# Imputing Compound Activities Based on Sparse and Noisy Data

**Thomas Whitehead*, Benedict Irwin†, Peter Hunt†, Matthew Segall†, Gareth Conduit***

***Intellegens †Optibrium. Email: matt@optibrium.com**

# The Challenges of Applying Deep Learning to Drug Discovery Data

- Application of conventional deep learning to traditional QSAR modelling offers little advantage

  - Robert Sheridan (Merck) reported an average improvement in $R^2$ of 0.04 over random forests across 30 representative QSAR data sets*

- Challenges

  - Compound bioactivity/property data is very sparse

  - 'Big data' in pharma is not very big! $O(10^6)$ compounds and $O(10^7)$ experimental data points

  - Biological data is noisy. ~0.3-0.5 log unit experimental variability

- How can we learn from these experimental data to make better predictions for compound bioactivities and properties?

*AI in Chemical Research, Switzerland, Sept.9 2018

# 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, …

intellegens.ai

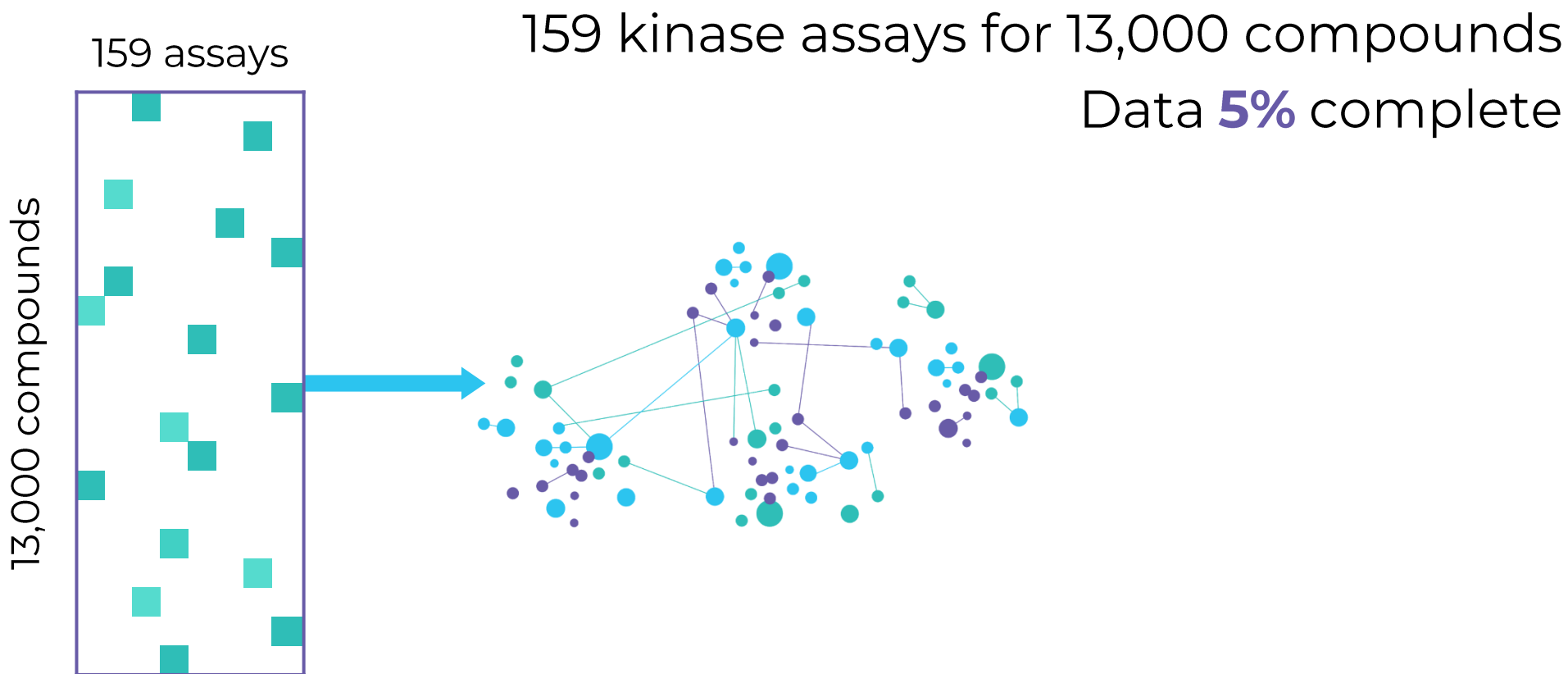# Deep learning



**Inputs** → **Outputs**

intellegens.ai

# Alchemite™ deep learning



Inputs → Outputs

intellegens.ai

# Novartis dataset to benchmark machine learning

159 assays

