



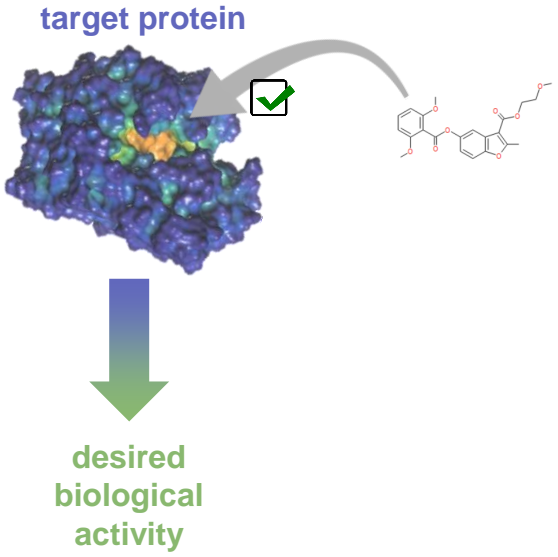
Building Chemogenomics Models From a Large-Scale Public Dataset and Applying them to Industrial Datasets



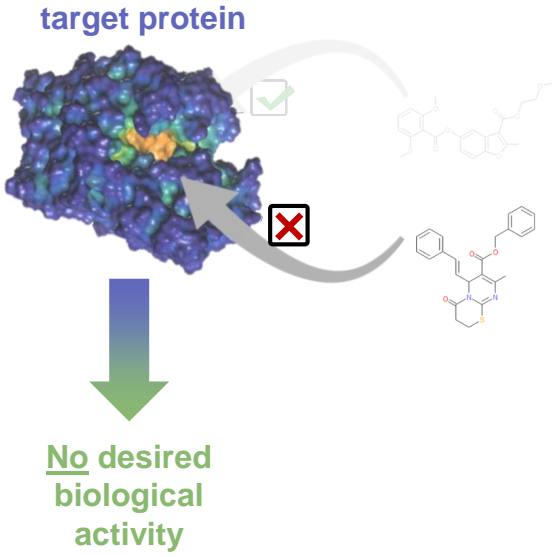
*Funded by the European Union's Horizon 2020 Research and Innovation program.
Grant Agreement no. 654168*



Chemogenomics modelling

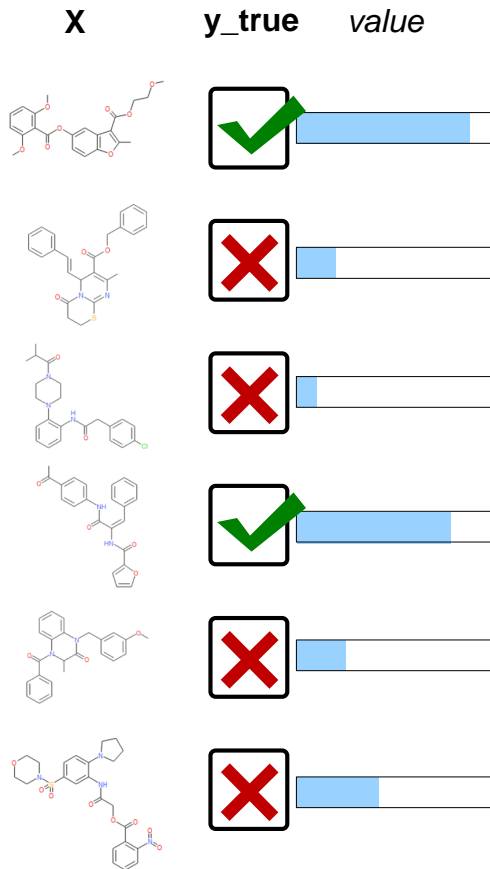
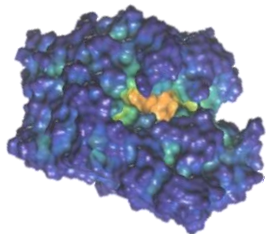


Chemogenomics modelling



Chemogenomics modelling

target protein



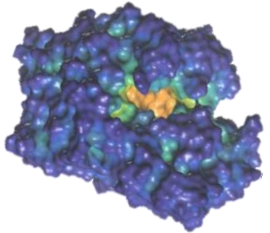
{ Molecule ; Activity }

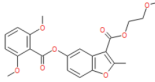
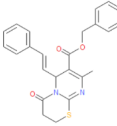
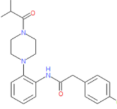
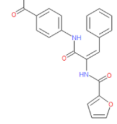
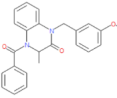
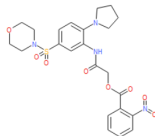
X \Rightarrow *y_true*

Classification: *y_true* is binary
Regression: *y_true* is continuous

Chemogenomics modelling

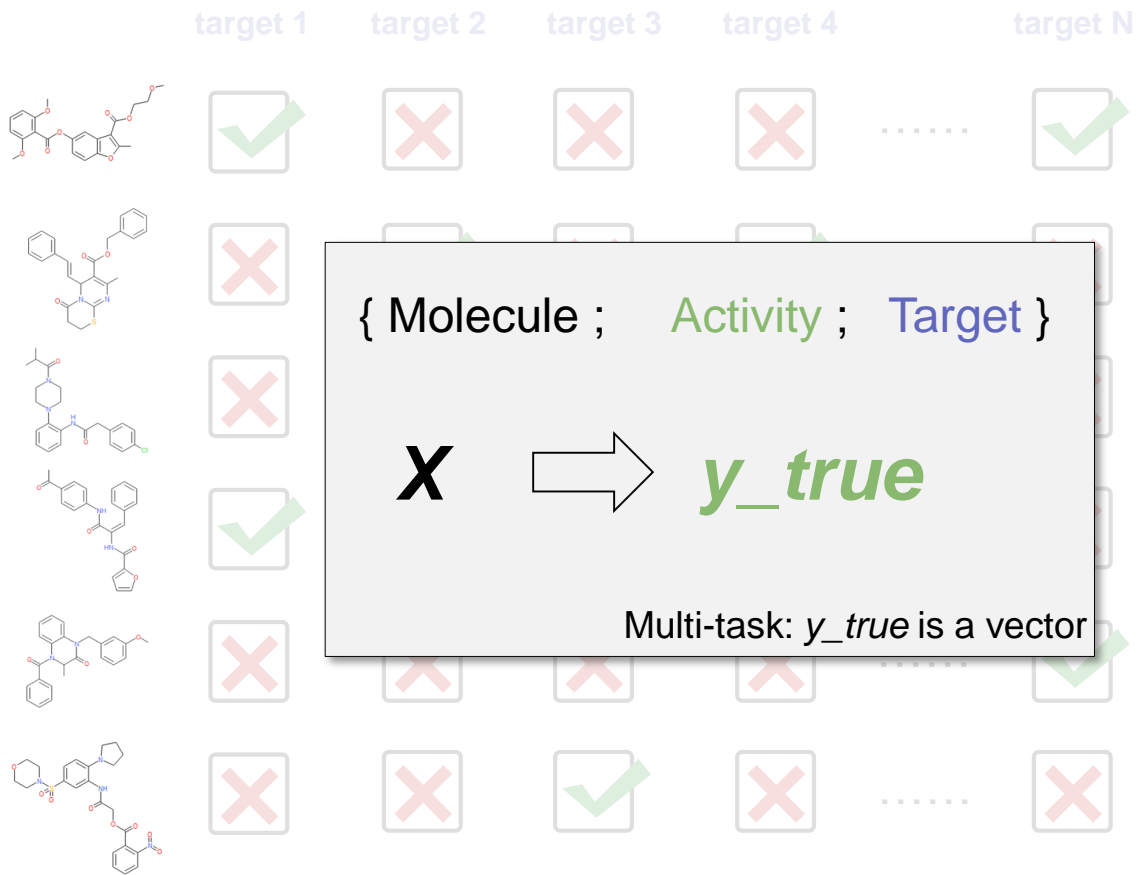
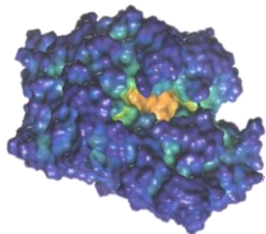
target protein



	target 1	target 2	target 3	target 4	target N
	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>
	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>
	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

Chemogenomics modelling

target protein



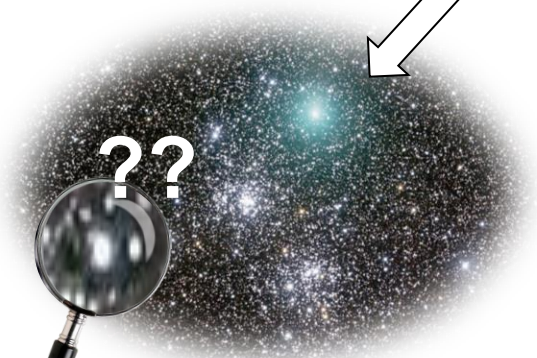
Aim of chemogenomics: Accelerate Drug Discovery

Drug Discovery Pipeline



Aim of chemogenomics: Accelerate Drug Discovery

Drug Discovery Pipeline

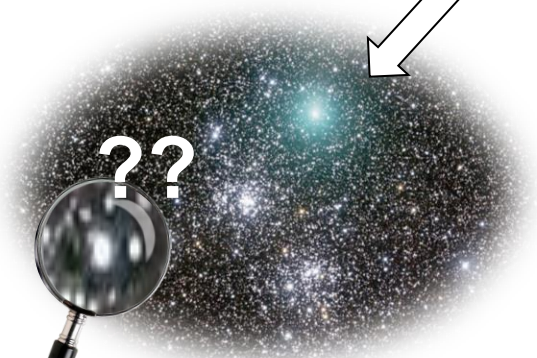


10^{60} molecules

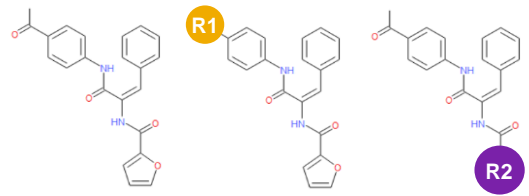
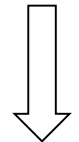
→ Predict activity of compounds from uncharted chemical space

Aim of chemogenomics: Accelerate Drug Discovery

Drug Discovery Pipeline



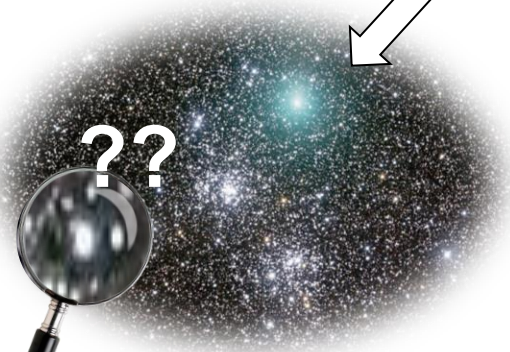
10^{60} molecules



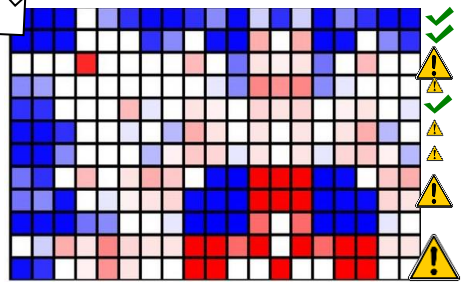
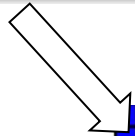
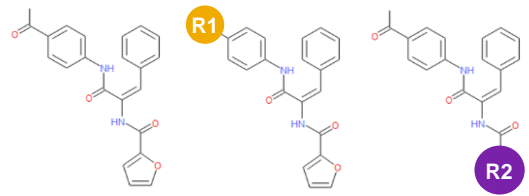
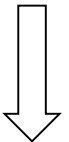
- Predict activity of compounds from uncharted chemical space
- Be sensitive enough to handle subtle chemical changes

Aim of chemogenomics: Accelerate Drug Discovery

Drug Discovery Pipeline



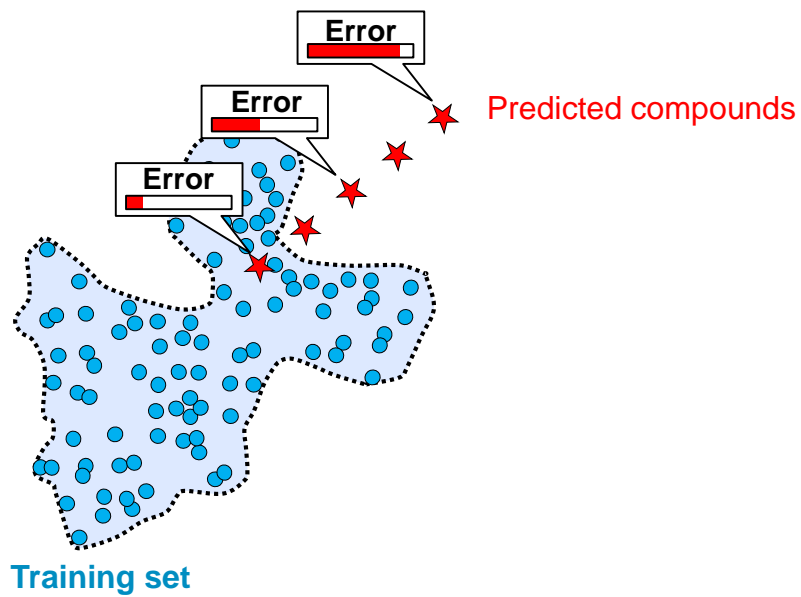
10^{60} molecules



Safety profiling panels

- Predict activity of compounds from uncharted chemical space
- Be sensitive enough to handle subtle chemical changes
- Be accurate across targets (MOA elucidation, safety)

Challenge: Application Domain



Challenge: Application Domain

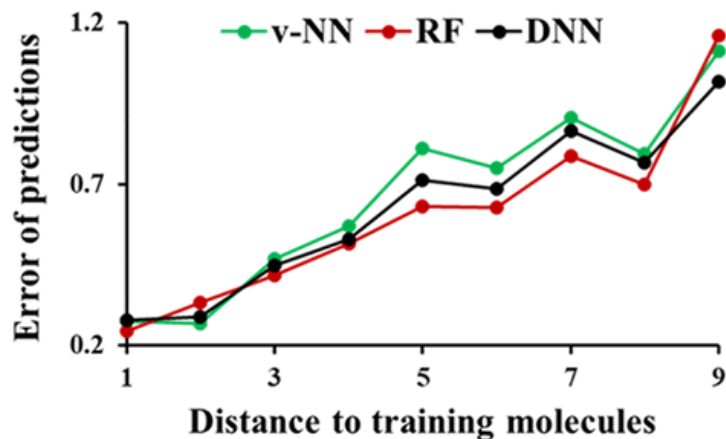
Dissecting Machine-Learning Prediction of Molecular Activity: Is an Applicability Domain Needed for Quantitative Structure–Activity Relationship Models Based on Deep Neural Networks?

