1    PhD7Faster 2.0: predicting clones propagating faster from the Ph.D.-7 phage display library

2    by coupling PseAAC and tripeptide composition

3

4    Bifang He[1,2], Heng Chen[1], Jian Huang[1]

5

6    [1] School of Medicine, Guizhou University, Guiyang 550025, China

7    [2] Center for Informational Biology, University of Electronic Science and Technology of

8    China, Chengdu 611731, China

9

10    Corresponding Author:

11    Jian Huang

12    No. 2006, Xiyuan Ave, West Hi-Tech Zone, University of Electronic Science and Technology

13    of China, Chengdu 611731, China

14    Email address: hj@uestc.edu.cn

15

16

**ABSTRACT**

Selection from phage display libraries empowers isolation of high-affinity ligands for various

targets. However, this method also identifies propagation-related target-unrelated peptides

(PrTUPs). These false positive hits appear because of their amplification advantages. In this

report, we present PhD7Faster 2.0 for predicting fast-propagating clones from the Ph.D.-7

phage display library, which was developed based on support vector machine (SVM). Feature

selection was performed against PseAAC and tripeptide composition using the incremental

feature selection method. Ten-fold cross-validation results show that PhD7Faster 2.0 succeeds

a decent performance with the accuracy of 81.84%, the Matthews correlation coefficient

(MCC) of 0.64 and the area under the ROC curve (AUC) of 0.90. The permutation test with

1000 shuffles resulted in $p <0.001$. We implemented PhD7Faster 2.0 into a publicly accessible

web tool (http://i.uestc.edu.cn/sarotup3/cgi-bin/PhD7Faster.pl) and constructed standalone

graphical user interface (GUI) and command-line versions for different systems. The

standalone PhD7Faster 2.0 is able to detect PrTUPs within small datasets as well as large-

scale datasets. This makes PhD7Faster 2.0 an enhanced and powerful tool for scanning and

reporting faster-growing clones from the Ph.D.-7 phage display library.

## INTRODUCTION

Phage display is a high throughput and powerful screening methodology for identifying

ligands for myriad target types, ranging from molecules (microRNA, protein, polysaccharide)

(He et al., 2013; Zhang et al., 2017) to inorganic (gold) (Causa et al., 2013), organic (epoxy)

(Swaminathan et al., 2013) and biological (tissue, organ) materials (Hung et al., 2018). Large

libraries of phage-displayed peptides or proteins consist of millions to billions of variant

members, which can be iteratively selected and amplified in a process referred to as

biopanning (Pande et al., 2010). Recently, next generation sequencing technologies have been

coupled with phage display, which have substantially contributed to the analysis of output

from combinatorial libraries and allowed for even faster and more robust discovery of novel

ligands (Christiansen et al., 2015; Matochko et al., 2014; Ngubane et al., 2013; Rentero

Rebollo et al., 2014; t Hoen et al., 2012). The ever-increasing utility and versatility makes

phage display a powerful tool in multiple research areas, such as materials science,

biotechnology, pharmacology, cell biology and diagnostics (Martins et al., 2016).

However, the phage display methodology is notorious for the enrichment of target-

unrelated peptides (TUPs) (Menendez et al., 2005). Therefore, biopanning results are a

mixture of true target binders and TUPs (Vodnik et al., 2011). These false positive TUPs have

no actual affinity towards the target of interest and can fall into two categories: selection- and

propagation-related TUPs (SrTUPs and PrTUPs) (Thomas et al., 2010). The SrTUPs can bind

to other components (plates, beads) of the screening system other than the desired target and

thus creep into the output of phage display. The PrTUPs sneak into the biopanning results due

to their propagation advantages, which allow them to outcompete clones with lower growth

55    rates (Brammer et al., 2008; Matochko et al., 2014; Nguyen et al., 2014; Thomas et al., 2010;

56    Zade et al., 2017; Zygiel et al., 2017). Apparently, these TUPs may misdirect ligand discovery

57    through biopanning and should be distinguished from actual target-binding peptides

58    (Bakhshinejad et al., 2016). Therefore, the diagnosis of TUPs is as crucial as the identification

59    of target binders.

60          Although several experimental strategies have been proposed to decrease TUP

61    isolation during biopanning and differentiate between TUPs and true binders post-biopanning

62    (Nguyen et al., 2014; Thomas et al., 2010; Vodnik et al., 2011), TUP analysis has benefitted

63    considerably from computational approaches. Databases (BDB (He et al., 2016a; He et al.,

64    2018; Huang et al., 2012; Ru et al., 2010), PepBank (Shtatland et al., 2007)) and

65    bioinformatics tools (He et al., 2016b; Huang et al., 2010; Li et al., 2017; Mandava et al.,

66    2004; Ru et al., 2014) have been widely employed to report both SrTUPs and PrTUPs.

