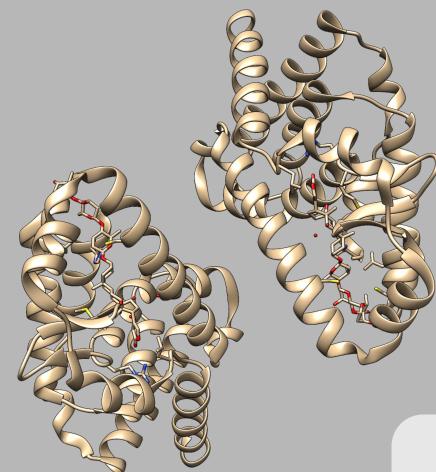


# Biological activity of nuclear receptors with RDKit

Rafał A. Bachorz

Institute of Medical Biology of  
Polish Academy of Sciences



# Agenda

- Introduction
- Nuclear receptors
- How to measure the biological activity?
- The QSAR workflow
- ZINC20 screening
- Summary

# Introduction

# Nuclear Receptors

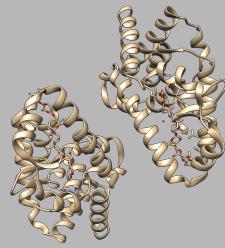
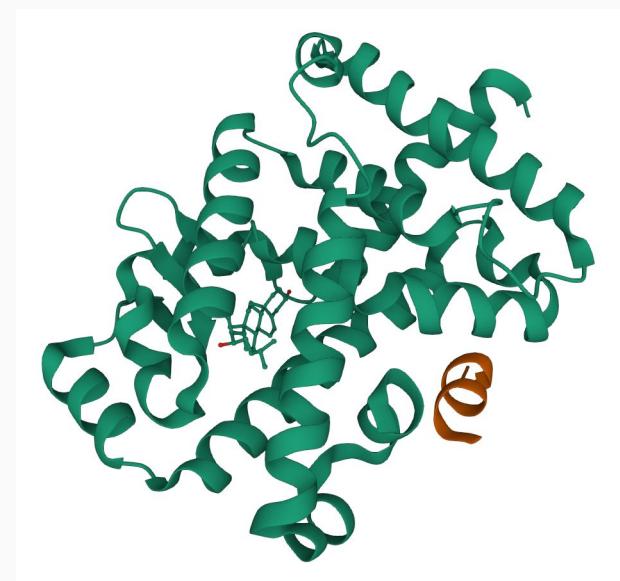


Open-Source Cheminformatics  
and Machine Learning



# Nuclear Receptors (NR)

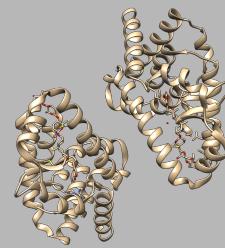
- Family of proteins
- RORs: the retinoic acid-related orphan receptors ( $\alpha$ ,  $\beta$ ,  $\gamma$ )
- Many low-weight compounds exert biological activity by binding to NRs
- The activity of NRs depends on the conformation change
- The conformation change can be initiated by binding small molecule to the protein moiety
- Binding of a ligand functions as a switch that induces a conformational switch
- Linked to many human diseases like: atherosclerosis, osteoporosis, autoimmunological disorders, obesity, asthma, and cancer



# How to measure the biological activity?

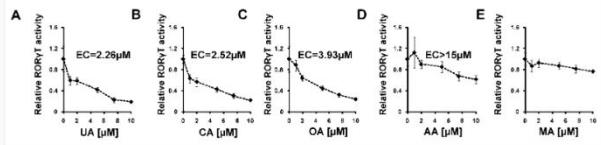


# Transfection and Luciferase Assay

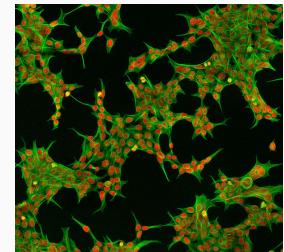


- The reporter vector: it turns into the luminescence protein after transfection and expression
- The GAL4-DBD ROR $\gamma$  fusion construct (containing human ROR $\gamma$  ligand binding domain)
- The HEK293 cell line
- Cotransfection with reporter vector and GAL4-DBD ROR $\gamma$  fusion construct
- Treatment of cotransfected cells with increasing concentrations of compounds
- Luminescence measurement reflects the inhibitory or activatory potential of considered compounds

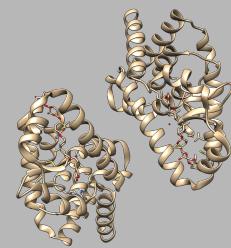
**Figure 5.** Effect of ursolic acid analogs on ROR $\gamma$ -dependent transcription in the HEK293 cell line. HEK293 cells were cotransfected with the pGL4.35[luc2P/9XGAL-4UAS/Hygro], GAL4-DBD ROR $\gamma$ , and pCMVSEAP vectors. Twenty-four hours later, the cells were treated with increasing concentrations of ursolic (**A**), corosolic (**B**), oleanolic (**C**), asiatic (**D**), and maslinic acids (**E**) for another 48 h. After that time, the cells were lysed, and luciferase activity was measured. Luciferase results are standardized for the transfection efficiency control, which was SEAP. Mean  $\pm$  SD,  $n = 3$ . EC50 values were calculated using AAT Bioquest (Sunnyvale, CA, USA) EC50 calculator, (<https://www.aatbio.com/tools/ec50-calculator/>, accessed on 6 December 2021).



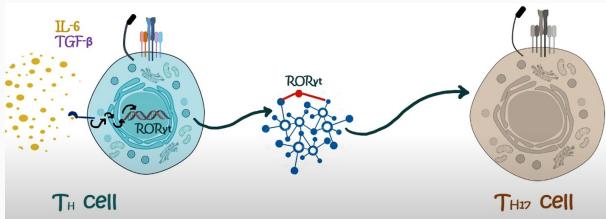
Pastwińska J, Karaś K, Salkowska A, Karwaciak I, Chałasięwicz K, Wojtczak BA, Bachorz RA, Ratajewski M. Identification of Corosolic and Oleanolic Acids as Molecules Antagonizing the Human ROR $\gamma$ T Nuclear Receptor Using the Calculated Fingerprints of the Molecular Similarity. *Int J Mol Sci.* 2022 Feb 8;23(3):1906. doi: 10.3390/ijms23031906. PMID: 35163824; PMCID: PMC8837092.



