

2 **DomSVR: domain boundary prediction with support vector**
3 **regression from sequence information alone**

4 Peng Chen · Chunmei Liu · Legand Burge · Jinyan Li ·
5 Mahmood Mohammad · William Southerland ·
6 Clay Gloster · Bing Wang

7 Received: 23 September 2009 / Accepted: 25 January 2010
8 © Springer-Verlag 2010

9 **Abstract** Protein domains are structural and fundamental
10 functional units of proteins. The information of protein
11 domain boundaries is helpful in understanding the evolu-
12 tion, structures and functions of proteins, and also plays an
13 important role in protein classification. In this paper, we
14 propose a support vector regression-based method to
15 address the problem of protein domain boundary identifi-
16 cation based on novel input profiles extracted from
17 AAindex database. As a result, our method achieves an
18 average sensitivity of ~36.5% and an average specificity
19 of ~81% for multi-domain protein chains, which is overall

better than the performance of published approaches to 20
identify domain boundary. As our method used sequence 21
information alone, our method is simpler and faster. 22
23

Keywords Domain boundary prediction · 24
Support vector regression · AAindex · 25
Principal component analysis 26

Introduction 27

Protein domains are importantly independent units of 28
protein tertiary structures and have been studied exten- 29
sively in recent decades. Edelman et al. studied the struc- 30
tures of immunoglobulins and first proposed some 31
important hypothesizes on domain structures (Edelman 32
1973; Porter 1973). Wetlaufer (1973) subsequently pro- 33
posed the concept of domain and defined domains as stable, 34
compact, and autonomously folding structures of proteins 35
based on a thorough investigation of immunoglobulins. A 36
domain can span an entire polypeptide chain or be a sub- 37
unit of a chain which can be folding into a stable tertiary 38
structure independently (Levitt and Chothia 1976). 39

Typically, most domains have a single continuous 40
polypeptide segment, while a few others consist of several 41
discontinuous segments. Furthermore, many protein chains 42
consist of more than one structural domains, all of them 43
form independently compact structures (Wetlaufer 1973). 44
Moreover, it is observed that a large protein may get its 45
optimal protein folding by domain formation, when giving 46
an observed random distribution of hydrophobic residues in 47
large proteins (George and Heringa 2002a, b). Actually, 48
each domain contains an individual hydrophobic core that 49
is built from secondary structures (Zhou et al. 1999). 50
Residues in hydrophobic core are more conserved than 51

A1 P. Chen (✉) · C. Liu · L. Burge
A2 Department of Systems and Computer Science,
A3 Howard University, 2400 Sixth Street, NW,
A4 Washington, DC 20059, USA
A5 e-mail: pchen1978@gmail.com

A6 P. Chen · J. Li
A7 Bioinformatics Research Center, School of Computer
A8 Engineering, Nanyang Technological University,
A9 Singapore 639798, Singapore

A10 M. Mohammad
A11 Department of Mathematics, Howard University,
A12 2400 Sixth Street, NW, Washington, DC 20059, USA

A13 W. Southerland
A14 Department of Biochemistry, Howard University,
A15 2400 Sixth Street, NW, Washington, DC 20059, USA

A16 C. Gloster
A17 Department of Electrical and Computer Engineering,
A18 Howard University, 2400 Sixth Street, NW,
A19 Washington, DC 20059, USA

A20 B. Wang
A21 School of Electrical Engineering and Information,
A22 Anhui University of Technology, Hudong Road 59,
A23 Ma'anshan 243002, Anhui, People's Republic of China

Author Proof

residues at the surface in a protein family unless the latter are involved in the functions of the protein (Zhou et al. 1999).

Previous works on the prediction of protein domain boundaries are roughly classified into two categories: template-based methods (Altschul et al. 1997; Cheng et al. 2006; Gewehr and Zimmer 2005; Marchler-Bauer et al. 2007; Marsden et al. 2002; Orengo et al. 1997) and ab initio methods (Copley et al. 2002; Dumontier et al. 2005; Galzitskaya and Melnik 2003; George and Heringa 2002b; Nagarajan and Yona 2004; Sikder and Zomaya 2006; Sim et al. 2005; Suyama and Ohara 2003). Template-based methods aim to predict domain boundaries using sequence alignment (Marchler-Bauer et al. 2007), secondary structure alignment (Cheng et al. 2006; Marsden et al. 2002), or other profile alignments. They align target profiles against profiles in a domain database. Among template-based methods, conserved domain database (CDD) (Marchler-Bauer et al. 2007) locates residues in domain boundaries using a search tool, reverse position-specific BLAST (RPS-BLAST). With CDD method, firstly, query sequences are compared to databases of position-specific score matrices (PSSMs). Secondly, *E* values are obtained in much the same way as in the PSI-BLAST application (Altschul et al. 1997). Overlapping domain hits are finally obtained by the sort of the *E* values. DomSSEA (Marsden et al. 2002) predicts domain boundaries by aligning the predicted secondary structures of target sequences against a database of observed secondary structures of chains that have known domain boundaries (Orengo et al. 1997). SSEPDomain method predicts domains with the alignment information of secondary structures and profile–profile as well as pattern searches (Gewehr and Zimmer 2005).

Most ab initio methods aim to identify protein domain boundaries based on the information of the properties of residues in protein chains using various machine learning techniques. Among them, CHOPnet addresses some issues in domain annotation with evolutionary information, amino acid composition, and amino acid flexibility (Copley et al. 2002); SnapDRAGON predicts domain boundaries using a distance geometry-based folding technique with a 3D domain assignment algorithm (George and Heringa 2002b); Galzitskaya and Melnik (2003) propose a simple approach to identify domain boundaries in proteins using side chain entropy of a residue region; DomCut's method predicts inter-domain linker regions using amino acid sequence information (Suyama and Ohara 2003); Nagarajan and Yona (2004) propose a neural network-based method to detect domain structure of a protein, which uses the information from multiple sequence alignments analysis, position-specific properties of amino acids, and predicted secondary structures; PRODO (Sim et al. 2005) uses a

neural network method with information from position-specific scoring matrix (PSSM) generated by PSI-BLAST (Altschul et al. 1997); Armadillo aims to predict domain boundaries by converting protein sequences to smoothed numeric profiles based on domain linker propensity index (DLI) from amino acids' composition (Dumontier et al. 2005); Dovidchenko et al. (2007) propose a simple and fast method with the use of a minimal number of amino acid sequence alone; DomainDiscovery detects domain boundaries by the use of support vector machines with sequence information including a PSSM, secondary structure, solvent accessibility information and inter-domain linker index (Sikder and Zomaya 2006); DOMpro applies recursive neural network to predict domain boundaries with evolutionary information, solvent evolutionary information, solvent accessibility information, and secondary structure (Cheng et al. 2006); Ye et al. (2007) present a Back-Propagation (BP) neural network approach to predict the domain boundaries with various property profiles; recently, Yoo et al. (2008) develop a new improved general regression network (IGRN) model to detect domain boundaries using a PSSM, secondary structure, information, and inter-domain linker index.

However, the accuracy of predicting multi-domain boundaries is considerably less than 40% in spite of great development on domain boundary prediction in the past years by the use of a large number of machine learners. Therefore, novel machine learning-based approaches should be developed to accurately identify protein domain boundaries.

Most previous work in the prediction of domain boundaries has been on the so-called “classification problem”. In this case, residues are assigned to one of two states, domain boundary or non-domain boundary, with arbitrary cutoff thresholds. However, the selection of thresholds is neither objective nor optimal, and the decomposition of residues into two classes decreases the prediction accuracy. To overcome such disadvantages, we predict domain boundary value for each residue. That is, our method predicts a series of real values representing residues in a protein sequence (also regarded as the boundary profile). In this paper, we develop an accurate, fast, and reliable ab initio protein domain boundary predictor, named as DomSVR, by the use of support vector regression (SVR) starting from protein sequence alone. The method just uses profiles extracted from AAindex database (Kawashima et al. 2008). Our proposed method DomSVR achieves an average sensitivity of ~36.5% and an average specificity of ~81% for multi-domain protein chains, which is overall better than the performance of published approaches to identify domain boundary. As our method used sequence information alone, our method is simpler and faster.

158 **Methods**

159 Dataset preparation

160 Our model is trained and tested on the dataset extracted from
 161 DOMpro method (Cheng et al. 2006). In this paper, we only
 162 consider proteins with more than one domain. Finally, 354
 163 multi-domain proteins are used to evaluate our proposed
 164 method of protein domain boundary prediction. In the
 165 dataset, sequence identity of each two protein chains is less
 166 than 25%. Moreover, all protein chains contain more than 40
 167 amino acid residues. The dataset consists of 282 two-domain
 168 chains, 50 three-domain chains, and 22 chains having more
 169 than three domains. The dataset can be found at our website:
 170 <http://mail.ustc.edu.cn/~bigeagle/DomSVR/index.htm>.