Ruifeng Liu¹, Hao Wang¹, Kyle P. Glover⁵, Michael G. Feasel⁵, and Anders Wallqvist¹

¹ Department of Defense, Biotechnology High Performance Computing Software Applications Institute, Telemedicine and Advanced Technology Research Center, U.S. Army Medical Research and Materiel Command, Fort Detrick, Maryland 21702, United States

⁵ U.S. Army–Edgewood Chemical Biological Center, Aberdeen Proving Ground, Maryland 21010, United States

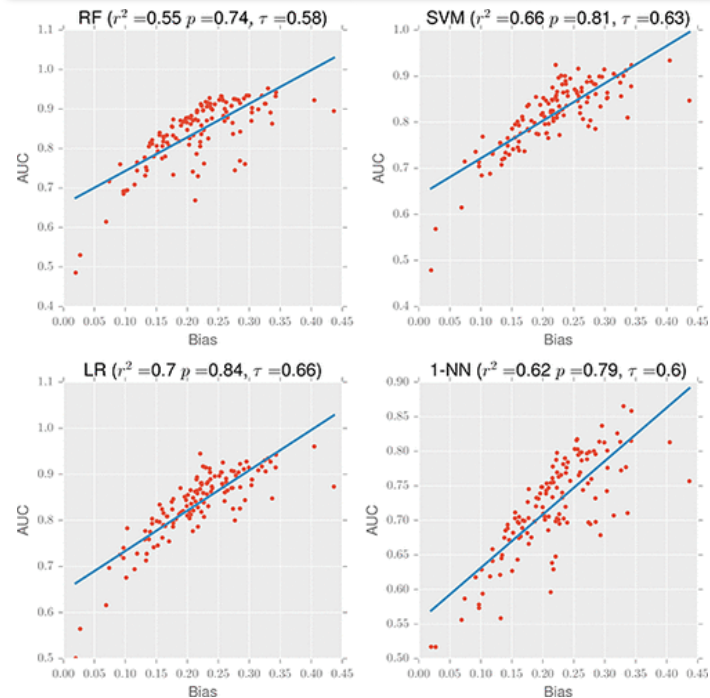
Distance to training molecules is the most important determinant of prediction accuracy



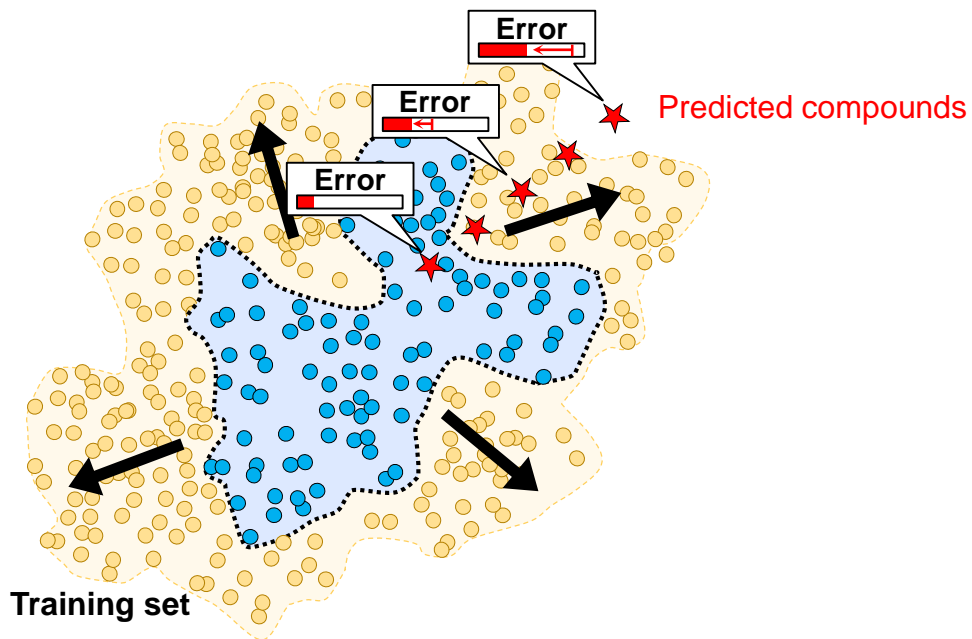
Most Ligand-Based Classification Benchmarks Reward Memorization Rather than Generalization

Izhar Wallach^{*} and Abraham Heifets^{*}

Atomwise Inc., 221 Main Street, Suite 1350, San Francisco, California 94105, United States



Challenge: Application Domain

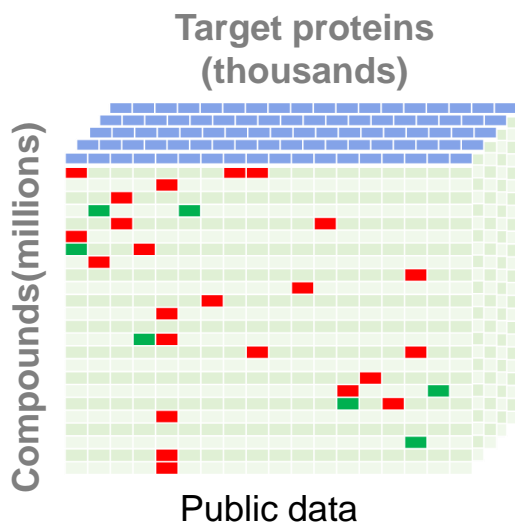


One Possible solution:

- Expand training set
- Integrate many examples of (active, inactive)
- public ML benchmark dataset



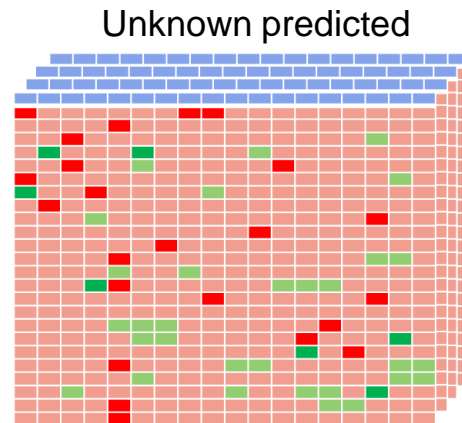
Exa-Scale Compound Activity Prediction Engine



Machine learning
algorithms



HPC



**But also,
predict new compounds!**




Chemogenomics Data



AstraZeneca
Janssen
IOFA consult

High Performance Computing



intel
IT4Innovations national supercomputing center
imec

Machine Learning Algorithms



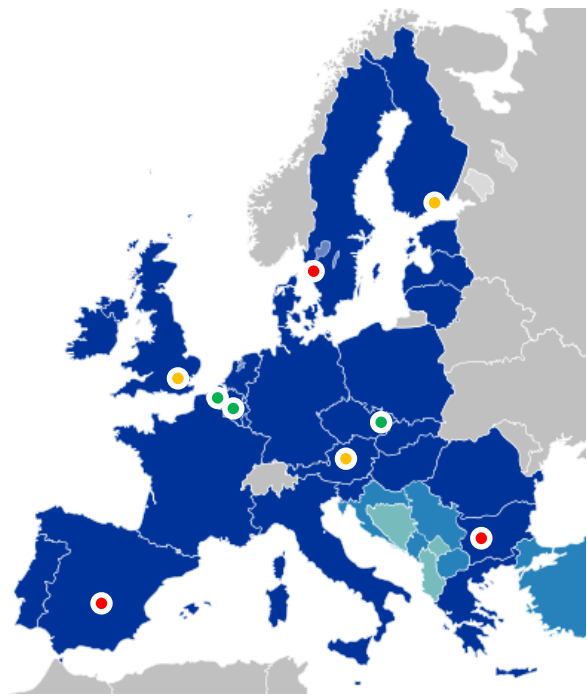
JKU
JOHANNES KEPLER UNIVERSITÄT LINZ
ROYAL HOLLOWAY UNIVERSITY OF LONDON
A? Aalto-yliopisto

9 partners / 8 countries

- Compound activity
- Data preparation
- Benchmarking

- Cluster deployment
- Programming model

- Deep learning
- Matrix factorisation
- Conformal prediction





Overall Workflow

1. Dataset collection
2. Dataset split
3. Descriptors
4. Performance metrics
5. Algorithms
6. Hyperparameter selection
7. Retrospective model evaluation
8. Prospective model evaluation



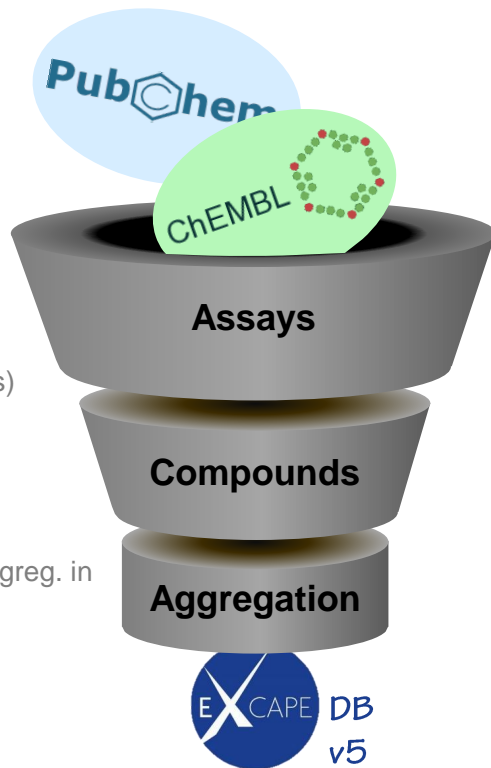
Dataset choice: ChEMBL + PubChem = *ExCAPE-DB*

AstraZeneca
Jiangming Sun
Hongming Chen

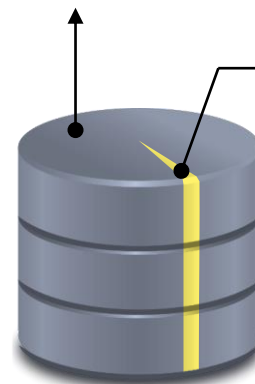
Janssen
Felip Golib
Vladimir Chupakhin

IOFA
consult
Nina Jeliaskova

- ✓ Single protein assays
- ✓ Human, mouse, rat
- ✓ Activity $\leq 10 \mu\text{M}$
- ✓ Inactive (flags from hts)
- ✓ Structure standization
- ✓ MW < 1000
- ✓ Heavy atoms > 12
- ✓ Compound activity aggreg. in pXC50



INACTIVES
Data pts 69,517,737
(719,192 cmpds)



ACTIVES
Data pts 1,332,426
(593,156 cmpds)

Table 1. Table 1: Summary of the dataset		
Category	Count	Percentage
Total	70,850,153	100%
Inactives	69,517,737	98.1%
Actives	1,332,426	1.9%
Unique Compounds	998,131	1.4%
Target Proteins	1667	0.002%

TOTAL
70,850,153 data points
998,131 unique compounds
1667 target proteins





Chemogenomics data: *ExCAPE-DB*

<https://solr.ideaconsult.net/search/excape/>

DATABASE

Open Access



ExCAPE-DB: an integrated large scale dataset facilitating Big Data analysis in chemogenomics

Jiangming Sun¹, Nina Jeliakova², Vladimir Chupakhin³, Jose-Felipe Golib-Dzib⁴, Ola Engkvist¹, Lars Carlsson¹, Jorg Wegner³, Hugo Ceulemans³, Ivan Georgiev², Vedrin Jeliakov², Nikolay Kochev^{2,5}, Thomas J. Ashby⁶ and Hongming Chen^{1*}



Home Help ▾

ExCAPE-DB: ExCAPE chemogenomics database

Free-text Similarity Substructure

▶ Data sources (998131)

▶ Species (998131)

▶ Orthologous group (997555)

▶ Gene symbol (998131)

▶ Entrez ID (997425)

▶ Results (998131)

Enter free text phrase...



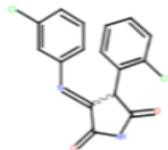
Hits list

Selection

Download

No filters selected!

< 1 2 3 ... 99813 99814 > displaying 11 to 20 of 998131



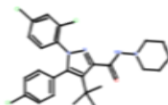
chembl20 **CHEMBL156506** (AGUJQIBUGCOXKU-UYBDAZJANA-N)

TOX.chembl20 N pXC50 = 3.71 [GSK3A] [CHEMBL:71130]

[more](#)

[Chemical structure](#)

[Add to Selection](#)



chembl20 **CHEMBL511907** (AGUKQCZXHXAYCS-SREBMQDQNA-N)

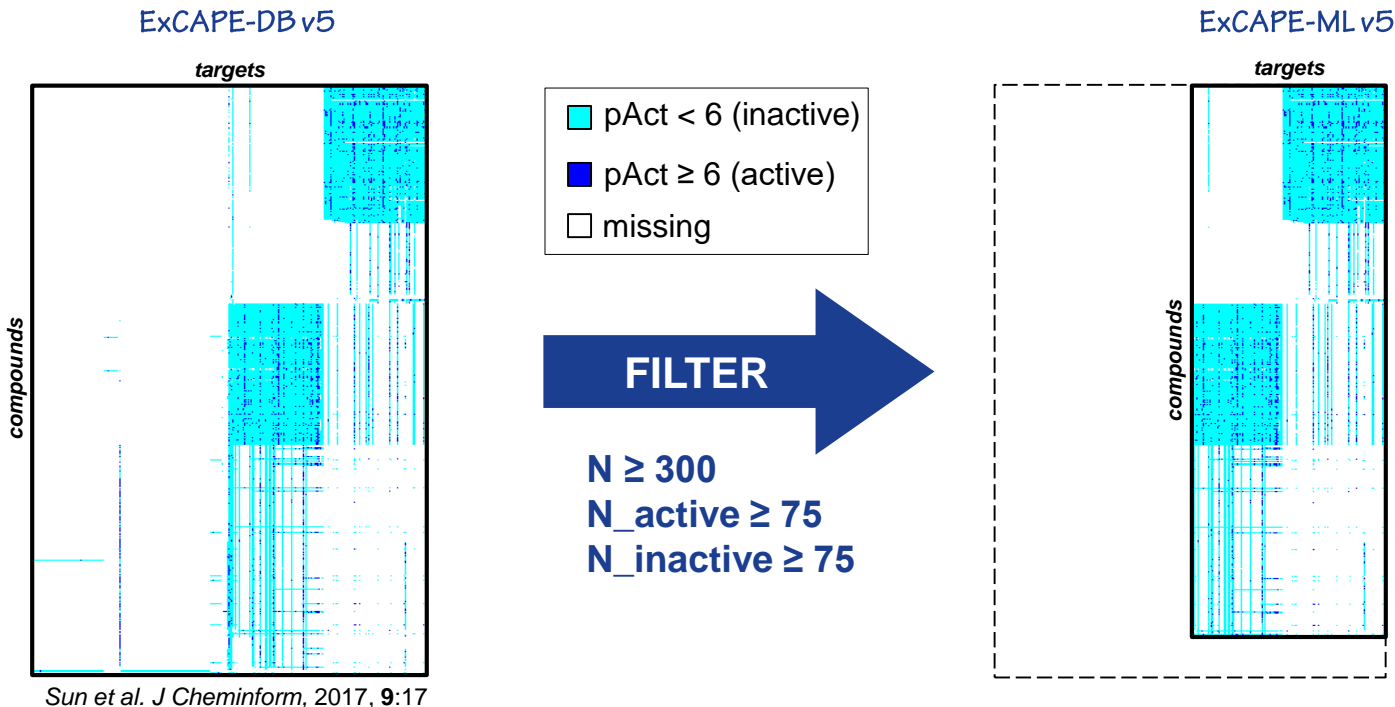
TOX.chembl20 A pXC50 = 8.21 [CNR1] [CHEMBL:500314]

