



# GRADUATION THESIS

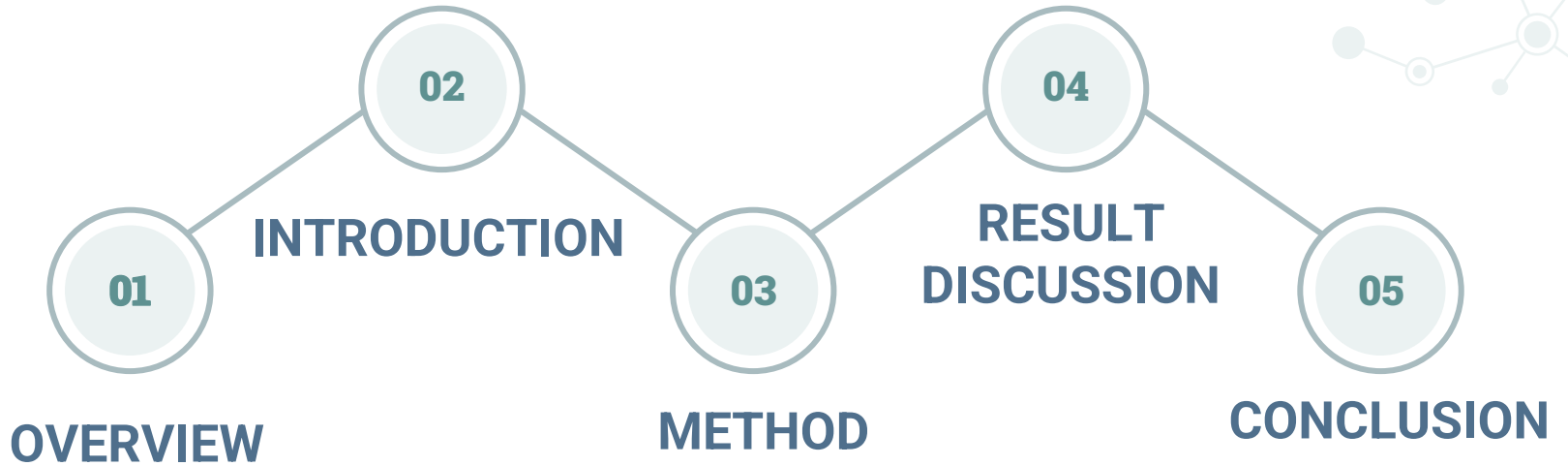
Department of Organic Chemistry

## *IN SILICO* MODELING FOR PREDICTION OF POTENTIAL HIV-1 INTEGRASE INHIBITORS

Presenter: Phan Tieu Long  
Supervisor: Assoc.Prof.Truong Ngoc Tuyen



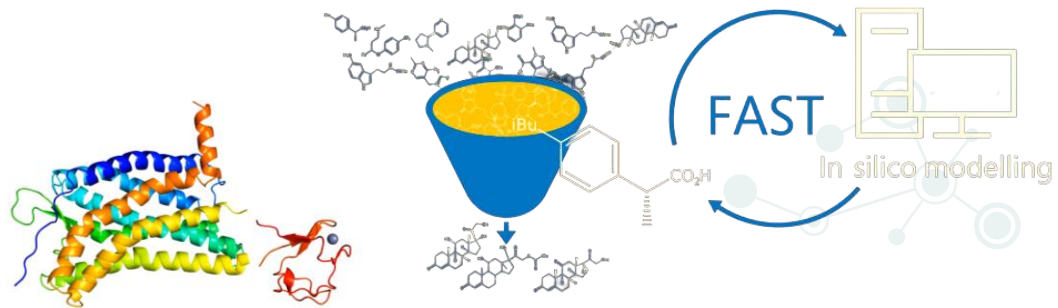
# CONTENT



# OVERVIEW



# OVERVIEW





# OVERVIEW



MOLECULAR DOCKING

PHAMACOPHORE

QSAR

VIRTUAL  
SCREENING





# OVERVIEW



01

Pharmacophore

02

QSAR  
classification

03

QSAR  
regression

04

Docking

05

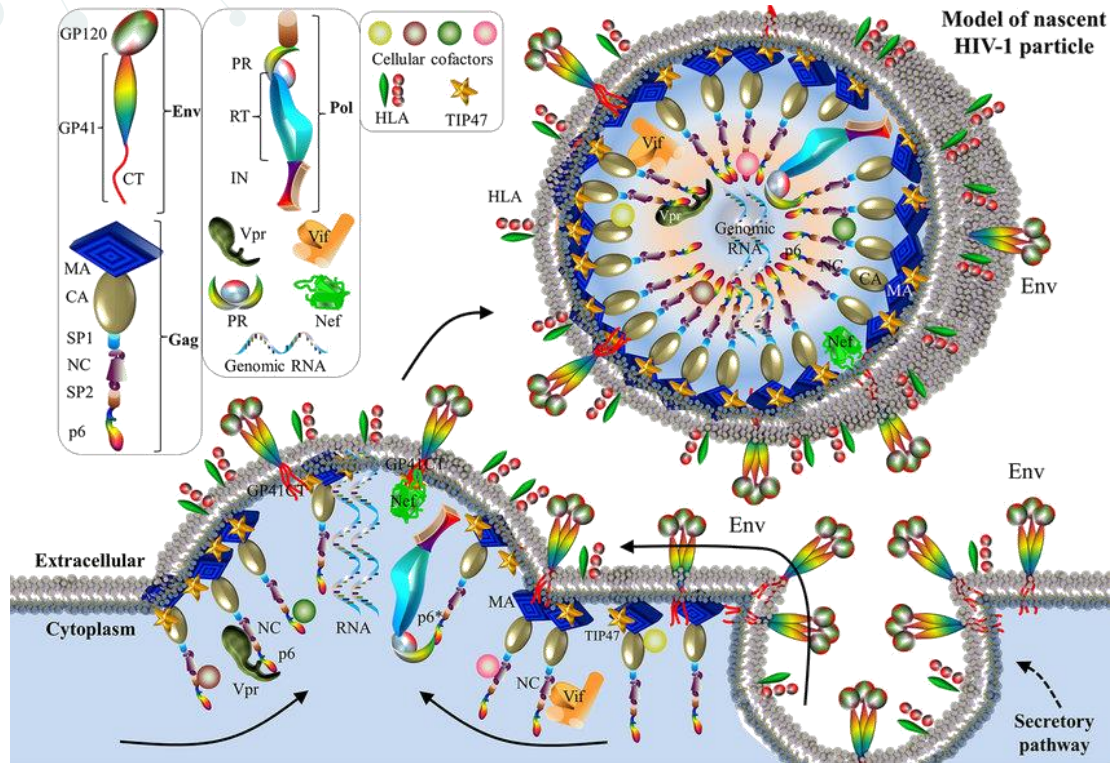
Screening





# INTRODUCTION

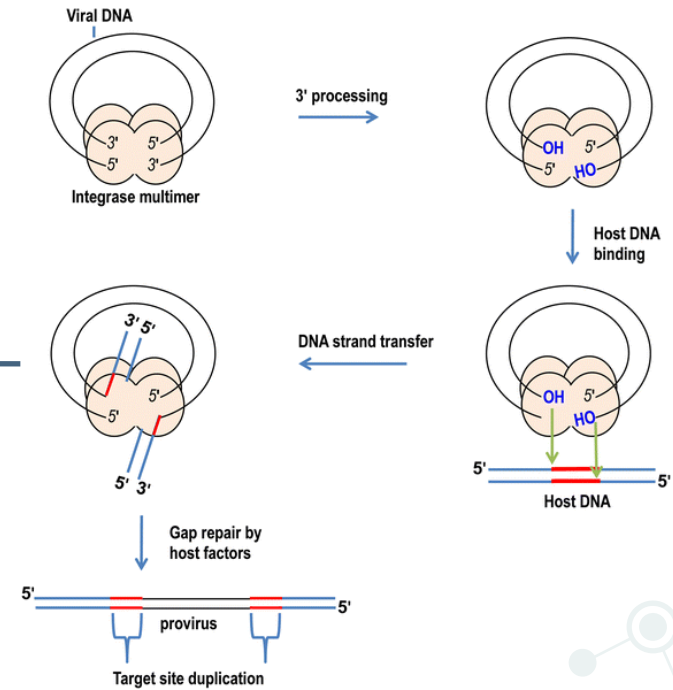
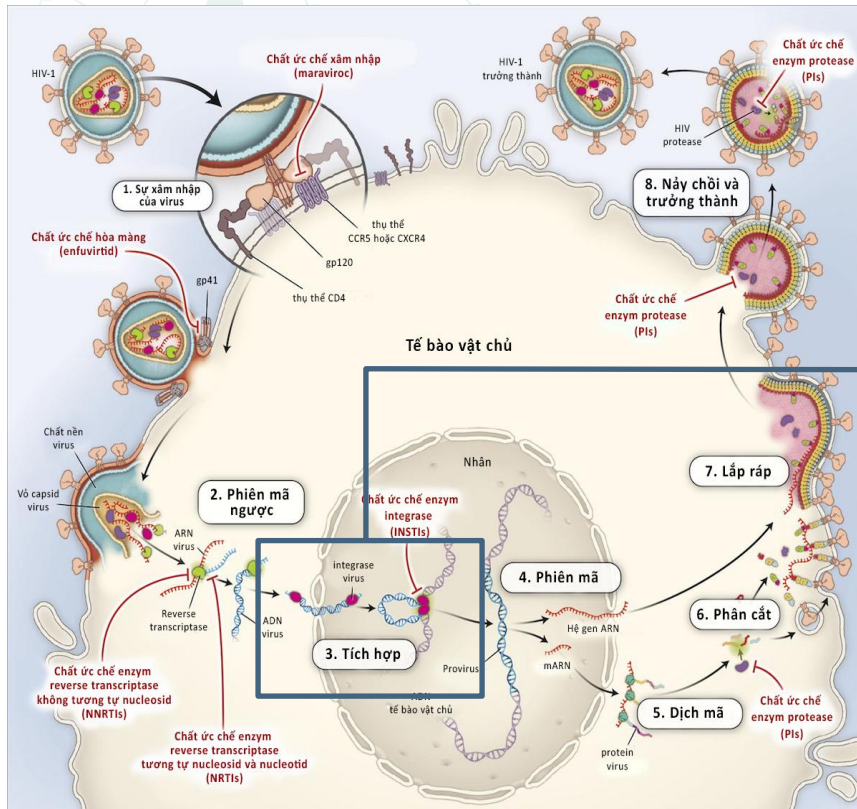
# INTRODUCTION



Li G, De Clercq E. HIV Genome-Wide Protein Associations: a Review of 30 Years of Research. *Microbiology and molecular biology reviews* : MMBR. 2016;80(3):679-731. doi:10.1128/mmr.00065-15



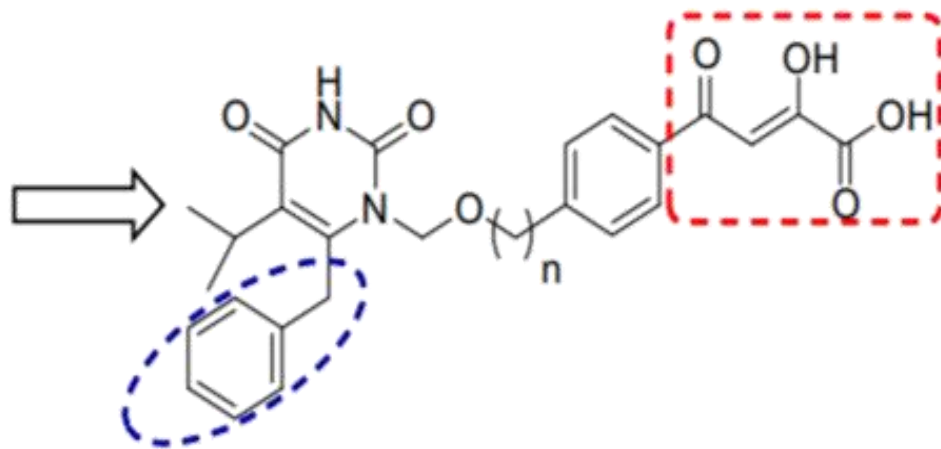
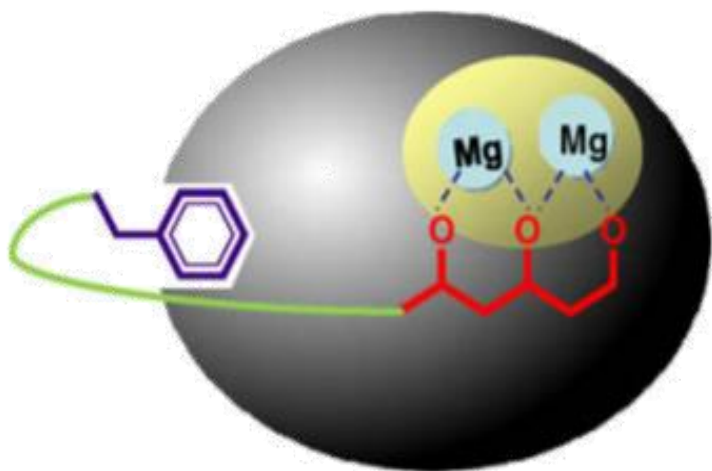
# OVERVIEW



## INTEGRATION MECHANISM

## OVERVIEW

## SAR OF INSTIS



Wang Z., Tang J., Salomon C. E. et al. (2010), "Pharmacophore and structure-activity relationships of integrase inhibition within a dual inhibitor scaffold of HIV reverse transcriptase and integrase", *Bioorganic & Medicinal Chemistry*. 18 (12), pp. 4202-4211

