

# **Additional file 1: Fig. S1-S5**

## **Systematic Prediction of Degrons and E3 Ubiquitin Ligase Binding via Deep Learning**

Chao Hou<sup>1,2</sup>, Yuxuan Li<sup>1,2</sup>, Mengyao Wang<sup>3,4</sup>, Hong Wu<sup>3,4,5</sup>, Tingting Li<sup>1,2,\*</sup>

<sup>1</sup>Department of Biomedical Informatics, School of Basic Medical Sciences, Peking University Health Science Center, Beijing 100191, China.

<sup>2</sup>Key Laboratory for Neuroscience, Ministry of Education/National Health Commission of China, Peking University, Beijing 100191, China.

<sup>3</sup>The MOE Key Laboratory of Cell Proliferation and Differentiation, School of Life Sciences, Peking University, Beijing 100871, China.

<sup>4</sup>Peking-Tsinghua Center for Life Sciences, Beijing, China.

<sup>5</sup>Peking University Institute of Hematology, National Clinical Research Center for Hematologic Disease, Beijing, China.

\*Corresponding Author: [litt@hsc.pku.edu.cn](mailto:litt@hsc.pku.edu.cn).

## Supplementary Figure legends

**Fig. S1 Architecture of the model and comparison with other predictors.** **a** Distribution of length of degrons and expanded degrons. **b** Detailed architecture of the TAPE BERT-based degron prediction model. **c** Performances of the BERT-based model, Motif\_RF and MoRFchibi in predicting degrons from motif matches, curves of the BERT-based model, Motif\_RF were averages in five-fold cross-validation, MoRFchibi was directly evaluated on all positive and negative motif matches. **d** Averaged receiver operating characteristic curves of the BERT-based and one-hot models. **e** Evaluation scores of the BERT-based model and one-hot model in predicting five clusters of degrons, AUC and recall were measured at positive to negative ratio equal to that in the training set, negative samples were AAs randomly selected from proteins possessing test degrons. **f** Averaged precision-recall curves of BERT-based models at different positive to negative ratios. **g** Detailed architecture of the one-hot degron prediction model. **f** Averaged precision-recall curves of one-hot models at different positive to negative ratios.

**Fig. S2 Cut-off determination, comparison with GPS experiment and PTMs in degrons.** **a** FDR and recall of Degpred at different cut-offs with positive: negative=1:20, negative samples were AAs randomly selected from proteins possessing degrons. **b** Distribution of terminus located known and Degpred degrons in GPS N-end and C-end experiments. Known and Degpred degrons located within 23 AAs distant from N-end and C-end were compared, consistent with the peptide length used in GPS experiments. The smaller the number of the bin, the stronger the ability of peptides of the bin to induce reporter protein degradation. **c** Enrichment of PTM sites in Degpred degrons, the bars show the fold change of ratios of modified AAs in Degpred degrons and random peptides from the human proteome. P-values were calculated using Fischer's exact test.

**Fig. S3 Statistics of collected and predicted ESI datasets and motifs for HECT E3s.** **a** Number of substrates for each E3 (left) and number of E3s for each substrate (right) in our collected ESI dataset. **b** Seqlogos of generated motifs of four HECT E3s: WWP1, WWP2, SMURF2, NEDD4L. **c** Precision and recall rates for 55 E3s. Precisions were calculated at positive: negative=1:20. Recall rates measured the percentage of known substrates that can be predicted using calculated motifs. The raw numbers of this figure were provided in Supplementary Table 3. **d** Comparison of ChenESI and generated motifs in predicting substrates of SPOP and FZR1 in manually collected ESIs of Ubibrowser2.0. ChenESI predicted 1186 SPOP substrates and 2020 FZR1 substrates, we selected top 890 SPOP substrates and top 1179 FZR1 substrates of Degpred to get the same FDR with ChenESI. **e** Number of substrates for each E3 (upper) and number of E3s for each substrate (lower) in the predicted ESI dataset.

