# Learning protein fitness models from evolutionary and assay-labeled data

(CSE 590C WI 22 - Alyssa La Fleur)

UNIVERSITY *of* WASHINGTON

**W**

# Protein fitness prediction

> **Protein "fitness":** any protein property (stability, enzyme activity, binding strength, etc.)
> **Predicting fitness for protein sequences:** assist with design, potential pathogenicity prediction
  – Pathogenicity prediction task != Fitness prediction task

# Protein fitness prediction

> **Protein "fitness":** any protein property (stability, enzyme activity, binding strength, etc.)
> **Predicting fitness for protein sequences:** assist with design, potential pathogenicity prediction
>    –   Pathogenicity prediction task != Fitness prediction task

**This paper evaluates existing fitness prediction methods, and proposes a new one**

W

# Two main ML strategies

1. **Evolutionary models**
   a. Get a sequence alignment for your target protein
   b. Model the probability density of these sequences
   c. Predict mutant fitness using the probability density model

# Two main ML strategies

1.  **Evolutionary models**
    a.  Get a sequence alignment for your target protein
    b.  Model the probability density of these sequences
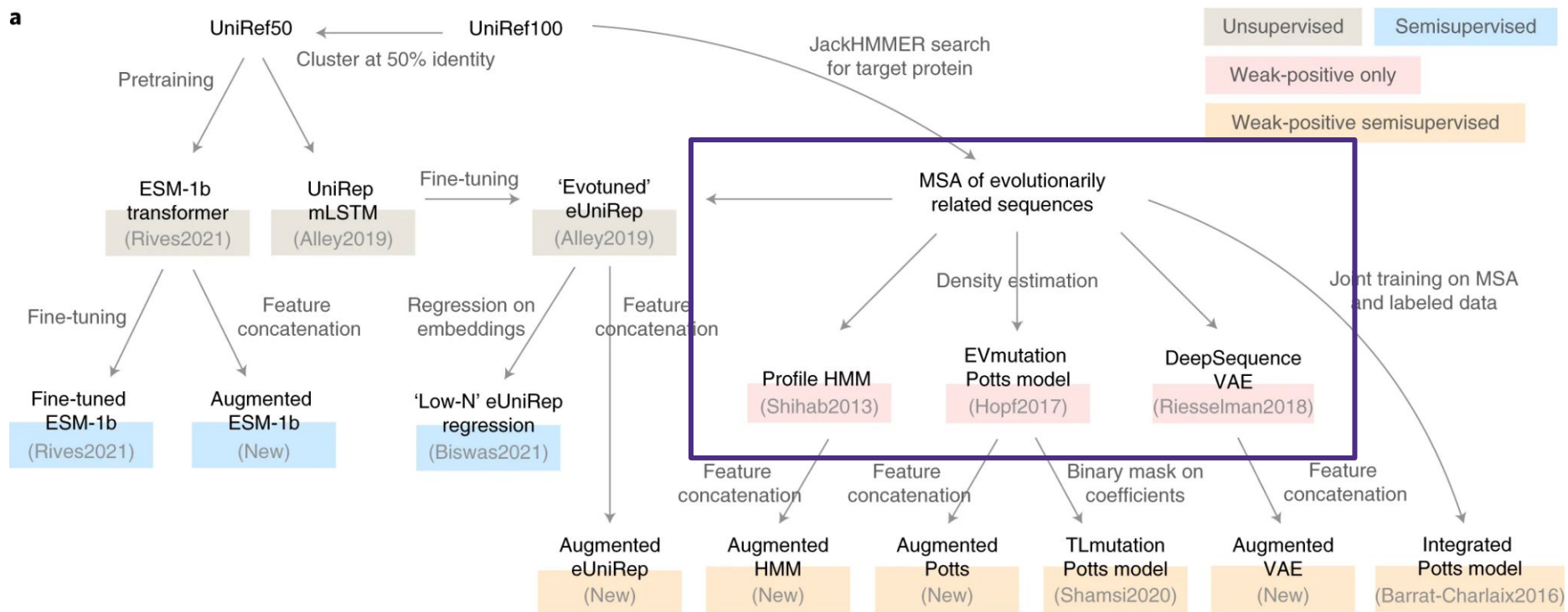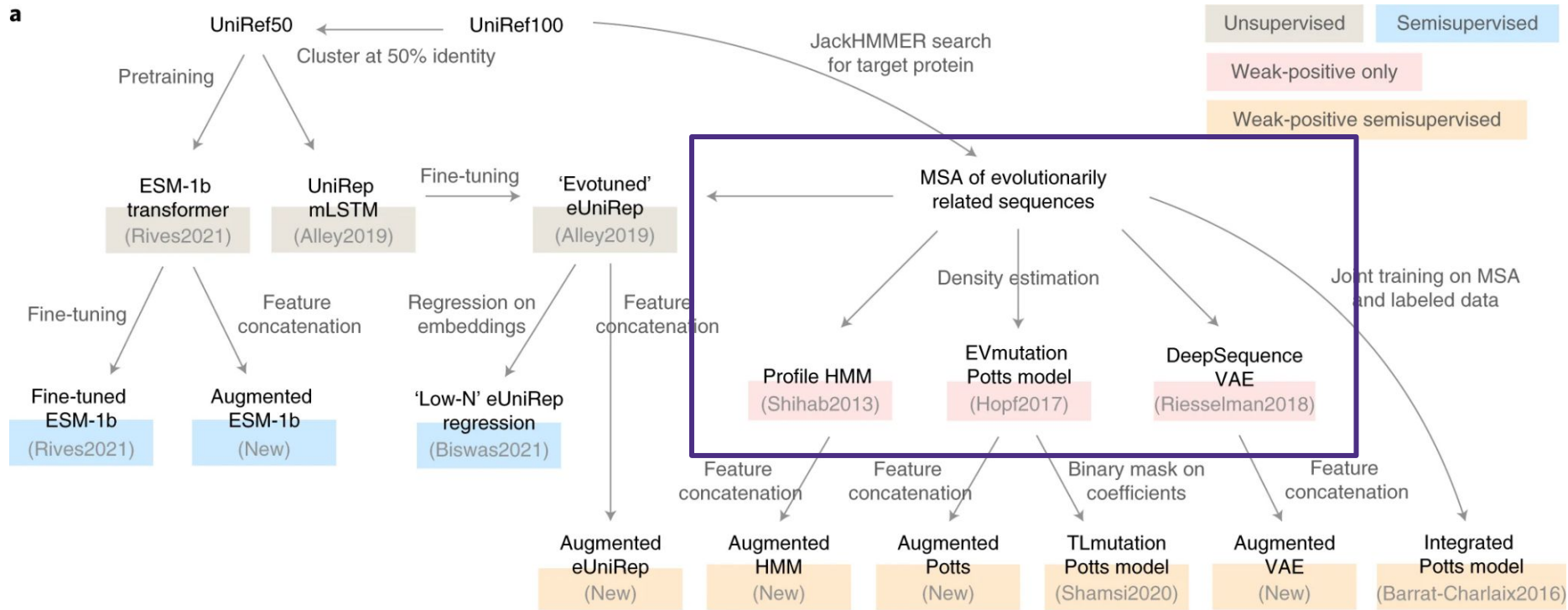    c.  Predict mutant fitness using the probability density model

*'Weak positive' learning - these approaches assume that evolutionary related sequences have similar functions to the target*

**W**

**a**

Unsupervised | Semisupervised

Weak-positive only

Weak-positive semisupervised

UniRef50 ← UniRef100
Cluster at 50% identity

JackHMMER search for target protein

Pretraining

ESM-1b transformer (Rives2021)

UniRep mLSTM (Alley2019)

Fine-tuning

'Evotuned' eUniRep (Alley2019)

MSA of evolutionarily related sequences

Density estimation

Joint training on MSA and labeled data

Fine-tuning

Feature concatenation

Regression on embeddings

Feature concatenation

Fine-tuned ESM-1b (Rives2021)

Augmented ESM-1b (New)

'Low-N' eUniRep regression (Biswas2021)

Profile HMM (Shihab2013)

EVmutation Potts model (Hopf2017)

DeepSequence VAE (Riesselman2018)

Feature concatenation

Feature concatenation

Binary mask on coefficients

Feature concatenation

Augmented eUniRep (New)

Augmented HMM (New)

Augmented Potts (New)

TLmutation Potts model (Shamsi2020)

Augmented VAE (New)

Integrated Potts model (Barrat-Charlaix2016)