# OVERVIEW

# FDA APPROVED DRUGS

## Cabotegravir 2021



## Bictegravir 2018



## Dolutegravir 2013



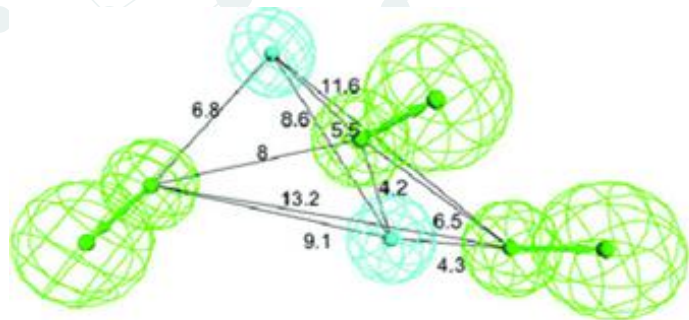
## Elvitegravir 2012



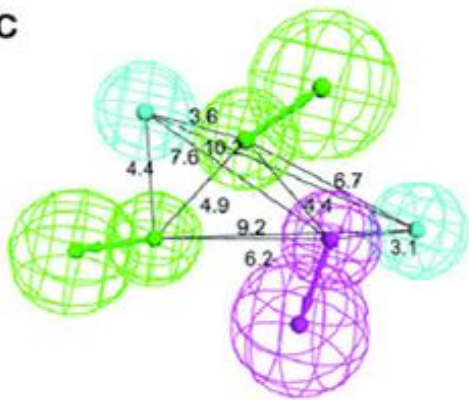
## Raltegravir 2007



# OVERVIEW

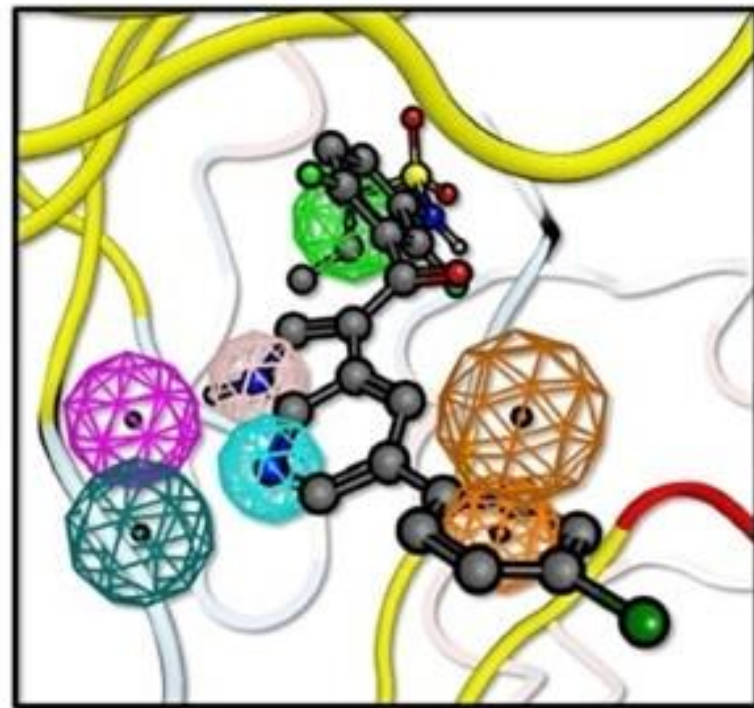


C



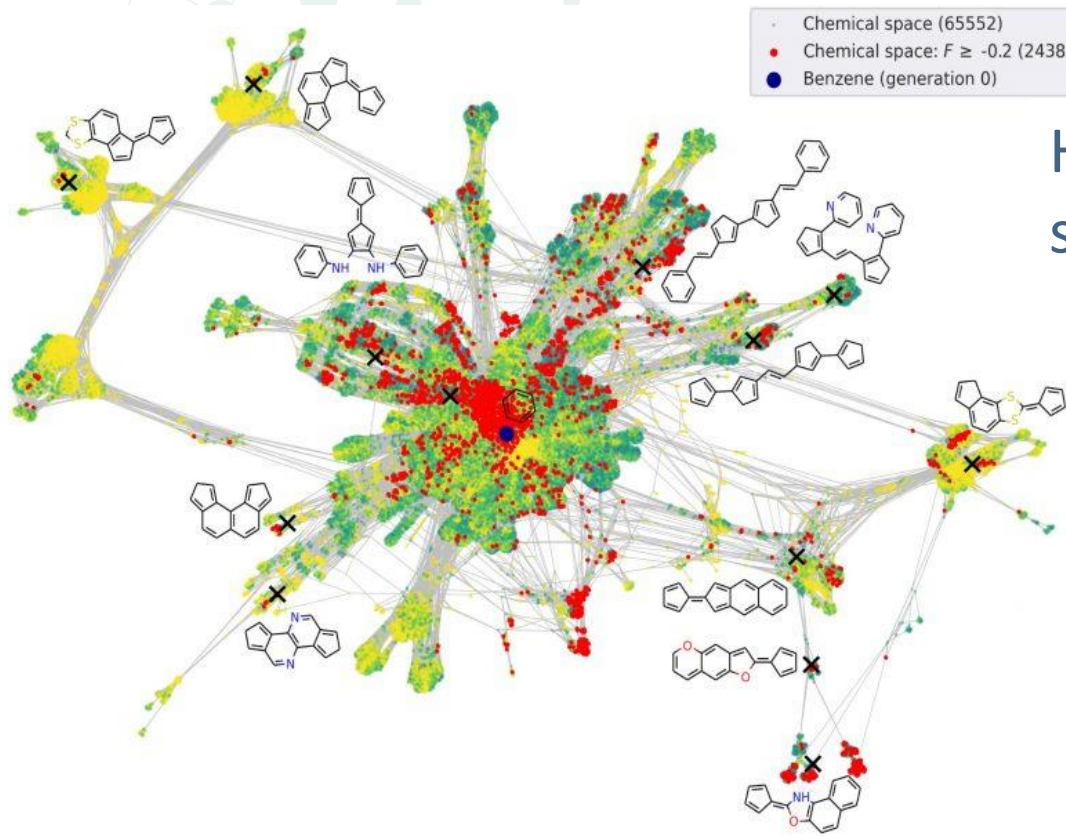
Ligand-based

# PHARMACOPHORE



Protein-based

## OVERVIEW



## PHARMACOPHORE

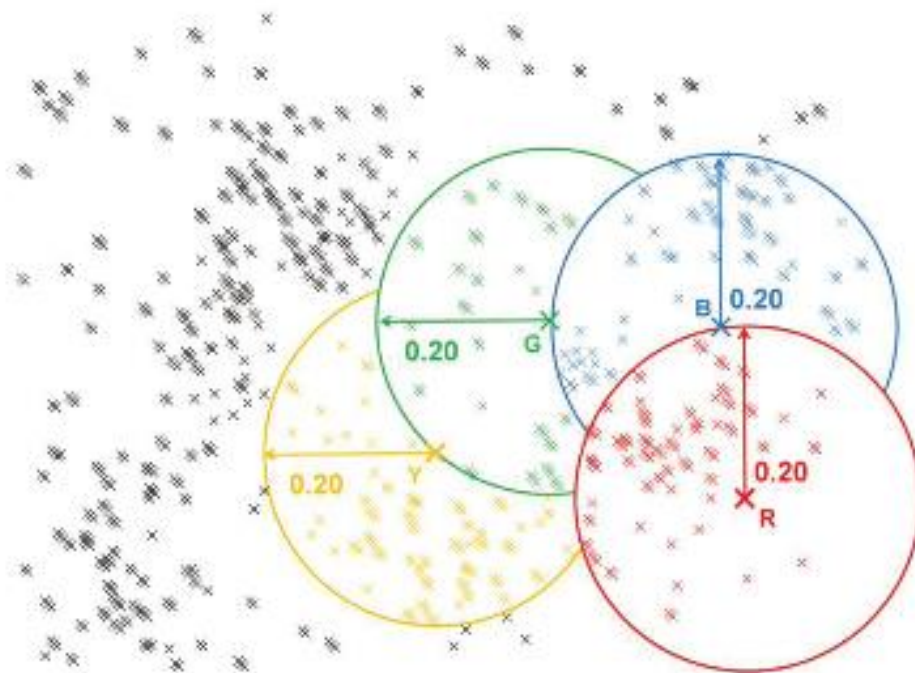
How can we select diverse subset for training model?



## OVERVIEW

# PHARMACOPHORE

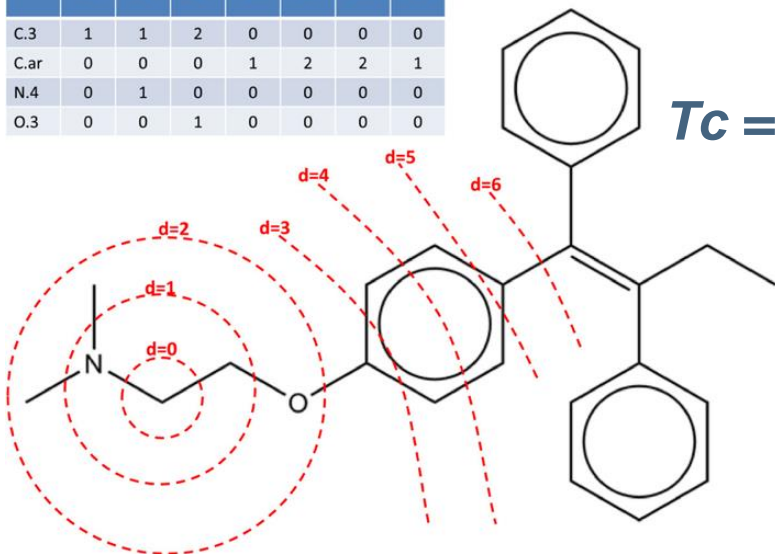
## Butina algorithm



- Calculate similarity matrix
- Select centroid
- "Single Pass" technique

# OVERVIEW

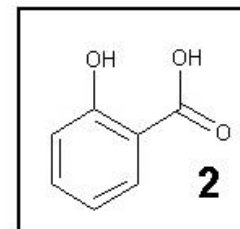
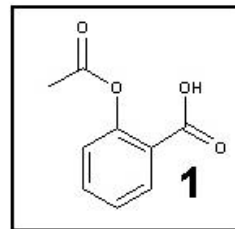
	d=0	d=1	d=2	d=3	d=4	d=5	d=6
C.3	1	1	2	0	0	0	0
C.ar	0	0	0	1	2	2	1
N.4	0	1	0	0	0	0	0
O.3	0	0	1	0	0	0	0



$$Tc = \frac{c}{a+b-c} = 0,6$$

# PHARMACOPHORE

## Molecular fingerprint



1	1	1	0	1	1	0	1	0
2	1	1	0	1	0	0	0	0

Rogers D, Hahn M. Extended-connectivity fingerprints. *Journal of chemical information and modeling*. 2010;50(5):742-54. doi:10.1021/ci100050t

Bender A, Glen RC. Molecular similarity: a key technique in molecular informatics. *Organic & biomolecular chemistry*. 2004;2(22):3204-18. doi:10.1039/b409813g

## OVERVIEW

# PHARMACOPHORE Model optimization

### Common Hits Approach

RPM1

**Hit-list :**  
Molecule A  
Molecule C  
Molecule E

RPM2

**Hit-list :**  
Molecule C  
Molecule D

RPM3

**Hit-list :**  
Molecule A  
Molecule C  
Molecule E

**Hit-list : Score :**

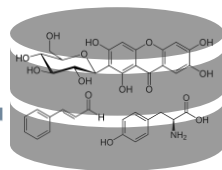
Molecule C	<b>3</b>
Molecule A	<b>2</b>
Molecule E	<b>2</b>
Molecule D	<b>1</b>

- Change query threshold
- Change overlap ratio
- Add exclusion volume
- Change ligand's shape
- **Common Hits Approach**



# OVERVIEW

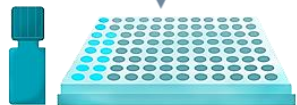
Database



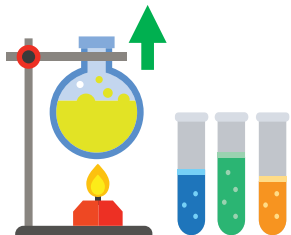
# QSAR MODEL

**Lead compound**

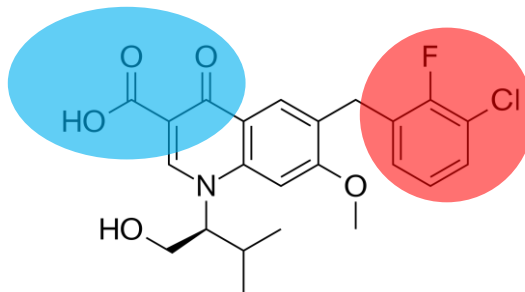
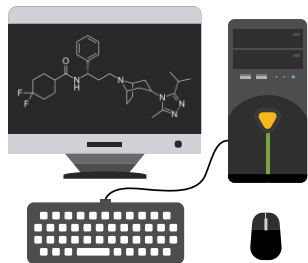
*in vitro*



