

Supplement

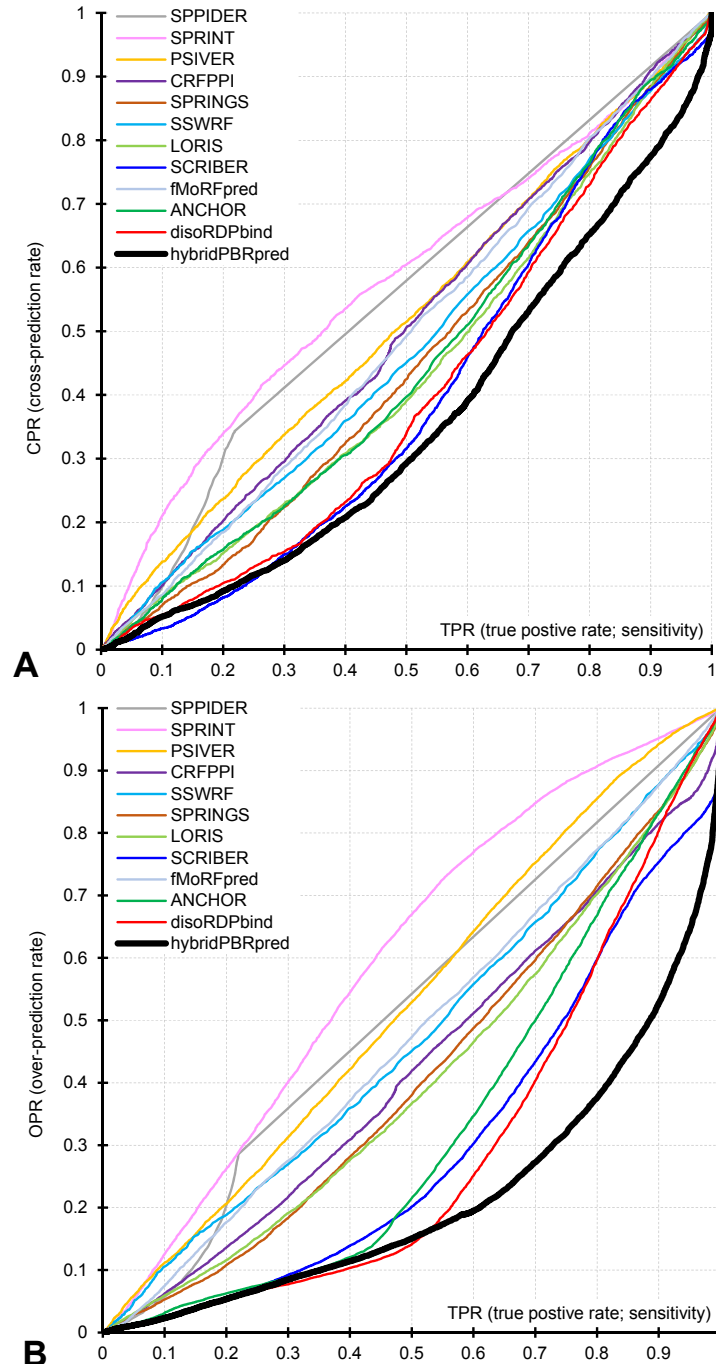
Prediction of protein-binding residues: dichotomy of sequence-based methods developed using structured complexes vs. disordered proteins