171 Creation of amino acid physicochemical profiles
 172 for inputs of SVR predictor

173 In this work, we encode input vectors of SVR predictor
 174 using amino acid profiles extracted from AAindex database
 175 (Kawashima et al. 2008). First, we need to assign physical
 176 and chemical properties to amino acid residues. Vectors of
 177 suitable amino acid physicochemical properties will then
 178 be created and be used for the domain boundary assign-
 179 ment. The physicochemical properties of amino acid resi-
 180 dues include inter-residue contact energy, secondary
 181 structure, residue charge, and other properties. In addition,
 182 the simple forms of the vectors make the entire algorithm
 183 robust, fast, and easy to apply.

184 The AAindex database contains a large number of
 185 experimental indexes, representing a large variety of
 186 physicochemical and biological properties of the amino
 187 acids. The AAindex1 section of the amino acid index
 188 database collects published indices together with the result
 189 of cluster analysis using the correlation coefficient as the
 190 distance between two indices (Kawashima et al. 2008). The
 191 section currently contains 544 indices, excluding all
 192 empirically derived propensities of amino acids. Taking all

193 these 544 amino acid properties as input features for a
 194 predictor may cause over-fitting. In order to distinguish and
 195 separate significant data and then construct our profile
 196 vectors, we applied principal component analysis (PCA)
 197 (Jolliffe 2002) on these properties. PCA is often used to
 198 reduce the dimensionality of a given dataset to lower
 199 dimensions for analysis. It can then produce a new set of
 200 principal components, which account for the top largest
 201 variations of the original data. PCA takes linear combina-
 202 tions of the data complying with the rule that the first
 203 principal component accounts for the maximum variation,
 204 the second principal component accounts for the next
 205 maximum variation which is subject to being orthogonal to
 206 the first one, the third one has the third maximum variation
 207 subject to being orthogonal to the first two, and so on.
 208 Nineteen principal components were created which account
 209 for 99.99% of the variance in the AAindex1 dataset. Among
 210 those components, the top four components account for
 211 93.78% of the experimental data variation. Using only four
 212 principal component vectors as shown in Table 1, the entire
 213 original dataset of properties is described with an approxi-
 214 mate 6.22% loss of variation. Thus, the dimensionality of
 215 the original data is significantly reduced. The first principal
 216 component, PrinComp1, which solely accounts for 55% of
 217 the data variation, has a strong correlation to inter-residue
 218 contact energy property (Miyazawa and Jernigan 1999).
 219 The second component, PrinComp2, is correlated to sec-
 220 ondary structure propensities of amino acids (Munoz and
 221 Serrano 1994). The third component, PrinComp3, is cor-
 222 related to entire chain composition of amino acids (Fukuchi
 223 and Nishikawa 2001). Finally, PrinComp4 is mainly cor-
 224 related to conformational and nucleation properties of
 225 individual amino acids (Rackovsky and Scheraga 1982).

226 For protein chain with L residues, in the case of Prin-
 227 Comp1 profile, each residue is encoded as the central
 228 residue in a sliding window with nine residues along the
 229 peptide chain. Then, the central residue is represented by a
 230 1×9 vector, and the value for each element of the vector
 231 corresponds to specific amino acid type in PrinComp1.

Table 1 The top four principal component profiles and the variation account rates

Profile	A/R	N/D	C/Q	E/G	H/I	L/K	M/F	P/S	T/W	Y/V	Rate (%)
PrinComp1	-81.9	-280.1	460.7	-257	-19.5	220.1	316.3	-262.3	-44.8	125	51.01
	-269.3	-134.5	-277.7	-260.5	271.5	-350.2	408.9	-262.3	467.8	229.8	
PrinComp2	357.2	-66.5	102.8	-101.2	-257.5	203.3	-140.9	-30.5	112.7	-178.2	25.45
	-276	-20.2	-209.1	377	74.5	-77.8	-30.3	189.5	-270	241.1	
PrinComp3	-55.8	-86	214.8	150.6	-155	-67.1	-71	-187.4	155.4	-76.3	10.09
	-18.4	243	-16.3	-105.7	-82	212.6	-35.8	-44.1	10.6	13.8	
PrinComp4	-26.8	55.2	209.6	-3.4	98.9	-179.4	100.1	151.9	31.1	-104.4	7.23
	-137.6	95.3	51.3	48.7	-58.3	-185.6	-67.8	28.5	-28.3	-78.9	

Each principal component profile needs to be equalized by normalized itself when applying to create input vectors for SVR predictor

232 Therefore, the protein chain is represented by a $L \times 9$
 233 matrix which corresponds to a real value vector $L \times 1$,
 234 where each residue is assigned to a real value that measures
 235 the sequence distance between the residue and the central
 236 residue of its closest domain boundary.

237 The outputs of SVR predictor

238 The identification of domain boundaries for each protein
 239 chain can be viewed as a binary regression problem. Each
 240 residue along the polypeptide chains is encoded by AA-
 241 index amino acid profiles and assigned a real target value.
 242 Following the conventions used in prior work (Cheng et al.
 243 2006; Liu and Rost 2004; Marsden et al. 2002), suppose
 244 that residues within more than 20 continuous amino acids
 245 of a domain boundary are regarded as domain boundary
 246 residues, and non-domain boundary residues otherwise.

247 Actually SVR is particularly suitable for solving such
 248 regression problem. Assigned real value to a residue as
 249 target can be more efficient and effective than the assign-
 250 ment of classification value 1 or 0 as target. In this work, a
 251 residue is assigned to a domain boundary (DB) value,
 252 which measures the residue distancing away from its
 253 closest domain boundary in sequence. The assignment for
 254 residue i is shown in the following form:

$$DB_i = \begin{cases} \frac{cb_m - |i - cb_m|}{cb_m} & \text{if } i \text{ in boundary} \\ -\frac{|i - r_{end}|}{r_{end} - r_{start}} & \text{if } i \text{ in non-boundary near} \\ & \text{the N-termini} \\ -\frac{|i - r_{start}|}{r_{end} - r_{start}} & \text{if } i \text{ in non-boundary near} \\ & \text{the C-termini} \\ -\frac{cnb_n - |i - cnb_n|}{cnb_n} & \text{Otherwise} \end{cases}, \quad (1)$$

256 where DB_i denotes the DB value for residue i , cb_m indicates
 257 the sequence position of central residue m in domain
 258 boundary cb if cb existed, cnb_n means the sequence position
 259 of central residue n in non-boundary cnb , while r_{start} and
 260 r_{end} stand for the sequence positions of the starting and the
 261 end residues in the non-boundary sequence, respectively.

262 The form of Eq. 1 is a triangular distribution with
 263 respect to residue position in primary sequence. Central
 264 residue in domain boundary is assigned to a bigger value,
 265 while the more far away from the boundary the more small
 266 value the residue is assigned to. Finally, the target vector
 267 DB also needs to be normalized to equalize itself.

268 For each residue in protein chains, in summary, vector to
 269 be input into SVR is represented as an array X_i , where each
 270 element in the array corresponds to amino acid type of each
 271 AAindex profile, while the corresponding target DB_i is
 272 another real value which is assigned by Eq. 1 in terms of
 273 the sequence distance between residue i and its closest
 274 domain boundary. Similar to most other machine learners,
 275 DomSVR method aims to learn the mapping from the input

array X onto the corresponding target array DB . Suppose
 that O is an output array of SVR, DomSVR is trained to
 make the output O as close as possible to the target DB .

Approach

Support vector regression aims to apply support vector
 machine to regression problems by introducing an alternative
 loss function. Likely as SVM approach (Chen et al. 2007),
 linear regression of SVR is performed in a high-dimensional
 feature space mapped from complex data with a non-linear
 mapping (Gunn 1998). With SVR, a ϵ -insensitive loss func-
 tion is used where only errors greater than a predefined
 parameter ϵ are considered in the loss function. Readers can
 refer to (Drucker et al. 1996; Gunn 1998) for more details.

Consider the problem of learning a set of data, (X_i, DB_i) ,
 such that $X_i \in \mathbb{R}^n$ is an input vector which characterizes a
 residue along protein chains, and $DB_i \in \mathbb{R}$ is a real target
 value which represents its associated boundary value mea-
 suring the separation between the residue i and the closest
 domain boundary in sequence, with a linear function,

$$f(X) = \langle w, X \rangle + b. \quad (2)$$

The optimized parameters w and b can be obtained by
 minimizing the following objective function:

$$\emptyset(w, \xi) = \frac{1}{2} \|w\|^2 + C \sum_i (\xi_i^- + \xi_i^+). \quad (3)$$

where C is a regularization constant that balances training
 errors and model complexity, and ξ^- and ξ^+ are slack
 variables representing upper and lower constraints which
 used to measure the deviation of samples outside the ϵ -
 insensitive zone.

