
Imputation of assay bioactivity data using deep learning



Intellegens



Tom Whitehead

Ben Irwin, Peter Hunt, Matt Segall, Gareth Conduit



Unique deep learning algorithm

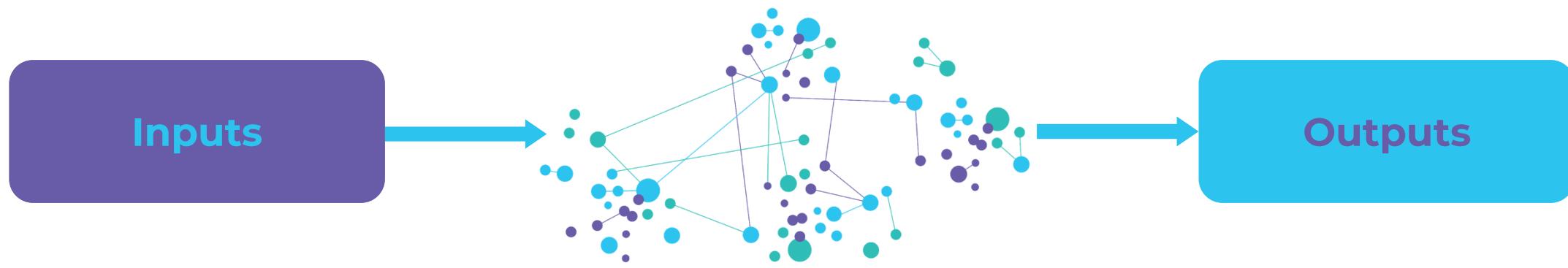
Utilise chemical descriptors, assay bioactivities, and simulations **in combination**

Understand and exploit **uncertainties** and noise to improve confidence in predictions

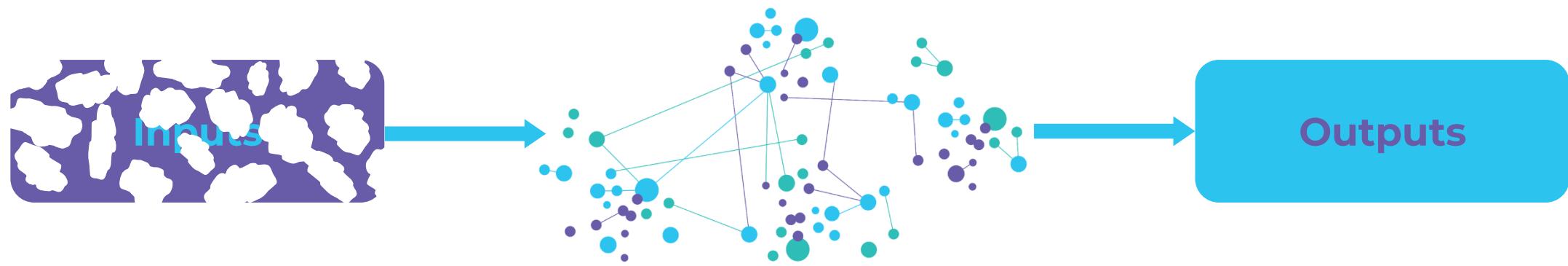
Broadly applicable algorithm with **proven** applications in drug design, materials discovery, patient analytics, ...