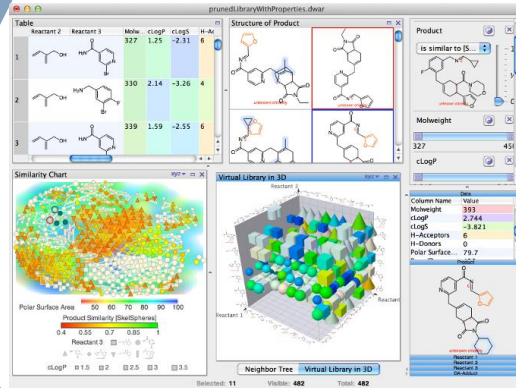
Synthesis



Design  
*in silico*



**QSAR**



# OVERVIEW

## mordred-descriptor/ documentation

<http://mordred-descriptor.github.io/documentation/master>

1 Contributor 0 Issues 0 Stars 0 Forks



and Machine Learning

## phi-grib/PaDEL- descriptor-ws

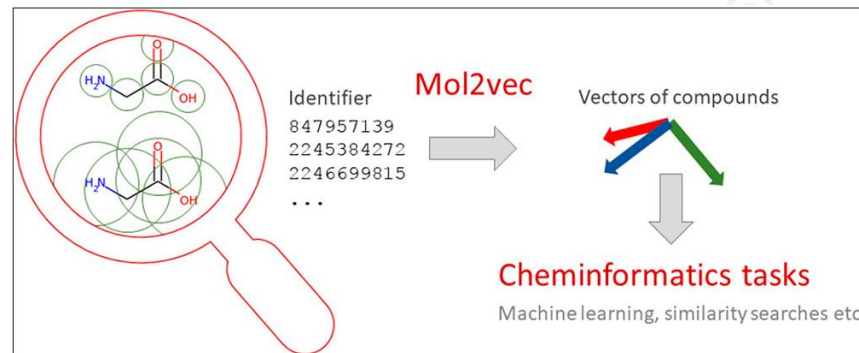
PaDEL ws descriptors engine

1 Contributor 1 Issue 11 Stars 6 Forks



# QSAR MODEL

## Molecular representation

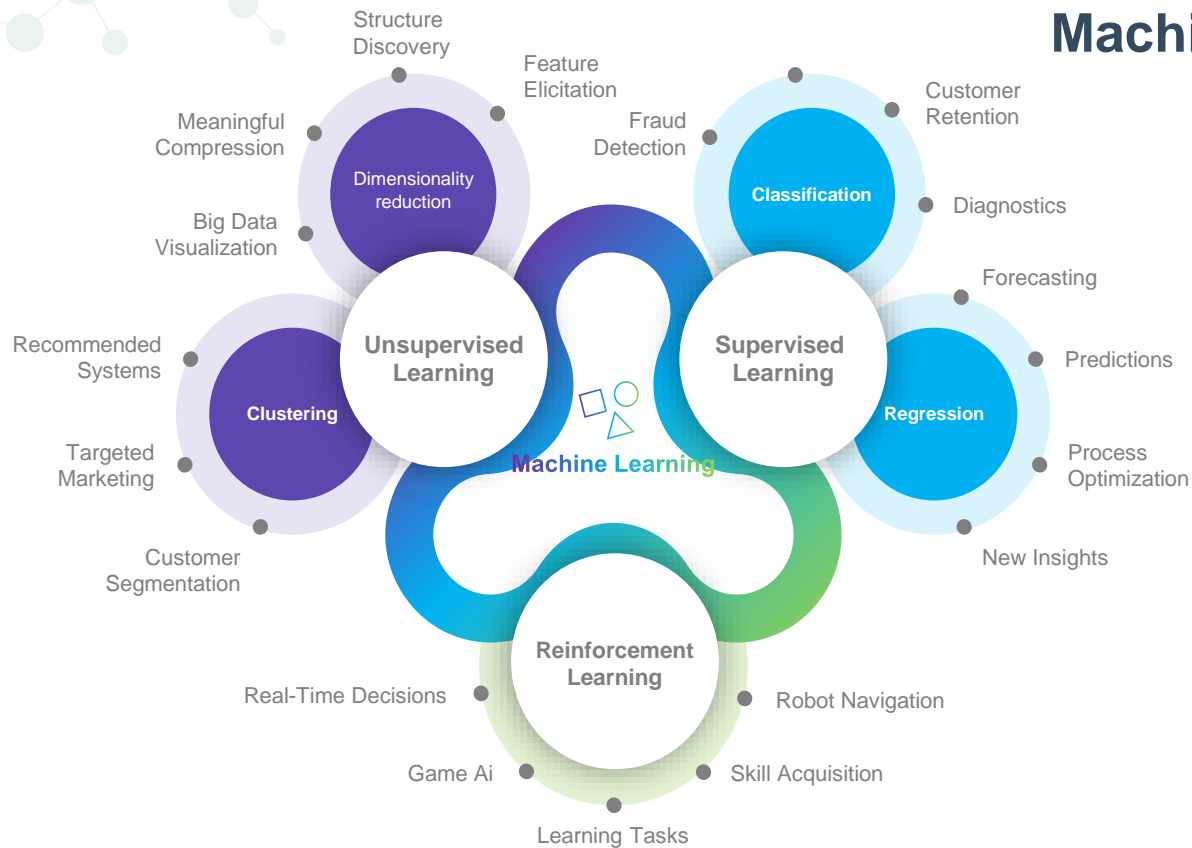


Open-Source Cheminformatics  
and Machine Learning

# OVERVIEW

# QSAR MODEL

## Machine Learning



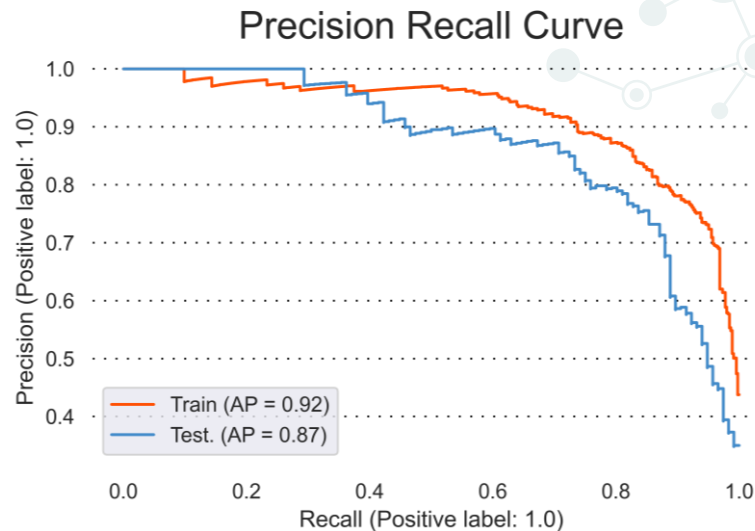
# OVERVIEW

# QSAR MODEL

## Evaluation metric - Classification

		Predicted Class		
		Positive	Negative	
Actual Class	Positive	True Positive (TP)	False Negative (FN) Type II Error	<b>Sensitivity</b> $\frac{TP}{(TP + FN)}$
	Negative	False Positive (FP) Type I Error	True Negative (TN)	<b>Specificity</b> $\frac{TN}{(TN + FP)}$
		<b>Precision</b> $\frac{TP}{(TP + FP)}$	<b>Negative Predictive Value</b> $\frac{TN}{(TN + FN)}$	<b>Accuracy</b> $\frac{TP + TN}{(TP + TN + FP + FN)}$

$$F_1 = \frac{2(\text{Precision} \times \text{Recall})}{\text{Precision} + \text{Recall}}$$



$$AP = \sum_n (R_n - R_{n-1}) P_n$$

## OVERVIEW

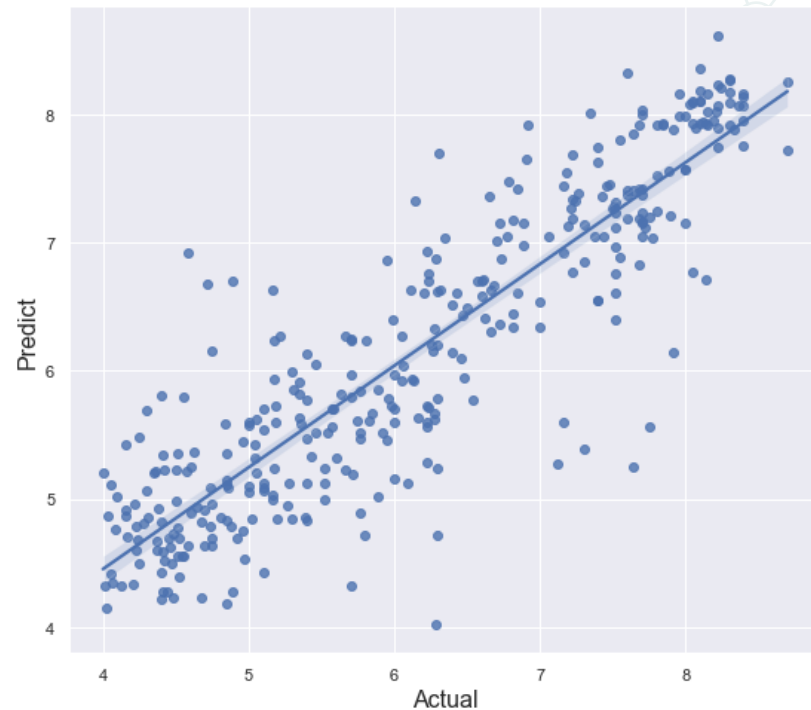
$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y}_i)^2}$$

$$\text{RMSE} = \frac{1}{n} \times \sqrt{\|\mathbf{y} - \hat{\mathbf{y}}\|_2^2}$$

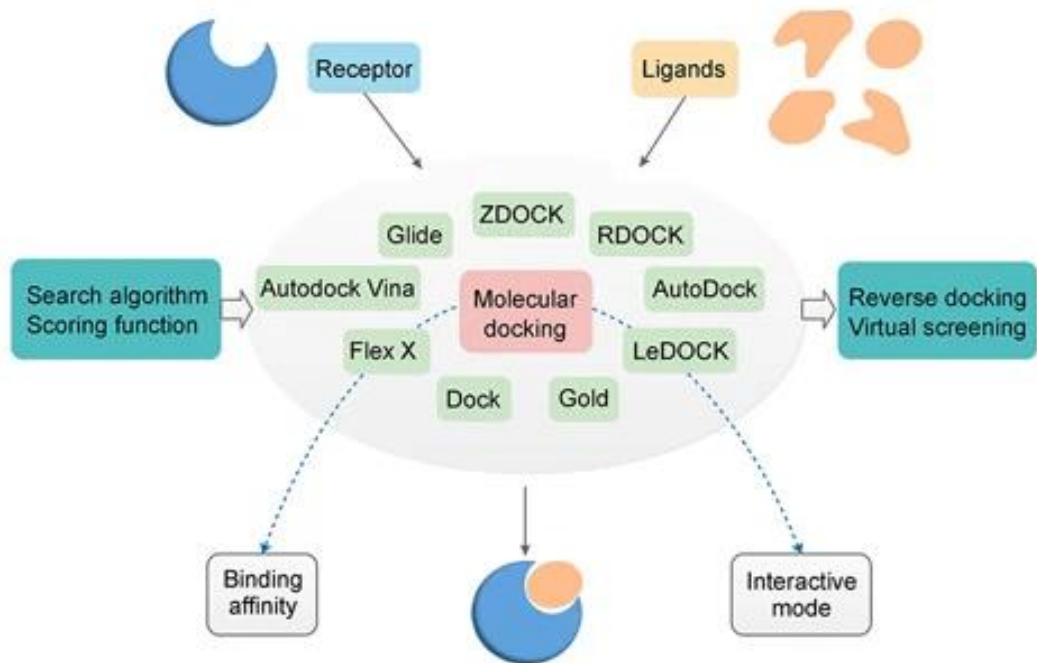
$$\text{MAE} = \frac{1}{n} \times \|\mathbf{y} - \hat{\mathbf{y}}\|_1$$

## QSAR MODEL

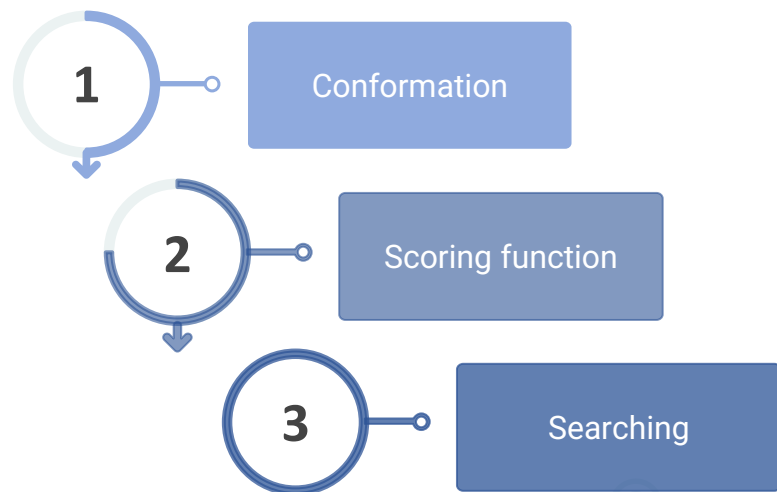
### Evaluation metric - Regression



# OVERVIEW



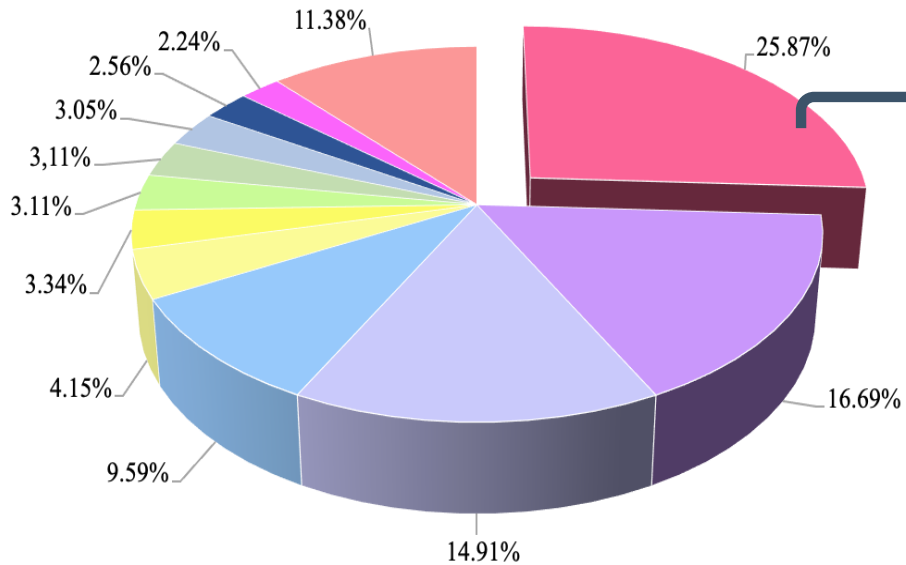
# MOLECULAR DOCKING



# OVERVIEW

# MOLECULAR DOCKING

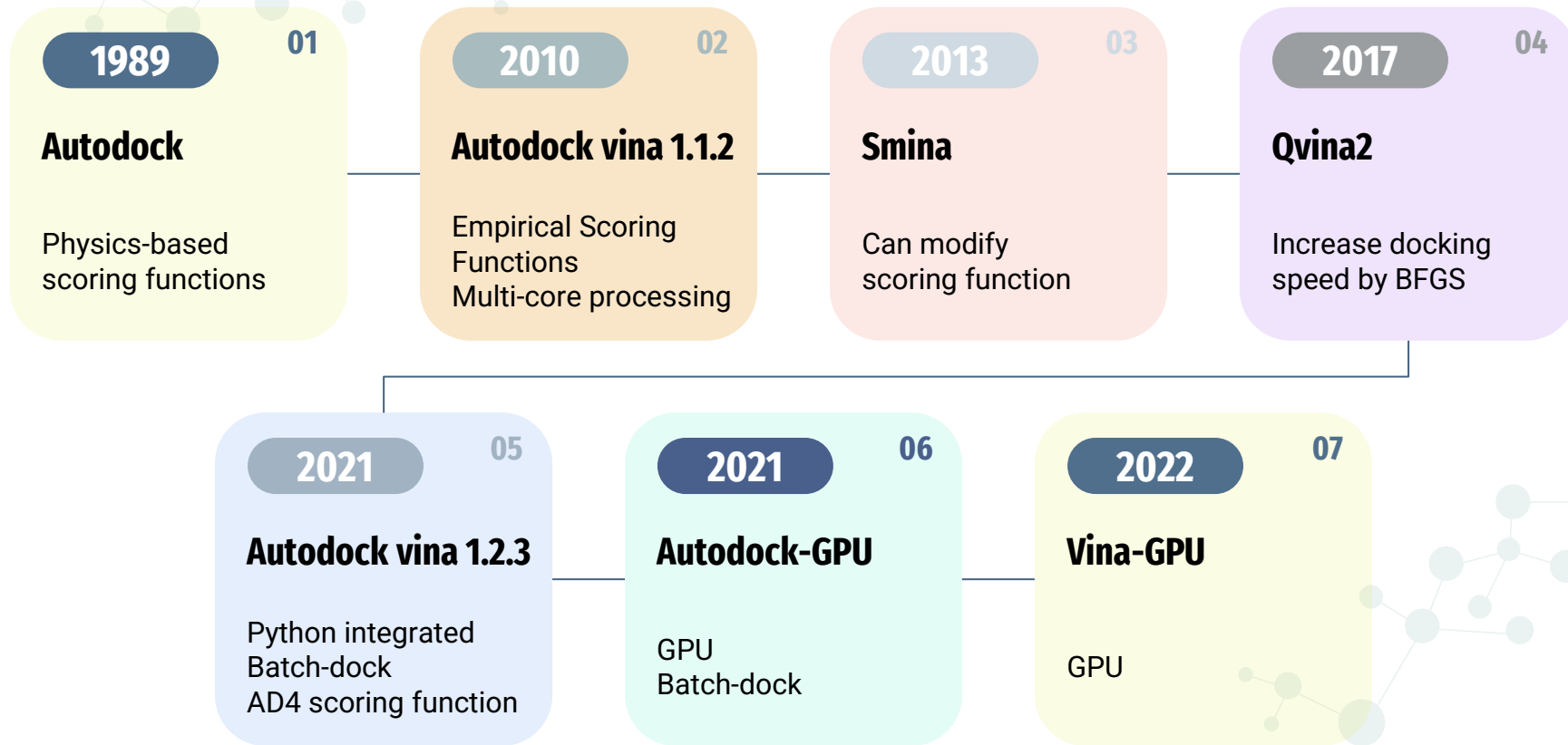
- Autodock
- Surflex-Dock
- HADDOCK
- GOLD
- FITTED
- ICM
- Glide
- Autodock Vina
- LigandFit
- FlexX, Flex-Ensemble
- MOE
- Others



Autodock had highest citation in ISI Web of Science (2005)