[Chemical structure](#)

[Add to Selection](#)



Machine learning dataset: *ExCAPE-ML*



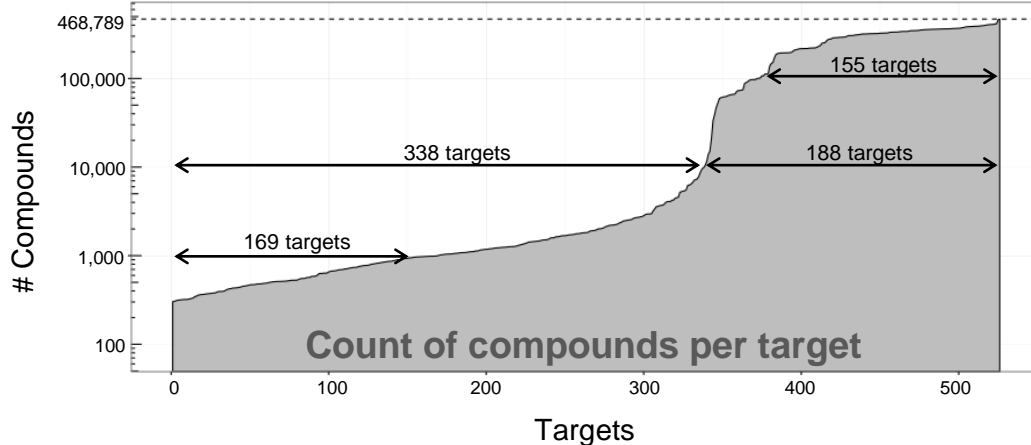
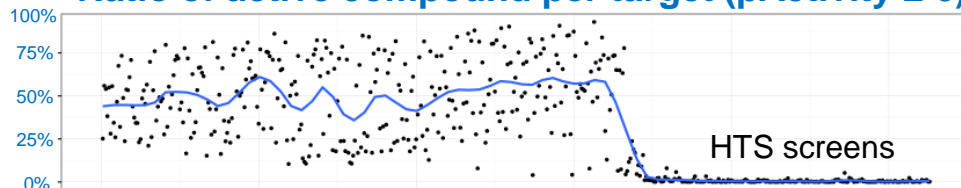
70,850,153 - data points (A:I = 69.5M: 1.3M)
998,131 - unique compounds
1667 - target proteins
4.3% - density

49,516,318 - data points (A:I = 48.8M: 0.5M)
955,386 - unique compounds
526 - target proteins
9.8% - density



Target dataset sizes and active ratio

Ratio of active compound per target (pActivity ≥ 6)

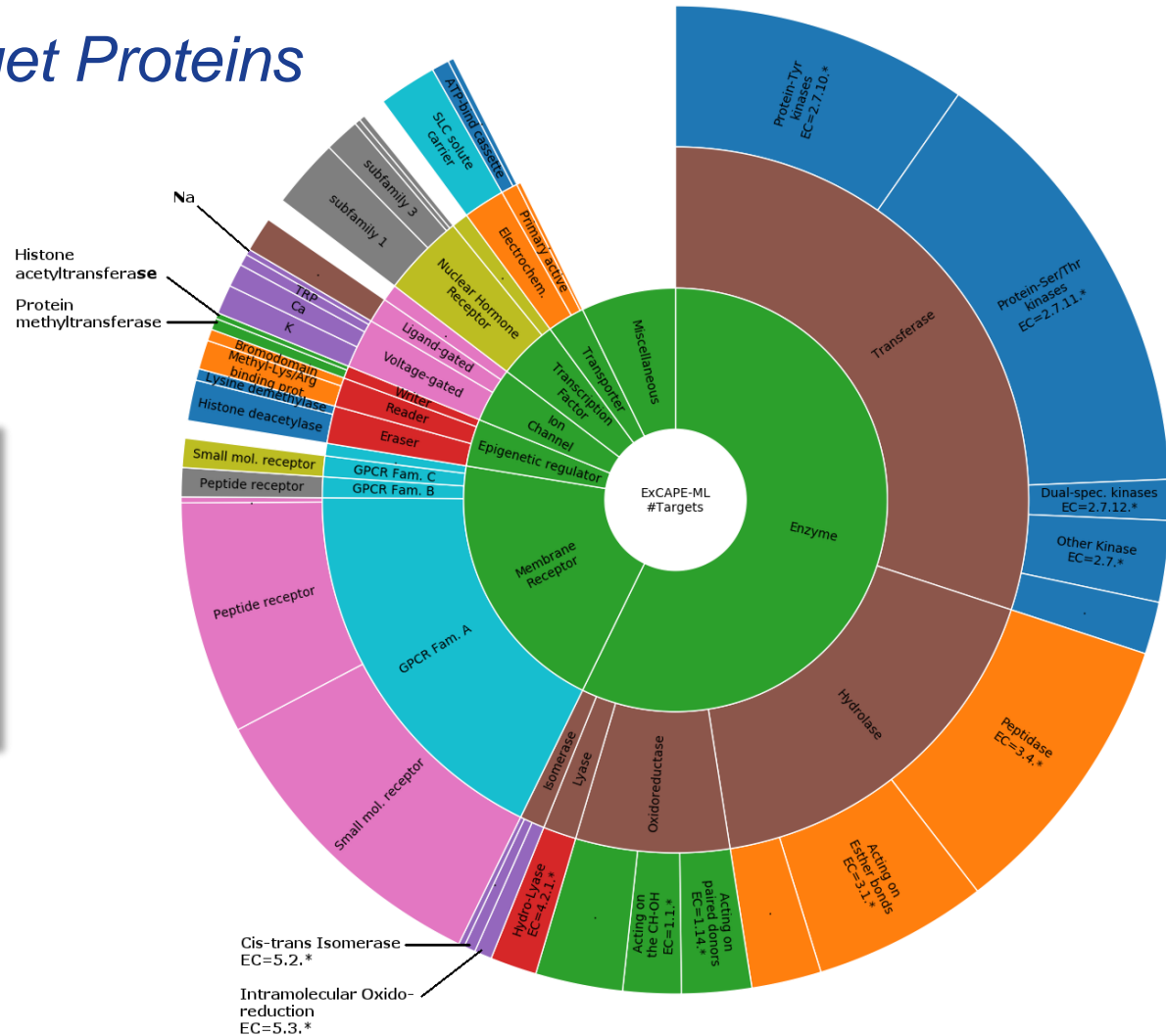


- 188 targets tested $\geq 10k$
- 155 targets tested $\geq 100k$
- 321 targets with $\geq 15\%$ actives have ~ 1500 data points on average.
- 205 targets with $< 15\%$ actives have $\sim 238,000$ data points on average.

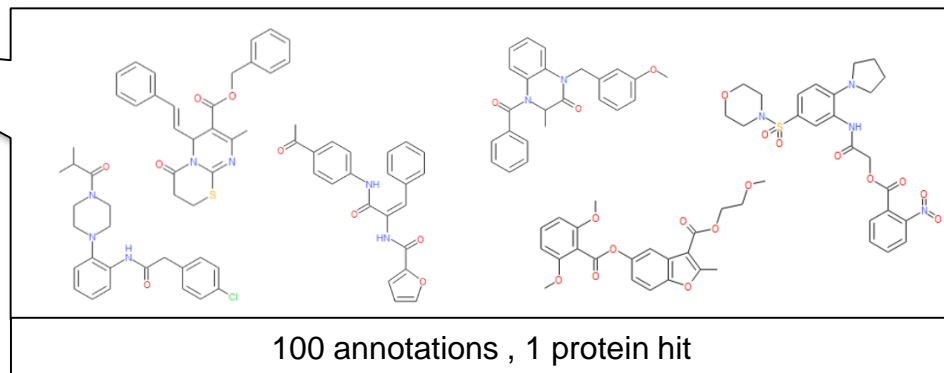
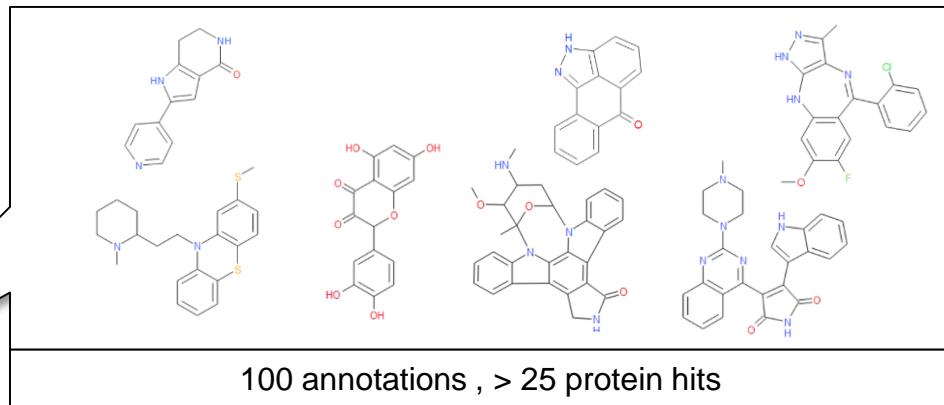
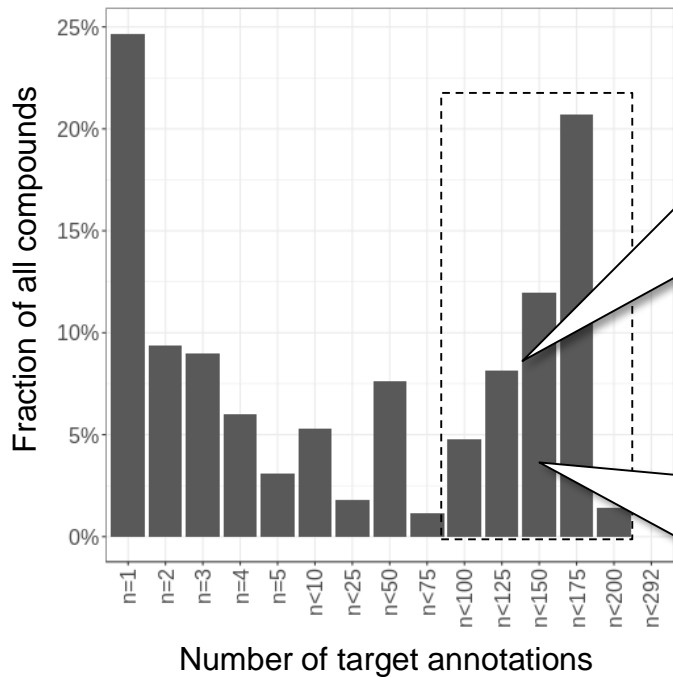
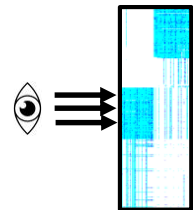


ExCAPE-ML: Target Proteins

- 301 Enzymes
- 107 Membrane receptors
- 24 Transcription factors
- 22 Ion channels
- 19 Epigenetic regulators
- 15 Transporters



ExCAPE Compounds annotations in ExCAPE-ML



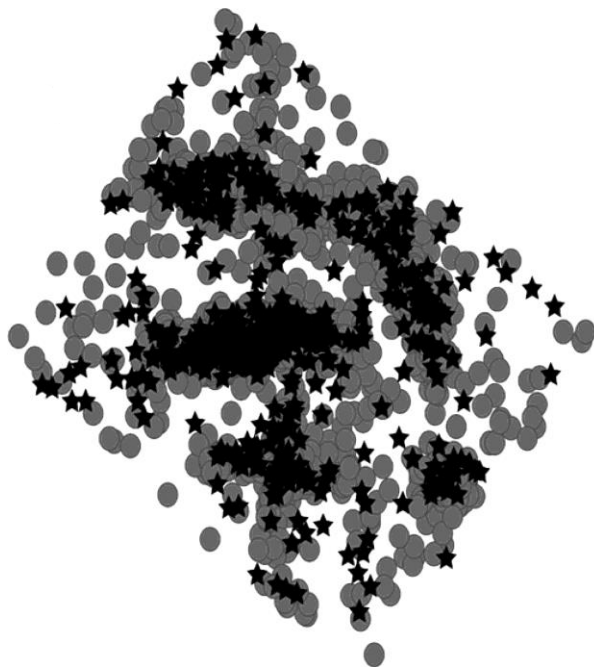


Overall Workflow

1. Dataset collection
- 2. Dataset split**
3. Descriptors
4. Performance metrics
5. Algorithms
6. Hyperparameter selection
7. Retrospective model evaluation
8. Prospective model evaluation



Chemical Series Bias



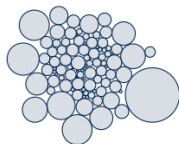
Random split

(Over optimistic performance)

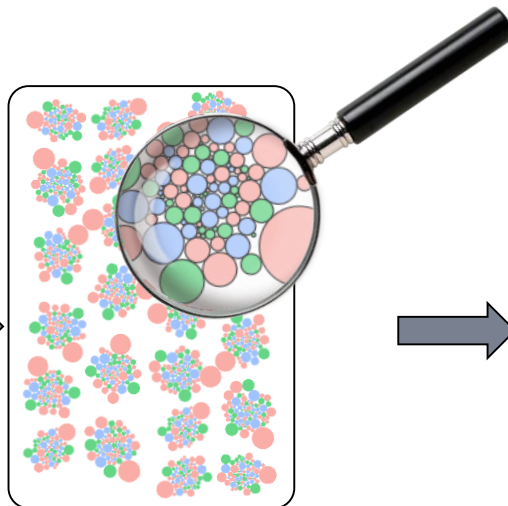
EXCAPE ML Chemical Series Bias: *a Realistic Dataset Split*



955K Compounds
526 Targets

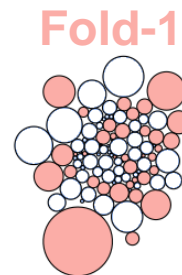


Cluster molecules
with ECFP6

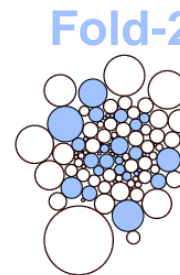


Search 3 folds with:

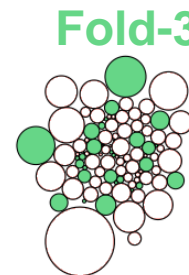
- 526 targets in each fold
- ~ even distribution of compounds