13,000 compounds

159 kinase assays for 13,000 compounds

Data **5%** complete

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

**intellegens.ai**

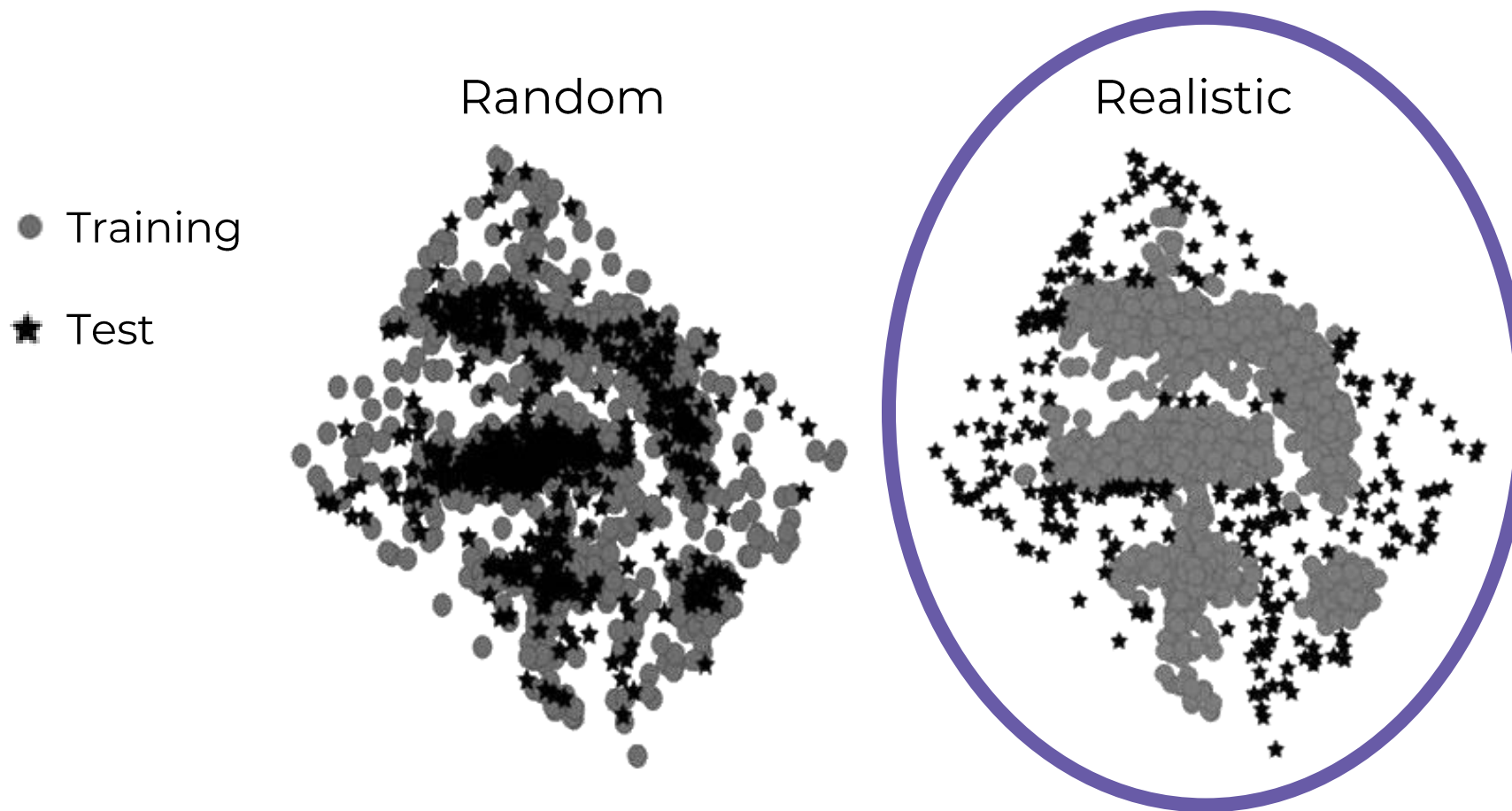# Novartis dataset distribution

Random

- ● Training
- ✶ Test



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

intellegens.ai

# Novartis dataset is realistically distributed



Random

Realistic

- ● Training

- ✸ Test

Data from ChEMBL
Martin, Polyakov, Tian, and Perez, J. Chem. Inf. Model. 57, 2077 (2017)
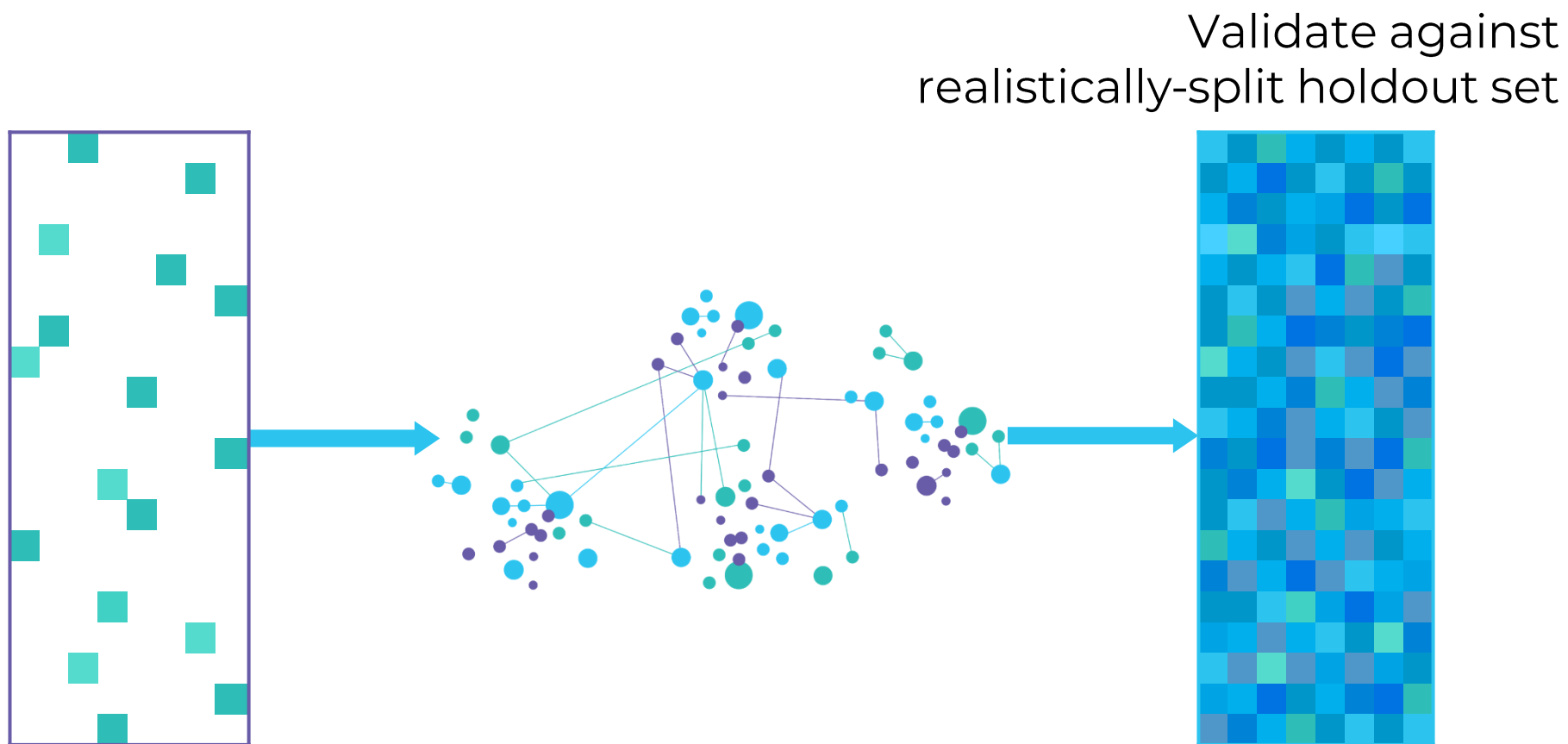
intellegens.ai

