

Prediction of the Effect of Single Amino Acid Protein Variants Using Deep Mutational Scanning Data

University of Bologna — Master Thesis in Bioinformatics

Pierotti Saul

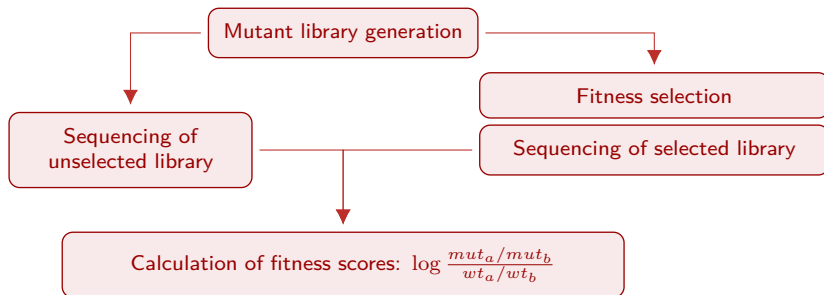
Internal Advisor: Prof. Pietro Di Lena

External Advisor: Prof. Arne Elofsson (Stockholm University)

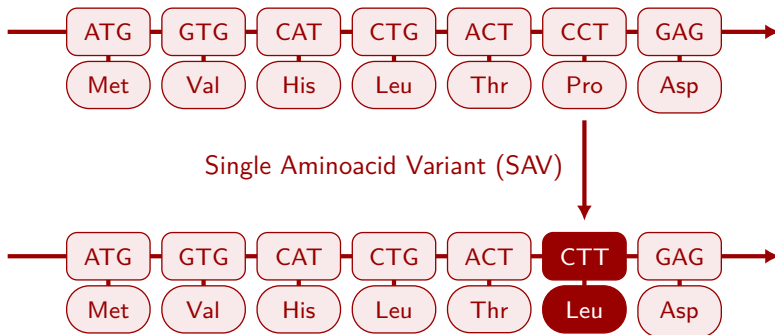
July 19, 2021

Deep Mutational Scanning

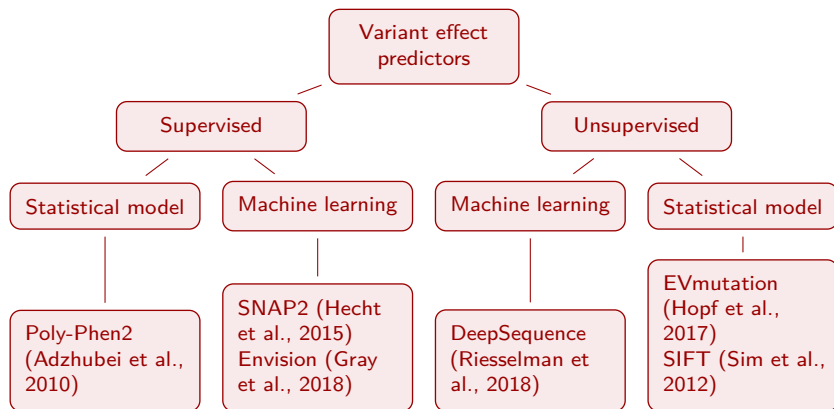
High-throughput technique for obtaining fitness information on a large number of mutations



