

Idea

Neural Networks have shown to be a viable tool for prediction of bioactivities. However, bioactivity datasets tend to be very small and imbalanced.

The problems:

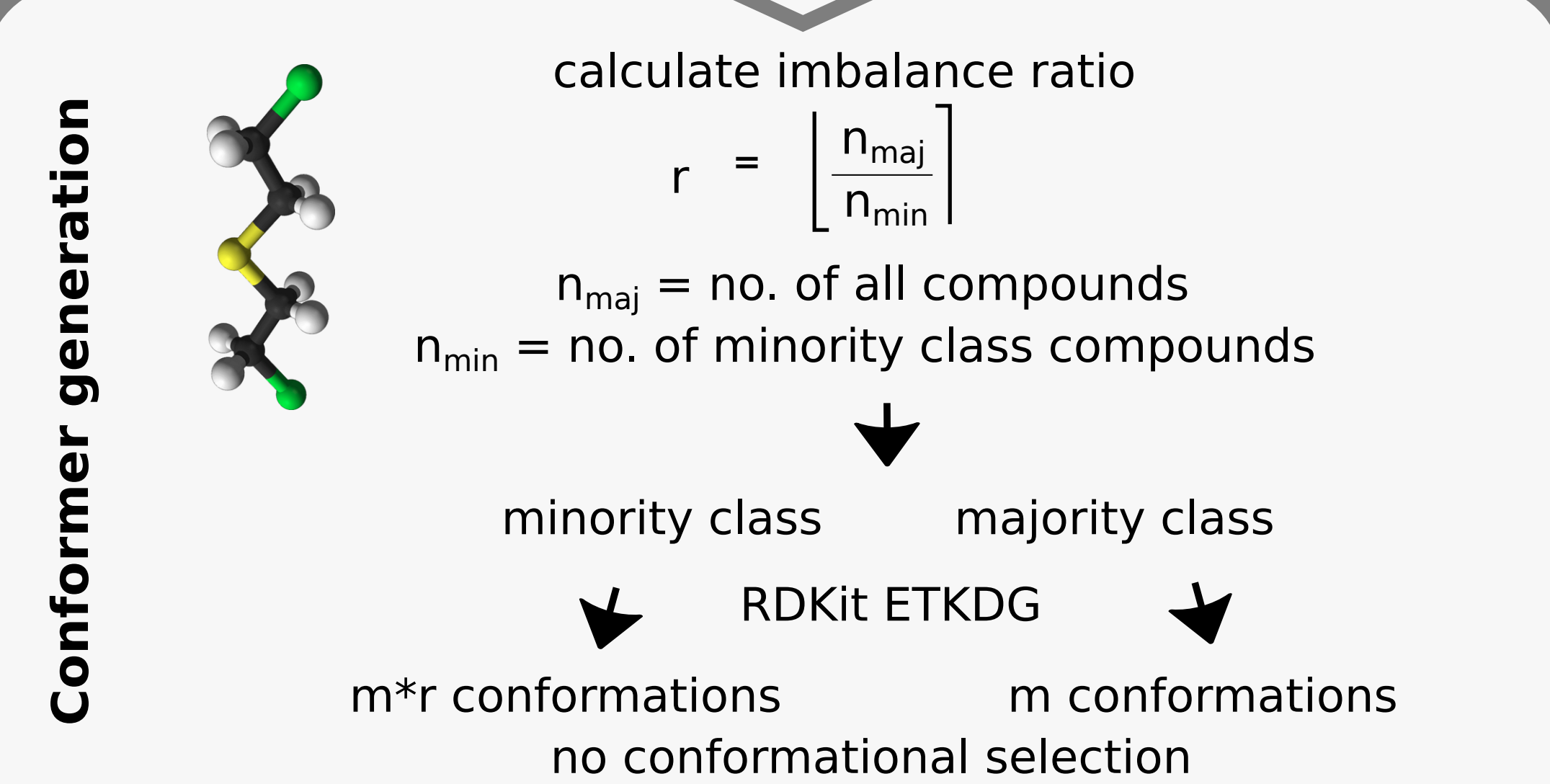
