



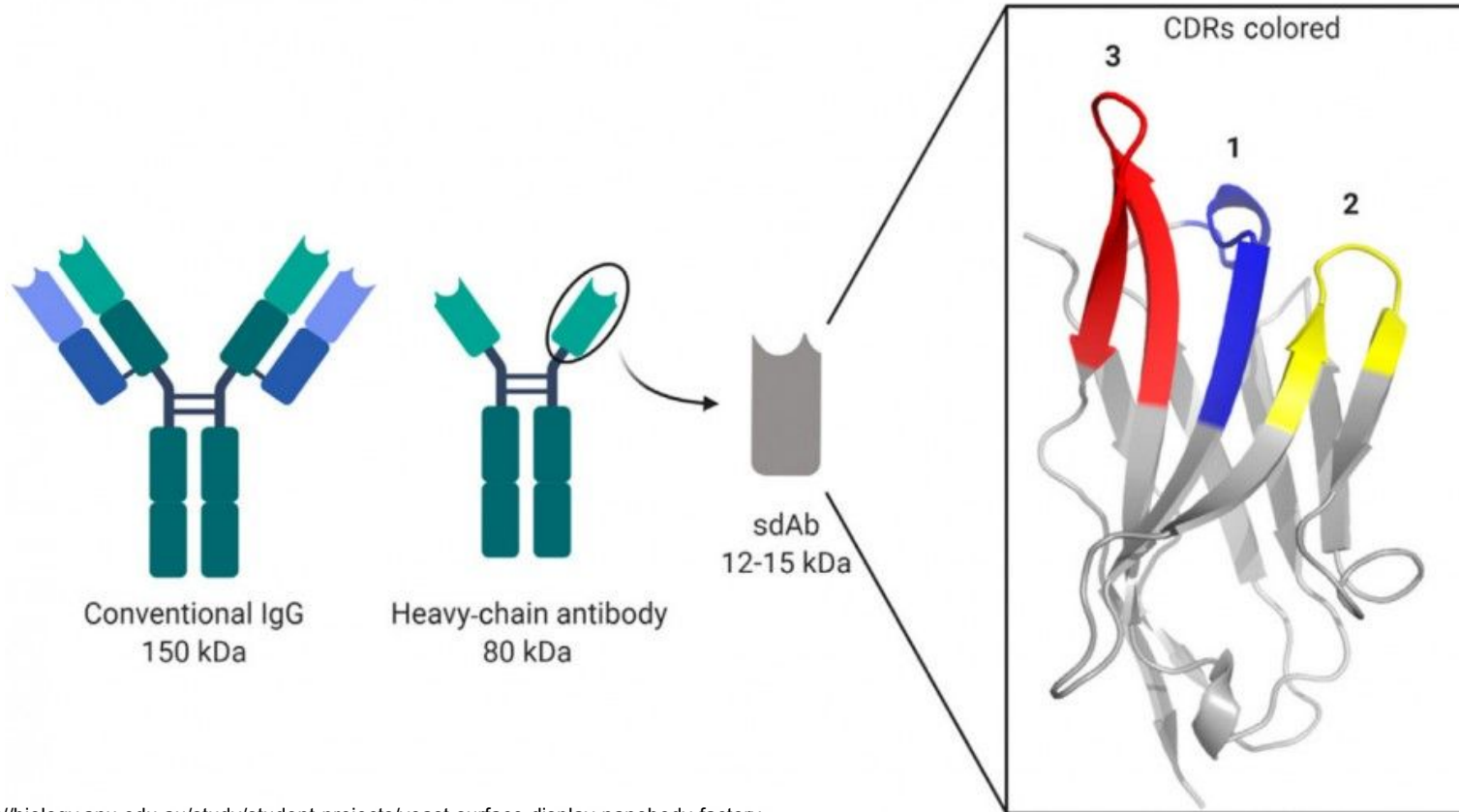
Model validation strategy is  
crucial to properly estimate the  
real-life applicability of machine  
learning models