#molecules:
332,657



#molecules:
308,169

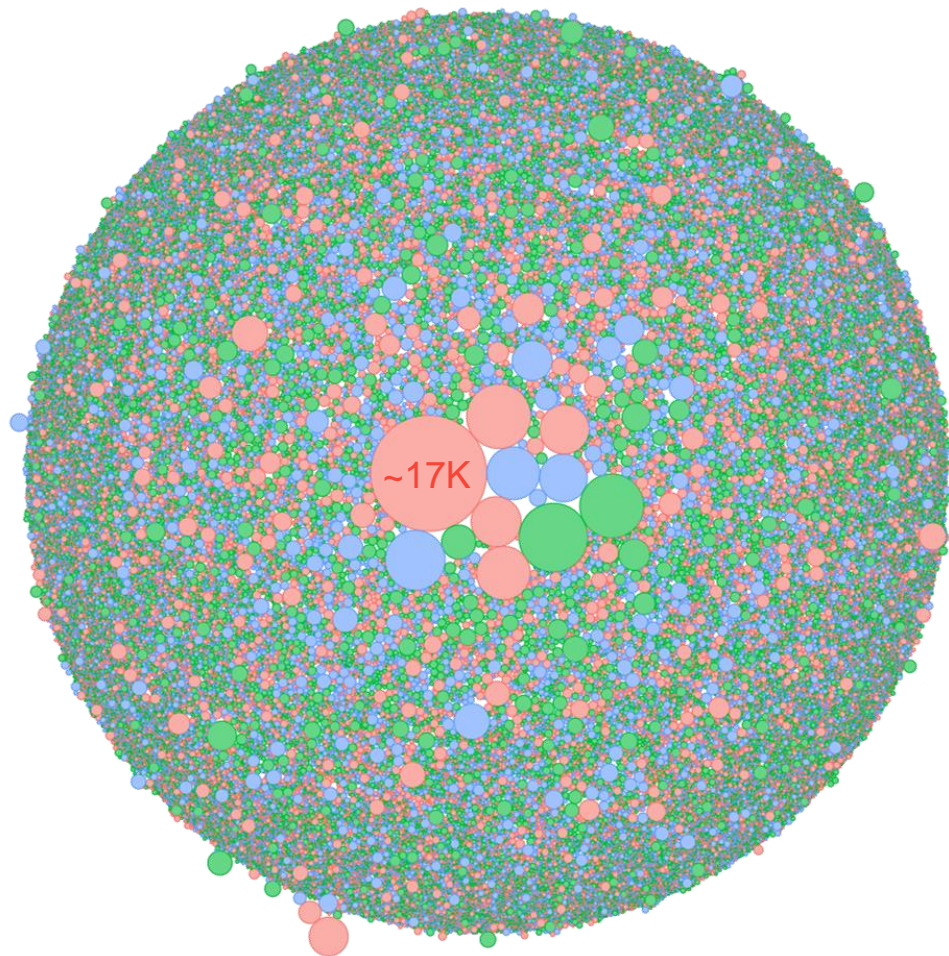


#molecules:
314,560



Cluster sizes

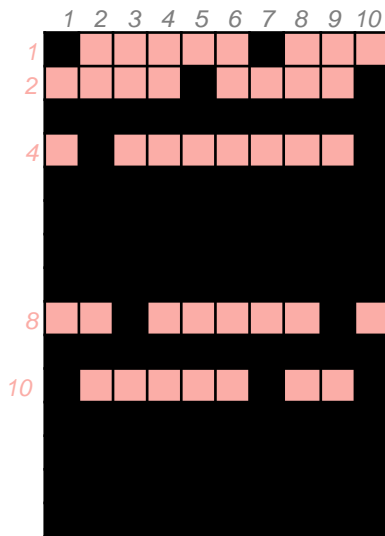
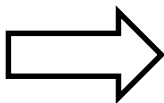
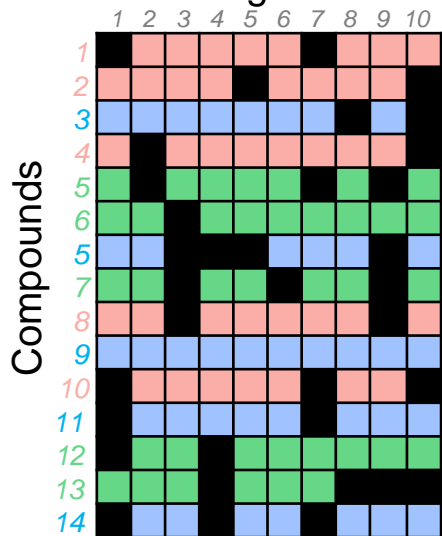
Fold	Clusters		Singletons	Total
	#cluster	#molec	#molec	#molec
1	20,614	320,451	12,206	332,657
2	21,204	295,557	12,612	308,169
3	20,592	302,330	12,230	314,560



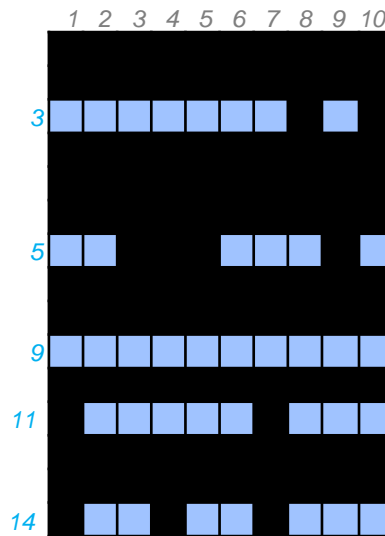
■ Fold-1 ■ Fold-2 ■ Fold-3

Chemogenomics Matrix [C x P]

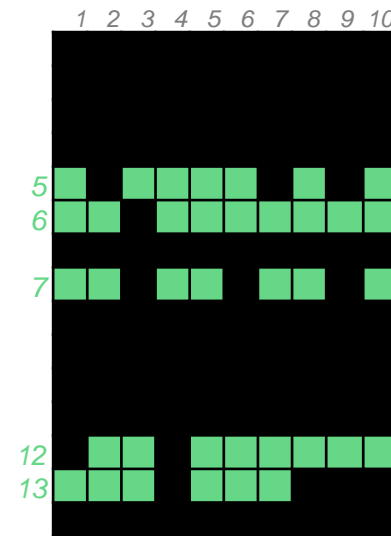
Targets



Fold-1



Fold-2



Fold-3

-  Fold-1
-  Fold-2
-  Fold-3
-  "Empty"

SciPy sparse matrices



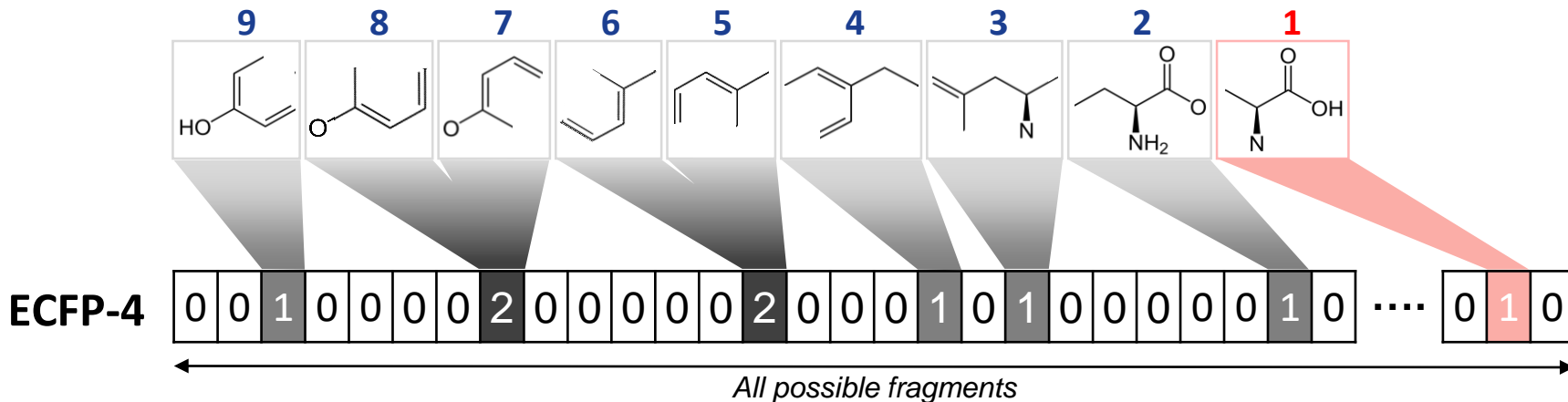
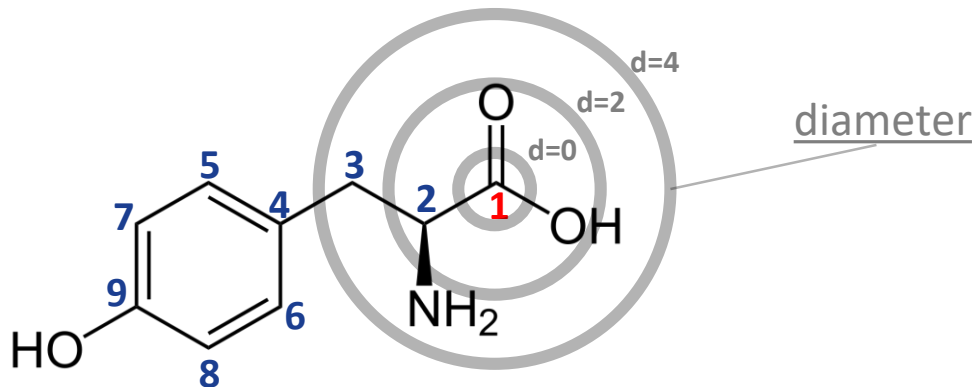


Overall Workflow

1. Dataset collection
2. Dataset split
- 3. Descriptors**
4. Performance metrics
5. Algorithms
6. Hyperparameter selection
7. Retrospective model evaluation
8. Prospective model evaluation

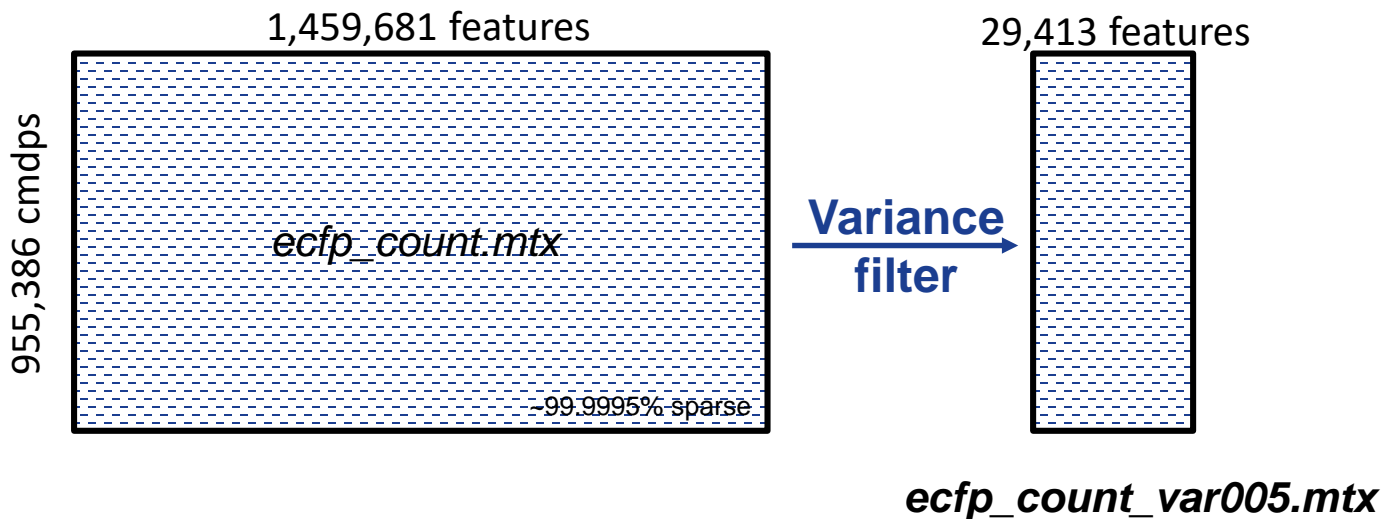


Molecular Features: *Extended Connectivity Fingerprint*





Molecular Features: *ECFP-6*



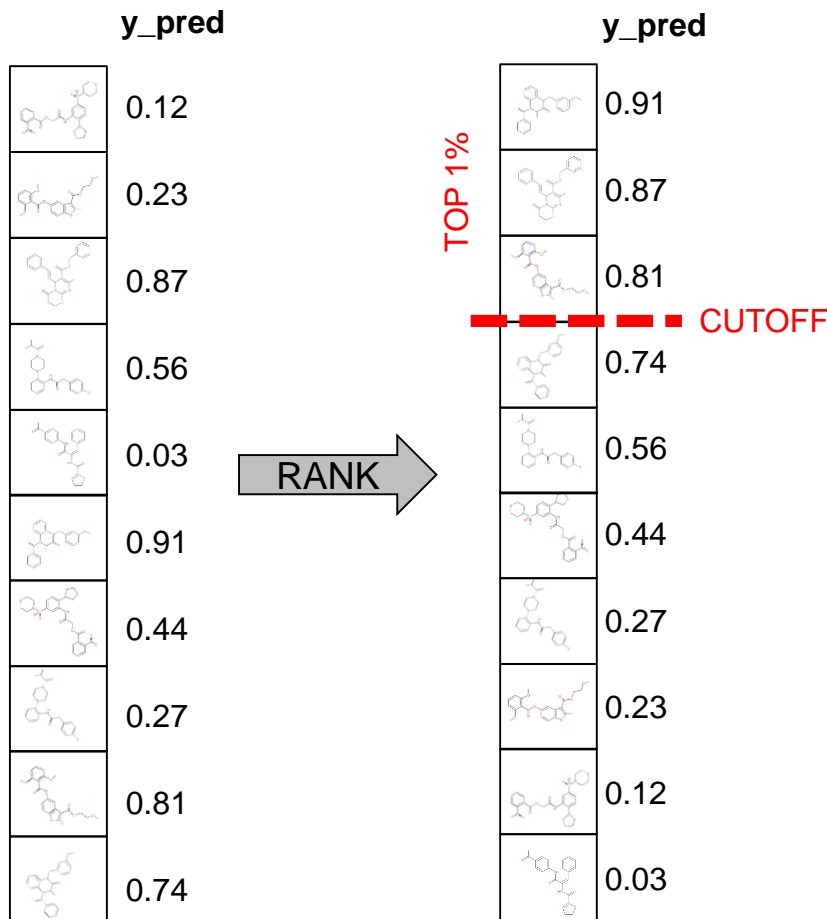


Overall Workflow

1. Dataset collection
2. Dataset split
3. Descriptors
- 4. Performance metrics**
5. Algorithms
6. Hyperparameter selection
7. Retrospective model evaluation
8. Prospective model evaluation