**Fig. S4 Half-lives of proteins with different disorder fractions and degron density.** Distribution of protein half-lives in four non-dividing cell types, proteins were divided into three groups with disorder fractions 0-10%, 10-30%, 30-100%, numbers of proteins were shown on the right. The

half-lives used were replicate 1 of four cell types in original paper, and the results of replicate 2 were identical (data not shown). Kernel density estimate plots showing the distribution of  $\log_{10}(\text{half-life} + 1)$ .

**Fig. S5 Functional analysis of degron-related mutations and E3s.** **a** Percentage of Degpred degon AAs in driver regions and other regions. P-value was calculated using Fischer's exact test. **b** Distribution of  $-\log(\text{P-value})$  calculated by SMDeg and FMDEg of motif matches overlapped with Degpred degons and not overlapped. SMDeg probes the over-representation of missense mutations in degons with respect to the number inferred from the distribution of all missense mutations observed in the protein. FMDEg computes the deviation of the average functional impact of missense mutations in degons from the expected impact. P-values were calculated using two side T-test. **c** Enriched function groups of short-lived proteins. **d** Percentage of driver mutations in degon-related regions and other regions of short-lived proteins and the other proteins in four non-dividing cell types. P-values were calculated using Fischer's exact test. **e** Average number of TCGA mutations occurring in degon-related regions bound by different E3s. **f** Percentage of driver mutations in degon-related regions bound by different E3s. **g** Percentage of mutations that alter the charge, hydrophobicity, phosphorylation sites, MoRF regions, predicted protein binding residues or lysine residues.