**One limitation is that alignment depth may vary - some targets may only have hundreds of usable sequences in their alignment**
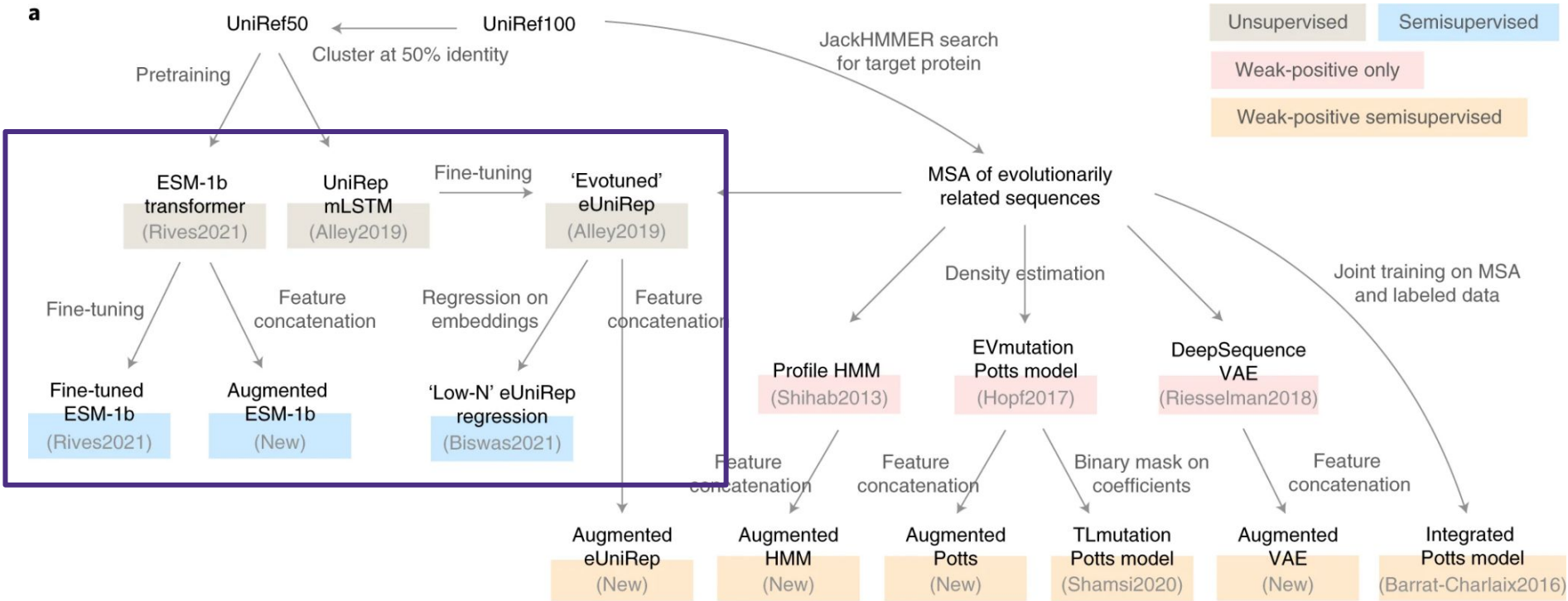
# Two main ML strategies

2. **Supervised regression models**
   a. Models range from simple (linear regression) to complex (CNN, LSTM, Transformers, etc.)
   b. **Semi-supervised:** Supervised regression models can also be trained using unsupervised NLP model protein representations

**W**

**Can be limited by number of mutants in training set, coverage of positions by mutation (few positions vs. many)**

# A combined strategy

> **Weak-positive semi-supervised learning:** learning a distribution of sequences using alignments, with supervised learning on labelled sequences
> **Their 'baseline' augmentation combined approach (Had max performance in 15/19 test sets)**

W

**Can be limited by number of mutants in training set, coverage of positions by mutation (few positions vs. many)**
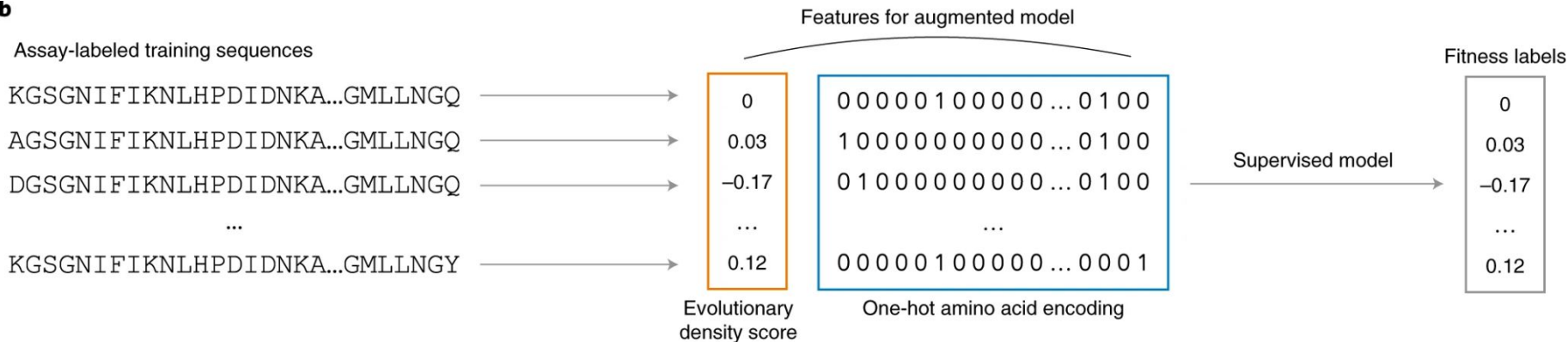
# Potts models

> **TLMutation Potts models SUMMARY GOES HERE**

W

# Augmentation combined approach

> **Sequence log-likelihoods from a sequence density model & one-hot encoded protein sequences**
> **Supervised model is ridge regression (L2)**



**b**

Assay-labeled training sequences

KGSGNIFIKNLHPDIDNKA...GMLLNGQ → 0

AGSGNIFIKNLHPDIDNKA...GMLLNGQ → 0.03

DGSGNIFIKNLHPDIDNKA...GMLLNGQ → −0.17

... 

KGSGNIFIKNLHPDIDNKA...GMLLNGY → 0.12

Features for augmented model

Evolutionary density score:
0
0.03
−0.17
…
0.12

One-hot amino acid encoding:
0 0 0 0 0 1 0 0 0 0 0 ... 0 1 0 0
1 0 0 0 0 0 0 0 0 0 0 ... 0 1 0 0
0 1 0 0 0 0 0 0 0 0 0 ... 0 1 0 0
…
0 0 0 0 0 1 0 0 0 0 0 ... 0 0 0 1

Supervised model →

Fitness labels:
0
0.03
−0.17
…
0.12

# Deep mutation scanning (DMS) datasets

> They used 19 of the DMS datasets from EVMutation (one of the competitor models they compared against) + a GFP fluorescence data set
> All had mutations throughout a domain or whole protein
> 16/19 had sequences one missense mutation away from WT (single mutants)
> Only evaluated mutants at positions with < 30% gaps in the MSAs generated

W

# Dataset splits

> **20% test**, varying sizes of training sets
> **80/20 train/test**, five fold cross-validation on the 80% where computationally feasible (For comparing to methods like TLMutation, ESM1-b)

**20 random seems were used for data partitioning for each approach**

W

# Ranking metrics

> **Spearman rank correlation coefficient:**

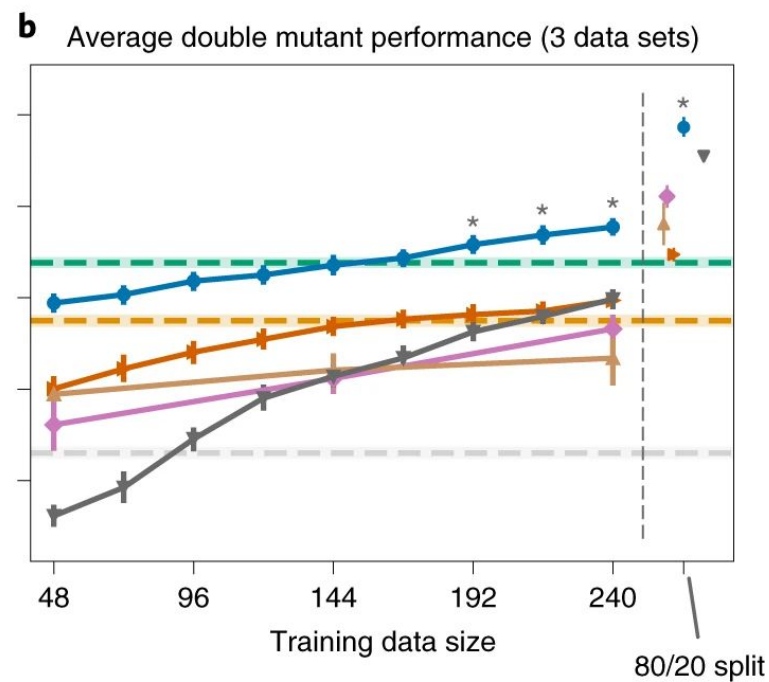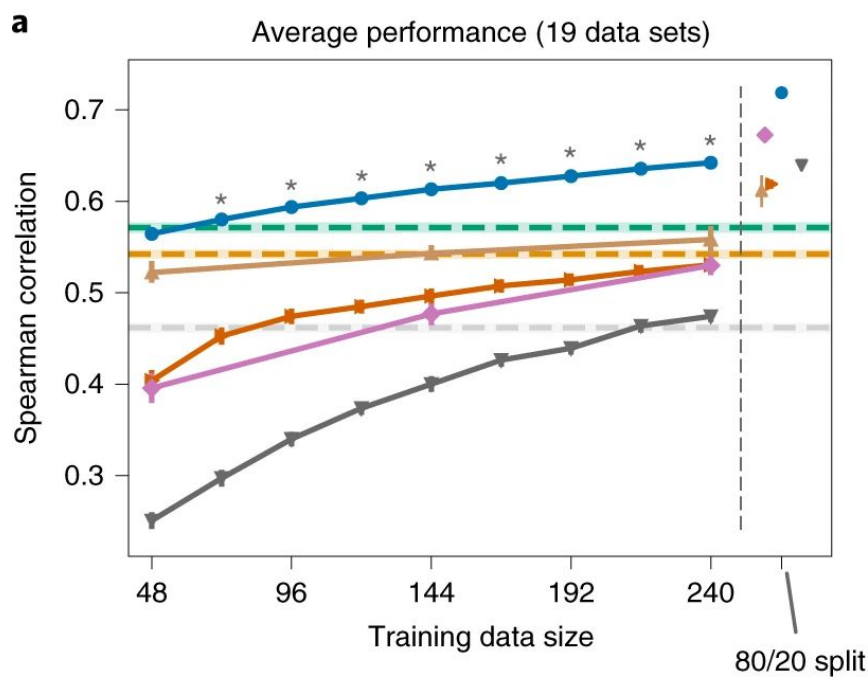> **Normalized discounted cumulative gain (NDCG):** From information retrieval, similar to a weighted Spearman rank which focuses on high value agreement

