

Benchmarking machine learning methods for pattern prediction and recovery in antibody sequences

Rahmad Akbar¹, Cédric R. Weber³, Igor Snapkov¹, Daniel Heinesen¹, Edvard Aksnes¹, Zixuan Liu¹, Milena Pavlovic^{1,2}, Geir Kjetil Sandve², Sai T. Reddy³, and Victor Greiff¹.



UiO University of Oslo

ETH

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

¹ UiO, Department of Immunology, Computational and Systems Immunology, Oslo, Norway

² UiO, Institute of Informatics, Biomedical Informatics, Oslo, Norway

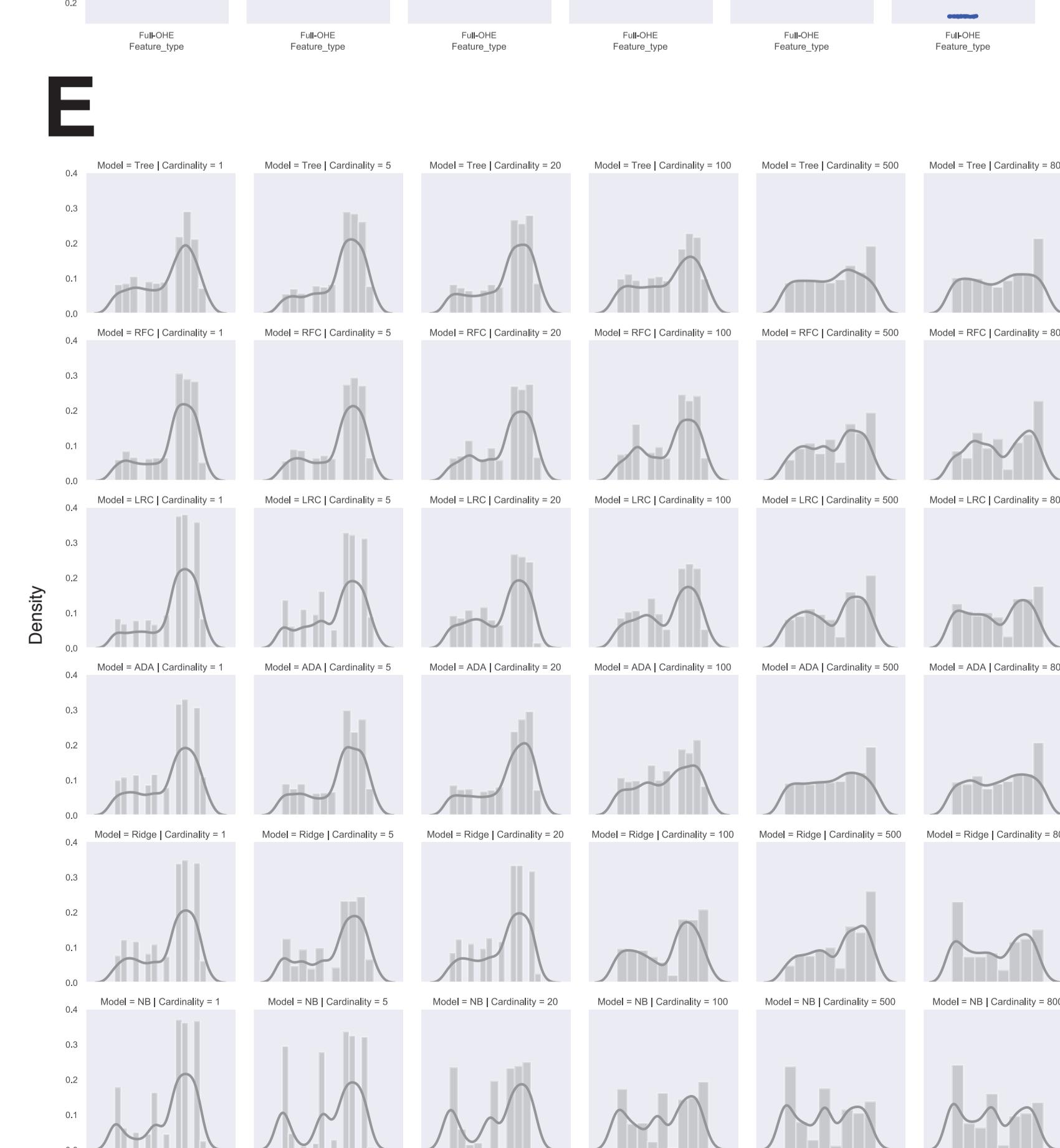
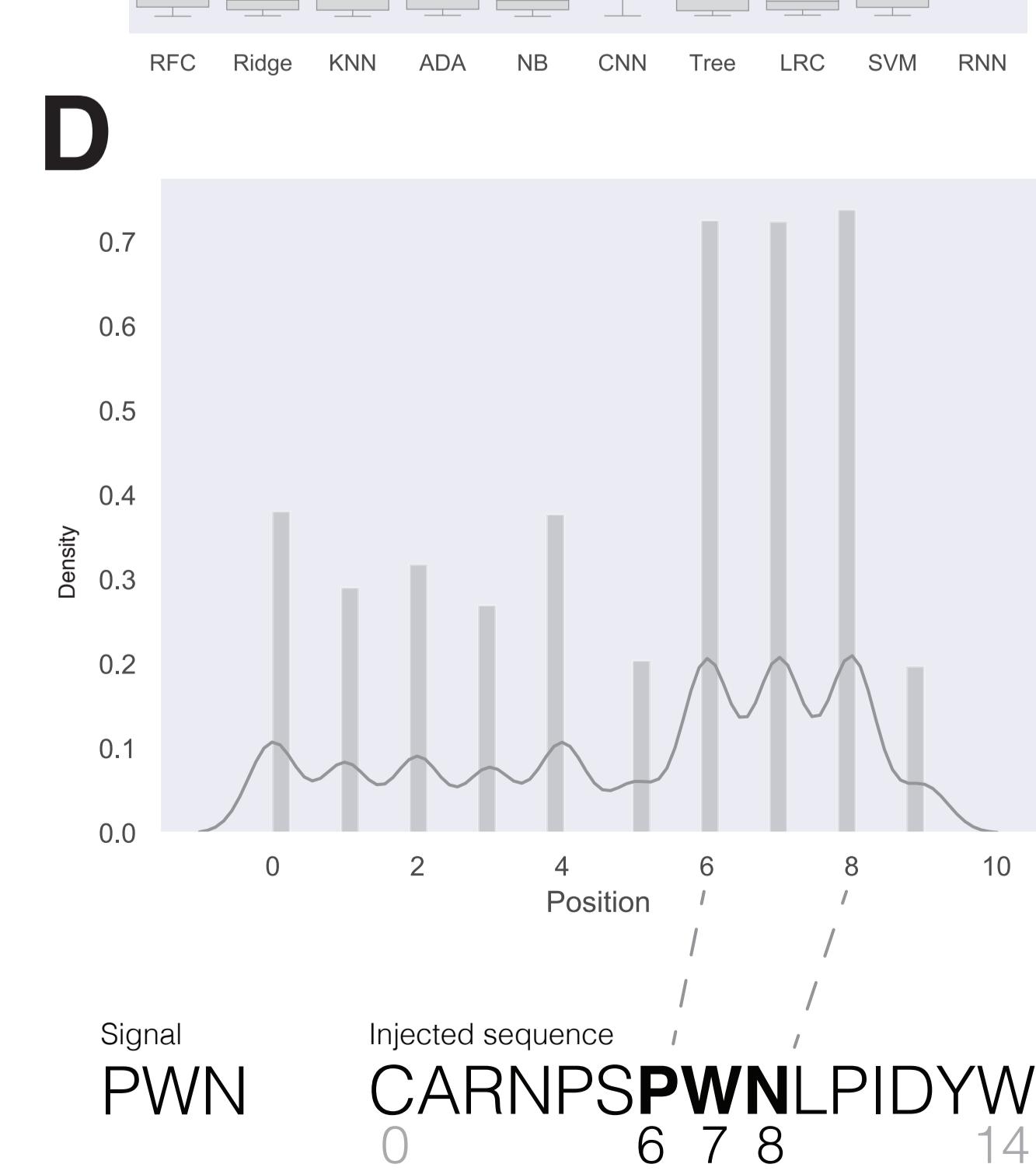