67    Searching against databases for biopanning data can uncover whether query peptides have

68    been isolated by many different targets. If so, query sequences are potential SrTUPs and

69    PrTUPs due to lack of target specificity. For example, the peptide HAIYPRH (a typical

70    PrTUP) has been identified by 23 completely different targets according to results of

71    searching the BDB database. The phage displayed the peptide was later verified to have a

72    propagation advantage owing to mutations in the regulatory region of the phage genome

73    (Brammer et al., 2008). HWGMWSY (a SrTUP) has been isolated by 10 completely different

74    targets according to records in the BDB database. The peptide was proved to be a plastic

75    binder (Vodnik et al., 2012), which resulted in this peptide repeatedly appearing in multiple

76    reported screening experiments. SABinder (He et al., 2016b) and PSBinder (Li et al., 2017)

77 have been designed for predicting streptavidin- and polystyrene surface-binding peptides,

78 respectively, as they are commonly known SrTUPs. The INFO tool in the RELIC suite

79 enables PrTUPs detection based on information content (Mandava et al., 2004), whereas

80 PhD7Faster (PhD7Faster 1.0) based on support vector machine (SVM) allows the prediction

81 of clones with amplification advantages from the popular commercial Ph.D.-7 phage display

82 library (Ru et al., 2014). However, PhD7Faster 1.0 can be improved in the following three

83 aspects. Firstly, the positive training dataset of PhD7Faster 1.0 was selected based on the copy

84 number of a peptide (15 or higher) after one round of amplification without consideration of

85 the corresponding copy number in the naïve Ph.D.-7 library. Secondly, only dipeptide

86 composition was employed to develop the classifier. Currently many reports have

87 demonstrated that predictors developed by combining pseudo amino acid composition

88 (PseAAC) (Chou, 2001; Chou, 2005) and tripeptide composition can achieve decent

89 predictive performances (Liao et al., 2011; Zhu et al., 2015). Thirdly, PhD7Faster 1.0 is

90 unable to process large datasets (e. g., next-generation sequencing data).

91      In this study, we develop a new predictor for identifying clones propagating faster from

92 the Ph.D.-7 phage display library. The SVM algorithm was employed to model the predictor

93 with the optimal feature subset after feature selection. The constructed SVM-based classifier

94 obtained an accuracy of 81.84% in the ten-fold cross-validation. The predictor was further

95 implemented into a web tool, called PhD7Faster 2.0, which is freely available at

96 http://i.uestc.edu.cn/sarotup3/cgi-bin/PhD7Faster.pl. We also developed the standalone

97 version of PhD7Faster 2.0 that enables the analysis of PrTUPs within large-scale datasets.

98 **DATA & METHODS**

## Benchmark Datasets

The dataset used to develop the predictor was acquired from (Matochko et al., 2014). Derda and coworkers employed high-throughput sequencing technology to characterize both the naïve Ph.D.-7 phage display library and the same library after one round of amplification. By comparing the abundance of each peptide before and after amplification using Bioconductor package edgeR, 770 unique peptides were identified with significantly higher growth rate (parasitic sequences) (Matochko et al., 2014), which were collected into the positive training dataset. The negative dataset was composed of those peptides with the copy number of one in the amplified Ph.D.-7 phage display library. The datasets were then processed as follows: (i) peptide sequences containing ambiguous residues (such as "X", "B" and "Z") were excluded; (ii) sequences within 2 Hamming distance (h=2, the Hamming distance between two strings of equal length is the minimum number of substitutions required to change one string into the other.) were removed. Finally, 749 peptides were retained in the positive dataset. To match the size of the positive dataset, we randomly selected 749 peptides from the negative dataset. No overlapping was found between the negative and positive datasets. Finally, the benchmark dataset was composed of 749 fast-growing peptides and 749 regular-growing peptides (See positive.fasta and negative.fasta in Supplementary Data).

## PseAAC and tripeptide composition

Extracting a set of informative features is a standard and important procedure for developing predictors. Chou initially formulated the PseAAC (Chou, 2001; Chou, 2005), which consists of more than 20 discrete numbers, where the top 20 represent the classical amino acid composition (AAC) of a protein sequence whereas the additional parameters incorporate

121 some sequence-order information. PseAAC and tripeptide composition have been widely used

122 in protein prediction related research (Chou, 2011; Lin et al., 2013). Here, they were

123 employed to encode each peptide in the benchmark dataset.

124 Given a peptide P with L amino acid residues:

$$P = [R_1 R_2 R_3 R_4 R_5 R_6 R_7 \dots R_L] \tag{1}$$

126 where $R_i$ ($i = 1, 2, 3 \dots$ L) is the residue at the $i$-th sequence position. Accordingly, any

127 sequence like the peptide P of Equation (1) can be presented using a set of feature vectors

128 with $8000 + n\lambda$ dimensions.

$$P = [P_1, P_2, \dots, P_{8000}, P_{8000+1}, \dots, P_{8000+n\lambda}] \tag{2}$$

130 where the first 8000 numbers $P_1, P_2, \dots, P_{8000}$ reflect the effect of the conventional tripeptide

131 composition; the remaining $n\lambda$ elements $P_{8000+1}, P_{8000+2}, \dots, P_{8000+n\lambda}$ reflect the amphipathic