**W**

# One hot linear model

# Low-N Training Predictions



**Legend:**
- Augmented EVmutation Potts (MSA + labeled)
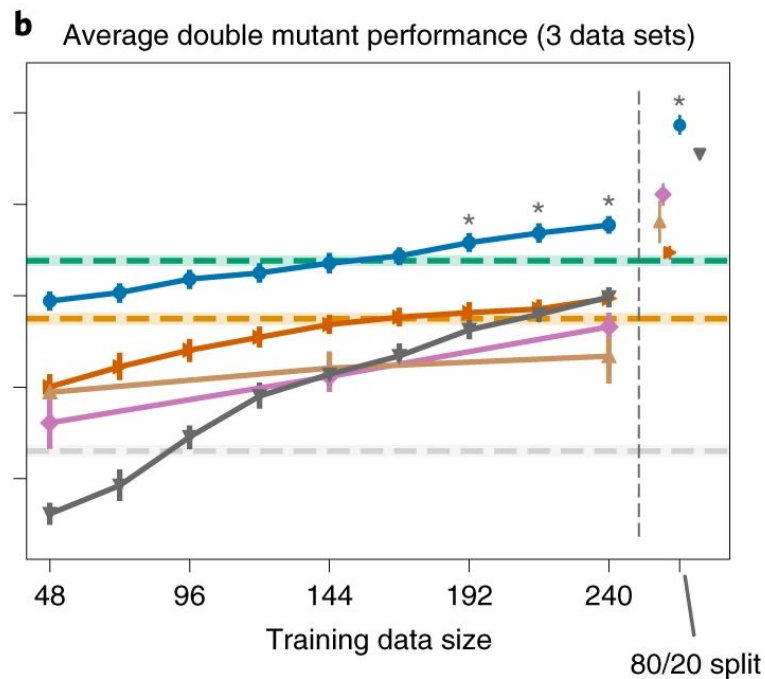- Profile HMM (MSA only)
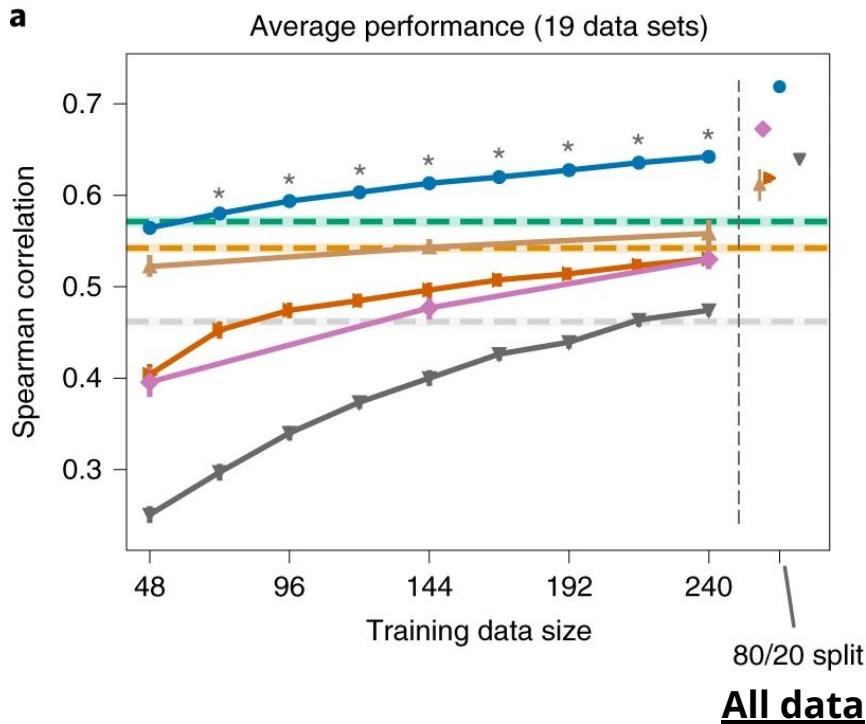- EVmutation Potts (Hopf2017) (MSA only)
- DeepSequence VAE (Riesselman2018) (MSA only)
- eUniRep regression (Biswas2021) (UniRef50 + MSA + labeled)
- Fine-tuned transformer (Rives2021) (UniRef50 + labeled, no MSA)
- Integrated Potts (Barrat-Charlaix2016) (MSA + labeled)
- Linear model with one-hot encoding (labeled only)

**a** Average performance (19 data sets)

**b** Average double mutant performance (3 data sets)

Axes: Spearman correlation vs Training data size (48, 96, 144, 192, 240), with 80/20 split

# Low-N Training Predictions



Legend:
- Augmented EVmutation Potts (MSA + labeled)
- Profile HMM (MSA only)
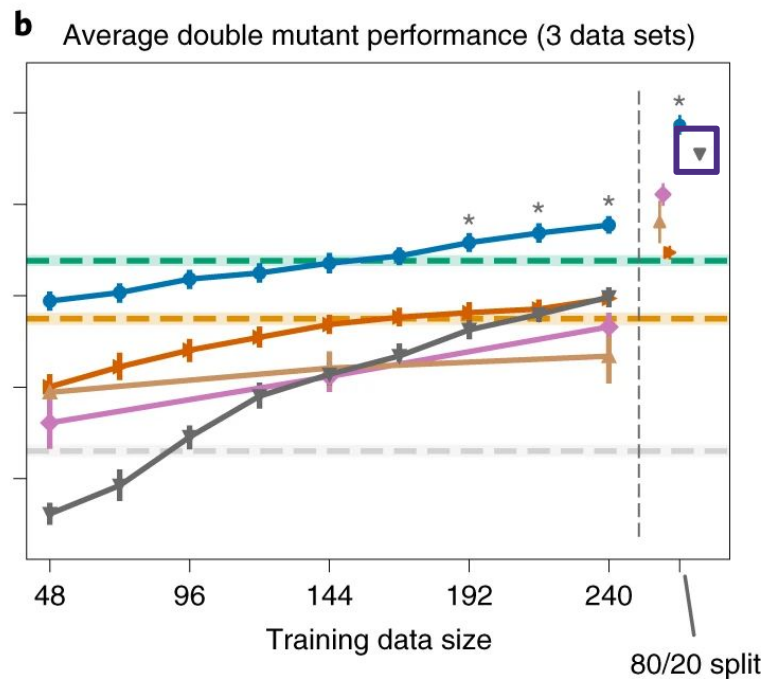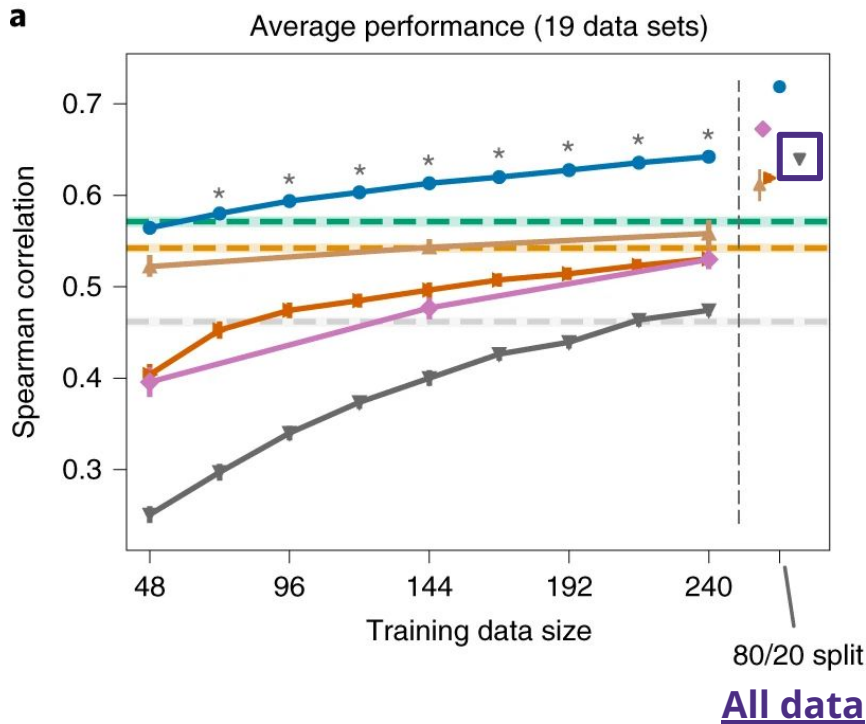- EVmutation Potts (Hopf2017) (MSA only)
- DeepSequence VAE (Riesselman2018) (MSA only)
- eUniRep regression (Biswas2021) (UniRef50 + MSA + labeled)
- Fine-tuned transformer (Rives2021) (UniRef50 + labeled, no MSA)
- Integrated Potts (Barrat-Charlaix2016) (MSA + labeled)
- Linear model with one-hot encoding (labeled only)

**a** Average performance (19 data sets)

**b** Average double mutant performance (3 data sets)

**All data**

# Low-N Training Predictions



Legend:
- Augmented EVmutation Potts (MSA + labeled)
- Profile HMM (MSA only)
- EVmutation Potts (Hopf2017) (MSA only)
- DeepSequence VAE (Riesselman2018) (MSA only)
- eUniRep regression (Biswas2021) (UniRef50 + MSA + labeled)
- Fine-tuned transformer (Rives2021) (UniRef50 + labeled, no MSA)
- Integrated Potts (Barrat-Charlaix2016) (MSA + labeled)
- Linear model with one-hot encoding (labeled only)