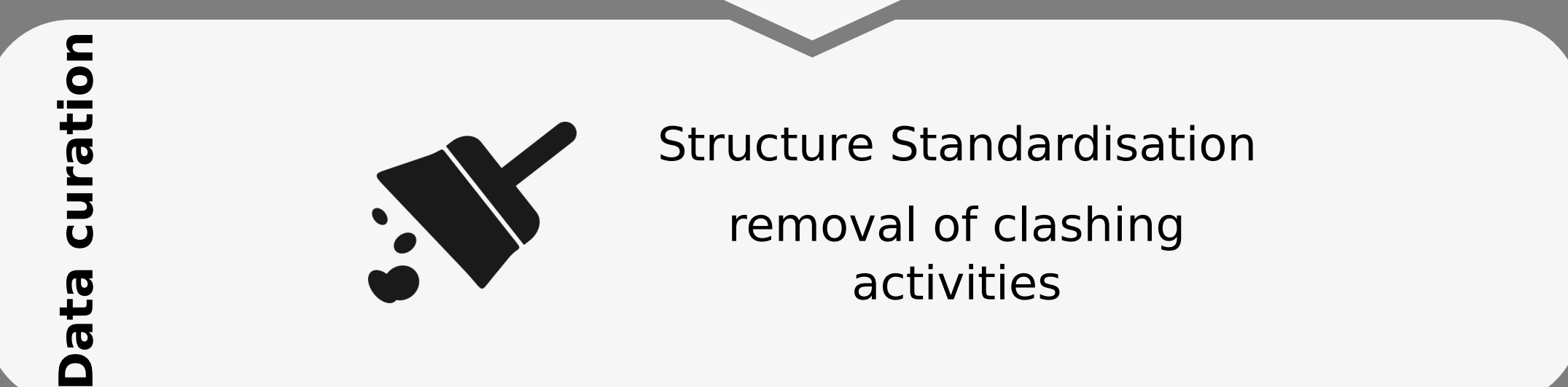
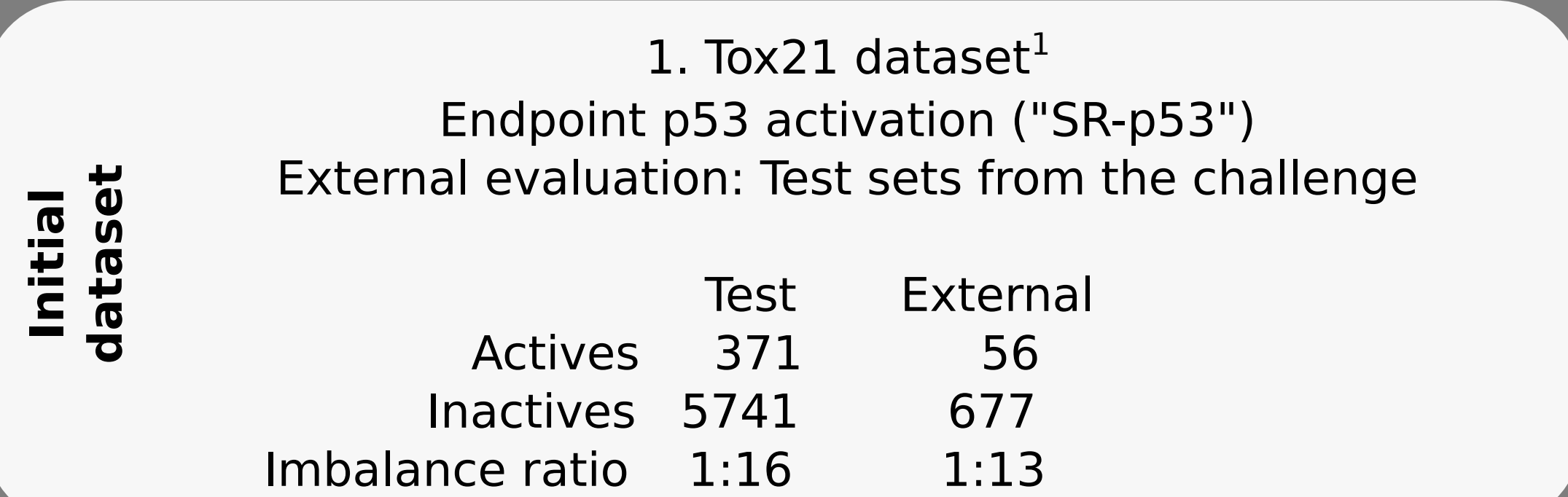
- (i) Neural networks overfit easily with small datasets.
- (ii) Neural networks tend to ignore the minority class for imbalanced datasets.