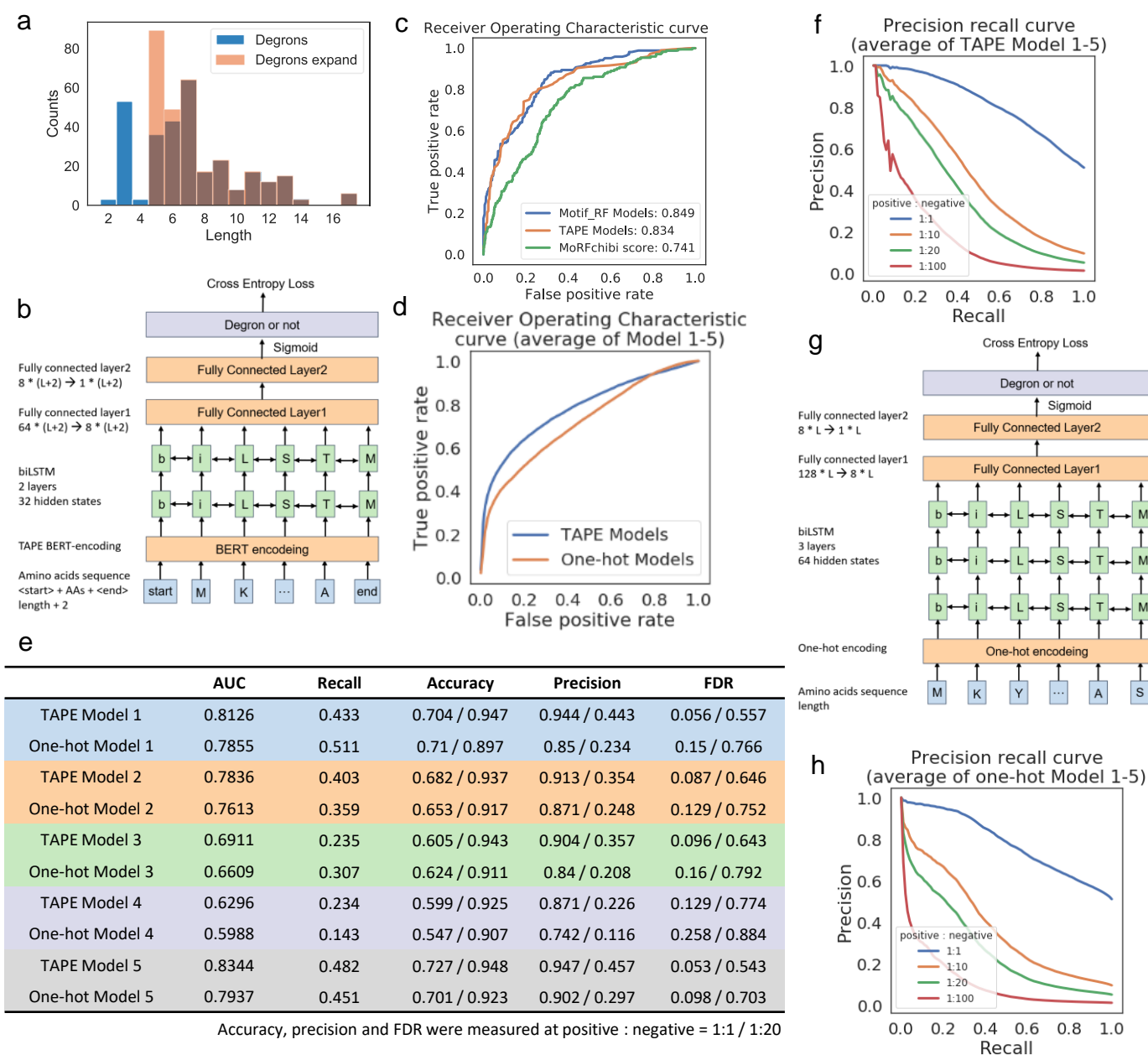


Fig. S1

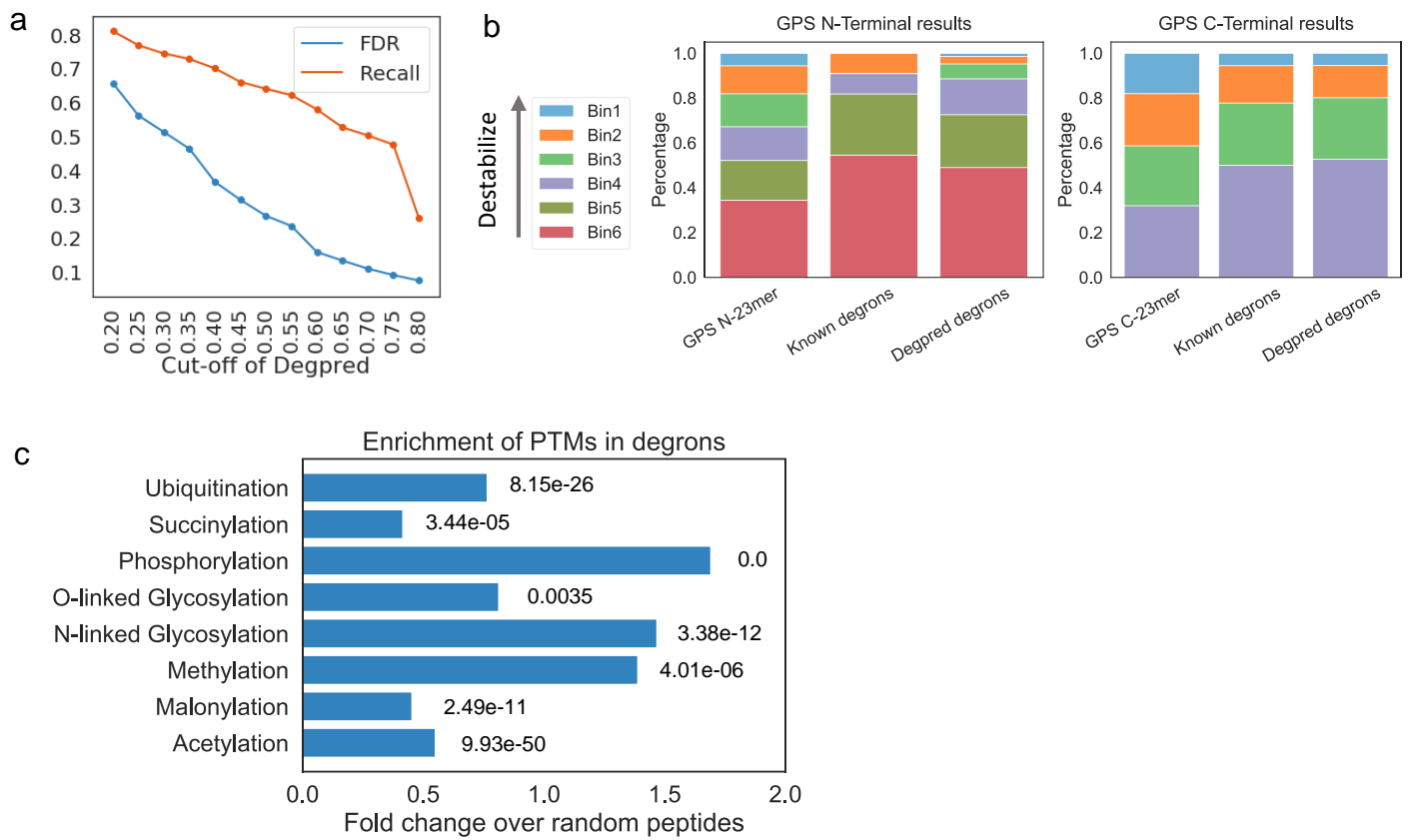


