf, we also obtained the rf-SDRs from the most highly contributing attributes. The analysis of the selected superfamilies showed that the rf-SDRs included many experimentally verified SDRs. Moreover, we showed that the rf-SDRs reflected the mechanisms of functional diversification within each superfamily; the rf-SDRs both indicate a general degree of functional diversity (as measured by the proportion of ASRs to be selected as rf-SDRs) and the specific characteristics of each superfamily represented by the conservations of each residue in a superfamily. Thus, EFPrf is a useful tool for predicting detailed enzyme functions and the rf-SDRs are a good resource for determining SDRs by experimental and computational methods and understanding functional diversity in a superfamily.

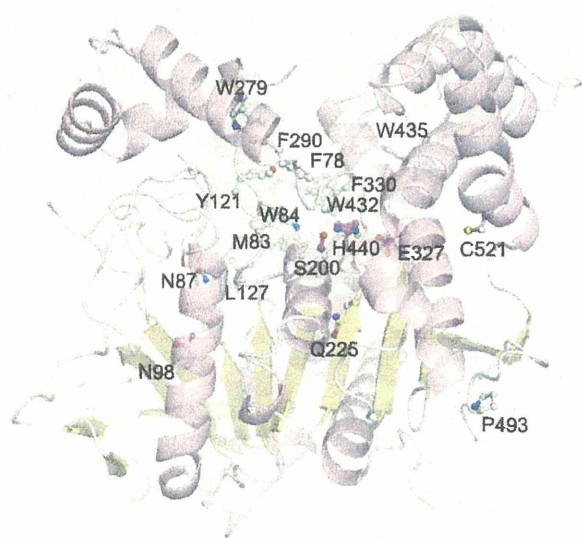
In this paper, we examined individual domain sequences pre-assigned to a CATH superfamily for validating EFPrf. In practice, enzyme sequences often consist of multiple domains and in the future, we will develop a method for combining prediction results for the individual domains of a query sequence and producing an overall function prediction. In recent years, many methods have been proposed for predicting protein functions described by GO terms [13]. Our method can be extended to GO term prediction and may be efficient in the low sequence similarity region, where GO terms are also difficult to predict [24,77].

## Materials and Methods

### Dataset preparation

Figure 2 shows an outline of the dataset construction. From the UniProtKB/Swiss-Prot database [39] (release 2010\_06), we selected the enzyme sequences that: i) had been annotated with complete four-digit EC numbers, ii) were not fragment sequences and iii) had domains assigned to CATH [38] superfamilies in the Gene3D database [40]. A total of 332,021 enzyme domain sequences were obtained. In the following, an enzyme sequence refers to a protein domain sequence thus created, which was associated with a single CATH superfamily. The domain sequences were treated as independent sequences, although some of these were obtained from single multi-domain proteins. In order to obtain structural information, the 72,993 enzymes in the CATH database (ver. 3.3) were added to the 332,021 enzyme sequences. In each enzyme (as distinguished by the four-digit EC number) in each superfamily, all these sequences were clustered at a 95% sequence identity cutoff by using blastclust [78]. Also for each enzyme, a single representative structure was selected as the CATH S-level representative structure with the longest sequence length and the highest resolution. In the 95%-identity cluster that included the representative structure, the corresponding sequence was considered the representative of the cluster and in the other 95%-identity clusters, the longest sequence was selected as the representative. After the removal of redundancy, 201,708 sequences remained.

In the remaining sequences, a predictor was constructed for an enzyme if: 1) the enzyme belonged to a superfamily that contained at least one other enzyme in it, 2) the enzyme had a representative



**Figure 9. The rf-SDRs for acetylcholine esterase (AChE, EC 3.1.1.7, CATH domain: 1w76B00) in  $\alpha/\beta$ -hydrolase superfamily (CATH 3.40.50.1820).** The rf-SDRs are represented by balls and sticks, where carbon atoms are colored white, nitrogen atoms are blue, oxygen atoms are red and sulfur atoms are yellow. The active site gorge is partially represented by green surface. At the bottom of the active site gorge, the catalytic triads, which are not selected to be the rf-SDRs, are represented by balls and sticks and colored magenta. Many rf-SDRs are positioned around the catalytic gorge region. doi:10.1371/journal.pone.0084623.g009

structure and ten or more sequences and 3) a total of ten or more sequences were available for the other enzymes as negative data in the superfamily. We randomly selected 80% of the sequences from a given enzyme and 80% of the sequences from the other enzymes in the superfamily for training. The remaining 20% of the sequences were used as a test dataset. A total of 1121 enzymes over 306 CATH homologous superfamilies were selected for benchmarking.

### Calculations of attributes for classifiers

In addition to the BLAST [14,15] bit score, we used two types of scores as attributes: the scores calculated by using a full-length sequence and the scores at the functionally important positions in the alignment of a query sequence to a representative structure. The functionally important positions were defined to be the active sites, ligand binding sites and conserved site residues. In the following sections, we describe the selection of these positions and the score calculations.

**Determination of the alignment positions used for attribute calculations.** i) Active site and ligand binding residue positions from the literature and structural information: We obtained the literature information about active site residues from the Enzyme Catalytic-Mechanism Database (EzCatDB, ver. 20100722) [79] and the Catalytic Site Atlas (CSA, ver. 2.2.12) [45] database. All annotations in the EzCatDB and the original, hand-annotated entries derived from the primary literature in the CSA were used.

Ligand (substrate, cofactor, intermediate, products and their analogues) information in the Protein Data Bank (PDB) [80] was obtained from the EzCatDB and PROCOGNATE (ver. 1.6) [81] databases. All annotations in the EzCatDB and the cognate ligand entries with similarity scores higher than 0.5 in PROCOGNATE were used. Ligand binding residues were defined from complex structures by using LIGPLOT [82]. The residues that interacted with the ligands through both hydrogen bonds and hydrophobic interactions were considered as ligand binding residues. Ligand assignments to obsolete PDB entries were ignored.