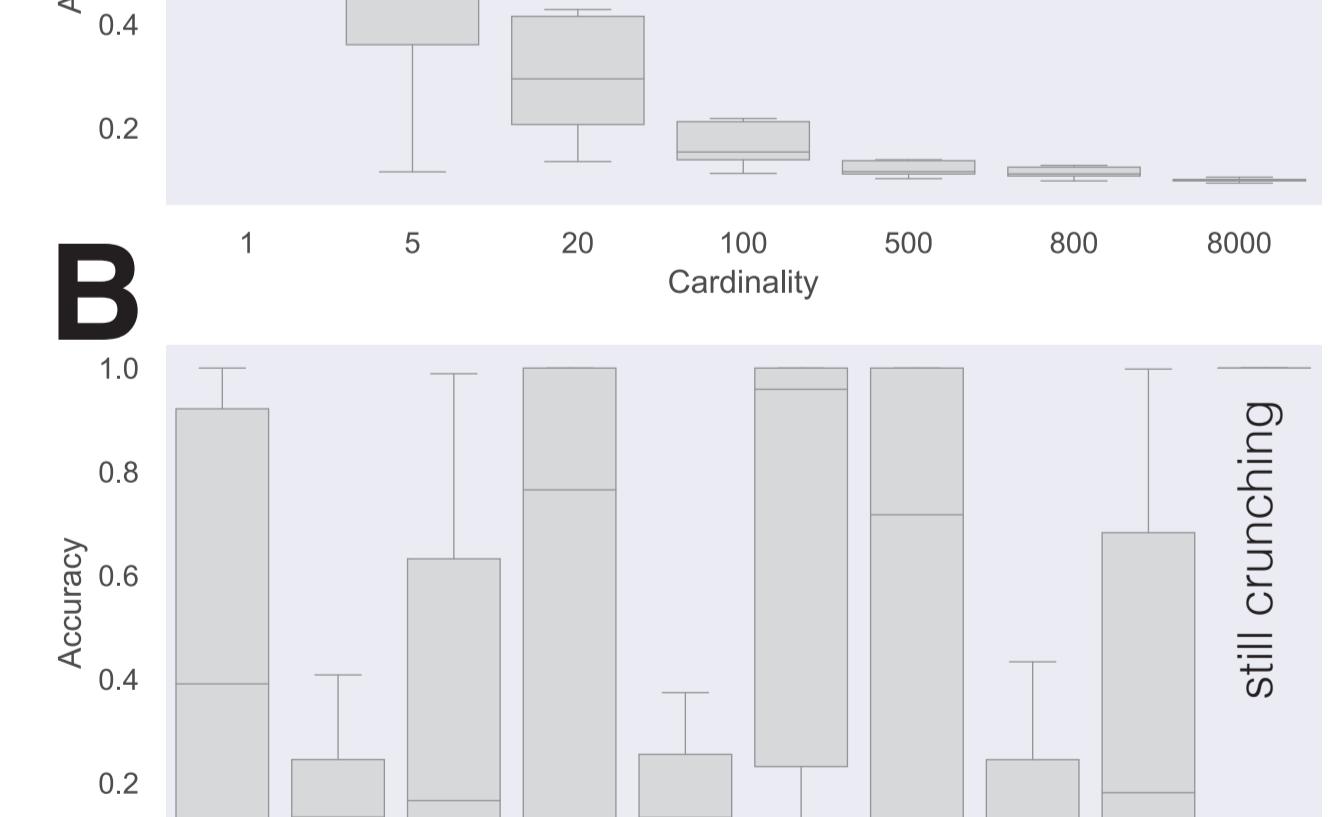
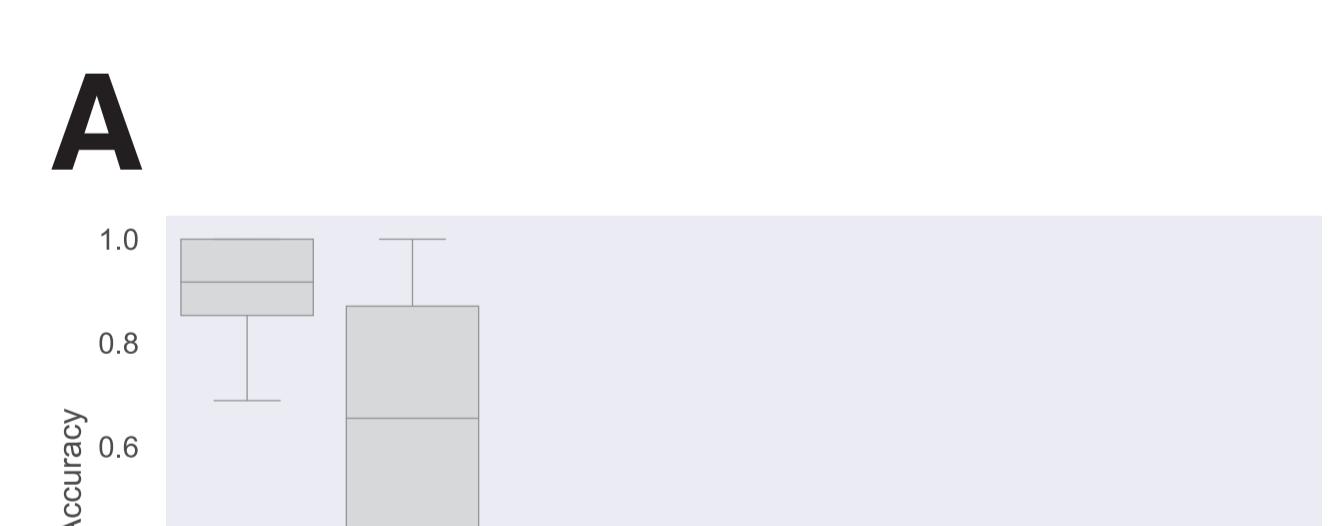
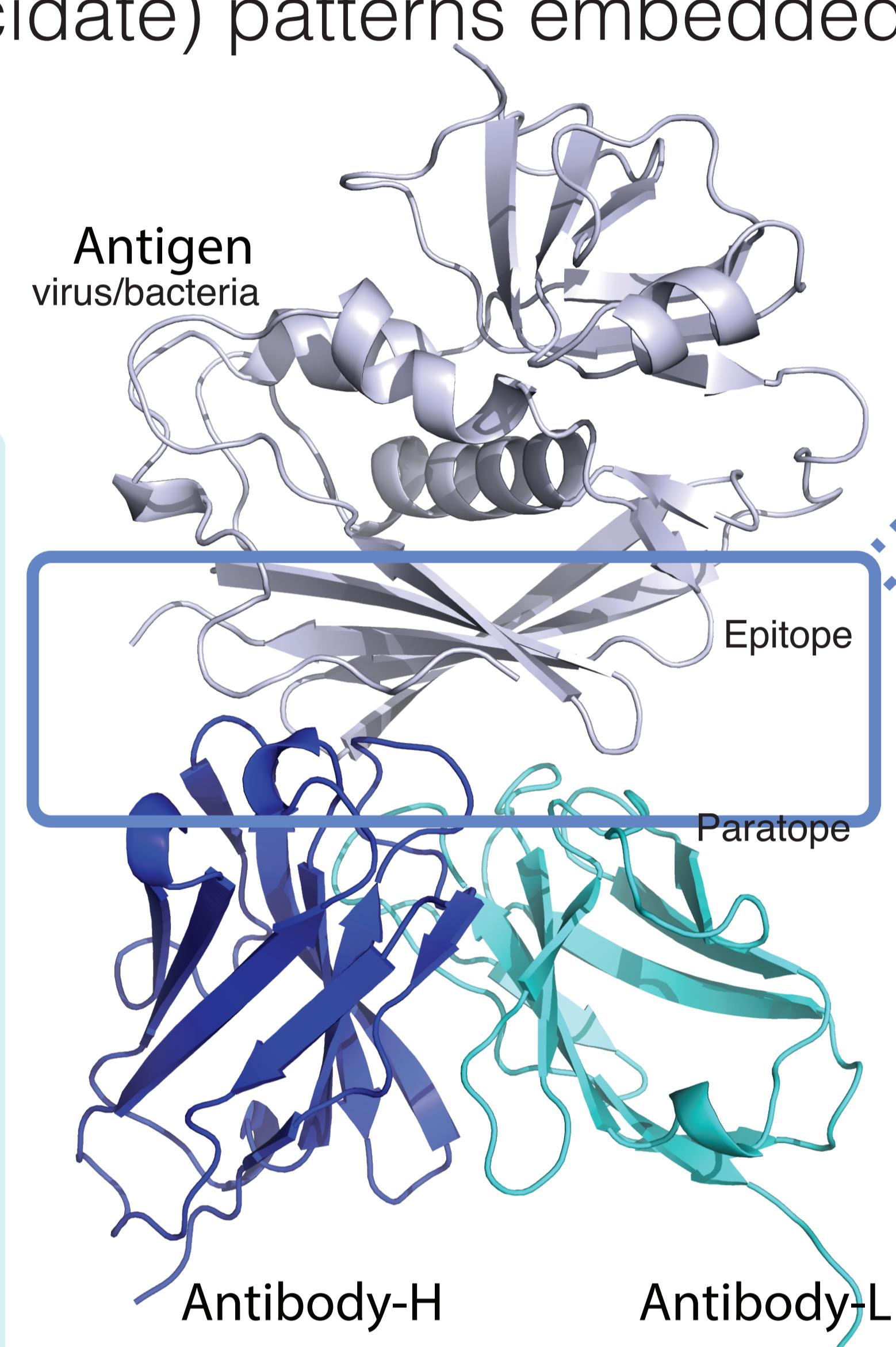
³ ETH Zürich, D-BSSE, Mattenstrasse 26, 4058, Basel, Switzerland

