

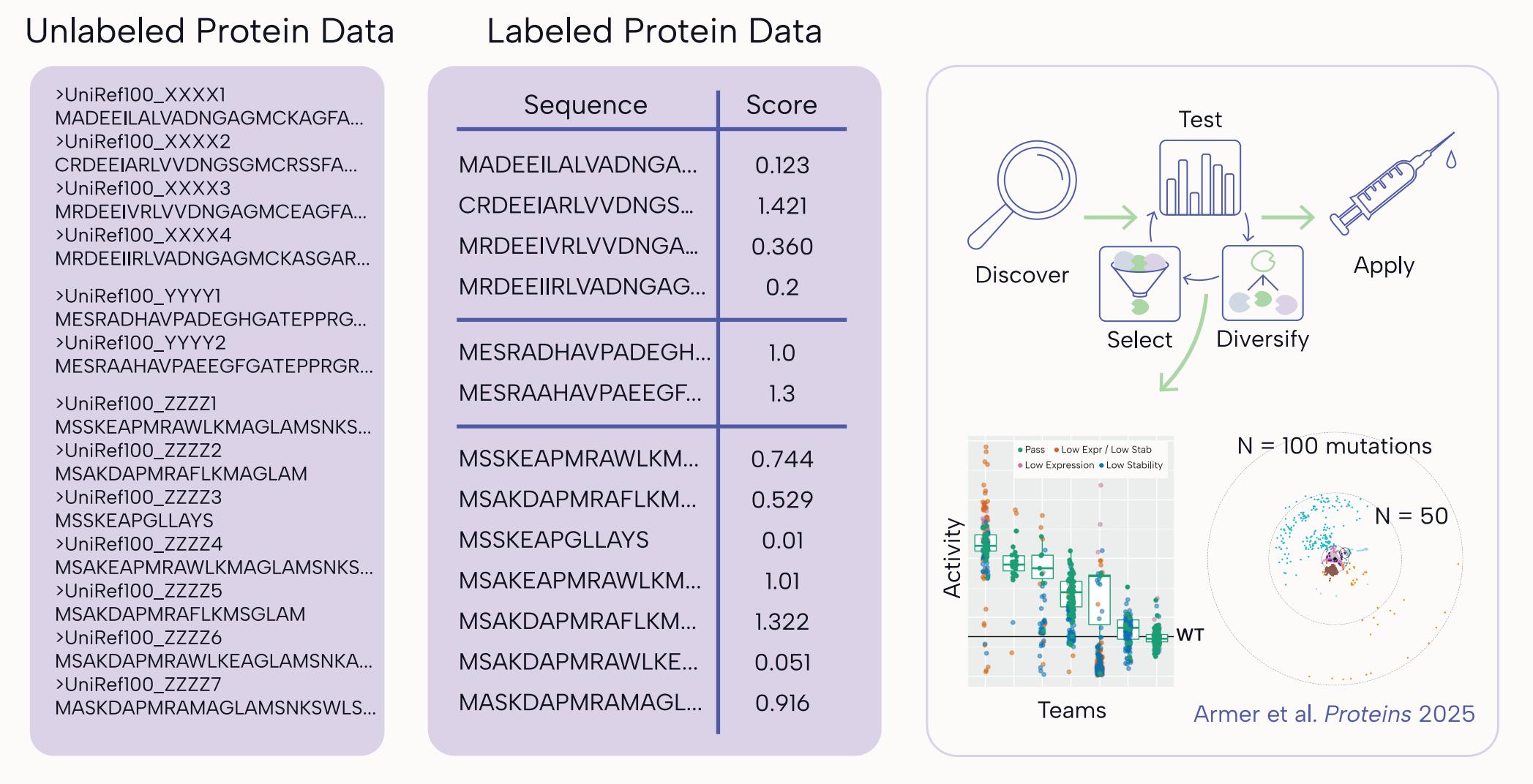


# Scaling and Data Saturation in Protein Language Models