**a** Average performance (19 data sets)

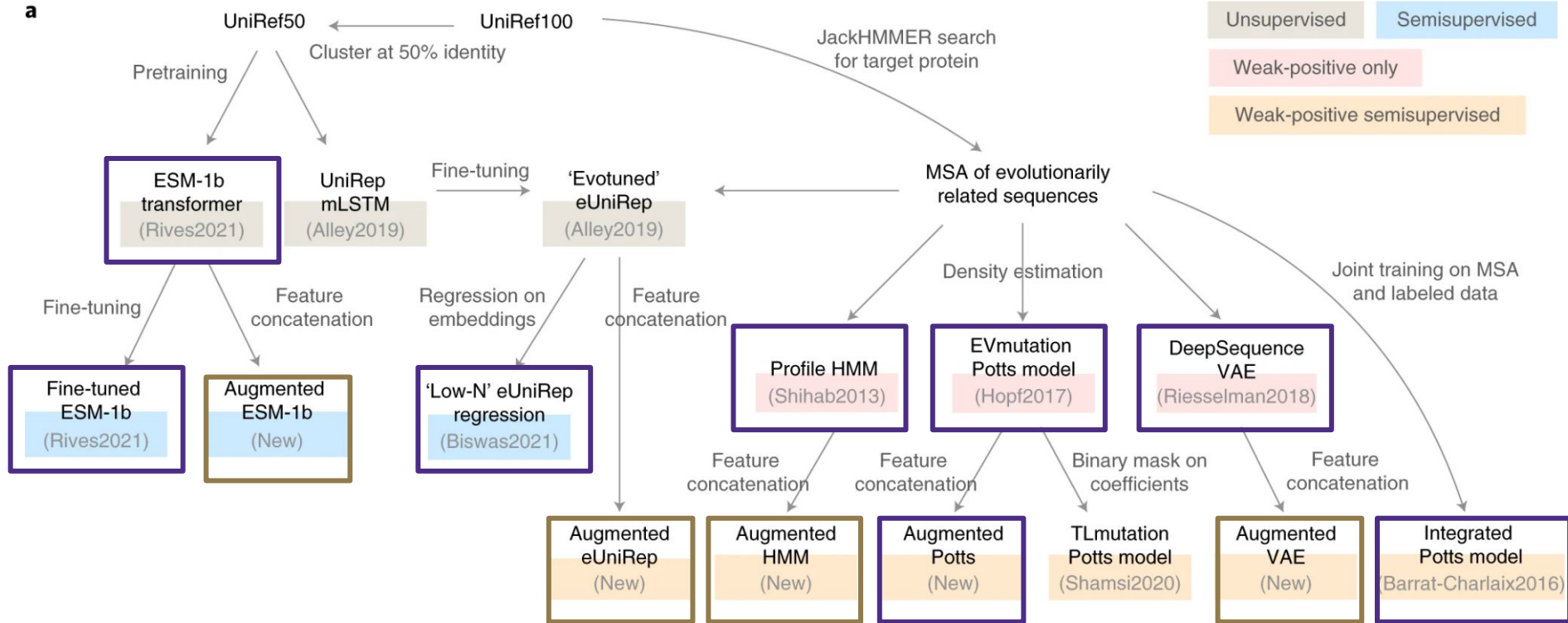**b** Average double mutant performance (3 data sets)
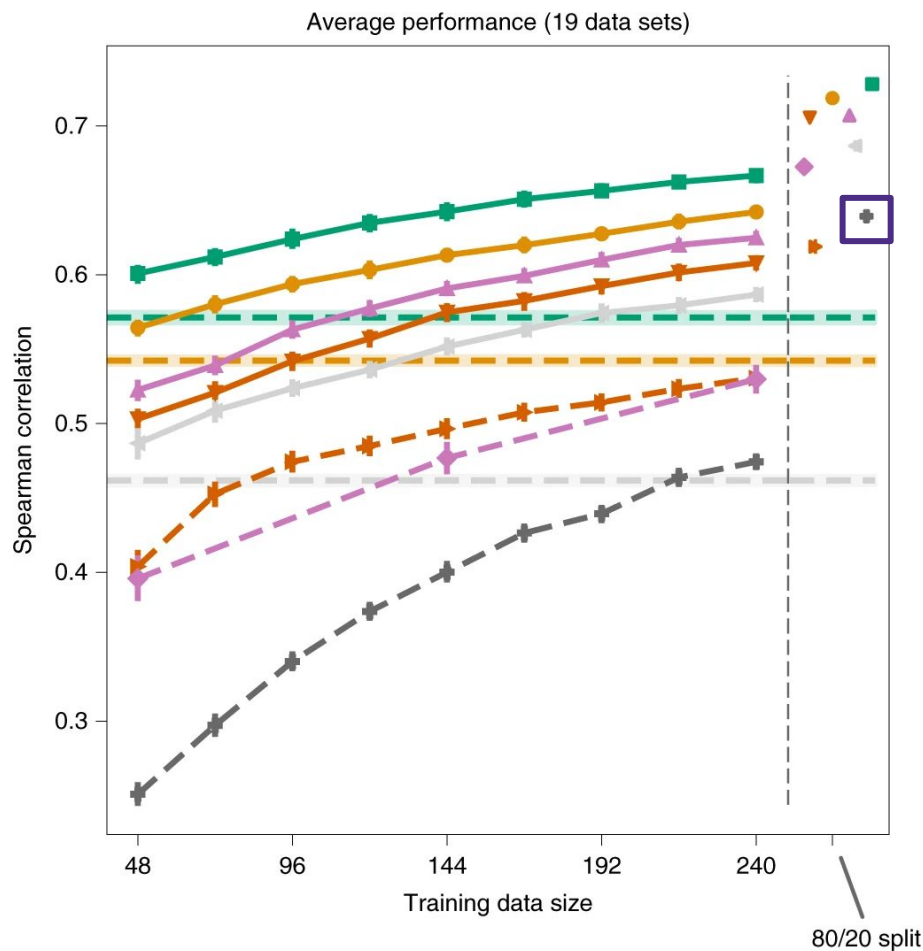
**All data**