Motivation

Six out of ten top selling drugs are antibodies. Hence, accurate prediction of antibody specificity from the antibody sequence alone is of paramount importance for the conception of next-generation antibody therapeutics. Can machine learning help in such predictions? To what extent machine learning algorithms recover (elucidate) patterns embedded in antibody sequences?

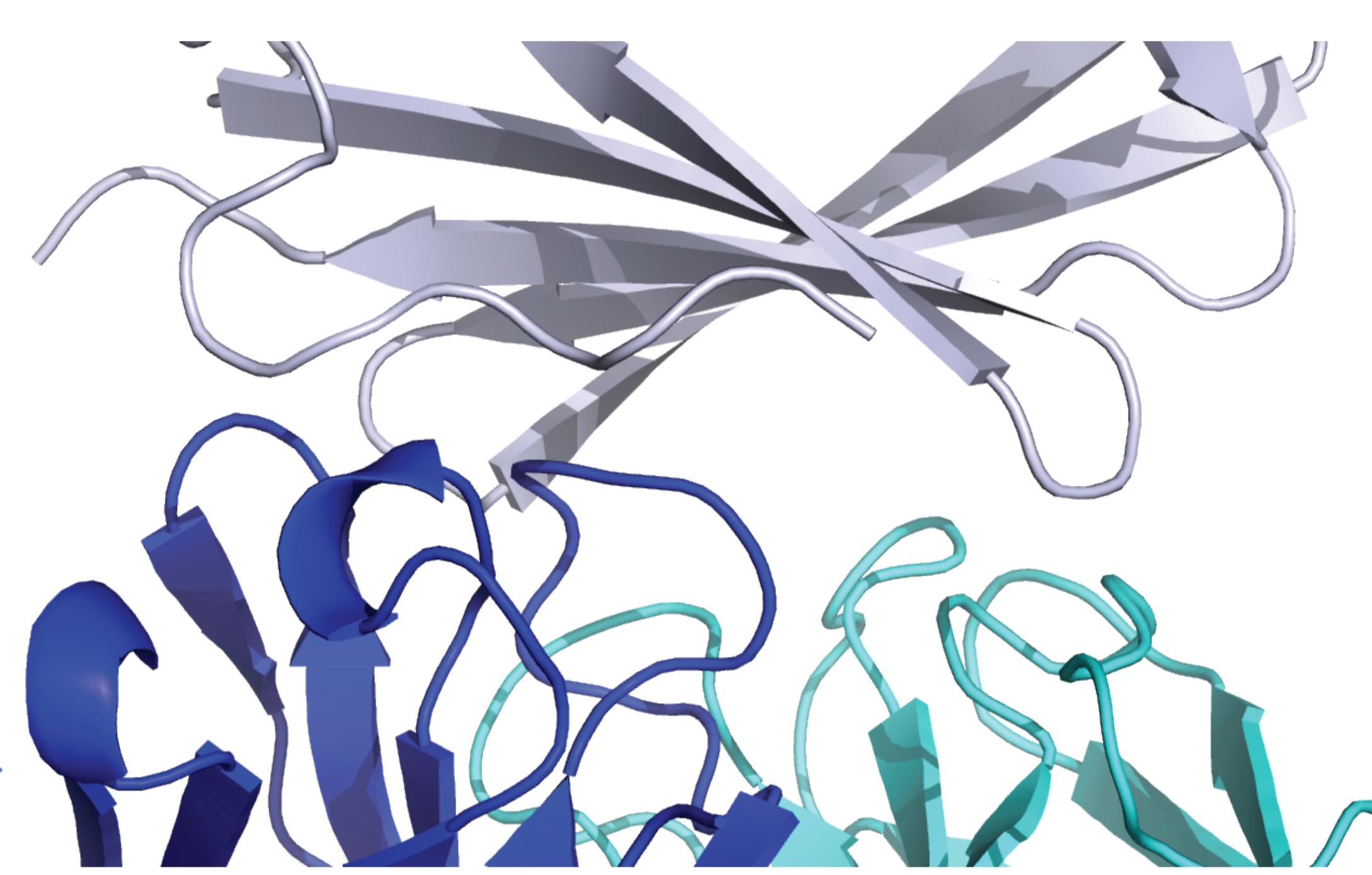
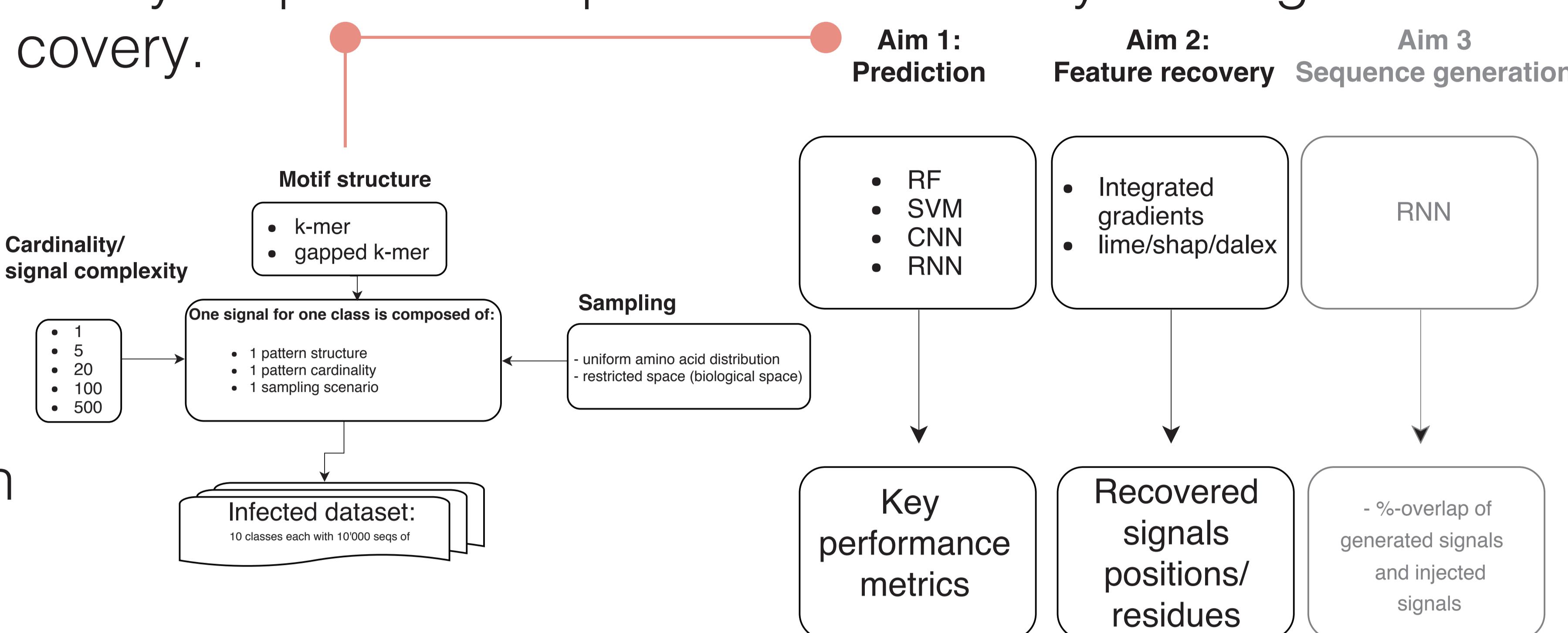
Results