# OVERVIEW

# MOLECULAR DOCKING





## OVERVIEW

### Active compounds

Active reference for docking model

### Decoy

Inactive reference for docking model

### Metrics

AUC-ROC; G-mean, TPR

### Meaning

Performance of docking softwares

01

02

03

04

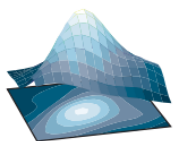
## MOLECULAR DOCKING

### Retrospective control



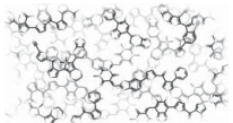
# Computational drug discovery: three schemes

## Functional space



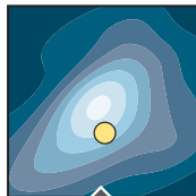
Desired properties (redox potential, solubility, toxicity)

## Chemical space

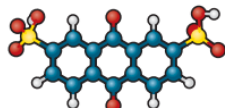


(Drug-like, photovoltaics, polymers, dyes)

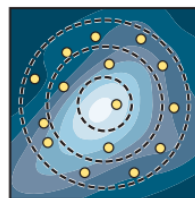
## Simulation



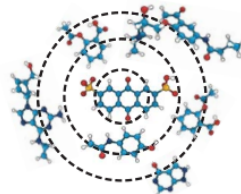
Experiment or simulation (Schrödinger equation)



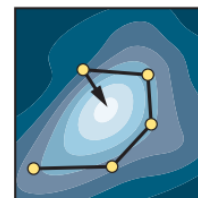
## Virtual screening



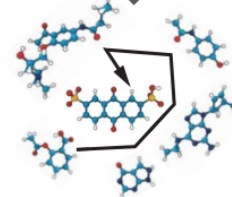
High-throughput virtual screening (e.g., with 3 filtering stages)



## De novo drug design



Optimization, evolutionary strategies, generative models (VAE, GAN, RL)



**DATA**

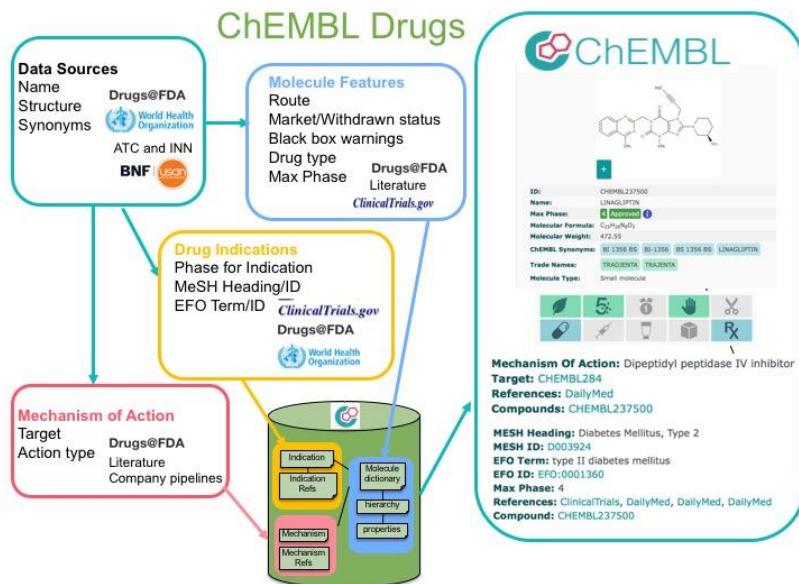


# DATA

## DATA FOR DEVELOPING MODELS

## DATA FOR SCREENING

3105 chất

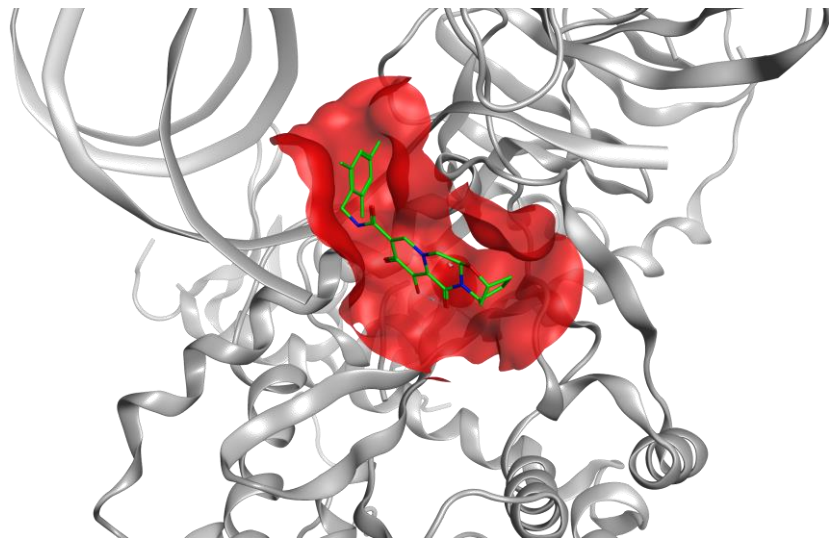


Internal database

2016 compounds

DATA

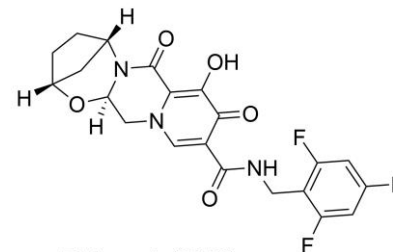
PROTEIN



Cryo-em

2,9 Å

PDB ID: 6PUW



Bictegravir (2018)