# Additional augmented models

> Augmented other models than just Potts
> Note that the transformer (not eUniRep) is the only method not using any evolutionary data
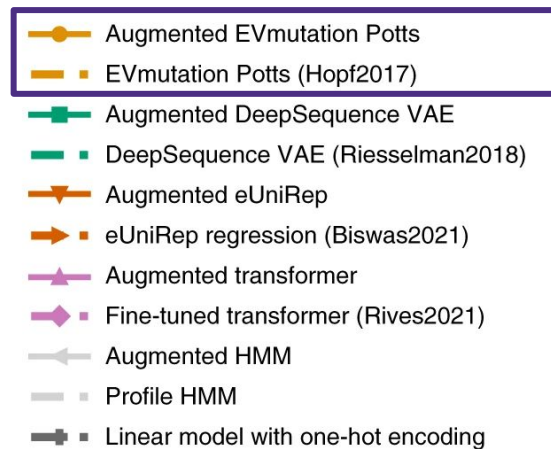> Augmented models outperformed non-augmented model, regardless of training set size
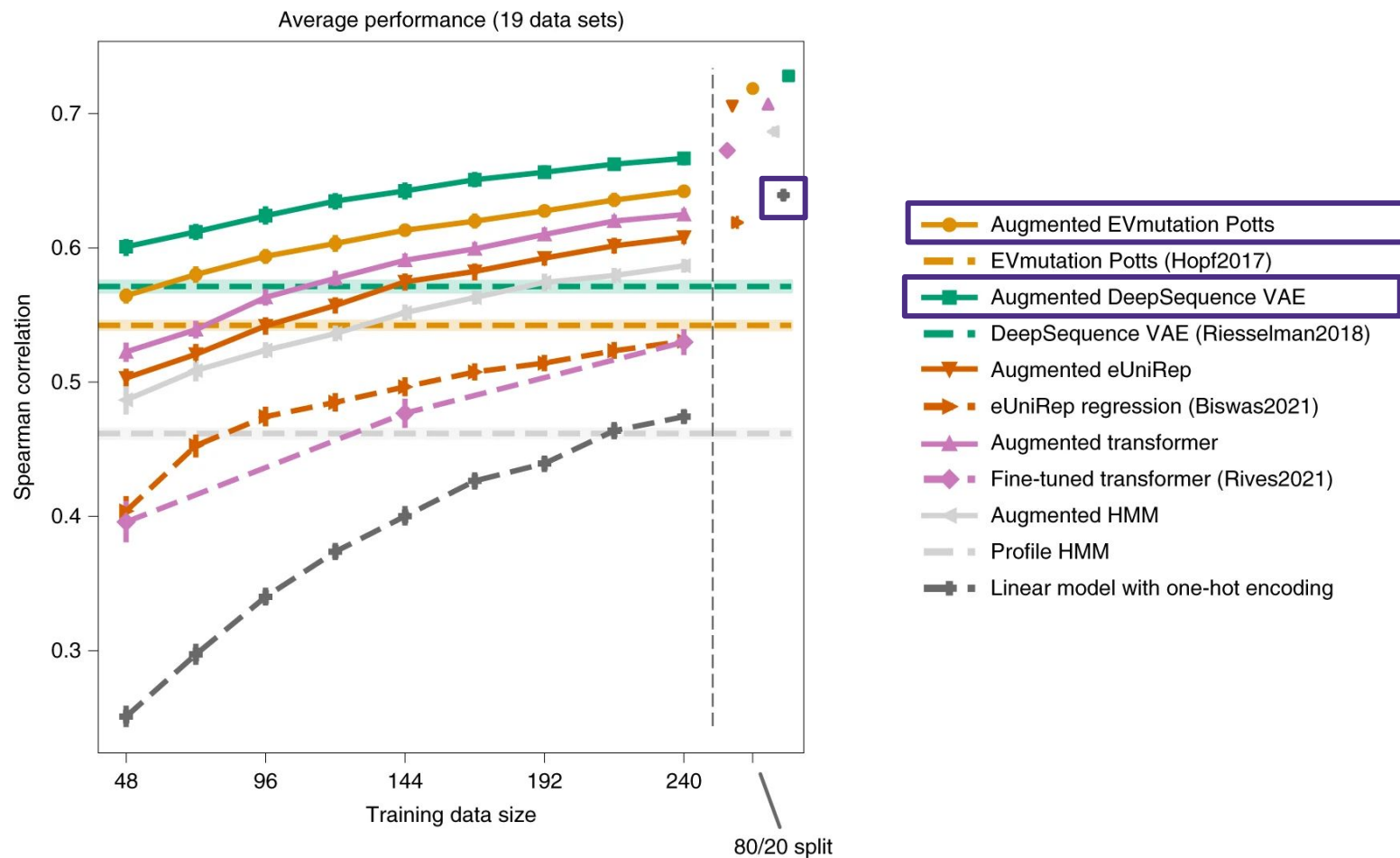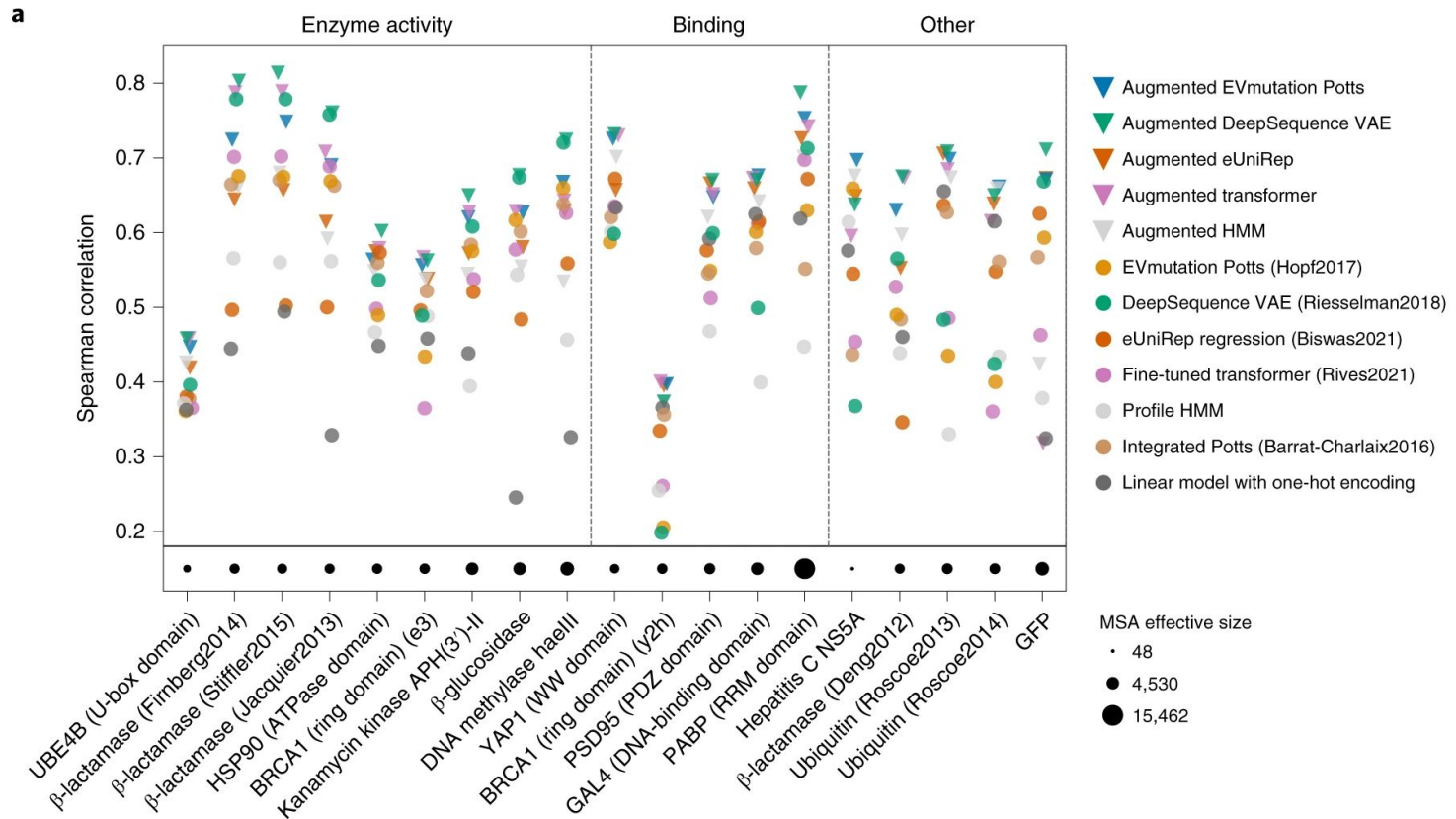
One hot linear model

a

UniRef50 ← UniRef100

Cluster at 50% identity

JackHMMER search for target protein

Unsupervised   Semisupervised

Weak-positive only

Weak-positive semisupervised

Pretraining

ESM-1b transformer (Rives2021)

UniRep mLSTM (Alley2019)

Fine-tuning

'Evotuned' eUniRep (Alley2019)

MSA of evolutionarily related sequences

Fine-tuning

Feature concatenation

Regression on embeddings

Feature concatenation

Density estimation

Joint training on MSA and labeled data

Fine-tuned ESM-1b (Rives2021)

Augmented ESM-1b (New)

'Low-N' eUniRep regression (Biswas2021)

Profile HMM (Shihab2013)

EVmutation Potts model (Hopf2017)

DeepSequence VAE (Riesselman2018)

Feature concatenation

Feature concatenation

Binary mask on coefficients

Feature concatenation

Augmented eUniRep (New)

Augmented HMM (New)

Augmented Potts (New)

TLmutation Potts model (Shamsi2020)

Augmented VAE (New)

Integrated Potts model (Barrat-Charlaix2016)

W

Average performance (19 data sets)

Solid and dashed lines of the same color are the aug. and non-aug. versions, respectively

Legend:
- Augmented EVmutation Potts
- EVmutation Potts (Hopf2017)
- Augmented DeepSequence VAE
- DeepSequence VAE (Riesselman2018)
- Augmented eUniRep
- eUniRep regression (Biswas2021)
- Augmented transformer
- Fine-tuned transformer (Rives2021)
- Augmented HMM
- Profile HMM
- Linear model with one-hot encoding

Spearman correlation

Training data size

80/20 split

Average performance (19 data sets)

Spearman correlation vs. Training data size

Legend:
- Augmented EVmutation Potts
- EVmutation Potts (Hopf2017)
- Augmented DeepSequence VAE
- DeepSequence VAE (Riesselman2018)
- Augmented eUniRep
- eUniRep regression (Biswas2021)
- Augmented transformer
- Fine-tuned transformer (Rives2021)
- Augmented HMM
- Profile HMM
- Linear model with one-hot encoding

80/20 split

# Performance w/ train N = 240



**Augmented DeepSequence VAE was the best (esp. enzyme activity)**

# Maximal Spearman values w/ train N = 240



| | Enzyme activity | Binding | Other |
|---|---|---|---|
| Augmented EVmutation Potts | 2/9 | 4/5 | 3/5 |
| Augmented DeepSequence VAE | 9/9 | 5/5 | 4/5 |
| Augmented eUniRep | 2/9 | 3/5 | 2/5 |
| Augmented transformer | 3/9 | 4/5 | 2/5 |
| Augmented HMM | 1/9 | 2/5 | 2/5 |
| EVmutation Potts (Hopf2017) | | | |
| DeepSequence VAE (Riesselman2018) | 3/9 | | |
| eUniRep regression (Biswas2021) | | | |
| Fine-tuned transformer (Rives2021) | | | |
| Profile HMM | | | |
| Integrated Potts (Barrat-Charlaix2016) | 1/9 | 1/5 | |
| Linear model with one-hot encoding | | 1/5 | |

**Augmented DeepSequence VAE was the best (esp. enzyme activity)**

# Performance w/ train N = 240

> **Models with evolutionary data had better Spearman correlation w/ larger effective MSA size**

> **Relative model ranking appeared to not relate to MSA size**

# Effect of reducing MSA size

> **Chose largest MSA data set (poly(A)-binding protein) and decreased effective size**
> **Examined aug. Potts model peformance**

# Single to higher order mutant prediction

> Trained the model on only single mutant data - tested on single, double, triple, and quadruple mutants

> Should capture how much epistasis contributes to the fitness landscape, and if/how much models capture it

> Only 3 datasets

W

Test on: single mutants | Double mutants (extrapolation) | Triple mutants (extrapolation) | Quadruple mutants (extrapolation)