Proposed solution:

Augmentation helps neural network training for image recognition, can we use that for molecules?

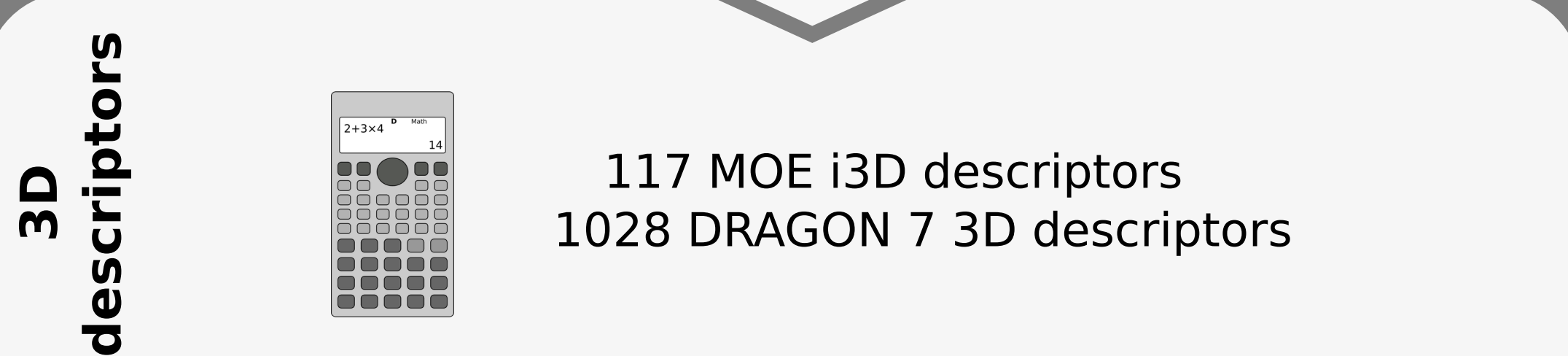
Methods

COVER Workflow

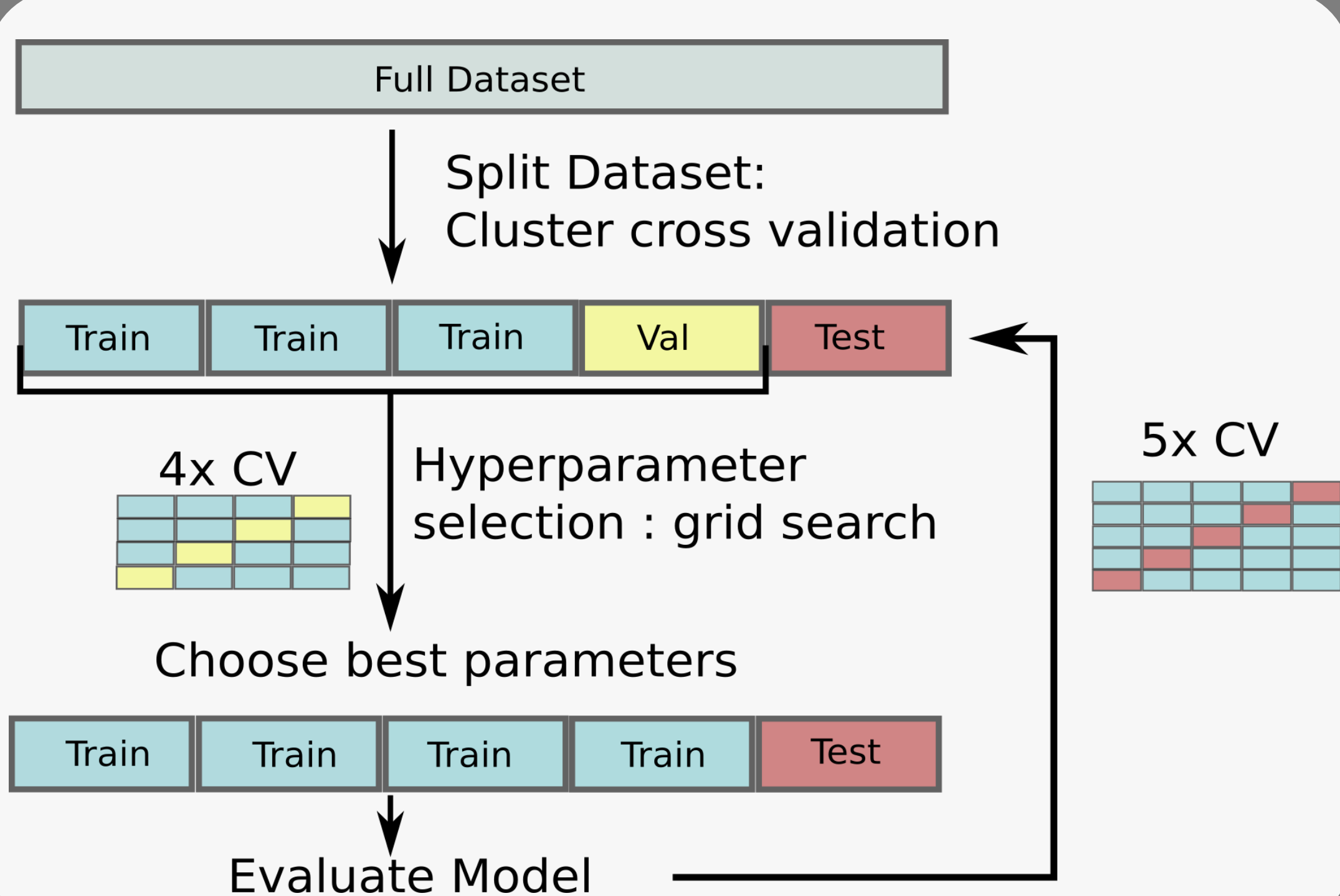


Training datasets

dataset	no. of conformations per		no. of molecules		
	inactive	active	inactive	active	overall
1-1	1	1	5502	341	5843
1-16	1	16	5502	5428	10930
2-2	2	2	11001	680	11681
2-32	2	32	11001	10865	21866
5-5	5	5	27504	1698	29202
5-80	5	80	27504	27145	54649