We defined active site and ligand binding positions of each enzyme as the alignment positions, which were used by at least one PDB entry corresponding to that enzyme as an active site or a ligand-binding site, respectively. The position used as both active and ligand binding sites was defined to be an active site residue (ASR) position. The ASRs and ligand binding residues (LBRs) were mapped on to the representative structure for the calculation of attributes based on a multiple structural alignment, generated by MUSTANG [83], between the available complex structures and the representative.

ii) Conserved amino acid residue positions: For each enzyme in the training dataset, a multiple sequence alignment was generated by clustalw [84] and this alignment was aligned to the representative structure by FUGUE [41]. FUGUE performs sequence-structure comparison by utilizing environment-specific substitution tables (ESSTs). An ESST-based structural profile was calculated for the representative structure of each enzyme. To examine amino acid conservation, the entropy  $S_k$  for each alignment position  $k$  was calculated as

$$S_k = -\sum_{i=1}^{21} P^i \log P^i,$$

where  $i$  represents 20 types of amino acids plus a gap and  $P^i$  is the fraction of amino acid type  $i$  at this position. The top 10% conserved residue positions (CBRs) in one enzyme were selected

for the calculation of attributes. The positions where the fraction of the gap was above 20% were excluded from the entropy calculation. If the positions selected as CBRs were already defined as ASRs or LBRs, those positions were defined to be ASRs or LBRs.

Position-specific scoring matrices (PSSMs) [43] were also calculated from the multiple sequence alignments. The PSSM scores at the  $i$ th alignment positions were given by

$$P_{ij} = \sum_{k=1}^{20} W_{ki} \text{sim}(k,j),$$

where  $i$  is the alignment position,  $j$  and  $k$  are the amino acid types and  $\text{sim}(k,j)$  is the score in the BLOSUM 62 matrix between amino acid types  $j$  and  $k$  [42]. The logarithmic weight  $W_{ki}$  was defined, depending on occurrences of amino acid type  $k$  at position  $i$ , as

$$W_{ki} = \frac{\ln \left[ 1 - \left( \frac{\sum_{n=1}^N \delta_{ki}}{N+1} \right) \right]}{\ln \left( \frac{1}{N+1} \right)}, \delta_{ki} = \begin{cases} 1, & \text{amino acid type is } k \\ 0, & \text{if amino acid type is not } k \end{cases}$$

where  $N$  is the number of sequences in the alignment.

**Calculation of scores.** Given a query sequence, a BLAST search was performed against the sequences in the training dataset for each enzyme in each superfamily. The bit score for the top hit was used as an attribute for the predictors (see below). In the training mode, the bit score for the top hit, except for its own sequence, was used.

The other attributes were calculated based on an alignment between the query sequence and the representative structure by using three different scoring matrices: BLOSUM62, ESSTs and PSSMs. The latter two matrices were specific to each enzyme, as described in the previous section. The full-length sequence scores and the scores at ASRs, LBRs and CBRs were calculated.

### Construction of predictors and evaluation of performance

Decision trees were constructed by C4.5 [85] algorithms implemented in WEKA, a data mining software tool in Java (ver. 3.6.5) [86], with default parameters. Forests of decision trees were constructed by the random forests [31] algorithm implemented in R (ver. 2.15.1), a language and environment for statistical computing [87]. The default value was used for the number of attributes to split on at each node ( $\text{floor}(\sqrt{n})$ , where  $n$  is the number of input attributes), since the number of attributes was different for each enzyme. The number of trees constructed for each classifier was set to be 500, by comparing averaged out-of-bag (OOB) error rates obtained from the models with 250, 500 and 750 trees (data not shown). In construction of random forest for each enzyme, the importance score for each attribute was calculated. We selected the top  $3 * \text{floor}(\sqrt{n})$  ranked attributes as highly contributing attributes, analyzed their properties and defined the associated residues as random forest-derived specificity determining residues (rf-SDRs).

In order to evaluate prediction performance in regions where sequence identities between test and training sequences are low, we calculated the maximal test to training sequence identity (MTTSI) following Arakaki *et al.* [4] (see the reference for the detailed definition of MTTSI). Table S12 shows the number of positive and negative sequences in each MTTSI bin of the test set.

Given a predictor for enzyme EC *a.a.a.a*, a set of prediction results were obtained (by using the test sequences) and these results were divided into eight bins according to their MTTSI values. Then for each bin, precision = TP/(TP+FP) and recall = TP/(TP+FN) were calculated, where TP is the number of true positives, FP is the number of false positives and FN is the number of false negatives. Finally, these precision and recall values were averaged over all the enzymes, for which it was possible to define the performance measure (i.e., (TP+FP) > 0 for precision and (TP+FN) > 0 for recall within a bin).

### Functional entropy of a superfamily

For classifying superfamilies at the EC third-digit level, we defined the functional entropy  $S_{func}$  for each superfamily as follows:

$$P_{a.b.c} = \frac{n_{a.b.c}}{N},$$

$$S_{func} = \sum_{a.b.c} -P_{a.b.c} \ln P_{a.b.c}$$

where  $n_{a.b.c}$  is the number of predictors that share the first three digits of their EC numbers (*a.b.c*) and  $N$  is the total number of predictors in the superfamily. Using the functional entropy, superfamilies were classified into three groups: highly diverged ( $1.5 \leq S_{func}$ ), moderately diverged ( $0.5 \leq S_{func} < 1.5$ ) and least diverged ( $0 \leq S_{func} < 0.5$ ). The cutoff values were determined such that the occurrences of distinct EC numbers at the third-digit level within each superfamily approximately corresponded to one, two to four, and more than four, respectively (data not shown).