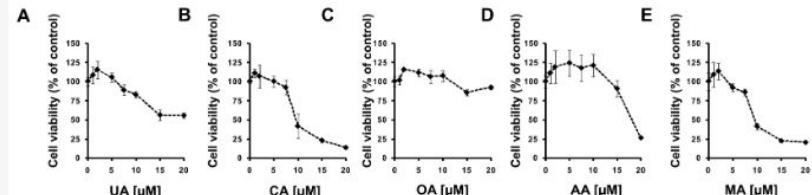
# Cytotoxicity towards Th17



- ROR $\gamma$ T are considered as master regulator of Th17 differentiation
- Cytotoxicity of considered compounds towards Th17 cells
- The CD4+ cells are subject to Th17 differentiation in the presence of increasing concentrations of selected compounds
- CellTiter-Glo® Luminescent Cell Viability Assay
- No noticeable cytotoxicity up to 7.5  $\mu$ M



**Figure 6.** (A–E) Effect of ursolic acid analogs on CD4+ lymphocyte viability. CD4+ cells were isolated from buffy coats of healthy donors and subjected to Th17 polarization in the presence of increasing concentrations of ursolic, corosolic, oleanolic, asiatic, and maslinic acids for 5 days. Then, cell viability was determined using the CellTiter-Glo® Luminescent Cell Viability Assay. Mean  $\pm$  SD,  $n = 3$ , compared with control cells.

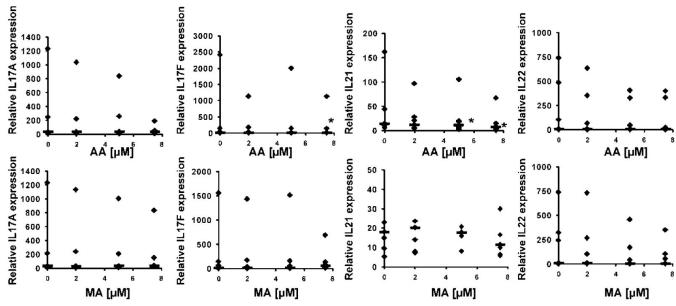


Pastwińska J, Karaś K, Salkowska A, Karwaciak I, Chałasiakiewicz K, Wojtczak BA, Bachorz RA, Ratajewski M. Identification of Corosolic and Oleanolic Acids as Molecules Antagonizing the Human ROR $\gamma$ T Nuclear Receptor Using the Calculated Fingerprints of the Molecular Similarity. *Int J Mol Sci.* 2022 Feb 8;23(3):1906. doi: 10.3390/ijms23031906. PMID: 35163824; PMCID: PMC8837092.



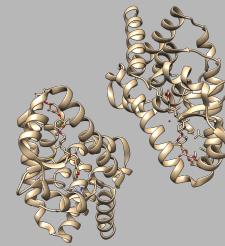
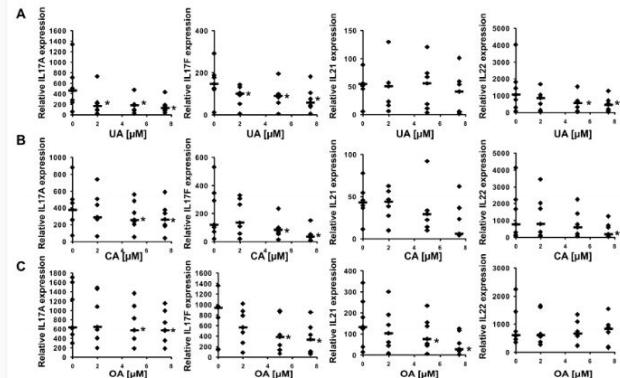
# CD4+ → Th17 differentiation

- The CD4+ cells are isolated from buffy coats
- They are differentiated into Th17 lymphocytes (protocol by Wilson et al.) in the presence of the compounds being verified
- Expression of ROR $\gamma$ T, IL17A, IL17F, IL21, IL22 and APOD was measured using real time RT-PCR methodology



Pastwińska J, Karaś K, Salkowska A, Kawaciak I, Chałasięwicz K, Wojtczak BA, Bachorz RA, Ratajewski M. Identification of Corosolic and Oleanolic Acids as Molecules Antagonizing the Human ROR $\gamma$ T Nuclear Receptor Using the Calculated Fingerprints of the Molecular Similarity. *Int J Mol Sci.* 2022 Feb 8;23(3):1906. doi: 10.3390/ijms23031906. PMID: 35163824; PMCID: PMC8837092.

**Figure 7.** Effect of ursolic acid analogs on the expression of selected genes in human Th17 cells. Human naive CD4+ cells were treated with increasing concentrations of ursolic (A), corosolic (B), and oleanolic (C) acids and cultured under Th17-polarizing conditions for 5 days. Then, cells were collected for RNA extraction. The expression of the *IL17A*, *IL17F*, *IL21*, and *IL22* genes was determined by real-time RT-PCR. The results were normalized to the housekeeping genes *Hprt1*, *HMBS*, and *RPL13A*. An asterisk indicates a statistically significant difference at  $p < 0.05$  compared with control cells. The data are presented as statistical dot plots with the median value (bars) from seven independent cultures ( $n = 7$ ).



# The QSAR workflow

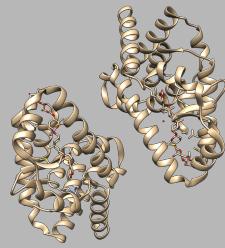
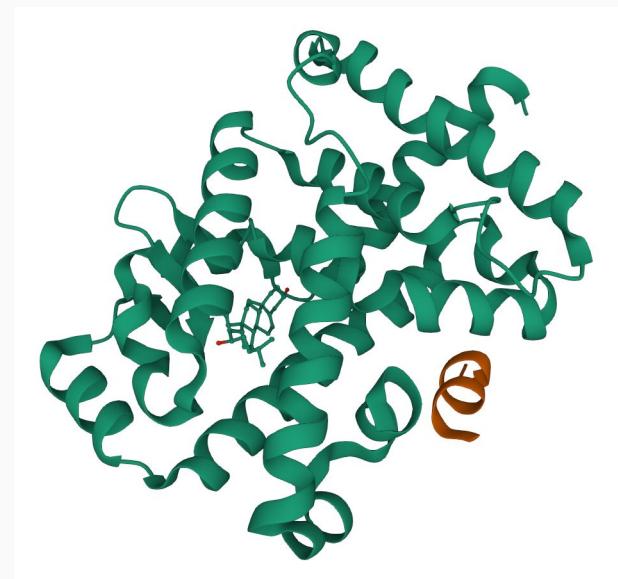


Open-Source Cheminformatics  
and Machine Learning



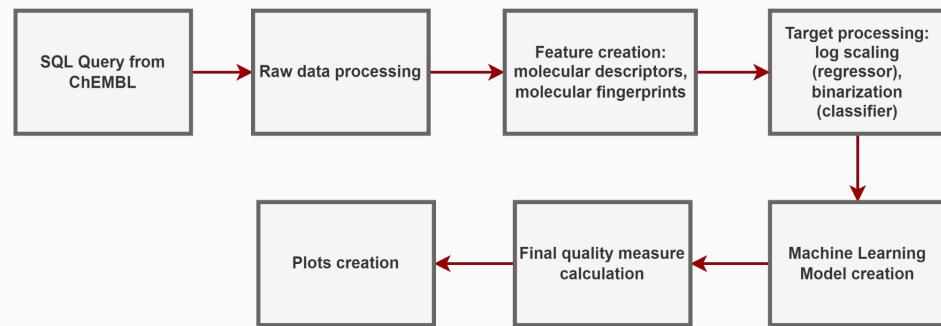
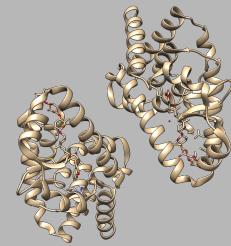
# Modeling of biological activity