Deep learning

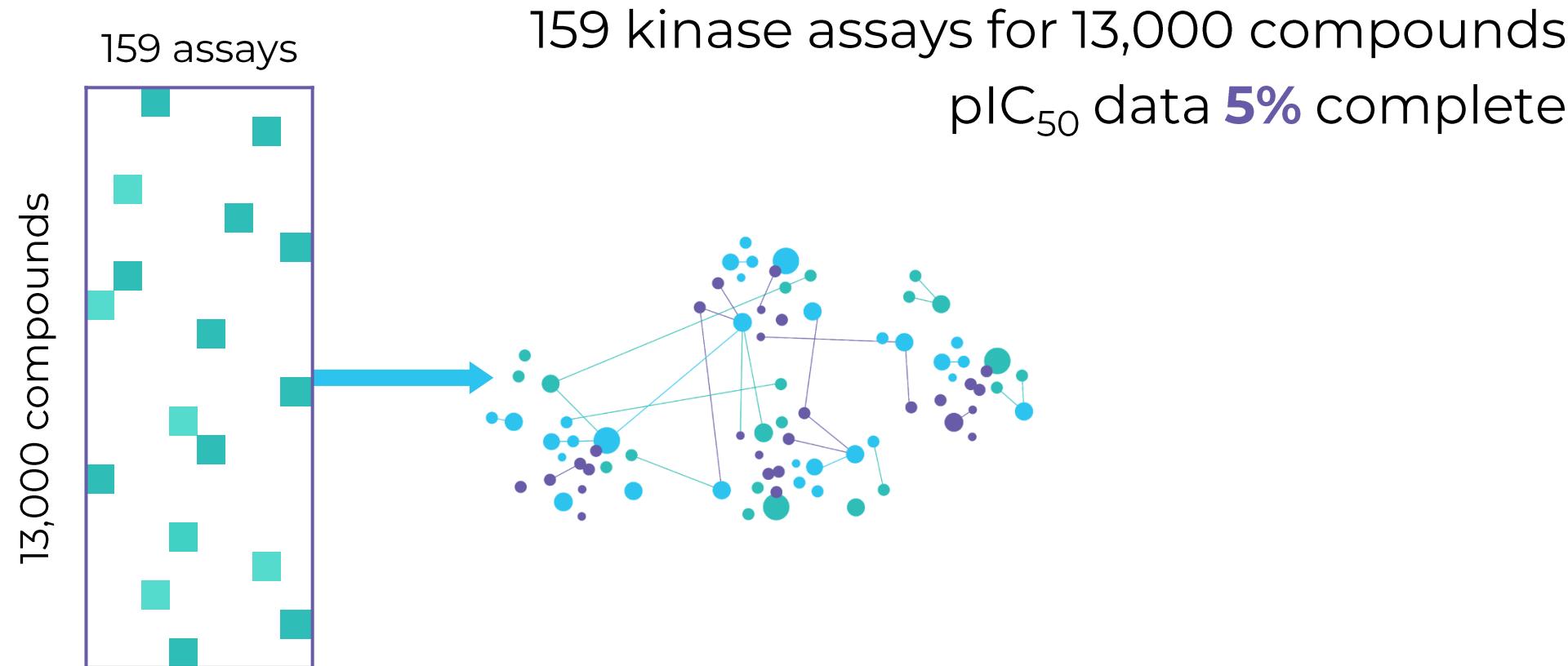


Alchemite™ deep learning





Novartis dataset to benchmark machine learning



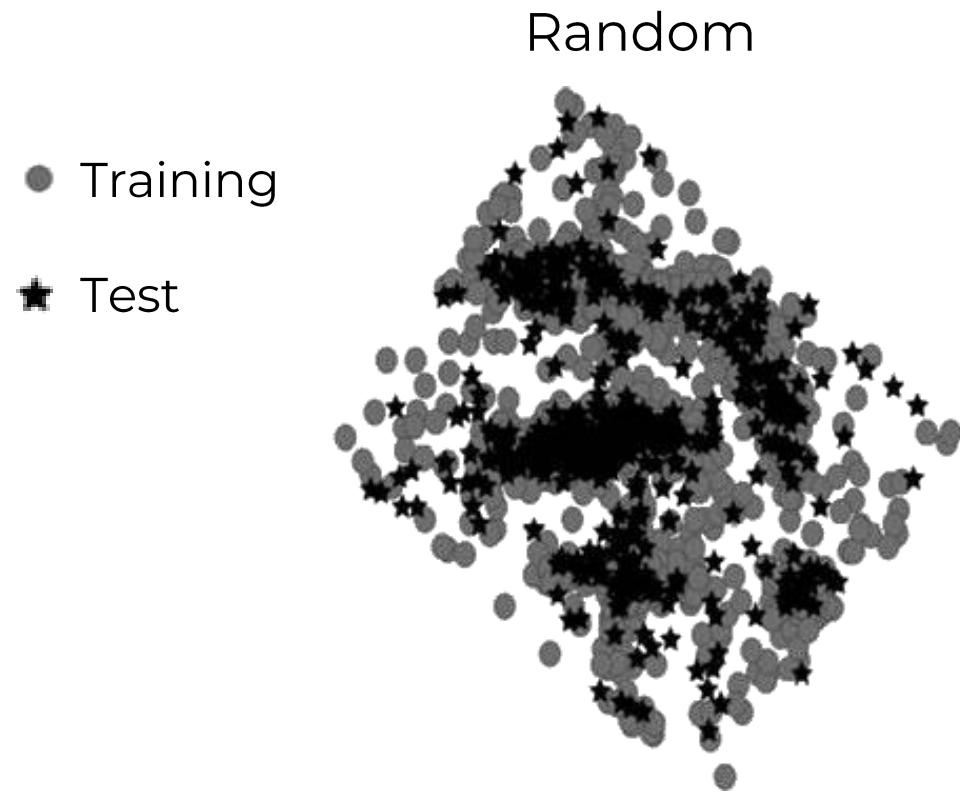
Data from ChEMBL

Martin, Polyakov, Tian, and Perez, J. Chem. Inf. Model. 57, 2077 (2017)