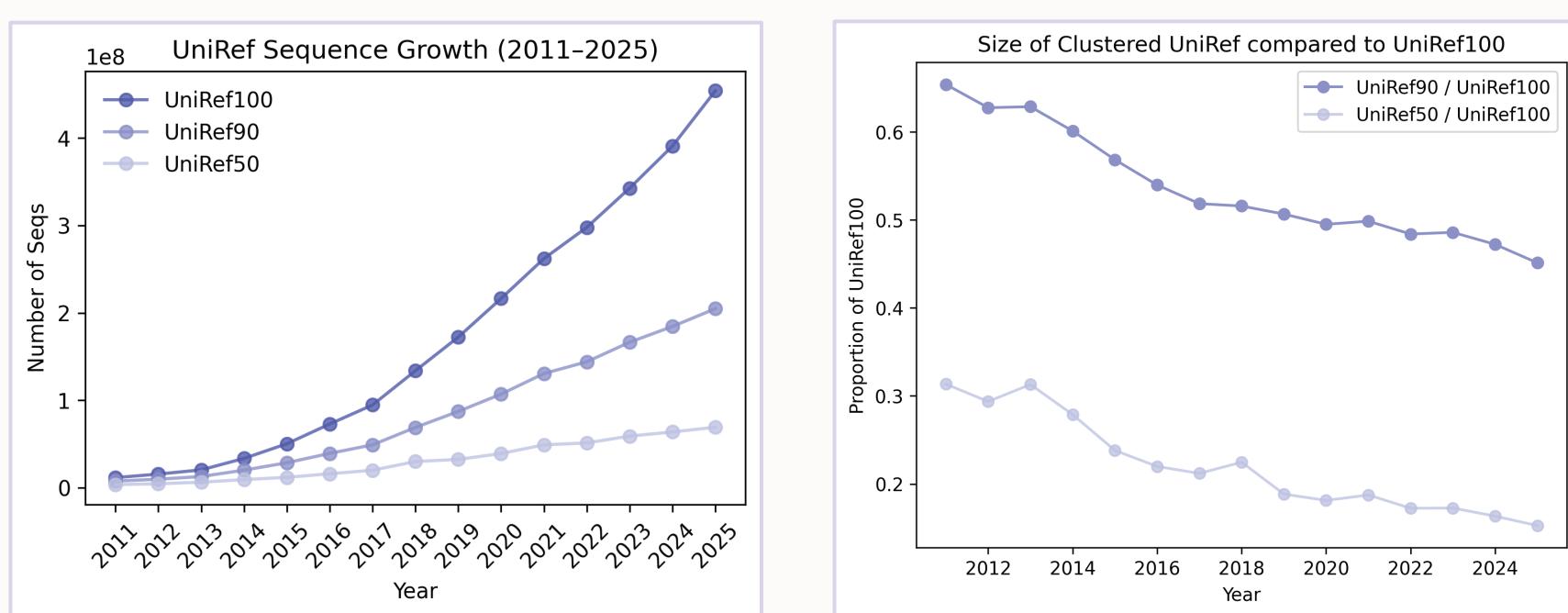
Aviv Spinner<sup>1</sup>, Erica DeBenedictis<sup>1</sup>, Corey M. Hudson<sup>1</sup>