**METHOD**



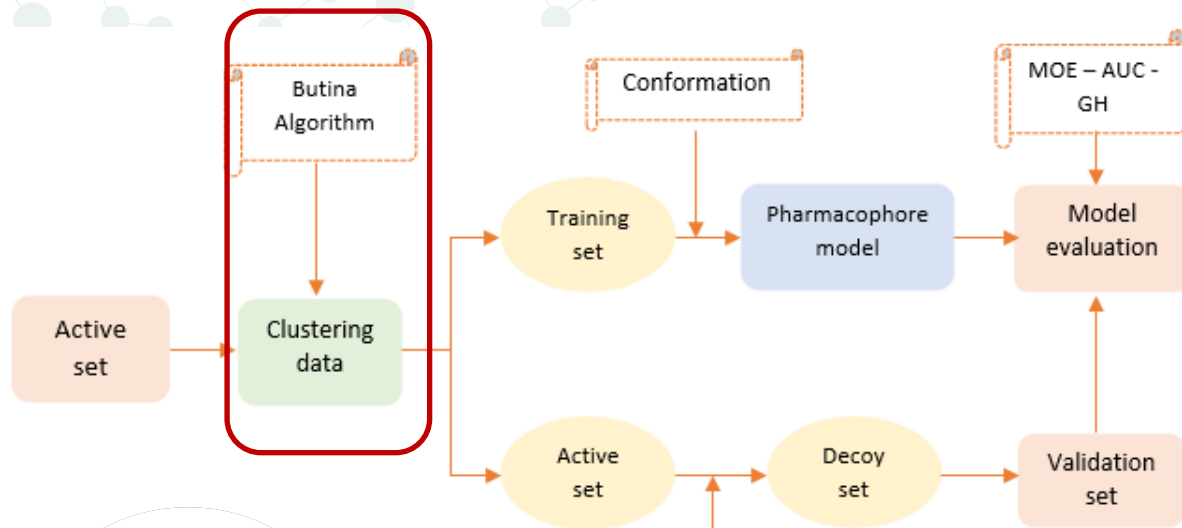
Method



01

PHARMACOPHORE

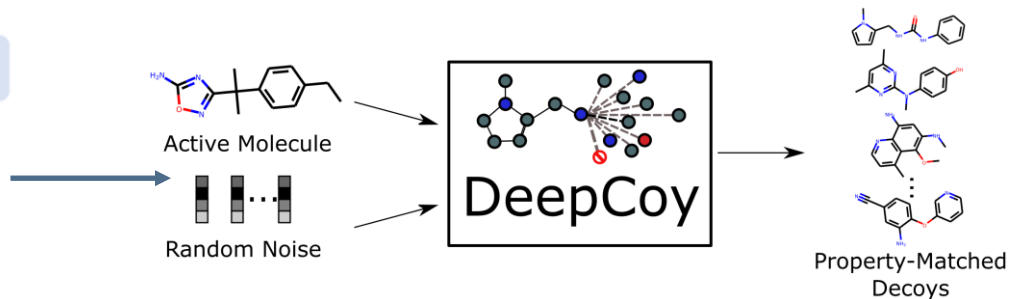
# METHOD



# PHARMACOPHORE



Chemical  
Computing  
Group

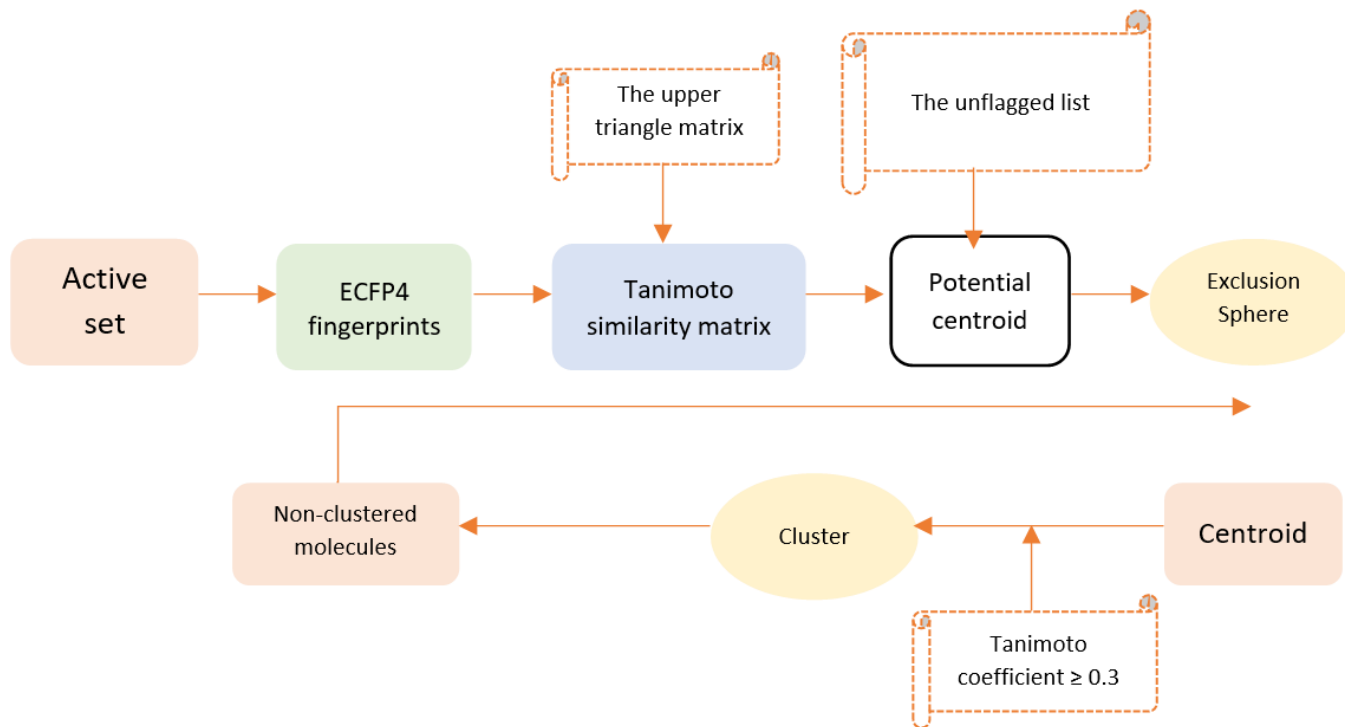




# METHOD

# PHARMACOPHORE

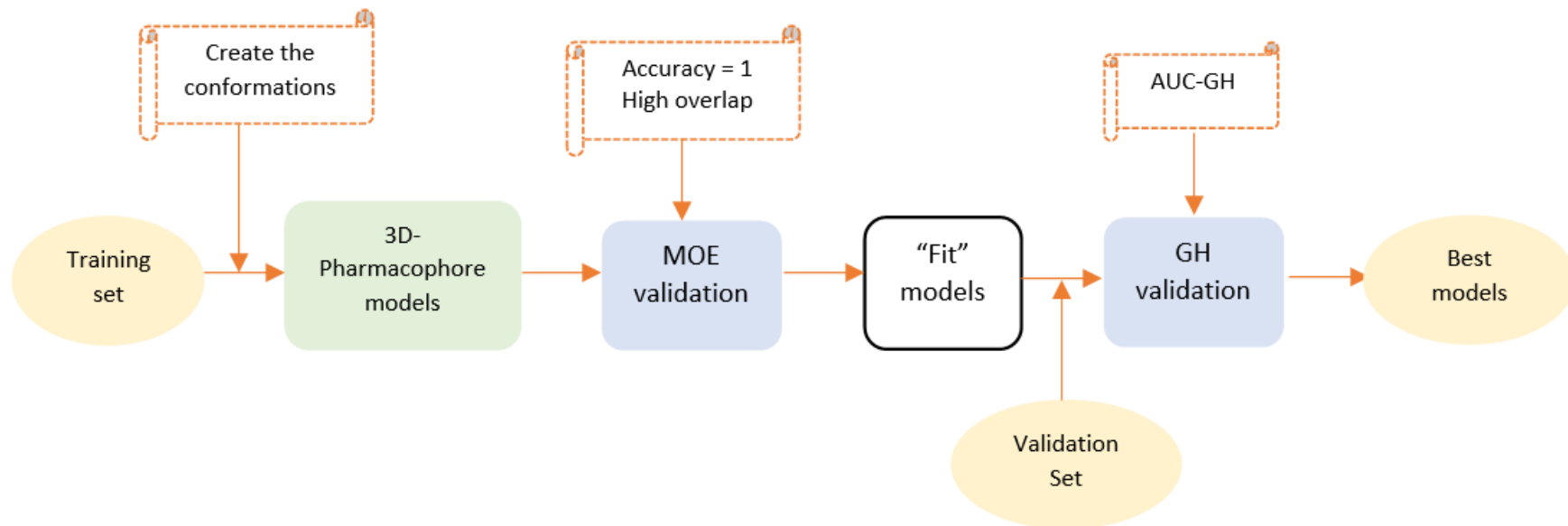
## Molecules clustering



# METHOD

# PHARMACOPHORE

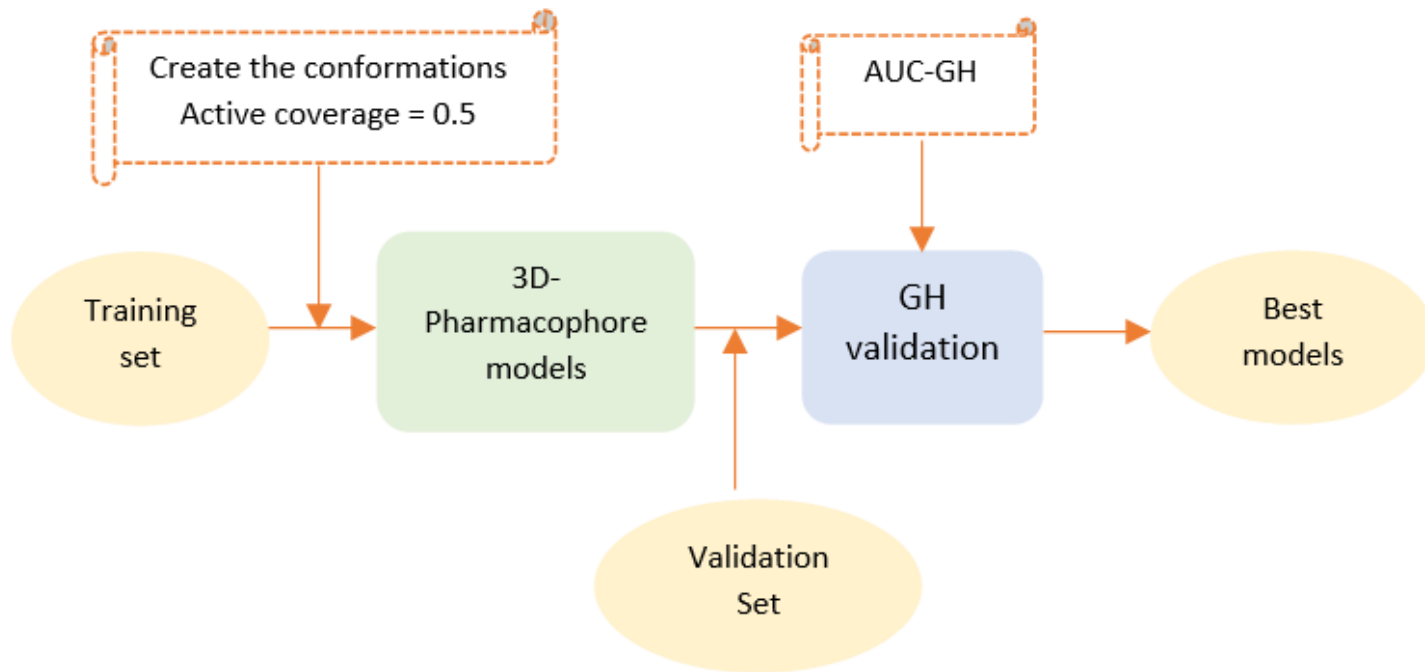
## Local Search



# METHOD

# PHARMACOPHORE

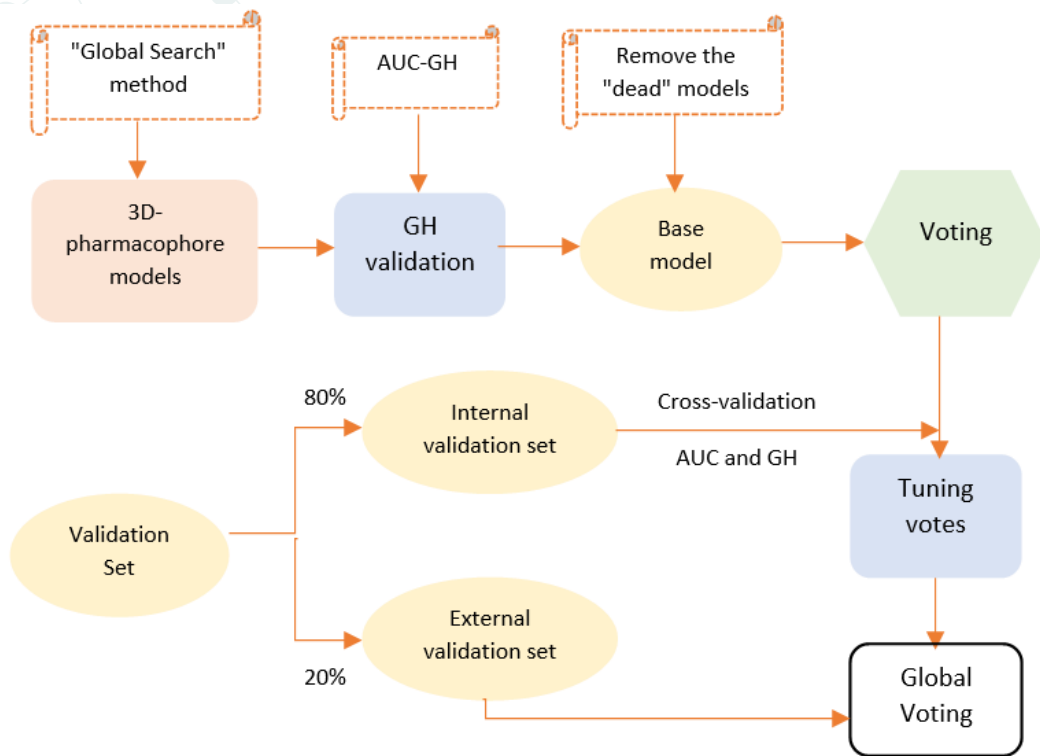
## Global Search



# METHOD

# PHARMACOPHORE

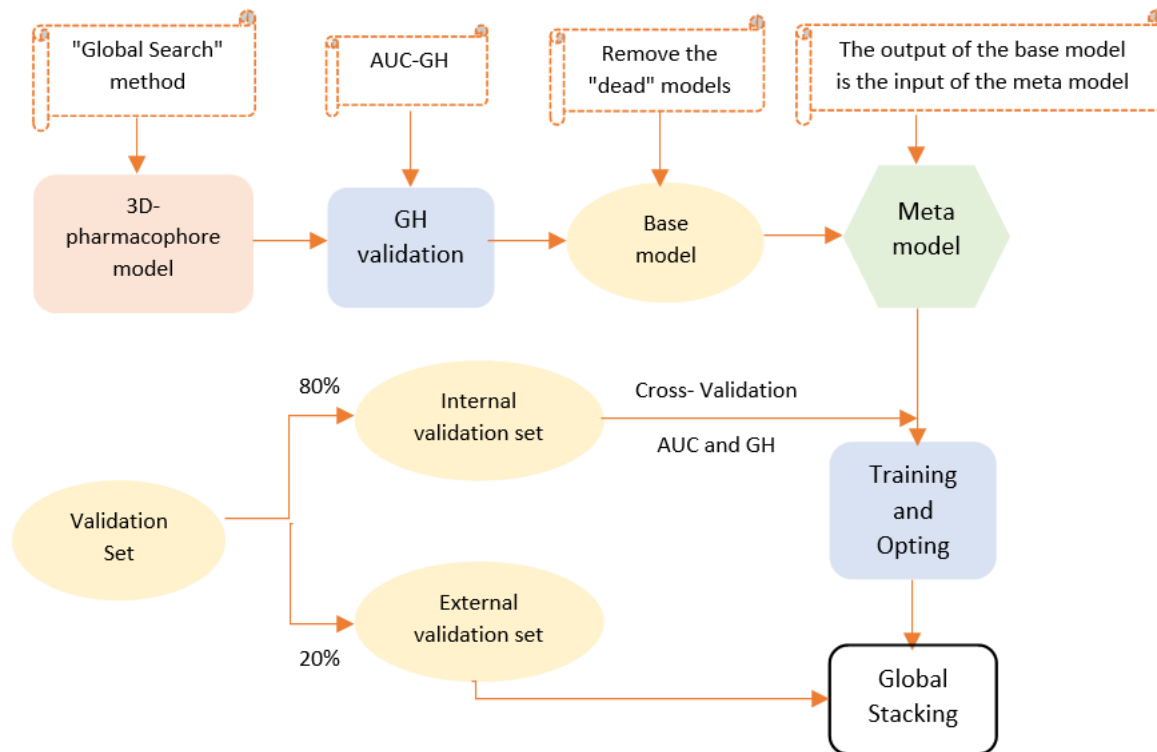
## Voting



# METHOD

# MÔ HÌNH PHARMACOPHORE

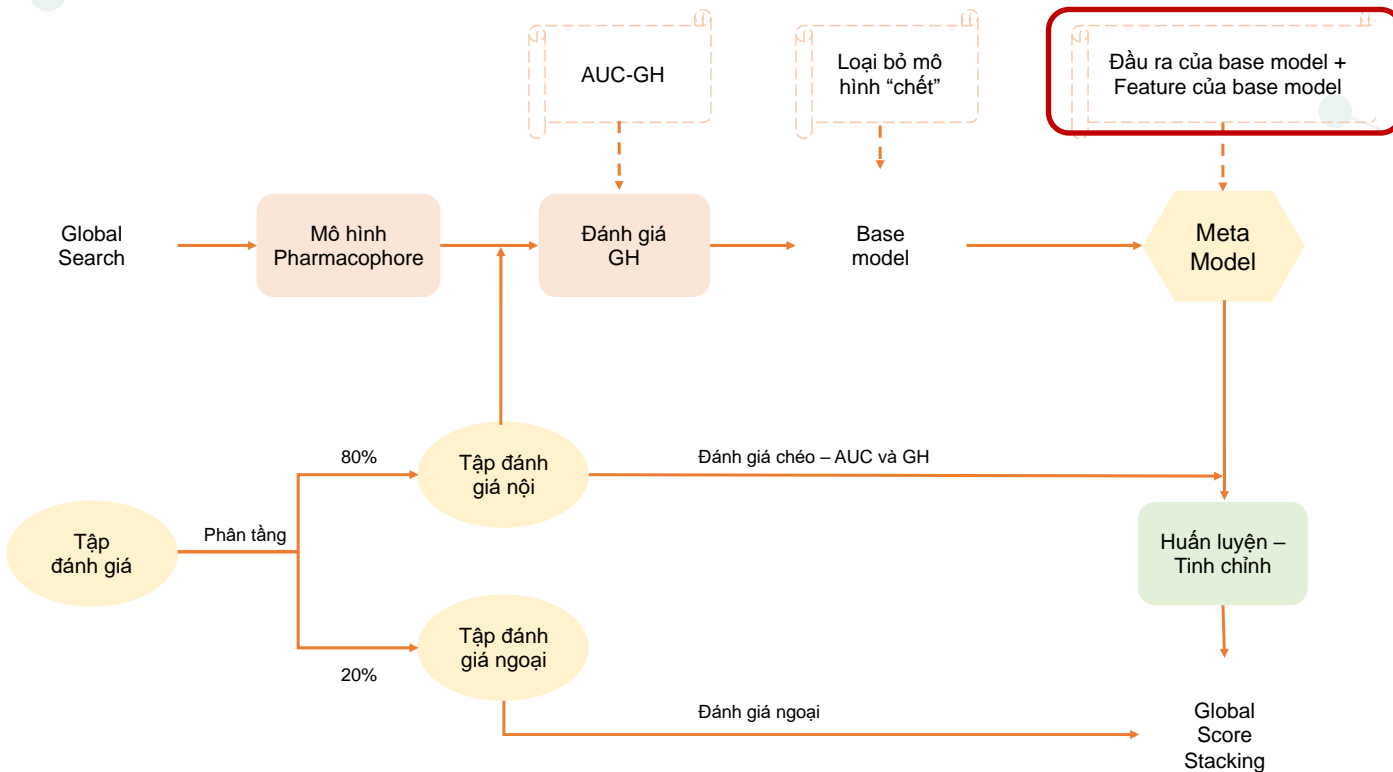
## Stacking



# METHOD

# PHARMACOPHORE

## Score Stacking





## METHOD

### Model

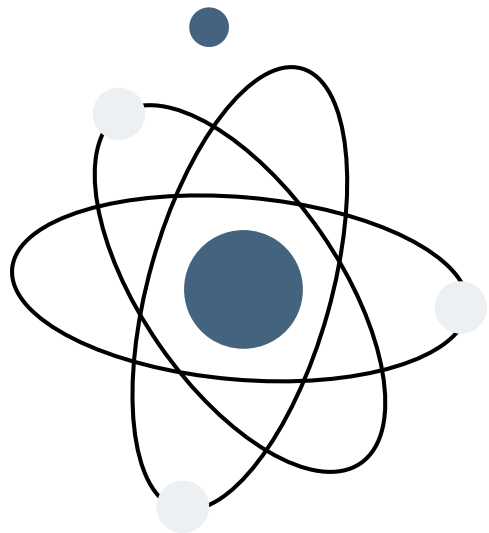
1. Local Search
2. Global Search
3. Voting
4. Stacking
5. Score Stacking
- 6. Feature Score Stacking**
- 7. Optimized Feature Score Stacking**

## PHARMACOPHORE

## EVALUATION

### Evaluation method

- **Cross validation:** post hoc  
Wilcoxon
- **External validation:** generalization
- **Metric:**  
AUC > GH > F1

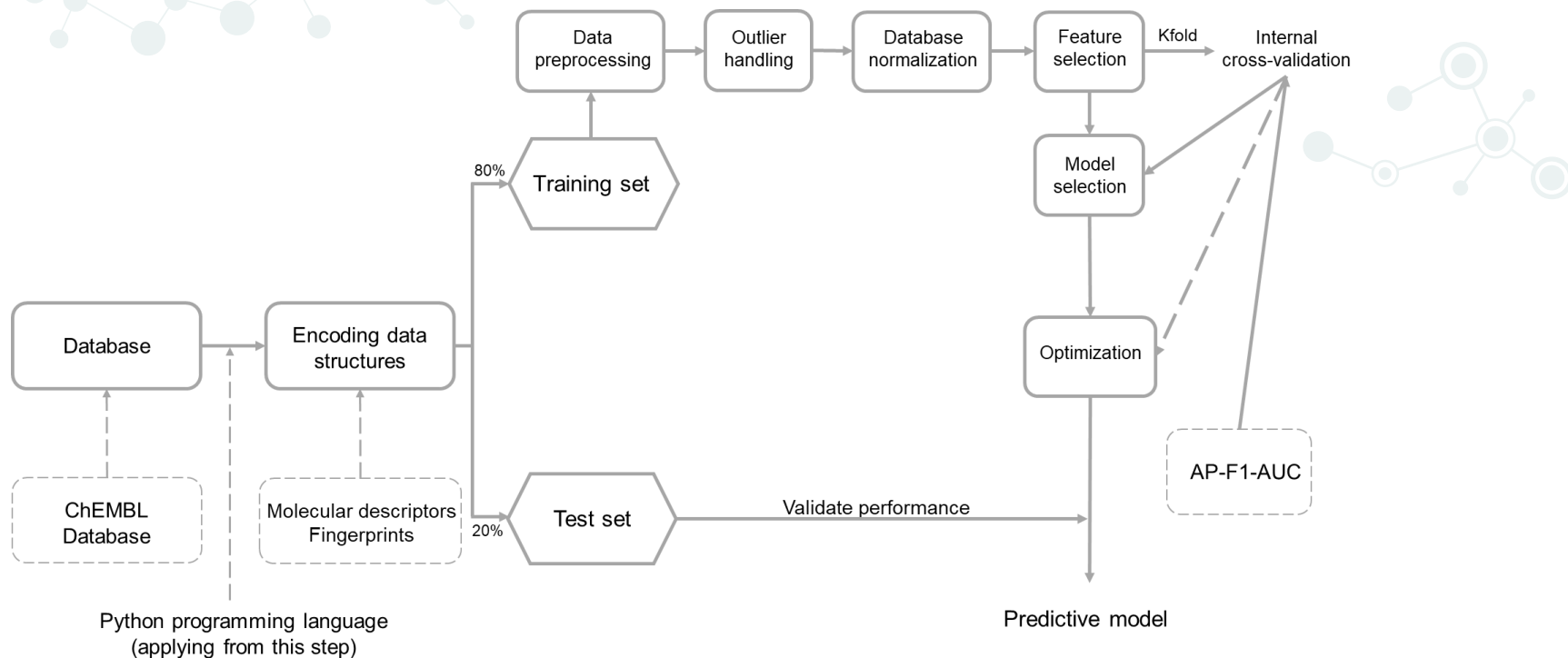
A dark blue background with a faint, light blue molecular structure graphic consisting of interconnected circles and lines, representing atoms and bonds. A vertical white line descends from the top center, ending in a white circle with a dark blue border containing the number "02".