intellegens.ai



Novartis dataset distribution



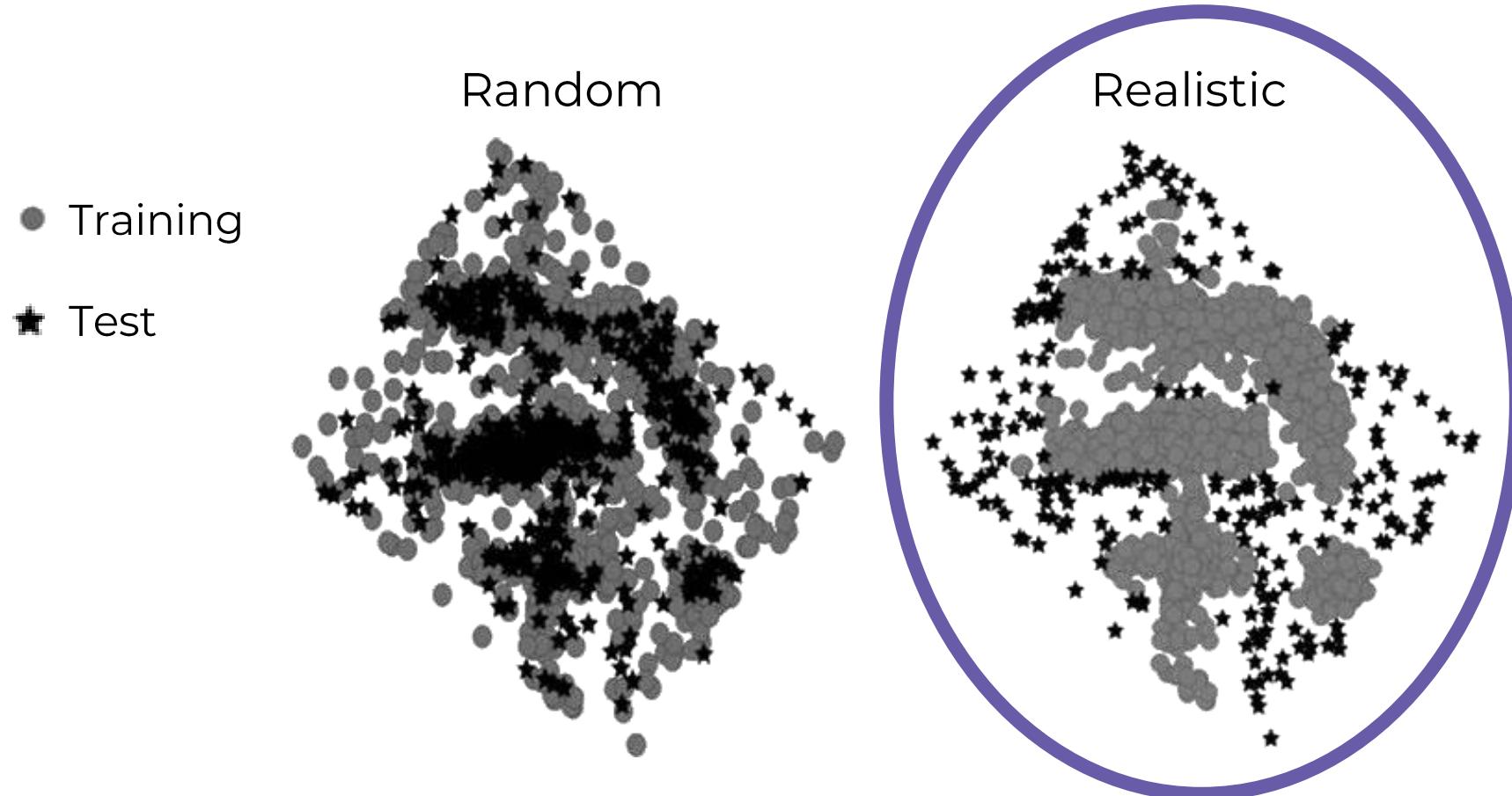
Data from ChEMBL

Martin, Polyakov, Tian, and Perez, J. Chem. Inf. Model. 57, 2077 (2017)

intellegens.ai



Novartis dataset is realistically distributed



Data from ChEMBL

Martin, Polyakov, Tian, and Perez, J. Chem. Inf. Model. 57, 2077 (2017)

intellegens.ai



Accuracy metrics

Coefficient of Determination, R^2

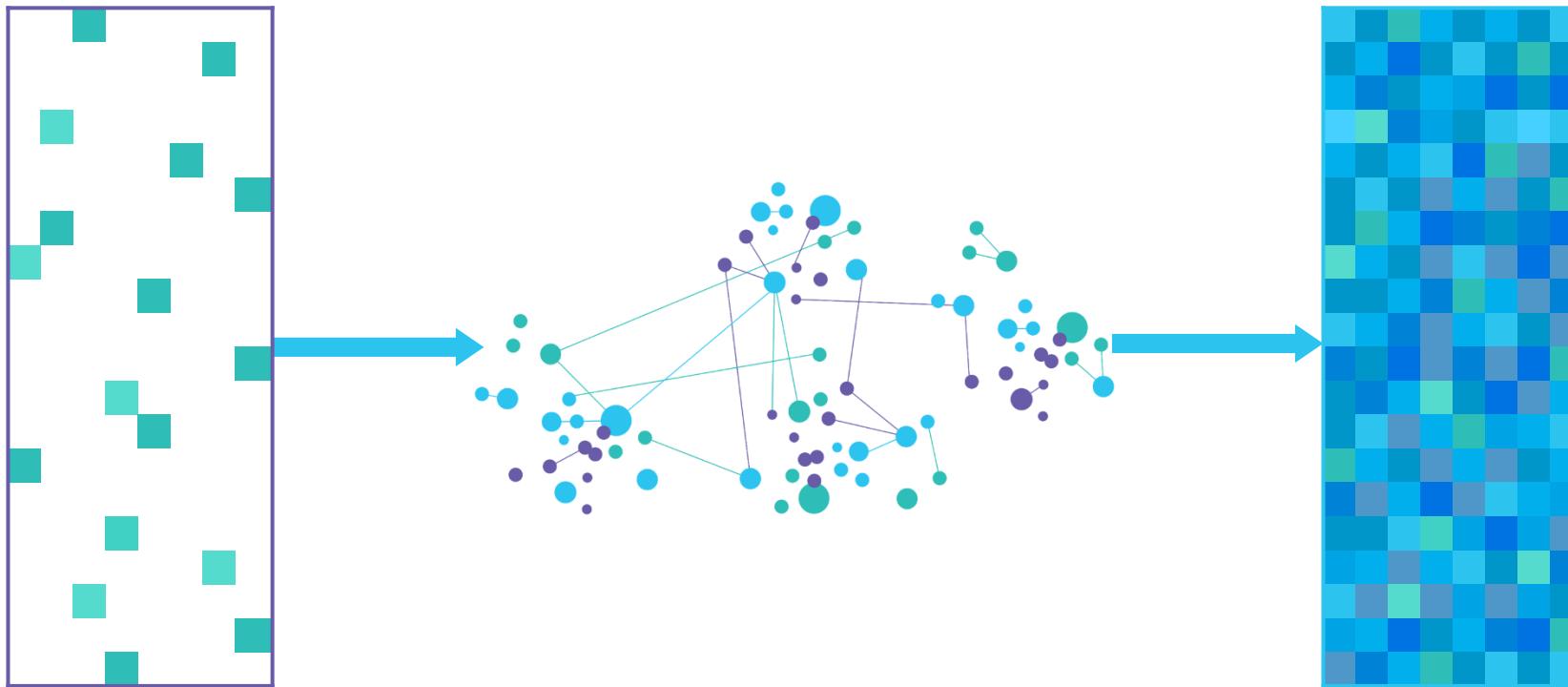
Root Mean Square Error, RMSE

Measure R^2 and RMSE per assay against realistic test set,
then report mean across assays