### Supporting Information

**Figure S1** Distribution of the number of enzyme predictors constructed in a superfamily. The region between 20 to 70 is expanded and represented in the figure. Fifteen superfamilies contained more than ten enzyme predictors and the largest superfamily was NAD(P)-binding Rossmann-like domain superfamily (CATH 3.40.50.720) with 65 predictors. (EPS)

**Figure S2** Distribution of the active site residues (ASRs) and ligand binding residues (LBRs) in all superfamilies. The white bars represent the ASRs and the light gray bars represent the LBRs. (EPS)

**Figure S3** Distributions of fractions of the rf-SDRs in active site residues (ASRs, A) and ligand binding residues (LBRs, B), observed in the superfamilies with low, medium and high degrees of functional diversity classified at the fourth-digit level of EC

numbers. The top and bottom of a box indicate 75th and 25th percentiles and the horizontal line in a box represents the median value. The top and bottom whiskers represent 90th and 10th percentiles.

(EPS)

**Table S1** Number of predictors in each CATH homologous superfamily. (XLSX)

**Table S2** Precision and recall of enzymes in each MTTSI bin. (DOCX)

**Table S3** Prediction performance of each predictor. (XLSX)

**Table S4** List of the rf-SDRs. (XLSX)

**Table S5** Differences of scoring matrices selected in the rf-SDRs. (DOCX)

**Table S6** Classifications of superfamilies at the third- and fourth-digit levels of EC numbers. (XLSX)

**Table S7** Averaged prediction performance for different classes of functional diversity at the third-digit level of EC numbers. (DOCX)

**Table S8** Averaged prediction performance for different classes of functional diversity at the fourth-digit level of EC numbers. (DOCX)

**Table S9** The average proportion of ASRs/LBRs to be selected as rf-SDRs for different classes of functional diversity at the third-digit level of EC numbers. (DOCX)

**Table S10** The average proportion of ASRs/LBRs to be selected as rf-SDRs for different classes of functional diversity at the fourth-digit level of EC numbers. (DOCX)

**Table S11** The number of rf-SDRs in ASRs, LBRs and CSRs. (DOCX)

**Table S12** The number of positive and negative queries in each MTTSI bin. (DOCX)

### Author Contributions

Conceived and designed the experiments: CN KM. Performed the experiments: CN. Analyzed the data: CN KM. Contributed reagents/materials/analysis tools: CN NN KM. Wrote the paper: CN KM.

### References

- Voet D, Voet JG (1990) Biochemistry: John Wiley and Sons, New York.
- Webb EC, NC-IUBMB (1992) Enzyme Nomenclature 1992, Recommendations of the Nomenclature Committee of the International Union of Biochemistry and Molecular Biology on the Nomenclature and Classification of Enzymes. San Diego, California: Academic Press.
- Wass MN, Barton G, Sternberg MJ (2012) CombFunc: predicting protein function using heterogeneous data sources. *Nucleic Acids Res* 40: W466–470.
- Arakaki AK, Huang Y, Skolnick J (2009) EFICAz2: enzyme function inference by a combined approach enhanced by machine learning. *BMC Bioinformatics* 10: 107.
- Kumar N, Skolnick J (2012) EFICAz2.5: application of a high-precision enzyme function predictor to 396 proteomes. *Bioinformatics* 28: 2687–2688.
- Redfern OC, Dessailly BH, Dallman TJ, Sillitoe I, Orengo CA (2009) FLORA: a novel method to predict protein function from structure in diverse superfamilies. *PLoS Comput Biol* 5: e1000485.
- Kumar C, Choudhary A (2012) A top-down approach to classify enzyme functional classes and sub-classes using random forest. *EURASIP J Bioinform Syst Biol* 2012: 1.
- Bray T, Doig AJ, Warwicker J (2009) Sequence and structural features of enzymes and their active sites by EC class. *J Mol Biol* 386: 1423–1436.
- Shen HB, Chou KC (2007) EzyPred: a top-down approach for predicting enzyme functional classes and subclasses. *Biochem Biophys Res Commun* 364: 53–59.
- Dobson PD, Doig AJ (2003) Distinguishing enzyme structures from non-enzymes without alignments. *J Mol Biol* 330: 771–783.
- Laskowski RA, Watson JD, Thornton JM (2005) ProFunc: a server for predicting protein function from 3D structure. *Nucleic Acids Res* 33: W89–93.
- Hulo N, Bairoch A, Bulliard V, Cerutti L, Cuče BA, et al. (2008) The 20 years of PROSITE. *Nucleic Acids Res* 36: D245–249.

