

Supplementary materials for the following manuscript:

Sheng Wang, Zhen Li, Yizhou Yu and Jinbo Xu. Folding membrane proteins by deep transfer learning.

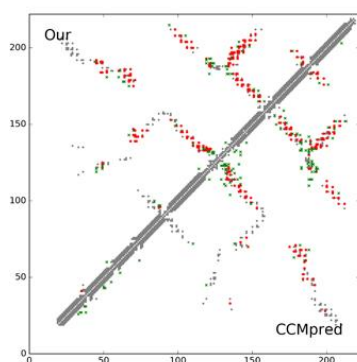
Table S1. A list of 510 non-redundant membrane proteins with solved structures in PDB, related to *Table 1 and STAR Methods section "Data for model parameter optimization and test"*.

| | | | | | | | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1a0sP | 1pw4A | 2evuA | 2lnlA | 2wpvB | 3cn5A | 3kvnA | 3udcA | 4chvA | 4il3A | 4or2A | 4wgvA | 5c8jI |
| 1ar1B | 1q16C | 2flcX | 2lomA | 2wsc1 | 3cx5C | 3l11A | 3ug9A | 4cskA | 4in5H | 4p6vB | 4wmzA | 5cfbA |
| 1bccE | 1q90A | 2f93B | 2loqA | 2wsc3 | 3d31C | 3lnmB | 3ukmA | 4czbA | 4in5L | 4p6vC | 4x5mA | 5ctgA |
| 1bctA | 1q90B | 2f95B | 2lorA | 2wscF | 3dd1A | 3lw54 | 3um7A | 4d5bA | 4j05A | 4p6vD | 4xk83 | 5d0yA |
| 1bhaA | 1qcrD | 2fynB | 2losA | 2wscG | 3dhwA | 3lw5H | 3uq7A | 4d6tD | 4j72A | 4p6vE | 4xnkA | 5dirA |
| 1c17M | 1qd6C | 2ge4A | 2lotA | 2wscH | 3dinE | 3m71A | 3ux4A | 4d6tG | 4j7cI | 4p6vF | 4xnvA | 5doqA |
| 1e7pC | 1qleC | 2gfpA | 2lp1A | 2wscK | 3dl8C | 3mk7A | 3v2wA | 4d6tJ | 4jkvA | 4p79A | 4xu4A | 5doqB |
| 1ehkB | 1rh5B | 2gr7A | 2m0qA | 2wscL | 3dl8E | 3mk7B | 3v5sA | 4d6uD | 4k1cA | 4pgrA | 4xxjA | 5ee7A |
| 1fftB | 1rh5C | 2gr8A | 2m20A | 2wswA | 3dwoX | 3mk7C | 3vmqA | 4djiA | 4kjrA | 4phzA | 4xydB | 5ek0A |
| 1fftC | 1rwtA | 2h8aA | 2m67A | 2wwbB | 3dwwA | 3mktA | 3vouA | 4dojA | 4knfA | 4pirA | 4y25A | 5ekeA |
| 1fw2A | 1s51B | 2h8pC | 2m6bA | 2wwbC | 3dzmA | 3mp7A | 3vr8C | 4dveA | 4kppA | 4px7A | 4y28G | 5eulE |
| 1fx8A | 1s51E | 2hdfA | 2m7gA | 2x4mA | 3effK | 3mp7B | 3vr8D | 4dxwA | 4kt0F | 4q2eA | 4y28K | 5ezmA |
| 1gzmA | 1s51X | 2ibzG | 2m8rA | 2xq2A | 3eh3A | 3njtA | 3vwiA | 4e1tA | 4kt0K | 4qncA | 4y28L | 5f1cA |
| 1h2sB | 1sqqK | 2ibzI | 2mafA | 2xutA | 3ejzA | 3nymA | 3wdoA | 4ea3A | 4ky0A | 4qndA | 4y7jA | 5fn2B |
| 1h6s1 | 1t16A | 2iubA | 2mfrA | 2y5yA | 3emnX | 3o0rB | 3wmfA | 4ezcA | 4l6rA | 4qtnA | 4ymkA | 5gaeh |
| 1iz1A | 1tlwA | 2j58A | 2mgvA | 2y69D | 3emoA | 3o7pA | 3wmm1 | 4f35A | 4l6v6 | 4quvA | 4ymsC | 5gaqA |
| 1iz1C | 1tqqA | 2j7aC | 2mm8A | 2y69G | 3fhhA | 3ohnA | 3wmmM | 4f41A | 4l6v8 | 4r1iA | 4ytpC | 5garO |
| 1jb0K | 1uunA | 2jafA | 2mmuA | 2y69I | 3fidA | 3orgA | 3wo7A | 4fqeA | 4ltoA | 4rdqA | 4ytpD | 5hk1A |