- Generic framework for QSAR studies
- QSAR approach
- The 2D (potentially also 3D) molecular descriptors and molecular fingerprints as features
- Fully configurable data processing workflow
- Implemented in the Python ecosystem
- Only open source libraries (RDKit and Mordred on the Cheminformatics site)
- Web application written with Jinja templates
- The source code will be published



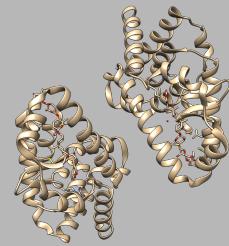


# Modeling of biological activity: workflow

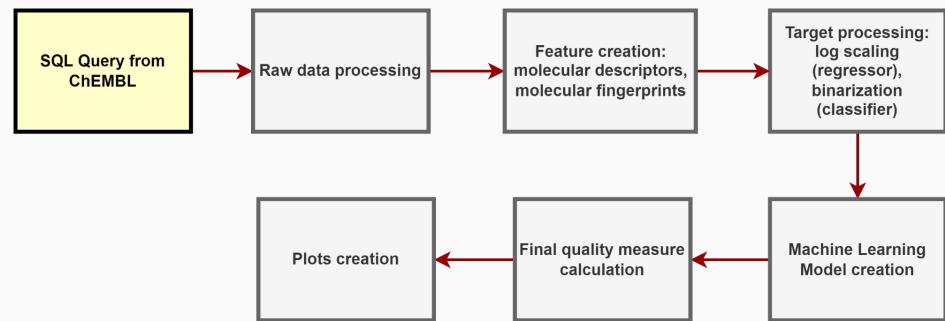




# Modeling of biological activity: workflow



- Literally the SQL query
- Local instance of the ChEMBL database (dockerized setup)
- The query brings the raw data
- Can be carried out for any target (ChEMBL ID needed)



```
1 def create_query(receptor):  
2     query = "select * from (\`  
3         select * from activities as act \  
4             inner join assays as assays on act.assay_id=assays.assay_id \  
5             left join docs as docs on act.doc_id=docs.doc_id \  
6             join target_dictionary as tardi on assays.tid=tardi.tid \  
7             inner join compound_structures as cs on act.molregno=cs.molregno \  
8             inner join compound_properties as cp on act.molregno=cp.molregno \  
9             left join molecule_dictionary as md on act.molregno=md.molregno \  
10            left join ligand_eff as le on act.activity_id=le.activity_id \  
11            where tardi.chembl_id in ('"+receptor+"')) as combined"  
12     return query
```

The ChEMBL database of molecules in a Docker environment

Rafal A. Bachors, PhD

Data Sources, Head of Advanced Analytics Team at Pfizer, Head of Molecular Biology at Polish Academy of Sciences

More from Medium

Python 3.6 with Docker that keep you happy all day

Iran Karimi

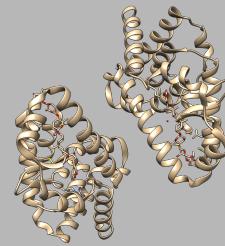
Small Python script that'll keep you happy all day

These books

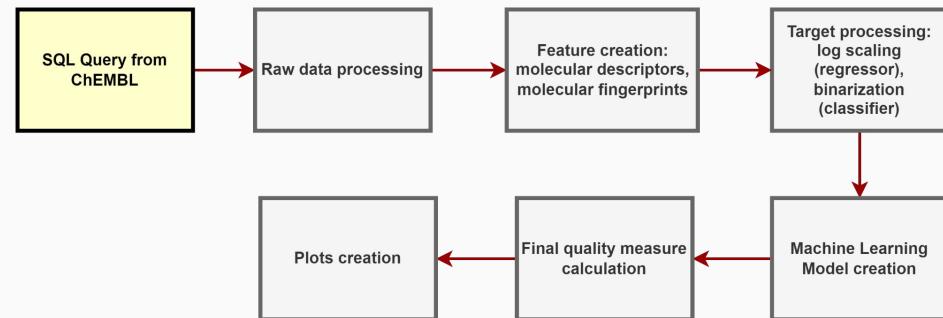
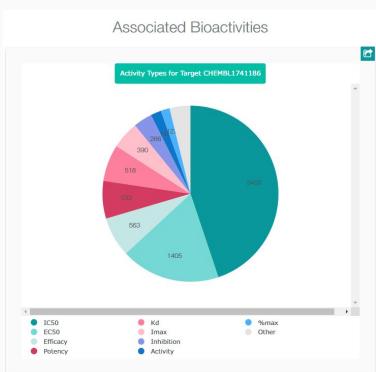
The 9 books that completely changed the way I see the world



# Modeling of biological activity: workflow

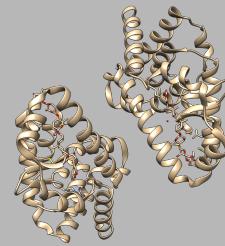


- Literally the SQL query
- Local instance of the ChEMBL database (dockerized setup)
- The query brings the raw data
- Can be carried out for any target (ChEMBL ID needed)

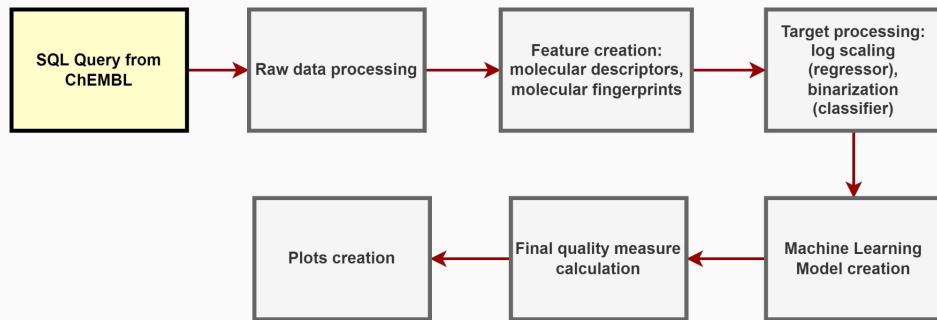
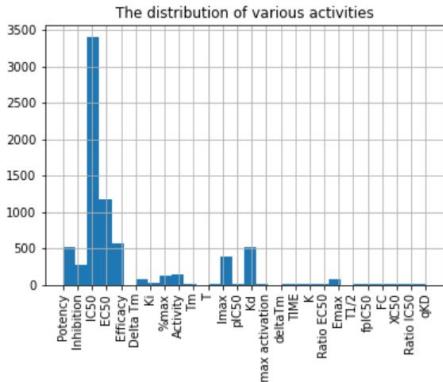




# Modeling of biological activity: workflow



- Literally the SQL query
- Local instance of the ChEMBL database (dockerized setup)
- The query brings the raw data
- Can be carried out for any target (ChEMBL ID needed)