GFP | PABP-RRM | UBE4B

Spearman correlation

Training data size

TS = 613 | TS = 3,662 | TS = 2,026 | TS = 844
TS = 1,188 | TS = 36,522
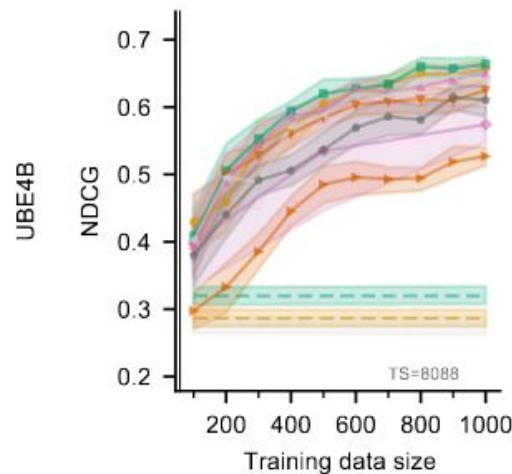TS = 613 | TS = 21,969 | TS = 8,088 | TS = 1,404
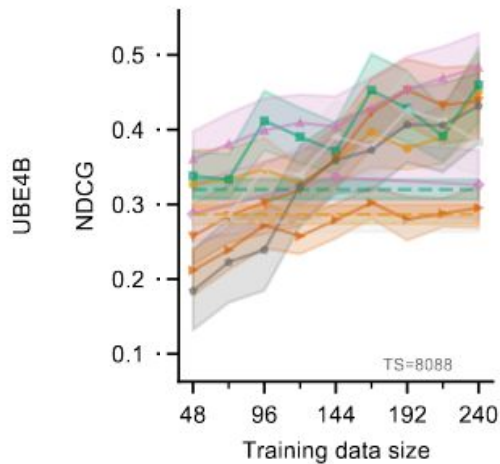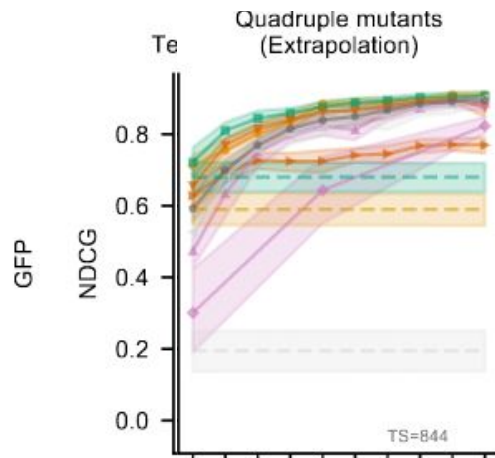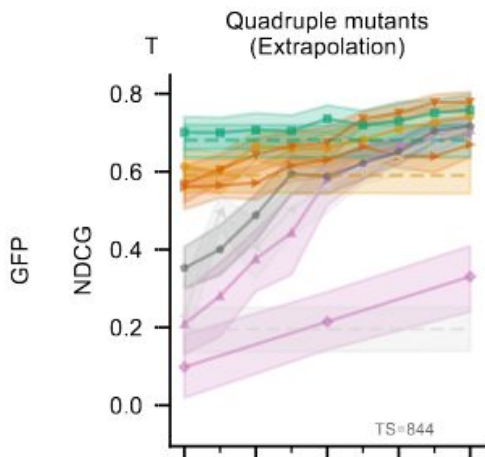
Legend:
- Augmented EVmutation Potts
- EVmutation Potts (Hopf2017)
- Augmented DeepSequence VAE
- DeepSequence VAE (Riesselman2018)
- Augmented eUniRep
- eUniRep regression (Biswas2021)
- Augmented transformer
- Fine-tuned transformer (Rives2021)
- Augmented HMM
- Profile HMM
- Linear model with one-hot encoding

# Single to higher order mutant prediction

> **Poor performance on ubiquitination factor E4B (UBE4B) may be due to evolutionary data not providing much relevant information to assayed value**

Test on: Single mutants | Double mutants | Triple mutants (Extrapolation) | Quadruple mutants (Extrapolation)
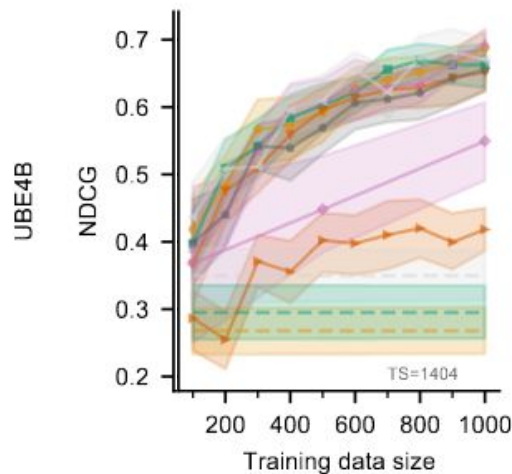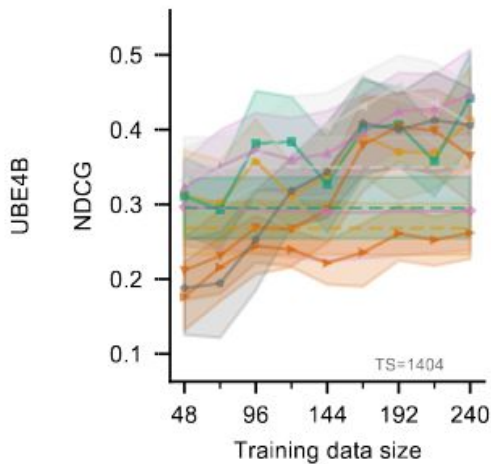
Legend:
- Augmented EVmutation Potts
- EVmutation Potts [Hopf2017]
- Augmented DeepSequence VAE
- DeepSequence VAE [Riesselman2018]
- Augmented eUniRep
- eUniRep regression [Biswas2021]
- Augmented Transformer
- Finetuned Transformer [Rives2021]
- Augmented HMM
- Profile HMM
- Linear model with one-hot encoding

TS=613, TS=3662, TS=2026, TS=844
TS=1188, TS=36522
TS=613, TS=21969, TS=8088, TS=1404

Left: single
Right: single and double

**Left: single**
**Right: single and double**

Legend:
- Augmented EVmutation Potts
- EVmutation Potts [Hopf2017]
- Augmented DeepSequence VAE
- DeepSequence VAE [Riesselman2018]
- Augmented eUniRep
- eUniRep regression [Biswas2021]
- Augmented Transformer
- Finetuned Transformer [Rives2021]
- Augmented HMM
- Profile HMM
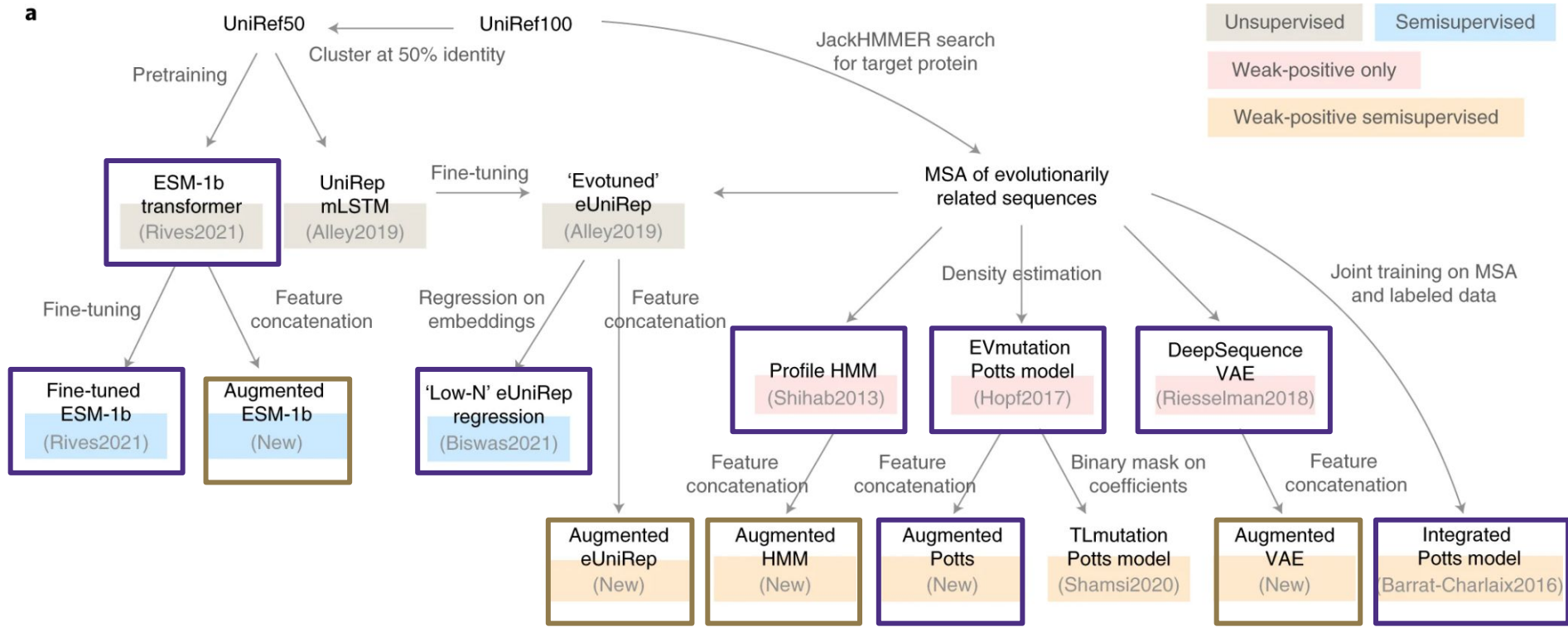- Linear model with one-hot encoding

**Left: single**
**Right: single and double**

Triple mutants (Extrapolation)