Virtual Screening: *prioritize active compounds*



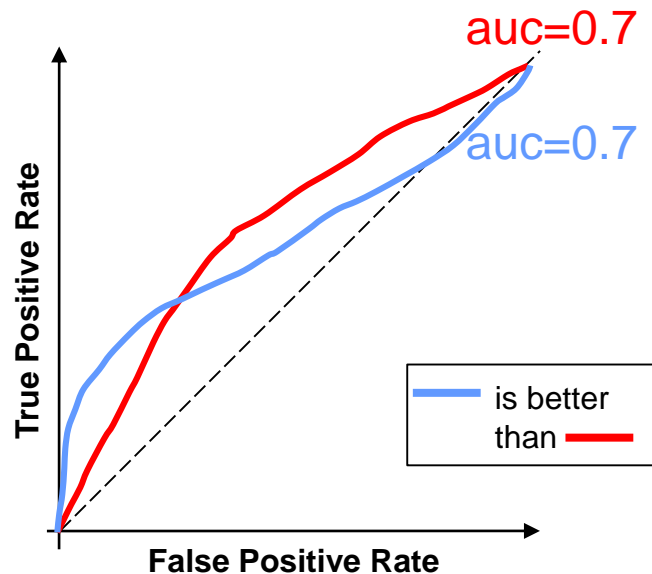
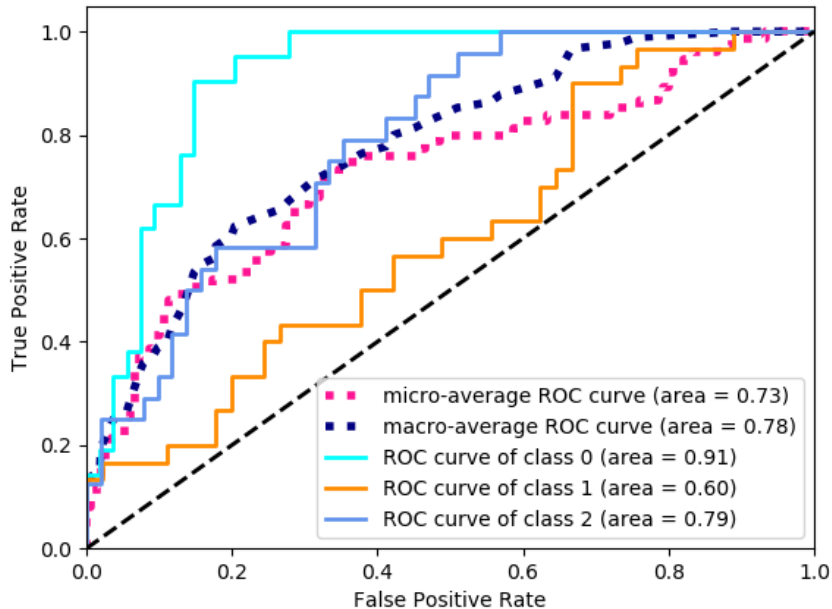
Active compound
enrichment
VS
HTS overall
hit rate



Performance metrics : ROC-AUC

$$ROCAUC = \int_{-\infty}^{+\infty} TPR(T)FPR'(T)dT$$

TPR : True Positive Rate | *FPR* : False Positive Rate





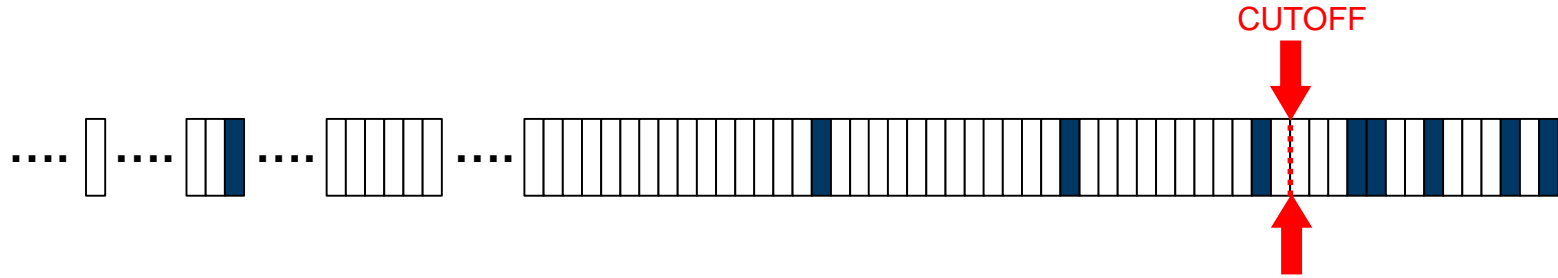
Performance metrics: *Confusion Matrix*

$$\mathbf{Precision} = \frac{TP}{TP+FP}$$

$$\mathbf{Recall} = \frac{TP}{TP+FN}$$

$$\mathbf{F1-score} = 2 \times \frac{\textit{precision} \cdot \textit{recall}}{\textit{precision} + \textit{recall}}$$

TP : True Positive | *FP* : False Positive | *FN* : False Negative | *TN* : True Negative





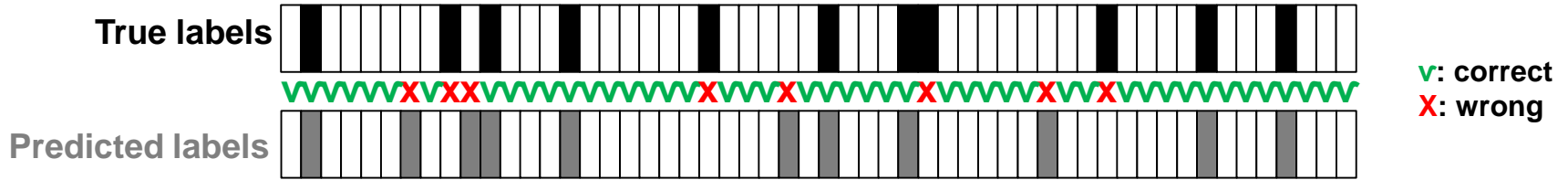
Performance metrics: *Cohen-Kappa*

$$\kappa = \frac{p_0 - p_e}{1 - p_e}$$

Inter-classification agreement

p_0 : Agreement between two classifiers (accuracy)

p_e : Expected agreement by chance of two classifiers





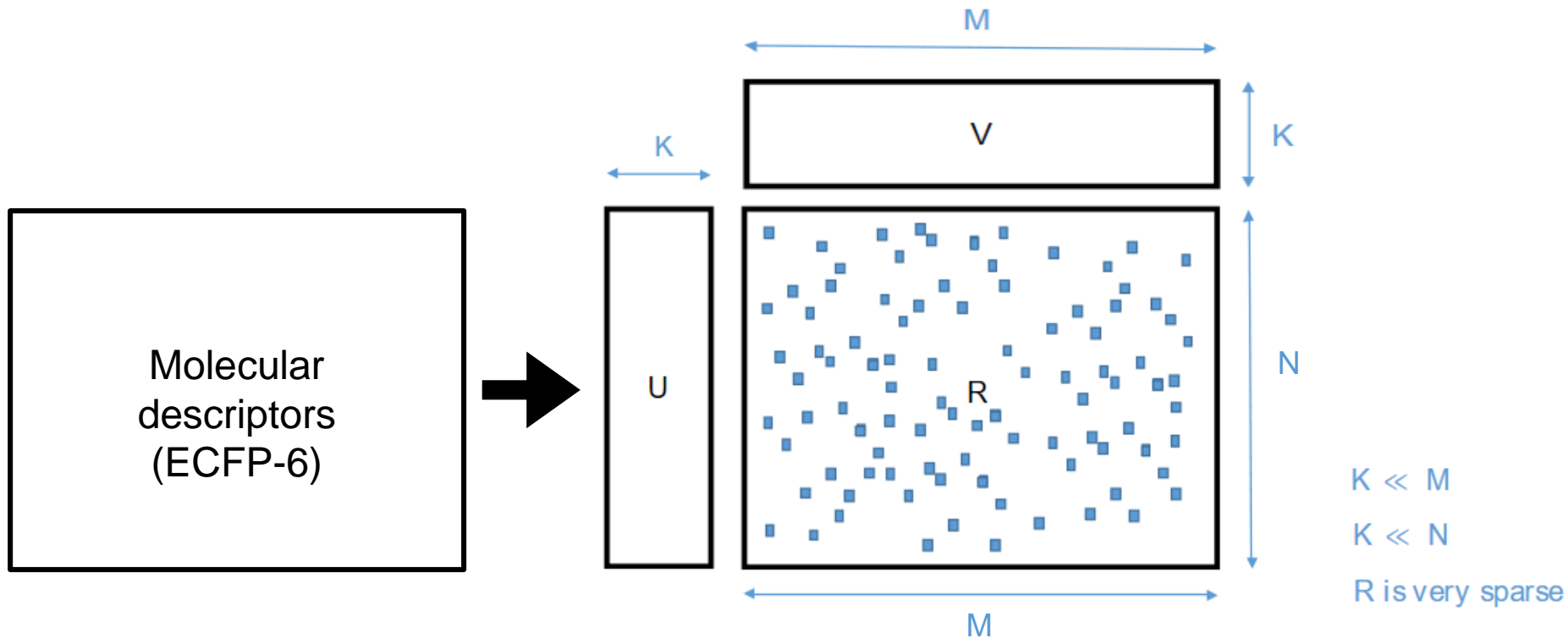
Overall Workflow

1. Dataset collection
2. Dataset split
3. Descriptors
4. Performance metrics
- 5. Algorithms**
6. Hyperparameter selection
7. Retrospective model evaluation
8. Prospective model evaluation



Matrix Factorization : *SMURFF* / *Macau*

umec Tom Vander
Tom Ashby

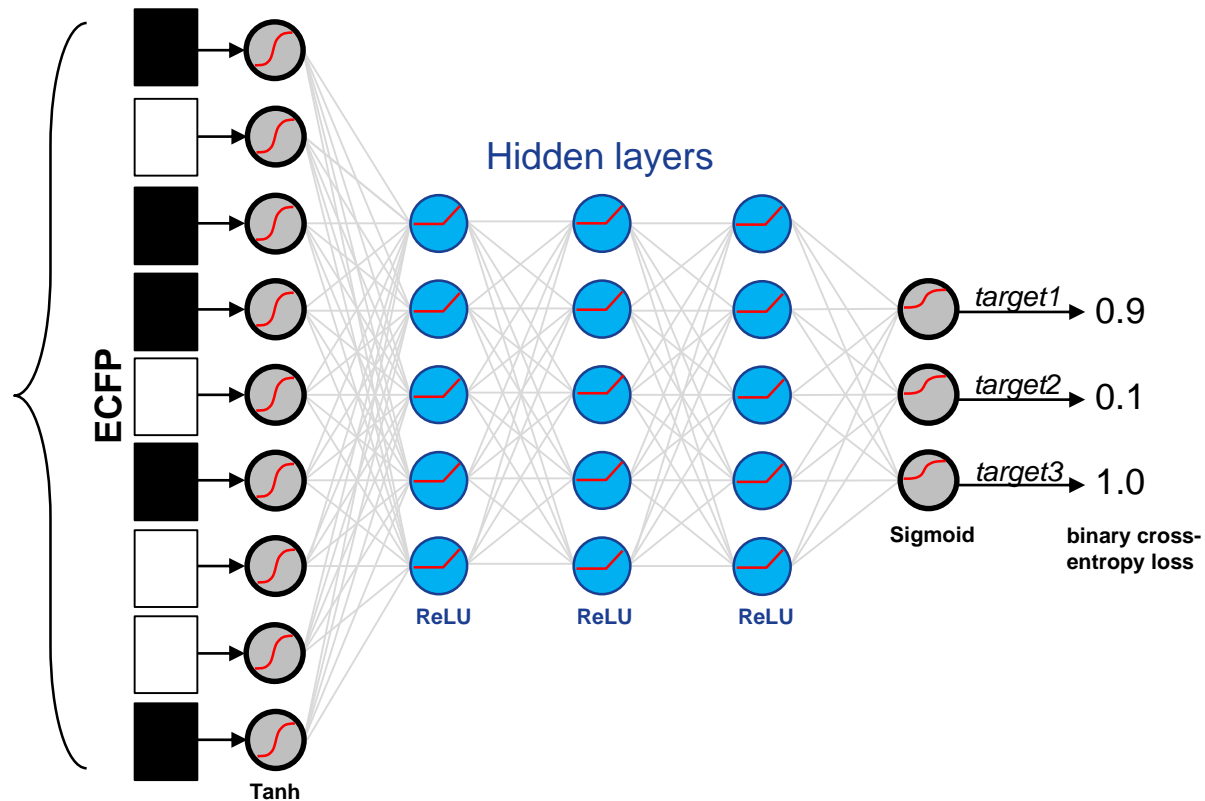
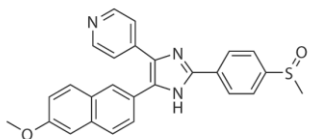




Deep Learning algorithm: *ExNET*

 JKU *Andreas Mayr*

 *Yves Vandriessche*





Overall Workflow

1. Dataset collection
2. Dataset split
3. Descriptors
4. Performance metrics
5. Algorithms
- 6. Hyperparameter selection**
7. Retrospective model evaluation
8. Prospective model evaluation



Hyperparameters & algorithms

SMURFF: MACAU

Parameters	Considered values
Dim latent space D	{8, 16, 32, 64}
Precision obs. TP	{1.0, 5.0, 10.0}
Precision feat. BP	{1.0, 5.0, 10.0}
# Samples	{100, 200, ..., 1700}

XGboost

Parameters	Considered values
Objective	binary:logistic
Booster	gbtree
Learning rate	0.05
Scale positive weight	{1, 5, 10}
Nbr estimators	{50, 100, 200}
Maximum depth	{5, 10}

Random Forest

Parameters	Considered values
Max features	auto
Class weight	balanced
Nbr estimators	{125, 256, 512}

SVM: LIBLINEAR

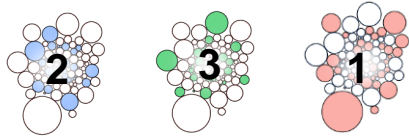
Parameters	Considered values
kernel	linear
C	{ 300, 100, 30,10, 1, 0.1, 0.05, 0.01, 0.001 }
penalty_loss_dual	("l2", "squared_hinge", True) ("l2", "squared_hinge", False) ("l2", "hinge", True) ("l1", "squared_hinge", True)
Sampling	none
Feature scaling	MaxAbsScaler

ExNET

Parameters	Considered values
Architecture	[1024x1024x1024], [2048x2048], [2048x2048x2048], [4096x4096x4096]
Learning Rate	{0.1, 0.01}
Input dropout	{0.0, 0.2}
Dropout	0.5
Momentum	{0.0, 0.4}
Activation Fct	ReLU
Minibatch size	64



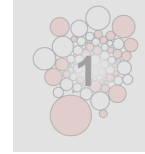
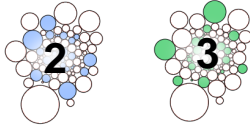
Model evaluation workflow: *Nested Cross-Validation*





Model evaluation workflow: *Nested Cross-Validation*

INNER FOLDS
Hyperparameter search



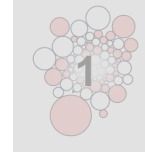
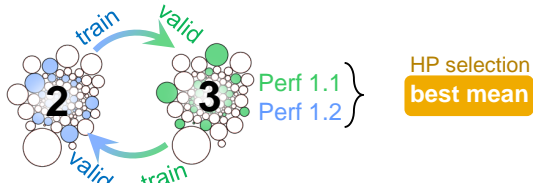
1. Leave one out