# QSAR Model



# METHOD

# MÔ HÌNH QSAR



[Source code](#)

# METHOD

# MÔ HÌNH QSAR

## Thông số mô tả

2D-Mordred  
RDKit  
Mol2vec

## Dấu vân đường dẫn

RDK5  
RDK6  
RDK7

## Dấu vân tay từ điển

MACCS  
PubChem  
Avalon

## Dấu vân tay tròn

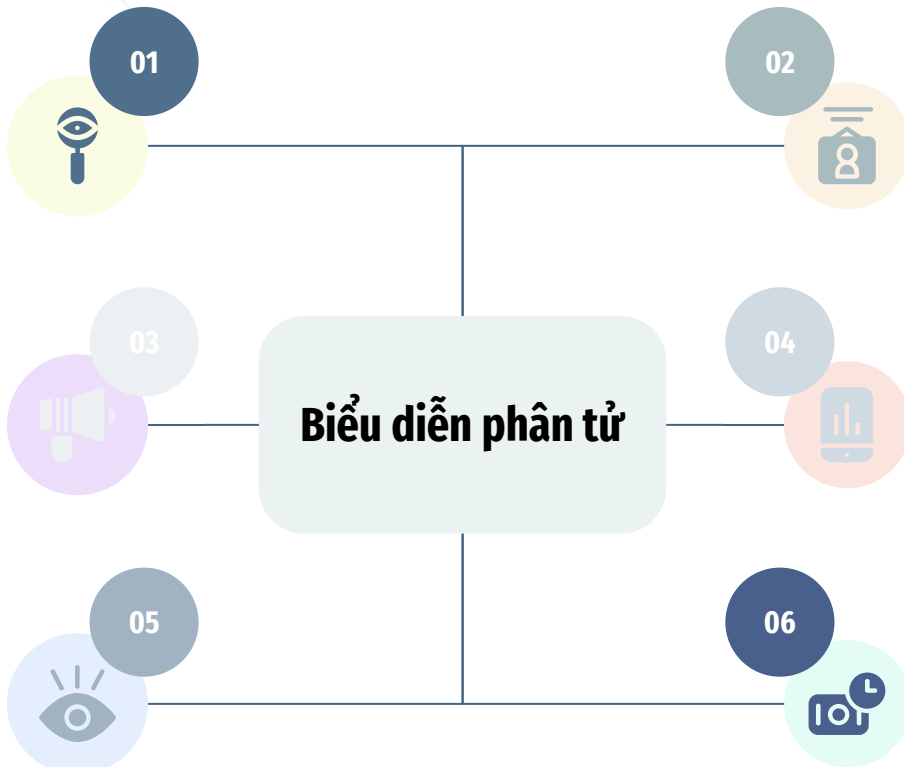
ECFP2  
ECFP4  
ECFP6

## Dấu vân tay PH4

Cats2D  
Pharmacophore Gobbi

## Dấu vân tay khác

map4  
secfp



# METHOD

QSAR

Data-centric



**01**

**02**

**03**

**04**

**05**

**06**

Xây dựng  
cơ sở dữ liệu

Khai phá  
dữ liệu

Lựa chọn bộ  
dữ liệu tối ưu

Lựa chọn  
đặc trưng

Lựa chọn  
mô hình

Tinh chỉnh  
mô hình



PHƯƠNG PHÁP

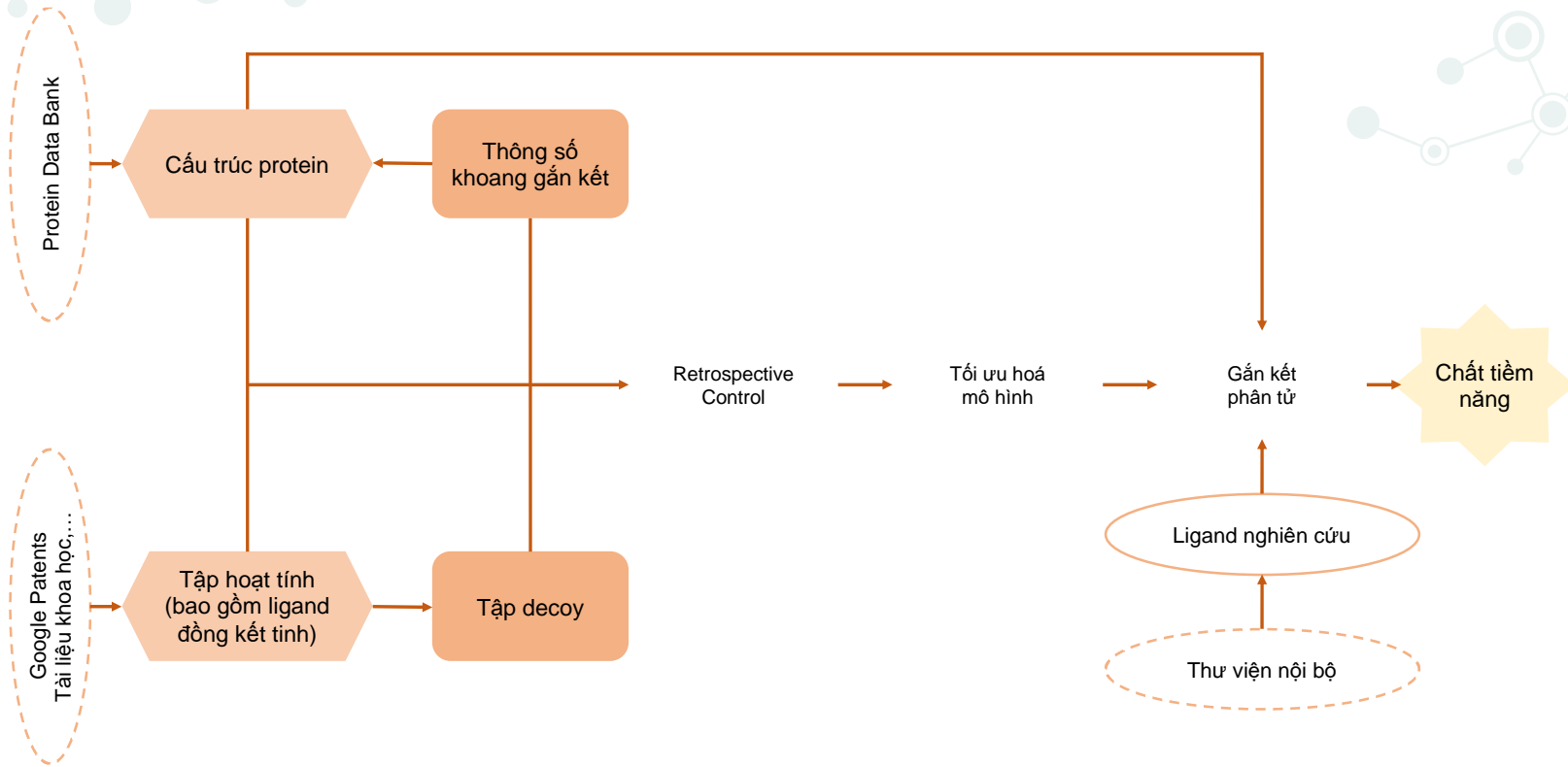


03

# MOLECULAR DOCKING

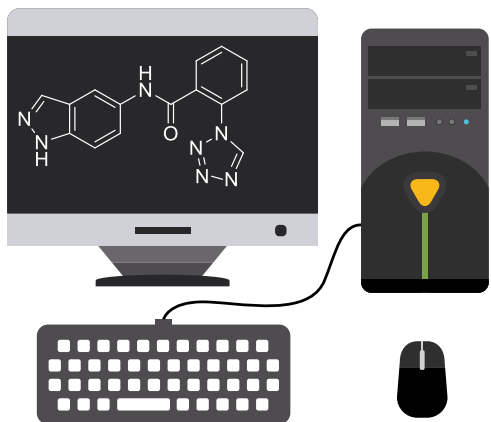
# METHOD

# MOLECULAR DOCKING





METHOD

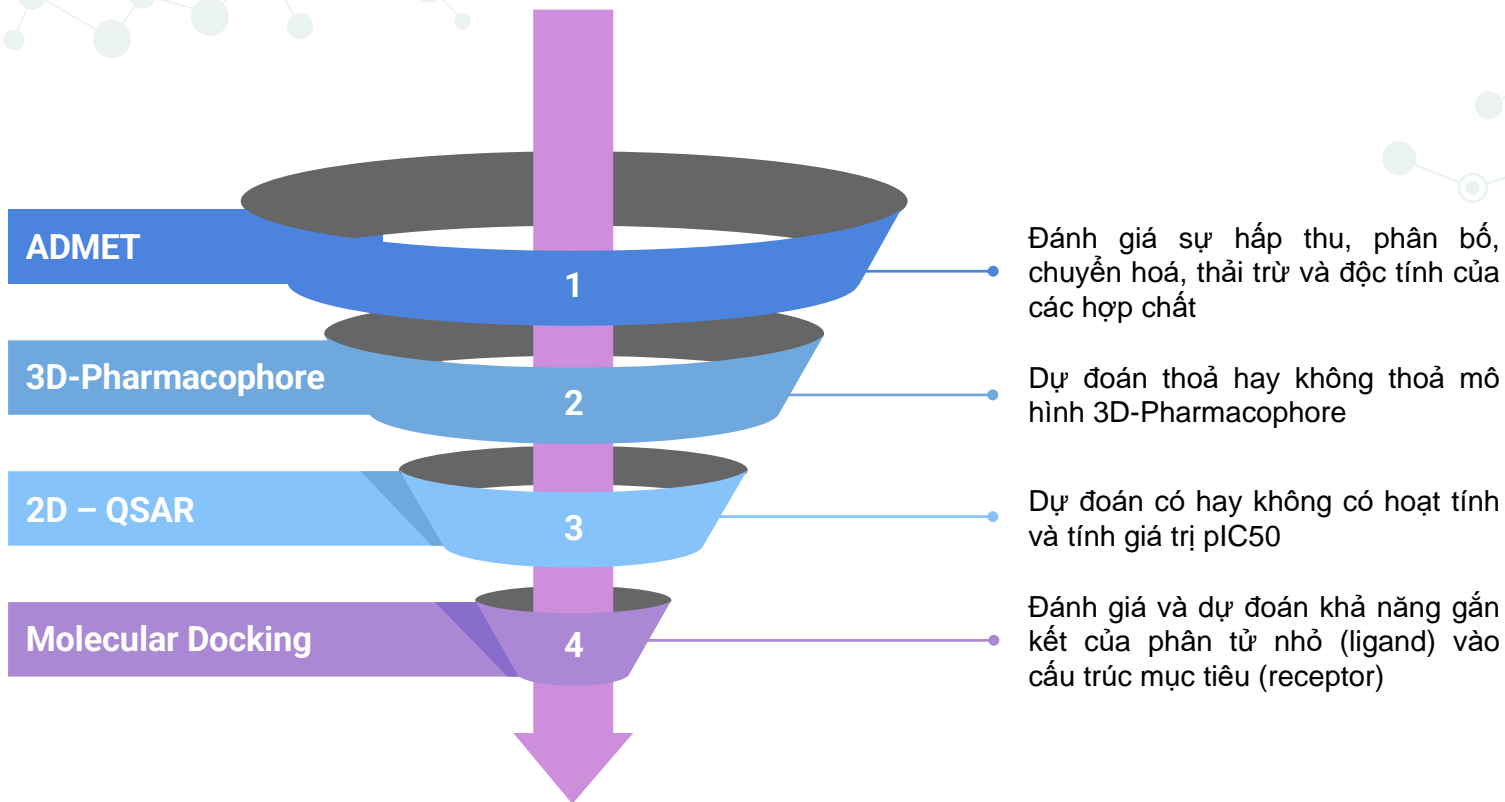


04

# VIRTUAL SCREENING

## METHOD

## VIRTUAL SCREENING





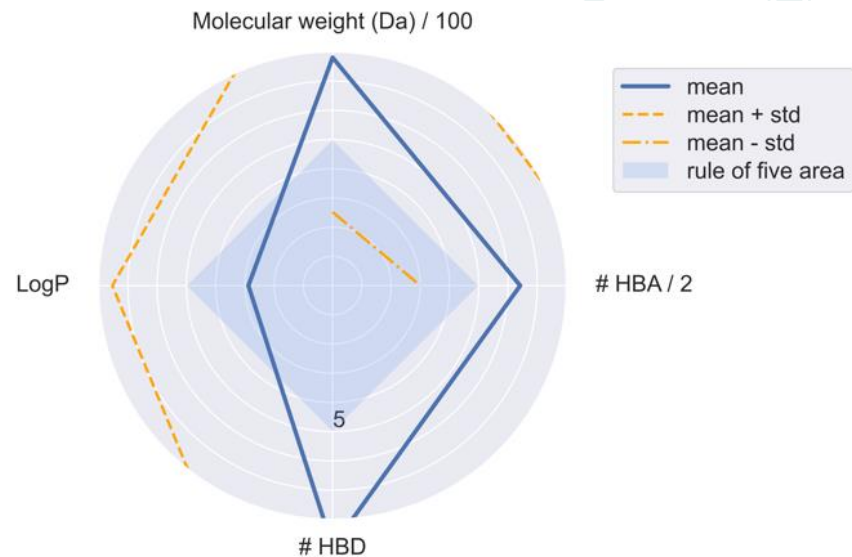
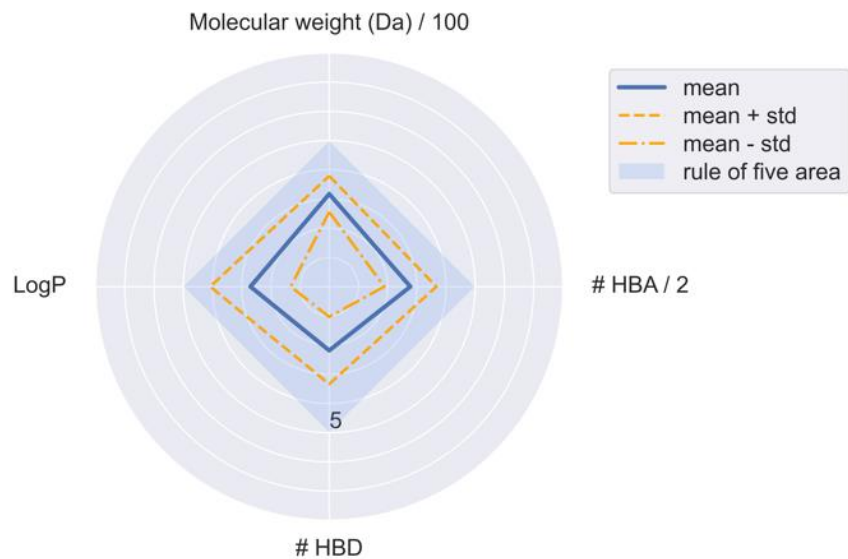
# RESULT- DISCUSSION