Earth



# Nanobodies are small antibodies



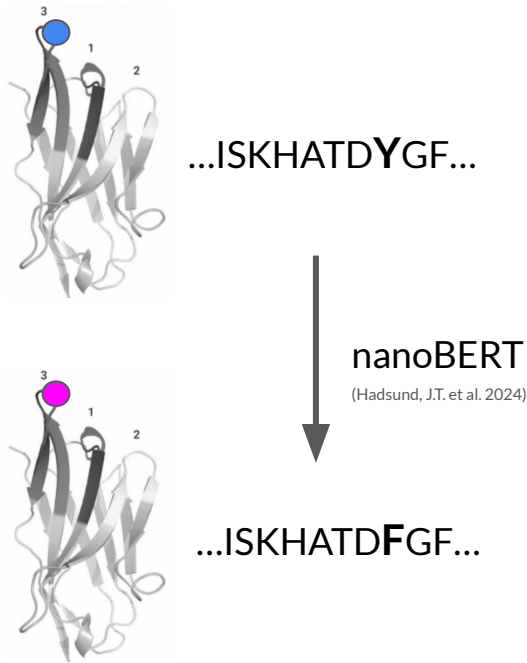
Can we find alternative nanobodies that are more stable than wild type?

# Suggested pipeline:



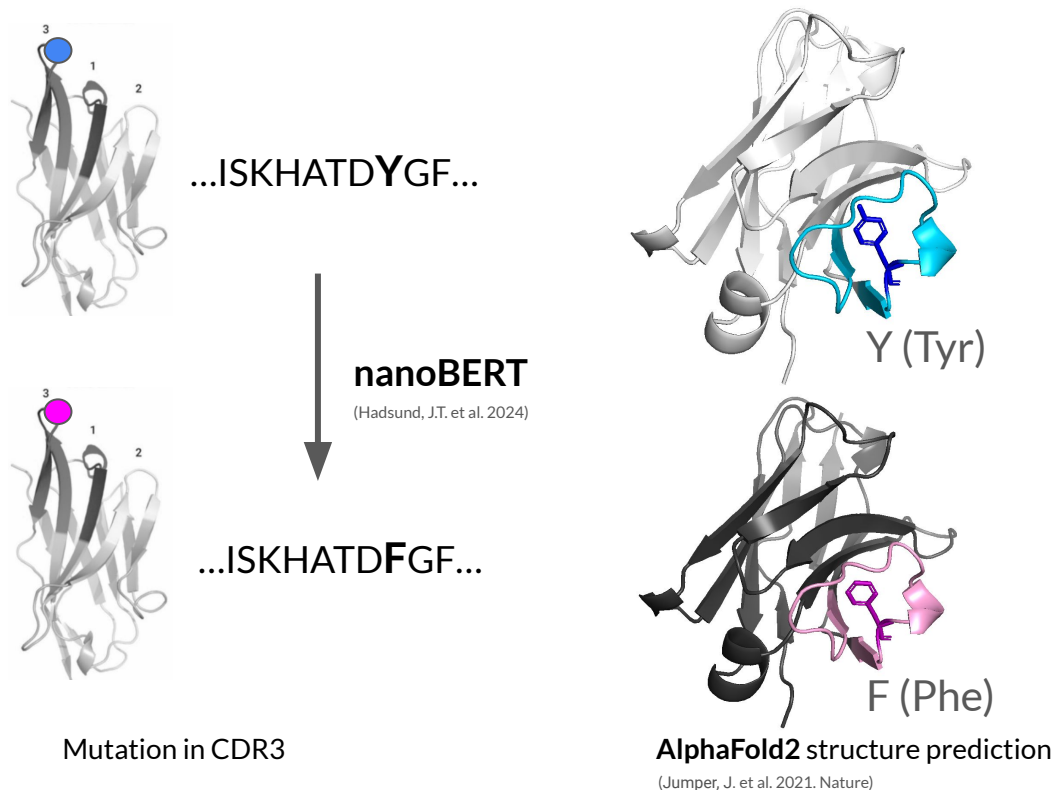
...ISKHATD**Y**GF...