13. Radiwojac P, Clark WT, Oron TR, Schnoes AM, Wittkop T, et al. (2013) A large-scale evaluation of computational protein function prediction. *Nat Methods* 10: 221–227.
14. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990) Basic local alignment search tool. *J Mol Biol* 215: 403–410.
15. Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, et al. (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res* 25: 3389–3402.
16. Todd AE, Orengo CA, Thornton JM (2001) Evolution of function in protein superfamilies, from a structural perspective. *J Mol Biol* 307: 1113–1143.
17. Tian W, Skolnick J (2003) How well is enzyme function conserved as a function of pairwise sequence identity? *J Mol Biol* 333: 863–882.
18. Addou S, Rentzsch R, Lec D, Orengo CA (2009) Domain-based and family-specific sequence identity thresholds increase the levels of reliable protein function transfer. *J Mol Biol* 387: 416–430.
19. Rost B, Liu J, Nair R, Wrzeszczynski KO, Ofra Y (2003) Automatic prediction of protein function. *Cell Mol Life Sci* 60: 2637–2650.
20. Bannert C, Wellie A, Aus dem Spring C, Schomburg D (2010) BrEFS: a flexible and automatic protocol to compute enzyme-specific sequence profiles for functional annotation. *BMC Bioinformatics* 11: 589.
21. Claudel-Renard C, Chevalet C, Faraut T, Kahn D (2003) Enzyme-specific profiles for genome annotation: PRIAM. *Nucleic Acids Res* 31: 6633–6639.
22. Nagao C, Nagano N, Mizuguchi K (2010) Relationships between functional subclasses and information contained in active-site and ligand-binding residues in diverse superfamilies. *Proteins* 78: 2369–2384.
23. George RA, Spriggs RV, Bartlett GJ, Gutteridge A, MacArthur MW, et al. (2005) Effective function annotation through catalytic residue conservation. *Proc Natl Acad Sci U S A* 102: 12299–12304.
24. Wass MN, Sternberg MJ (2008) ConFunc—functional annotation in the twilight zone. *Bioinformatics* 24: 798–806.
25. Kristensen DM, Ward RM, Lisewski AM, Erdin S, Chen BY, et al. (2008) Prediction of enzyme function based on 3D templates of evolutionarily important amino acids. *BMC Bioinformatics* 9: 17.
26. Tian W, Arakaki AK, Skolnick J (2004) EFICAZ: a comprehensive approach for accurate genome-scale enzyme function inference. *Nucleic Acids Res* 32: 6226–6239.
27. Capra JA, Singh M (2008) Characterization and prediction of residues determining protein functional specificity. *Bioinformatics* 24: 1473–1480.
28. Kalinina OV, Novichkov PS, Mironov AA, Gelfand MS, Rakhmaninova AB (2004) SDPpred: a tool for prediction of amino acid residues that determine differences in functional specificity of homologous proteins. *Nucleic Acids Res* 32: W424–428.
29. Lichtarge O, Bourne HR, Cohen FE (1996) An evolutionary trace method defines binding surfaces common to protein families. *J Mol Biol* 257: 342–358.
30. Addington Ta Fau - Mertz RW, Mertz Rw Fau - Siegel JB, Siegel Jb Fau - Thompson JM, Thompson Jm Fau - Fisher AJ, Fisher Aj Fau - Filkov V, et al. Janus: prediction and ranking of mutations required for functional interconversion of enzymes.
31. Breiman L (2001) Random Forests. *Machine Learning Journal*: 5–32.
32. Diaz-Uriarte R, Alvarez de Andres S (2006) Gene selection and classification of microarray data using random forest. *BMC Bioinformatics* 7: 3.
33. Svetnik V, Liaw A, Tong C, Culberson JC, Sheridan RP, et al. (2003) Random forest: a classification and regression tool for compound classification and QSAR modeling. *J Chem Inf Comput Sci* 43: 1947–1958.
34. Lec BJ, Shin MS, Oh YJ, Oh HS, Ryu KH (2009) Identification of protein functions using a machine-learning approach based on sequence-derived properties. *Proteome Sci* 7: 27.
35. Chen XW, Liu M (2005) Prediction of protein-protein interactions using random decision forest framework. *Bioinformatics* 21: 4394–4400.
36. Cai CZ, Han LY, Ji ZL, Chen YZ (2004) Enzyme family classification by support vector machines. *Proteins* 55: 66–76.
37. Syed U, Yona G (2009) Enzyme function prediction with interpretable models. *Methods Mol Biol* 541: 373–420.
38. Orengo CA, Michie AD, Jones S, Jones DT, Swindells MB, et al. (1997) CATH—a hierarchical classification of protein domain structures. *Structure* 5: 1093–1108.
39. UniProt Consortium (2012) Reorganizing the protein space at the Universal Protein Resource (UniProt). *Nucleic Acids Res* 40: D71–75.
40. Lees J, Yeats C, Redfern O, Clegg A, Orengo C (2010) Gene3D: merging structure and function for a Thousand genomes. *Nucleic Acids Res* 38: D296–300.
41. Shi J, Blundell TL, Mizuguchi K (2001) FUGUE: sequence-structure homology recognition using environment-specific substitution tables and structure-dependent gap penalties. *J Mol Biol* 310: 243–257.
42. Henikoff S, Henikoff JG (1992) Amino acid substitution matrices from protein blocks. *Proc Natl Acad Sci U S A* 89: 10915–10919.
43. Gribskov M, Luthy R, Eisenberg D (1990) Profile analysis. *Methods Enzymol* 183: 146–159.
44. Nagao C, Izako N, Soga S, Khan SH, Kawabata S, et al. (2012) Computational design, construction, and characterization of a set of specificity determining residues in protein-protein interactions. *Proteins* 80: 2426–2436.
45. Porter CT, Bartlett GJ, Thornton JM (2004) The Catalytic Site Atlas: a resource of catalytic sites and residues identified in enzymes using structural data. *Nucleic Acids Res* 32: D129–133.
46. Gutteridge A, Thornton JM (2005) Understanding nature's catalytic toolkit. *Trends Biochem Sci* 30: 622–629.
47. Henrissat B, Davies G (1997) Structural and sequence-based classification of glycoside hydrolases. *Curr Opin Struct Biol* 7: 637–644.
48. Nagano N, Porter CT, Thornton JM (2001) The (betaalpha)8 glycosidases: sequence and structure analyses suggest distant evolutionary relationships. *Protein Eng* 14: 845–853.
49. Cantarel BL, Coutinho PM, Rancurel C, Bernard T, Lombard V, et al. (2009) The Carbohydrate-Active EnZymes database (CAZy): an expert resource for Glycogenomics. *Nucleic Acids Res* 37: D233–238.
50. Dominguez R, Souchon H, Spinelli S, Dauter Z, Wilson KS, et al. (1995) A common protein fold and similar active site in two distinct families of beta-glycanases. *Nat Struct Biol* 2: 569–576.
51. Ducros V, Czjzek M, Belaich A, Gaudin C, Fierobe HP, et al. (1995) Crystal structure of the catalytic domain of a bacterial cellulase belonging to family 5. *Structure* 3: 939–949.
52. Dominguez R, Souchon H, Lascombe M, Alzari PM (1996) The crystal structure of a family 5 endoglucanase mutant in complexed and uncomplexed forms reveals an induced fit activation mechanism. *J Mol Biol* 257: 1042–1051.
53. Bortoli-German I, Haich J, Chippaux M, Barras F (1995) Informational suppression to investigate structural functional and evolutionary aspects of the *Erwinia chrysanthemi* cellulase EGZ. *J Mol Biol* 246: 82–94.
54. Navas J, Beguin P (1992) Site-directed mutagenesis of conserved residues of *Clostridium thermocellum* endoglucanase CelC. *Biochem Biophys Res Commun* 189: 807–812.
55. Belaich A, Fierobe HP, Baty D, Busetta B, Bagnara-Tardif C, et al. (1992) The catalytic domain of endoglucanase A from *Clostridium cellulolyticum*: effects of arginine 79 and histidine 122 mutations on catalysis. *J Bacteriol* 174: 4677–4682.
56. Zolotnitsky G, Cogan U, Adir N, Solomon V, Shoham G, et al. (2004) Mapping glycosidic hydrolase substrate subsites by isothermal titration calorimetry. *Proc Natl Acad Sci U S A* 101: 11275–11280.
57. Charnock SJ, Lakey JH, Virden R, Hughes N, Sinnott ML, et al. (1997) Key residues in subsite F play a critical role in the activity of *Pseudomonas fluorescens* subspecies cellulosa xylanase A against xylooligosaccharides but not against highly polymeric substrates such as xylan. *J Biol Chem* 272: 2942–2951.
58. Erskine PT, Norton E, Cooper JB, Lambert R, Coker A, et al. (1999) X-ray structure of 5-aminolevulinic acid dehydratase from *Escherichia coli* complexed with the inhibitor levulinic acid at 2.0 Å resolution. *Biochemistry* 38: 4266–4276.
59. Lang D, Thoma R, Henn-Sax M, Sterner R, Wilmanns M (2000) Structural evidence for evolution of the beta/alpha barrel scaffold by gene duplication and fusion. *Science* 289: 1546–1550.
60. Nagano N, Orengo CA, Thornton JM (2002) One fold with many functions: the evolutionary relationships between TIM barrel families based on their sequences, structures and functions. *J Mol Biol* 321: 741–765.
61. Sharma V, Grubmeyer C, Sacchetti JC (1998) Crystal structure of quinolinic acid phosphoribosyltransferase from *Mycobacterium tuberculosis*: a potential TB drug target. *Structure* 6: 1587–1599.
62. Fujimoto Z, Kaneko S, Momma M, Kobayashi H, Mizuno H (2003) Crystal structure of rice alpha-galactosidase complexed with D-galactose. *J Biol Chem* 278: 20313–20318.
63. List F, Sterner R, Wilmanns M (2011) Related (betaalpha)8-barrel proteins in histidine and tryptophan biosynthesis: a paradigm to study enzyme evolution. *ChemBiochem* 12: 1487–1494.
64. Reisinger B, Bocola M, List F, Claren J, Rajendran C, et al. (2012) A sugar isomerization reaction established on various (betaalpha)8-barrel scaffolds is based on substrate-assisted catalysis. *Protein Eng Des Sci* 25: 751–760.
65. Noda-Garcia L, Camacho-Zarco AR, Medina-Ruiz S, Gaytan P, Carrillo-Trippa M, et al. (2013) Evolution of Substrate Specificity in a Recipient's Enzyme Following Horizontal Gene Transfer. *Mol Biol Evol* 30: 2024–2034.
66. Henn-Sax M, Thoma R, Schmidt S, Hennig M, Kirschner K, et al. (2002) Two (betaalpha)8-barrel enzymes of histidine and tryptophan biosynthesis have similar reaction mechanisms and common strategies for protecting their labile substrates. *Biochemistry* 41: 12032–12042.
67. Duc AV, Kuper J, Gerloff A, von Kries JP, Wilmanns M (2011) Bisubstrate specificity in histidine/tryptophan biosynthesis isomerase from *Mycobacterium tuberculosis* by active site metamorphosis. *Proc Natl Acad Sci U S A* 108: 3554–3559.
68. Nakanishi T, Nakatsu T, Matsuoka M, Sakata K, Kato H (2005) Crystal structures of pyruvate phosphate dikinase from maize revealed an alternative conformation in the swiveling-domain motion. *Biochemistry* 44: 1136–1144.
69. Nardini M, Dijkstra BW (1999) Alpha/beta hydrolase fold enzymes: the family keeps growing. *Curr Opin Struct Biol* 9: 732–737.
70. Holmquist M (2000) Alpha/Beta-hydrolase fold enzymes: structures, functions and mechanisms. *Curr Protein Sci* 1: 209–235.
71. Sussman JL, Harel M, Frolow F, Oefner C, Goldman A, et al. (1991) Atomic structure of acetylcholinesterase from *Torpedo californica*: a prototypic acetylcholine-binding protein. *Science* 253: 872–879.
72. Harel M, Kryger G, Rosenberry TL, Mallerer WD, Lewis T, et al. (2000) Three-dimensional structures of *Trophilata melanogaster* acetylcholinesterase and of its complexes with two potent inhibitors. *Protein Sci* 9: 1063–1072.
73. Greenblatt HM, Guillou C, Guenard D, Argaman A, Boti S, et al. (2004) The complex of a bivalent derivative of galanthamine with torpedotoxin acetylcholinesterase displays drastic deformation of the active-site gorge:

- implications for structure-based drug design. *J Am Chem Soc* 126: 15405–15411.
74. Bourne Y, Taylor P, Radic Z, Marchot P (2003) Structural insights into ligand interactions at the acetylcholinesterase peripheral anionic site. *EMBO J* 22: 1–12.
  75. Harel M, Schalk I, Ehret-Sabatier L, Bouet F, Goeldner M, et al. (1993) Quaternary ligand binding to aromatic residues in the active-site gorge of acetylcholinesterase. *Proc Natl Acad Sci U S A* 90: 9031–9035.
  76. Vellom DC, Radic Z, Li Y, Pickering NA, Camp S, et al. (1993) Amino acid residues controlling acetylcholinesterase and butyrylcholinesterase specificity. *Biochemistry* 32: 12–17.
  77. Erdin S, Venner E, Lisewski AM, Lichtarge O (2013) Function prediction from networks of local evolutionary similarity in protein structure. *BMC Bioinformatics* 14 Suppl 3: S6.
  78. Dondoshansky I, Wolf Y (2002) Blastclust (NCBI Software Development Toolkit) Bethesda: NCBI.
  79. Nagano N (2005) EzCatDB: the Enzyme Catalytic-mechanism Database. *Nucleic Acids Res* 33: D407–412.
  80. Berman H, Henrick K, Nakamura H (2003) Announcing the worldwide Protein Data Bank. *Nat Struct Biol* 10: 980.
  81. Bashton M, Nobeli I, Thornton JM (2008) PROCOGNATE: a cognate ligand domain mapping for enzymes. *Nucleic Acids Res* 36: D618–622.
  82. Wallace AC, Laskowski RA, Thornton JM (1995) LIGPLOT: a program to generate schematic diagrams of protein-ligand interactions. *Protein Eng* 8: 127–134.
  83. Konagurthu AS, Whisstock JC, Stuckey PJ, Lesk AM (2006) MUSTANG: a multiple structural alignment algorithm. *Proteins* 64: 559–574.
  84. Thompson JD, Higgins DG, Gibson TJ (1994) CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res* 22: 4673–4680.
  85. Quinlan JR (1993) C4.5: Programs for Machine Learning. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc.
  86. Mark H, Eibe F, Geoffrey H, Bernhard P, Peter R, et al. (2009) The WEKA data mining software: an update. *SIGKDD Explor Newsl* 11: 10–18.
  87. R Development Core Team (2008) R: A language and environment for statistical computing. Vienna, Austria: R Foundation for Statistical Computing.