Rafal A. Bachors, PhD  
Data Scientist, Head of Advanced Analytics Team at PSL Polaris, Head of Molecular Modeling at Institute of Chemical Biology of Polish Academy of Sciences  
bio profile

The ChEMBL database of molecules in a Docker environment

Rafal A. Bachors, PhD  
Data Scientist, Head of Advanced Analytics Team at PSL Polaris, Head of Molecular Modeling at Institute of Chemical Biology of Polish Academy of Sciences  
bio profile

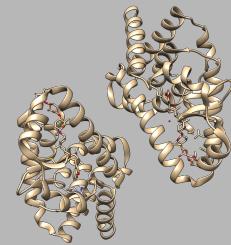
Docker Mechanics

Search

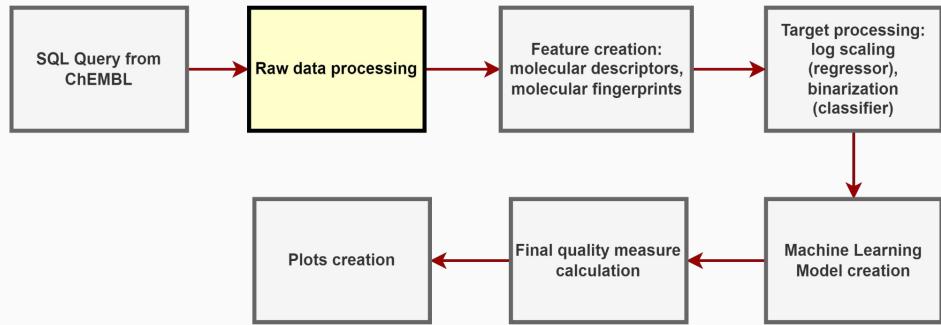
More from Medium

Python 3.6 with Numpy that's...  
Ian Kaiser  
Small Python script that'll keep you from being a dog  
Paul Braren  
The 9 books that completely changed the way I see the world

# Modeling of biological activity: workflow

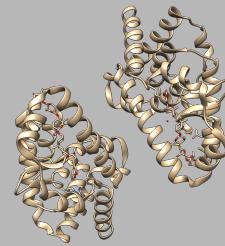


- Raw data sometimes confusing
  - Multiple entries for the same molecule
  - Sometimes contradictory
  - The processing strategy
    - depending on std value the observation can be removed
    - for the remaining the aggregation strategy (mean, median)

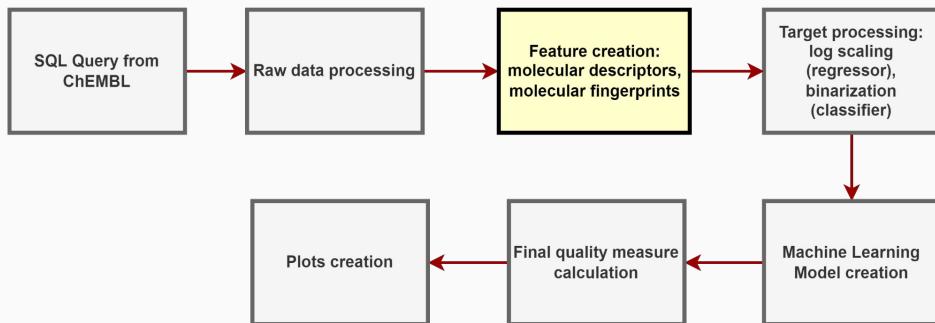
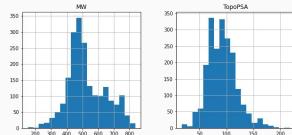
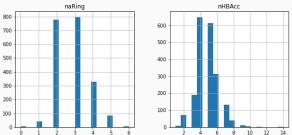
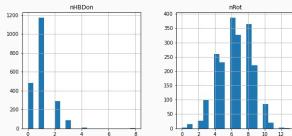
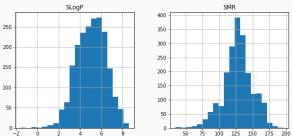




# Modeling of biological activity: workflow

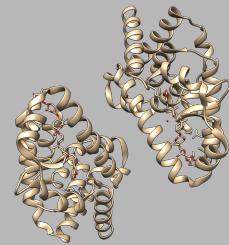


- The molecular descriptors with Mordred/RDKit
- The molecular fingerprints with RDKit
- The features can be transformed with PCA
  - Important for fingerprints taken by neural networks (sparsity)

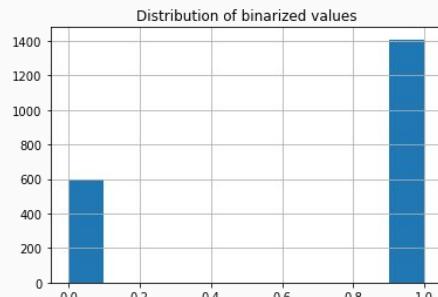
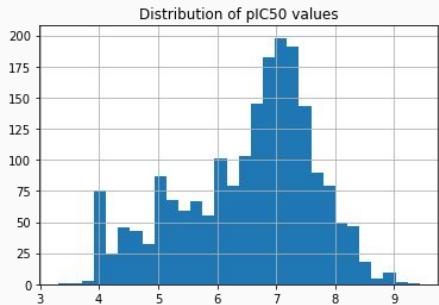
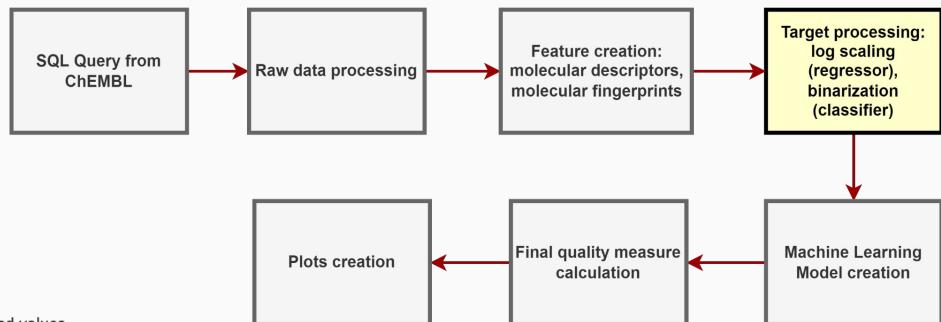




# Modeling of biological activity: workflow

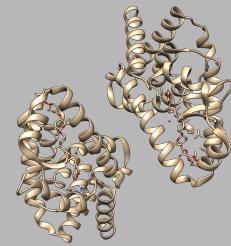


- The logarithmic transformation ( $\text{IC}_{50} \rightarrow \text{pIC}_{50}$ ), usually done for regressors
- The binarization with certain threshold (here 1000 nM), classifiers