# RESULT

2445  
compounds





RESULT

1

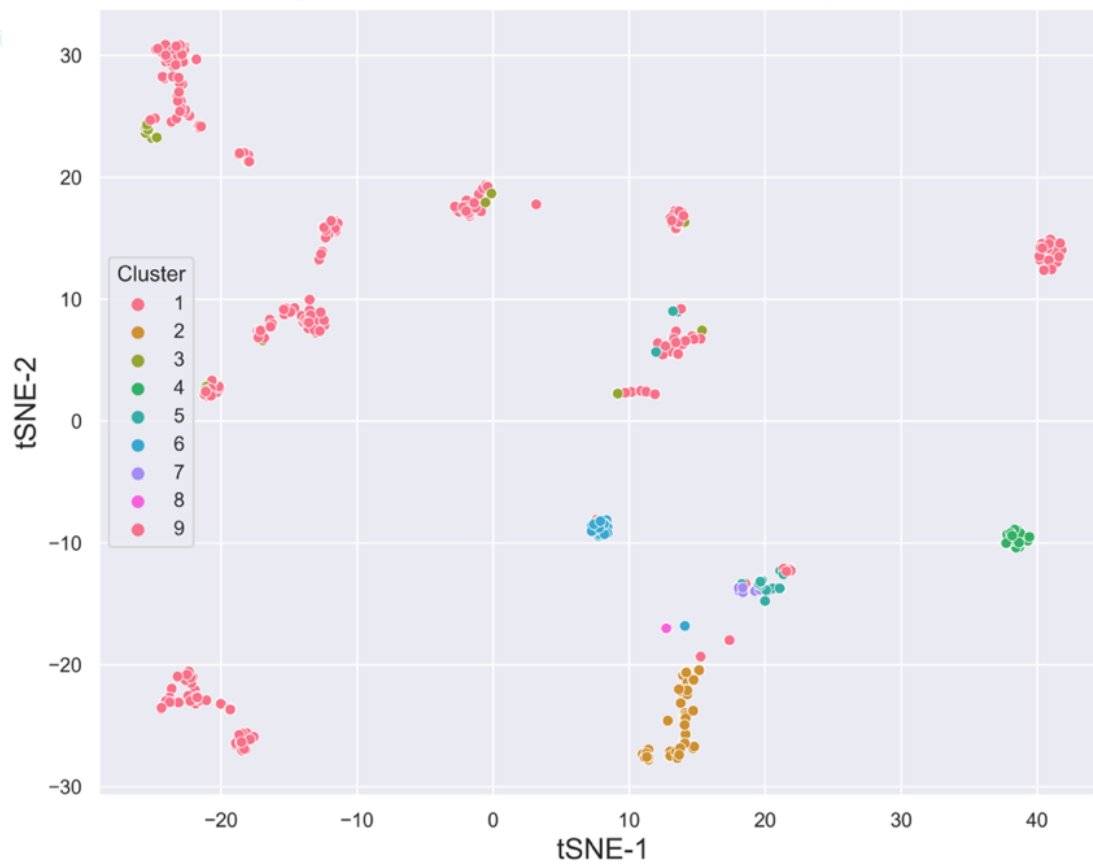
# Pharmacophore

Ligand-based pharmacophore





# RESULT



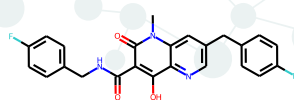
## Pharmacophore

### Clustering

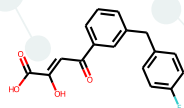




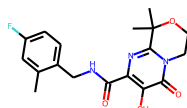
# RESULT



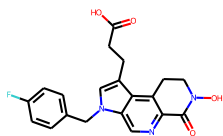
CHEMBL1914556



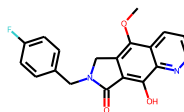
CHEMBL584360



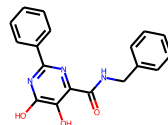
CHEMBL4126686



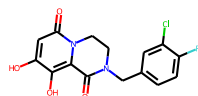
CHEMBL1773405



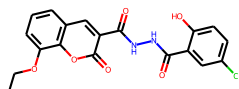
CHEMBL209440



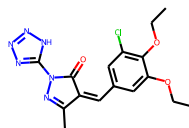
CHEMBL385951



CHEMBL237727



CHEMBL4463247



CHEMBL429327

# Pharmacophore

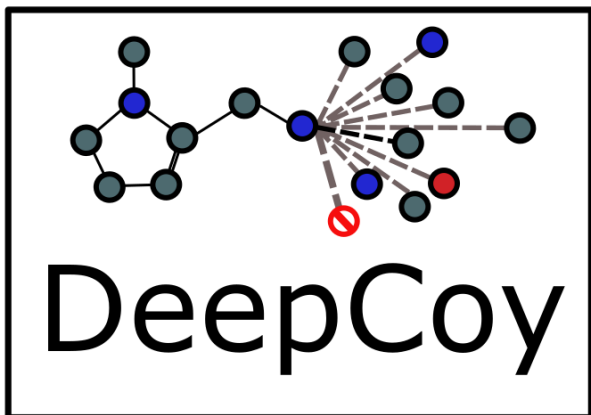




RESULT

Pharmacophore

Decoy



Metric	Values
AUC-ROC 1NN	0,533
AUC-ROC RF	0,702
DOE score	0,065
Doppelganger score mean	0,271
Doppelganger score max	0,421



## RESULT

Pharmacophore

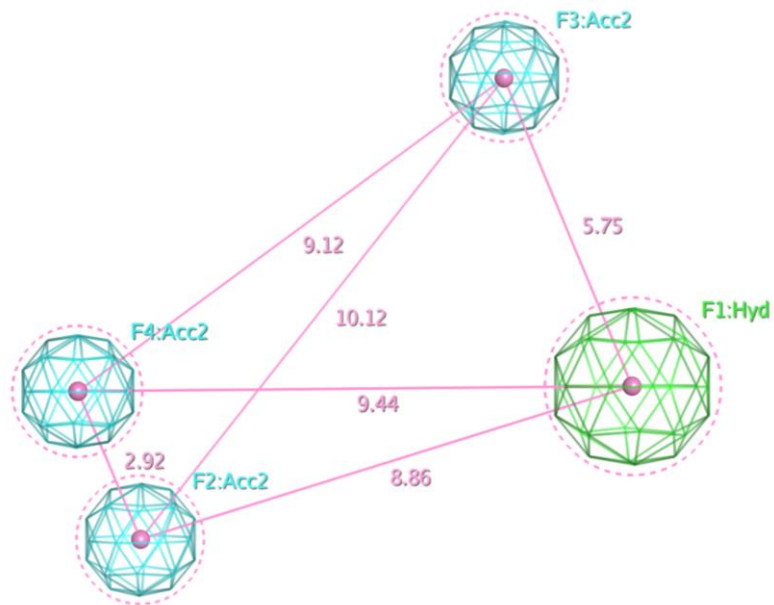
Local search

Model	$S_e$	$S_p$	$P_r$	AUC	GH
Haaa_7	0,814	0,573	0,037	0,794	0,133
Haaa_3	0,765	0,482	0,029	0,632	0,103
Haaa_4	0,8	0,6	0,039	0,748	0,138
Haaa_6	0,403	0,631	0,021	0,489	0,074
....					



# RESULT

## Haaa\_7 model



Pharmacophore

Local search



 RESULT

Pharmacophore

Global Search

12 models have AUC &gt; 0,7

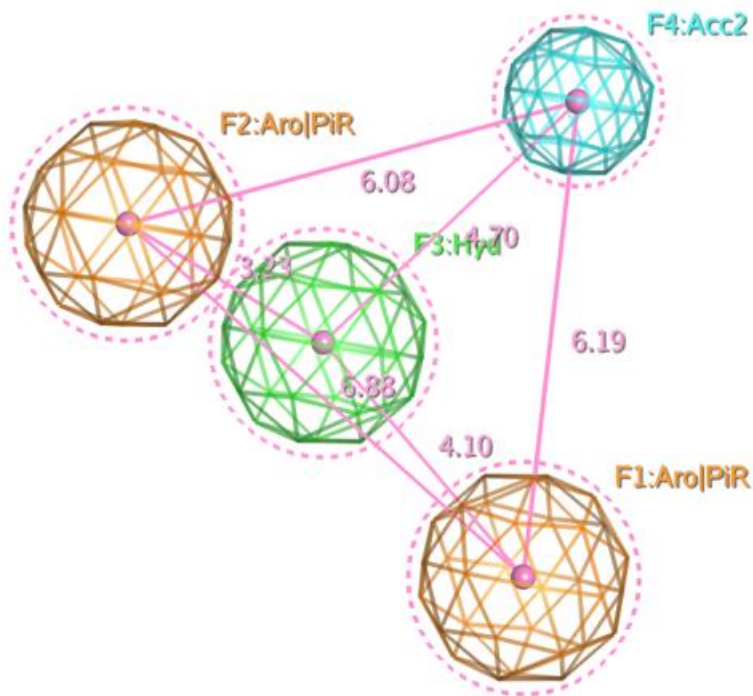
Model	S <sub>e</sub>	S <sub>p</sub>	P <sub>r</sub>	AUC	GH
RRHa_4	0,834	0,894	0,137	0,87	0,278
Haaa_7	0,814	0,573	0,037	0,794	0,133
RRHa_1	0,761	0,812	0,075	0,791	0,200
RHaa_7	0,846	0,546	0,036	0,775	0,130
....					





## RESULT

### RRHa\_4 model



## Pharmacophore

## Global search



- **Cover 7**
- F1 and F2 – Aromatic (Aro:  $r=1,4\text{\AA}$ )
- F3 is hydrophobic (Hyd:  $r=1,4\text{\AA}$ )
- F4 hydro acceptor (Acc2:  $r=1,0\text{\AA}$ )