In this work, we adopt an ϵ -insensitive loss function,

$$L_\epsilon(DB) = \begin{cases} 0 & \text{if } |f(X) - DB| - \epsilon \\ |f(X) - DB| - \epsilon & \text{Otherwise} \end{cases}. \quad (4)$$

To solve the optimization problem, therefore, two
 Lagrange multipliers α_i and α_i^* are applied and the solution
 is given by

$$\begin{aligned} \text{Maximize} & \quad -\frac{1}{2} \sum_{i=1}^L \sum_{j=1}^L (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) Q_{ij} \\ & \quad + \sum_{i=1}^L \alpha_i (DB_i - \epsilon) - \alpha_i^* (DB_i + \epsilon) \\ \text{subject to} & \quad 0 \leq \alpha_i, \alpha_i^* \leq C, \quad i = 1, \dots, L \\ & \quad \text{and } \sum_{i=1}^L (\alpha_i - \alpha_i^*) = 0. \end{aligned} \quad (5)$$

where $Q_{ij} = K(x_i, x_j) \equiv \emptyset(x_i)^T \emptyset(x_j)$.

Finally the decision function is

$$\sum_{i=1}^L (\alpha_i - \alpha_i^*) K(X_i, X) + b. \quad (6)$$

Once the Lagrange multipliers α_i and α_i^* and the bias b
 are determined from the training data, Eq. 6 can be applied

Author Proof

315 to predict the domain boundary values for a test protein
 316 chain.

317 As a result, our model infers the domain boundary
 318 regions from predictions of domain boundary values for a
 319 test protein chain. The larger the prediction value is, the
 320 more possible the corresponding residue is belonging to
 321 domain boundary. In this work, a series of continuous
 322 residues are considered to be in domain boundary if the
 323 residue amount is more than 20 and their DB values are
 324 larger than other neighboring ones. At the same time, a
 325 series of continuous residues with bigger DB values are
 326 ignored if the residue amount is less than 5. Moreover, two
 327 inferred boundary regions that separate less than 10 resi-
 328 dues should be merged into one region. The test chain is
 329 then cut into domain regions linked by boundary region
 330 (regions).

331 Evaluation measures

332 To evaluate our method, three measurements are used to
 333 evaluate the performance of the predictor: criteria of sensi-
 334 tivity (Sen), specificity (Spec), and accuracy (Acc) (Baldi
 335 et al. 2000; Saini and Fischer 2005). They are defined as
 336 follows:

$$\text{Sen} = \frac{\text{TP}}{\text{TP} + \text{FN}}, \quad \text{Spec} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \quad \text{Acc} = \frac{\text{TP} + \text{TN}}{N_{\text{total}}} \quad (7)$$

338 where TP denotes the number of true positives (residues in
 339 domain boundaries), FP denotes the number of false posi-
 340 tives, TN stands for the number of true negatives (residues
 341 in non-domain boundaries), and N_{total} stands for the num-
 342 ber of total residues.

343 When assessing predictor with respect to domain
 344 boundary, evaluation is based on the above measures of
 345 Sen and Spec and, for the assessment with respect to
 346 domain number, measure of accuracy is the ratio of the
 347 number of chains whose domain number was predicted
 348 correctly to that of total protein chains.

349 Results

350 Domain boundary distribution

351 In this work, there are total 354 protein chains, each of
 352 which contains more than one domain. Figure 1 shows the
 353 distribution of sequence positions of residues at the center
 354 of domain boundaries. Most domain boundaries are far
 355 from the start and the end of the protein sequences. The
 356 distribution is helpful for limiting random noise of outputs
 357 from domain boundary prediction methods and further
 358 improves the identification rate of domain residues.

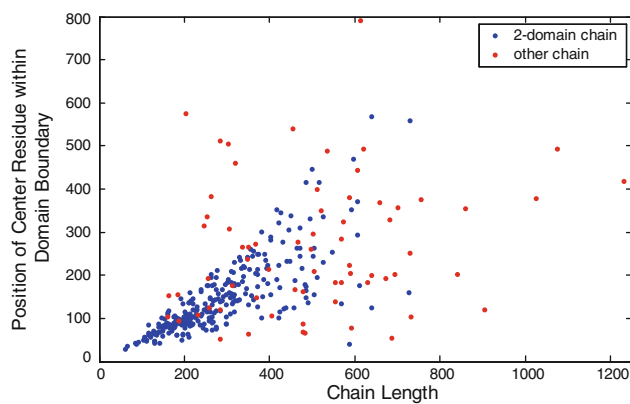


Fig. 1 Distribution of sequence positions of residues at the center of domain boundaries. Blue dot denotes two-domain chain while red dot stands for protein chain containing more than two domains

359 Figure 2 shows chain length distributions of multi-
 360 domain chains in the non-redundant set. From Fig. 2, the
 361 length distributions of multi-domain chains are not dis-
 362 crete, which has implications in domain prediction. As
 363 chain length increases, the likelihood of the chain having a
 364 multi-domain conformation almost increases. Most two-
 365 domain chains contain 100–200 amino acids. Most of
 366 three-domain chains contain 200–700 amino acids. Fur-
 367 thermore, chains containing more than 800 amino acid
 368 residues always have four or more domains.

369 The output from domain boundary predictor is quite
 370 noisy. To limit random noises that come from false positive
 371 hits and false negative hits, smoothing technique is used to
 372 correct the random fluctuation of outputs for neighboring
 373 residues (Goodall 1990). The smoothing technique is
 374 accomplished by averaging over a window around each
 375 residue position. For instance, Fig. 3 shows a case study of
 376 prediction for protein chain PDB:1qu6A, where each resi-
 377 due is assigned a state (boundary/not boundary) by a cutoff
 378 threshold at 0.5 to the output of model. A residue will be
 379 assigned to 1 (boundary state) when the corresponding
 380 output is larger than the threshold and, 0 (not boundary
 381 state) otherwise. After smoothing the outputs for each
 382 residue, the center of the domain boundary was predicted at
 383 residue 80 and the domain number was also correctly
 384 predicted. Figure 3 also illustrates how smoothing tech-
 385 nique helps reducing noises found in the raw outputs from
 386 the model. It is evident from Fig. 3 that the domain
 387 boundary threshold used to define the two classes (domain
 388 boundary and non-domain boundary) strongly affects the
 389 absolute classification results.

390 Performance of the PCA profiles

391 Figures 4, 5, 6, 7, and 8 show the ROC analysis of protein
 392 chains in CATH according to class membership, with the
 393 top four principal components being used as property

Author Proof

Fig. 2 Chain length distributions as observed in the CATH representative set used in this study. Intervals were calculated with a width of 100 residues. The domain frequencies were used to calculate probabilities of predicted domain sizes

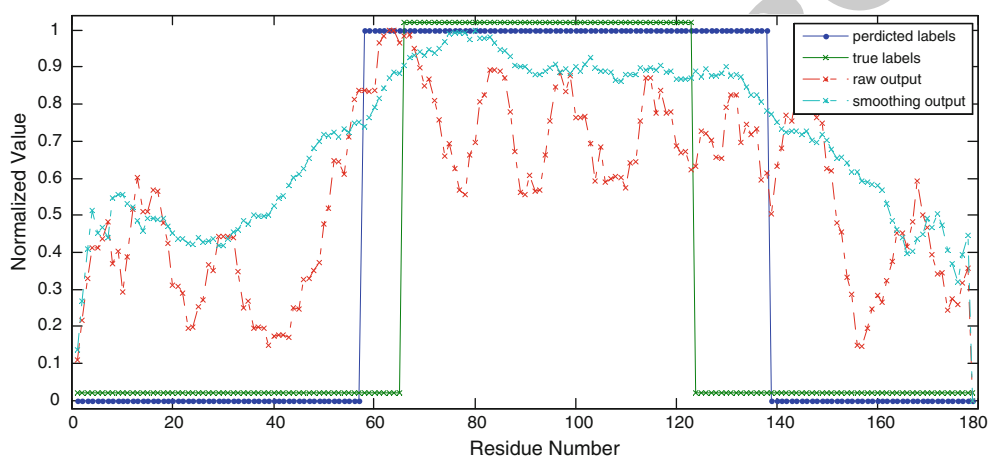
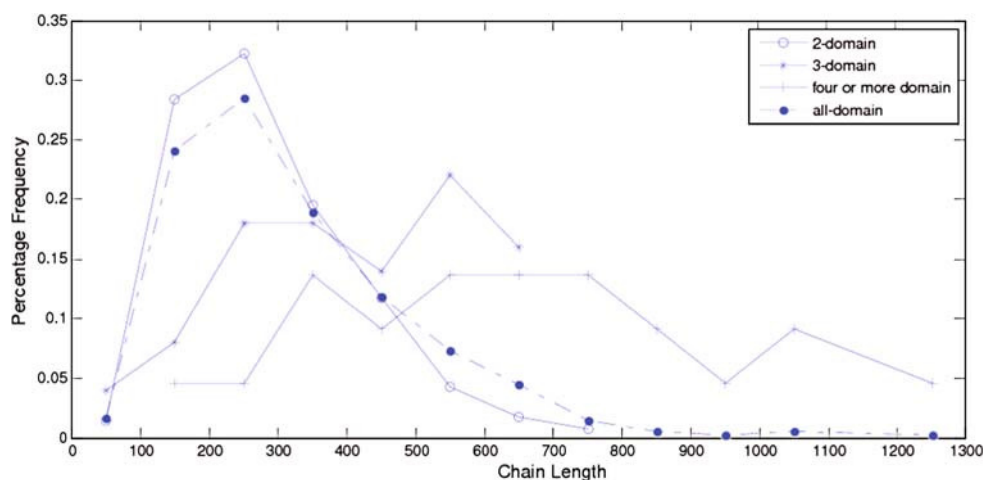