Aim: impute missing assay values

Validate against
realistically-split holdout set



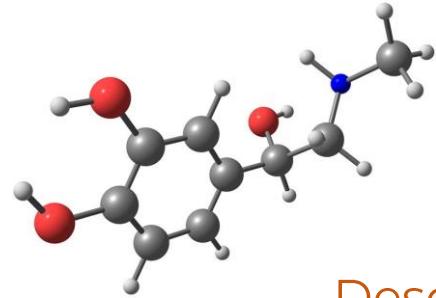
Data from ChEMBL

Martin, Polyakov, Tian, and Perez, J. Chem. Inf. Model. 57, 2077 (2017)

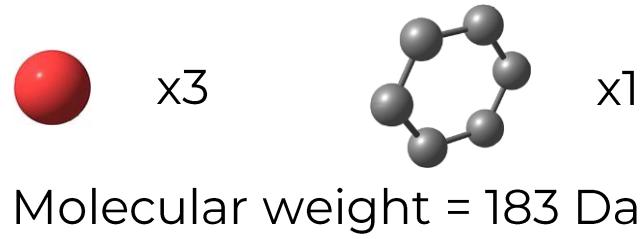
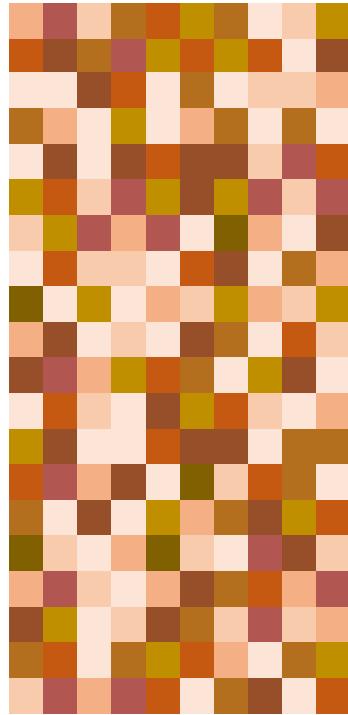
intellegens.ai



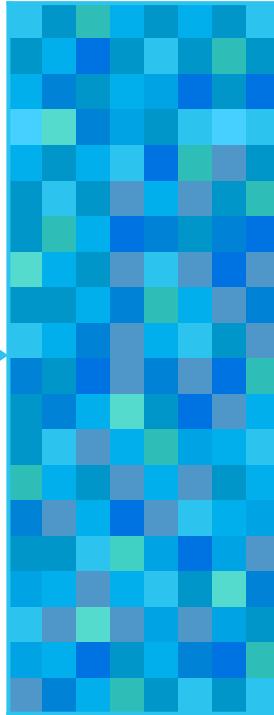
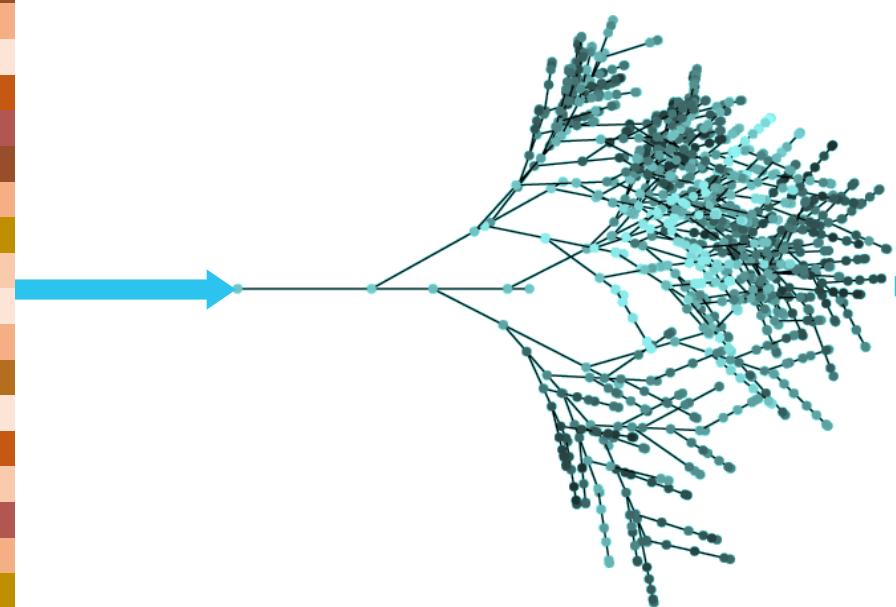
Random forest regression