# Suggested pipeline: find functional mutated sequences

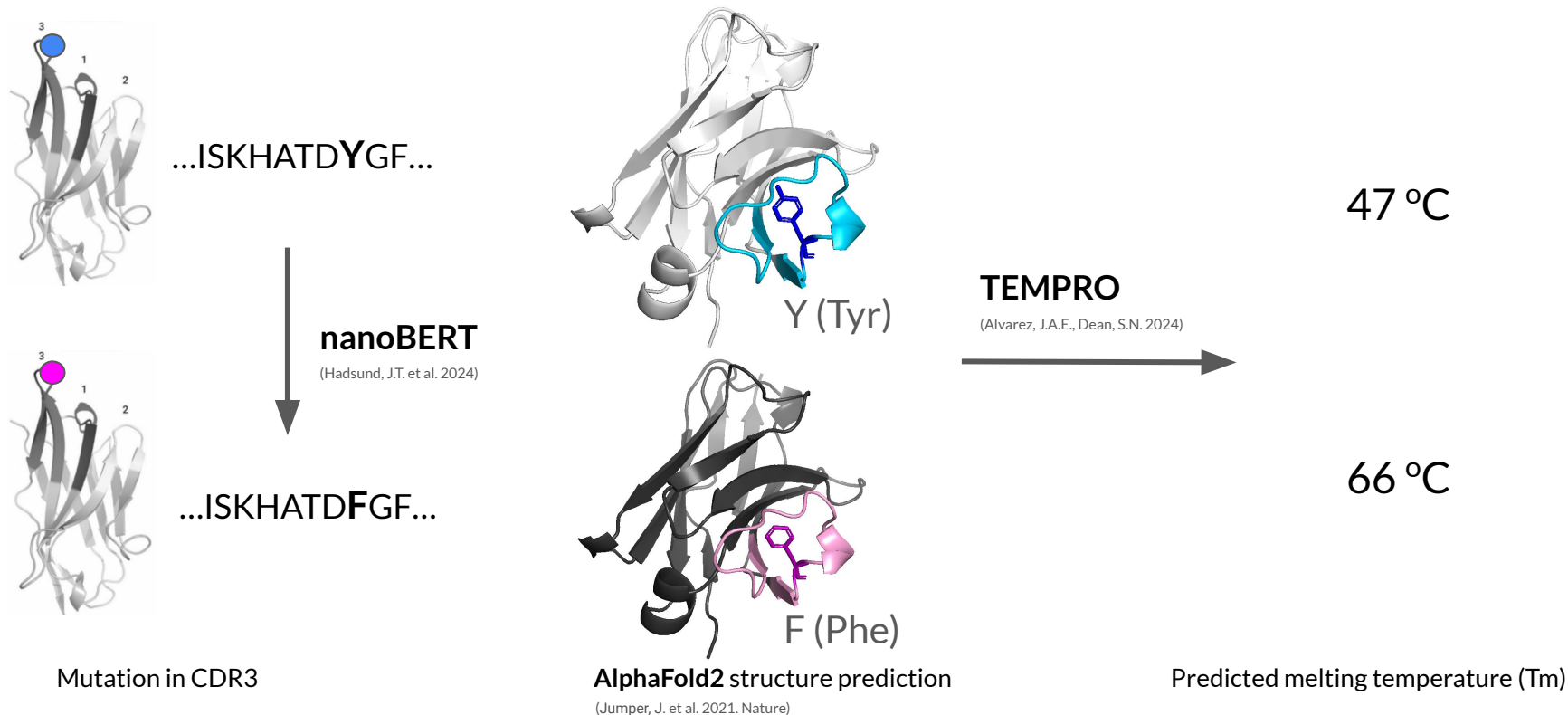


Mutation in CDR3

# Suggested pipeline: predict 3D structure



# Suggested pipeline: predict melting temperature



But do we trust this model?



# Sensitivity analyses and uncertainty estimates are crucial to trust model predictions.

dataset: NbThermo (Valdés-Tresanco et.al. 2023)

TEMPRO (Alvarez, J.A.E., Dean, S.N. 2024)

splitting randomly

80% training

20% validation and testing

splitting randomly

80% training and validation

20% testing

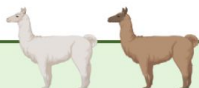
80% training

20% validation

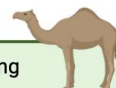
splitting by species



80% training and validation



20% testing



80% training

20% validation

splitting by sequence homology

80% training and validation

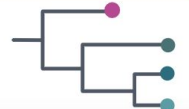


20% testing



80% training

20% validation



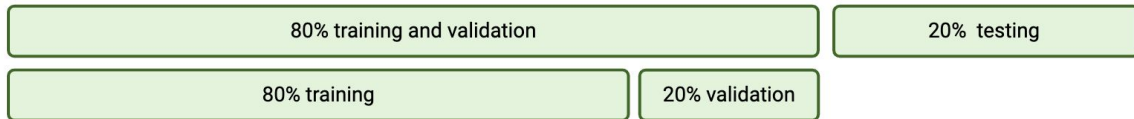
# Sensitivity analyses and uncertainty estimates are crucial to trust model predictions.

