

Table S5 Evaluation of our method with respect to comprehensive interaction prediction

<sup>1</sup> dataset	<sup>2</sup> neg.	<sup>3</sup> firsts	<sup>4</sup> P10275	<sup>4</sup> P11229	<sup>4</sup> P35367	<sup>5</sup> rec <sub>0.5</sub> (%)	<sup>5</sup> rec <sub>0.95</sub> (%)	<sup>6</sup> evaluation
(A) one-layer SVM								
<i>mlt</i>	16	—	714	1408	1187	100	98.97	82.50
<i>mlt</i>	14	—	709	1820	1634	100	*97.94	*79.02
<i>max</i>	16	—	4073	5956	6964	82.47	*56.70	*47.51
random	14	—	1896.7(±53.6)	10627.3(±648.9)	10204.0(±640.7)	100	99.66(±1.09)	69.20(±0.57)
<i>random</i>	16	—	1869.3(±136.1)	10503.3(±1250.7)	9305.3(±517.8)	100	99.66(±1.09)	69.45(±0.32)
(B) two-layer SVM- <i>subpos</i>								
<i>mlt</i>	14	10	177	535	451	96.91	93.81	75.56
<i>mlt</i>	14	11	205	671	491	96.91	91.75	73.54
<i>mlt</i>	14	9	239	513	403	95.88	91.75	73.87
<i>mlt</i>	14	8	290	456	363	88.66	82.47	66.58
<i>mlt</i>	12	10	224	561	612	95.88	92.78	73.25
<i>mlt</i>	16	10	162	466	415	94.85	89.69	73.47
<i>min</i>	14	10	2525	6098	3326	97.94	96.91	69.52
<i>mle</i>	14	10	168	526	599	97.94	92.78	74.79
<i>max</i>	14	10	32	386	191	92.78	*85.57	*72.27
<i>random</i>	14	10	848.3(±345.0)	1531.7(±628.9)	988.0(±411.4)	96.56(±2.89)	81.10(±19.44)	66.44(±7.82)
(C) two-layer SVM- <i>allpos</i>								
<i>max</i>	16	9	28	231	129	100	97.94	82.92
<i>max</i>	16	10	29	238	131	100	98.97	82.73
<i>max</i>	16	8	29	243	133	100	96.91	82.09
<i>max</i>	14	9	29	243	129	100	96.91	82.00
<i>mle</i>	16	9	28	267	140	100	100	80.99
<i>mlt</i>	16	9	67	248	141	100	100	80.72
<i>random</i>	16	9	74.7(±42.6)	255.3(±32.2)	146.7(±8.3)	100	100	80.67(±0.93)
(D) only compound SVM <sup>7</sup>								
—	—	—	640	1791	838	86.60	71.13	59.66
(E) similarity search <sup>8</sup>								
—	—	—	1869	1816	1580	—	—	—

<sup>1</sup> refers to negative data expansion rules (details are provided in Sec. 1.3 in Supplementary Materials). “random” indicates that three types of random pairs comprising a protein and a drug are used as negatives. The 95% confidence intervals are shown.

<sup>2</sup>: the number of negatives (= 1,750×*x*).

<sup>3</sup>: the number of the first-layer SVM models utilized for the construction of the second-layer SVM model.

<sup>4</sup>: target proteins whose ligands were predicted on the basis of 109,841 compounds. The number of predicted binding compounds is shown.

<sup>5</sup>: rec<sub>*x*</sub> is the recall rate (=TP/(TP+FN)) at the threshold *x*, ranging from 0 to 1. 0.5 is the threshold following the definition of SVM. TP: true positives, FN: false negatives.

<sup>6</sup>:

$$\text{evaluation} = 100 \times \left( \frac{1}{2} \left[ \text{rec}_{0.5} + \frac{\text{rec}_{0.95} + \text{prec}_{0.95}}{2\{1 + (1 - \text{rec}_{0.95})(1 - \text{prec}_{0.95})\}} \right] - \frac{\text{total \# of predicted positives} - \text{\# of known positives}}{\text{total \# of prediction targets} - \text{\# of known positives}} \right)$$

Here, prec<sub>*x*</sub> is the precision (=TP/(TP+FP)) at the threshold *x*. FP: false positives.

<sup>7</sup>: SVM model in which chemical compounds binding to each target protein were treated as positives and all other compounds in the DrugBank dataset were regarded as negatives.

<sup>8</sup>: A chemical compound *i* was predicted as a binding ligand of a protein  $\alpha$  by using the similarity method if  $\text{pred}_{\text{sim}}(i) = \max_{j \in A} |I \cap J| / |I \cup J| \geq 0.9$ , where *A* represents the known binding ligands of the protein  $\alpha$ , and *I* (or *J*) represents a set of substructures considered in calculating the feature vector of the chemical compounds.

\*: the threshold was set to 0.9 instead of 0.95 for the calculation of “evaluation”.