## Control of adhesion of human induced pluripotent stem cells to plasma-patterned polydimethylsiloxane coated with vitronectin and $\gamma$ -globulin

Ryotaro Yamada,<sup>1</sup> Koji Hattori,<sup>2</sup> Saoko Tachikawa,<sup>1</sup> Motohiro Tagaya,<sup>3</sup> Toru Sasaki,<sup>4</sup> Shinji Sugiura,<sup>2</sup> Toshiyuki Kanamori,<sup>2</sup> and Kiyoshi Ohnuma<sup>1,5,\*</sup>

Department of Bioengineering, Nagaoka University of Technology, 1603-1 Kamitomioka, Nagaoka, Niigata 940-2188, Japan,<sup>1</sup> Research Center for Stem Cell Engineering, National Institute of Advanced Industrial Science and Technology (AIST), Tsukuba Central 5th, 1-1-1 Higashi, Tsukuba, 5 Ibaraki 305-8565, Japan,<sup>2</sup> Department of Materials Science and Technology, Nagaoka University of Technology, 1603-1 Kamitomioka, Nagaoka, Niigata 940-2188, Japan,<sup>3</sup> Department of Electrical Engineering, Nagaoka University of Technology, 1603-1 Kamitomioka, Nagaoka, Niigata 940-2188, Japan,<sup>4</sup> and Top Runner Incubation Center for Academia-Industry Fusion, Nagaoka University of Technology, 1603-1 Kamitomioka, Nagaoka, Niigata 940-2188, Japan<sup>5</sup>

Received 19 December 2013; accepted 11 February 2014  
Available online xxx

Human induced pluripotent stem cells (hiPSCs) are a promising source of cells for medical applications. Recently, the development of polydimethylsiloxane (PDMS) microdevices to control the microenvironment of hiPSCs has been extensively studied. PDMS surfaces are often treated with low-pressure air plasma to facilitate protein adsorption and cell adhesion. However, undefined molecules present in the serum and extracellular matrix used to culture cells complicate the study of cell adhesion. Here, we studied the effects of vitronectin and  $\gamma$ -globulin on hiPSC adhesion to plasma-treated and untreated PDMS surfaces under defined culture conditions. We chose these proteins because they have opposite properties: vitronectin mediates hiPSC attachment to hydrophilic siliceous surfaces, whereas  $\gamma$ -globulin is adsorbed by hydrophobic surfaces and does not mediate cell adhesion. Immunostaining showed that, when applied separately, vitronectin and  $\gamma$ -globulin were adsorbed by both plasma-treated and untreated PDMS surfaces. In contrast, when PDMS surfaces were exposed to a mixture of the two proteins, vitronectin was preferentially adsorbed onto plasma-treated surfaces, whereas  $\gamma$ -globulin was adsorbed onto untreated surfaces. Human iPSCs adhered to the vitronectin-rich plasma-treated surfaces but not to the  $\gamma$ -globulin-rich untreated surfaces. On the basis of these results, we used perforated masks to prepare plasma-patterned PDMS substrates, which were then used to pattern hiPSCs. The patterned hiPSCs expressed undifferentiated-cell markers and did not escape from the patterned area for at least 7 days. The patterned PDMS could be stored for up to 6 days before hiPSCs were plated. We believe that our results will be useful for the development of hiPSC microdevices.

© 2014, The Society for Biotechnology, Japan. All rights reserved.

[Key words: Polydimethylsiloxane; Low-pressure air plasma; Microenvironment control; Microdevice; Serum-free culture; Feeder-free culture; iPSC cells; Cell adhesion; Competitive adsorption]

Human pluripotent stem cells (hPSCs), including both human embryonic stem cells and human induced pluripotent stem cells (hiPSCs), exhibit infinite self-renewal capacity and pluripotency (1–3). Because hiPSCs and embryonic stem cells generated by somatic cell nuclear transfer contain the donor's genetic information, medical applications of autologous stem cells offer the hope of rejection-free transplantation of tissues and patient-specific drug screening (2,3).