The Align Foundation | alignbio.org

**1** Biological data is expensive and extremely useful for modeling sequence to function relationships.



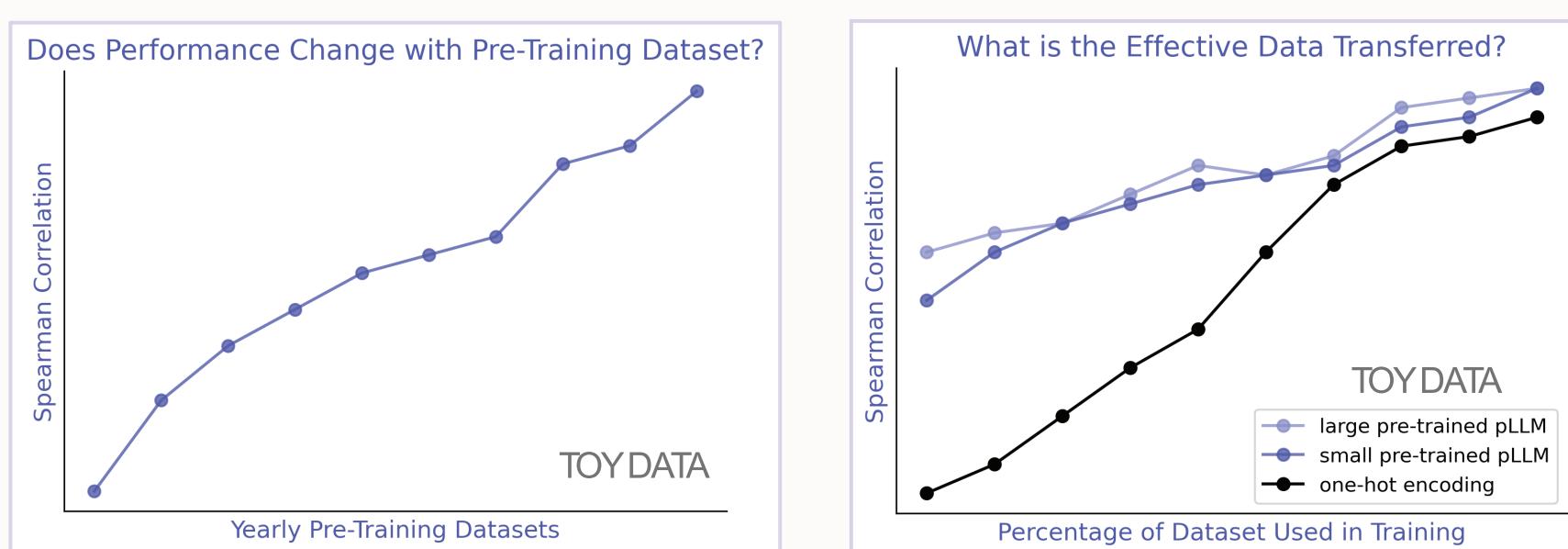
**2** Biological data is growing at an ever-increasing rate...