Fig. 3 Comparison of raw and smoothing outputs from SVR model for protein chain 1qu6A. The protein chain has 179 residues and contains two domains lined by a domain boundary. The center of the domain boundary is at residue 94. The two types of outputs are

normalized to the range [0, 1]. The two *square curves* denote the two kinds of residue labels. One is true labels describing residues' states (boundaries/not boundaries); the other is predicted labels

394 descriptors. Based on CATH architecture, protein chains in
395 our dataset are classified into four classes, i.e., mainly
396 alpha, mainly beta, alpha and beta, and fewer secondary
397 structure (SS). If all domains of a protein chain belong to
398 one CATH class, the chain is classified into the same class.
399 Inversely, if domains of a protein chain belong to different
400 CATH classes, the chain is classified into class "Others".

401 It is clearly shown that all the four profiles behave
402 similar in their predictive ability. The average accuracy
403 increases with the increase of the threshold, and all pre-
404 dictors reach high accuracy near the value of 0.7 for all
405 protein classes. However, many key differences of their
406 performance should be noted. An increase of the cutoff
407 threshold positively affects performance of the domain
408 boundaries prediction. The tradeoff for the increase of the
409 sensitivity is the dramatic decrease of the specificity for
410 almost all the four principal component profiles, as illus-
411 trated in Figs. 4, 5, 6, 7, and 8. In other words, from Eq. 7,

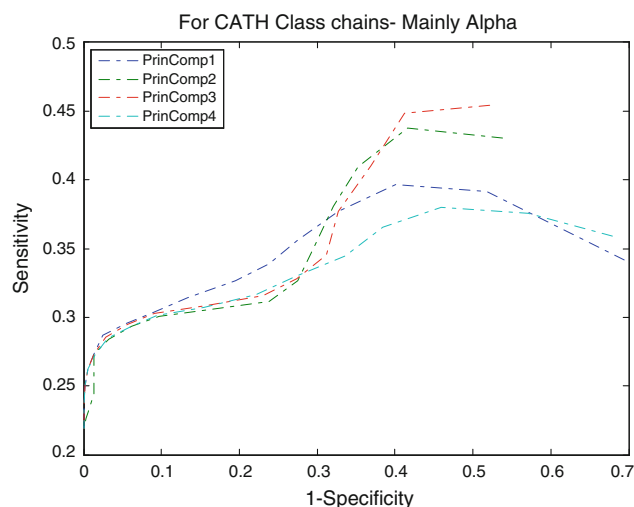


Fig. 4 ROC analysis for mainly alpha proteins with respect to threshold

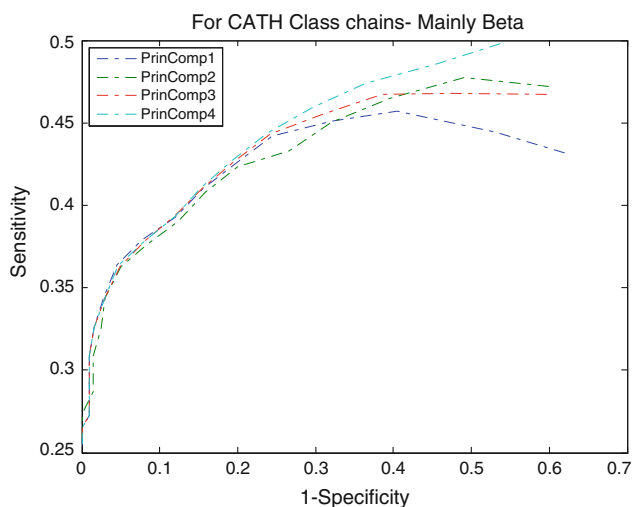


Fig. 5 ROC analysis for mainly beta proteins with respect to threshold

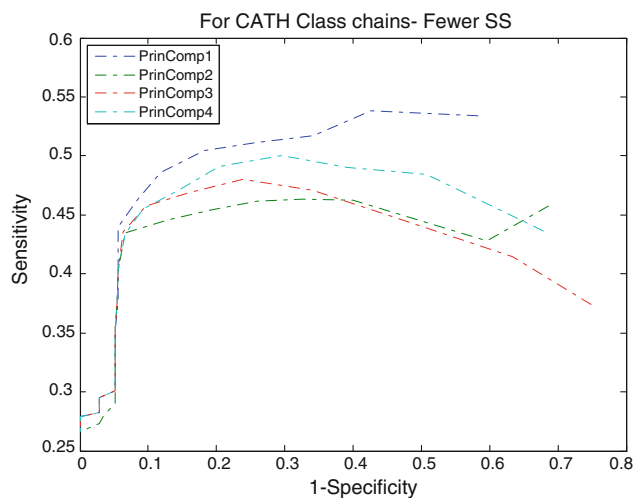


Fig. 7 ROC analysis for fewer secondary structures proteins with respect to threshold

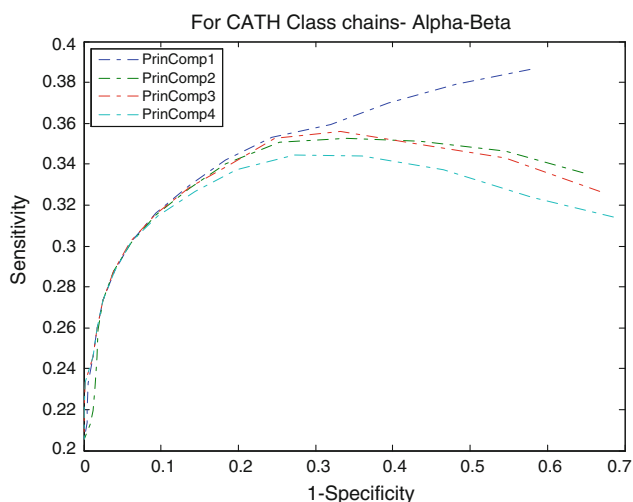


Fig. 6 ROC analysis for alpha-beta proteins with respect to threshold

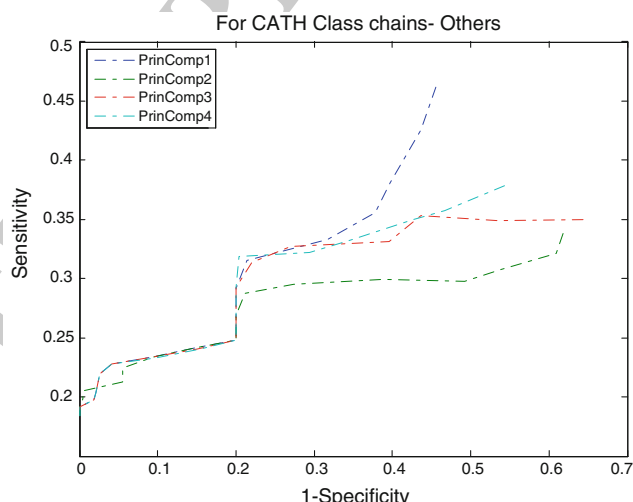


Fig. 8 ROC analysis for other proteins with respect to threshold

412 the decrease of false domain boundary residues leads to the
 413 dramatic increase of false domain residues. In general,
 414 however, the decrease of the specificity (the same as the
 415 increase of the $1 - \text{specificity}$ being shown in the figures)
 416 will lead to the decrease of the sensitivity starting from a
 417 point in ROC curve. The point for mainly alpha proteins is
 418 near specificity 0.55 (i.e., $1 - \text{specificity} = 0.45$), 0.6 for
 419 mainly beta proteins, 0.75 for alpha-beta proteins, and 0.7
 420 for fewer SS proteins.

421 From Fig. 5 we can observe that for the set of mainly
 422 alpha proteins, PrinComp1 provides good predictions
 423 compared to other three profiles. This could be an indication
 424 that inter-residue contact energy is very important. Predic-
 425 tions using the first profile are also important for fewer SS
 426 proteins. Furthermore, predictions from PrinComp4 are

427 important for mainly beta proteins but show poor prediction
 428 for alpha-beta proteins and all alpha proteins. PrinComp2
 429 shows a much lower prediction performance for fewer SS
 430 proteins and other proteins.