dataset: NbThermo (Valdés-Tresanco et.al. 2023)

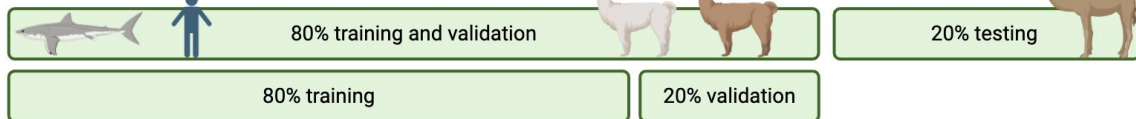
TEMPRO (Alvarez, J.A.E., Dean, S.N. 2024) **splitting randomly**



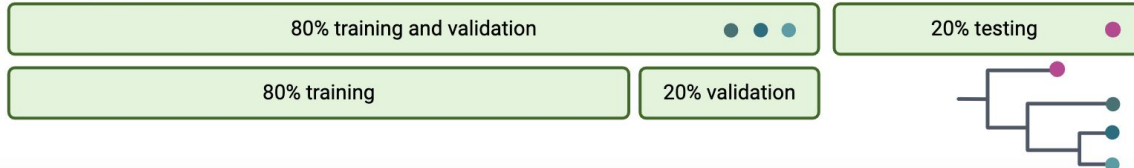
**splitting randomly**



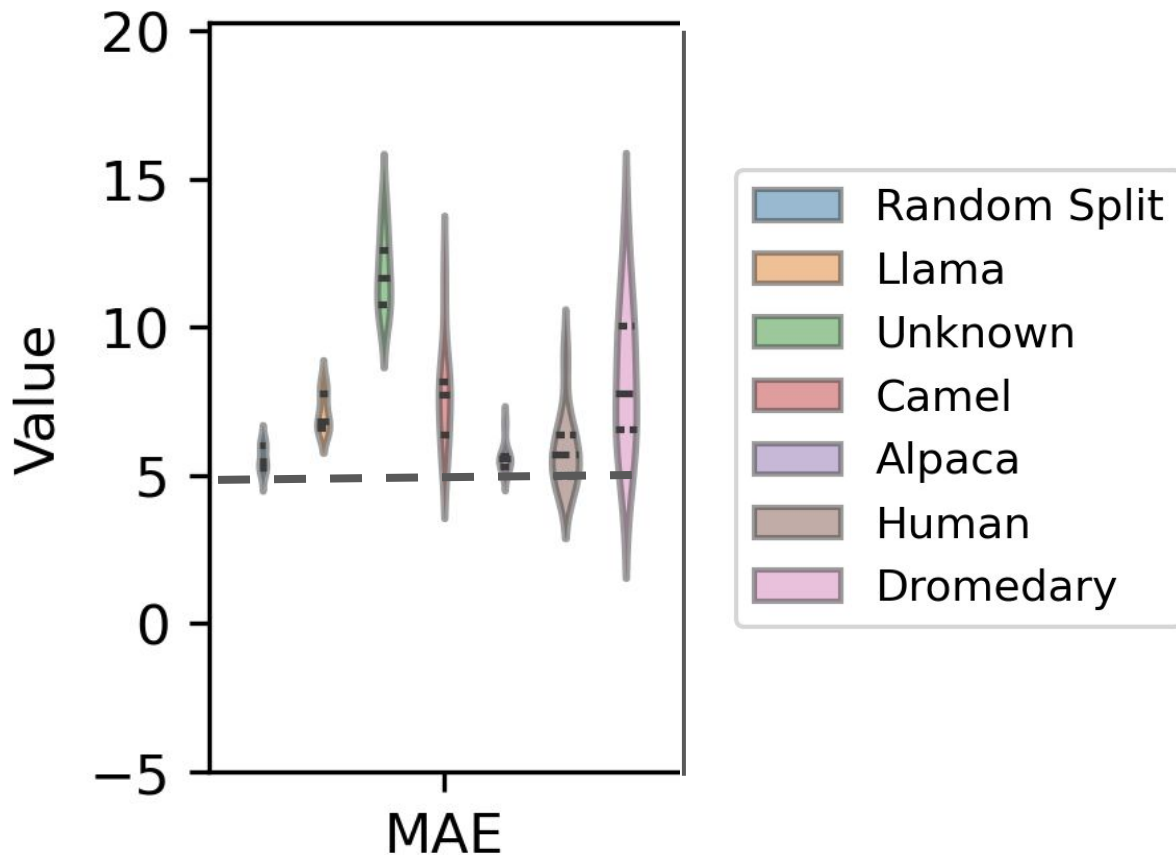
**splitting by species**



**splitting by sequence homology**

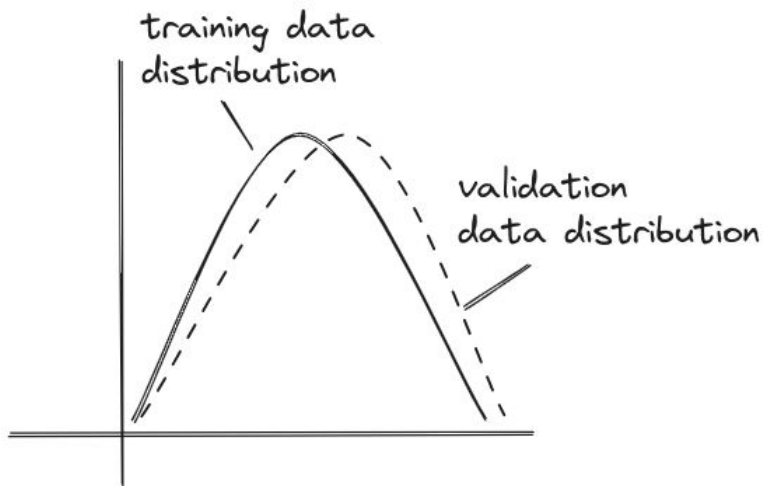