| 1k24A | 1uynX | 2jlnA | 2mn6A | 2y69J | 3g67A | 3oufA | 3wvfA | 4fuvA | 4m58A | 4rfsS | 4z34A | 5i1mV |
| 1kf6C | 1vclA | 2jo1A | 2mpnA | 2y69K | 3gi8C | 3p5nA | 3wxvA | 4g1uA | 4m64A | 4ri2A | 4z3nA | 5i20A |
| 1kf6D | 1vf5B | 2jp3A | 2mxbA | 2y69L | 3hd6A | 3pjsK | 3x29A | 4g7vS | 4mbsA | 4rjwA | 4z7fA | 5i32A |
| 1kqfB | 1vf5D | 2k01A | 2n4xA | 2y69M | 3hw9A | 3pjaA | 3x2rA | 4g80I | 4meeA | 4rl8A | 4zp0A | 5i6cA |
| 1kqfC | 1wrgA | 2k21A | 2n61A | 2yevB | 3iyzA | 3pwhA | 3x3bA | 4gbyA | 4mndA | 4rl9A | 4zr0A | 5i6zA |
| 1kzuA | 1xioA | 2k73A | 2n7qA | 2yevC | 3iz1A | 3q7kA | 3ze3A | 4gd3A | 4mqSA | 4rlcA | 4zr1A | 5id3A |
| 1lghA | 1x14A | 2k9pA | 2nmrA | 2yiuA | 3j08A | 3qe7A | 3zevA | 4gx5A | 4mt4A | 4rngA | 4zw9A | 5iofA |
| 1m56B | 1yc9A | 2kluA | 2nq2A | 2ynkA | 3j1zP | 3qnqA | 3zjzA | 4gycB | 4n74A | 4rp8A | 5a1sA | 5irxA |
| 1m56D | 1yewC | 2kogA | 2nr9A | 2z73A | 3j9tR | 3qraA | 3zk1A | 4h33A | 4n75A | 4ryiA | 5a40A | 5ivaA |
| 1m57A | 1yq3C | 2ks9A | 2nrgA | 2ziyA | 3jbrE | 3rbzA | 3zuxA | 4he8A | 4njnA | 4s0vA | 5a63C | 5iwsA |
| 1mm4A | 1yq3D | 2ksdA | 2o01F | 2zjsE | 3jcuD | 3rgwS | 4a2nB | 4he8C | 4nppA | 4tkrA | 5a63D | 5ixmB |
| 1mprA | 1zrtE | 2kseA | 2oarA | 2zxeB | 3jcuH | 3rkoA | 4atvA | 4he8D | 4ntjA | 4tq3A | 5a6eB | 5jagA |
| 1n71A | 1zzaA | 2ksfA | 2pnoA | 2zxeG | 3jcuK | 3rkoB | 4aw6A | 4hkrA | 4nykA | 4tquM | 5abbZ | |
| 1nekC | 2a01A | 2ksrA | 2q67A | 3a2sX | 3jcuR | 3rkoC | 4b4aA | 4hqjE | 4o6mA | 4tquN | 5araT | |
| 1nekD | 2a9hA | 2kyhA | 2q7mA | 3a7kA | 3jcuS | 3rkoD | 4bemJ | 4httA | 4o6yA | 4twkA | 5araW | |
| 1o5wA | 2akhA | 2l35A | 2qomA | 3anzA | 3jcuW | 3rkoF | 4bgnA | 4huqS | 4o9pA | 4u15A | 5awwG | |
| 1occD | 2akhB | 2l8sA | 2r6gF | 3b4rA | 3jcuX | 3rkoG | 4bog3 | 4huqT | 4o9pB | 4u4tA | 5awwY | |
| 1oedC | 2bg9A | 2lckA | 2r6gG | 3b5dA | 3jcuZ | 3s0xA | 4bpmA | 4hw9A | 4o9uB | 4u91A | 5awzA | |
| 1orsC | 2bl2A | 2lhfA | 2vpwC | 3b9wA | 3jycA | 3sljA | 4bwzA | 4hycA | 4od4A | 4uc1A | 5aymA | |
| 1p49A | 2cpbA | 2lkgA | 2w1pA | 3bryA | 3k3fA | 3sybA | 4c9jA | 4hyoA | 4ogqC | 4us3A | 5azbA | |
| 1p4tA | 2d57A | 2llyA | 2wjqa | 3chxB | 3k3jA | 3tjja | 4cadC | 4hzuS | 4oh3A | 4v1fA | 5bwkE | |
| 1p7bA | 2ervA | 2lmeA | 2wpdJ | 3chxC | 3kp9A | 3tx3A | 4cfgA | 4iffA | 4oo9A | 4wd7A | 5c6oA | |