Fig. S2

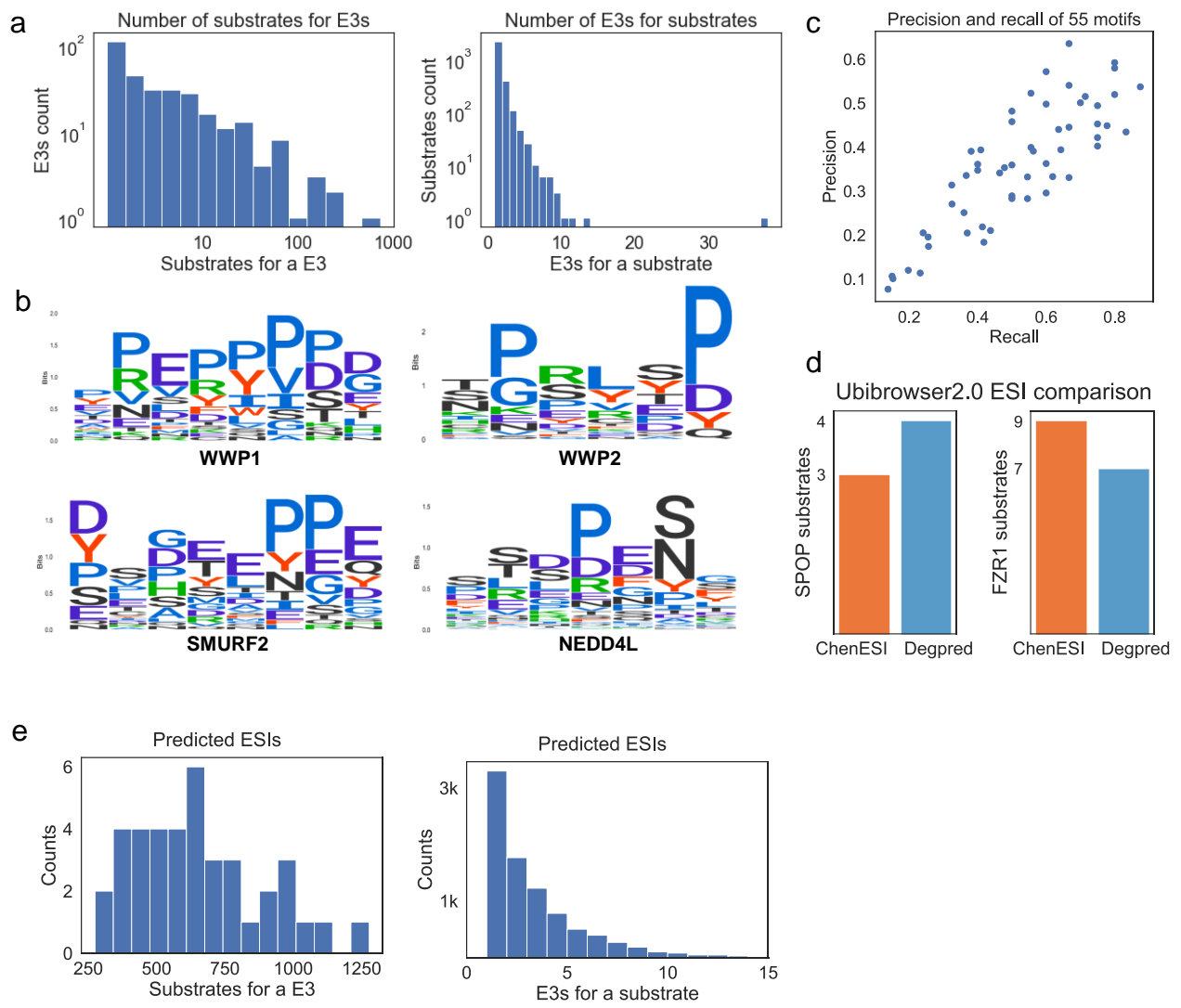


Fig. S3