431 It has also been observed that the sensitivities of predic-
 432 tions from PrinComp2 are the same as those from
 433 PrinComp3 for mainly alpha, mainly beta, and alpha-beta
 434 proteins in CATH. The specificities of predictions from
 435 PrinComp2 are the same as those from PrinComp1 for
 436 mainly alpha, mainly beta, and alpha-beta proteins in
 437 CATH. More importantly, all the four profiles show good
 438 predictions for mainly beta proteins compared to other
 439 proteins in CATH. The fewer SS proteins also show the
 440 same results although containing fewer numbers of
 441 proteins.

Author Proof

442 Performance with respect to protein classes

443 Tables 2 and 3 show the performance comparisons of the
 444 model on protein chains in our dataset classified by CATH
 445 and SCOP architectures, respectively. In the case of CATH
 446 architecture, protein chains are classified into seven classes
 447 in terms of the composition of secondary structure (SS), i.e.,
 448 all alpha, all beta, alpha/beta, alpha + beta, multi-domain
 449 proteins, membrane and cell surface proteins, and small
 450 proteins. In this work, similar to the above discussion, all
 451 domains of a protein chain belonging to one SCOP class have
 452 the chain to be classified into the class. Inversely, all domains
 453 of a protein chain belonging to different SCOP classes may
 454 make the chain being classified into class "Other".

455 When being classified by SCOP, small protein chains,
 456 although having six members, show the best performance.
 457 The overall sensitivity and the accuracy are around 0.666
 458 and 0.75 from all the four profiles. However, all beta
 459 proteins and alpha + beta proteins have the second best
 460 sensitivities and accuracies. Proteins in other classes have
 461 sensitivity and specificity of 0.413 and 1 from all the four
 462 profiles, respectively. It has also been observed that the
 463 sensitivities of predictions from PrinComp2 tend to be the
 464 same as those from PrinComp3 and PrinComp4 for all
 465 alpha, all beta, alpha/beta, alpha + beta proteins when
 466 being classified by SCOP database.

467 As a result, the PrinComp1 profile shows a good pre-
 468 diction for all proteins compared to the other three profiles.
 469 Moreover, predictions from PrinComp3 are very similar to
 470 those from PrinComp4. The reason behind the similarity of
 471 the predictions between PrinComp3 and PrinComp4 is that
 472 even though the two profiles are correlated to entire chain
 473 composition of amino acids and conformational properties
 474 of individual amino acids, they may also share other
 475 physicochemical properties from the original 544 proper-
 476 ties set in AAindex1 database. In general, using all the four
 477 principal components leads to higher prediction accuracy.

478 Not all protein chains demonstrate similar behavior in
 479 the domain boundary prediction. It is noted that for some
 480 chains such as 1tf3A and 1dx5L, DomSVR predicts a very
 481 few number of false positives and false negatives, which
 482 lead to higher sensitivity and specificity performance. For
 483 protein chains such as 1hf2B, 1cfb0, and 1jr3E, our method
 484 make bad predictions, close to zeros for sensitivities and
 485 specificities with all the four profiles.

486 The important conclusion from these figures and tables
 487 is that PrinComp1, which as stated above is related to inter-
 488 residue contact energy, provides the most reliable predic-
 489 tion. This is due to the fact that in general PrinComp1 has
 490 the largest domain boundaries of predictions compared to
 491 the other three profiles. The average sensitivity of predic-
 492 tions over all protein chains is 0.365 for PrinComp1, 0.356

Table 2 Comparison of protein chains classified by CATH (%)

SS	No.	PrinComp1			PrinComp2			PrinComp3			PrinComp4		
		Sen	Spec	Acc	Sen	Spec	Acc	Sen	Spec	Acc	Sen	Spec	Acc
Mainly alpha	40	32.9	76.2	63.7	31.7	72.5	62.6	31.8	72.6	62.8	32	729	62.8
Mainly beta	95	41.6	80.1	68.1	41.4	80	67.9	41.6	80.6	68.3	41.7	80.8	68.3
Alpha + beta	194	33.2	81.6	65.4	33	81.6	65.2	33	81.1	65.1	32.7	80.3	64.9
Fewer SS	9	47.6	88.1	72.4	44.4	80	69.6	45.9	83.8	70.9	46.1	85.5	71.2
Others	16	30.6	78.6	64.8	28.5	72.7	63.7	30.4	78	64.8	30.9	79.7	65.1

Table 3 Comparison of protein chains classified by SCOP (%)

SS	No.	PrinComp1			PrinComp2			PrinComp3			PrinComp4		
		Sen	Spec	Acc	Sen	Spec	Acc	Sen	Spec	Acc	Sen	Spec	Acc
All alpha	6	34.9	72.3	63.6	33.2	67.6	62.1	33.3	67.8	62.3	33.3	67.8	62
All beta	36	37.5	83.6	67.3	37	82.9	67	37.5	84	67.4	37.9	84.9	67.8
Alpha/beta	80	28.4	84.1	64.2	27.6	82.4	63.7	27.7	82.6	63.5	27.3	81.5	63.1
Alpha + beta	85	34.5	78.3	65.6	34.4	78.9	65.5	34.6	78.6	65.7	34.5	78.1	65.6
Multi-domain	101	29.5	89.7	65.5	28.5	87.2	64.2	27.9	84.4	64.2	29.3	88.8	65.5
Membrane and cell	10	30.5	84.5	65	30.9	85.7	65.3	30.8	85.7	65.2	29.6	82.7	64.1
Small proteins	8	66.6	74.4	75	66.6	73.9	75	66.8	74.8	75.3	66.4	74.5	74.9
Others	28	41.3	100	72.6	41.3	100	72.6	41.3	100	72.6	41.3	100	72.6
Total	354	36.5	80.8	66.3	35.6	80	65.8	35.9	80	66	35.8	80	65.9

493 for PrinComp2, 0.359 for PrinComp3, and 0.358 for Prin-
494 Comp4; the average specificity of predictions for all pro-
495 tein chains is 0.808 for PrinComp1 and 0.8 over all other
496 three profiles.

497 Accuracy for different chains comparison with other
498 methods

499 Our DomSVR method aims to predict domain boundaries
500 for protein chains containing more than one domain.
501 However, it is also suitable for the identification of one-
502 domain protein chain. To make the comparison with other
503 methods, we trained DomSVR predictor on our dataset
504 integrating with other 963 one-domain chains, and then
505 evaluated it both with respect to one-domain chains and
506 multi-domain chains on CAFASP-4 and CASP7 bench-
507 mark datasets. The experiments on one-domain proteins
508 were similar to those on multi-domain proteins. The dataset
509 of one-domain chains is also available at our website:
510 <http://mail.ustc.edu.cn/~bigeagle/DomSVR/index.htm>.

511 The detailed comparison with other similar methods is
512 shown in Table 4 based on the PrinComp1 profile. Table 4
513 shows 13 previous predictors evaluated in the Critical
514 Assessment of Fully Automated Structure Prediction 4
515 (CAFASP-4) (Saini and Fischer 2005), where some stati-
516 stical data are extracted from DOMpro paper (Cheng
517 et al. 2006). The evaluation dataset of CAFASP-4 consists
518 of 41 one-domain CASP6 targets and 17 two-domain
519 CASP6 targets (58 targets in total). The targets in CA-
520 FASP-4 dataset are divided into two main divisions:

521 homology modeling and fold recognition targets. Twenty
522 one-domain chains and 7 two-domains chains are homol-
523 ogy modeling targets, and 21 one-domain chains and 10
524 two-domain chains are fold recognition targets. In the
525 CAFASP-4, seven predictors belong to the category of
526 template-based methods, which have an advantage due to
527 this evaluation set contains only comparative modeling and
528 fold recognition targets (no new fold targets). Our method
529 achieves higher sensitivity and specificity than other
530 ab initio predictors when averaging over all of the targets.
531 Moreover, in spite of our model outperforms even better
532 than some template-based methods such as ADDA, Inter-
533 ProScan, and Dompred-Domssea, it performs worse than
534 other template-based methods such as Dopro, SSEPD-
535 domain, and Robetta-Ginzu.