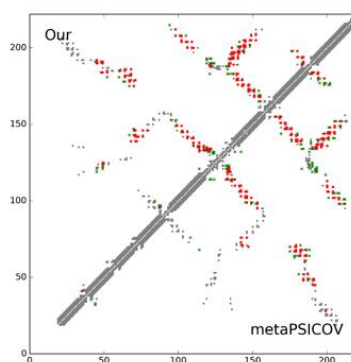
Figure S1. Case study of one CAMEO target 5jkiA, related to section “Blind test in CAMEO” and Figure 4. (A) The long- and medium-range contact prediction accuracy of our methods, MetaPSICOV, CCMpred, and EVfold (Web Server). (B-D) The overlap between top L predicted all-range contacts and the native contact map. A grey, red and green dot represents a native contact, a correct prediction and a wrong prediction, respectively. (E) The superimposition between our predicted model (in red) and the native structure (in green).

| | Long range accuracy | | | | Medium range accuracy | | | |
|-------------|---------------------|-------|-------|-------|-----------------------|-------|-------|-------|
| | L | L/2 | L/5 | L/10 | L | L/2 | L/5 | L/10 |
| Our method | 0.658 | 0.883 | 1.000 | 1.000 | 0.185 | 0.351 | 0.659 | 0.864 |
| MetaPSICOV | 0.554 | 0.820 | 0.977 | 1.000 | 0.158 | 0.279 | 0.523 | 0.727 |
| CCMpred | 0.495 | 0.703 | 0.773 | 0.818 | 0.131 | 0.207 | 0.477 | 0.682 |
| EVfold(web) | 0.514 | 0.712 | 0.773 | 0.841 | 0.126 | 0.207 | 0.432 | 0.727 |

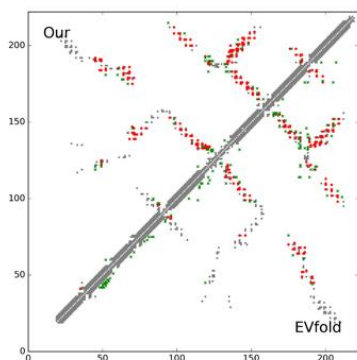
(A)