132 sequence-order pattern. These features are calculated through the following equations:

$$P_u = \begin{cases} \dfrac{f_u}{\sum_{i=1}^{8000} f_i + w \sum_{j=1}^{n\lambda} \tau_j} & (1 \le u \le 8000) \\[4mm] \dfrac{w\tau_u}{\sum_{i=1}^{8000} f_i + w \sum_{j=1}^{n\lambda} \tau_j} & (8000 + 1 \le u \le 8000 + n\lambda) \end{cases} \tag{3}$$

134 where $f_i$ ($i = 1, 2, 3, \dots, 8000$) are the normalized occurrence frequencies of the 8000

135 tripeptides in peptide P; w is the weight factor for the sequence-order effect; $\tau_j$ ($j = 1, 2, \dots,$

136 $n\lambda$) is the $j$-tier sequence-correlated factor as formulated by:

137

$$
\begin{cases}
\tau_1 = \dfrac{1}{L-1}\displaystyle\sum_{i=1}^{L-1} H^1_{i,i+1} \\[2ex]
\tau_2 = \dfrac{1}{L-1}\displaystyle\sum_{i=1}^{L-1} H^2_{i,i+1} \\[1ex]
\quad \cdots \\[1ex]
\tau_n = \dfrac{1}{L-1}\displaystyle\sum_{i=1}^{L-1} H^n_{i,i+1} \\[2ex]
\tau_{n+1} = \dfrac{1}{L-2}\displaystyle\sum_{i=1}^{L-2} H^1_{i,i+2} \\[2ex]
\tau_{n+2} = \dfrac{1}{L-2}\displaystyle\sum_{i=1}^{L-2} H^2_{i,i+2} \\[1ex]
\quad \cdots \\[1ex]
\tau_{2n} = \dfrac{1}{L-2}\displaystyle\sum_{i=1}^{L-2} H^n_{i,i+2} \\[1ex]
\quad \cdots \\[1ex]
\tau_{n\lambda-1} = \dfrac{1}{L-\lambda}\displaystyle\sum_{i=1}^{L-\lambda} H^{n-1}_{i,i+\lambda} \\[2ex]
\tau_{n\lambda} = \dfrac{1}{L-\lambda}\displaystyle\sum_{i=1}^{L-\lambda} H^n_{i,i+\lambda}
\end{cases}
\tag{4}
$$

where $H^n_{i,j}$ is the physicochemical property correlation function and can be computed

according to the following equation:

$$
H^n_{i,j} = h^n(\mathrm{R}_i).h^n(\mathrm{R}_j)
\tag{5}
$$

where $h^n(\mathrm{R}_i)$ and $h^n(\mathrm{R}_j)$ are the values of the n-th type of physicochemical property of $\mathrm{R}_i$

and $\mathrm{R}_j$ in Equation (1), respectively. It is noteworthy that before substituting the values of all

physicochemical properties into Equation (5), they were undergone a standard conversion as

described below:

$$
h^k(\mathrm{R}_i) = \frac{h^k_0(\mathrm{R}_i) - \sum_{\alpha=1}^{20} h^k_0(\mathrm{R}_\alpha)/20}{\sqrt{\sum_{u=1}^{20}[h^k_0(\mathrm{R}_i) - \sum_{\alpha=1}^{20} h^k_0(\mathrm{R}_\alpha)/20]^2}}
\tag{6}
$$

where $\mathrm{R}_i$ (i = 1, 2, . . . , 20) denotes the 20-standard amino acid in the alphabetical order of

their single-letter codes. $h^k_0(\mathrm{R}_i)$ is the initial value of the k-th type of physicochemical

149   property for amino acid residue $R_i$. Nine kinds of physicochemical properties, namely

150   hydrophobicity, hydrophilicity, mass, pK1, pK2, pI, rigidity, flexibility and irreplaceability,

151   were considered in this report.

## Feature Selection

153   Generally, not all features make an equal contribution to the prediction system. A part of

154   features make significant contributions, while some others make less important contributions

155   (Zhao et al., 2016). Feature selection, thus, is a critical step to reduce feature dimensionality

156   and build a highly effective prediction model (Su et al., 2018; Tang et al., 2016). In this work,

157   the fselect.py program in the LIBSVM 3.23 package was applied to evaluate each feature's

158   significance to the classification system (Chang et al., 2011). As a consequence, each feature

159   corresponds to an F-score. The greater F-score implies the larger importance of the

160   corresponding feature to the prediction model. We rearranged all features by F-scores in

161   descending order. The incremental feature selection (IFS) strategy was then utilized to

162   determine the optimal feature subset (He et al., 2016b; Li et al., 2017), which can produce the

163   maximal accuracy. Feature selection was conducted as follows: (i) investigating the accuracy

164   of the first feature subset which included the feature with the largest F-score; (ii) examining

165   the accuracy of the second feature subset that was generated by appending the feature with the

166   second largest F-score; (iii) iterating the second step from the larger F-score to the smaller F

167   score until all candidate features were added. The best feature subset with the highest

168   accuracy can be finally obtained.

## Support Vector Machine

170   The SVM is a powerful supervised learning method, which has been widely applied in

171    classification (He et al., 2016b; Kang et al., 2018; Li et al., 2017; Ru et al., 2014) and

172    regression analysis. In this study, we utilized the LIBSVM 3.23 program (Chang & Lin, 2011)

173    that could be freely available for download from http://www.csie.ntu.edu.tw/~cjlin/libsvm/.