536 Table 5 shows the performance comparison of the 14
537 domain boundary predictors, random predictor, and our
538 DomSVR predictor with PrinComp1 profile on the selected
539 CASP7 dataset. Currently, the dataset contains 95 peptide
540 chains where some chains were removed by assessors of
541 CASP7. It consists of 62 one-domain chains, 30 two-
542 domain chains, 2 three-domain chains and 1 four-domain
543 chain. In this work, we made comparison of our method
544 and 14 predictors in the CASP7 assessment by evaluated
545 on one-domain chains, two-domain chains, and even chains
546 containing more than two domains. All the prediction
547 data for the 14 predictors are created from CASP7
548 <http://www.predictioncenter.org/casp7/>. In Table 5, the
549 accuracy is calculated as the ratio of the number of chains
550 with correctly predicted domain number to that of chains

Table 4 Performance comparison with other methods on CAFASP-4 benchmark dataset

Predictor	1-D ^a		2-D		All-D	
	Sen	Spec	Sen	Spec	Sen	Spec
DomSVR ^b	0.8	0.9	0.34	0.78	0.67	0.87
ADDA (Heger and Holm 2003) ^b	0.85	0.73	0.18	0.33	0.66	0.67
Armadillo ^b	0.1	1	0.24	0.18	0.14	0.31
Biozon (Nagarajan and Yona 2004) ^b	0.1	1	0.35	0.19	0.17	0.29
Dompred-DPS (Bryson et al. 2005) ^b	0.68	0.78	0.47	0.5	0.62	0.69
DOMpro ^b	0.85	0.76	0.35	0.5	0.71	0.71
Globplot (Linding et al. 2003) ^b	0.83	0.71	0.18	0.6	0.64	0.7
Mateo (Lexa and Valle 2003) ^b	0.51	0.78	0.12	0.15	0.4	0.58
Dompred-Domssea (Marsden et al. 2002)	0.8	0.75	0.29	0.63	0.66	0.73
Dopro (von Ohlsen et al. 2004)	0.85	0.88	0.53	0.64	0.76	0.81
InterProScan (Zdobnov et al. 2001)	0.93	0.75	0.24	0.67	0.72	0.74
Robetta-Ginzu (Chivian et al. 2003)	0.8	0.92	0.53	0.69	0.72	0.86
Robetta-Rosettadom	0.83	0.94	0.71	0.75	0.79	0.88
SSEPDdomain (Gewehr et al. 2005)	0.93	0.84	0.47	0.73	0.79	0.82

^a 1-D denotes that each tested protein chain is a 1-domain one, 2-D denotes that each tested protein chain contains more than one domain, while All-D stands for all tested protein chains

^b Ab initio method

Table 5 Performance comparison with other methods on CASP7 benchmark dataset (%)

Predictor	1-D	2-D	3-D ^a	All-D
DomSVR ^b	82.26 (51/62)	46.67 (14/30)	33.33 (1/3) ^c	69.47 (66/95)
chop ^b	53.66 (22/41)	28.57 (6/21)	0 (0/3)	43.08 (28/65)
chop_homo ^b	58.33 (21/36)	36.36 (8/22)	0 (0/3)	47.54 (29/61)
DomFOLD ^b	97.96 (48/49)	20.69 (6/29)	0 (0/3)	66.67 (54/81)
DPS ^b	78.95 (30/38)	42.31 (11/26)	0 (0/3)	61.19 (41/67)
Distill ^b	77.42 (48/62)	46.67 (14/30)	33.33 (1/3)	66.32 (63/95)
NN_PUT_lab	77.59 (45/58)	10.34 (3/29)	33.33 (1/3)	54.44 (49/90)
BAKER-ROSETTADOM	88.52 (54/61)	80 (24/30)	0 (0/3)	82.98 (78/94)
DomSSEA	97.44 (38/39)	30.77 (8/26)	33.33 (1/3)	69.12 (47/68)
FOLDpro	98.36 (60/61)	76.67 (23/30)	33.33 (1/3)	89.36 (84/94)
HHpred1	96 (48/50)	14.29 (4/28)	33.33 (1/3)	65.43 (53/81)
HHpred3	94.12 (48/51)	17.24 (5/29)	33.33 (1/3)	65.06 (54/83)
Ma-OPUS-DOM	87.8 (36/41)	76.92 (20/26)	33.33 (1/3)	81.43 (57/70)
Robetta-Ginzu	83.61 (51/61)	86.67 (26/30)	33.33 (1/3)	82.98 (78/94)
Meta-DP	97.56 (40/41)	14.81 (4/27)	0 (0/3)	61.97 (44/71)
Random predictor	65.21 (40.43/62)	31.51 (9.45/30)	3.17 (0.0951/3)	52.61 (49.98/95)

^a “1-D”, “2-D”, and “3-D” denote that each tested protein chain is a 1-domain one, 2-domain one, and chain with three or more domains, respectively. In addition “All-D” stands for all tested protein chains

^b Ab initio method

^c The numbers in parentheses denote correctly predicted chains and the amount of chains used to the prediction

551 for one-domain, two-domain, three-domain, or all-domain
 552 category. In this case, template-based predictors outper-
 553 form ab initio-based predictors due to the advantage of
 554 containing similar fold targets in their template set. Sta-
 555 tistically, our method performs better than other ab initio-
 556 based predictors and even better than some template-based
 557 predictors, such as HHpred1, HHpred2, and DomSSEA. In
 558 addition, our method also makes better prediction than a
 559 meta predictor, Meta-DP, which integrated several pre-
 560 dictors in order to obtain better predictions than the use of
 561 single predictor (Saini and Fischer 2005).

562 One important aspect should be noted that split-domain
 563 in chain involved in CAFASP-4 and CASP7 datasets is
 564 treated as one single domain due to the complex domain
 565 topology. For the CAFASP-4 database, there are five such
 566 targets, T0226, T0248, T0268, T0279, and T0280. In the
 567 case of target T0226, predictors Robetta-Rosettadom,
 568 Biozon, and DOMpro make correct predictions of domain
 569 number but predict the domain boundary between the first
 570 split of the split-domain and another domain as non-
 571 boundary. Our method makes a similar prediction as
 572 DOMpro predictor. Other predictors in CADASP-4 make
 573 wrong predictions of domain number for target T0226. For
 574 other four targets, all predictors perform similar. For the
 575 CASP7 dataset, there are 18 such targets containing 17
 576 two-domain chains and 1 three-domain chains. Some
 577 methods in CASP7 identify split-domain as two or more
 578 domains and some other ones correctly predict one split of

579 the domain. Table 4 demonstrates prediction performance
 580 excluding the targets having split-domain on CAFASP-4
 581 dataset, while Table 5 shows prediction performance
 582 involving in 18 split-domain targets on CASP7 dataset. We
 583 evaluate the predictors on the condition that split-domain in
 584 one chain is treated as one domain. Performance of each
 585 method is varied with and without involving these split-
 586 domain targets, and the comparison excluding such targets
 587 is shown in Table 6. Note that no method can make correct
 588 predictions for three-domain chains and, additionally, in
 589 Tables 5 and 6 all predictions for the 1 four-domain chain
 590 are not correct.

591 However, predictions may be changed if the evaluation
 592 is with respect to both domain boundary and domain
 593 number, but not with respect to domain number alone.
 594 Suppose that a chain is correctly predicted if its domain
 595 number was predicted correctly and the predicted domain
 596 boundaries distance from the true boundaries less than ± 20
 597 residues in primary sequence. In this case, accuracies of
 598 our method are 82.26, 40, 33.33, and 67.37% for one-
 599 domain, two-domain, three-domain, and all-domain cate-
 600 gories, respectively, which are a little less than the case of
 601 those with respect to domain number alone. In detail, the
 602 predictions of domain boundaries for targets T0330 and
 603 T0379 are wrong although the predictions of domain
 604 number were correct by our model. Target T0330 consists
 605 of two domains: one domain is split into two so-called
 606 split-domains containing residues from SER2 to LYS16

Table 6 Performance comparison with other methods on CASP7 benchmark dataset excluding chains having split-domain (%)

Predictor	1-D	2-D	3-D ^a	All-D
DomSVR ^b	82.26 (51/62)	53.85 (7/13)	0 (0/2) ^c	75.32 (57/77)
chop ^b	53.66 (22/41)	22.22 (2/9)	0 (0/2)	46.15 (24/52)
chop_homo ^b	58.33 (21/36)	33.33 (3/9)	0 (0/2)	51.06 (24/47)
DomFOLD ^b	97.96 (48/49)	25 (3/12)	0 (0/2)	80.96 (51/63)
DPS ^b	78.95 (30/38)	60 (6/10)	0 (0/2)	72 (36/50)
Distill ^b	77.42 (48/62)	46.15 (6/13)	0 (0/2)	70.13 (54/77)
NN_PUT_lab	77.59 (45/58)	16.67 (2/12)	0 (0/2)	65.28 (47/72)
BAKER-ROSETTADOM	88.52 (54/61)	53.85 (7/13)	0 (0/2)	80.26 (61/76)
DomSSEA	97.44 (38/39)	40 (4/10)	0 (0/2)	82.35 (42/51)
FOLDpro	98.36 (60/61)	69.23 (9/13)	0 (0/2)	90.79 (69/76)
HHpred1	96 (48/50)	9.09 (1/11)	0 (0/2)	77.78 (49/63)
HHpred3	94.12 (48/51)	16.67 (2/12)	0 (0/2)	76.92 (50/65)
Ma-OPUS-DOM	87.8 (36/41)	60 (6/10)	0 (0/2)	79.25 (42/53)
Robetta-Ginzu	83.61 (51/61)	69.23 (9/13)	0 (0/2)	78.95 (60/76)
Meta-DP	97.56 (40/41)	30 (3/10)	0 (0/2)	81.13 (43/53)
Random predictor	80.54 (49.92/62)	16.98 (2.21/13)	1.25 (0.025/2)	67.75 (52.17/95)

^a “1-D”, “2-D”, “3-D”, and “All-D” are the same as in Table 5

^b Ab initio method

^c The numbers in parentheses denote correctly predicted chains and the amount of chains used to the prediction

607 and from THR92 to THR229, while the other one is located
608 from VA117 to ILE91. As a result, the predicted domain
609 boundary is located from residue LEU115 to residue
610 ILE154. Actually, some residues of the target were missed
611 in the structure-determined experiments, and the target
612 structure also contains several “non-standard” groups. All
613 of these make the prediction of domain boundary hard. In
614 the case of target T0339, it also consists of two domains:
615 one domain is split into two split-domains containing res-
616 idues from MSE1 to LEU16 and from LEU84 to GLN207,
617 while the other one is located from ASN17 to PHE83.
618 Containing “non-standard” groups and missed residues
619 makes the same effect on the prediction of domain
620 boundary as the Target T0330.

621 To make sure the prediction is accurate, a random pre-
622 dictor was constructed and the prediction performance
623 based on CASP7 dataset is appended to the last row of
624 Tables 5 and 6. In the case of evaluation on CASP7, the
625 random predictor was constructed in the same form of
626 CASP7 dataset which consists of 62 one-domain chains, 30
627 two-domain chains, and three chains having three or more
628 domains. To better simulate the real random sampling test,
629 we ran the random predictor 10,000 times and one average
630 accuracy of 52.61% was achieved. From the Table 5, most
631 of methods outperform the random one except for predic-
632 tors “chop” and “chop_homo”. In the case of evaluation
633 on CASP7 without chains having split-domain, random
634 predictor was created and ran in the same way. The dataset

635 consists of 62 one-domain chains, 13 two-domain chains,
636 and two chains with three or more domains. The last row of
637 Table 6 can be seen on average accuracy of 67.75% for
638 random predictor. From Table 6, predictors “chop”,
639 “chop_homo”, and NN_PUT_lab perform worse than
640 random predictor.

641 Moreover, we assess both template-based and ab initio
642 predictors on the CASP7 dataset, respectively. Figure 9,
643 respectively, illustrates domain number comparison of such
644 two categories of predictors, our model, and random pre-
645 dictor, with and without split-domain chains. The overall
646 accuracies of domain number prediction for the template-
647 based and ab initio predictors are 72.53 and 56.96%,
648 respectively; while the accuracies are respectively 79.19
649 and 64.06% if excluding split-domain chains.

650 As discussed above, it can be found that our SVR
651 model outperforms other predictors despite of obtaining a
652 lower accuracy for three-domain chains, probably due to
653 the small number of three-domain chains in CASP7
654 dataset. Actually, more one-domain chains and less chains
655 with two or more domains may make the prediction over-
656 estimated. In addition, the small number of chains in
657 CAFASP-4 and CASP7 datasets may also aggravate the
658 trend. Therefore, the evaluation based on a small size of
659 dataset cannot fully reflect the advantages and disadvan-
660 taged of these methods. As a result, larger benchmark
661 dataset is more desirable to compare these similar meth-
662 ods in the future.

Fig. 9 Performance comparison based on CASP7 dataset. No yellow bar is shown in the right graph for template-based, ab initio, and DomSVR predictors, since the prediction accuracies for three-domain chains are zeros

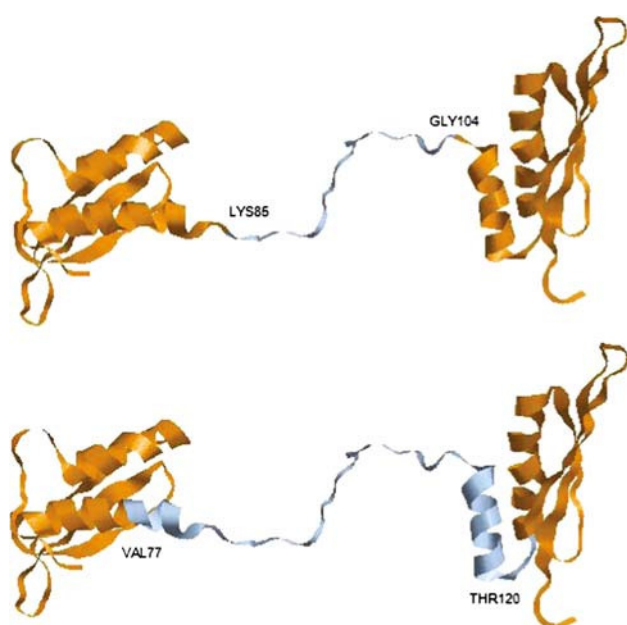
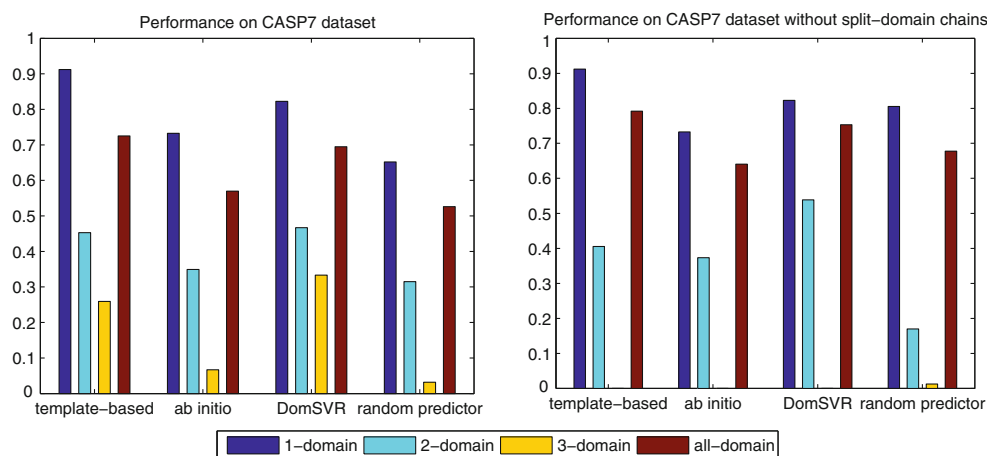


Fig. 10 Comparison of natural versus predicted domain boundaries for protein chain Iqu6_A. The chain is colored in gold and the domain boundary (true or predicted) are colored in blue. **a** True domain boundary for protein chain Iqu6A, **b** Predicted domain boundary for protein chain Iqu6A

663 A case study of domain boundary prediction

664 In order to illustrate the prediction of domain boundaries
 665 directly, protein chain Iqu6A (the same protein discussed
 666 as Fig. 3) is taken as a case of domain boundary prediction
 667 and shown in Fig. 10. The protein chain has 179 residues
 668 and consists of two double-stranded RNA (dsRNA)-binding
 669 domains linked by a domain boundary ranging from
 670 residue LYS85 to GLY104 (shown in Fig. 10). The protein
 671 Iqu6, categorized as kinase PKR (protein kinase RNA-
 672 regulated), is an interferon-induced enzyme that plays a
 673 key role in the control of viral infections and cellular

homeostasis (Nanduri et al. 1998). Protein kinase PKR is
 674 activated by a distinct mechanism that involves dsRNA
 675 binding in its N-terminal region in an RNA sequence-
 676 independent fashion. The structure of dsRNA-binding
 677 domain exhibits a dumb-bell shape comprising two tandem
 678 linked dsRNA-binding motifs both with an alpha-beta-
 679 beta-beta-alpha fold. The structure may reveal a highly
 680 conserved RNA-binding site on each dsRNA-binding motif
 681 and suggests a novel mode of protein-RNA recognition.
 682 The central linker between the two dsRNA-binding motifs
 683 is highly flexible, which may enable the two motifs to wrap
 684 around the RNA duplex for cooperative and high-affinity
 685 binding and advance the overall change of PKR conforma-
 686 tion and its activation (Nanduri et al. 1998). The domain
 687 boundary prediction for this protein chain is demonstrated
 688 in Fig. 10. In this case, our approach predicted the domain
 689 boundary actually but a little extension to several residues,
 690 ranging from residue VAL77 to residue THR120. 691

692 Conclusions

693 In this paper, we addressed the problem of domain
 694 boundaries prediction from sequence information alone.
 695 Amino acid residue profiles were taken from AAindex
 696 database using PCA technique to extract necessary physico-
 697 chemical properties. The profiles were then used to train
 698 and test our predictor by the form of input vectors. As a
 699 result, our method achieves a sensitivity of 36.5% and a
 700 specificity of 80.8%. Our method is also evaluated on two
 701 datasets: the CAFASP-4 dataset and the CASP7 benchmark
 702 dataset. On the CAFASP-4 test dataset, our method per-
 703 forms better than the template-based method InterProScan
 704 and comparably to all other template-based methods with
 705 respect to specificities. Moreover, our method performs
 706 significantly better than all other ab initio methods for
 707 domain boundary prediction. On the CASP7 test dataset,

708 our method is able to outperform all the other ab initio
709 methods for two-domain protein chains and slightly worse
710 than some other methods for one-domain protein chains.
711 However, the overall accuracy of our model is the best. It
712 should be noted that the purpose of the comparison is just
713 to estimate the current state-of-the-art of domain boundary
714 prediction instead of ranking these methods, because pre-
715 dictors used different scales of protein set from the CA-
716 FASP-4 and CASP7 datasets to evaluate themselves.