I Considered Only Single Amino Acid Variants

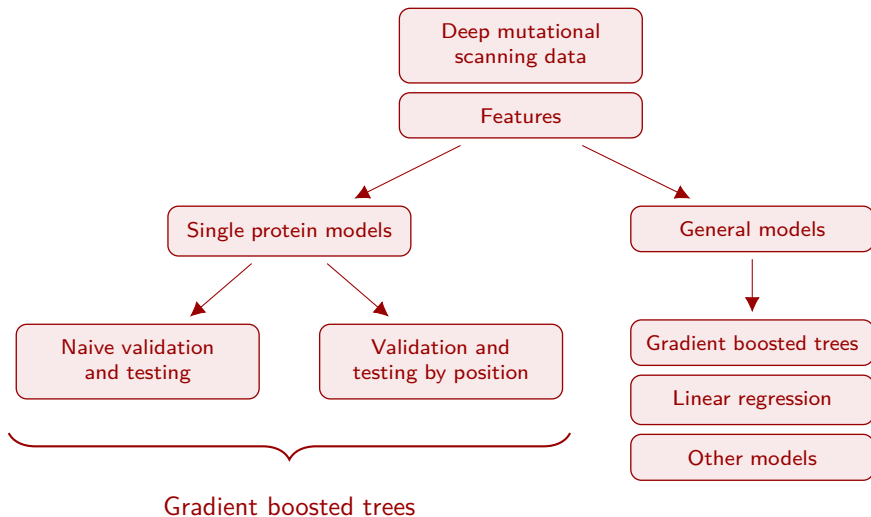


Some Notable Variant Effect Predictors

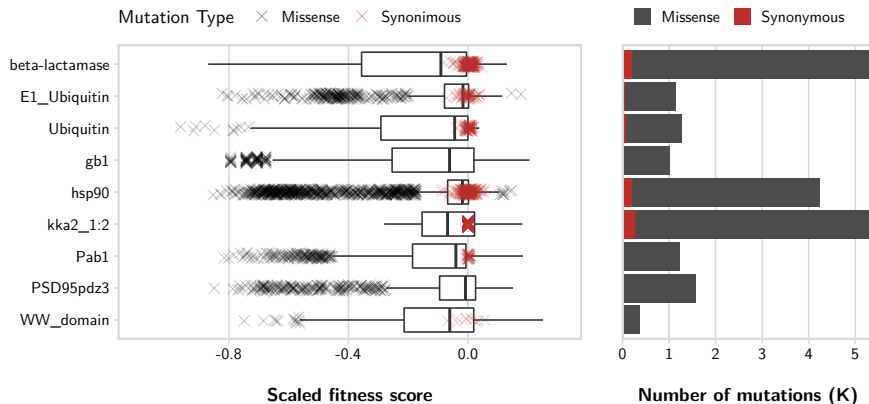


Among these predictors only Envision was trained on deep mutational scanning data

Structure of the Project

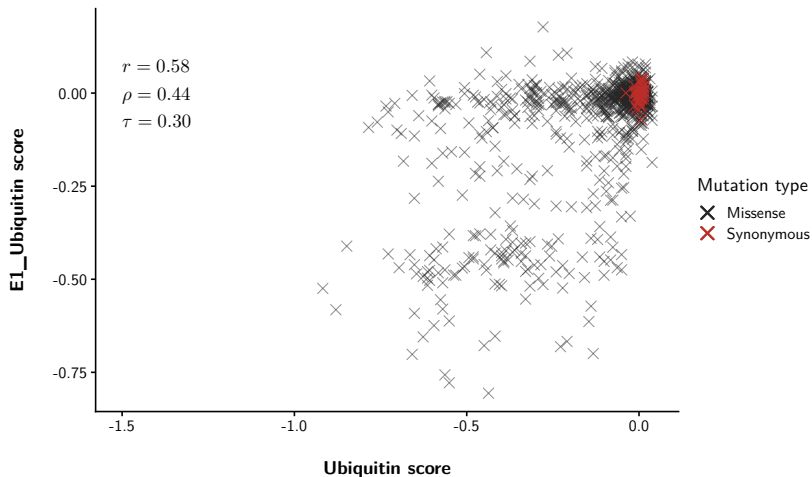


I Used the Training Dataset of Envision (Gray et al., 2018)



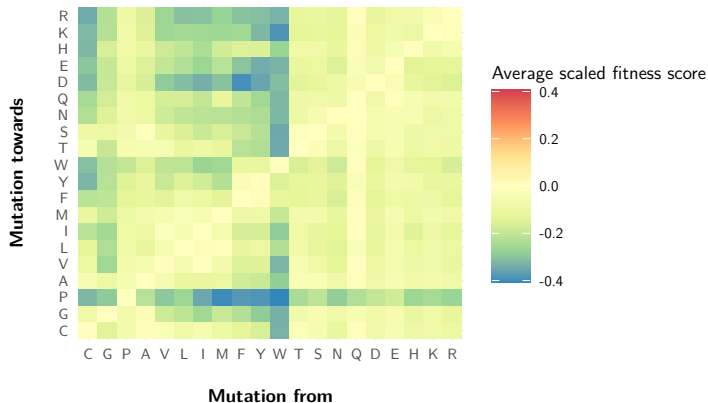
Poor Correlation among Experimental Results

Two independent deep mutational scanning experiments on Ubiquitin are present in the training dataset. Their correlation is low.

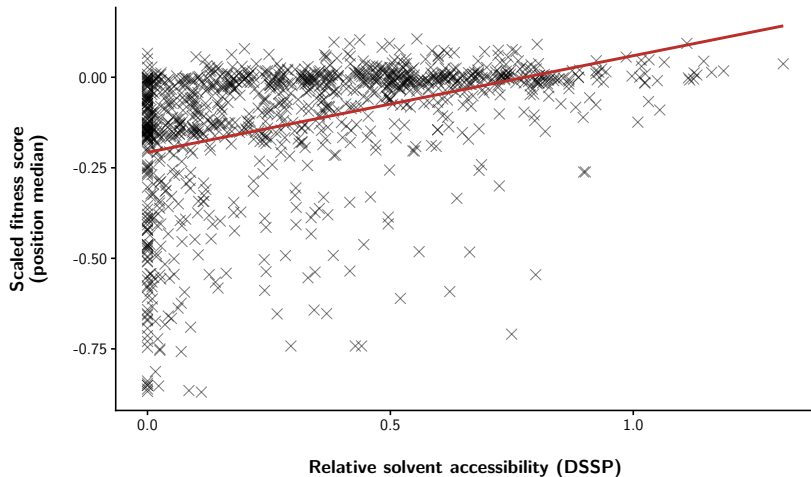


The Identities of the Mutated Residues Are Important

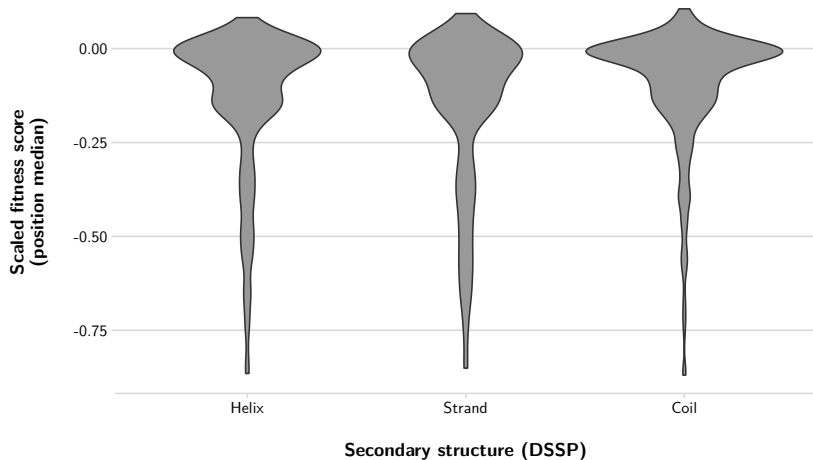
Mutations from polar residues are less detrimental than other mutations. This effect disappears when filtering by relative solvent accessibility (not shown here).



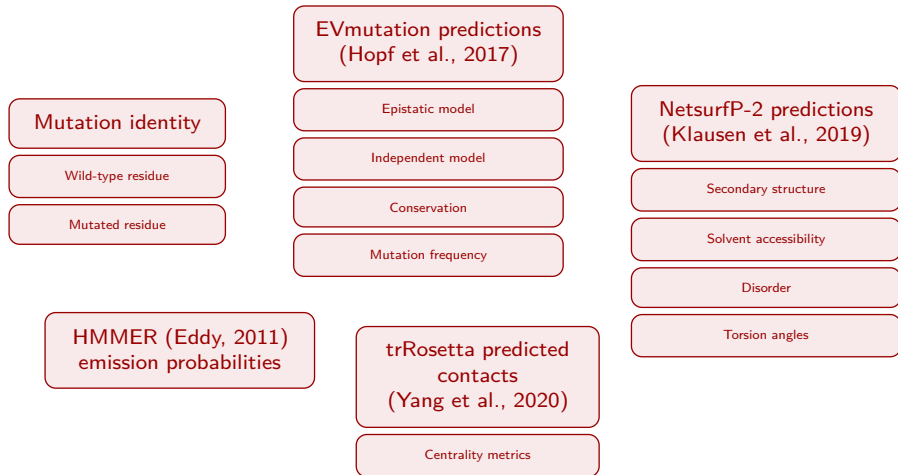
Buried Residues Are More Conserved



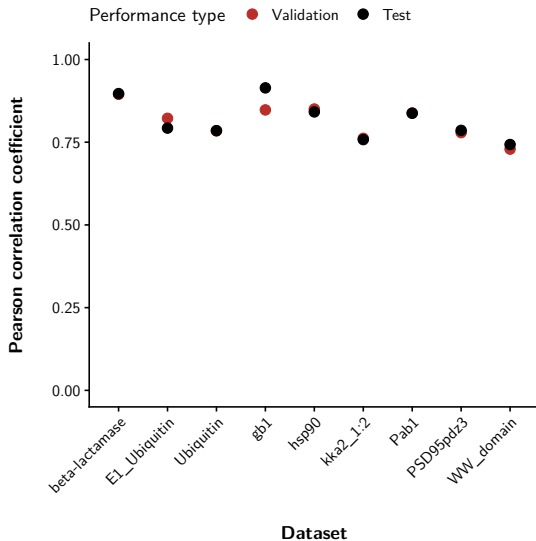
Secondary Structure Is Not Critical



I Did Not Use Structural Information

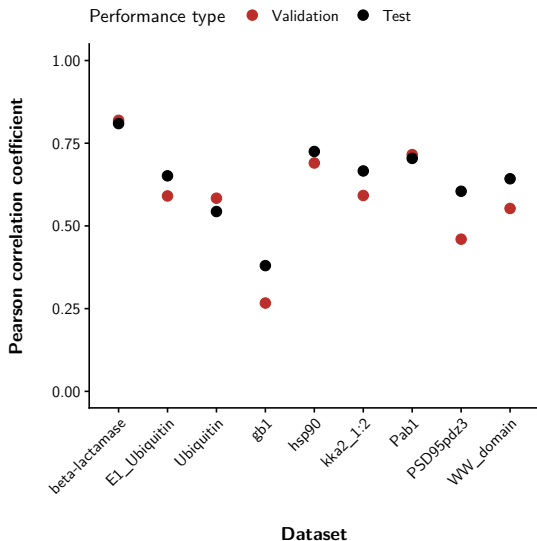


Single Protein Models with Naive Testing



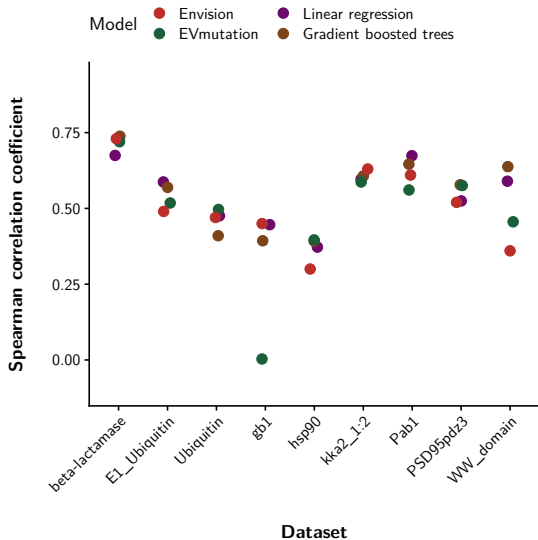
- ▶ Half of the mutations in a protein randomly set aside for testing
- ▶ Hyperparameters optimized in the remaining half with 5-fold cross-validation
- ▶ Good results but likely overfitting the testing set

Single Protein Models with Testing by Position



- ▶ Half of the mutations in a protein set aside for testing but avoiding mutations in the same protein position to end up in different splits
- ▶ Hyperparameters optimized in the remaining half with 5-fold cross-validation
- ▶ Performance more realistic

Leave-One-Protein-Out (LOPO) Models



- ▶ For the left-out protein, half of the mutations used for testing
- ▶ Spearman correlation coefficient used for evaluation
- ▶ Small difference among gradient boosted trees and linear regression
- ▶ Performances comparable to those of Envision

Discussion

Complex models do not improve much on linear regression

Unsupervised models perform similarly to supervised models

There is strong variability between datasets

How validation and testing are performed is crucial

Performances on par with other predictors can be reached without structural features

Future Directions

Unsupervised models seem promising
and may be worth exploring more

Training on more deep
mutational scanning studies

Tuning the
set of features

Using residue contacts in a graph
convolutional neural network

Finding a better normalization strategy
for the scores from different experiments

Bibliography I

- Adzhubei, I. A., Schmidt, S., Peshkin, L., Ramensky, V. E., Gerasimova, A., Bork, P., Kondrashov, A. S. & Sunyaev, S. R. (2010). A method and server for predicting damaging missense mutations. *Nature Methods*, 7(4), 248–249. <https://doi.org/10.1038/nmeth0410-248>
- Eddy, S. R. (2011). Accelerated profile HMM searches (W. R. Pearson, Ed.). *PLoS Computational Biology*, 7(10), e1002195. <https://doi.org/10.1371/journal.pcbi.1002195>
- Gray, V. E., Hause, R. J., Luebeck, J., Shendure, J. & Fowler, D. M. (2018). Quantitative missense variant effect prediction using large-scale mutagenesis data. *Cell Systems*, 6(1), 116–124.e3. <https://doi.org/10.1016/j.cels.2017.11.003>
- Hecht, M., Bromberg, Y. & Rost, B. (2015). Better prediction of functional effects for sequence variants. *BMC Genomics*, 16(S8). <https://doi.org/10.1186/1471-2164-16-s8-s1>
- Hopf, T., Ingraham, J., Poelwijk, F., Schärfe, C., Springer, M., Sander, C. & Marks, D. (2017). Mutation effects predicted from sequence co-variation. *Nature Biotechnology*, 35(2), 128–135. <https://doi.org/10.1038/nbt.3769>

Bibliography II

- Klausen, M. S., Jespersen, M. C., Nielsen, H., Jensen, K. K., Jurtz, V. I., Søndersby, C. K., Sommer, M. O. A., Winther, O., Nielsen, M., Petersen, B. & Marcatili, P. (2019). NetSurfP-2.0: Improved prediction of protein structural features by integrated deep learning. *Proteins: Structure, Function, and Bioinformatics*, 87(6), 520–527. <https://doi.org/10.1002/prot.25674>
- Riesselman, A. J., Ingraham, J. B. & Marks, D. S. (2018). Deep generative models of genetic variation capture the effects of mutations. *Nature Methods*, 15(10), 816–822. <https://doi.org/10.1038/s41592-018-0138-4>
- Sim, N.-I., Kumar, P., Hu, J., Henikoff, S., Schneider, G. & Ng, P. C. (2012). SIFT web server: Predicting effects of amino acid substitutions on proteins. *Nucleic Acids Research*, 40(W1), W452–W457. <https://doi.org/10.1093/nar/gks539>
- Yang, J., Anishchenko, I., Park, H., Peng, Z., Ovchinnikov, S. & Baker, D. (2020). Improved protein structure prediction using predicted interresidue orientations. *Proceedings of the National Academy of Sciences*, 117(3), 1496–1503. <https://doi.org/10.1073/pnas.1914677117>