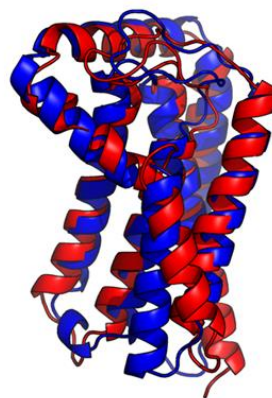
(B)



(C)



(D)

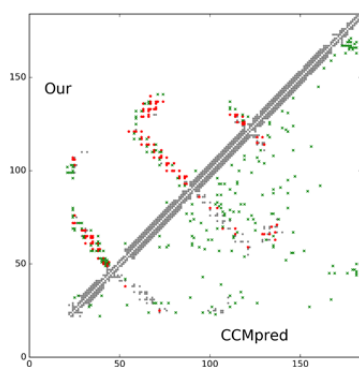


(E)

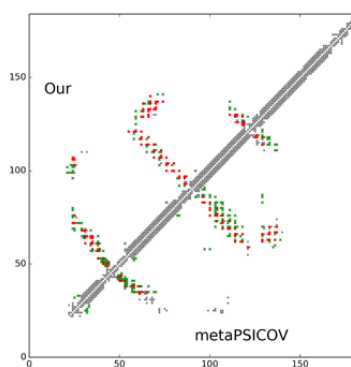
Figure S2. Case study of one CAMEO target 510wA, related to section “*Blind test in CAMEO*” and *Figure 4*. (A) The long- and medium-range contact prediction accuracy of our methods, MetaPSICOV and CCMpred. (B-C) The overlap between top L predicted all-range contacts and the native contact map. A grey, red and green dot represents a native contact, a correct prediction and a wrong prediction, respectively. (D) The superimposition between our predicted model (in red) and the native structure (in green).

| | Long range accuracy | | | | Medium range accuracy | | | |
|------------|---------------------|-------|-------|-------|-----------------------|-------|-------|-------|
| | L | L/2 | L/5 | L/10 | L | L/2 | L/5 | L/10 |
| Our method | 0.397 | 0.674 | 0.889 | 1.000 | 0.103 | 0.207 | 0.444 | 0.778 |
| metaPSICOV | 0.250 | 0.391 | 0.528 | 0.722 | 0.098 | 0.163 | 0.278 | 0.389 |
| CCMpred | 0.087 | 0.109 | 0.222 | 0.333 | 0.016 | 0.033 | 0.056 | 0.056 |

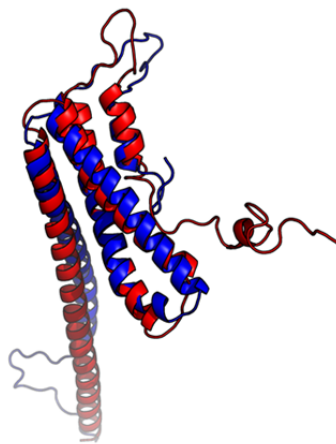
(A)



(B)



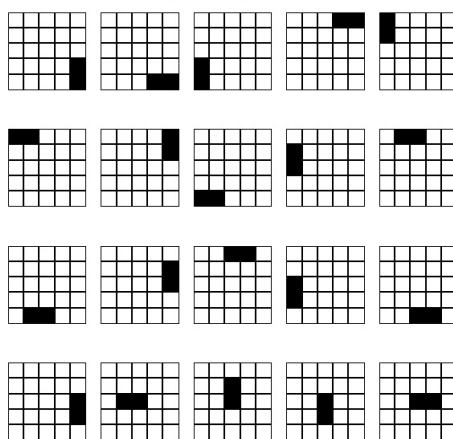
(C)



(D)

Supplementary Figure 3. Top 20 long-range 5x5 contact occurrence patterns, related to section “*Why does deep transfer learning work?*” and STAR method.