174    We chose the radial basis function (RBF) kernel as the kernel function. The optimal kernel

175    width parameter γ and penalty constant C were selected by using the parameter selection tool

176    in the LIBSVM 3.23 (Chang & Lin, 2011).

177    **Performance Evaluation**

178    The ten-fold cross-validation was adopted to evaluate the predictive model in this study. Four

179    commonly-used parameters, including sensitivity (Sn), specificity (Sp), accuracy (Acc) and

180    Matthews correlation coefficient (MCC), were employed to investigate the performance of the

181    constructed model. These measures were expressed as follows:

182
$$\text{Sn} = \frac{TP}{TP + FN} \tag{7}$$

183
$$\text{Sp} = \frac{TN}{FP + TN} \tag{8}$$

184
$$\text{Acc} = \frac{TP + TN}{TP + FN + FP + TN} \tag{9}$$

185
$$\text{MCC} = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \tag{10}$$

186    where *TP* and *TN* denote the number of true positives and negatives, respectively. *FP* and *FN*

187    are the number of false positives and negatives, respectively. The area under the ROC curve

188    (AUC) was also calculated as a performance measure. The AUC ranges from zero to one. The

189    AUC of one represents a perfect prediction, 0.5 a random guess.

190        To estimate the statistical significance of the predictive accuracy, a permutation test with

191    1,000 shuffles was performed by exchanging the labels of the benchmark dataset. The ten-fold

192    cross-validation was then conducted against the label-rearranged dataset. Thus, each

193 permutation trial corresponds to an accuracy value. The $p$ value was calculated by the number

194 of permutations that the Acc produced by the permuted dataset was higher than Acc based on

195 the un-permuted dataset divided by the overall shuffle times. P values of <0.05 were referred

196 to as statistically significant.

## Standalone Version Implementation

198 The standalone version of PhD7Faster 2.0 was developed with open source Qt 5.7 under the

199 GPL & LGPLv3 licenses, which uses standard C++ for developing multiple-platform

200 applications. Both graphical user interface (GUI) and command-line versions of PhD7Faster

201 2.0 were implemented. We provided different versions for Windows and Linux systems with

202 little or no modification. All versions and source code are freely available at

203 http://i.uestc.edu.cn/sarotup3/download.html.

## RESULTS

> **Commented [TT1]:** This should also be a part of materials and methods

## Parameter Optimization

206 Two important parameters: $\lambda$ and $w$ in Equation (3) were necessary to be optimized before

207 building the model. To obtain the best parameters, multiple experiments were performed

208 according to the following standard:

$$\begin{cases} 1 \leq \lambda \leq 6 \text{ with step } \Delta = 1 \\ 0.05 \leq w \leq 0.70 \text{ with step } \Delta = 0.05 \end{cases} \tag{1}$$

210 Thus, a total of $6 \times 14 = 84$ individual combinations were obtained. Then, we used the ten-

211 fold cross-validation to investigate the accuracy of the model, which was built with SVM and

212 the feature set without feature selection. $\lambda=3$ and $w=0.15$ produced the highest accuracy,

213 which was considered as the best parameter combination.

## Performance of PhD7Faster 2.0

215 ~~The optimal feature subset with 644 features was determined through feature selection against~~

216 ~~8027 features including 8000 tripeptide features and 27 PseAAC features. The SVM-based~~

217 ~~model was then trained with the optimal feature set. The results from the ten-fold cross-~~

218 ~~validation showed that the Acc of the predictive model was 81.84 % with MCC of 0.64, Sn of~~

219 ~~84.51% and Sp of 79.17% when the threshold to distinguish between predicted positives and~~

220 ~~negatives ($tp$) was set to be 0.5. The ROC curve for model tuning is shown in Fig. 1, where~~

221 ~~the AUC is approximately 0.90. The permutation test resulted in a p-value of <0.001. The~~

222 ~~above results indicated that PhD7Faster 2.0 achieved a promising performance.~~

## Web and Standalone Versions of PhD7Faster 2.0

For the convenience of users, the SVM-based predictive model was implemented into a user-friendly web server, called PhD7Faster 2.0, which is freely available at http://i.uestc.edu.cn/sarotup3/cgi-bin/PhD7Faster.pl. The standalone GUI and command-line versions of PhD7Faster 2.0 for Windows and Linux systems were also provided. The interface as well as the utilization of the GUI version is remarkably similar to those of the web version (Fig. 2). A dataset with 20,000 peptides from the Ph.D.-7 phage display library was constructed (see testdataset.fasta in Supplementary Data). The standalone PhD7Faster 2.0 can complete analysis of the dataset within 60 seconds on a regular computer with Intel Core i3 Processor and 4GB RAM, which suggests that PhD7Faster 2.0 is highly efficient in processing massive datasets. PhD7Faster 2.0 was integrated into the SAROTUP 3.0 suite, which contains a series of computational tools to identify TUPs.