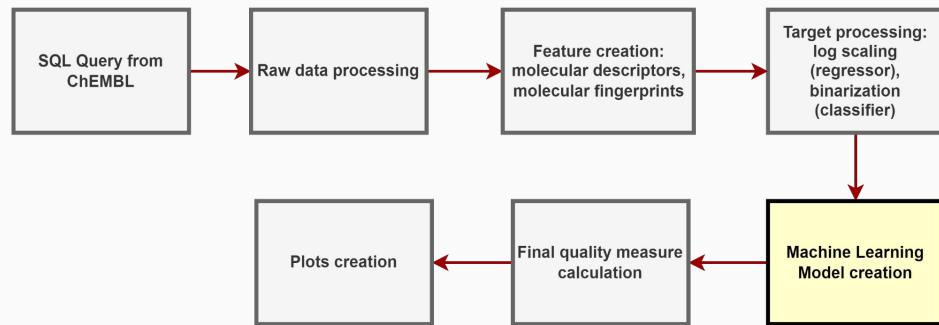




# Modeling of biological activity: workflow

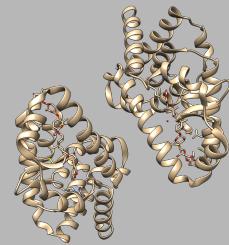


- The methodologies:
  - XGBoost
  - Multilayer perceptron
  - Random forest
  - Bagging
  - Support vector machines
  - Ridge
- Hyperparameters optimization with hyperopt
- Resulting models stored in a serialized form
- The serialized model is ready-to-deployment

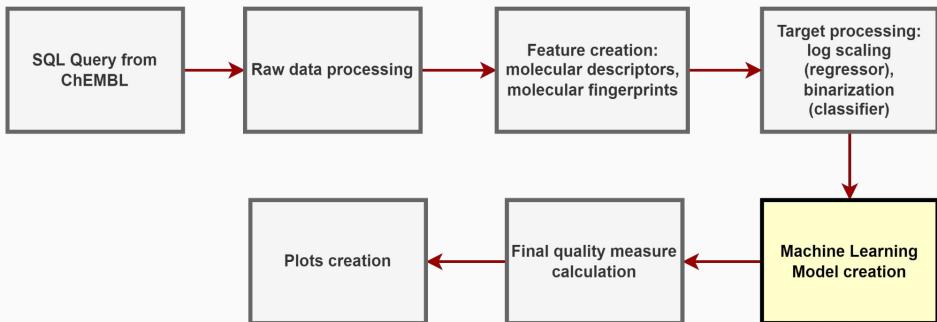




# Modeling of biological activity: workflow

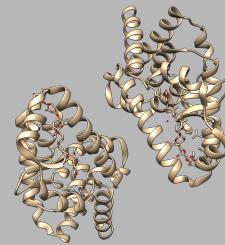


```
pipelines > xgboost > pipeline_configuration_xgboost_classification_Morgan_1024_pca_256_QED.json > ...  
1 {  
2     "feature_transform":  
3     {  
4         "Calculate_molecular_features":  
5         {  
6             "proceed": "yes",  
7             "calculate_fps": "yes",  
8             "calculate_descriptors": "yes",  
9             "fingerprint_size": 1024,  
10            "fp_type": "morgan",  
11            "descriptors_types": ["QED", "Lipinski"],  
12            "meta_descriptors_types": ["QED"],  
13            "label": "Molecular features"  
14        },  
15        "Scaling":  
16        {  
17            "proceed": "no",  
18            "scaling_features": ["SlogP", "SMR", "naRing", "nHAcc", "nHBDon", "nRot", "Pw", "TopoPSA", "QED"],  
19            "label": "scaling"  
20        },  
21        "PCA":  
22        {  
23            "proceed": "yes",  
24            "apply_to_md": "no",  
25            "apply_to_fp": "yes",  
26            "n_components": 256,  
27            "label": "PCA"  
28        }  
29    },  
30    "target_transform":  
31    {  
32        "Target_transform":  
33        {  
34            "proceed": "no"  
35        },  
36        "Target_binarization":  
37        {  
38            "proceed": "yes",  
39            "threshold": 1000,  
40            "label": "Target Binarization"  
41        },  
42        "Outlier_removal":  
43        {  
44            "proceed": "no",  
45            "factor": 1.5,  
46            "label": "Outlier removal"  
47        }  
48    }  
49}
```

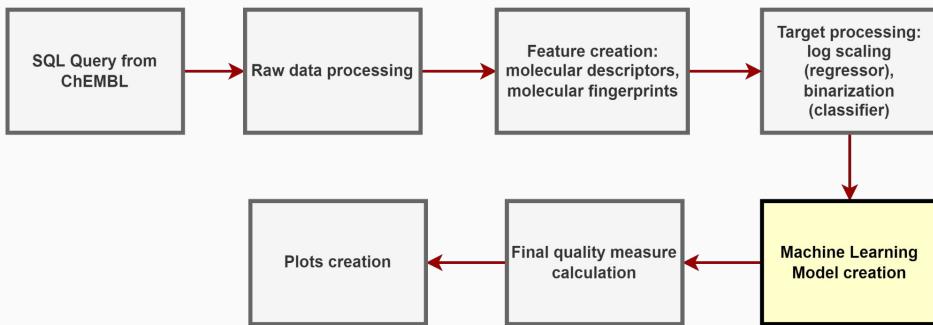




# Modeling of biological activity: workflow

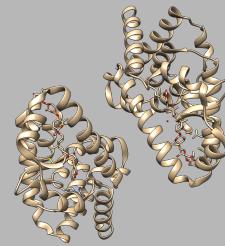


```
ng_configurations > mean_100 > bagging_classification_Morgan_1024_pca_512_QED.json > ...
{
    "prediction_methodology": "M_Bagging",
    "prediction_type": "classification",
    "GPU_support": true,
    "training_aux_data": {
        "run_name": "test",
        "n_outer": 1,
        "n_cv": 5,
        "development": false,
        "goal_function": "precision",
        "goal function_multiplier": -1.0,
        "threshold": 0.5,
        "track_model": false,
        "experiment": null,
        "comment": null,
        "max_evals": 200,
        "pipeline_file": "bagging/pipeline_configuration_bagging_classification_Morgan_1024_pca_512_QED.json"
    },
    "data_preparation": {
        "data_file": "ror_gamma_homo_sapiens_chembl_3.0_IC50.csv",
        "max_level_activity": 15,
        "std_threshold": 100,
        "strategy": "mean"
    },
    "model_storage": {
        "resulting_model": "model_bagging_classifier_Morgan_1024_pca_512_QED_mean_100.model"
    }
}
```