Jian Zhang¹, Sina Ghadermarzi², and Lukasz Kurgan^{2*}

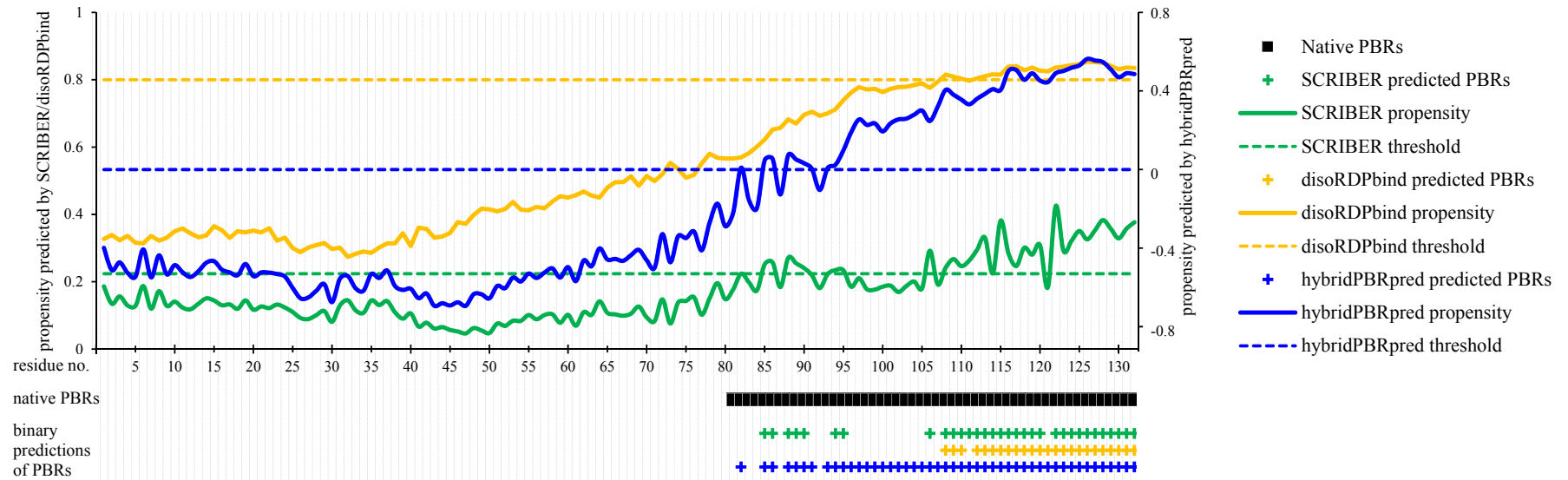
¹School of Computer and Information Technology, Xinyang Normal University, Xinyang, China, 464000;

²Department of Computer Science, Virginia Commonwealth University, Richmond, VA, USA, 23284

*corresponding author: phone: +1-804-827-3986; fax: +1-804-828-2771; email: lkurgan@vcu.edu



Supplementary Figure S1. The cross-prediction curves (panel A) and the over-prediction curves (panel B) computed on the entire benchmark dataset. The cross-predictions consider the residues that interact with all non-protein partners (nucleic acids and small ligands).



Supplementary Figure S2. Comparison of predictions generated by hybridPBRpred, SCRIBER and disoRDPbind for art v 1 protein, a pollen allergen from *Artemisia vulgaris* (UniProt ID: Q84ZX5; DisProt ID: DP00807). The solid lines at the top represent the putative propensities for protein binding generated by SCRIBER (in green), disoRDPbind (orange) and hybridPBRpred (blue). The propensities are converted into the binary predictions (PBRs vs. non-PBRs) using the predictor specific color-coded thresholds represented by the dashed horizontal lines. The binary predictions are shown using the color-coded horizontal lines at the bottom of the figure. The black markers denote the native annotations of the protein-binding residues.

Supplementary Table S1. The cross-prediction rates (AUCPC) and over-prediction rates (AUOPC) for the 11 predictors of PBRs and the new hybridPBRpred. We note that lower values of these metrics indicate stronger predictive performance. We test robustness of predictions across diverse protein sets by bootstrapping 10 sets of 50% randomly selected proteins from a given dataset; we report the corresponding averages \pm standard deviations. For each dataset, we evaluate statistical significance of differences between the predictor shown in bold font (the best overall predictor) and the other predictors based on the 10 bootstrapped tests; ++ and + mean that the best predictor is significantly better with p -value < 0.001 and p -value < 0.05 , respectively; = means that the difference is not significant (p -value ≥ 0.05); -- and - mean that the best predictor is significantly worse with p -value < 0.001 and p -value < 0.05 , respectively.