## RESULTS~~DISCUSSION~~

## Performance of PhD7Faster 2.0

The optimal feature subset with 644 features was determined through feature selection against 8027 features including 8000 tripeptide features and 27 PseAAC features. The SVM-based model was then trained with the optimal feature set. The results from the ten-fold cross-validation showed that the Acc of the predictive model was 81.84 % with MCC of 0.64, Sn of 84.51% and Sp of 79.17% when the threshold to distinguish between predicted positives and negatives ($tp$) was set to be 0.5. The ROC curve for model tuning is shown in Fig. 1, where the AUC is approximately 0.90. The permutation test resulted in a p-value of <0.001. The above results indicated that PhD7Faster 2.0 achieved a promising performance.

## Comparison between PhDFaster 2.0 and 1.0

Parasitic sequences were identified significantly enriched in the amplified Ph.D.-7 phage display library by differential enrichment analysis of naïve and amplified Ph.D.-7 phage display libraries (Matochko et al., 2014). These parasitic peptides were grouped into the positive dataset of PhD7Faster 2.0. However, the positive dataset of PhD7Faster 1.0 was constructed based on threshold in copy numbers after one round of amplification in one replicate of sequencing data, irrespective of copy numbers in the naïve Ph.D.-7 library. Peptides with high abundances in both the naïve and amplified Ph.D.-7 libraries may also be selected as fast-growing sequences. Therefore, the positive training dataset of PhD7Faster 2.0 is more reliable than that of PhD7Faster 1.0.

PhD7Faster 2.0 was developed based on the combination of PseAAC and tripeptide composition, whereas only dipeptide composition was employed to build PhD7Faster 1.0. We also tried to use dipeptide composition to encode each peptide in the training dataset of PhD7Faster 2.0, but only 64% accuracy was obtained in the ten-fold cross-validation after

259 feature selection. PseAAC coupled with tripeptide composition has been used in multiple

260 protein prediction fields, such as predicting the subcellular localization of mycobacterial

261 proteins (Zhu et al., 2015) and predicting apoptosis protein subcellular location (Liao et al.,

262 2011). They contain more sequence-order information than dipeptide composition and hence

263 can better reflect the feature of a peptide sequence. Thus, PhD7Faster 2.0 has 5% sensitivity,

264 2% accuracy and 0.04 MCC higher than PhD7Faster 1.0 (Table 1).

265     The standalone PhD7Faster 2.0 is empowered to identify PrTUPs within output of

266 conventional phage display as well as large next-generation sequencing data, whereas

267 PhD7Faster 1.0 can only work with small-scale data sets (several hundreds of peptides). This

268 important improvement makes PhD7Faster 2.0 as an enhanced and powerful tool for scanning

269 and reporting PrTUPs from the Ph.D.-7 phage display library. The emergence of PhD7Faster

270 2.0 highlights the significance of high throughput sequencing of different types of phage

271 display libraries and developing bioinformatics tools for identifying PrTUPs from these

272 libraries.

## PhD7Faster 2.0 cannot predict the censorship in the Ph.D. libraries

274 It is possible that some peptides are likely to be censored from being displayed on the phage

275 in the first place. The censorship of positively charged amino acids has been reported since

276 these residues suppress proper insertion of pIII into the inner membrane of *Escherichia coli*

277 (*E. coli*), thus decreasing efficiency of the assembly and extrusion of phage clones (Peters et

278 al., 1994). Makowski et al. also observed that peptides of α-helix or β-sheet conformations

279 were censored in Ph.D.-12 and Ph.D.-C7C libraries (Rodi et al., 2002). Plückthun and

280 coworkers have shown that maturely folded proteins are displayed poorly via the Sec

281 translocation pathway (Steiner et al., 2006). However, this censorship is a completely

282 different phenomenon from that of phage growing faster. Therefore, PhD7Faster 2.0 cannot be

283 able to predict this censorship.

284 **PhD7Faster 2.0 predict PrTUPs in the Ph.D.-7 library**

285 The PrTUPs have significantly higher proliferation rates than normal-growing phage and are

286 favored during the amplification steps. The proliferation advantage of some PrTUPs have

287 been verified to be intrinsic to mutations in the 5'-untranslated region (UTR) of gene II in

288 M13 phage (Brammer et al., 2008; Nguyen et al., 2014; Zygiel et al., 2017). Zygiel et al. have

289 also described the likelihood that these mutations compensate for the replication defect

290 afforded by the lacZα insert present in the M13 bacteriophage-based vector upon which the

291 Ph.D.-7 (and Ph.D.-12) library was based (Zygiel et al., 2017). Thus, the particular peptide

292 displayed (e.g., HAIYPRH, GKPMPPM, AKIDART) is merely a stowaway on a clone that

293 propagates fast due to its gene II 5'-UTR mutation(s). In these clones, the peptide itself is

294 completely arbitrary; it just happens to be the peptide displayed on a clone that picked up a

295 mutation prior to or during library construction. As these mutations in the phage genome are

296 unrelated to the displayed peptide, PhD7Faster 2.0 may not be able to predict this type of

297 PrTUPs in the Ph.D.-7 library. In addition, Smith et al. indicated that the enhanced

298 propagation rate of some PrTUPs may be due to the displayed peptide (Thomas et al., 2010),

299 and PhD7Faster 2.0 can be used to predict this type of PrTUPs in the Ph.D.-7 library.