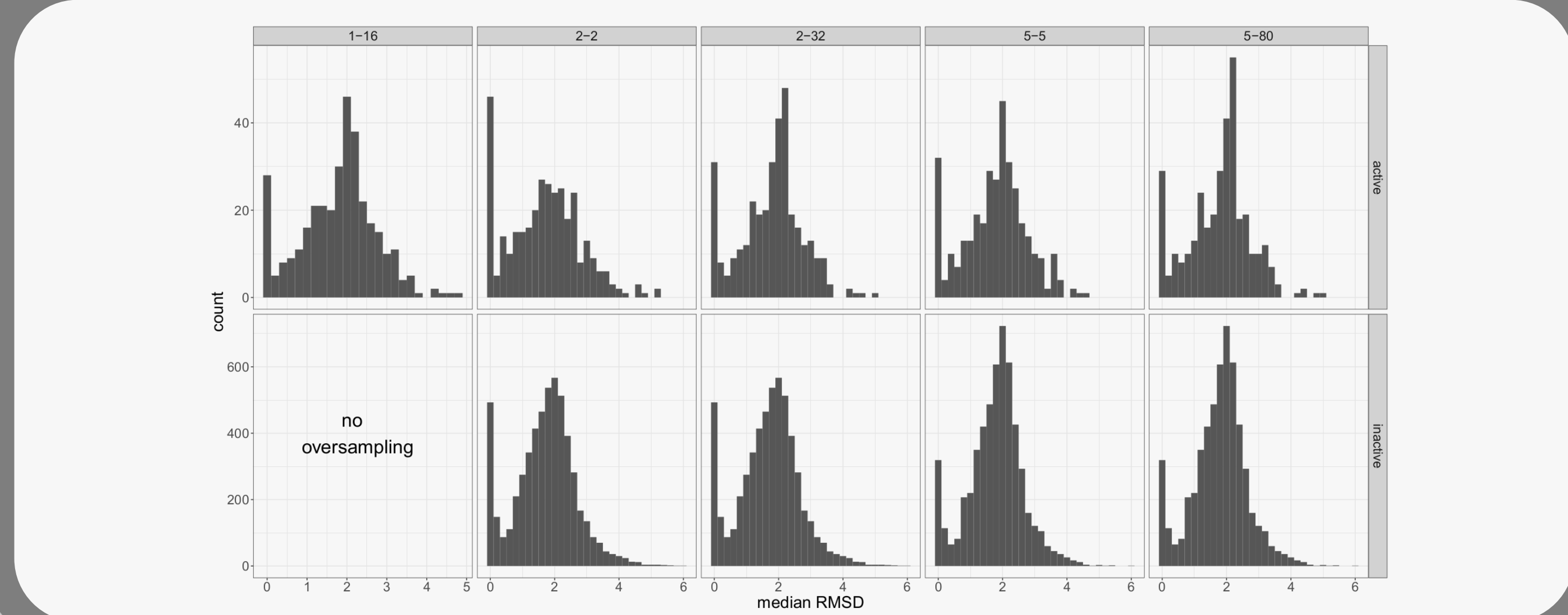


Training procedure

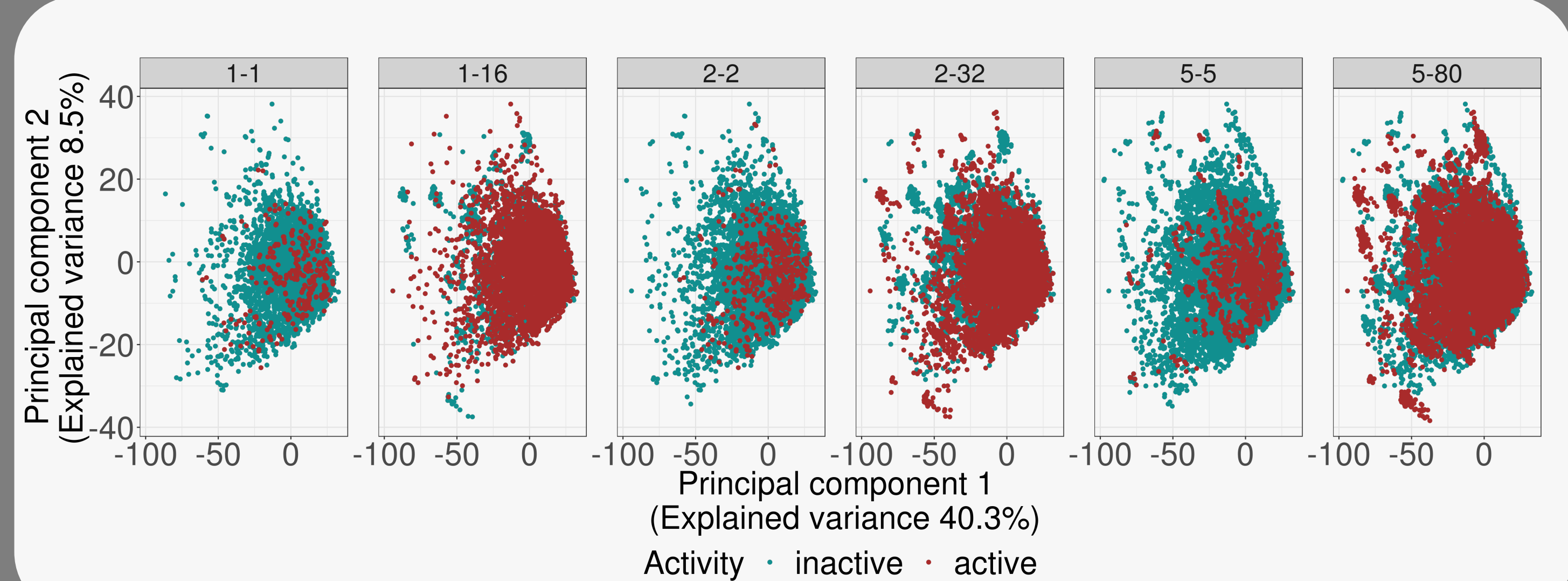


Results

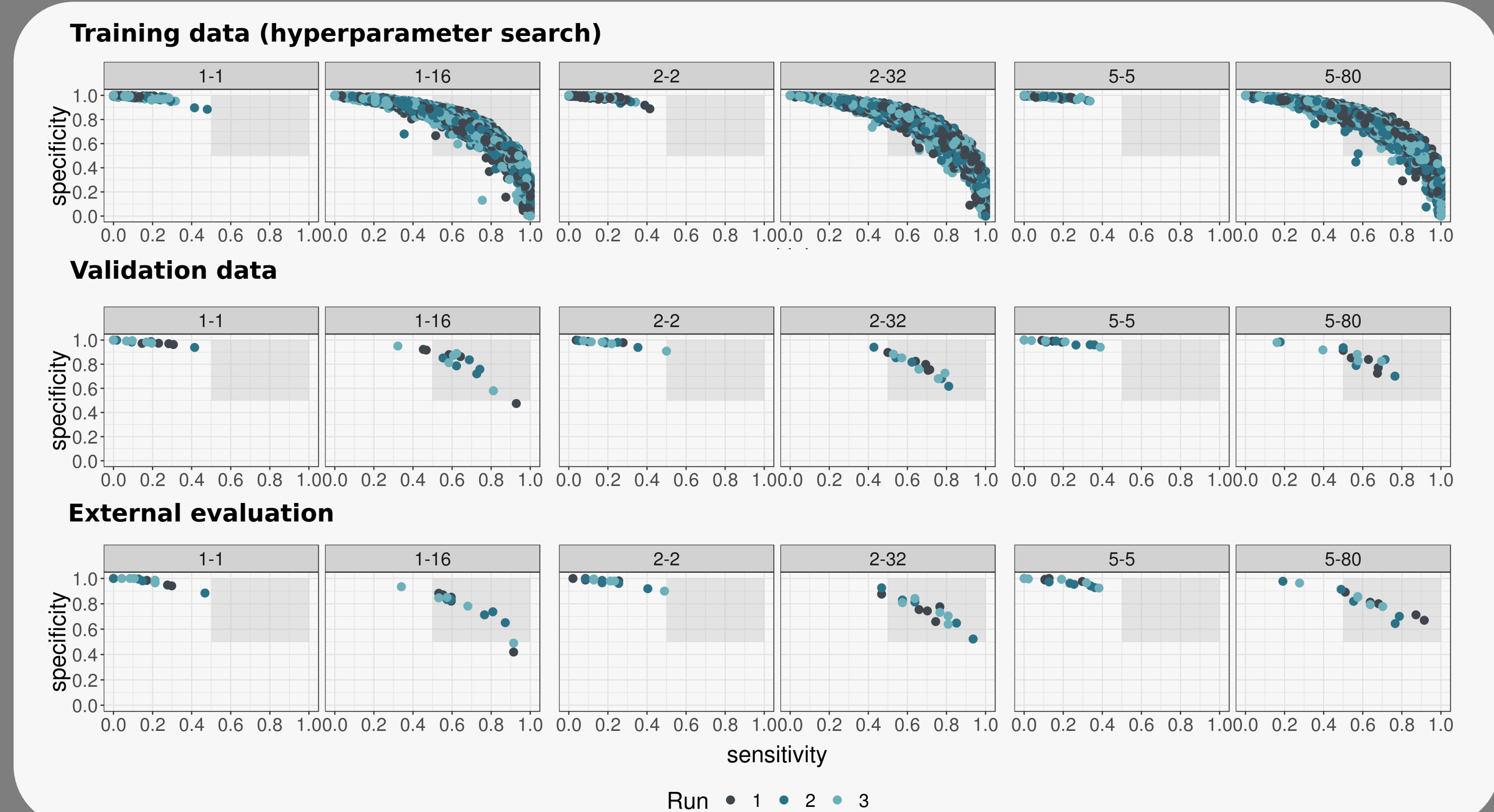
Comparison of median RMSDs



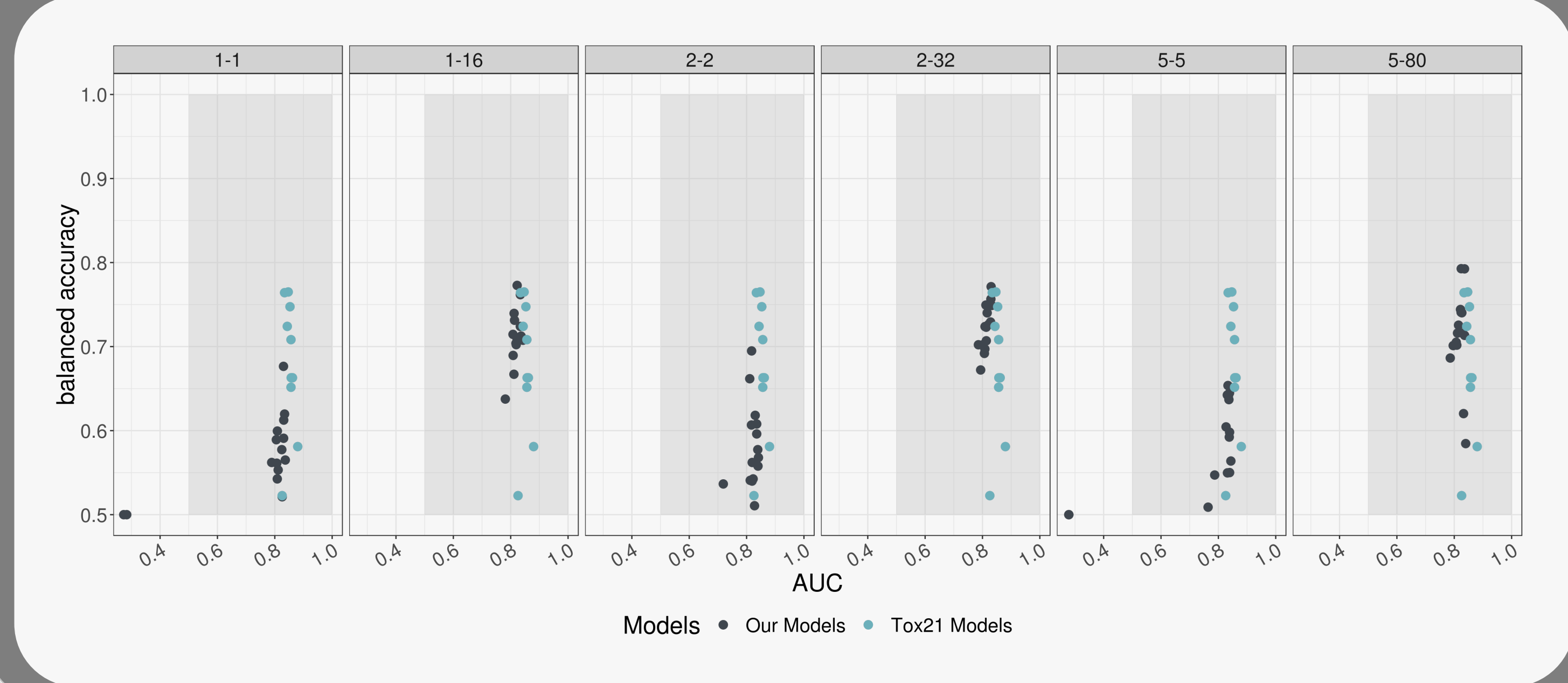
PCA of the oversampled datasets



Model performances



Comparison to Tox21 top 10 models



Conclusion

Did we find a solution?

Yes, COVER helps to increase the training performance, especially with respect to sensitivity

But:

Oversampling alone does not increase the performance.

Balancing is necessary to increase the performance of the models.

Acknowledgements

This project has received funding from the Innovative Medicines Initiative 2 Joint Undertaking under grant agreement No 777365 ("eTRANSafe"). This Joint Undertaking receives support from the European Union's Horizon 2020 research and innovation programme and EFPIA. It also has received funding from the Austrian Science Fund FWF (grant W1232).

References

¹<https://tripod.nih.gov/tox21/challenge/>