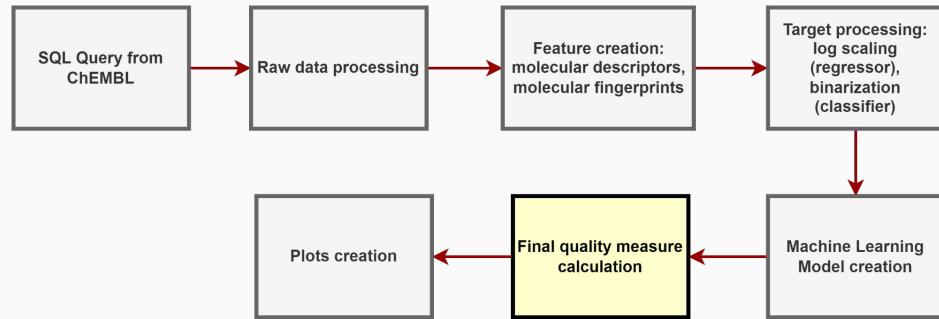




# Modeling of biological activity: workflow

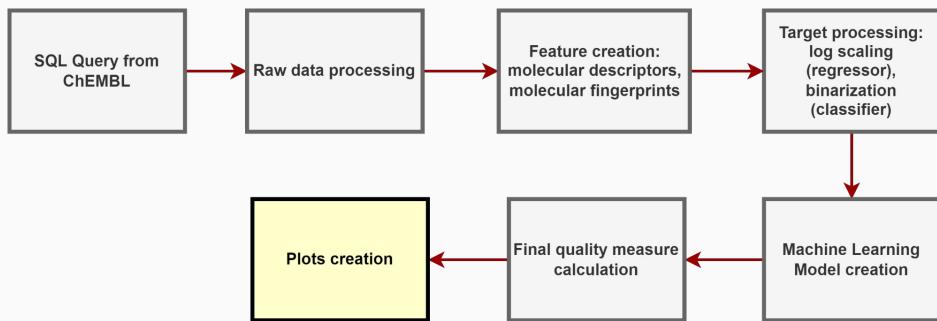
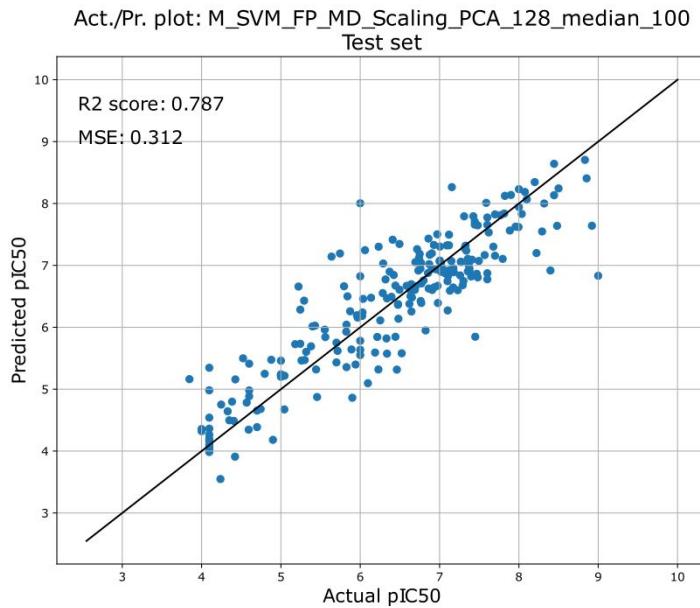
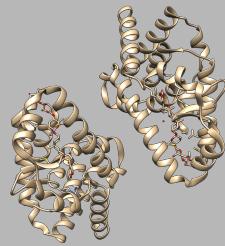


- The regressors
  - R<sup>2</sup> score
  - Mean squared (log) error
  - Mean absolute error
  - Mean absolute percentage error
- The classifiers
  - Precision
  - Recall
  - Accuracy
  - F1-score
  - ROC-AUC
  - Average precision
  - Matthews coefficient



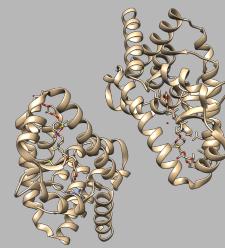


# Modeling of biological activity: workflow

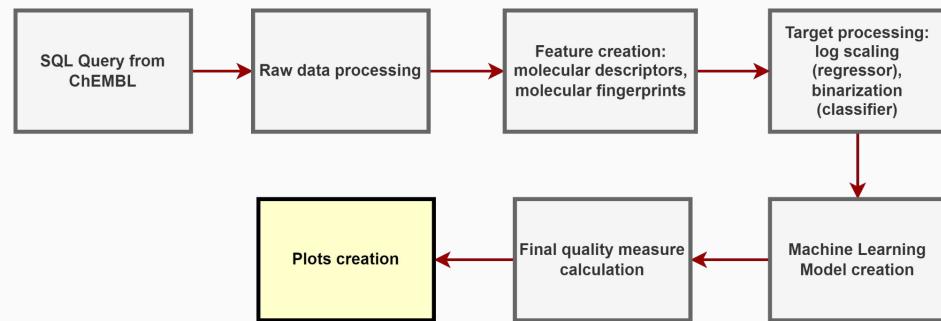
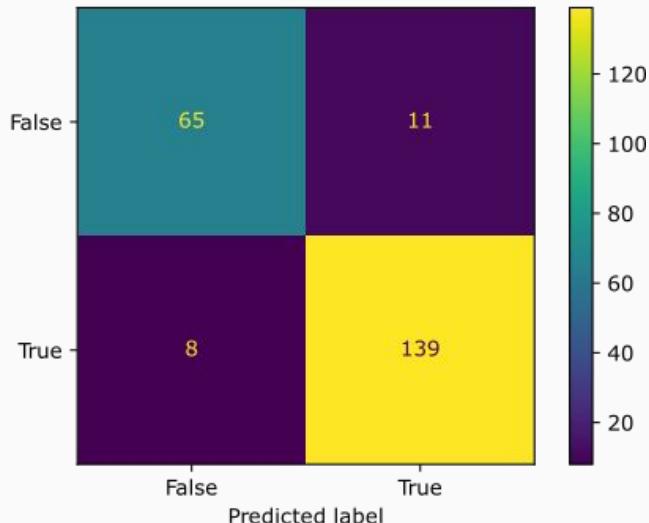




# Modeling of biological activity: workflow



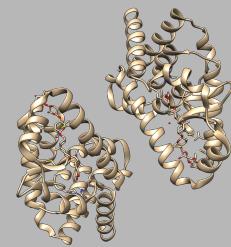
CM: M\_XGBoost\_FP\_NoScaling\_NoPCA\_mean\_100  
Test set





INSTYTUT  
BiolMed

# Modeling of biological activity: webapp



- Web application
  - Jinja templates
  - Dockerized setup
  - With load balancing (NGINX)
  - Used by experimental scientists