Descriptors

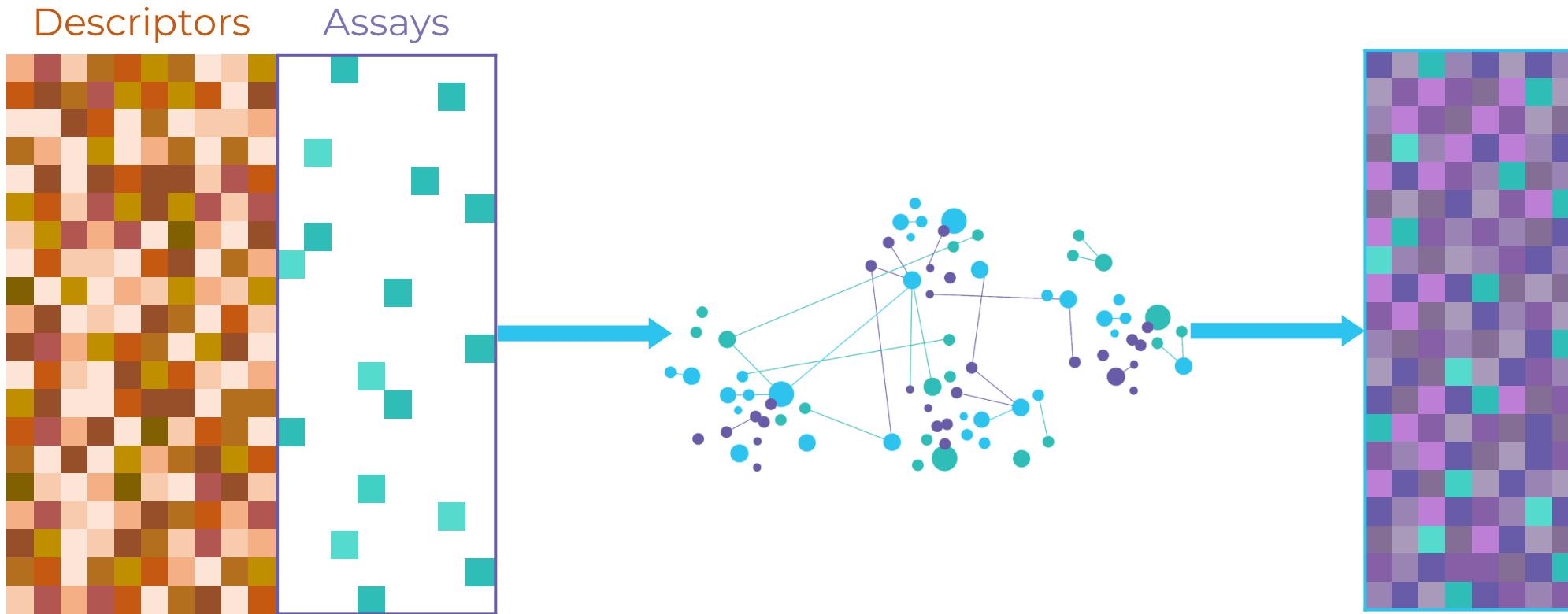


$R^2 = -0.19$
 $RMSE = 0.89$





Descriptors and bioactivity values





Deep learning predictions

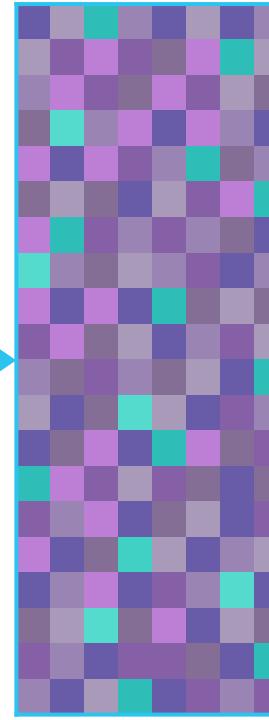
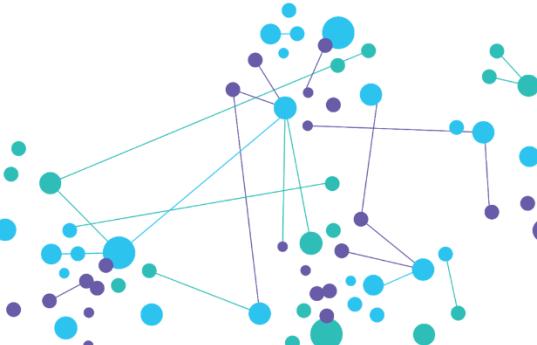
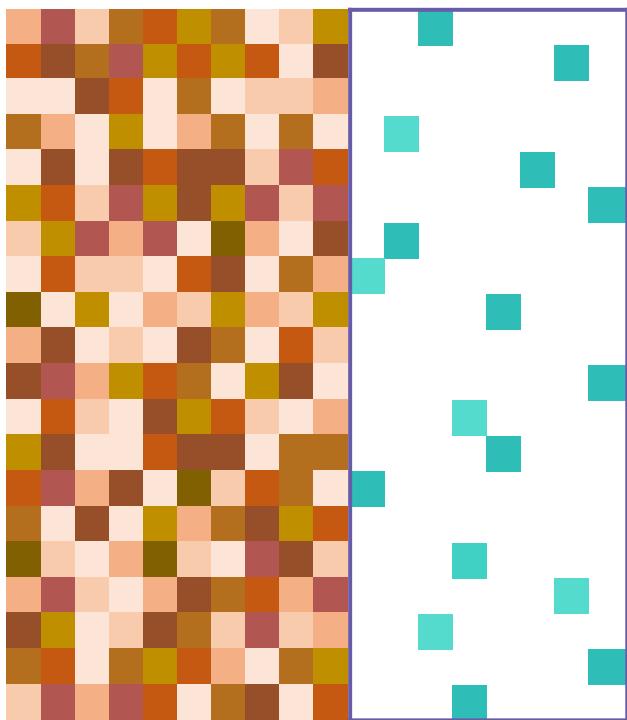
$R^2 = 0.46$