# Accuracy metrics

Coefficient of determination, $R^2$

Root Mean Square Error, RMSE

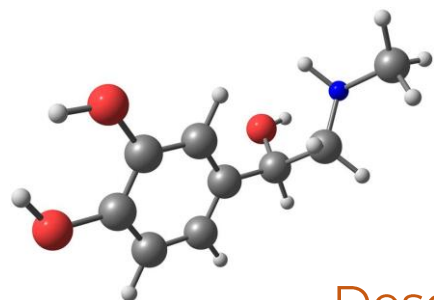Measure per assay against realistic test set, then report mean across assays

**intellegens.ai**

# 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)
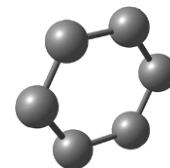
intellegens.ai

# Random forest regression



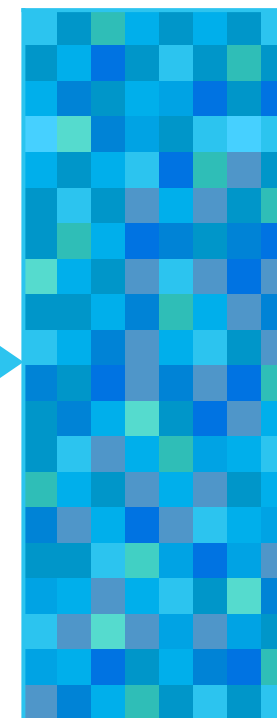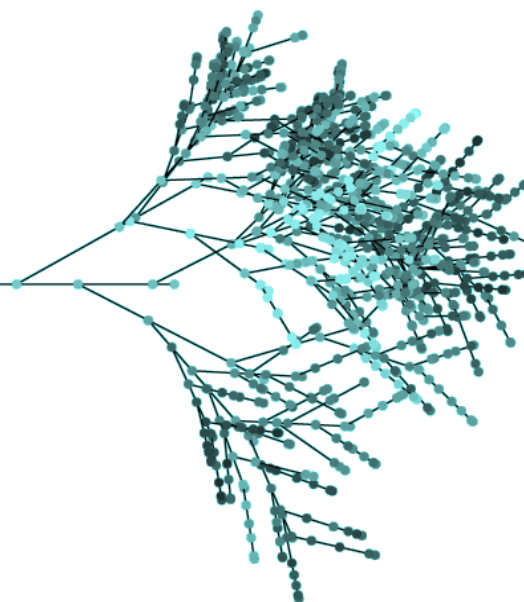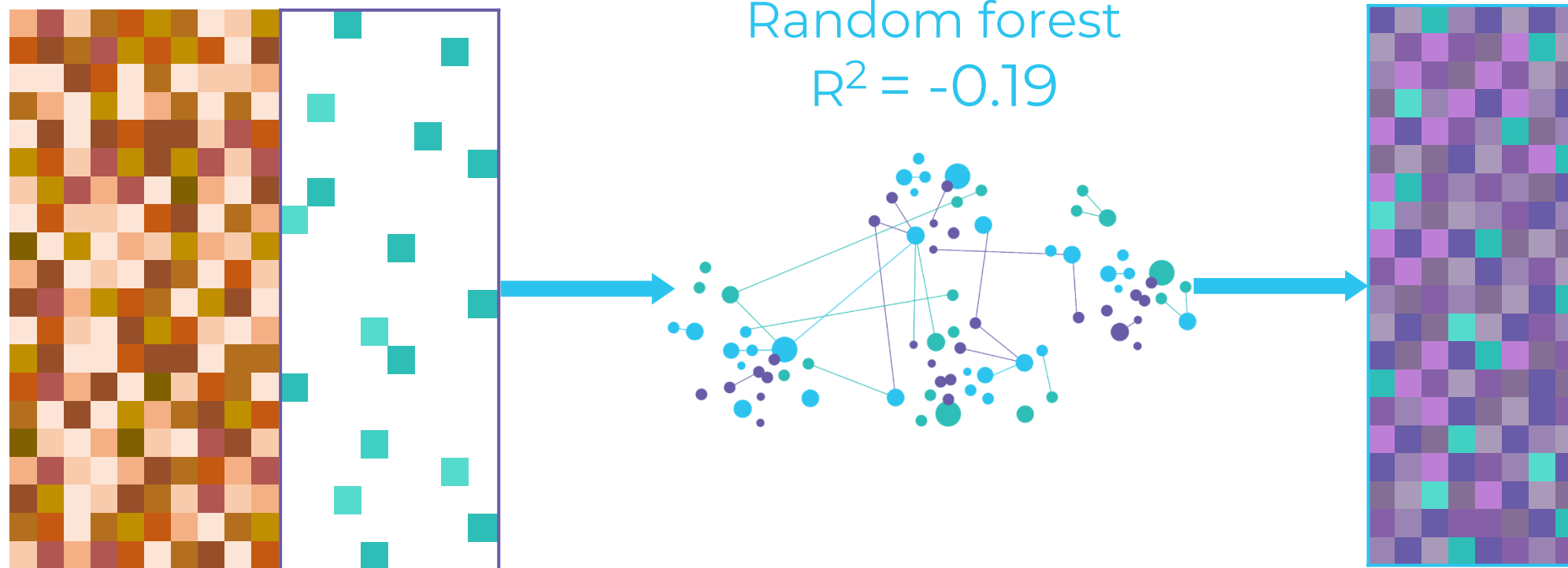x3    x1    x1