9627 non-MPs long-rang top-20 patterns with m contacts, $m \geq 2$



345 multi-pass MPs long-rang top-20 patterns with m contacts, $m \geq 2$

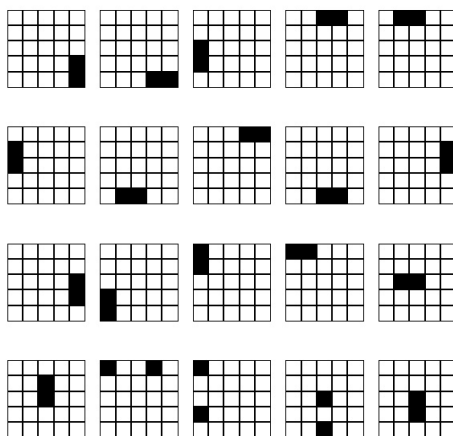


Table S2. Contact prediction accuracy on the 87 COMSAT test proteins, *related to STAR Method and Table I*. Only contacts between two transmembrane segments are evaluated. Acc, Cov, Sp and Mcc represent accuracy, coverage, specificity and Mathew correlation coefficient, respectively. Note that here our result is based upon the deep model trained without using any membrane proteins while COMSAT was trained by some membrane proteins. The result of COMSAT is taken from its paper. Following the COMSAT paper, the results of the latter three methods are calculated on top L_m predicted contacts where L_m is the length of transmembrane regions in a test protein.

| Accuracy of top L_m predicted contacts when a contact is defined by C_α - C_α distance less than 14Å. | | | | | | | | | | | | |
|---|----------|------|------|------|-----------|-------|------|------|-----------|-------|------|------|
| Method | ≥ 6 | | | | ≥ 12 | | | | ≥ 24 | | | |
| | Acc | Cov | Sp | Mcc | Acc | Cov | Sp | Mcc | Acc | Cov | Sp | Mcc |
| CCMpred | 0.63 | 0.07 | 0.98 | 0.15 | 0.61 | 0.07 | 0.98 | 0.15 | 0.57 | 0.08 | 0.97 | 0.14 |
| MetaPSICOV | 0.73 | 0.08 | 0.99 | 0.19 | 0.72 | 0.08 | 0.99 | 0.19 | 0.69 | 0.10 | 0.98 | 0.19 |
| Our Method | 0.86 | 0.10 | 0.99 | 0.23 | 0.85 | 0.10 | 0.99 | 0.23 | 0.82 | 0.12 | 0.98 | 0.24 |
| COMSAT | 0.65 | 0.05 | 0.99 | 0.11 | 0.63 | 0.054 | 0.99 | 0.11 | 0.61 | 0.052 | 0.99 | 0.10 |
| Accuracy of top L_m predicted contacts when a contact is defined by C_β - C_β distance less than 8Å. | | | | | | | | | | | | |
| Method | ≥ 6 | | | | ≥ 12 | | | | ≥ 24 | | | |
| | Acc | Cov | Sp | Mcc | Acc | Cov | Sp | Mcc | Acc | Cov | Sp | Mcc |
| CCMpred | 0.32 | 0.25 | 0.98 | 0.26 | 0.31 | 0.26 | 0.97 | 0.26 | 0.29 | 0.27 | 0.96 | 0.25 |
| MetaPSICOV | 0.39 | 0.31 | 0.98 | 0.32 | 0.38 | 0.31 | 0.98 | 0.32 | 0.35 | 0.33 | 0.96 | 0.31 |
| Our Method | 0.59 | 0.46 | 0.98 | 0.48 | 0.58 | 0.46 | 0.98 | 0.48 | 0.53 | 0.49 | 0.97 | 0.47 |
| COMSAT | 0.43 | 0.14 | 0.99 | 0.21 | 0.43 | 0.14 | 0.99 | 0.21 | 0.44 | 0.14 | 0.98 | 0.21 |