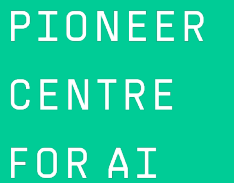# Kermut: Composite kernel regression for protein variant effects

**Peter** Mørch Groth*, **Mads** Herbert Kerrn*,
Lars Olsen, Jesper Salomon, and Wouter Boomsma

* Equal contributions

KØBENHAVNS UNIVERSITET
UNIVERSITY OF COPENHAGEN

novonesis

M·L·L·S

PIONEER CENTRE FOR AI

# Variant effect prediction

– Variant effects are measurable **changes in protein function caused by mutations** in the amino acid sequence.

– Predicting potentially **beneficial** or **deleterious** mutations is crucial for engineering and optimizing proteins, e.g., to increase activity and stability.

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Reference | V | F | A | H | P | E | T | L | 1.0 |
| Variant 1 | V | F | A | H | P | W | T | L | 0.2 |
| Variant 2 | V | F | A | H | A | E | T | L | 1.2 |
| Variant 3 | V | E | A | H | P | E | T | L | 1.5 |

– Central question:

*How can we predict variant effects given a reference protein and experimental data for a number of variants?*

# Desiderata

– **Supervised** model

  – Not all protein properties correlate with zero-shot fitness estimates

  – We want to learn from our data to guide exploration for protein engineering

– **Uncertainty quantification**

  – Valuable to quantify predictive uncertainties

  – Uncertainties should be **well-calibrated**

– **Leverage pre-trained models**

  – We're often working with few labeled sequences

# Model of choice: **Gaussian processes**

– Explicitly uses similarities between datapoints to reason about the function of interest

– Provides predictive uncertainties

– Fully specified by **two** components:

$$f(\mathbf{x}) \sim \mathcal{GP}(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$$
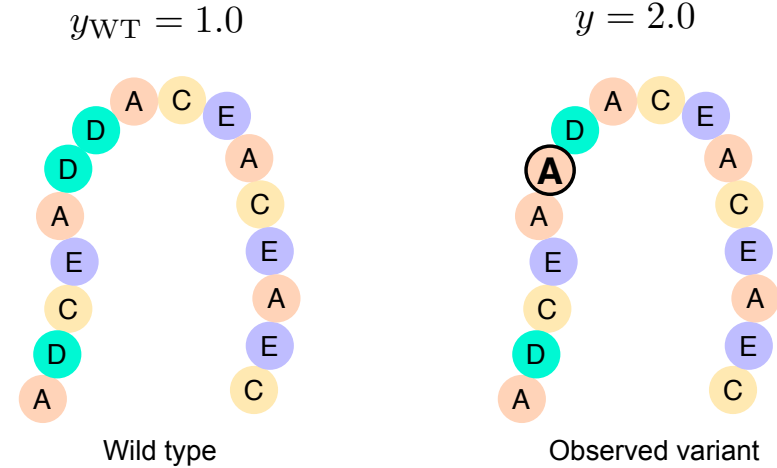
– A mean function

$$m(\mathbf{x})$$

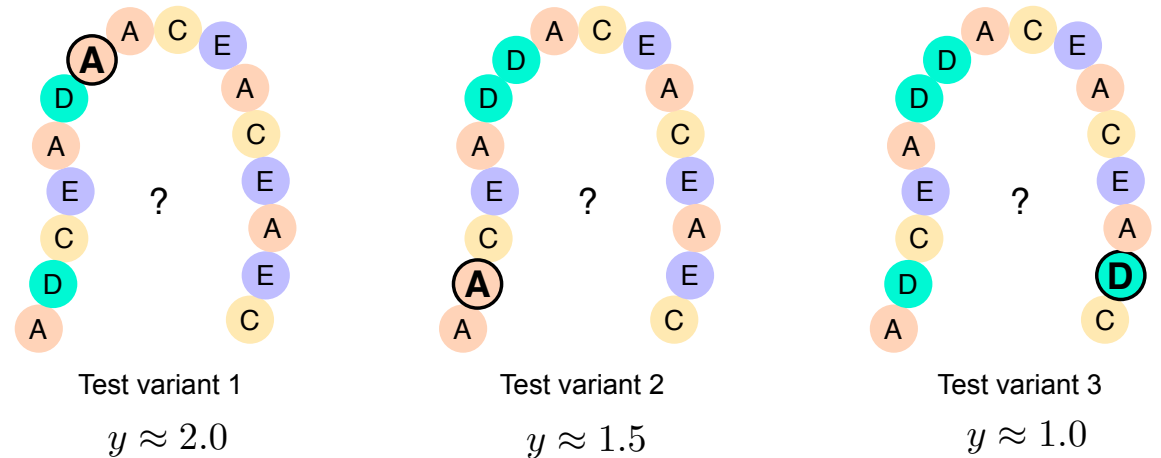– And a covariance function (also known as kernel function)