300 However, no direct evidence supports that displayed peptides allow the phage to propagate

301 faster, and the biological mechanism remains to be further examined.

302 **CONCLUSION**

303 In this report, we propose an SVM-based tool, PhD7Faster 2.0, for predicting clones growing

304 faster from the Ph.D.-7 phage display library. Ten-fold cross-validation results show that

305 PhD7Faster 2.0 achieves an accuracy of 81.84% with 0.64 MCC and 0.90 AUC. The

306 standalone version of the tool was also developed, which ~~is capable of predicting~~can predict

307 PrTUPs within both traditional biopanning data and next generation phage display data. We

308 also implemented a web-server for the proposed method, which can be freely accessible from

309 http://i.uestc.edu.cn/sarotup3/cgi-bin/PhD7Faster.pl.

## ACKNOWLEDGEMENTS

## REFERENCES

314 **Bakhshinejad B, Zade HM, Shekarabi HS, and Neman S. 2016.** Phage display biopanning and isolation of target-
315     unrelated peptides: in search of nonspecific binders hidden in a combinatorial library. *Amino Acids*
316     **48(12)**:2699-2716 DOI 10.1007/s00726-016-2329-6.
317 **Brammer LA, Bolduc B, Kass JL, Felice KM, Noren CJ, and Hall MF. 2008.** A target-unrelated peptide in an M13
318     phage display library traced to an advantageous mutation in the gene II ribosome-binding site. *Anal*
319     *Biochem* **373(1)**:88-98 DOI 10.1016/j.ab.2007.10.015.
320 **Causa F, Della Moglie R, Iaccino E, Mimmi S, Marasco D, Scognamiglio PL, Battista E, Palmieri C, Cosenza C,**
321     **Sanguigno L, Quinto I, Scala G, and Netti PA. 2013.** Evolutionary screening and adsorption behavior of
322     engineered M13 bacteriophage and derived dodecapeptide for selective decoration of gold interfaces.
323     *J Colloid Interface Sci* **389(1)**:220-229 DOI 10.1016/j.jcis.2012.08.046.
324 **Chang C-C, and Lin C-J. 2011.** LIBSVM: a library for support vector machines. *ACM transactions on intelligent*
325     *systems and technology (TIST)* **2(3)**:27 DOI 10.1145/1961189.1961199.
326 **Chou K-C. 2011.** Some remarks on protein attribute prediction and pseudo amino acid composition. *Journal of*
327     *theoretical biology* **273(1)**:236-247 DOI 10.1016/j.jtbi.2010.12.024.
328 **Chou KC. 2001.** Prediction of protein cellular attributes using pseudo-amino acid composition. *Proteins* **43(3)**:246-
329     255 DOI 10.1002/prot.1035.
330 **Chou KC. 2005.** Using amphiphilic pseudo amino acid composition to predict enzyme subfamily classes.
331     *Bioinformatics* **21(1)**:10-19 DOI 10.1093/bioinformatics/bth466.
332 **Christiansen A, Kringelum JV, Hansen CS, Bogh KL, Sullivan E, Patel J, Rigby NM, Eiwegger T, Szepfalusi Z, de**
333     **Masi F, Nielsen M, Lund O, and Dufva M. 2015.** High-throughput sequencing enhanced phage display
334     enables the identification of patient-specific epitope motifs in serum. *Sci Rep* **5**:12913 DOI
335     10.1038/srep12913.