**RORgamma Bioactivity Prediction**

**Input information:**  
 Enter the smiles code:  
CC(O)c1ccc(C(=O)Nc2ccc([Si](C)(C)C)c2)N(C)C(=O)c2ccco2)cc1

**Choose the prediction model:**  
 activity\_prediction\_20201026\_ap\_pickle

**The predictive model description:**  
 The model optimized towards the best average precision (20210726)

**Calculate**

**Results:**

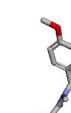
**The activity prediction**

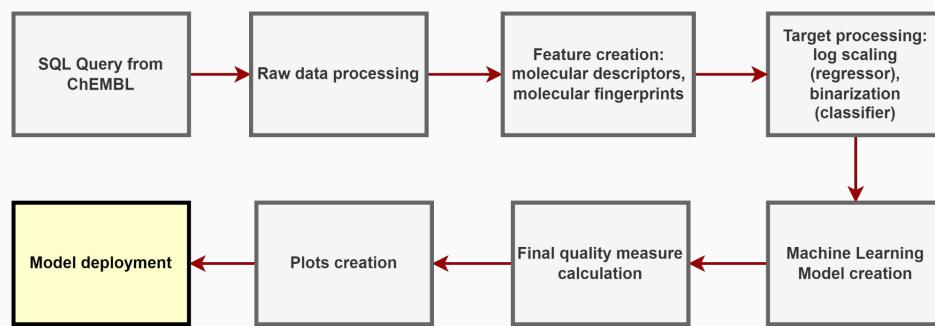
smiles code	activity
<chem>CC(O)c1ccc(C(=O)Nc2ccc([Si](C)(C)C)c2)N(C)C(=O)c2ccco2)cc1</chem>	0.973

**The selected molecular descriptors**

Molecular weight	436.182
Number of aromatic bonds	17
Number of atoms	59
Number of heavy atoms	31

**The graphical representation**

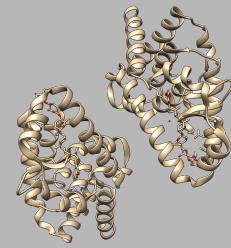




Manuscript in preparation...

# The ZINC20 screening

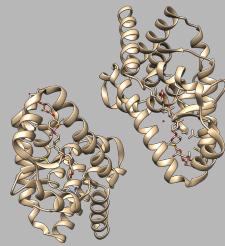
# The ZINC20 search



- The “active” ROR $\gamma$  ligands were taken as the reference species
  - “active”: IC50 < 1000 nM
  - 1704 molecules
  - The ZINC20 library was considered as a target
  - ca. 884 Ms of species
  - ca. 750 Ms of them are purchasable
  - The **chemfp** library for similarity search

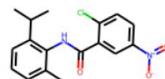
Molecular Weight (up to, Daltons)												Totals, by LogP
	200	250	300	325	350	375	400	425	450	500	>500	
-1	27,791	172,563	710,795	1,072,978	2,241,498	786,738	276,834	116,066	92,417	77,790	7,310	5,582,780
0	139,434	934,776	3,655,384	5,126,157	10,608,025	3,498,214	1,663,579	708,919	570,546	507,344	4,734	27,417,112
1	362,437	2,884,636	12,030,074	16,154,544	33,650,249	11,885,957	6,807,876	3,178,487	2,648,581	2,412,998	9,940	92,025,779
2	467,220	4,584,223	22,941,208	30,909,513	65,047,385	26,752,849	17,839,254	9,349,272	8,099,970	7,686,687	24,554	193,701,135
LogP (up to 2.5)	167,513	2,136,113	12,849,121	17,977,157	38,682,058	18,584,223	13,812,274	8,111,104	7,197,414	6,979,014	24,126	126,520,117
	90,548	1,570,772	11,037,383	16,282,627	34,831,558	19,940,391	16,037,132	10,339,743	9,362,233	9,118,717	37,422	128,648,526
3.5	36,748	929,872	7,920,574	12,499,662	27,380,104	18,703,024	16,485,194	11,784,160	10,774,472	10,693,411	58,791	117,257,012
4	9,017	369,565	4,332,131	6,472,808	10,487,856	13,034,155	14,329,253	11,683,208	10,891,665	11,003,975	86,862	82,699,695
4.5	993	86,613	1,814,492	3,457,942	6,367,225	8,853,064	10,320,054	9,945,353	9,486,869	9,825,079	117,980	60,275,064
5	150	13,393	536,018	1,405,708	3,168,584	4,995,850	6,471,525	7,025,034	6,976,742	7,325,833	144,297	38,063,134
5.5	39	1,097	22,854	103,521	376,905	927,395	1,670,856	2,195,160	2,588,702	3,052,048	787,762	11,706,339
i, by Weight	1,301,890	13,683,623	77,850,034	111,452,617	232,841,447	127,961,860	105,713,831	74,436,506	68,689,411	68,682,896	1,283,178	883,897,293
												Substances 1.9K Tranches

# The ZINC20 search

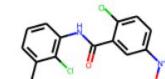


- The **chemfp** library was used to calculate the Morgan fingerprints
- The potential of entire content of the ZINC20 library was used
- The Tanimoto similarity coefficients were calculated between ZINC20 and “active” ROR $\gamma$  antagonists
- The similarity threshold: 0.7
- After removal of training set duplicates: 1673 potential ROR $\gamma$  antagonists
- By considering the commercial issues: three of them were initially selected to the experimental verification

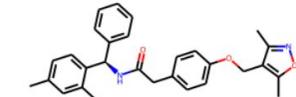
## Found species:



Classifier: 0.95  
Regressor: 922 nM

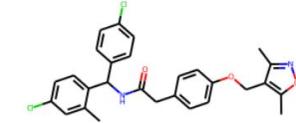
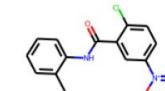
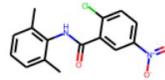


Classifier: 0.93  
Regressor: 889 nM

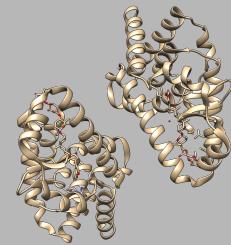


Classifier: 0.91  
Regressor: 37 nM

## The most similar analogues (from ChEMBL)

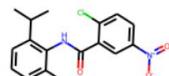


# The ZINC20 search

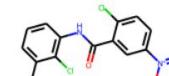


- The **chemfp** library was used to calculate the Morgan fingerprints
- The potential of entire content of the ZINC20 library was used
- The Tanimoto similarity coefficients were calculated between ZINC20 and “active” ROR $\gamma$  antagonists
- The similarity threshold: 0.7
- After removal of training set duplicates: 1673 potential ROR $\gamma$  antagonists
- By considering the commercial issues: three of them were initially selected to the experimental verification

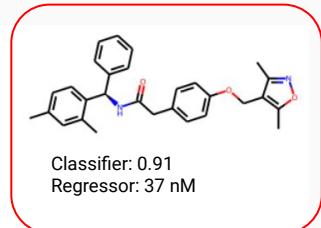
## Found species:



Classifier: 0.95  
Regressor: 922 nM

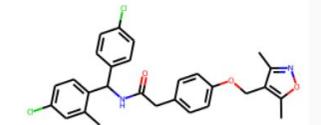
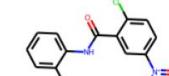
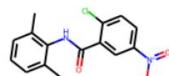


Classifier: 0.93  
Regressor: 889 nM



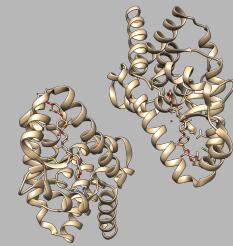
Classifier: 0.91  
Regressor: 37 nM

## The most similar analogues (from ChEMBL)



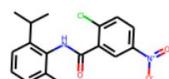
Already known potent ROR $\gamma$  antagonist  
not included in the ChEMBL database  
Measured IC50: 30 nM