Dataset	Type of Methods	Predictors	AUOPC	AUCPC residues binding non-protein partners (nucleic acids and small ligands)	AUCPC _{NA} nucleic acid-binding residues	AUCPC _{SL} small ligand binding residues	
Structure-annotated proteins	Trained on proteins annotated from structure	SPPIDER	0.486 \pm 0.012 ++	0.626 \pm 0.011 ++	0.605 \pm 0.039 ++	0.632 \pm 0.012 ++	
		SPRINT	0.431 \pm 0.014 ++	0.684 \pm 0.021 ++	0.537 \pm 0.040 ++	0.706 \pm 0.015 ++	
		PSIVER	0.441 \pm 0.010 ++	0.568 \pm 0.024 ++	0.663 \pm 0.109 ++	0.557 \pm 0.030 ++	
		CRF-PPI	0.320 \pm 0.010 ++	0.443 \pm 0.024 ++	0.479 \pm 0.044 ++	0.444 \pm 0.024 ++	
		SSWRF	0.329 \pm 0.012 ++	0.439 \pm 0.030 ++	0.500 \pm 0.059 ++	0.436 \pm 0.030 ++	
		SPRINGS	0.392 \pm 0.007 ++	0.522 \pm 0.013 ++	0.492 \pm 0.088 ++	0.527 \pm 0.012 ++	
		LORIS	0.363 \pm 0.010 ++	0.459 \pm 0.019 ++	0.461 \pm 0.083 ++	0.461 \pm 0.020 ++	
		SCRIBER	0.278\pm0.016	0.265\pm0.013	0.161\pm0.038	0.279\pm0.016	
	Trained on proteins annotated with disordered PBRs	fMoRFpred	0.498 \pm 0.010 ++	0.447 \pm 0.018 ++	0.437 \pm 0.038 ++	0.446 \pm 0.018 ++	
		ANCHOR	0.524 \pm 0.024 ++	0.497 \pm 0.046 ++	0.597 \pm 0.100 ++	0.485 \pm 0.038 ++	
		disoRDPbind	0.477 \pm 0.028 ++	0.465 \pm 0.040 ++	0.507 \pm 0.123 ++	0.458 \pm 0.035 ++	
		Hybrid predictor	hybridPBRpred	0.294 \pm 0.021 =	0.283 \pm 0.018 +	0.217 \pm 0.088 =	0.292 \pm 0.017 =
	Disorder-annotated proteins	Trained on proteins annotated from structure	SPPIDER	0.539 \pm 0.027 ++	0.508 \pm 0.031 ++	0.508 \pm 0.031 ++	0.537 \pm 0.076 ++
SPRINT			0.654 \pm 0.035 ++	0.484 \pm 0.041 ++	0.455 \pm 0.041 +	0.661 \pm 0.082 ++	
PSIVER			0.580 \pm 0.040 ++	0.478 \pm 0.050 ++	0.447 \pm 0.045 =	0.645 \pm 0.045 ++	
CRF-PPI			0.470 \pm 0.038 ++	0.503 \pm 0.051 ++	0.488 \pm 0.054 =	0.589 \pm 0.041 ++	
SSWRF			0.453 \pm 0.043 ++	0.465 \pm 0.069 +	0.449 \pm 0.076 =	0.566 \pm 0.042 ++	
SPRINGS			0.425 \pm 0.063 ++	0.389 \pm 0.070 =	0.363 \pm 0.063 =	0.515 \pm 0.068 ++	
LORIS			0.429 \pm 0.058 ++	0.402 \pm 0.073 =	0.377 \pm 0.071 =	0.527 \pm 0.060 ++	
		SCRIBER	0.272 \pm 0.051 +	0.454 \pm 0.090 +	0.464 \pm 0.094 =	0.422 \pm 0.095 +	
Trained on proteins annotated with disordered PBRs		fMoRFpred	0.458 \pm 0.013 ++	0.514 \pm 0.015 ++	0.521 \pm 0.013 ++	0.479 \pm 0.035 ++	
		ANCHOR	0.265 \pm 0.029 +	0.445 \pm 0.038 +	0.464 \pm 0.044 +	0.311 \pm 0.095 +	
		disoRDPbind	0.212\pm0.033	0.381\pm0.052	0.402\pm0.062	0.227\pm0.116	
		Hybrid predictor	hybridPBRpred	0.192 \pm 0.017 =	0.386 \pm 0.066 =	0.406 \pm 0.082 =	0.276 \pm 0.090 =
All proteins		Trained on proteins annotated from structure	SPPIDER	0.525 \pm 0.019 ++	0.546 \pm 0.027 ++	0.506 \pm 0.030 =	0.620 \pm 0.048 ++
	SPRINT		0.593 \pm 0.038 ++	0.537 \pm 0.044 ++	0.406 \pm 0.043 -	0.752 \pm 0.062 ++	
	PSIVER		0.536 \pm 0.030 ++	0.504 \pm 0.039 ++	0.432 \pm 0.038 =	0.625 \pm 0.037 ++	
	CRF-PPI		0.417 \pm 0.029 ++	0.483 \pm 0.041 ++	0.461 \pm 0.051 =	0.527 \pm 0.033 ++	
	SSWRF		0.406 \pm 0.030 ++	0.455 \pm 0.051 +	0.444 \pm 0.070 =	0.482 \pm 0.038 ++	
	SPRINGS		0.410 \pm 0.040 ++	0.429 \pm 0.048 +	0.382 \pm 0.052 -	0.504 \pm 0.038 ++	
	LORIS		0.403 \pm 0.036 ++	0.418 \pm 0.050 =	0.386 \pm 0.061 -	0.474 \pm 0.040 ++	
		SCRIBER	0.293 \pm 0.041 ++	0.400 \pm 0.061 =	0.425 \pm 0.068 =	0.366 \pm 0.060 ++	
	Trained on proteins annotated with disordered PBRs	fMoRFpred	0.472 \pm 0.006 ++	0.493 \pm 0.010 ++	0.526 \pm 0.009 +	0.441 \pm 0.009 ++	
		ANCHOR	0.341 \pm 0.035 ++	0.467 \pm 0.041 ++	0.565 \pm 0.042 ++	0.310 \pm 0.068 +	
		disoRDPbind	0.292 \pm 0.038 ++	0.422 \pm 0.049 +	0.525 \pm 0.033 +	0.272 \pm 0.028 +	
		Hybrid predictor	hybridPBRpred	0.211\pm0.023	0.376\pm0.037	0.469\pm0.057	0.228\pm0.047