336  **He B, Chai G, Duan Y, Yan Z, Qiu L, Zhang H, Liu Z, He Q, Han K, Ru B, Guo FB, Ding H, Lin H, Wang X, Rao N, Zhou**
337      **P, and Huang J. 2016a.** BDB: biopanning data bank. *Nucleic Acids Res* **44(D1)**:D1127-1132 DOI
338      10.1093/nar/gkv1100.
339  **He B, Jiang L, Duan Y, Chai G, Fang Y, Kang J, Yu M, Li N, Tang Z, Yao P, Wu P, Derda R, and Huang J. 2018.**
340      Biopanning data bank 2018: hugging next generation phage display. *Database (Oxford)* **2018**:1-8 DOI
341      10.1093/database/bay032.
342  **He B, Kang J, Ru B, Ding H, Zhou P, and Huang J. 2016b.** SABinder: A Web Service for Predicting Streptavidin-
343      Binding Peptides. *Biomed Res Int* **2016**:9175143 DOI 10.1155/2016/9175143.
344  **He B, Mao C, Ru B, Han H, Zhou P, and Huang J. 2013.** Epitope mapping of metuximab on CD147 using phage
345      display and molecular docking. *Comput Math Methods Med* **2013**:983829 DOI 10.1155/2013/983829.
346  **Huang J, Ru B, Li S, Lin H, and Guo FB. 2010.** SAROTUP: scanner and reporter of target-unrelated peptides. *J*
347      *Biomed Biotechnol* **2010**:101932 DOI 10.1155/2010/101932.
348  **Huang J, Ru B, Zhu P, Nie F, Yang J, Wang X, Dai P, Lin H, Guo FB, and Rao N. 2012.** MimoDB 2.0: a mimotope
349      database and beyond. *Nucleic Acids Res* **40(Database issue)**:D271-277 DOI 10.1093/nar/gkr922.
350  **Hung LY, Fu CY, Wang CH, Chuang YJ, Tsai YC, Lo YL, Hsu PH, Chang HY, Shiesh SC, Hsu KF, and Lee GB. 2018.**
351      Microfluidic platforms for rapid screening of cancer affinity reagents by using tissue samples.
352      *Biomicrofluidics* **12(5)**:054108 DOI 10.1063/1.5050451.
353  **Kang J, Fang Y, Yao P, Li N, Tang Q, and Huang J. 2018.** NeuroPP: A Tool for the Prediction of Neuropeptide
354      Precursors Based on Optimal Sequence Composition. *Interdiscip Sci* DOI 10.1007/s12539-018-0287-2.
355  **Li N, Kang J, Jiang L, He B, Lin H, and Huang J. 2017.** PSBinder: A Web Service for Predicting Polystyrene Surface-
356      Binding Peptides. *Biomed Res Int* **2017**:5761517 DOI 10.1155/2017/5761517.
357  **Liao B, Jiang JB, Zeng QG, and Zhu W. 2011.** Predicting apoptosis protein subcellular location with PseAAC by
358      incorporating    tripeptide    composition.    *Protein    Pept    Lett*    **18(11)**:1086-1092    DOI
359      10.2174/092986611797200931.
360  **Lin H, Ding C, Yuan L-F, Chen W, Ding H, Li Z-Q, Guo F-B, Huang J, and Rao N-N. 2013.** Predicting subchloroplast
361      locations of proteins based on the general form of Chou's pseudo amino acid composition: approached
362      from optimal tripeptide composition. *International Journal of Biomathematics* **6(02)**:1350003 DOI
363      10.1142/S1793524513500034.
364  **Mandava S, Makowski L, Devarapalli S, Uzubell J, and Rodi DJ. 2004.** RELIC--a bioinformatics server for
365      combinatorial peptide analysis and identification of protein-ligand interaction sites. *Proteomics*
366      **4(5)**:1439-1460 DOI 10.1002/pmic.200300680.
367  **Martins IM, Reis RL, and Azevedo HS. 2016.** Phage Display Technology in Biomaterials Engineering: Progress and
368      Opportunities for Applications in Regenerative Medicine. *ACS Chem Biol* **11(11)**:2962-2980 DOI
369      10.1021/acschembio.5b00717.
370  **Matochko WL, Cory Li S, Tang SK, and Derda R. 2014.** Prospective identification of parasitic sequences in phage
371      display screens. *Nucleic Acids Res* **42(3)**:1784-1798 DOI 10.1093/nar/gkt1104.
372  **Menendez A, and Scott JK. 2005.** The nature of target-unrelated peptides recovered in the screening of phage-
373      displayed random peptide libraries with antibodies. *Anal Biochem* **336(2)**:145-157 DOI
374      10.1016/j.ab.2004.09.048.
375  **Ngubane NA, Gresh L, Ioerger TR, Sacchettini JC, Zhang YJ, Rubin EJ, Pym A, and Khati M. 2013.** High-throughput
376      sequencing enhanced phage display identifies peptides that bind mycobacteria. *PLoS One* **8(11)**:e77844
377      DOI 10.1371/journal.pone.0077844.
378  **Nguyen KT, Adamkiewicz MA, Hebert LE, Zygiel EM, Boyle HR, Martone CM, Melendez-Rios CB, Noren KA,**
379      **Noren CJ, and Hall MF. 2014.** Identification and characterization of mutant clones with enhanced

380      propagation rates from phage-displayed peptide libraries. *Anal Biochem* **462**:35-43 DOI
381      10.1016/j.ab.2014.06.007.

382 **Pande J, Szewczyk MM, and Grover AK. 2010.** Phage display: concept, innovations, applications and future.
383      *Biotechnol Adv* **28(6)**:849-858 DOI 10.1016/j.biotechadv.2010.07.004.

384 **Peters EA, Schatz PJ, Johnson SS, and Dower WJ. 1994.** Membrane insertion defects caused by positive charges
385      in the early mature region of protein pIII of filamentous phage fd can be corrected by prlA suppressors.
386      *J Bacteriol* **176(14)**:4296-4305 DOI 10.1128/jb.176.14.4296-4305.1994.

387 **Rentero Rebollo I, Sabisz M, Baeriswyl V, and Heinis C. 2014.** Identification of target-binding peptide motifs by
388      high-throughput sequencing of phage-selected peptides. *Nucleic Acids Res* **42(22)**:e169 DOI
389      10.1093/nar/gku940.

390 **Rodi DJ, Soares AS, and Makowski L. 2002.** Quantitative assessment of peptide sequence diversity in M13
391      combinatorial peptide phage display libraries. *J Mol Biol* **322(5)**:1039-1052 DOI 10.1016/S0022-
392      2836(02)00844-6.