- A.** Accuracy decays as signals get more complex. **B.** Models with non-linear decision functions are superior to their linear counterparts. **C.** Longer signals (k-mers) yield higher accuracy. **D.** Signals were successfully recovered using local linear approximations. **E.** Signal recovery decays as a function of signal complexity.



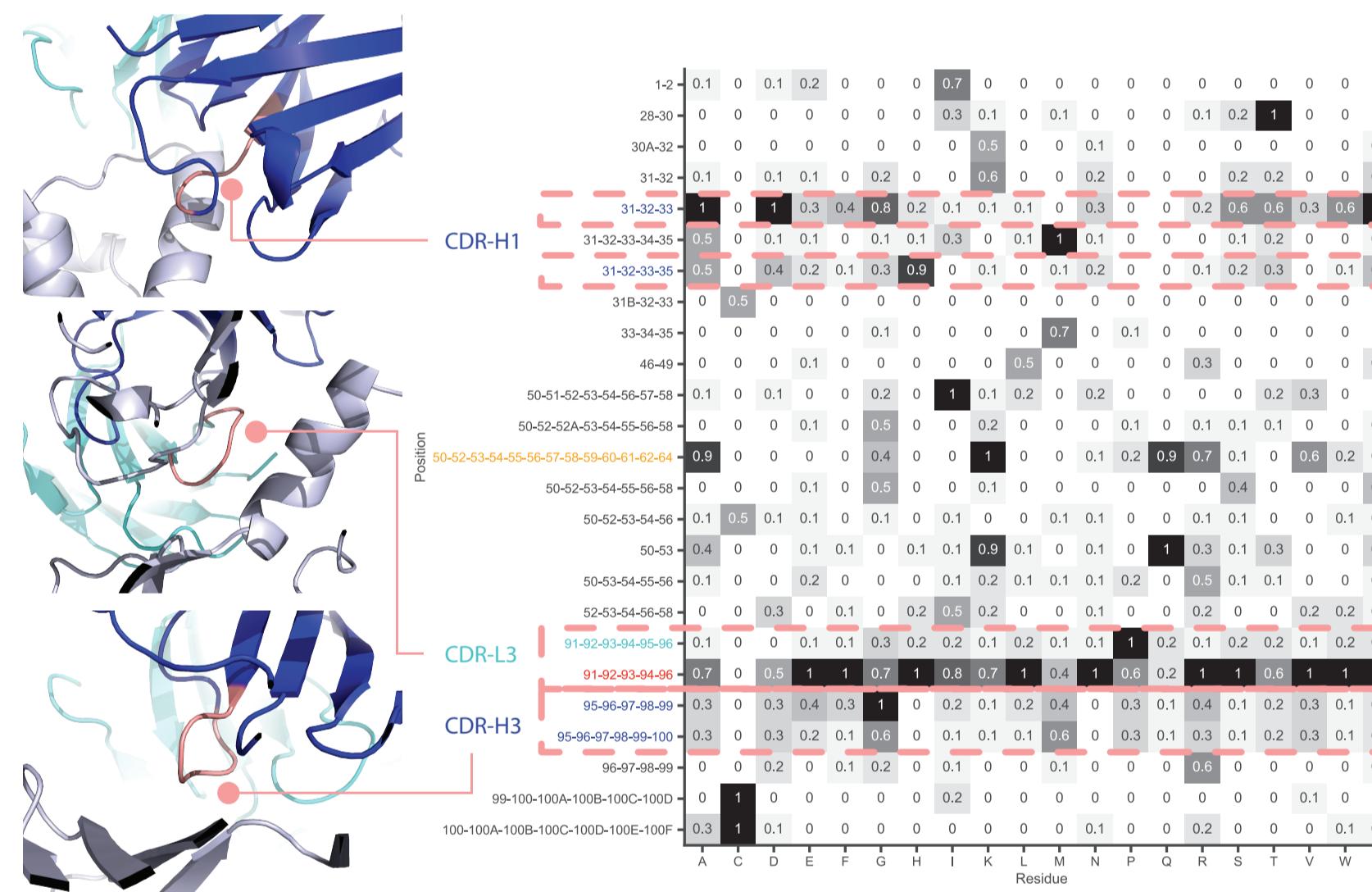
Approach and methods

Evaluate a set of machine learning models on synthetic antibody sequences for prediction accuracy and signal recovery.