Log(half-life) for proteins with different density of degrons

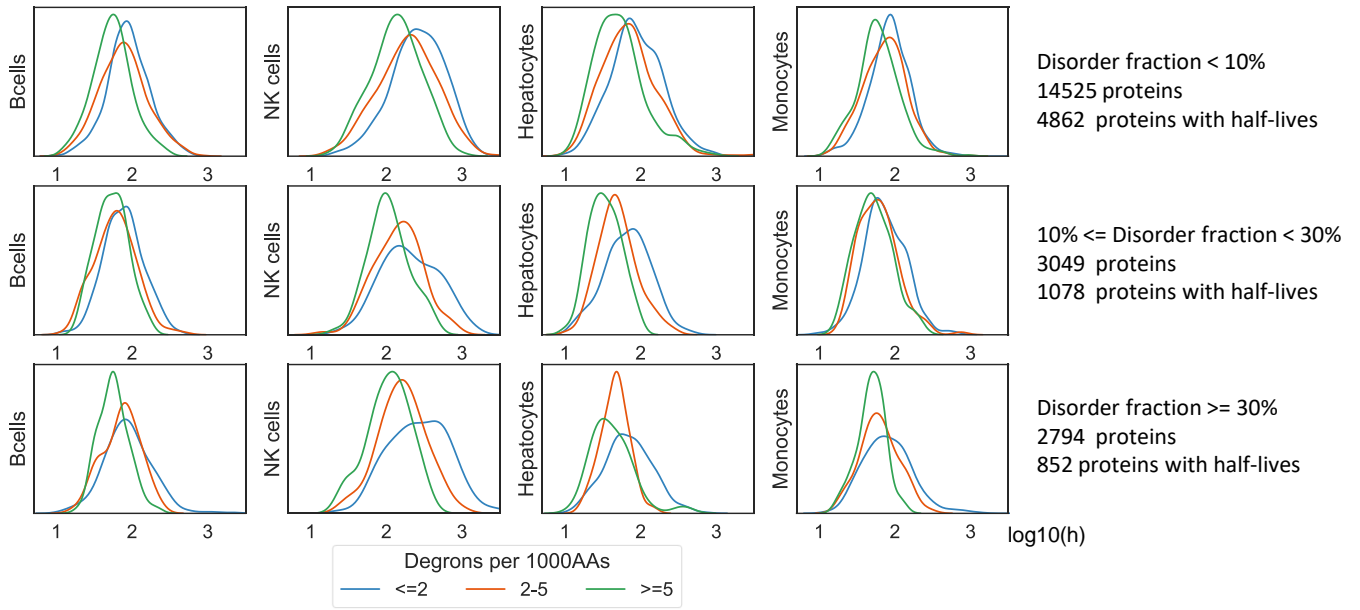