# Random split validation underestimates the error for novel species

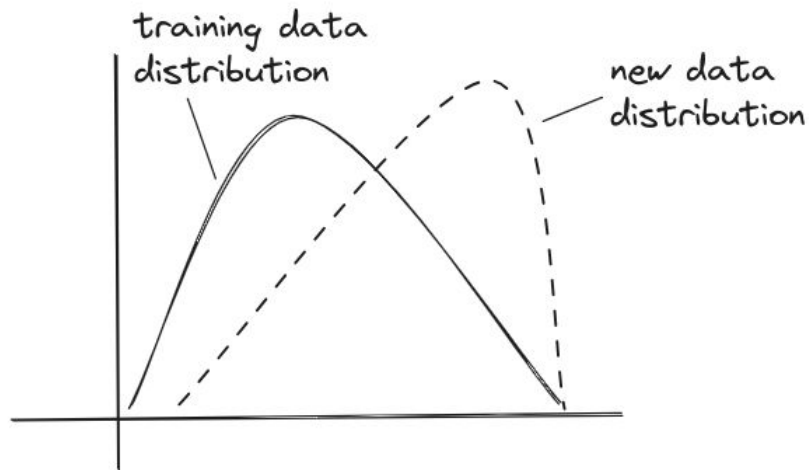


# Take-home message

Training: great performance



Inference: wrong predictions





# Thank you!

**Contributors:**

Finnja Becker

Madalina Giurgiu

Dr. Anne Hartebrodt

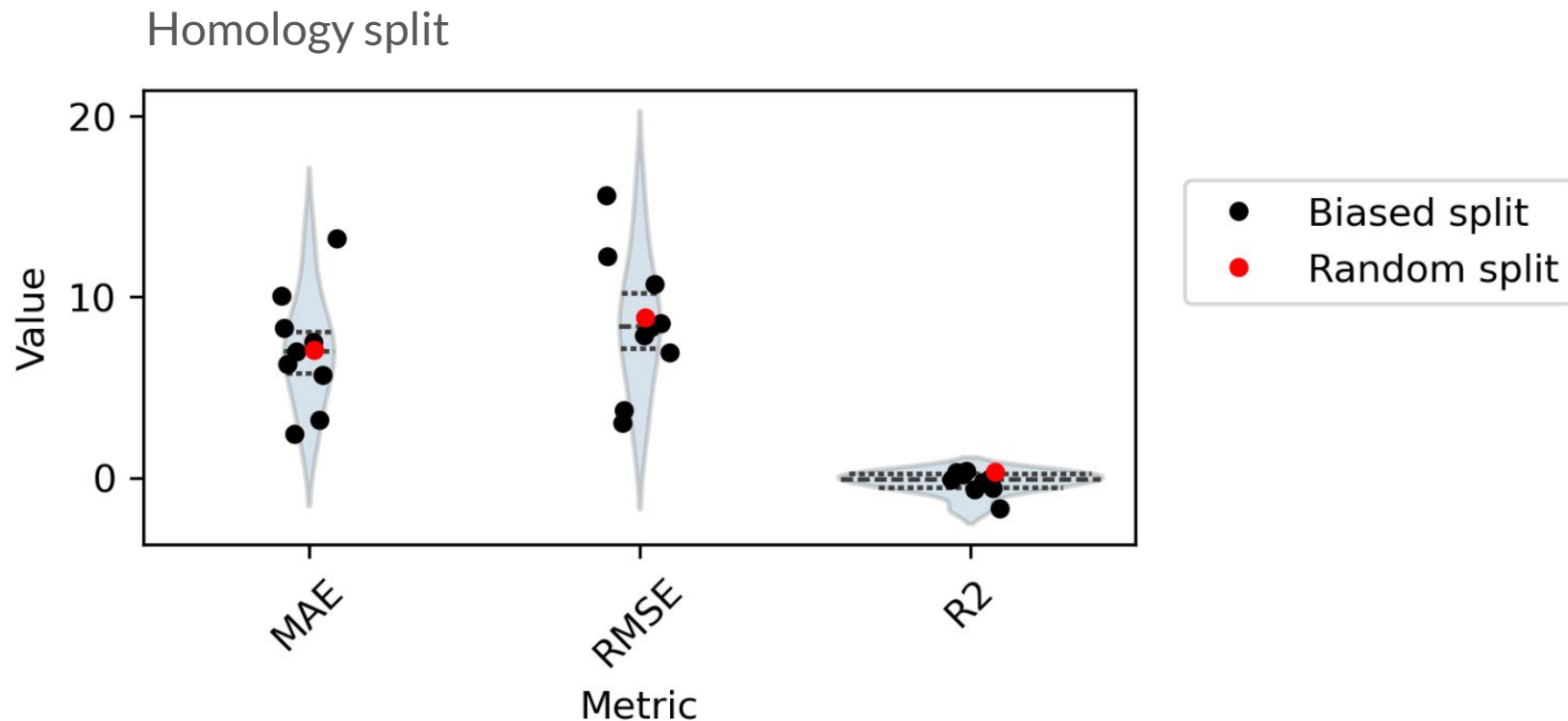