TS=2026

TS=8088

Legend:
- Augmented EVmutation Potts
- EVmutation Potts [Hopf2017]
- Augmented DeepSequence VAE
- DeepSequence VAE [Riesselman2018]
- Augmented eUniRep
- eUniRep regression [Biswas2021]
- Augmented Transformer
- Finetuned Transformer [Rives2021]
- Augmented HMM
- Profile HMM
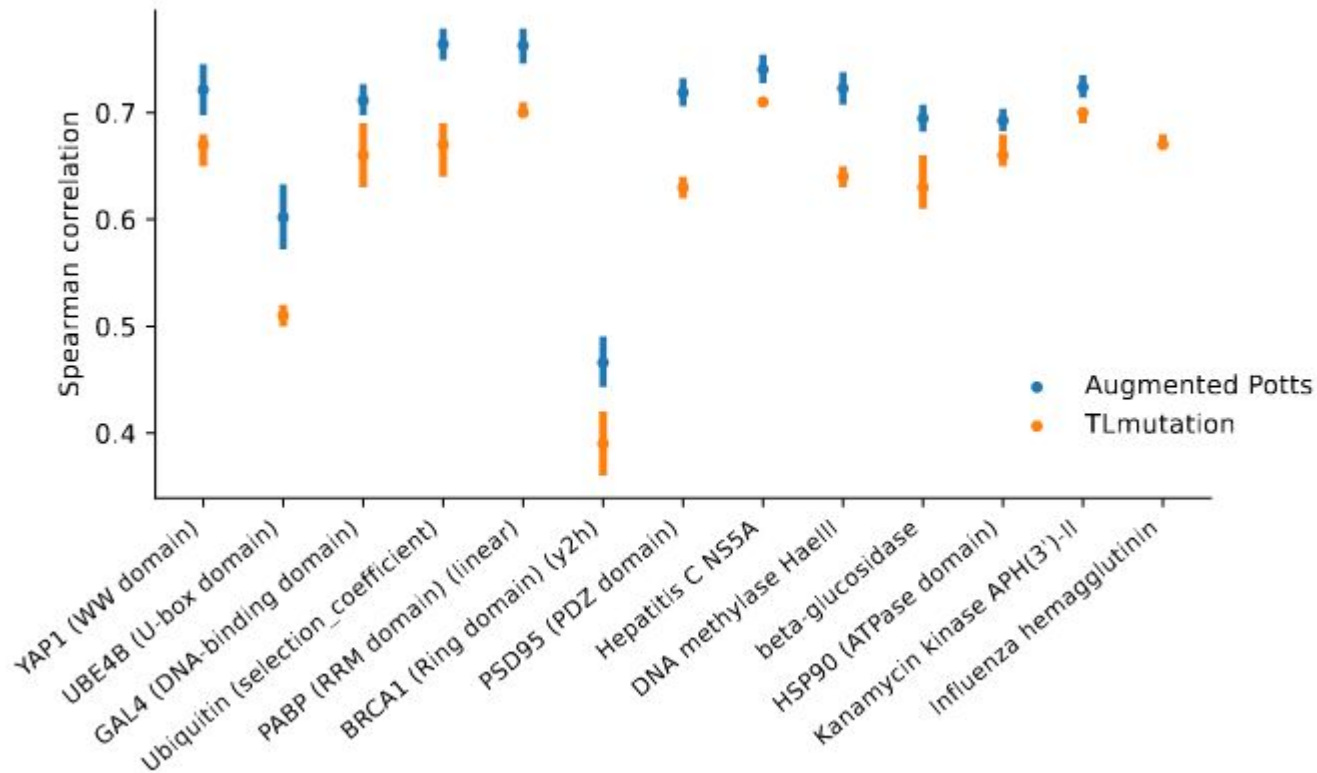- Linear model with one-hot encoding

Left: single
Right: single and double

Quadruple mutants (Extrapolation)

- Augmented EVmutation Potts
- EVmutation Potts [Hopf2017]
- Augmented DeepSequence VAE
- DeepSequence VAE [Riesselman2018]
- Augmented eUniRep
- eUniRep regression [Biswas2021]
- Augmented Transformer
- Finetuned Transformer [Rives2021]
- Augmented HMM
- Profile HMM
- Linear model with one-hot encoding

# TLMutation

> **Conceptually similar to the aug. Potts model (combining density model and supervised learning)**
> **Allows for zeroing out Potts model parameters with supervised learning - learns a mask.**
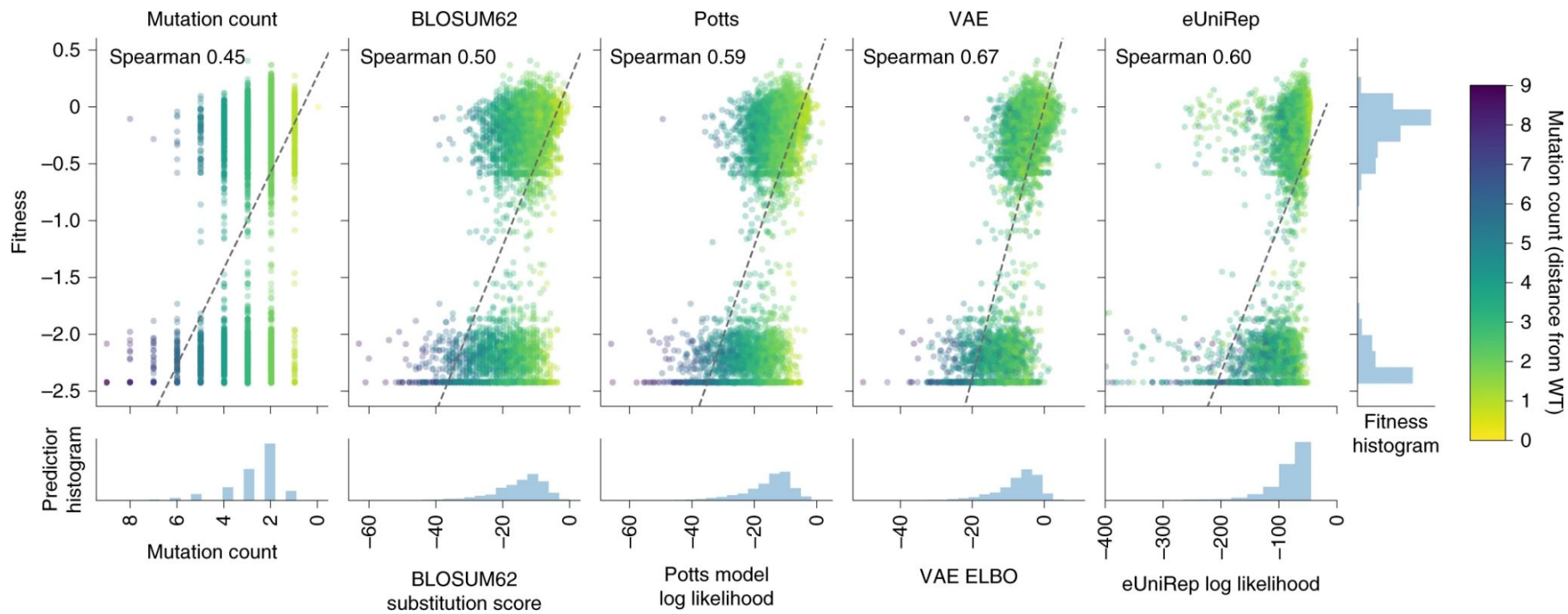> **More computationally expensive, worse than the aug Potts model**
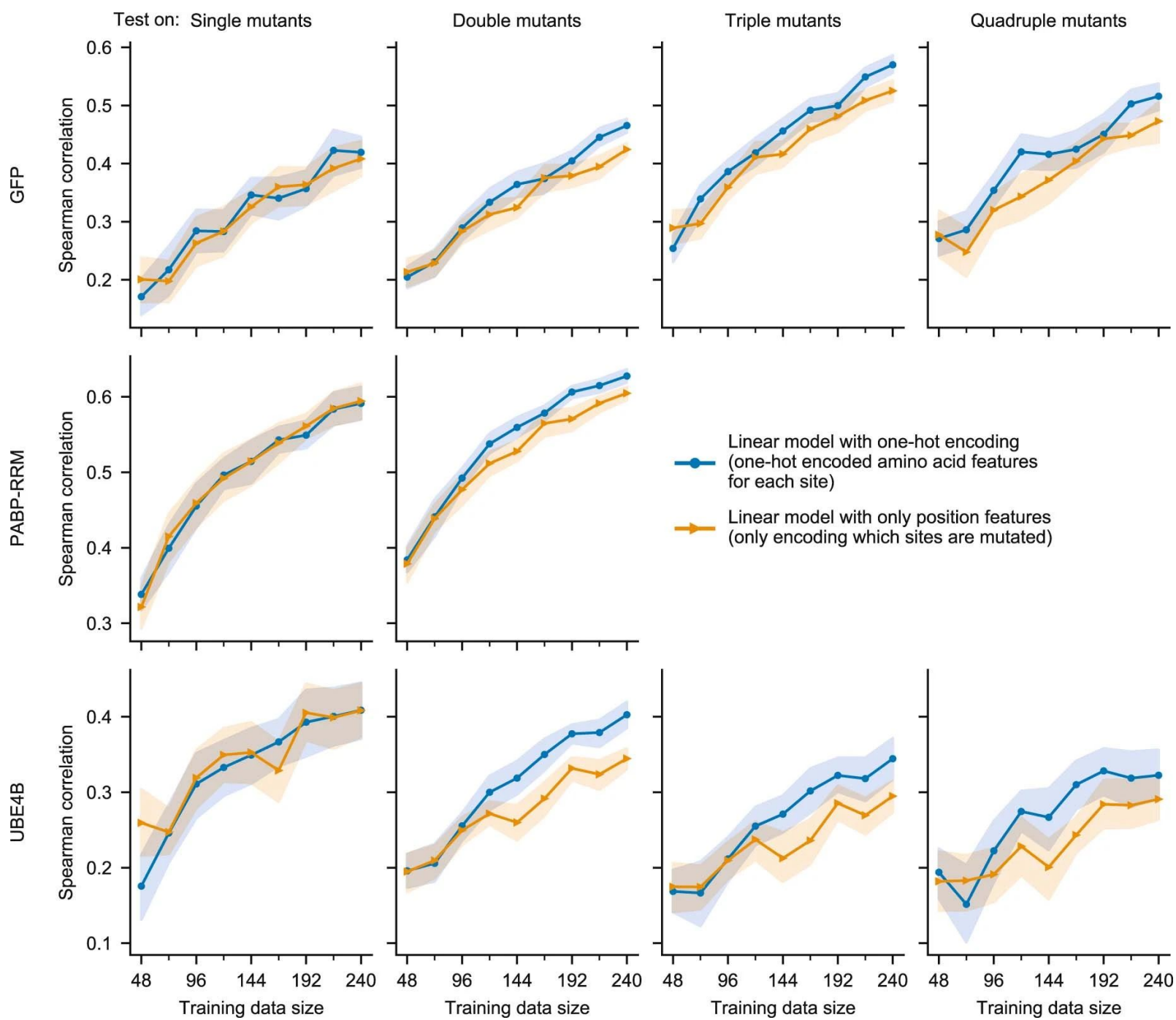
# TLMutation comparison