Model evaluation workflow: *Nested Cross-Validation*

INNER FOLDS

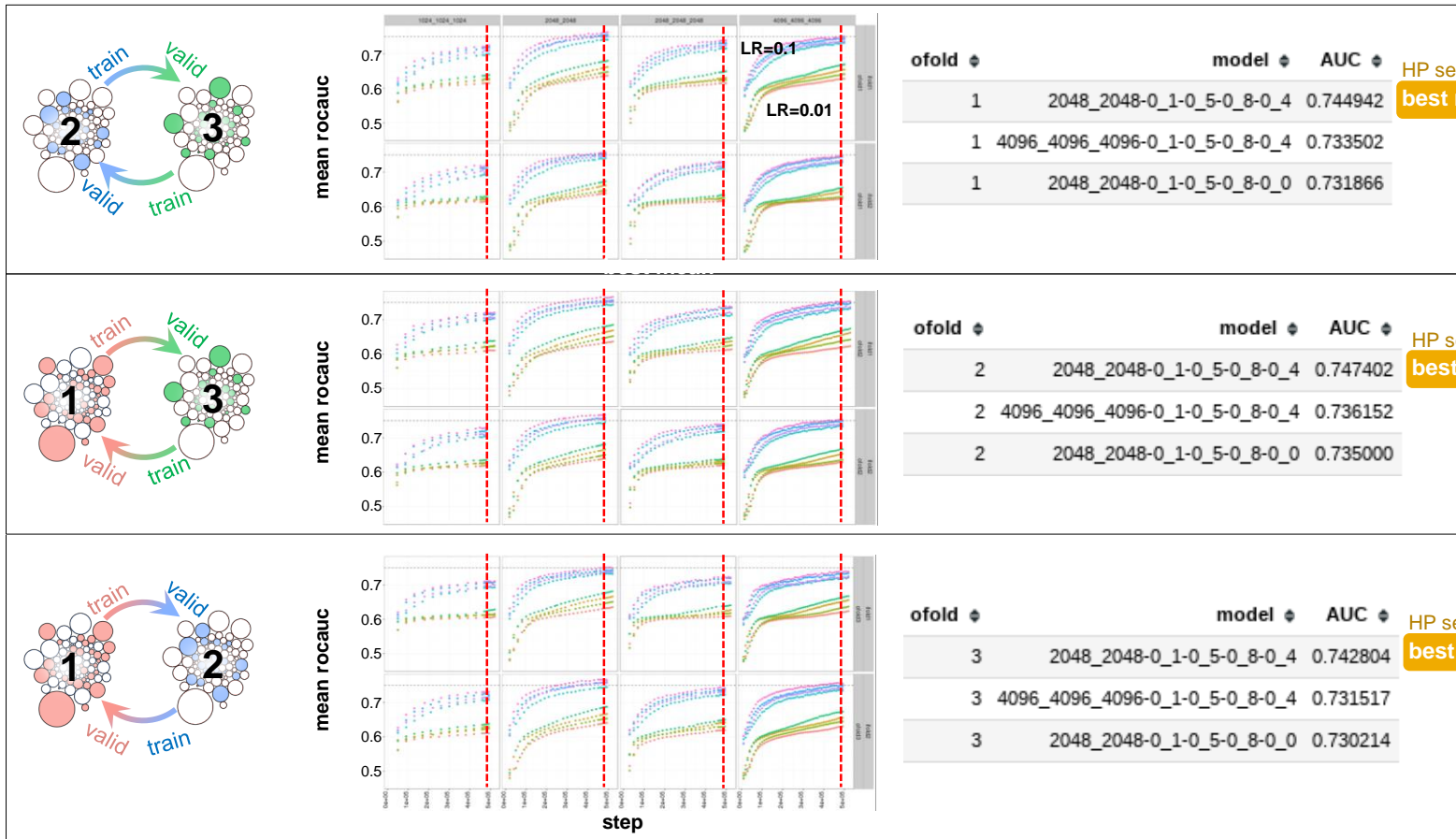
Hyperparameter search



1. Leave one out
2. Select hyperparam.



Hyperparameter Selection: *ExNET*



--- Last checkpoint in common (480k steps)

hyperSet

- 0_01-0_5-0_8-0_0
- 0_01-0_5-1_0-0_0
- 0_1-0_5-0_8-0_0
- 0_1-0_5-1_0-0_0
- 0_01-0_5-0_8-0_4
- 0_01-0_5-1_0-0_4
- 0_1-0_5-0_8-0_4
- 0_1-0_5-1_0-0_4





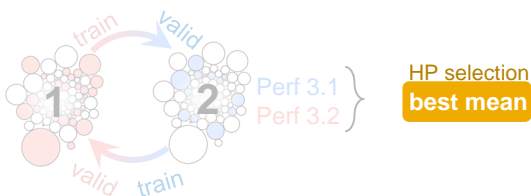
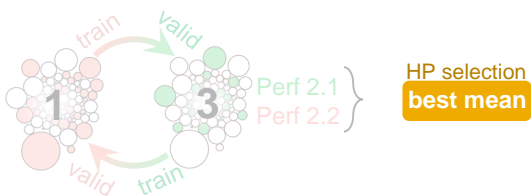
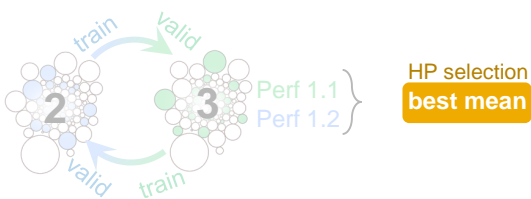
Overall Workflow

1. Dataset collection
2. Dataset split
3. Descriptors
4. Performance metrics
5. Algorithms
6. Hyperparameter selection
- 7. Retrospective model evaluation**
8. Prospective model evaluation

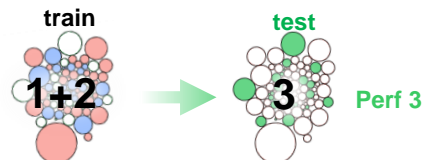
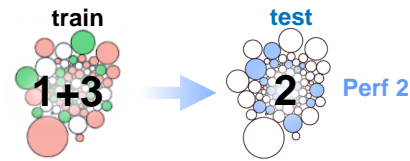
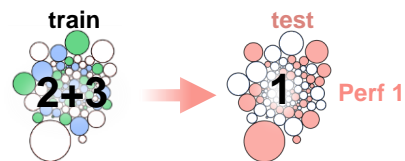


Model evaluation workflow: *Nested Cross-Validation*

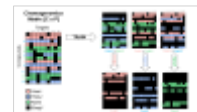
INNER FOLDS *Hyperparameter search*



OUTER FOLDS *Retrospective models testing*



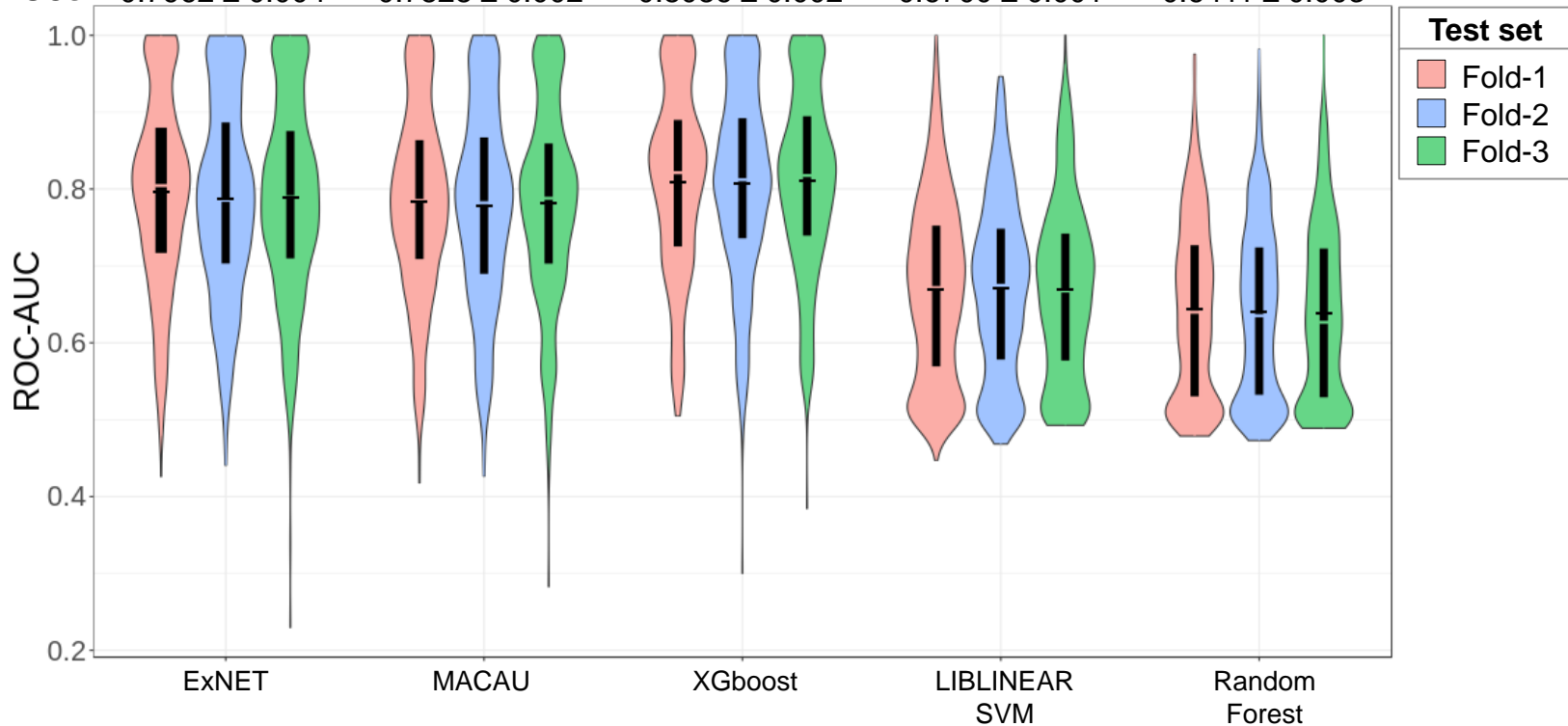
1. Leave one out
2. Select hyperparam.
3. Test models





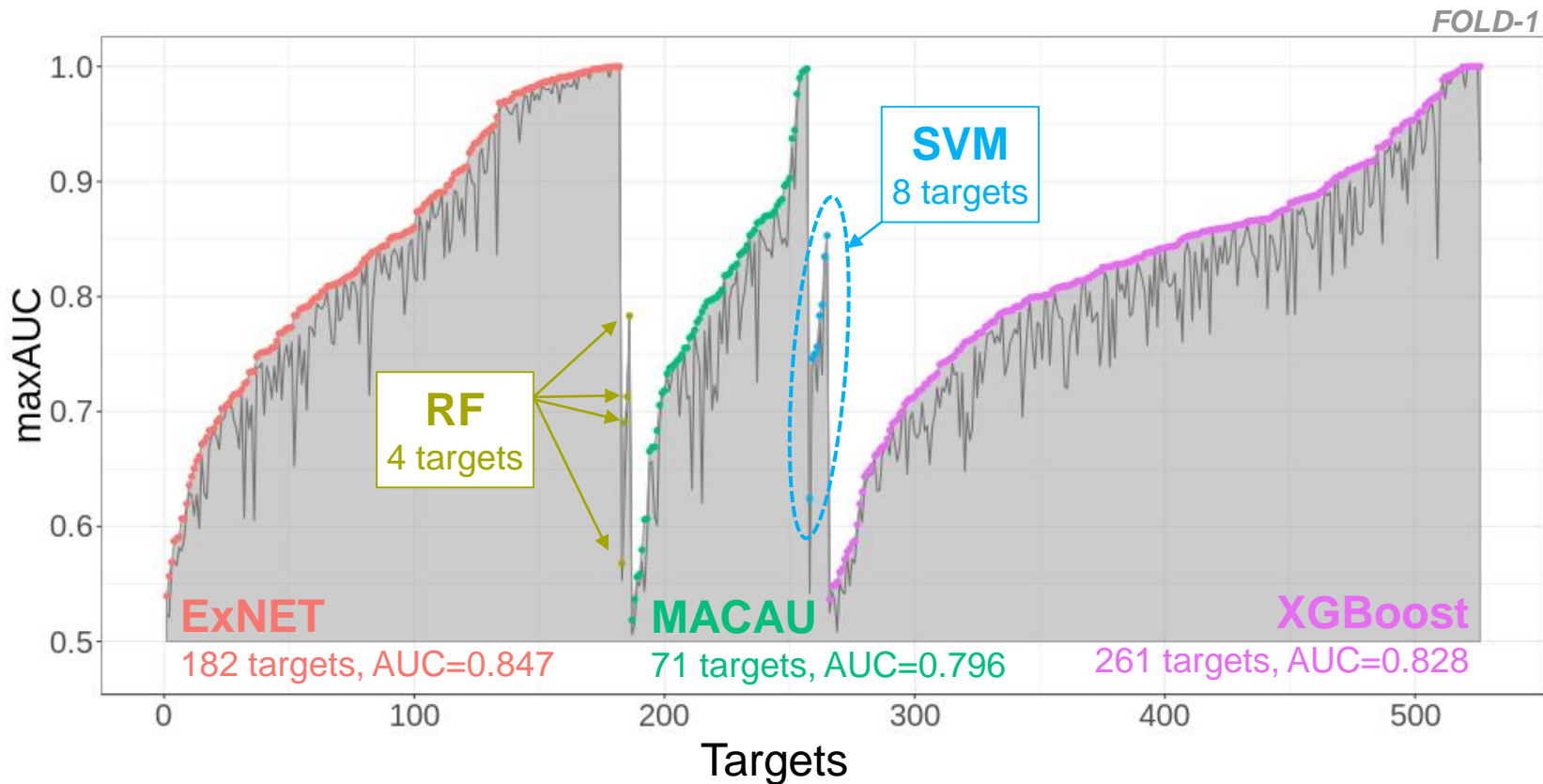
Test performances

Mean ROC-AUCs: 0.7962 ± 0.004 0.7828 ± 0.002 0.8086 ± 0.002 0.6700 ± 0.001 0.6411 ± 0.003