Fig. S4

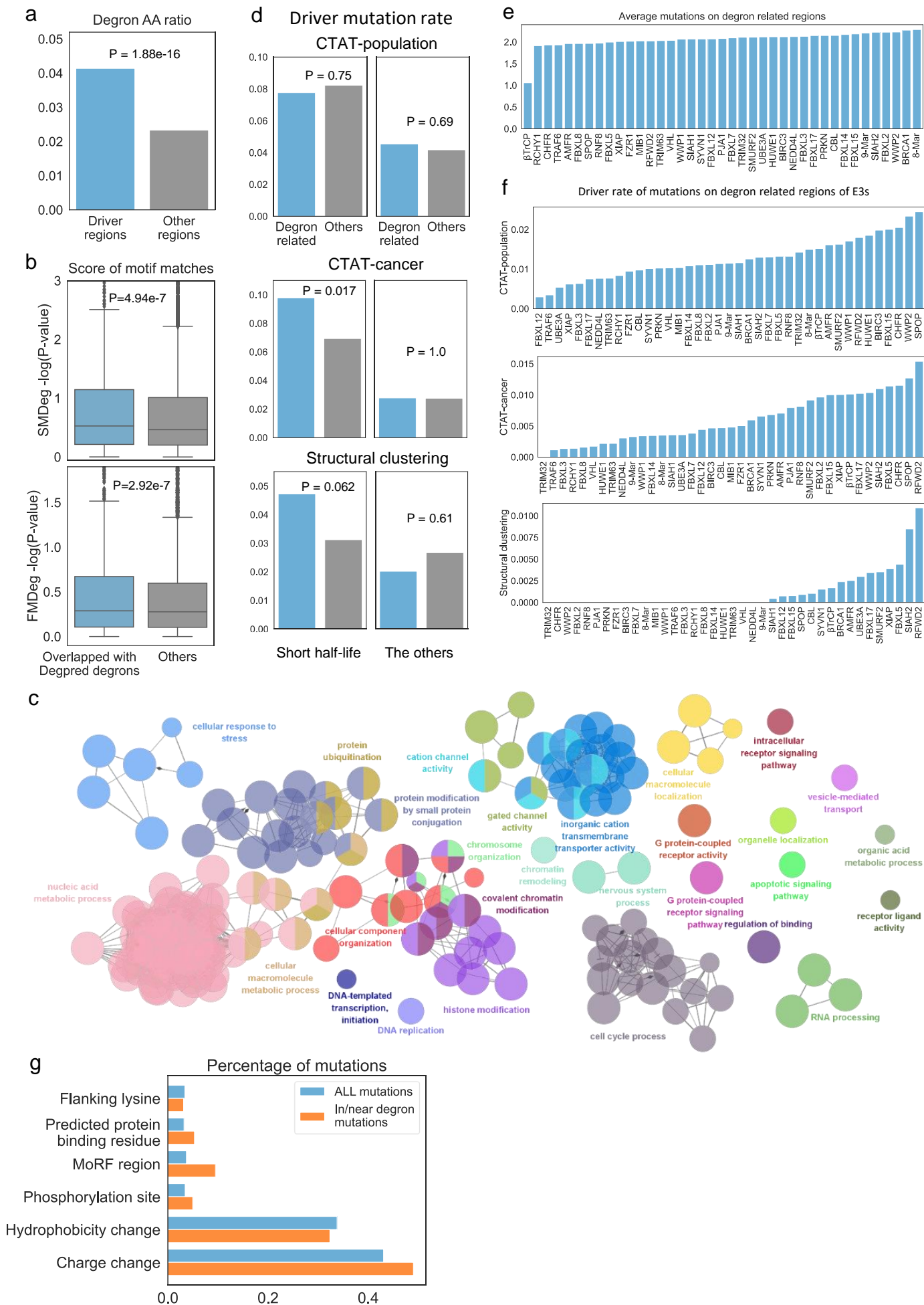


Fig. S5