$RMSE = 0.59$

Random forest

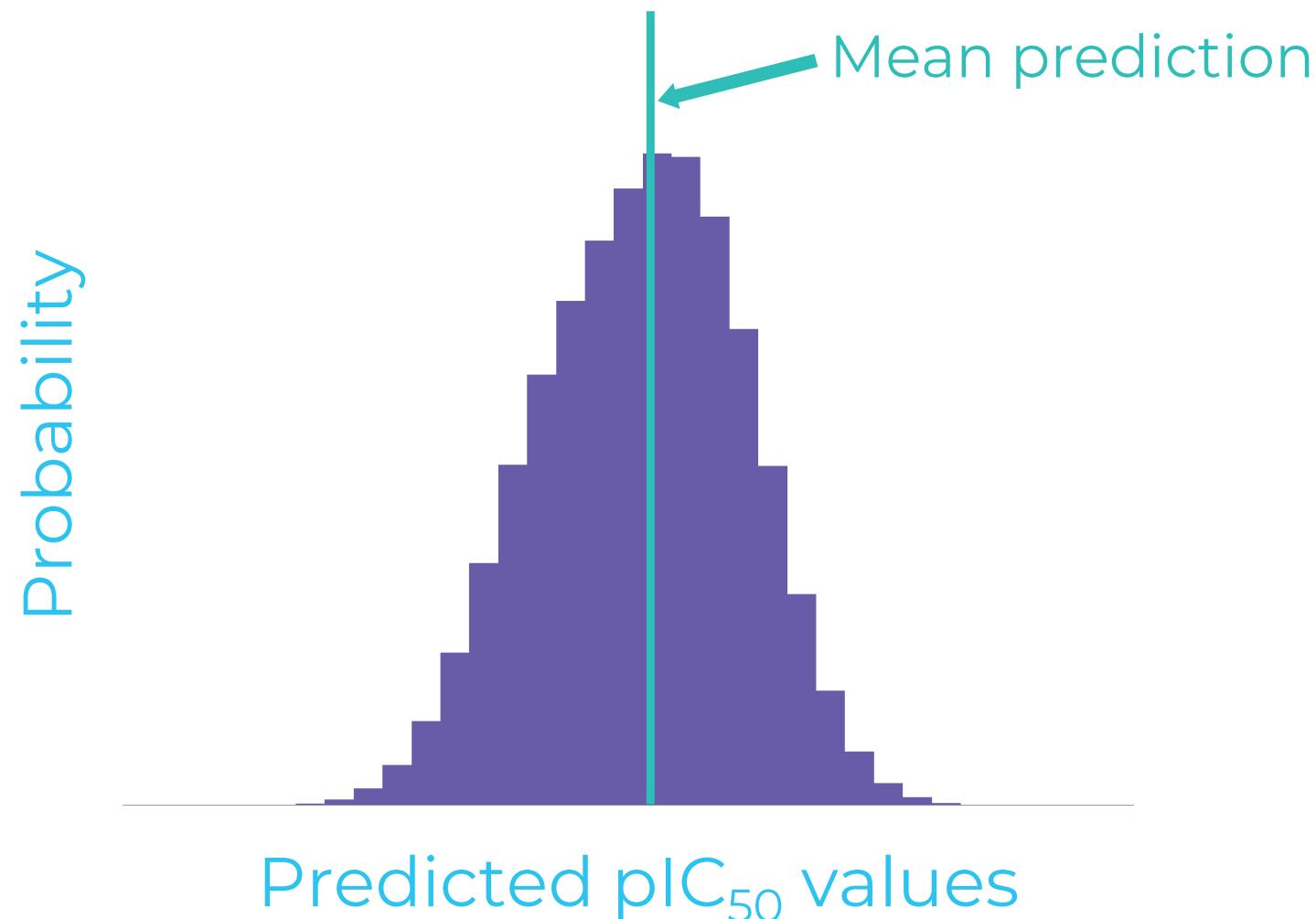
$R^2 = -0.19$