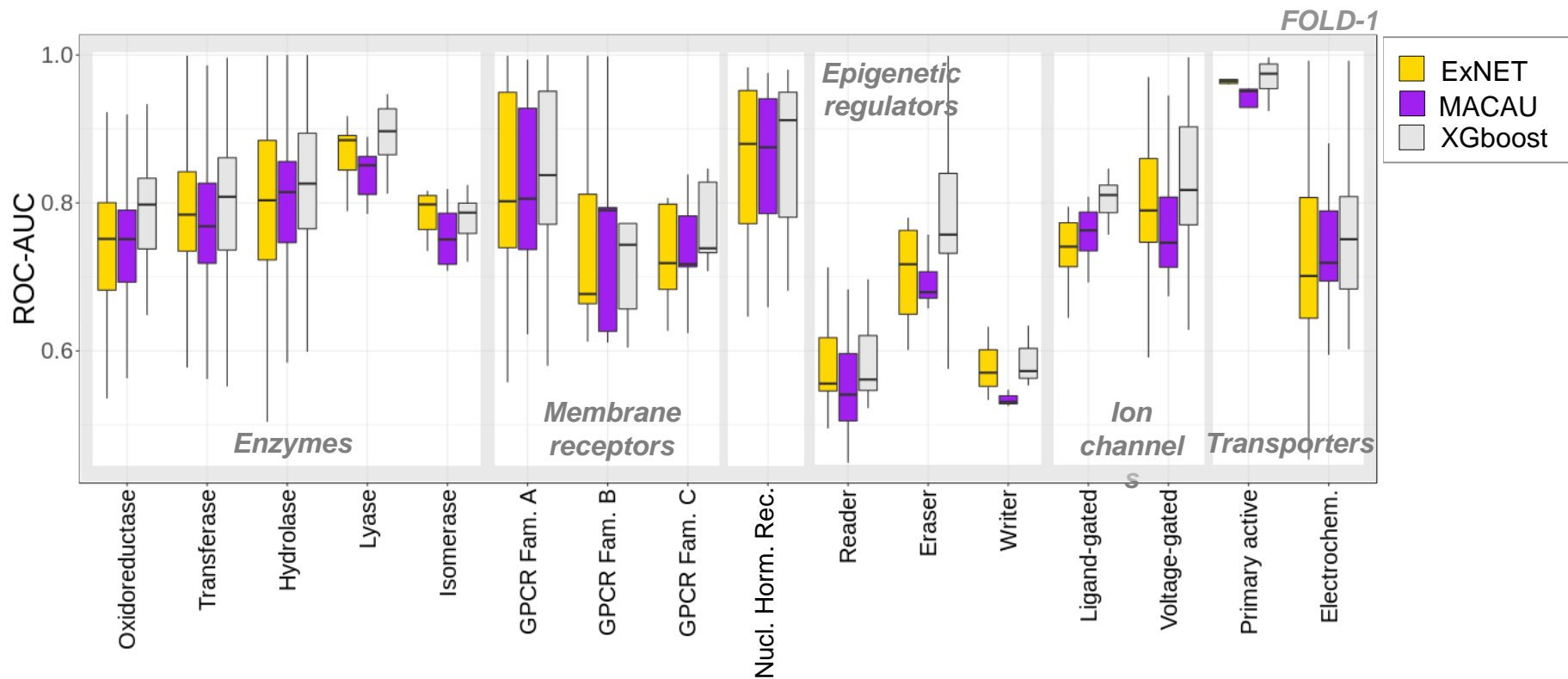


Winning algorithms



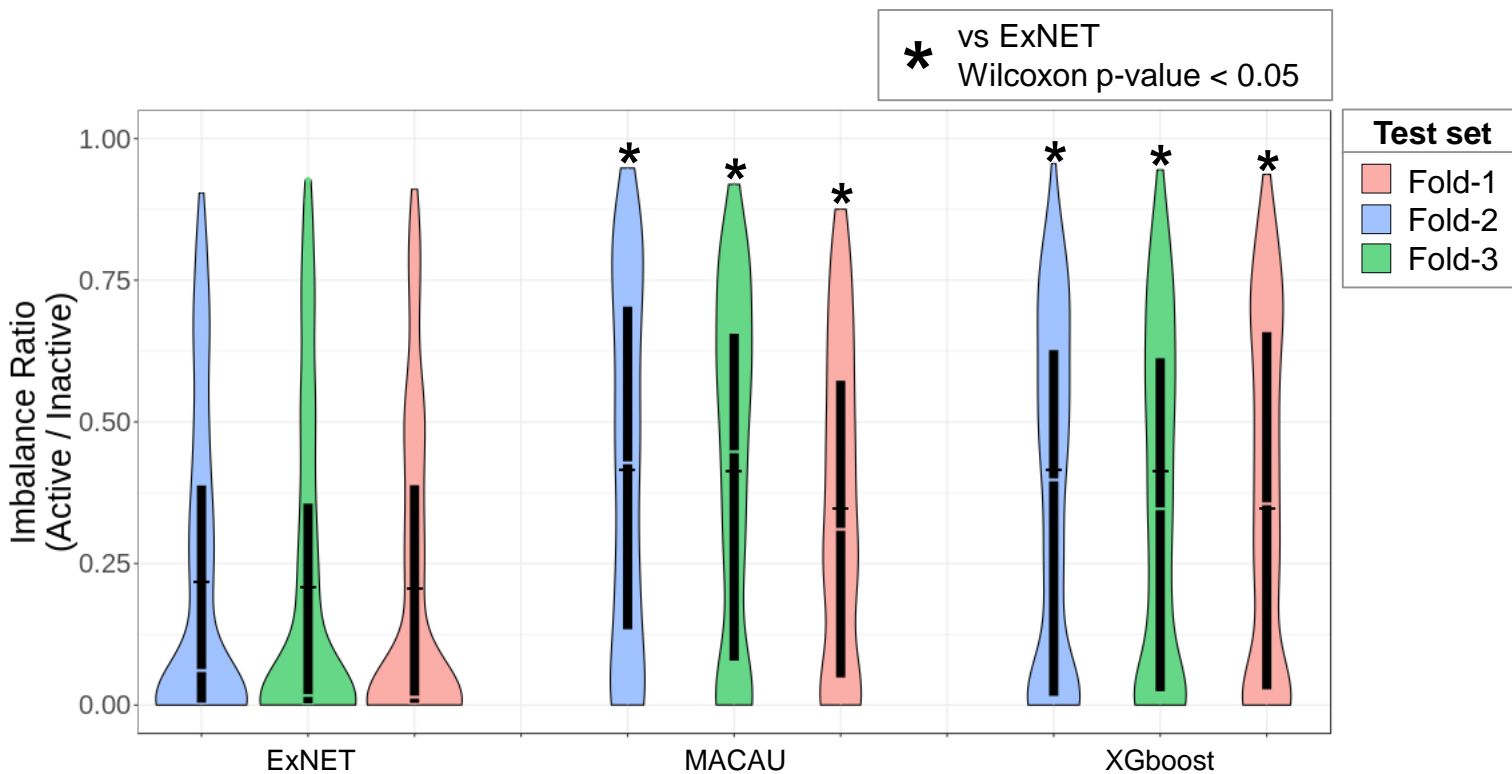


Performance per protein family





Imbalance Ratio of Target Datasets by Winning Algo.



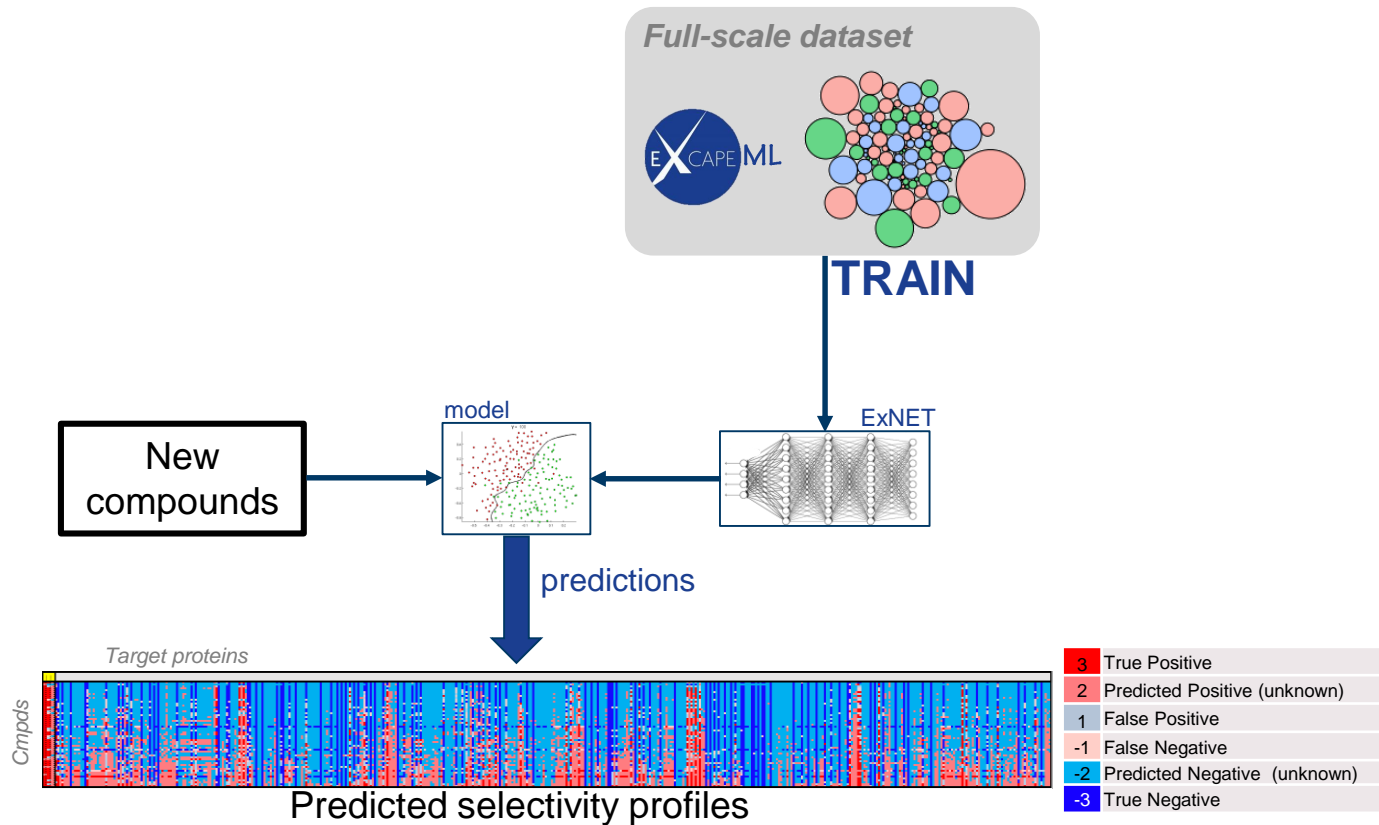


Overall Workflow

1. Dataset collection
2. Dataset split
3. Descriptors
4. Performance metrics
5. Algorithms
6. Hyperparameter selection
7. Retrospective model evaluation
- 8. Prospective model evaluation**

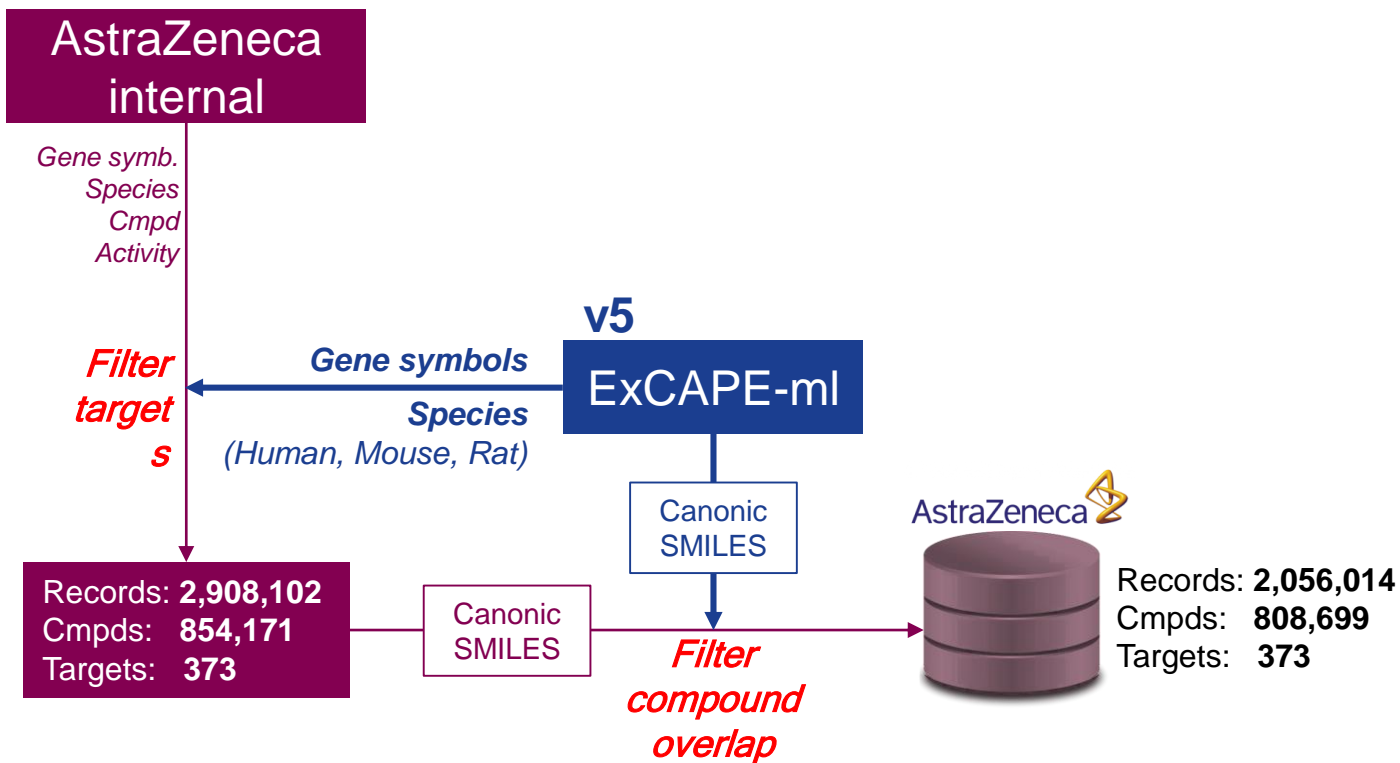


Prospective predictions





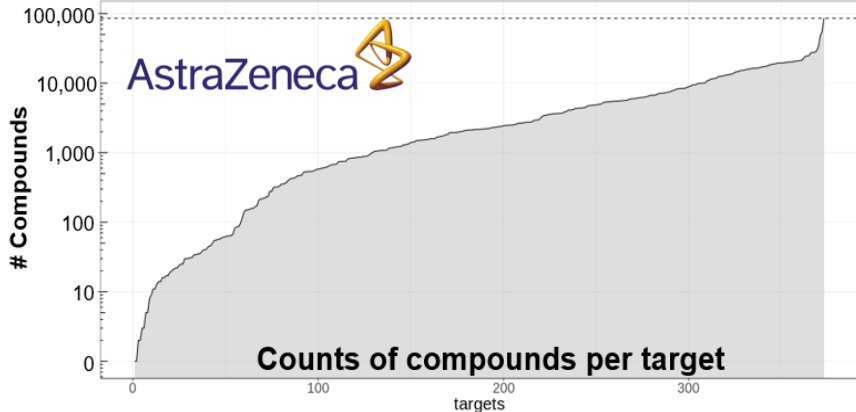
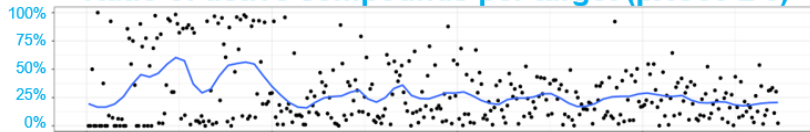
Dataset preparation: AZ