Project advisor: Dr. Ni Fang

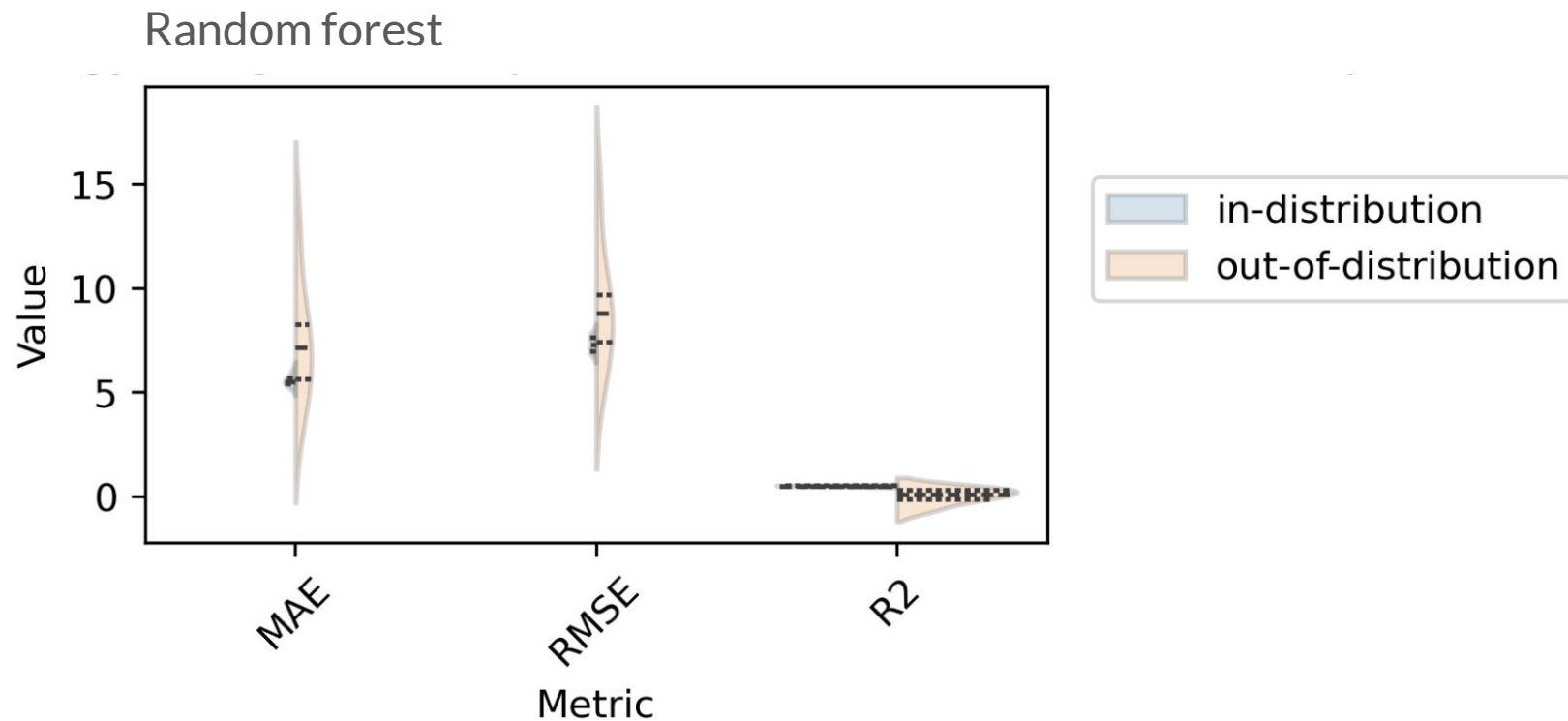
Github repository: <https://github.com/AnneHartebrodt/earth-ml-sensitivity>

# Backup

# Validation strategy underestimates the error for novel species and non-homologous sequences



# Baseline ML model accuracy

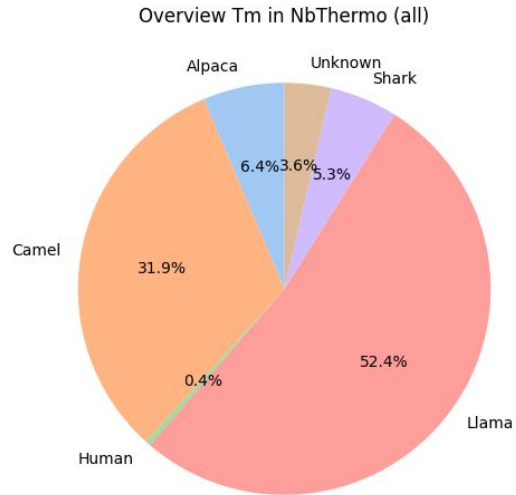




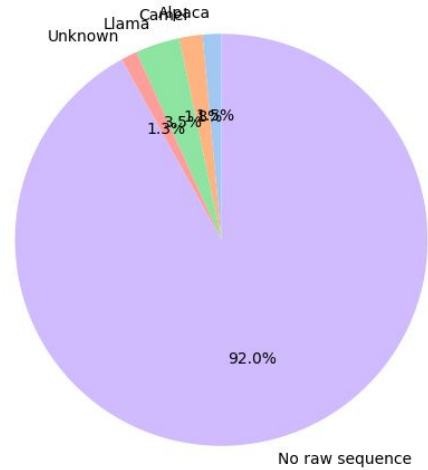
# Model improvement ideas

- Ensemble classifier
- More data
- Train different models (e.g. protein melting temp. models) on nanobody data
- Fine-tuning (use protein data to train larger model, fine-tune this model using the limited nanobody data)

# NbThermo



Overview Tm in NbThermo data with sequence



# NbThermo:

