

**Supplementary Table 1.** CAMEO (9) pairwise comparison of 3D modelling servers. Dataset from last 6-months (2014-07-18 to 2015-01-09). IntFOLD3-TS (server33) is the reference server compared against older IntFOLD versions and other public servers on common targets.

Server Name	Number of models in common subset	Avg. IDDT difference*	Avg. IDDT-C $\alpha$ difference**	Avg. IDDT binding sites difference***	Avg. GDT_HA difference*
Robetta	220	-4.4	-3.88	-0.02	-3.43
RaptorX	216	-1.14	-1.12	<b>0.01</b>	-0.67
IntFOLD2-TS	214	<b>0.17</b>	<b>0.07</b>	0	-0.05
RBO Aleph	186	<b>1.19</b>	<b>1.68</b>	<b>0.02</b>	<b>2.94</b>
M4T	120	<b>1.74</b>	<b>1.97</b>	<b>0.03</b>	<b>1.66</b>
IntFOLD-TS	208	<b>3.56</b>	<b>3.75</b>	<b>0.02</b>	<b>3.7</b>
SWISS-MODEL	229	<b>4.32</b>	<b>4.8</b>	-0.01	<b>1.53</b>
Princeton_TEMPLATE	215	<b>4.71</b>	<b>1.88</b>	<b>0.11</b>	<b>2.27</b>
HHpredB	206	<b>6.85</b>	-1.01	-0.05	-0.56
NaiveBlits	119	<b>7.87</b>	<b>9.08</b>	<b>0.01</b>	<b>3.99</b>
NaiveBLAST	203	<b>10.67</b>	<b>12.43</b>	<b>0.06</b>	<b>6.51</b>
Phyre2	172	<b>15.16</b>	<b>11.75</b>	<b>0.07</b>	<b>6.41</b>

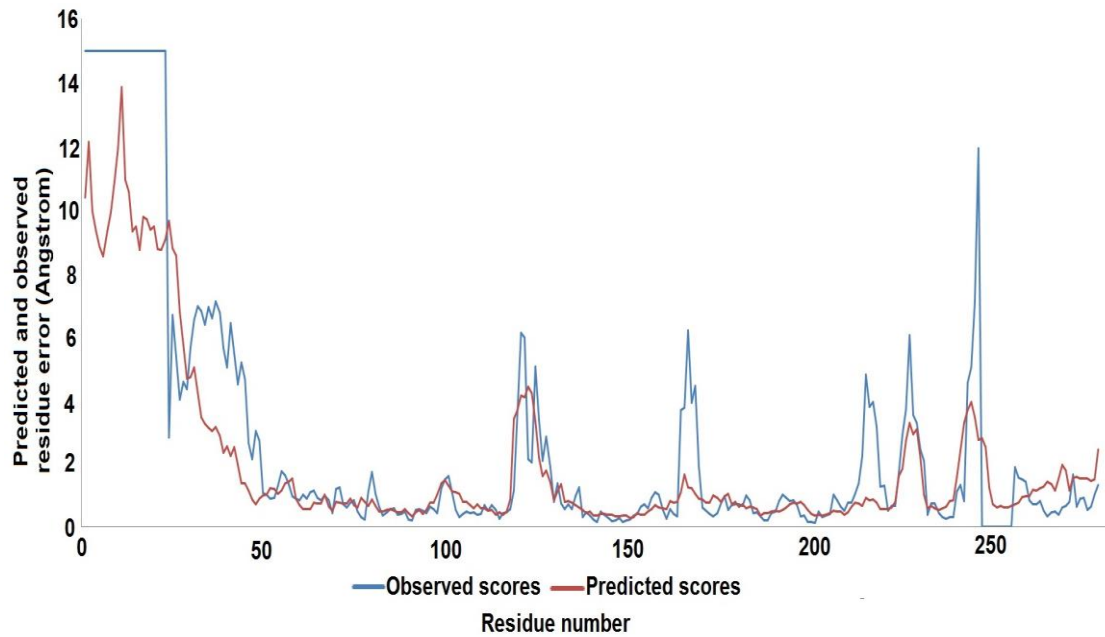
The IDDT score is the local Distance Difference Test on all atoms that assess the quality of the local atomic environment of a model. The IDDT acts as an incentive scheme for correctly predicted interatomic distances in a model at different threshold levels. If the IDDT score difference\* is below the given threshold level, the interaction between interatomic distances is considered to be preserved in the prediction. The final IDDT-all score is computed by averaging the fraction of correctly modelled interactions for the following four distance difference thresholds: 0.5, 1, 2, and 4 Å. This is the same as GDT\_HA which identifies collection of residues that are away from the target but not more than stated C $\alpha$  distance thresholds\*\* for varying superposition. The Avg. IDDT score\*\*\* is limited to those residues which form the binding site(s) on the respective target. Source of data and information on scoring methods: <http://www.cameo3d.org/>.

**Supplementary Table 2.** CAMEO (9) model quality estimation performance. Data from last 6-months (2014-07-18 to 2015-01-09). ModFOLD5 (server 9) and ModFOLD4 comparison with other methods.

Predictor Name	Number of models processed (out of 1935 submitted)	ROC		ROC normalised		PR		PR normalised	
		AUC <sub>0,1</sub>	AUC <sub>0,0.2</sub>	AUC <sub>0,1</sub>	AUC <sub>0,0.2</sub>	AUC <sub>0,1</sub>	AUC <sub>0.8,1</sub>	AUC <sub>0,1</sub>	AUC <sub>0.8,1</sub>
ModFOLD5 (Server 9)	1837	0.86	<b>0.59</b>	<b>0.82</b>	<b>0.56</b>	0.79	0.54	<b>0.75</b>	<b>0.51</b>
ModFOLD4	1745	<b>0.87</b>	<b>0.59</b>	0.78	0.53	<b>0.8</b>	<b>0.55</b>	0.72	0.49
Qmean 7.11	1927	0.82	0.5	0.81	0.5	0.74	0.47	0.74	0.47
ProQ2	1659	0.84	0.57	0.72	0.49	0.79	0.49	0.68	0.42
Verify3d smoothed	1935	0.71	0.34	0.71	0.34	0.62	0.4	0.62	0.4
Dfire v1.1	1935	0.66	0.25	0.66	0.25	0.53	0.39	0.53	0.39
Naive PSIBlast	1928	0.66	0.25	0.65	0.25	0.53	0.39	0.53	0.39

The ROC (Receiver Operating Characteristics) and PR (Precision and Recall) scores in bold indicate the highest prediction score of the ModFOLD5/4 comparing to other predictors. The AUC<sub>0,1</sub> (area under the curve) is the range from 0 to 1 (0 – perfect predictions and 0.5 – random prediction) that describes True Positive Rate (TPR) and False Positive Rate (FPR) computed for all positives thresholds. AUC<sub>0,0.2</sub> is the partial AUC of the ROC ‘trimmed’ threshold of 0.2 and scaled between 0 and 1. ROC<sub>normalised</sub> is the same scale as ROC but has been normalised. PR curve analysis is an alternative to ROC for task with large skew in the class distribution and it helps highlighting differences in the predictor performances that are not clear in the ROC. PR<sub>normalised</sub> is the same metric of the column PR with the values normalised by the submitted/received target ratio. Source of data and information on scoring methods: <http://www.cameo3d.org/>.

**Supplementary Figure 1.** Line plot showing the overlay of the predicted and observed per-residue errors in the top IntFOLD3-TS model for CASP target T0762.



**Supplementary Figure 2.** Scatter plot and correlation analysis of the predicted and observed per-residue errors in the top IntFOLD3-TS model for CASP target T0762. The plot shows a strong positive correlation between the observed and predicted residue scores: Spearman's rho (0.772), Kendall's tau B (0.588), and Pearson's R (0.917). The observed versus predicted correlation test indicated that there is a significant ( $P < 0.001$ ) positive correlation between the observed and predicted scores.