The development of new cell culture devices for patient-specific drug screening using hPSCs requires control of the microenvironment of the cells, including the spatiotemporal distribution of soluble factors, cell–cell interactions, and cell–substrate interactions; and microfabricated devices are increasingly being

developed for this purpose (4,5). Polydimethylsiloxane (PDMS) is one of the most popular biocompatible materials for such devices because this elastomer is non-toxic, chemically inert, transparent, and gas permeable (6). For the fabrication of microdevices, PDMS surfaces have often been modified by gas-phase processing methods including plasma treatment (in this paper, plasma refers to low-pressure air plasma, not blood plasma, unless otherwise stated), ultraviolet irradiation, chemical vapor deposition, and sputter coating of metal compounds (7). Plasma treatment is easy to carry out and is used for various purposes, including PDMS–PDMS and PDMS–glass bonding, cleaning PDMS surfaces, and facilitating the coating of surfaces with cell-adhesive extracellular matrix (ECM) proteins (8,9). Therefore, we frequently use plasma treatment in the fabrication of PDMS microdevices for cell culture (10,11).

One of the most fundamental requirements for PDMS microdevices for hPSC applications is that the cells adhere to the PDMS surface, because hPSCs form flat colonies on culture dishes and cannot maintain their pluripotency without adhesion (1,12). Although there have been many studies of adsorption of

\* Corresponding author at: Department of Bioengineering, Nagaoka University of Technology, 1603-1 Kamitomioka, Nagaoka, Niigata 940-2188, Japan. Tel./fax: +81 258 47 9454.

E-mail addresses: kohnuma@vos.nagaokaut.ac.jp, kyohnuma@gmail.com (K. Ohnuma).

cell-adhesive and non-cell-adhesive molecules, including ECM components, on biocompatible surfaces (13,14), the mechanism of adhesion remains to be revealed. The study of cell adhesion is complicated by the fact that the culture environment contains unknown amounts of various undefined molecules, including those in the ECM and in serum (15–17). Thus, investigation of how PDMS surface modifications and cell-adhesive and non-cell-adhesive proteins affect hPSC adhesion under defined culture conditions is urgently needed if medical applications of microfabricated devices for hPSCs are to be developed.

Here, we studied the effects of two proteins, vitronectin and  $\gamma$ -globulin, which is one of the most abundant protein in serum, on the adhesion of hiPSCs to plasma-treated and untreated PDMS surfaces under defined culture conditions. We chose vitronectin for three reasons: (i) hiPSCs adhere to tissue culture dishes coated with vitronectin (18,19); (ii) vitronectin in serum plays a major role in mediating adhesion of cells to the hydrophilic surface of glass, as reflected in the protein's name ("vitro" = "glass", "nectin" = "cell adhesion molecule") (20,21); and (iii) vitronectin may adsorb well on PDMS, because PDMS, like glass, is rich in Si–O bonds (6). In contrast,  $\gamma$ -globulin (an immunoglobulin) has the opposite adhesion properties and thus can be expected to block adsorption of vitronectin on PDMS for three reasons: (i) although adhesion of PSCs to polymers is mediated by integrins, cadherin, and glycans (16,19,22,23),  $\gamma$ -globulin has not been reported to mediate PSC adhesion; (ii)  $\gamma$ -globulin has a hydrophobic fragment crystallizable (Fc) region that is involved in adsorption on hydrophobic surfaces (7,24); and (iii) PDMS is rich in hydrophobic methylene groups (6).

We investigated the relationships between plasma treatment of PDMS surfaces, vitronectin and  $\gamma$ -globulin adsorption, and hiPSC adhesion under defined culture conditions using hESF9a, a serum- and feeder-free culture medium (25,26); this medium allowed us to study these relationships without masking by undefined factors derived from serum and feeder cells. We used the results of our initial investigations to pattern a PDMS surface with hiPSCs.

## MATERIALS AND METHODS

**Culture and subculture of hiPSCs** Two hiPSC cell lines, 201B7 (2) and 253G1 (27), were obtained from RIKEN BRC Cell Bank (Tsukuba, Japan) through the National BioResource Project for the Ministry of Education, Culture, Sports, Science and Technology, Japan. The 201B7 line was used unless otherwise stated. For all experiments, hiPSCs cultured in KSR-based medium on mouse embryonic fibroblast feeder cells were transferred to serum- and feeder-free culture conditions in hESF9a medium (11,25) on dishes coated with 2  $\mu$ g/mL fibronectin from bovine blood plasma (F-1141, Sigma–Aldrich, St. Louis, MO, USA) and were passaged at least once before use (11,26). For subculturing, the cells were detached from the culture dish by using 0.2–0.5 U/mL dispase (17105-041, Life Technologies, Grand Island, NY, USA) in hESF9a medium and replated in hESF9a medium with 5  $\mu$ M ROCK inhibitor (Y-27632, Wako Pure Chemical Industries, Osaka, Japan), which blocks dissociation-induced apoptosis of hiPSCs (12). The hESF9a medium was changed daily. For the adhesion experiments, hiPSCs were dissociated into single cells by incubation and trituration in 0.02% (w/w) ethylenediaminetetraacetic acid (EDTA) in PBS<sup>-/-</sup> and then plated in hESF9a solution with 5  $\mu$ M ROCK inhibitor.

**Preparation and plasma treatment of PDMS surfaces** PDMS prepolymer and curing agent (Sylgard 184, Dow Corning, Midland, MI, USA) were thoroughly mixed at a 10:1 weight ratio. To make PDMS sheets, we poured the mixture between two polyethylene terephthalate films separated with 0.5 mm rubber spacers and cured it in an oven at 120°C for 2 h. To make perforated masks, we perforated the sheet with 2-mm-diameter holes by using a hole punch. The resultant 0.5-mm-thick PDMS sheets and perforated masks were rinsed with ethanol and sterilized at 160°C for 2 h.