$RMSE = 0.89$



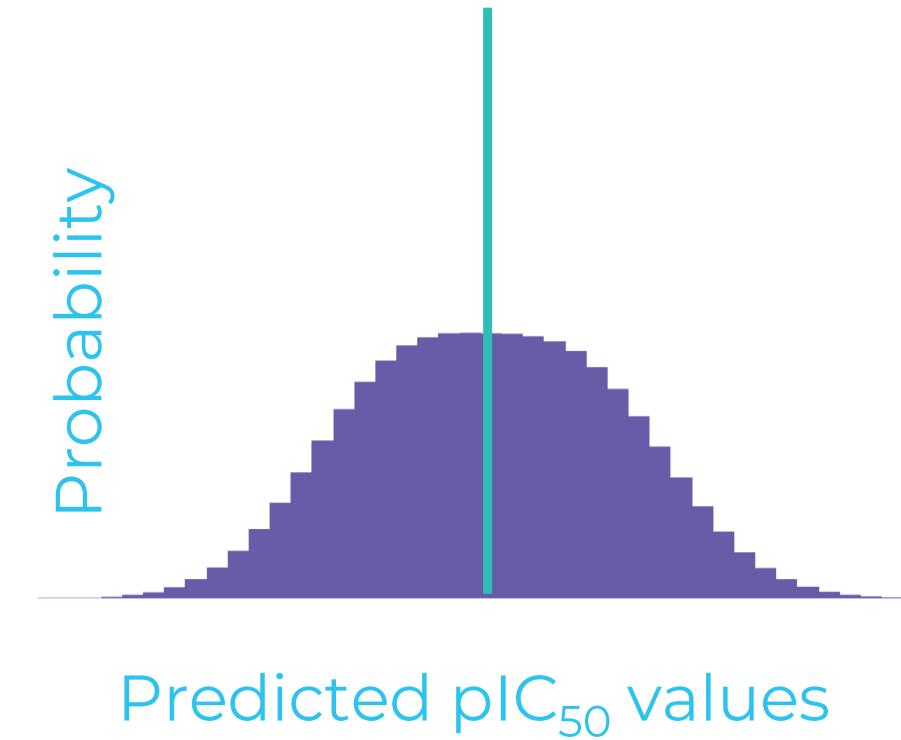
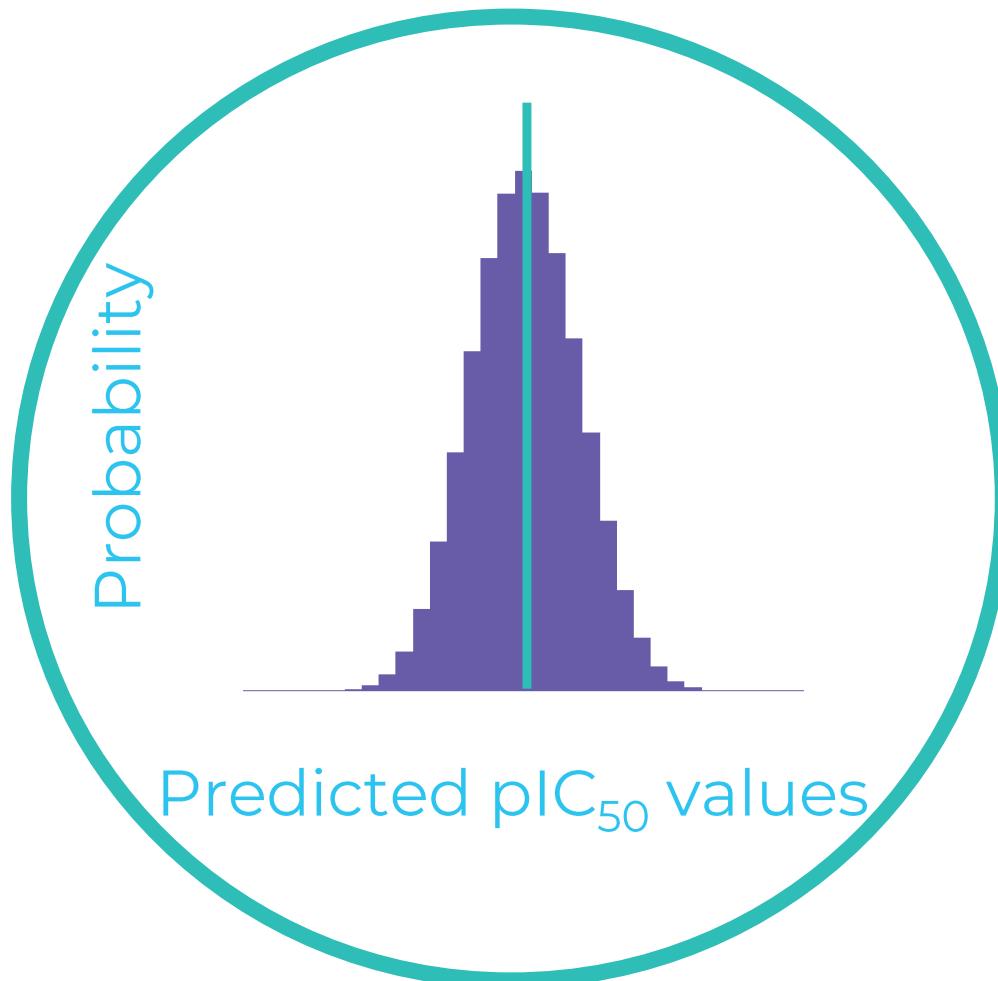