# The ZINC20 search

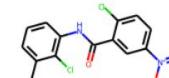


- The **chemfp** library was used to calculate the Morgan fingerprints
- The potential of entire content of the ZINC20 library was used
- The Tanimoto similarity coefficients were calculated between ZINC20 and “active” ROR $\gamma$  antagonists
- The similarity threshold: 0.7
- After removal of training set duplicates: 1673 potential ROR $\gamma$  antagonists
- By considering the commercial issues: three of them were initially selected to the experimental verification

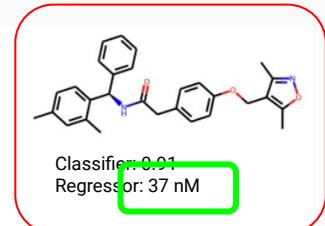
## Found species:



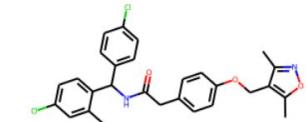
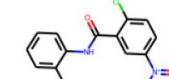
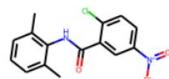
Classifier: 0.95  
Regressor: 922 nM



Classifier: 0.93  
Regressor: 889 nM



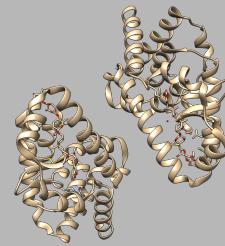
## The most similar analogues (from ChEMBL)



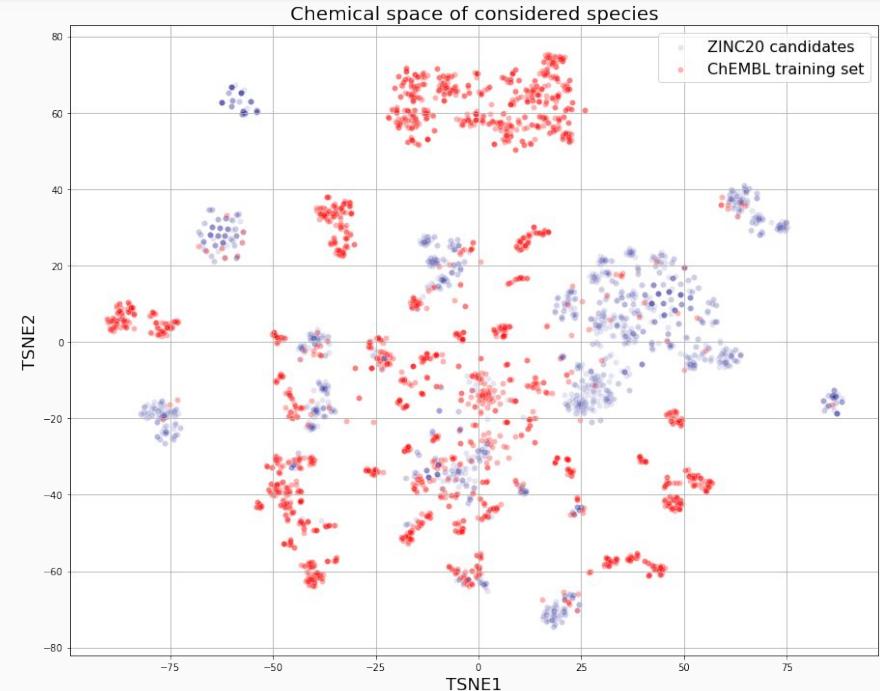
Already known potent ROR $\gamma$  antagonist  
not included in the ChEMBL database  
Measured IC50: **30 nM**



# The chemical space analysis

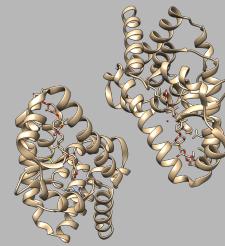


- The 2D representation of the chemical space
- The transformation:
  - 1024 Morgan fingerprint
  - PCA transformation into 1024 space of principal components
  - t-SNE into 2D space
- The area of applicability of the model
- Applied for the training set as well as for the set of potential inhibitors (ZINC20)
- Are the predictions trustworthy?
- The ML has the regression rather than extrapolation potential

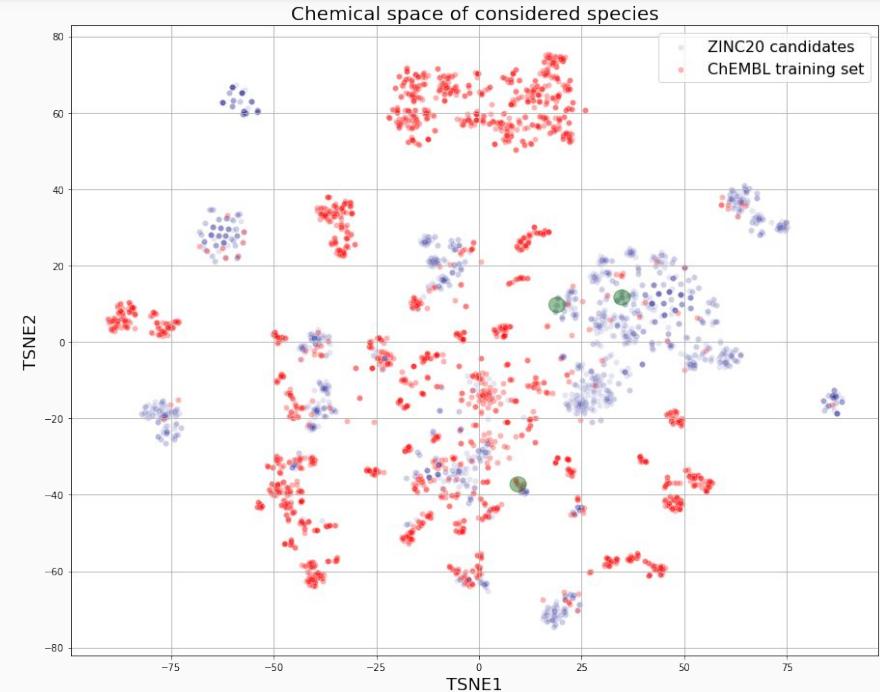




# The chemical space analysis



- The 2D representation of the chemical space
- The transformation:
  - 1024 Morgan fingerprint
  - PCA transformation into 1024 space of principal components
  - t-SNE into 2D space
- The area of applicability of the model
- Applied for the training set as well as for the set of potential inhibitors (ZINC20)
- Are the predictions trustworthy?
- The ML has the regression rather than extrapolation potential



# Summary

- Introduction
- Nuclear receptors
- How to measure the biological activity?
- The QSAR workflow
- ZINC20 screening
- Summary

# Acknowledgements

- **Thank you for your attention!**
- Marcin Ratajewski (for experiments and funds)
- Damian Nowak (for support with coding)
- Greg Landrum (for RDKit - a great library)
- Andrew Dalke (for chemfp)
- Polish National Science Centre Grant  
2019/33/B/NZ7/00795