717 In general, we are not only interested in the overall
718 performance of domain boundary prediction, but also
719 interested in how the prediction accuracy varies across
720 different protein classes by CATH and SCOP architectures.
721 Three hundred and fifty-four protein chains representing all
722 major classes from CATH and SCOP have been chosen for
723 training and testing our method. Mainly beta proteins and
724 fewer SS proteins achieve better prediction compared to
725 other proteins when classifying by CATH. When being
726 classified by SCOP, small proteins show the best sensitiv-
727 ities although containing six protein chains. However, all
728 beta proteins and alpha + beta proteins achieve the second
729 best sensitivities and accuracies. PrinComp1, having strong
730 correlation to inter-residue contact energy property, is the
731 one that the predictor achieves the most reliable results
732 from. The model also achieves very accurate predictions
733 from PrinComp2, PrinComp3, and PrinComp4, but the
734 number of correctly predicted domain boundary residues
735 from them is smaller than the model gets from PrinComp1.

736 The DomSVR algorithm described in this work gives
737 good results for most of proteins in our dataset taken from
738 PDB database. The successful application of SVR approach
739 in this study suggests that SVR can accurately describe the
740 relationship between primary sequence and domain
741 boundaries using amino acid information alone. The pre-
742 dicted domain boundaries can be used for classification of
743 proteins and understanding the evolutions, structures and
744 functions of proteins, which motivate us to improve the
745 algorithm and apply it to other protein chains. In future
746 work, we expect that the improved version of our predictor
747 can test more protein chains and reevaluate the chains that
748 have already been tested with our current predictor.

749 **Acknowledgments** This work was supported in part by grant 2 G12
750 RR003048 from the RCMI program, Division of Research Infra-
751 structure, National Center for Research Resources, NIH and the
752 Mordecai Wyatt Johnson program of Howard University. This work
753 was also supported in part by the Singapore MOE ARC Tier-2 funding
754 grant T208B2203 and the National Science Foundation of China
755 (No. 60803107). CL's work was supported by NSF (CCF-0845888).

756 References

757 Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W,
758 Lipman DJ (1997) Gapped BLAST and PSI-BLAST: a new

- generation of protein database search programs. *Nucleic Acids Res* 25:3389–3402 759
- Baldi P, Brunak S, Chauvin Y, Andersen CA, Nielsen H (2000) 760
Assessing the accuracy of prediction algorithms for classifica- 761
tion: an overview. *Bioinformatics* 16:412–424 762
- Chen P, Wang B, Wong HS, Huang D.S. (2007) Prediction of protein 763
B-factors using multi-class bounded SVM. *Protein Pept Lett* 764
14(2):185–190 765
- Cheng J, Sweredoski MJ, Baldi P (2006) DOMpro: protein domain 766
prediction using profiles, secondary structure, relative solvent 767
accessibility, and recursive neural networks. *Data Min Knowl* 768
Discov 13:1–10 769
- Copley RR, Doerksa T, Letunica I, Borcka P (2002) Protein domain 770
analysis in the era of complete genomes. *FEBS Lett* 513:129– 771
134 772
- Dovidchenko NV, Lobanov MY, Galzitskaya OV (2007) Prediction 773
of number and position of domain boundaries in multi-domain 774
proteins by use of amino acid sequence alone. *Curr Protein Pept* 775
Sci 8(2):189–195 776
- Drucker H, Burges CJC, Kaufman L, Smola AJ, Vapnik V (1996) 777
Support vector regression machines. In: *Proceedings of the* 778
NIPS, pp 155–161 779
- Dumontier M, Feldman R, Yao HJ, Hogue CWV (2005) Armadillo: 780
domain boundary prediction by amino acid composition. *J Mol* 781
Biol 350:1061–1073 782
- Edelman GM (1973) Antibody structure and molecular immunology. 783
Science 180:830–840 784
- Fukuchi S, Nishikawa K (2001) Protein surface amino acid compo- 785
sitions distinctively differ between thermophilic and mesophilic 786
bacteria. *J Mol Biol* 309:835–843 787
- Galzitskaya OV, Melnik BS (2003) Prediction of protein domain 788
boundaries from sequence alone. *Protein Sci* 12:696–701 789
- George RA, Heringa J (2002) Protein domain identification and 790
improved sequence similarity searching using PSI-BLAST. 791
Proteins: Struct Funct Gen 48:672–681 792
- George RA, Heringa J (2002) SNAPDRAGON: a new method to 793
predict protein structural domain boundaries from sequence data. 794
J Mol Biol 316:839–851 795
- Gewehr JE, Zimmer R (2005) SSEP-Domain: protein domain 796
prediction by alignment of secondary structure elements and 797
profiles. *Bioinformatics* 22:181–187 798
- Goodall C (1990) *Modern methods of data analysis*. Sage Publica- 799
tions, Newbury Park, CA 800
- Gunn SR (1998) *Support vector machines for classification and* 801
regression. Faculty of Engineering and Applied Science, Uni- 802
versity of Southampton 803
- Jolliffe IT (2002) *Principal component analysis*. Springer, NY. 804
- Kawashima S, Pokarowski P, Pokarowska M, Kolinski A, Katayama 805
T, Kanehisa M (2008) AAindex: amino acid index database, 806
progress report. *Nucleic Acids Res* 36:D202–D205 807
- Levitt M, Chothia C (1976) Structural patterns in globular proteins. 808
Nature 261:552–558 809
- Liu J, Rost B (2004) Sequence-based prediction of protein domains. 810
Nucleic Acids Res 32:3522–3530 811
- Marchler-Bauer A, Anderson JB, Derbyshire MK, DeWeese-Scott C 812
(2007) CDD: a conserved domain database for interactive 813
domain family analysis. *Nucleic Acids Res* 35:D237–240 814
- Marsden RL, McGuffin LJ, Jones DT (2002) Rapid protein domain 815
assignment from amino acid sequence using predicted secondary 816
structure. *Protein Sci* 11:2814–2824 817
- Miyazawa S, Jernigan RL (1999) Self-consistent estimation of inter- 818
residue protein contact energies based on an equilibrium mixture 819
approximation of residues. *Proteins* 34:49–68 820
- Munoz V, Serrano L (1994) Intrinsic secondary structure propensities 821
of the amino acids, using statistical phi–psi matrices: comparison 822
with experimental scale. *Proteins* 20:301–311 823
824

- 825 Nagarajan N, Yona G (2004) Automatic prediction of protein
826 domains from sequence information using a hybrid learning
827 system. *Bioinformatics* 20:1335–1360
- 828 Nanduri S, Carpick BW, Yang Y, Williams BR, Qin J (1998)
829 Structure of the double-stranded RNA-binding domain of the
830 protein kinase PKR reveals the molecular basis of its dsRNA-
831 mediated activation. *EMBO J* 17:5458–5465
- 832 Orengo CA, Michie AD, Jones DT, Swindells MB, Thornton JM
833 (1997) CATH: a hierarchic classification of protein domain
834 structures. *Structure* 5:1093–1108
- 835 Porter RR (1973) Structural studies of immunoglobulins. *Science*
836 180:713–716
- 837 Rackovsky S, Scheraga HA (1982) Differential geometry and
838 polymer conformation. 4. Conformational and nucleation prop-
839 erties of individual amino acids. *Macromolecules* 15:1340–1346
- 840 Saini HK, Fischer D (2005) Meta-DP: domain prediction meta server.
841 *Bioinformatics* 21:2917–2920
- 842 Sikder AR, Zomaya AY (2006) Improving the performance of
843 DomainDiscovery of protein domain boundary assignment using
844 inter-domain linker index. *BMC Bioinform* 7:S6
- Sim J, Kim SY, Lee J (2005) PRODO: prediction of protein domain
boundaries using neural networks. *Proteins* 59:627–632
- Suyama M, Ohara O (2003) DomCut: prediction of inter-domain
linker regions in amino acid sequences. *Bioinformatics* 19:673–
674
- Wetlaufer DB (1973) Nucleation, rapid folding, and globular
intrachain regions in proteins. *Proc Natl Acad Sci USA*
70:697–701
- Ye L, Liu T, Wu Z, Zhou R (2007) Sequence-based protein domain
boundary prediction using BP neural network with various
property profiles. *Proteins: Struct Funct Bioinform* 71:300–307
- Yoo PD, Sikder AR, Zhou BB, Zomaya AY (2008) Improved general
regression network for protein domain boundary prediction.
BMC Bioinform 9:S12
- Zhou Y, Vitkup D, Karplus M (1999) Native proteins are surface-
molten solids: application of the Lindemann criterion for the
solid versus liquid state. *J Mol Biol* 285:1371–1375

UNCORRECTED PROOF