Calculate probability distribution





Focus on most confident predictions



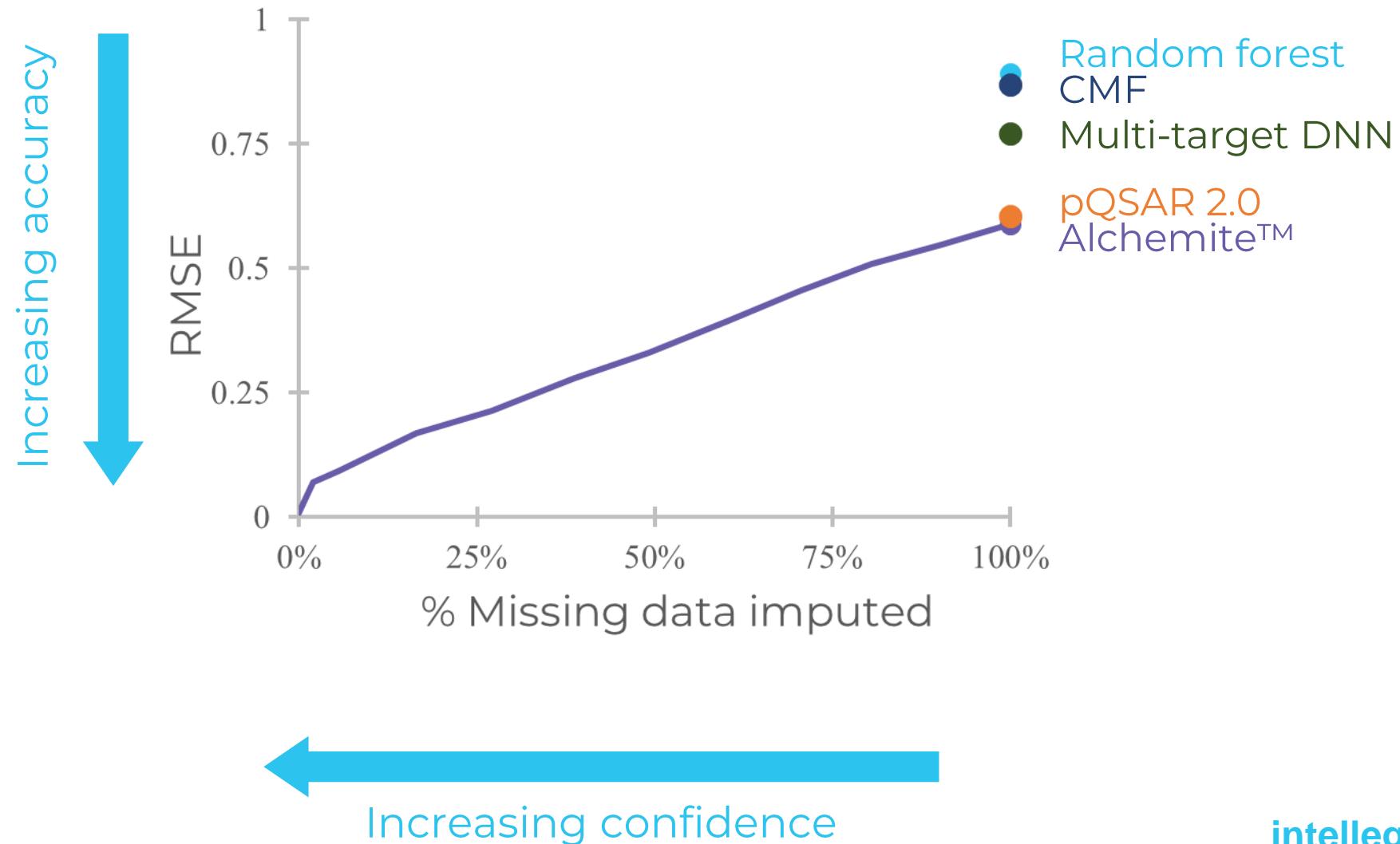


Reporting only most confident predictions





Reporting only most confident predictions





Reporting only most confident predictions





Real project application

Big Pharma company

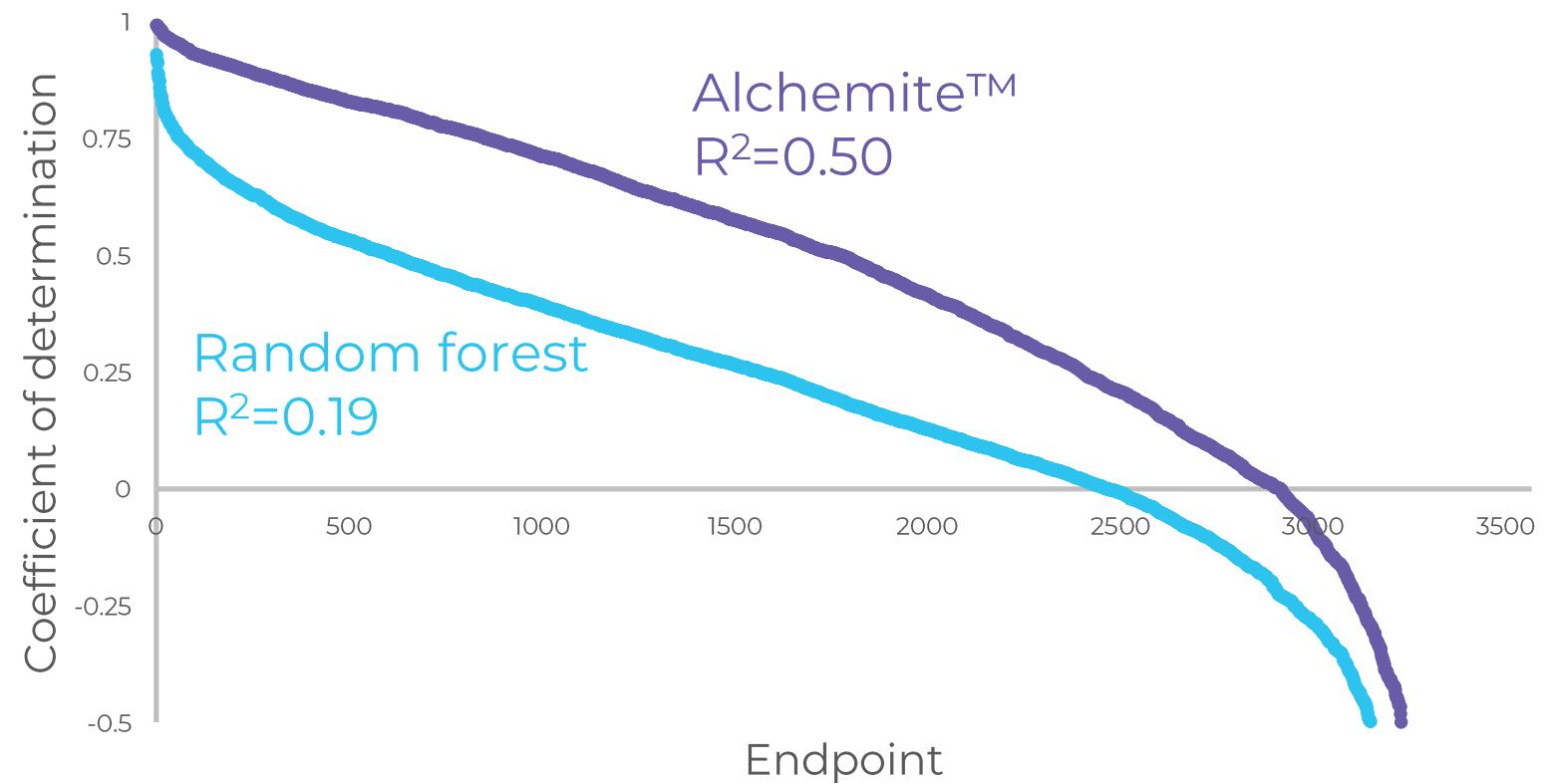
710,305 compounds

2,171 assays totalling **3,568** endpoints

Covering a **full range** of drug discovery assays, including compound activities and ADME properties



Results





Summary

Train across all endpoints simultaneously to capture
activity-activity correlations using sparse data as **input**

Understand and exploit **probability distribution** to focus on most confident results

Applicable to pharma-scale data sets **today** through Alchemite platform

For more details: [Whitehead et al., J. Chem. Inf. Model. 59, 1197 \(2019\)](#)
tom@intellegens.ai

Collaborate using Alchemite

Optibrium and **Intellegens** are working together on applying Alchemite to drug discovery and compound optimisation

If you'd like to discuss how to use Alchemite on your project or data, please contact

info@optibrium.com