Molecular weight = 183 Da

Descriptors

$R^2 = -0.19$

intellegens.ai

11

# Descriptors and bioactivity values



Descriptors          Assays

Whitehead *et al.* J. Chem. Inf. Model (2019) **59**(3), pp 1197–1204

intellegens.ai

# Deep learning predictions



$$R^2 = 0.46$$

Random forest
$$R^2 = -0.19$$

Whitehead *et al.* J. Chem. Inf. Model (2019) **59**(3), pp 1197–1204

# Comparison with other methods

| Method | $R^2$ | RMSE |
|---|---|---|
| **Alchemite** | **0.46*** | **0.59** |
| Profile QSAR 2.0 | 0.43 | 0.61 |
| Multi-target deep neural network (tensor-flow) | 0.11 | 0.77 |
| Collective matrix factorisation | -0.11 | 0.87 |
| Random forest | -0.19 | 0.89 |

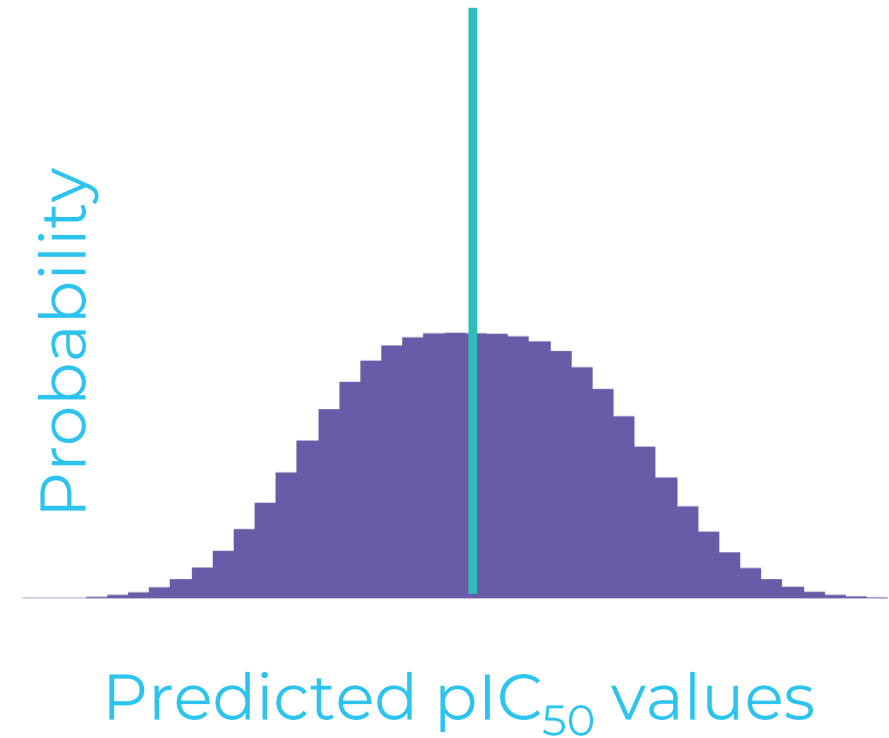* N.B. Improved over published results ($R^2$=0.44). Also 1000× faster to build model!