$$k(\mathbf{x}, \mathbf{x}') = \mathrm{cov}(f(\mathbf{x}), f(\mathbf{x}'))$$

# When are variants similar?

$y_{\text{WT}} = 1.0$      $y = 2.0$



Wild type      Observed variant

— Similarity given **local environments**

— Use an inverse-folding model to obtain **structure-conditioned amino acid distributions** at all sites

Given the wild type and observed variant, what can we say about the following variants?

1. Local environments should be similar
2. Individual mutations should be similar
3. Mutates sites should be physically close



Test variant 1      Test variant 2      Test variant 3

$y \approx 2.0$      $y \approx 1.5$      $y \approx 1.0$

# Kermut: a kernel for modeling mutation similarity

$$k(\mathbf{x}, \mathbf{x}') = \pi k_{\text{struct}}(\mathbf{x}, \mathbf{x}') + (1 - \pi)k_{\text{seq}}(\mathbf{x}, \mathbf{x}')$$

$$k_{\text{struct}}(\mathbf{x}, \mathbf{x}') = \lambda k_H(\mathbf{x}, \mathbf{x}') \cdot k_p(\mathbf{x}, \mathbf{x}') \cdot k_d(\mathbf{x}, \mathbf{x}')$$
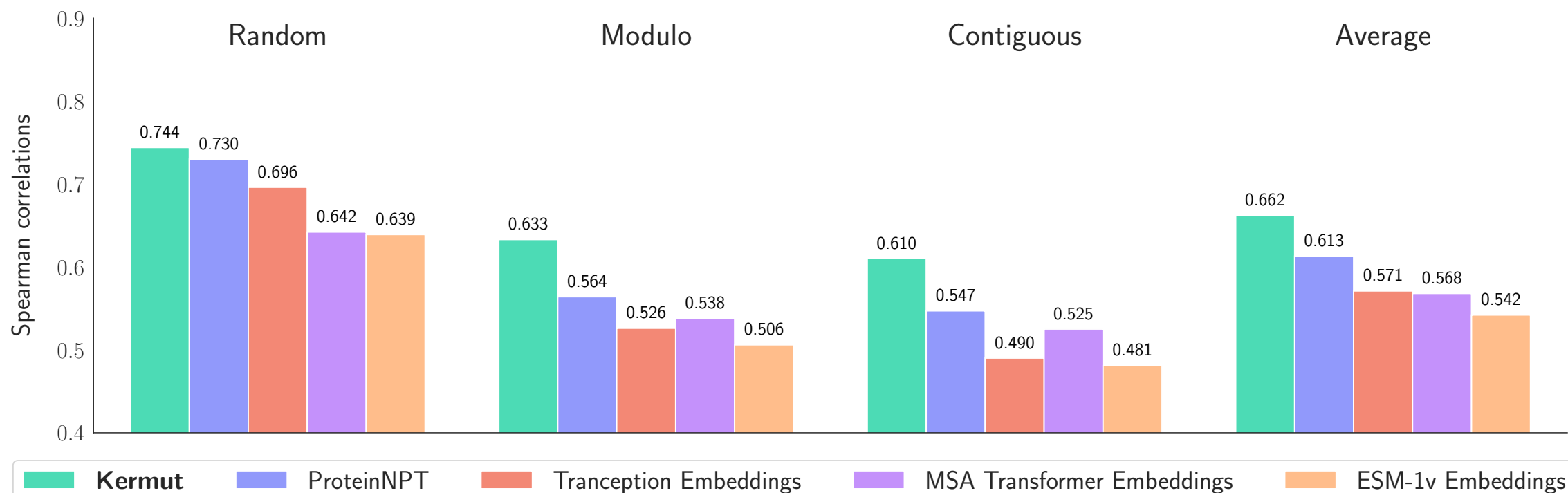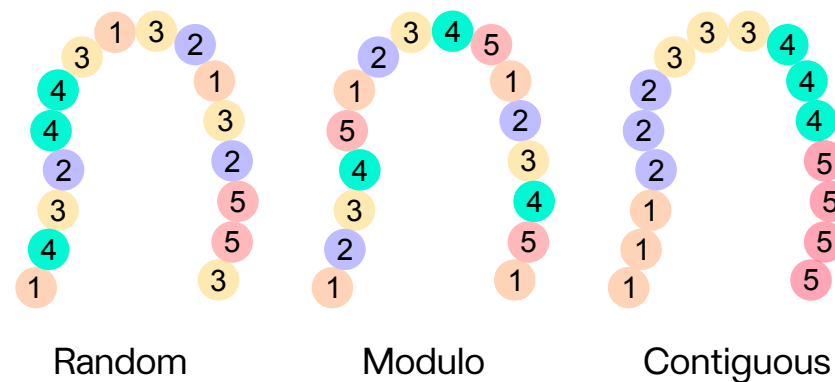
Site comparison   Mutation comparison   Distance comparison