393 **Ru B, Huang J, Dai P, Li S, Xia Z, Ding H, Lin H, Guo F, and Wang X. 2010.** MimoDB: a new repository for mimotope
394      data derived from phage display technology. *Molecules* **15(11)**:8279-8288 DOI
395      10.3390/molecules15118279.

396 **Ru B, t Hoen PA, Nie F, Lin H, Guo FB, and Huang J. 2014.** PhD7Faster: predicting clones propagating faster from
397      the Ph.D.-7 phage display peptide library. *J Bioinform Comput Biol* **12(1)**:1450005 DOI
398      10.1142/S021972001450005X.

399 **Shtatland T, Guettler D, Kossodo M, Pivovarov M, and Weissleder R. 2007.** PepBank--a database of peptides
400      based on sequence text mining and public peptide data sources. *BMC Bioinformatics* **8**:280 DOI
401      10.1186/1471-2105-8-280.

402 **Steiner D, Forrer P, Stumpp MT, and Pluckthun A. 2006.** Signal sequences directing cotranslational translocation
403      expand the range of proteins amenable to phage display. *Nat Biotechnol* **24(7)**:823-831 DOI
404      10.1038/nbt1218.

405 **Su ZD, Huang Y, Zhang ZY, Zhao YW, Wang D, Chen W, Chou KC, and Lin H. 2018.** iLoc-lncRNA: predict the
406      subcellular location of lncRNAs by incorporating octamer composition into general PseKNC.
407      *Bioinformatics* **34(24)**:4196-4204 DOI 10.1093/bioinformatics/bty508.

408 **Swaminathan S, and Cui Y. 2013.** Recognition of epoxy with phage displayed peptides. *Mater Sci Eng C Mater
409      Biol Appl* **33(5)**:3082-3084 DOI 10.1016/j.msec.2013.02.011.

410 **t Hoen PA, Jirka SM, Ten Broeke BR, Schultes EA, Aguilera B, Pang KH, Heemskerk H, Aartsma-Rus A, van
411      Ommen GJ, and den Dunnen JT. 2012.** Phage display screening without repetitious selection rounds.
412      *Anal Biochem* **421(2)**:622-631 DOI 10.1016/j.ab.2011.11.005.

413 **Tang H, Chen W, and Lin H. 2016.** Identification of immunoglobulins using Chou's pseudo amino acid composition
414      with feature selection technique. *Mol Biosyst* **12(4)**:1269-1275 DOI 10.1039/c5mb00883b.

415 **Thomas WD, Golomb M, and Smith GP. 2010.** Corruption of phage display libraries by target-unrelated clones:
416      diagnosis and countermeasures. *Anal Biochem* **407(2)**:237-240 DOI 10.1016/j.ab.2010.07.037.

417 **Vodnik M, Strukelj B, and Lunder M. 2012.** HWGMWSY, an unanticipated polystyrene binding peptide from
418      random phage display libraries. *Anal Biochem* **424(2)**:83-86 DOI 10.1016/j.ab.2012.02.013.

419 **Vodnik M, Zager U, Strukelj B, and Lunder M. 2011.** Phage display: selecting straws instead of a needle from a
420      haystack. *Molecules* **16(1)**:790-817 DOI 10.3390/molecules16010790.

421 **Zade HM, Keshavarz R, Shekarabi HSZ, and Bakhshinejad B. 2017.** Biased selection of propagation-related TUPs
422      from phage display peptide libraries. *Amino Acids* **49(8)**:1293-1308 DOI 10.1007/s00726-017-2452-z.

423 **Zhang Y, He B, Liu K, Ning L, Luo D, Xu K, Zhu W, Wu Z, Huang J, and Xu X. 2017.** A novel peptide specifically

424      binding to VEGF receptor suppresses angiogenesis in vitro and in vivo. *Signal Transduction and Targeted*
425      *Therapy* **2**:17010 DOI 10.1038/sigtrans.2017.10.

426      **Zhao YW, Lai HY, Tang H, Chen W, and Lin H. 2016.** Prediction of phosphothreonine sites in human proteins by
427      fusing different features. *Sci Rep* **6**:34817 DOI 10.1038/srep34817.

428      **Zhu PP, Li WC, Zhong ZJ, Deng EZ, Ding H, Chen W, and Lin H. 2015.** Predicting the subcellular localization of
429      mycobacterial proteins by incorporating the optimal tripeptides into the general form of pseudo amino
430      acid composition. *Mol Biosyst* **11(2)**:558-563 DOI 10.1039/c4mb00645c.

431      **Zygiel EM, Noren KA, Adamkiewicz MA, Aprile RJ, Bowditch HK, Carroll CL, Cerezo MAS, Dagher AM, Hebert**
432      **CR, Hebert LE, Mahame GM, Milne SC, Silvestri KM, Sutherland SE, Sylvia AM, Taveira CN,**
433      **VanValkenburgh DJ, Noren CJ, and Hall MF. 2017.** Various mutations compensate for a deleterious
434      lacZalpha insert in the replication enhancer of M13 bacteriophage. *PLoS One* **12(4)**:e0176421 DOI
435      10.1371/journal.pone.0176421.

436