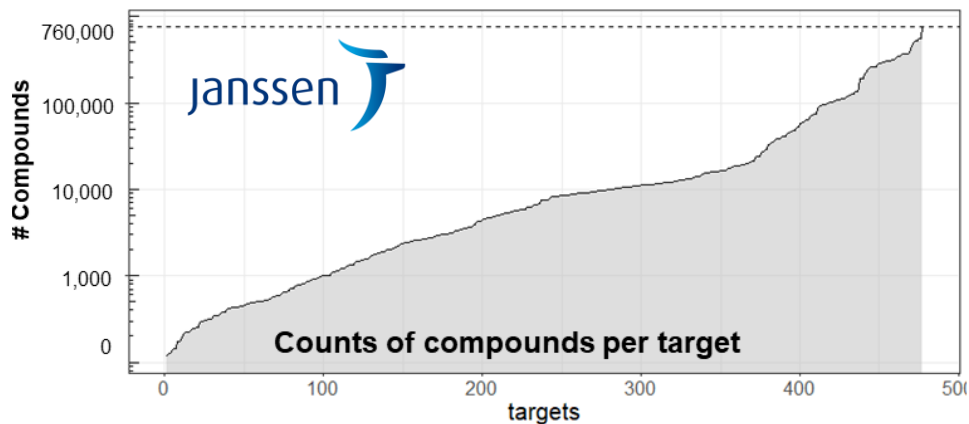
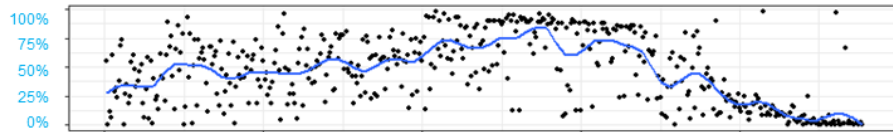
Industrial dataset profiles

Ratio of active compounds per target ($pXC50 \geq 6$)



~800K compounds
373 targets

Ratio of active compounds per target ($pXC50 \geq 6$)



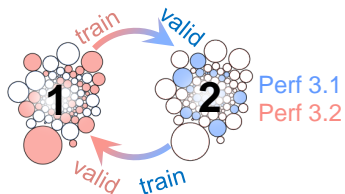
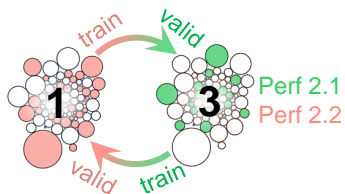
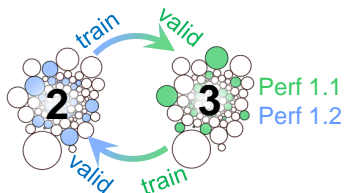
~1.7M compounds
467 targets



Application of Models to Industrial Datasets

INNER FOLDS

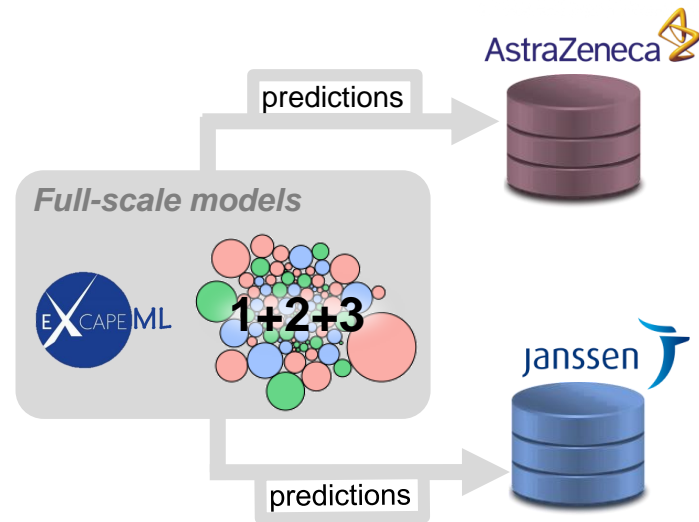
Hyperparameter search



HP selection
INNER FOLDS
best overall mean

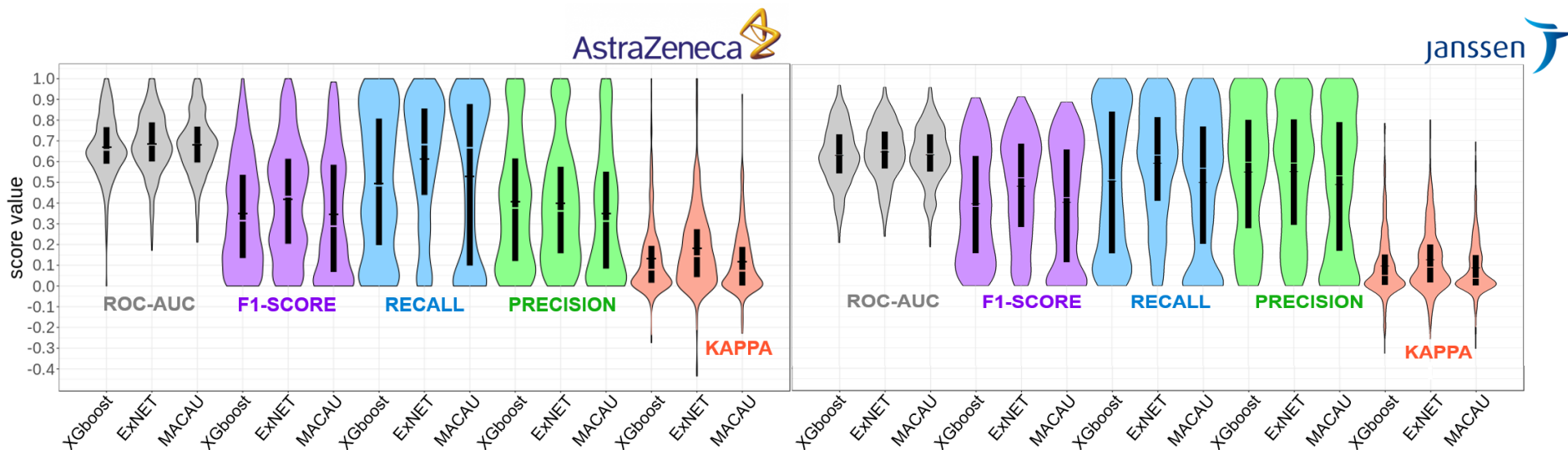
EXTERNAL TEST SETS

"Prospective" model testing



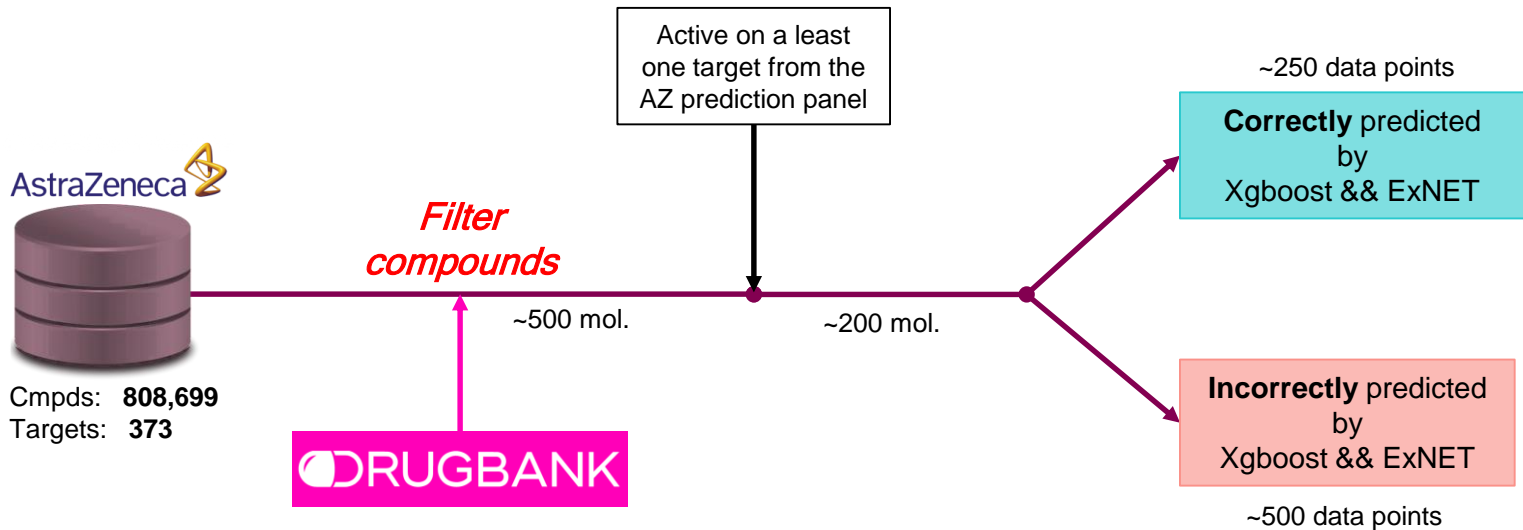


ExCAPE-ML full-scale model performance



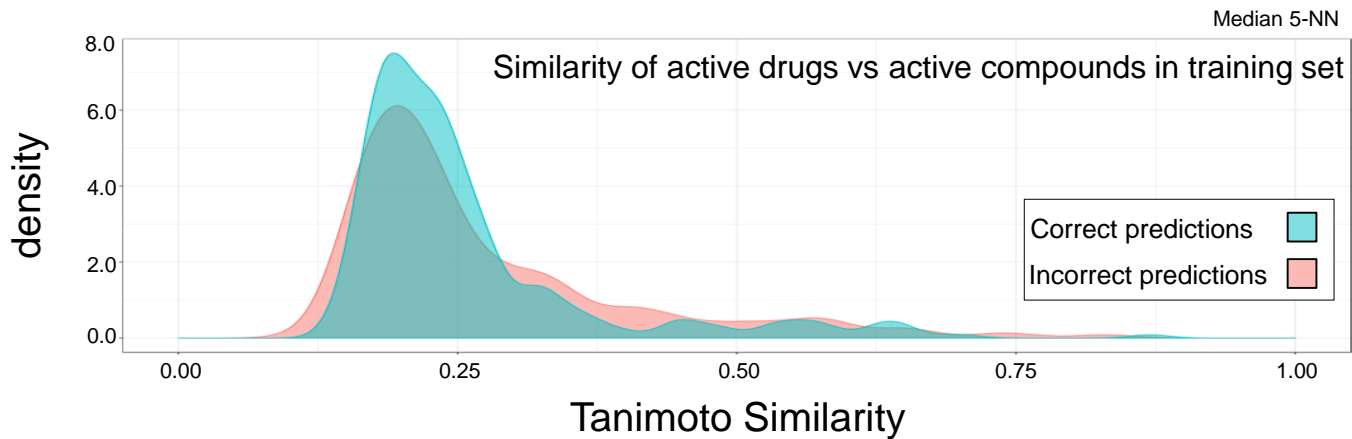


Investigating missed predictions...



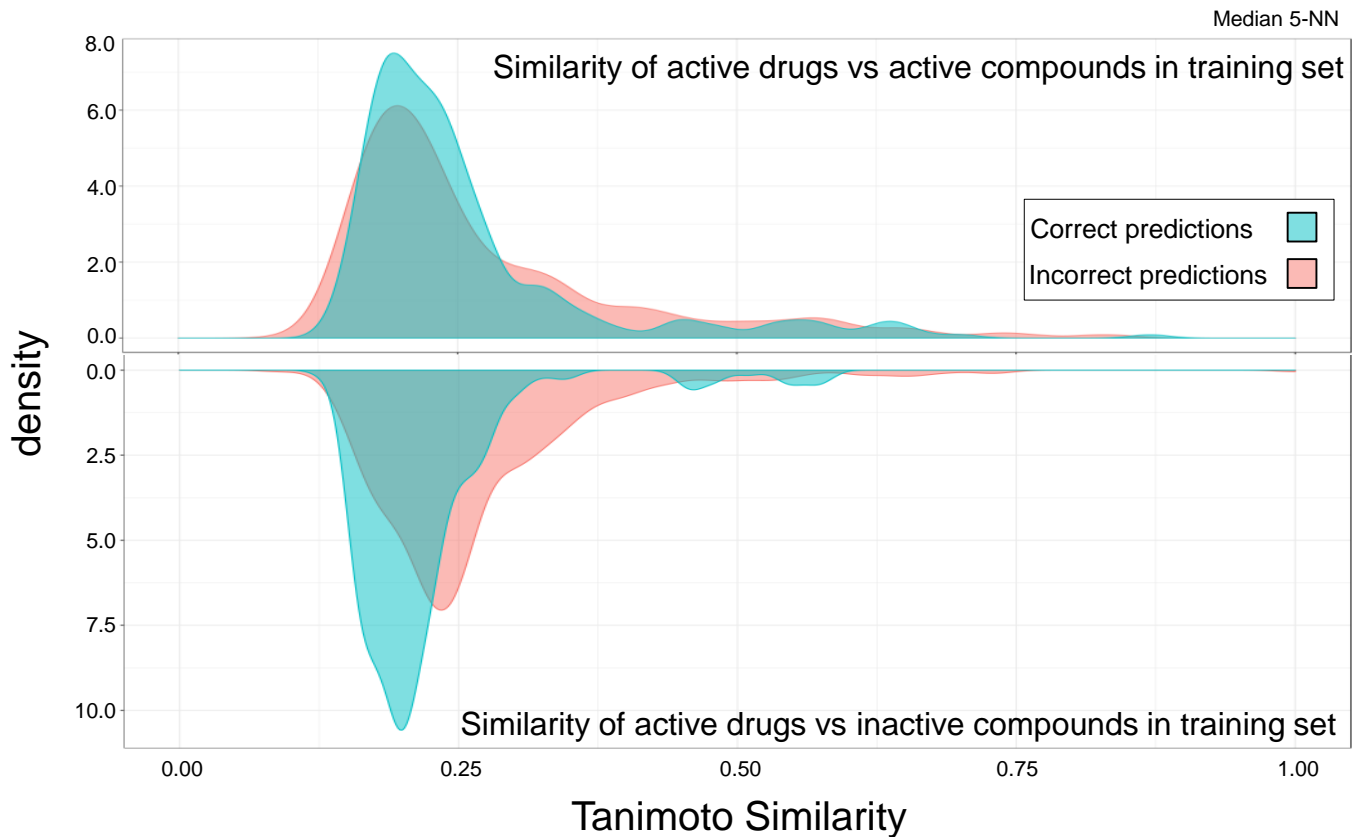


Investigating missed predictions...





Investigating missed predictions...





Acknowledgments

ExCAPE consortium



Jiangming Sun
Hongming Chen
Lars A Carlsson
Ernst Ahlberg Helgee
Ola Engkvist



Thanh Le Van
Felipe Golib
Vladimir Chupakhin



Nina Jeliaskova



Yves Vandriessche



Vojtech Cima



Tom Vander
Tom Ashby



Andreas Mayr
Günter Klambauer



Paolo Toccaceli



Xiangju Qin
Samuel Kaski