ISRGSSLDYWI RADGSDISVGR

But what kind of signal? Our pre-liminary analyses

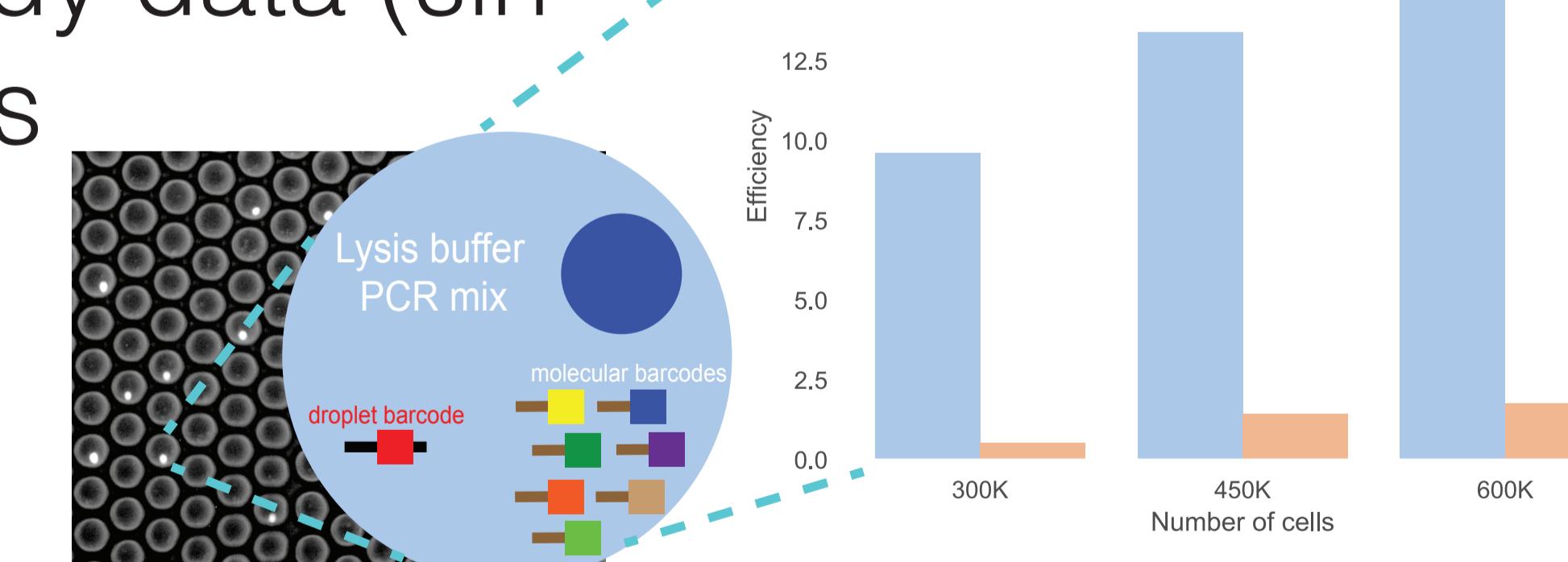


(i) Generate synthetic datasets with varying signal complexity. (ii) Feed the datasets to machine learning algorithms. (iii) Evaluate and recover key performance metrics and features.

But what kind of signal? Our preliminary analyses suggest that short and continuous (linear) residue-segments mediate the majority of antibody-antigen interactions.

Future works

Train and evaluate the models on single-cell resolution antibody data (single-cell immunogenomics and single antibody proteomics) for in-silico design of antibodies.



Concluding remarks

- Non-linear models are superior to models with linear decision boundaries (at least for this kind of datasets).
- Signal complexity dictates the performance of the models.
- Signal complexity dictates recovery.
- Interacting segments (signals) tend to be short and continuous.