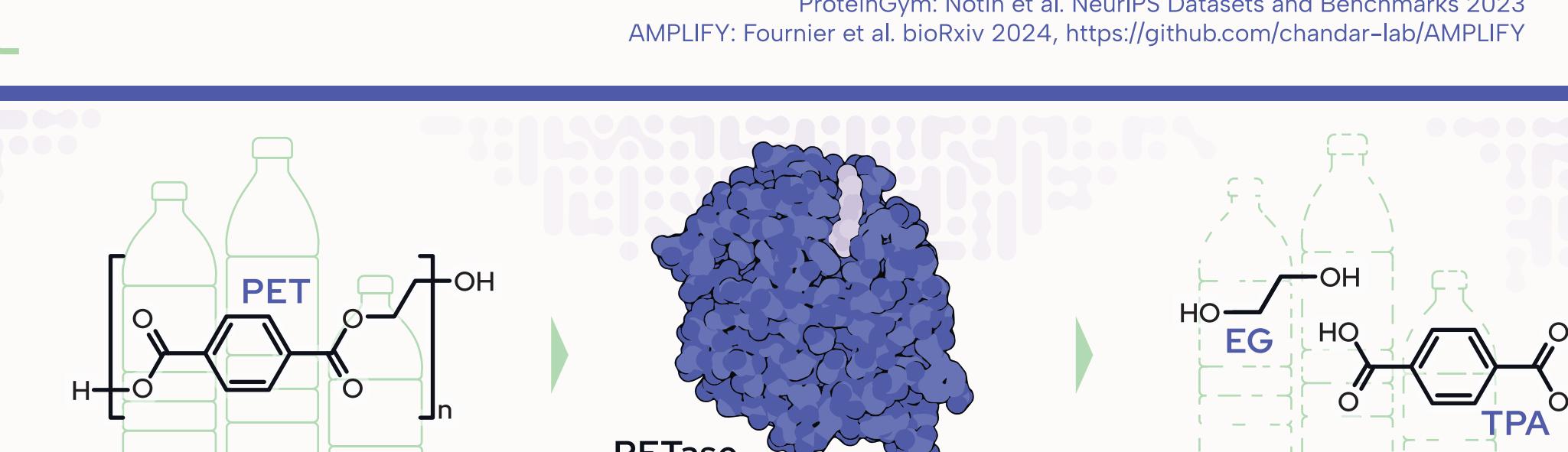
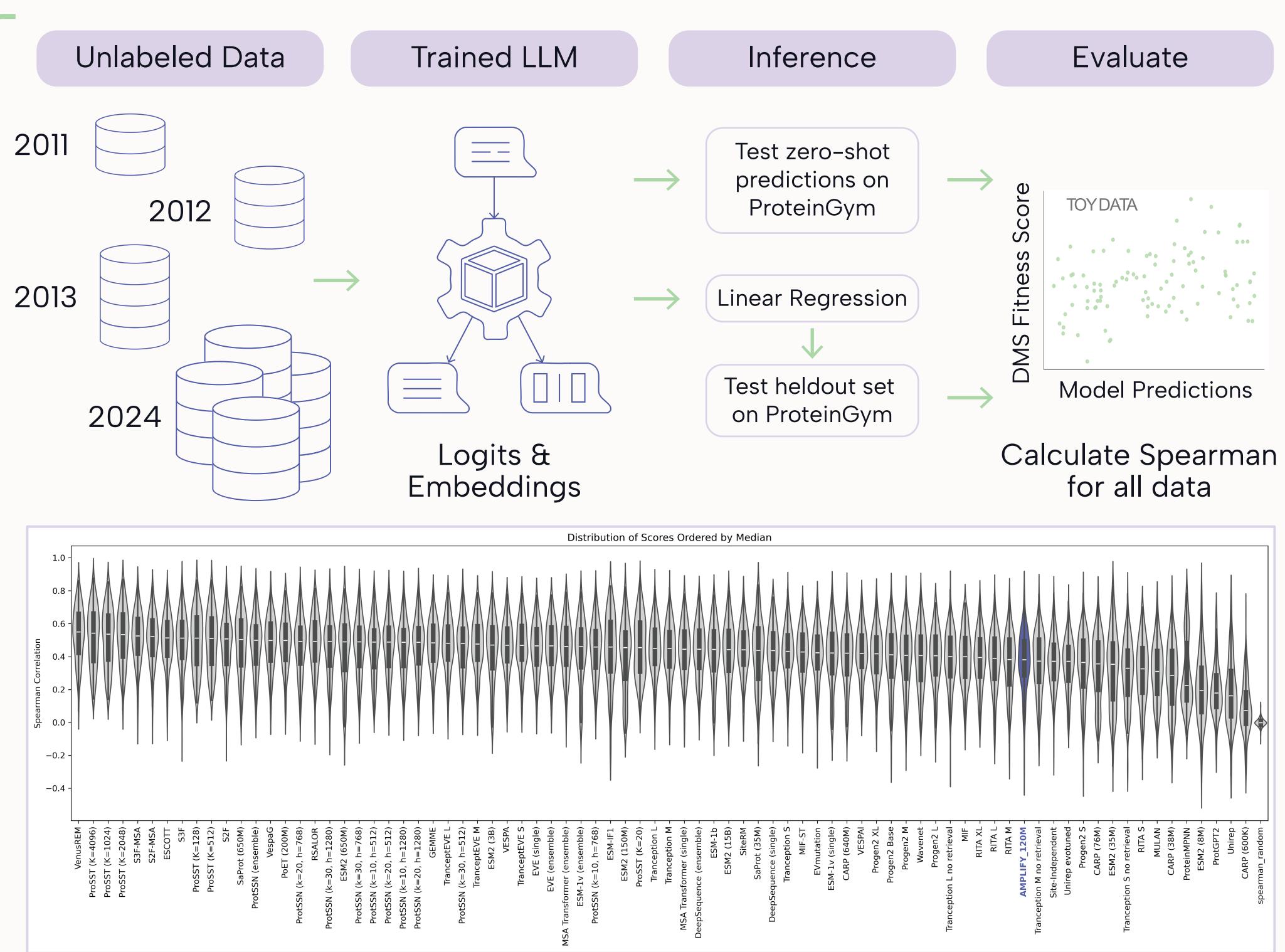
... and the data pose unique challenges

- > Redundancy and imbalance
- > Annotation sparsity
- > Noisy and heterogeneous data sources
- > Proteins have multiple functions