Supplementary Table S2. Predictive performance of several approaches used to combine the scores generated by the structure-based predictor (SCRIBER) and the disordered-based predictor (disoRDPbind) on the complete benchmark dataset. We test robustness of predictions across diverse protein sets by bootstrapping 10 sets of 50% randomly selected proteins from the benchmark dataset; we report the corresponding averages \pm standard deviations. The binary predictions were generated from the propensity scores using threshold that ensures that the number of predicted PBRs equals to the number of native PBRs, allowing direct comparison between predictors. We evaluate statistical significance of differences between the best overall method shown in bold font (the scheme used in hybridPBRpred) and the other three approaches based on the 10 bootstrapped tests; ++ and + mean that the best predictor is significantly better with p-value < 0.001 and p-value < 0.05, respectively; = means that the difference is not significant (p-value \geq 0.05); -- and - mean that the best predictor is significantly worse with p-value < 0.001 and p-value < 0.05, respectively.

Methodology to combine predictors	Sensitivity	Specificity	F1	MCC	AUC	AULCratio	AUPRC
Maximum of the normalized predictions	0.223 \pm 0.035 ++	0.873 \pm 0.022 ++	0.223 \pm 0.035 ++	0.096 \pm 0.031 ++	0.639 \pm 0.018 ++	1.977 \pm 0.422 ++	0.206 \pm 0.035 ++
Minimum of the normalized predictions	0.168 \pm 0.019 ++	0.863 \pm 0.025 ++	0.168 \pm 0.019 ++	0.031 \pm 0.038 ++	0.474 \pm 0.018 ++	1.744 \pm 0.424 ++	0.152 \pm 0.017 ++
Mean of the normalized predictions	0.254 \pm 0.037 ++	0.878 \pm 0.017 ++	0.254 \pm 0.037 ++	0.132 \pm 0.028 ++	0.646 \pm 0.029 ++	2.311 \pm 0.250 ++	0.225 \pm 0.036 ++
hybridPBRpred	0.567\pm0.053	0.812\pm0.017	0.418\pm0.053	0.309\pm0.049	0.779\pm0.023	3.314\pm0.490	0.322\pm0.062
Maximum of the normalized predictions for residues which at least one method predicts as PBRs; otherwise mean of the normalized predictions							