$$k_{\text{seq}}(\mathbf{x}, \mathbf{x}') = k_{\text{SE}}(f_1(\mathbf{x}), f_1(\mathbf{x}')) = k_{\text{SE}}(\mathbf{z}, \mathbf{z}') = \exp\left(-\frac{||\mathbf{z} - \mathbf{z}'||_2^2}{2\sigma^2}\right)$$
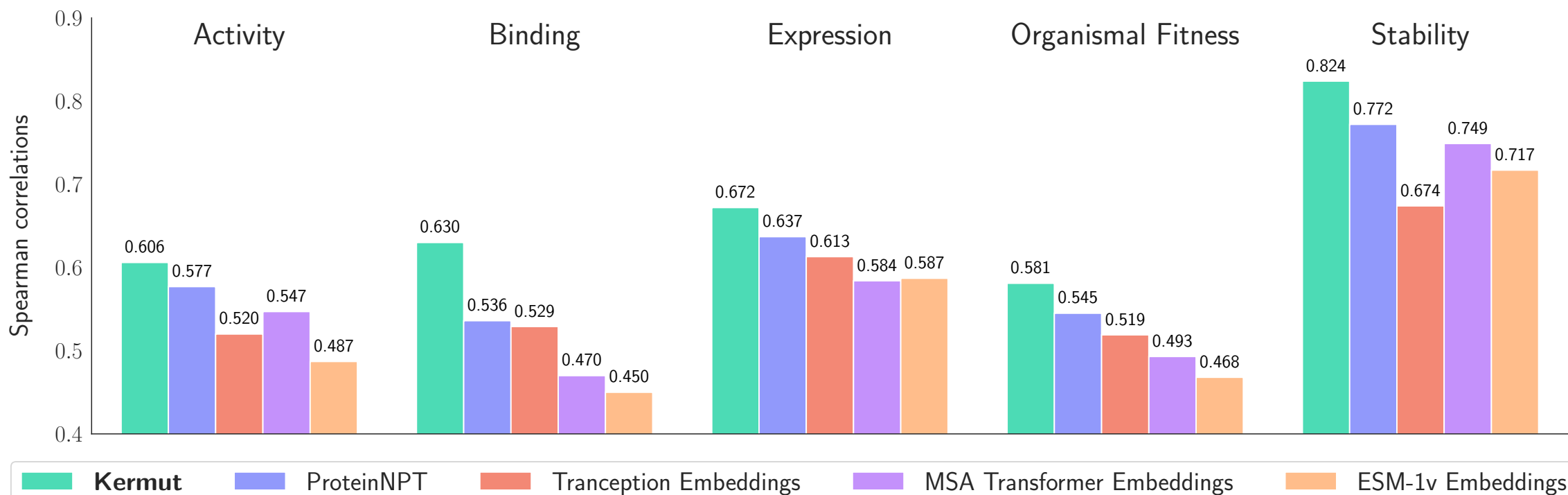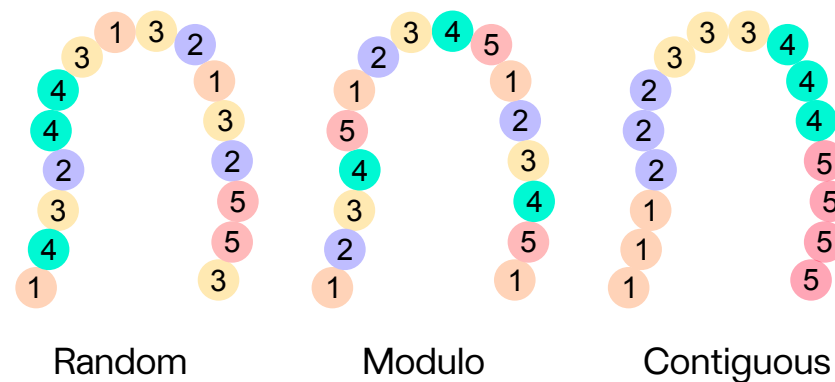
Embedding comparison