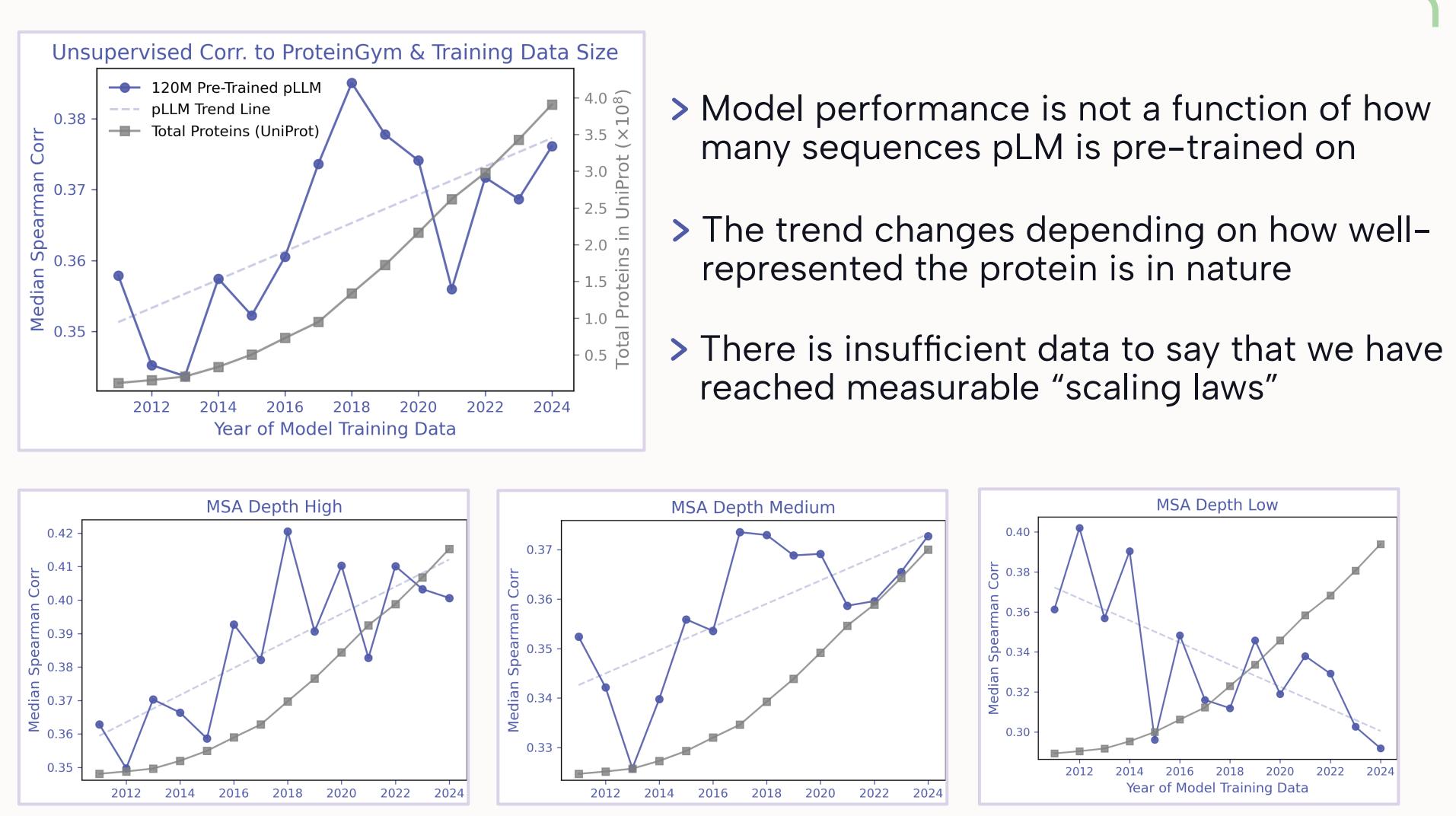
**3** How does this data impact the performance of protein language models?