PDMS sheets with or without a perforated mask were hydrophilized by treatment with a low-pressure air plasma for 60 s (YHS-R, Sakigake-Semiconductor Co., Kyoto, Japan) after 5 min under vacuum (ultimate vacuum, 2 Pa; TA150XA, Tasco, Osaka, Japan). Between 30 min and 1 h later (unless otherwise stated), the perforated mask was removed, if one was used, and the PDMS sheet surface was coated with 5.5 mg/mL rabbit  $\gamma$ -globulin (011-000-002, Jackson ImmunoResearch Laboratories, Inc., West Grove, PA, USA), 0.6  $\mu$ g/cm<sup>2</sup> human blood plasma-derived vitronectin (2349-VN, R&D Systems), or both, and the coated sheet was incubated overnight at 37°C.

**Contact angle measurement and Fourier transform infrared spectroscopy** The water contact angle of the PDMS surfaces was analyzed in air by the sessile drop method using a droplet of distilled water (2  $\mu$ L). Droplets were photographed with a digital camera (CX3, Ricoh, Tokyo, Japan), and the angles were estimated by half-angle contact methods using ImageJ software (NIH, Bethesda, MD, USA).

For Fourier transform infrared (FT-IR) spectroscopy, a PDMS thin film was prepared as reported before (28). Briefly, 0.1 g of a mixture of PDMS was dissolved in 30 mL of chloroform, and the solution was spin-coated (3000 rpm, 20 s) onto an oxidized Si(111) substrate. The resulting PDMS film was cured at 65°C for 12 h. The surface chemical bonding of the PDMS thin film was analyzed by FT-IR spectroscopy (Nicolet, Thermo Fisher Scientific). Spectra were accumulated from 32 scans at a resolution of 1.0 cm<sup>-1</sup> in transmittance mode in the wavenumber range between 4,000 and 400 cm<sup>-1</sup>.

**Cell attachment assay and immunostaining** Attached living cells were stained with 1  $\mu$ M calcein AM (Dojindo, Kumamoto, Japan), a fluorescent dye that can be transported into living cells, for 20 min at 37°C. For immunostaining of the surface of PDMS coated with proteins, the PDMS surface was rinsed with PBS containing 0.5 mM CaCl<sub>2</sub> and 0.5 mM MgCl<sub>2</sub> (PBS<sup>+/+</sup>), fixed in 4% formaldehyde (Sigma–Aldrich) with 0.5 mM MgCl<sub>2</sub> and 0.5 mM CaCl<sub>2</sub>, and reacted with primary antibodies overnight; the primary antibodies were then visualized with secondary antibodies (Table S1). The antibodies were diluted in PBS<sup>+/+</sup> containing 10 mg/mL bovine serum albumin. For immunocytochemistry, hiPSCs plated on conventional culture dishes were rinsed with PBS<sup>+/+</sup>, fixed in 4% formaldehyde with 0.5 mM MgCl<sub>2</sub> and 0.5 mM CaCl<sub>2</sub>, permeabilized, blocked with PBS<sup>+/+</sup> containing 0.2% Triton X-100 and 10 mg/mL bovine serum albumin, and reacted with primary antibodies, which were visualized with secondary antibodies (Table S1). The antibodies were diluted in PBS<sup>+/+</sup> containing 0.2% Triton X-100 and 10 mg/mL bovine serum albumin. Nuclei were stained with 0.4  $\mu$ M 4',6-diamidino-2-

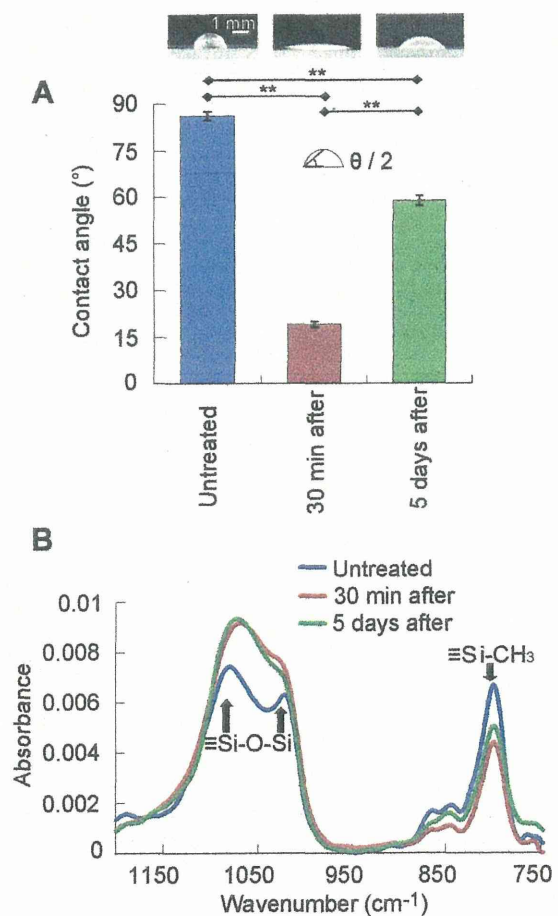


FIG. 1. Effects of plasma treatment of PDMS surfaces. (A) Contact angle ( $\theta$ ) of water. Data are means  $\pm$  SE ( $n = 5$ ). \*\*  $P < 1 \times 10^{-10}$ , Tukey's multiple comparison. Insets are photographs of 2  $\mu$ L water droplets. (B) FT-IR absorbance spectra of PDMS before plasma treatment (blue) and 30 min (red) or 5 days (green) after 60-s plasma treatment. The arrows represent functional groups whose absorbances were changed by plasma treatment. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)