# What about really simple models?

> They tried using edit distance from WT to predict fitness & found some correlation with GFP
> Non-aug. predictions were unimodal, but fitness values were bimodal
> Less correlation with UBE4B
> Tried just encoding position information +

**W**

# Summary

> Simple linear regression with one-hot encoded amino acid features and a evolutionary density feature from density models outperforms said density models

> Deep learning models may be used with these features instead - but this was not tested

> Aug. transformer could be used for small MSA proteins

# Discussion questions

> **Do you think using their augmented features with a more complicated regression model would lead to a better predictor?**
>  – **Would it be worth (presumably) trade-offs in requiring higher N training set sizes?**
>  – **Would it be worth it for protein design?**
> **Do you think their decision to remove TLMutation from their comparison figures throughout was fair to the assessment?**
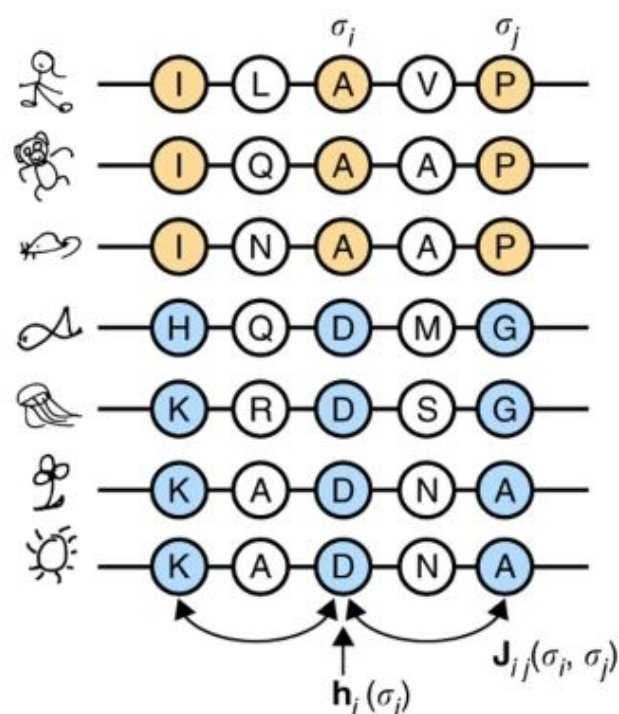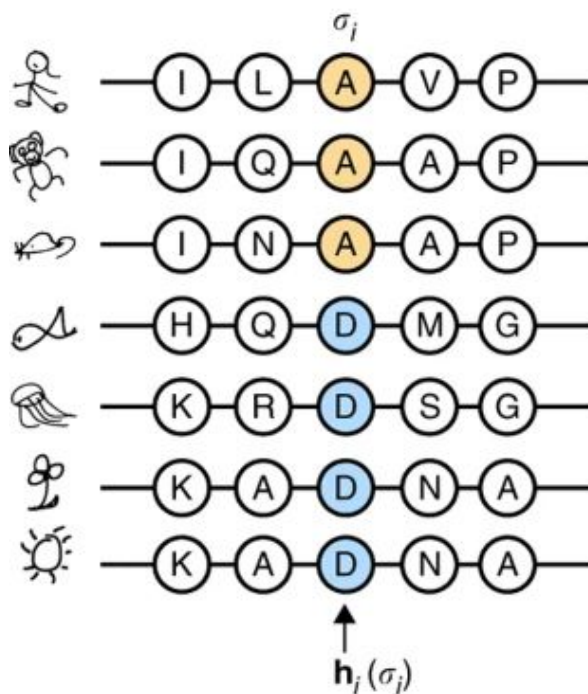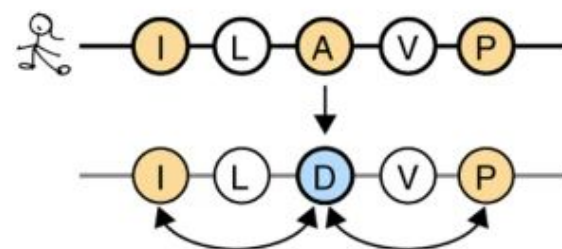
**W**

# Misc.

## Existing approaches

### Independent model

## Our approach (EVmutation)

### Epistatic model

## Model constraints on sequences

$\sigma_i$

$\sigma_i$ $\sigma_j$

$\mathbf{h}_i(\sigma_i)$

$\mathbf{h}_i(\sigma_i)$ $\mathbf{J}_{ij}(\sigma_i, \sigma_j)$

## Predict effects of mutations

✗ A → D wrongly predicted neutral ignoring sequence context

✓ A → D correctly predicted damaging needs couplings to other sites

| Protein | UniProt ID | Measurement | # Mutations before exclusion | # Mutations after excluding positions with $\geq$ 30% gaps | MSA size | Reference |
|---|---|---|---|---|---|---|
| β-glucosidase | (Custom sequence) | Enzyme activity | 3000 (1) | 2634 (1) | 28048 | Romero et al., PNAS, 2015 |
| β-lactamase | BLAT_ECOLX | Ampicillin resistance<br>Ampicillin resistance<br>Amoxicillin resistance<br>Stability | 5199 (1)<br>4997 (1)<br>990 (1)<br>4998 (1) | 4611 (1)<br>4807 (1)<br>951 (1)<br>4808 (1) | 8403 | Firnberg et al., Mol Biol Evol, 2014<br>Stiffler et al., Cell, 2015<br>Jacquier et al., PNAS 2013<br>Deng et al., JMB, 2012 |
| BRCA 1 (RING domain) | BRCA1_HUMAN | E3 ligase activity<br>BARD1 interaction | 4872 (1)<br>1748 (1) | 1382 (1)<br>1335 (1) | 25828 | Starita et al., Genetics, 2015 |
| PSD95 (PDZ domain) | DLG4_RAT | Peptide binding | 1578 (1) | 1578 (1) | 102410 | McLaughlin et al., Nature, 2012 |
| GAL4 (DNA-binding domain) | GAL4_YEAST | Transcriptional activity | 1196 (1) | 1123 (1) | 17521 | Kitzmann et al., Nature Methods, 2015 |
| HSP90 (ATPase domain) | HSP82_YEAST | ATPase activity | 4324 (1) | 4104 (1) | 15329 | Mishra et al., Cell Reports, 2016 |
| Kanamycin kinase APH(3')-II | KKA2_KLEPN | Kinase activity | 4582 (1) | 4385 (1) | 12861 | Melnikov et al., NAR, 2014 |
| DNA methylase HaeIII | (Custom sequence) | DNA methyltransferase activity | 1778 (1, filtered) | 1634 (1) | 14115 | Rockah-Shmuel et al., PLOS Comp Bio, 2015 |
| Poly(A)-binding protein (RRM domain) | PABP_YEAST | RNA binding | 1188 (1)<br>36522 (2) | 1188 (1)<br>36522 (2) | 152041 | Melamed et al., RNA, 2013 |
| Hepatitis C NS5A | POLG_HCVJF | Viral replication | 1632 (1) | 1632 (1) | 8106 | Qi et al., PLOS Pathogens, 2014 |
| Ubiquitin | RL401_YEAST | Growth<br>E1 reactivity | 1196 (1)<br>1366 (1) | 1161 (1)<br>1295 (1) | 21448 | Roscoe et al, JMB, 2013<br>Roscoe et al, JMB, 2014 |
| UBE4B (U-box domain) | UBE4B_MOUSE | Ligase activity | 91031 (1-9 mut.) | 613 (1)<br>21969 (2)<br>8088 (3)<br>1404 (4)<br>216 ($\geq$ 5) | 9172 | Starita et al., PNAS, 2013 |
| YAP1 (WW domain 1) | YAP1_HUMAN | Peptide binding | 563 (1) | 319 (1) | 40302 | Araya et al., PNAS, 2012 |
| Green Fluorescent Protein | (Custom sequence) | Fluorescence | 51715 (1-14 mut.) | 613 (1)<br>3662 (2)<br>2026 (3)<br>844 (4)<br>630 ($\geq$ 5) | 22535 | Sarkisyan et al., Nature, 2016 |