**4** We use a suite of AMPLIFY models trained on time points of UniRef100 and evaluated with ProteinGym

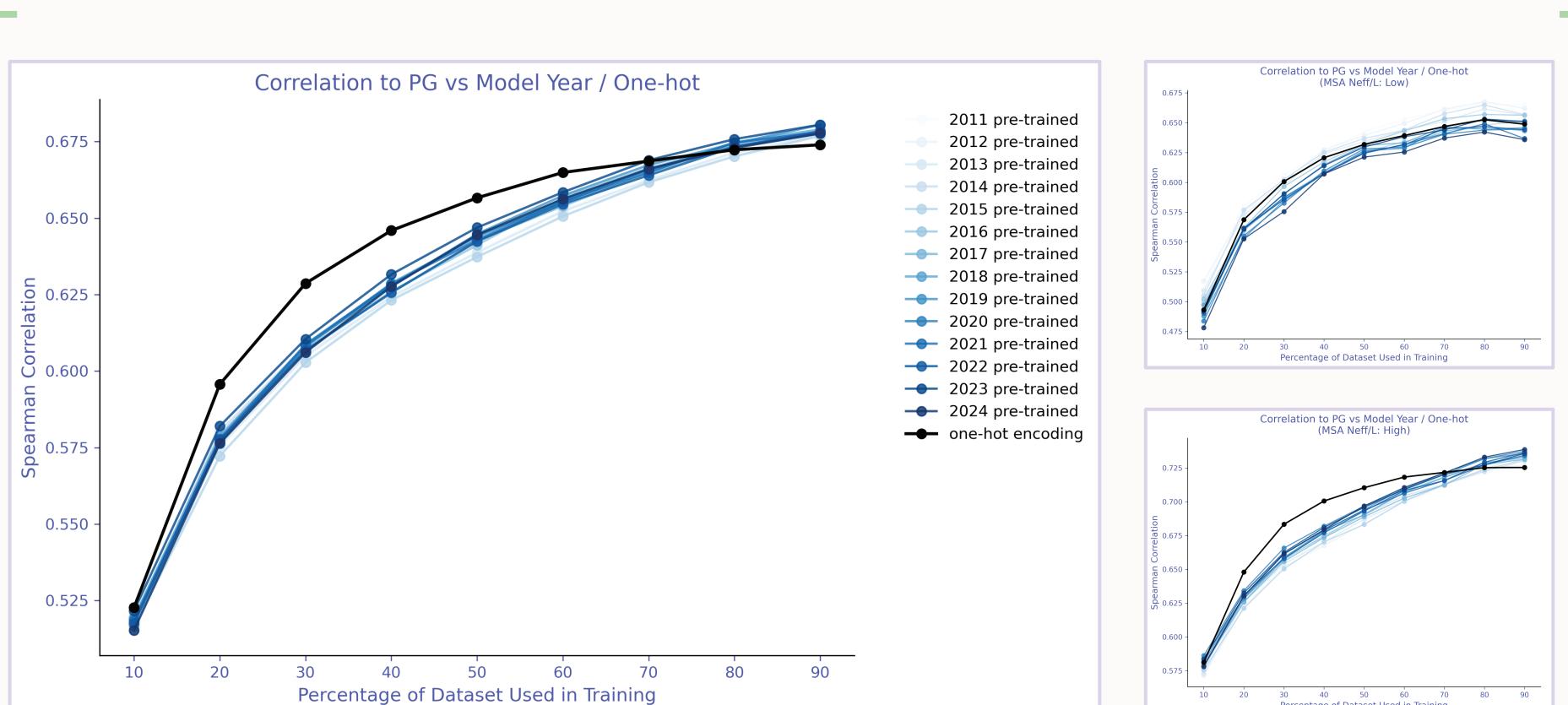


**5** More unlabeled sequence data does not yield monotonically better performance

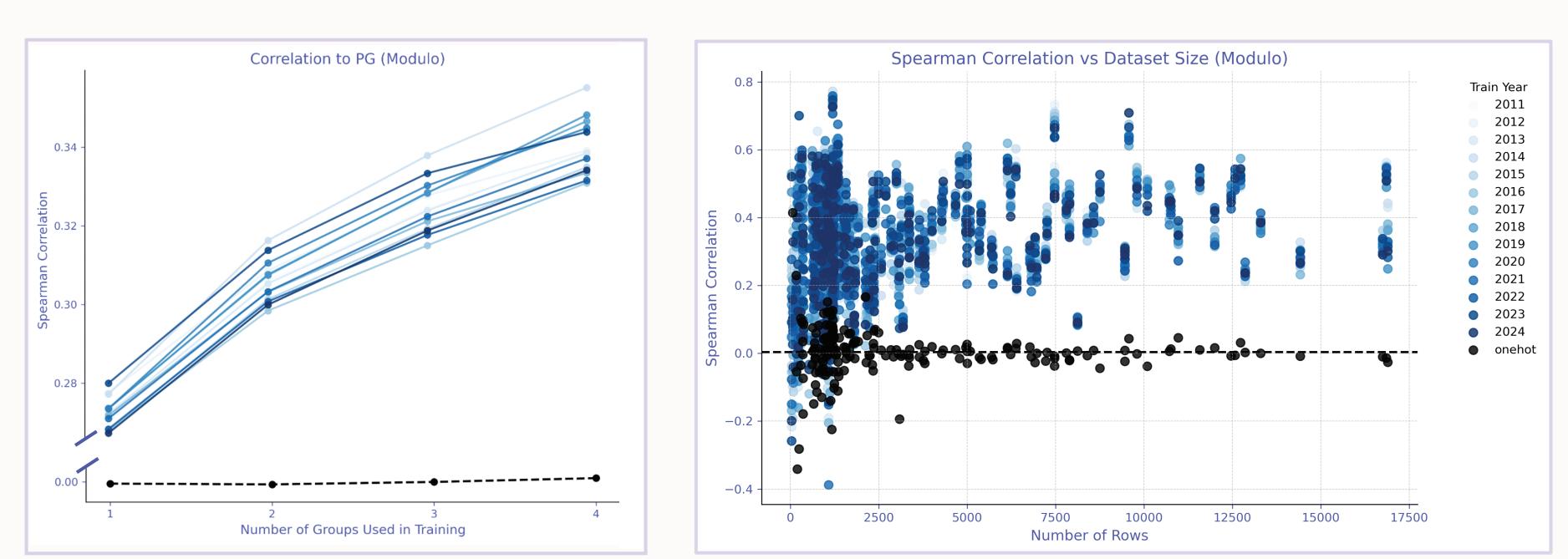


- > Model performance is not a function of how many sequences pLLM is pre-trained on
- > The trend changes depending on how well-represented the protein is in nature
- > There is insufficient data to say that we have reached measurable "scaling laws"

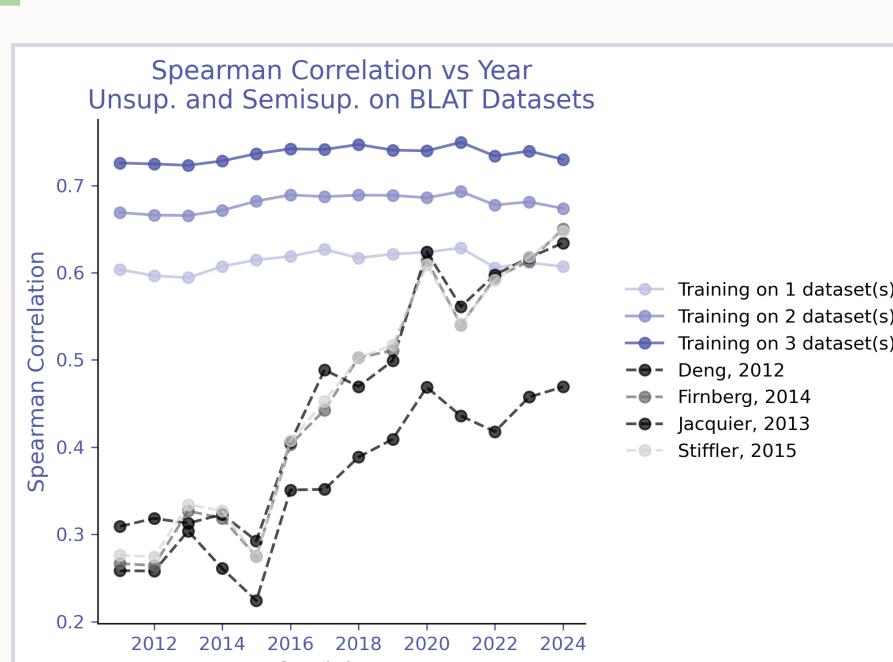
**6** Including labeled data boosts performance, one-hot encodings are helpful



- > In random splits, correlation with experimental data is strengthened by a little training data (i.e. 10% → 0.52, 0-shot → 0.38).
- > Correlation is overall worse in position-based splits, especially for one-hot encodings



**7** Deep dive on beta-lactamase suggests some emergent scaling laws



- > Grey traces = zero-shot predictions
- > Blue traces = train on dataset(s), test on heldout dataset(s)
- > Training on any one dataset gives similar performance as a model trained on all sequences in UniRef100 from 2020
- > Training on two or three datasets gives unbeatable performance

Note: these three experimental datasets are nearly identical

## Future Directions

> Train a suite of models to answer this question of scaling more directly  
**(Want to collaborate?)**

> Semi-supervised data splits \*between\* datasets to better understand transfer learning between different proteins

> Expanding semi-supervised models to include ProteinNPT, Kermut, etc.

## Align's 2025 Protein Engineering Tournament

Whose protein model reigns supreme?  
You generate. We test. **Science wins.**