Supplementary Material

Confidence intervals in prediction performances

Dataset name	Model	95 % C.I. (Pearson)	95 % C.I. (Spearman)
beta-lactamase	Naive	0.89 to 0.91	0.86 to 0.88
beta-lactamase	By position	0.79 to 0.83	0.75 to 0.78
beta-lactamase	LOPO	—	0.68 to 0.72
WW_domain	Naive	0.67 to 0.82	0.70 to 0.84
WW_domain	By position	0.57 to 0.73	0.52 to 0.72
WW_domain	LOPO	—	0.56 to 0.73
PSD95pdz3	Naive	0.74 to 0.83	0.73 to 0.80
PSD95pdz3	By position	0.55 to 0.67	0.54 to 0.65
PSD95pdz3	LOPO	—	0.53 to 0.63
kka2_1:2	Naive	0.74 to 0.78	0.72 to 0.76
kka2_1:2	By position	0.65 to 0.69	0.64 to 0.68
kka2_1:2	LOPO	—	0.60 to 0.64
hsp90	Naive	0.82 to 0.87	0.65 to 0.71
hsp90	By position	0.69 to 0.76	0.50 to 0.57
hsp90	LOPO	—	0.38 to 0.45
Ubiquitin	Naive	0.75 to 0.83	0.78 to 0.83
Ubiquitin	By position	0.49 to 0.60	0.46 to 0.57
Ubiquitin	LOPO	—	0.30 to 0.43
Pab1	Naive	0.80 to 0.87	0.75 to 0.82
Pab1	By position	0.65 to 0.76	0.64 to 0.73
Pab1	LOPO	—	0.60 to 0.70
E1_Ubiquitin	Naive	0.75 to 0.85	0.65 to 0.75
E1_Ubiquitin	By position	0.59 to 0.72	0.36 to 0.51
E1_Ubiquitin	LOPO	—	0.50 to 0.63
gb1	Naive	0.90 to 0.93	0.88 to 0.92
gb1	By position	0.31 to 0.46	0.29 to 0.45
gb1	LOPO	—	0.32 to 0.47

Supplementary Material

Statistical significance of performance differences. Bonferroni-corrected $\alpha = \frac{0.05}{27} = 0.00185185$. Starred values are significant.

Dataset name	Model 1	Model 2	<i>p</i> value
beta-lactamase	Linear regression	Gradient boosted trees	$1 \cdot 10^{-4}$ *
beta-lactamase	Linear regression	EVmutation	$1 \cdot 10^{-4}$ *
beta-lactamase	EVmutation	Gradient boosted trees	$1 \cdot 10^{-4}$ *
WW_domain	Linear regression	Gradient boosted trees	$1 \cdot 10^{-4}$ *
WW_domain	Linear regression	EVmutation	$1 \cdot 10^{-4}$ *
WW_domain	EVmutation	Gradient boosted trees	$1 \cdot 10^{-4}$ *
PSD95pdz3	Linear regression	Gradient boosted trees	$1 \cdot 10^{-4}$ *
PSD95pdz3	Linear regression	EVmutation	$1 \cdot 10^{-4}$ *
PSD95pdz3	EVmutation	Gradient boosted trees	0.51
kka2_1:2	Linear regression	Gradient boosted trees	0.00
kka2_1:2	Linear regression	EVmutation	0.01
kka2_1:2	EVmutation	Gradient boosted trees	$1 \cdot 10^{-4}$ *
hsp90	Linear regression	Gradient boosted trees	$1 \cdot 10^{-4}$ *
hsp90	Linear regression	EVmutation	$1 \cdot 10^{-4}$ *
hsp90	EVmutation	Gradient boosted trees	0.24

The table continues on the next slide.

Supplementary Material

Statistical significance of performance differences. Bonferroni-corrected $\alpha = \frac{0.05}{27} = 0.00185185$. Starred values are significant.

The table continues from the previous slide.

Dataset name	Model 1	Model 2	<i>p</i> value
Ubiquitin	Linear regression	Gradient boosted trees	$1 \cdot 10^{-4}$ *
Ubiquitin	Linear regression	EVmutation	$1 \cdot 10^{-4}$ *
Ubiquitin	EVmutation	Gradient boosted trees	$1 \cdot 10^{-4}$ *
Pab1	Linear regression	Gradient boosted trees	$1 \cdot 10^{-4}$ *
Pab1	Linear regression	EVmutation	$1 \cdot 10^{-4}$ *
Pab1	EVmutation	Gradient boosted trees	$1 \cdot 10^{-4}$ *
E1_Ubiquitin	Linear regression	Gradient boosted trees	$1 \cdot 10^{-4}$ *
E1_Ubiquitin	Linear regression	EVmutation	$1 \cdot 10^{-4}$ *
E1_Ubiquitin	EVmutation	Gradient boosted trees	$1 \cdot 10^{-4}$ *
gb1	Linear regression	Gradient boosted trees	$1 \cdot 10^{-4}$ *
gb1	Linear regression	EVmutation	$1 \cdot 10^{-4}$ *
gb1	EVmutation	Gradient boosted trees	$1 \cdot 10^{-4}$ *