Whitehead *et al*. J. Chem. Inf. Model (2019) **59**(3), pp 1197–1204

**intellegens.ai**

# Calculate probability distribution



Mean prediction

Probability

Predicted $pIC_{50}$ values

Whitehead *et al.* J. Chem. Inf. Model (2019) **59**(3), pp 1197–1204

# Focus on most confident predictions



intellegens.ai

# Reporting only most confident predictions



RMSE (y-axis) vs % Missing data imputed (x-axis, 0% to 100%)

Legend:
- Random forest
- CMF
- Multi-target DNN
- pQSAR 2.0
- Alchemite™

← Increasing confidence

intellegens.ai

# Reporting only most confident predictions



Increasing accuracy

RMSE

% Missing data imputed

Increasing confidence

Random forest
CMF
Multi-target DNN

pQSAR 2.0
Alchemite™

intellegens.ai

# Reporting only most confident predictions



- Random forest
- CMF
- Multi-target DNN

- pQSAR 2.0
- Alchemite™

Increasing accuracy

RMSE

% Missing data imputed

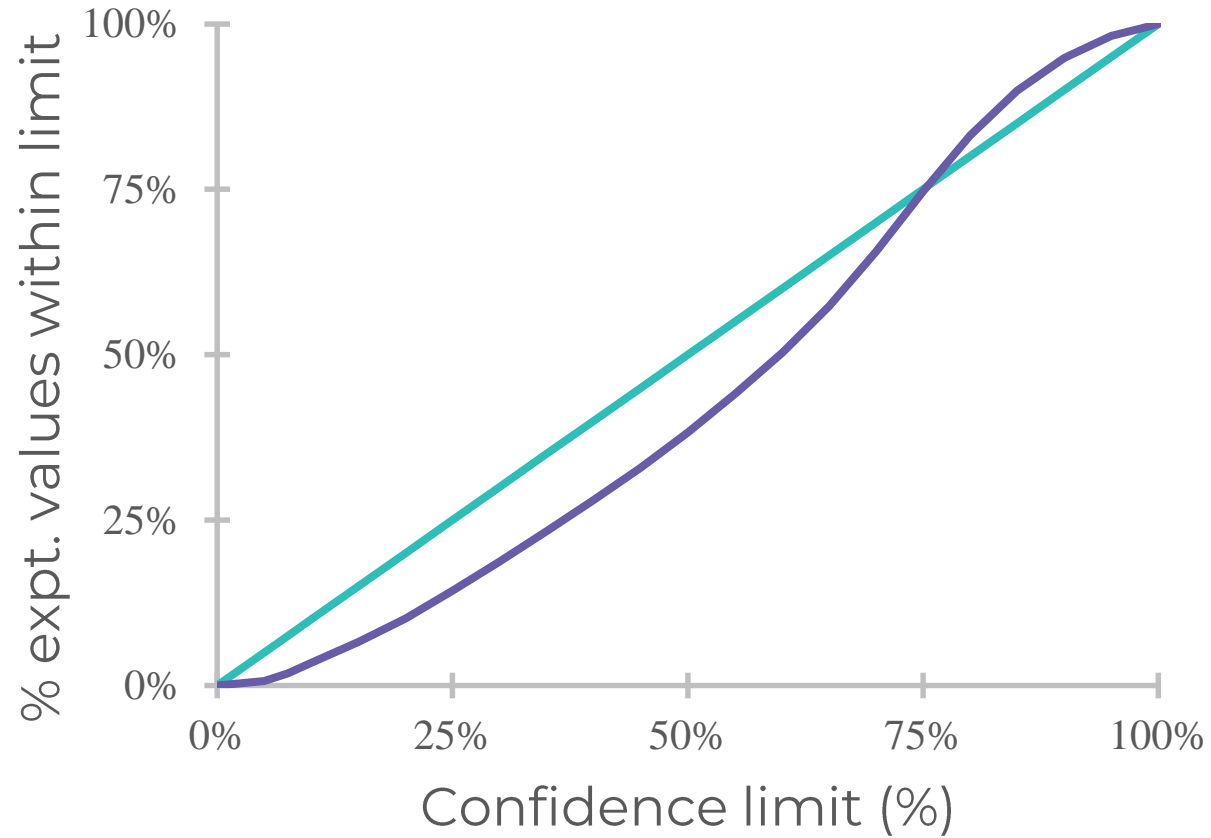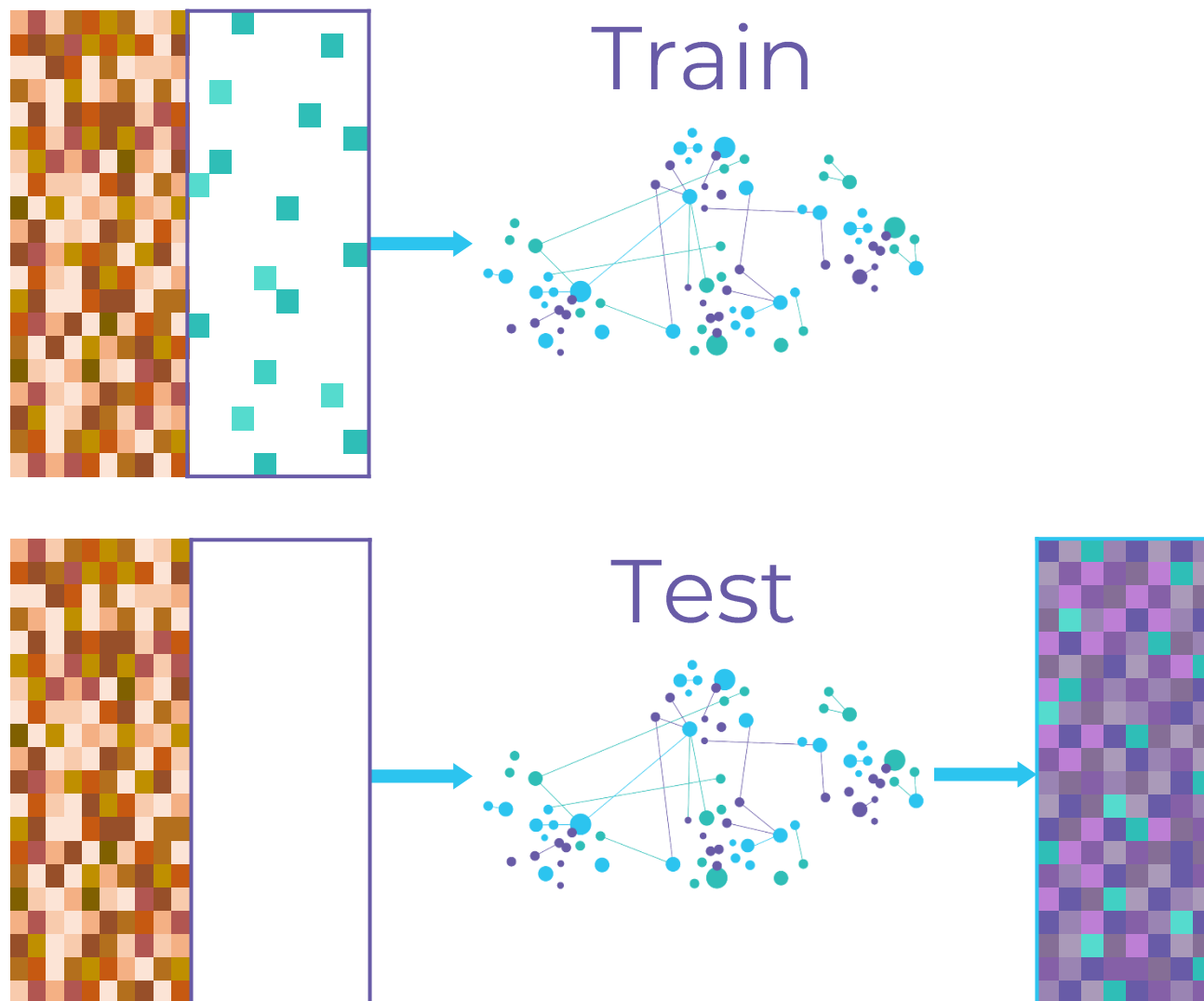Increasing confidence

intellegens.ai

# Absolute accuracy of uncertainties



N.B. Assumes normally distributed
errors e.g. 62% of results within 1 SD

intellegens.ai

# Application to virtual compounds



Train

Test

intellegens.ai

# Application to virtual compounds



Increasing accuracy

Virtual compound Alchemite™

RMSE

% Missing data imputed

Random forest
CMF
Multi-target DNN

pQSAR 2.0
Alchemite™

Increasing confidence

intellegens.ai

# Random forest confidence predictions



intellegens.ai

# Conclusions

- 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

- Impute results of missing assays to **high accuracy**

- **Broadly applicable** to other endpoints, e.g. physicochemical, ADME, tox…

- **Applicable** to pharma-scale data sets

- For more details: **Whitehead *et al.* J. Chem. Inf. Model (2019) 59(3), pp 1197–1204**
  - matt@optibrium.com