# Results

- Supervised ProteinGym benchmark
  - 217 DMS substitutions assays
  - 3 split schemes are defined for each assay with 5-fold CV in each



Random     Modulo     Contiguous

# Results per functional category

- Supervised ProteinGym benchmark
  - 217 DMS substitutions assays
  - 3 split schemes are defined for each assay with 5-fold CV in each



Random      Modulo      Contiguous



Figure: Spearman correlations across functional categories (Activity, Binding, Expression, Organismal Fitness, Stability) for Kermut, ProteinNPT, Tranception Embeddings, MSA Transformer Embeddings, and ESM-1v Embeddings.

| Category | Kermut | ProteinNPT | Tranception Embeddings | MSA Transformer Embeddings | ESM-1v Embeddings |
|---|---|---|---|---|---|
| Activity | 0.606 | 0.577 | 0.520 | 0.547 | 0.487 |
| Binding | 0.630 | 0.536 | 0.529 | 0.470 | 0.450 |
| Expression | 0.672 | 0.637 | 0.613 | 0.584 | 0.587 |
| Organismal Fitness | 0.581 | 0.545 | 0.519 | 0.493 | 0.468 |
| Stability | 0.824 | 0.772 | 0.674 | 0.749 | 0.717 |

# Speed

| Dataset | Kermut | ProteinNPT | $N$ | $L$ |
|---|---|---|---|---|
| BLAT_ECOLX | 111s | $\approx 32$h | 4996 | 286 |
| PA_I34A1 | 45s | $\approx 52$h | 1820 | 716 |
| TCRG1_MOUSE | 19s | $\approx 22$h | 621 | 37 |
| OPSD_HUMAN | 14s | $\approx 40$h | 165 | 348 |

# Conclusion

– Kermut achieves **state-of-the-art performance** for supervised variant effect prediction

– Provides **well-calibrated uncertainties** out-of-the-box

– Can be trained and evaluated orders of magnitude **faster** than competing methods

– Can easily be adapted for new pre-trained models

Limitations

– Does not support insertions and deletions

– Due to scaling, GPs scale cubically with number of datapoints[*]

– Structure kernel models multi-mutants linearly – only epistasis via sequence embeddings

– Extrapolation to higher order mutations is difficult and needs further analysis

*: Not a practical concern in most protein engineering campaigns.