RESULT

**Local Voting**



17 model with cover 9 vote



**Voting**

12 model (AUC > 0,7) vote



**Best Global Voting**

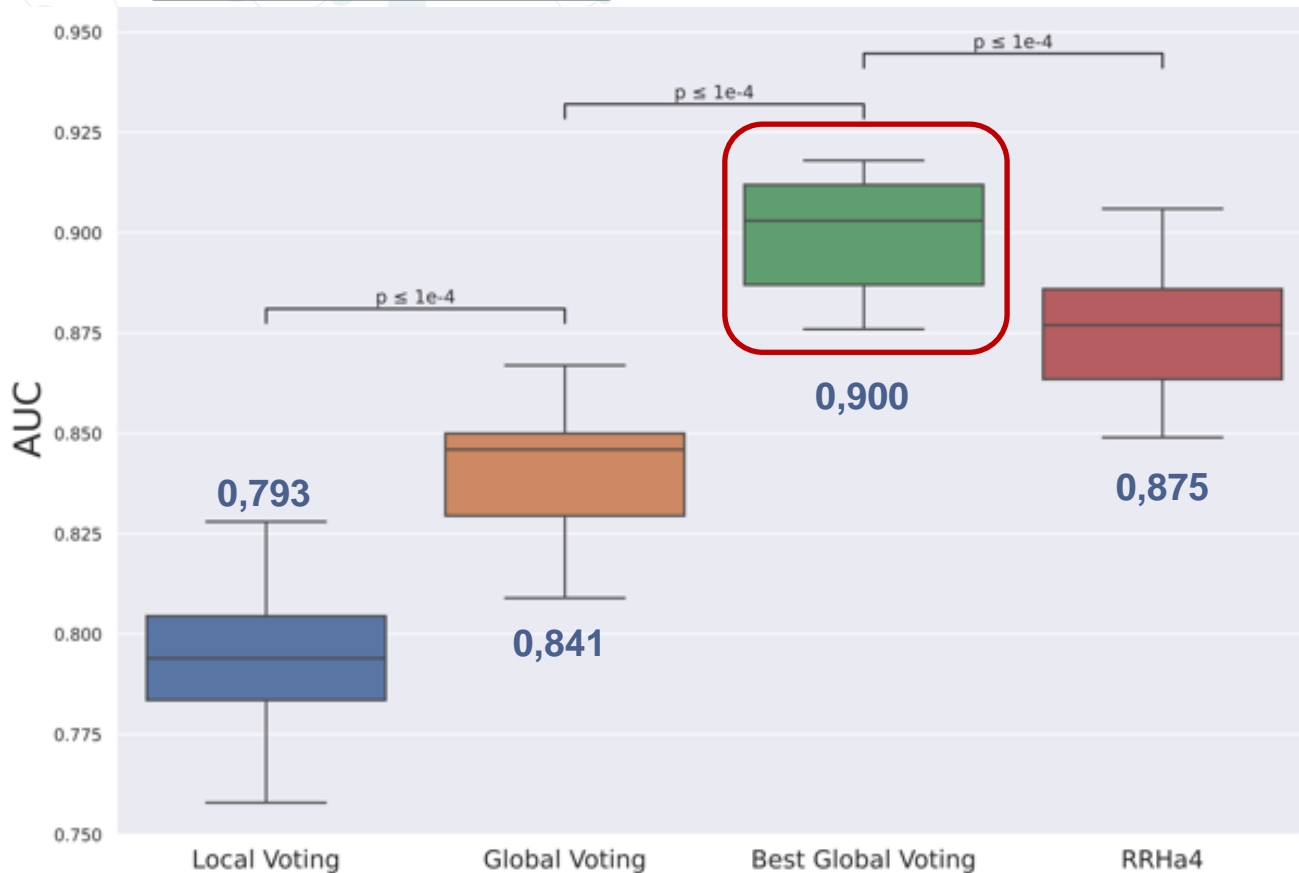
**Global Voting**



80 model vote



## RESULT



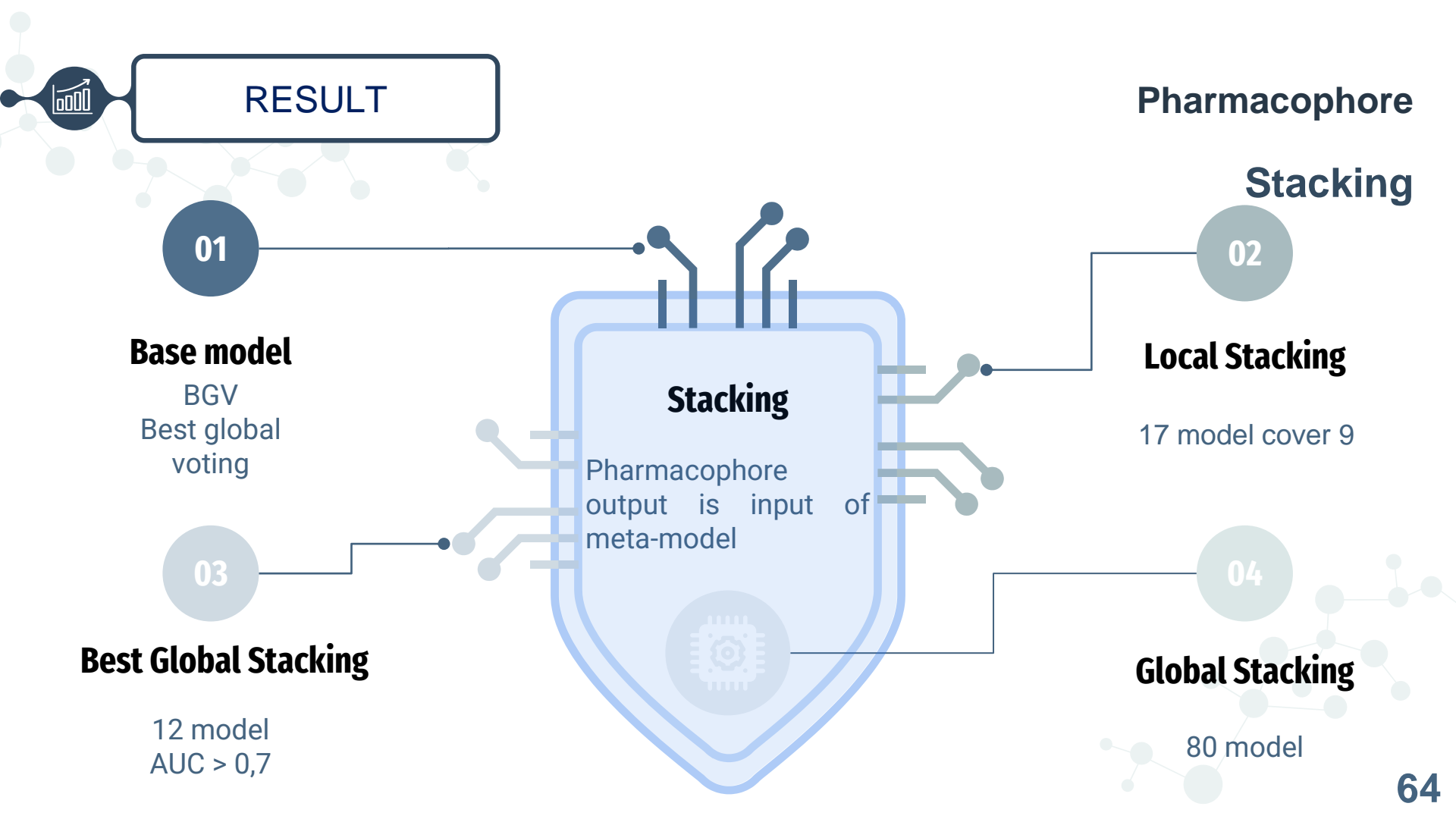
## Pharmacophore

### Voting



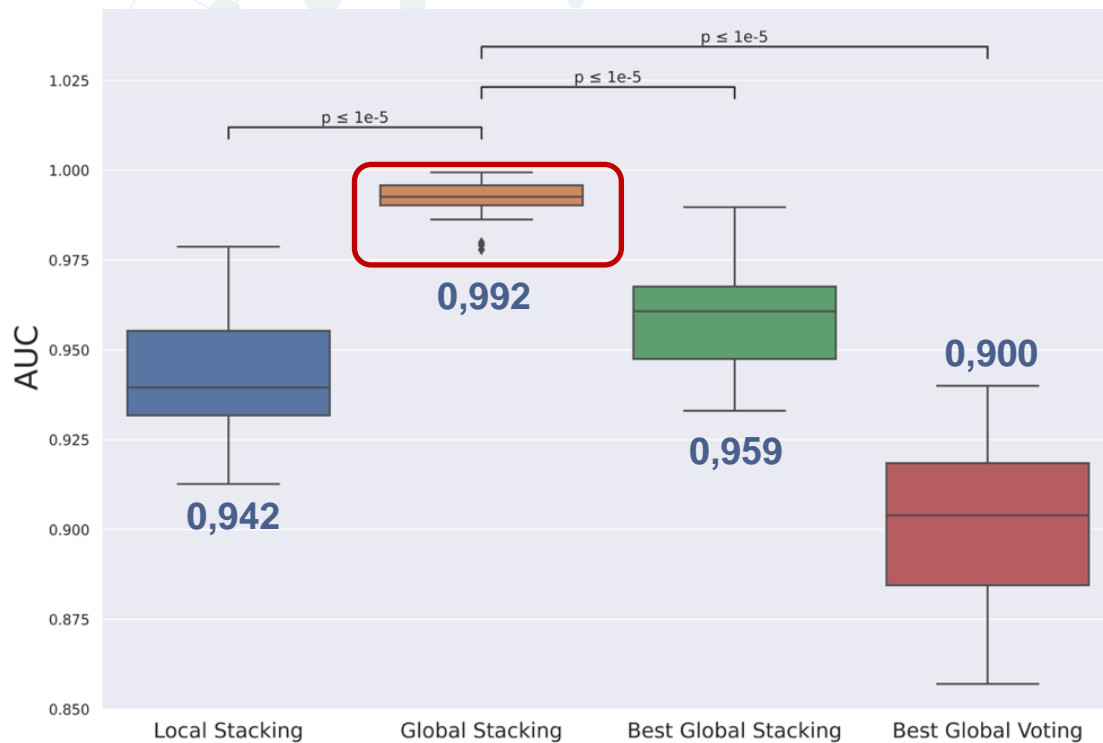
*Best global voting:*

- $AUC = 0,900 \pm 0,015$
- $GH = 0,334 \pm 0,049$





# RESULT



## Pharmacophore

## Stacking

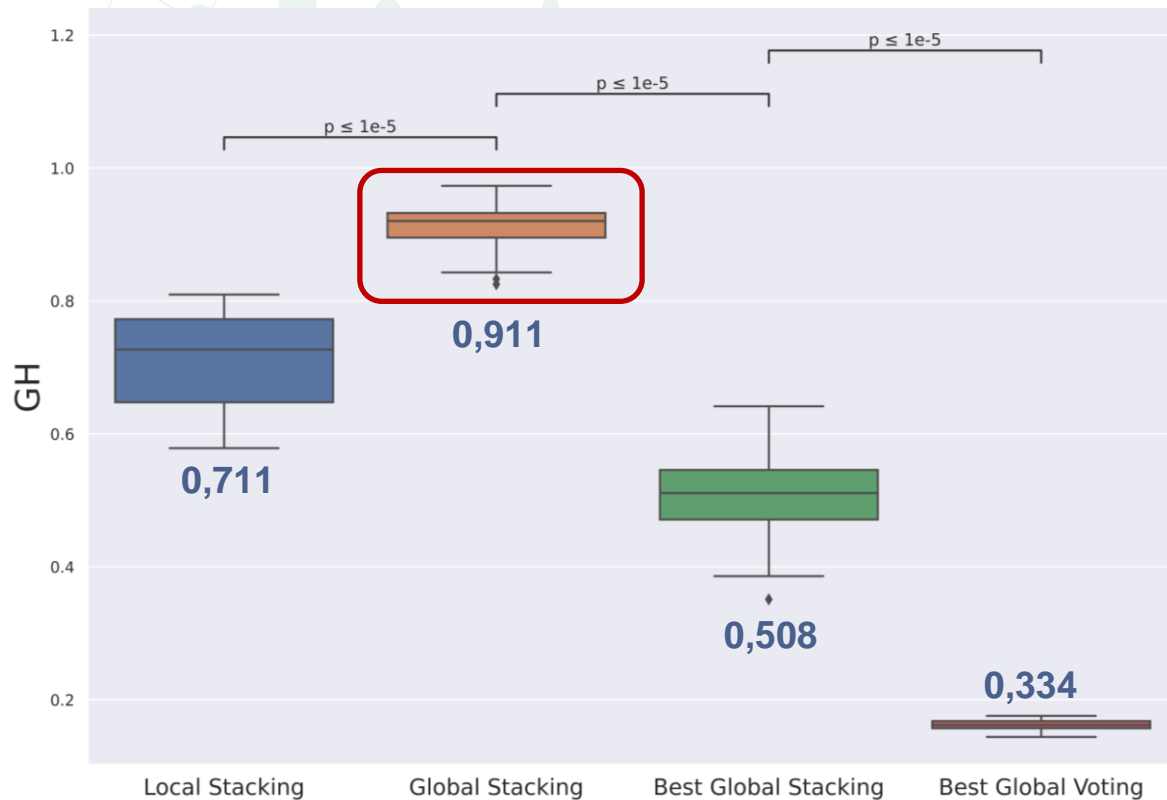


### Global stacking:

- $AUC = 0,992 \pm 0,017$
- $GH = 0,911 \pm 0,035$



## RESULT



## Pharmacophore

### Stacking



### Global stacking:

- $AUC = 0,992 \pm 0,017$
- $GH = 0,911 \pm 0,035$



# RESULT

01



Enhance data

**GSS model**

Using rescore from  
pharmacophore model

02



Feature selection

**F-GSS Model**

GSS feature selection  
XGBoost: n\_estimators =  
300 max\_depth = 4

03



Model optimization

**OF-GSS Model**

Optimize F-GSS

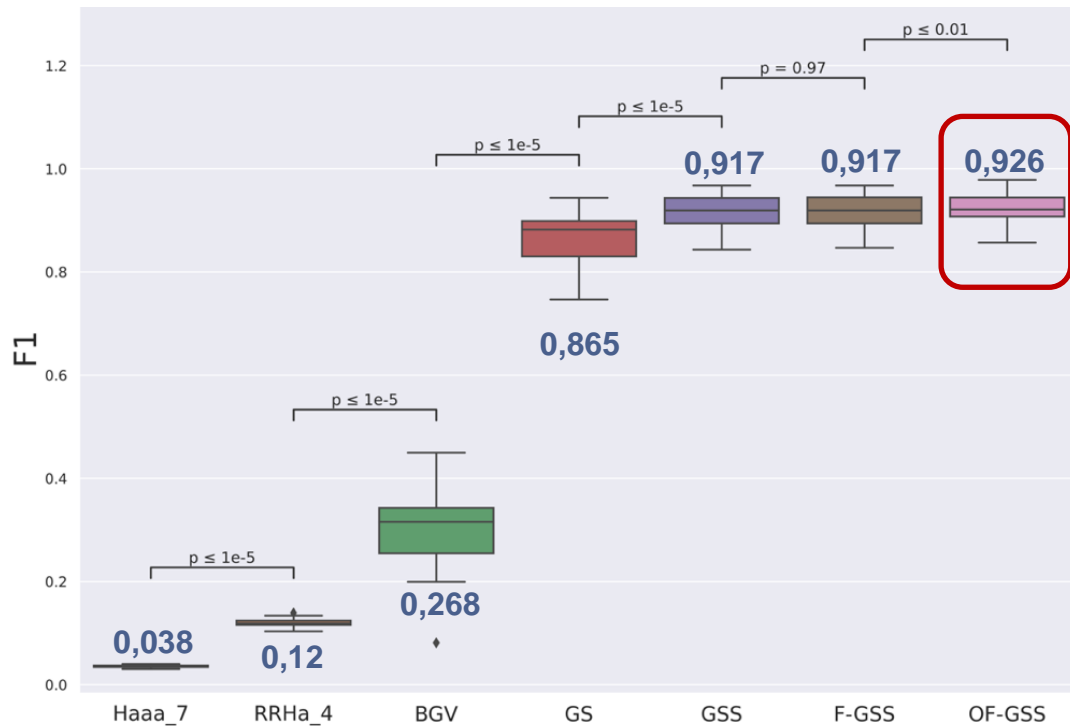
Pharmacophore  
Score Stacking



## DISCUSSION

## Pharmacophore

### Performance comparison

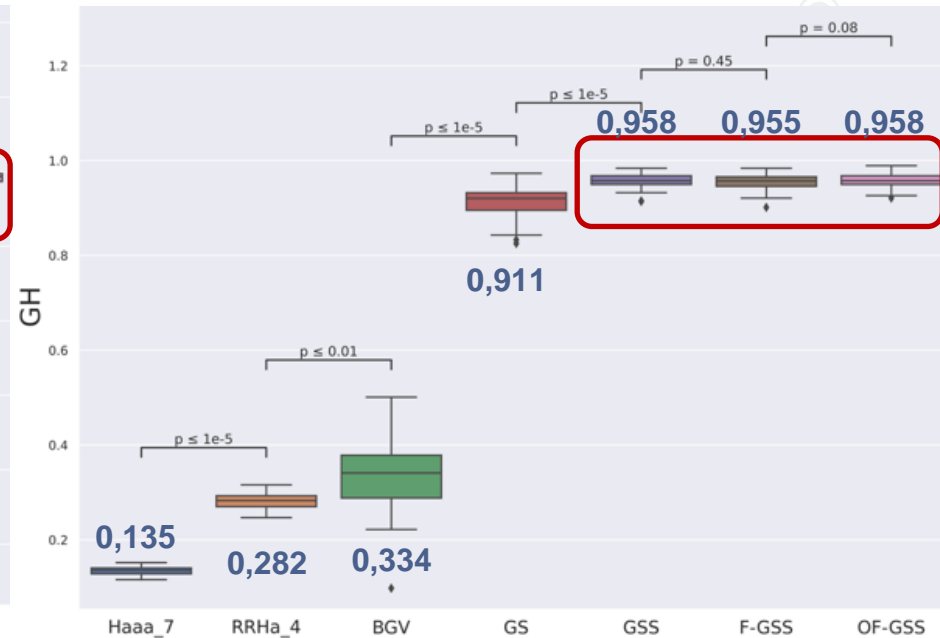
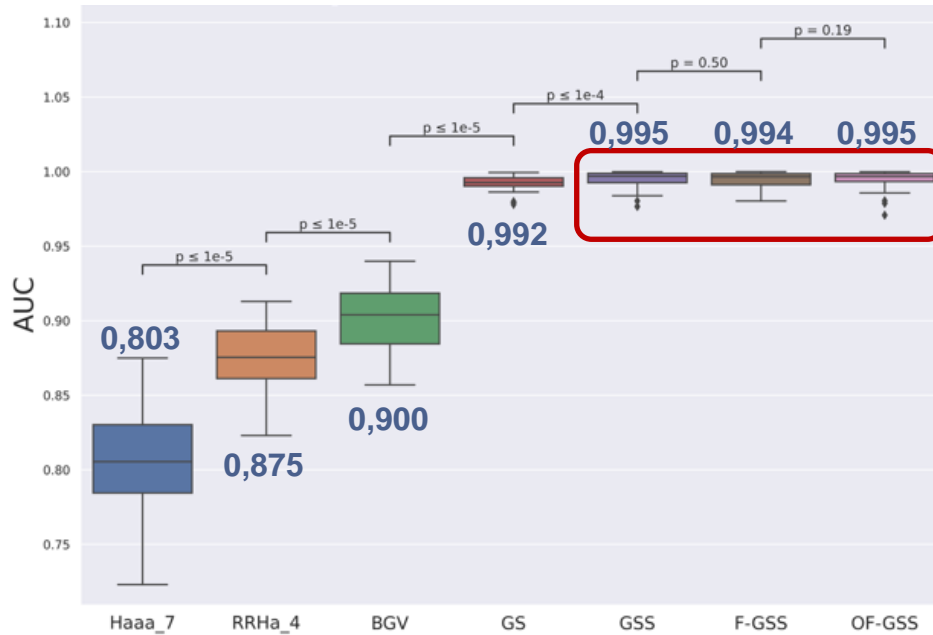


1. Local search: **Haaa\_7**
2. Global search: **Rhaa**
3. **BGV**: 12 model with AUC > 0,7
4. **GS**: stacking 80 model
5. **GSS**: using rescore
6. **F-GSS**: feature selection of GSS
7. **OF-GSS**: optimization of F-GSS

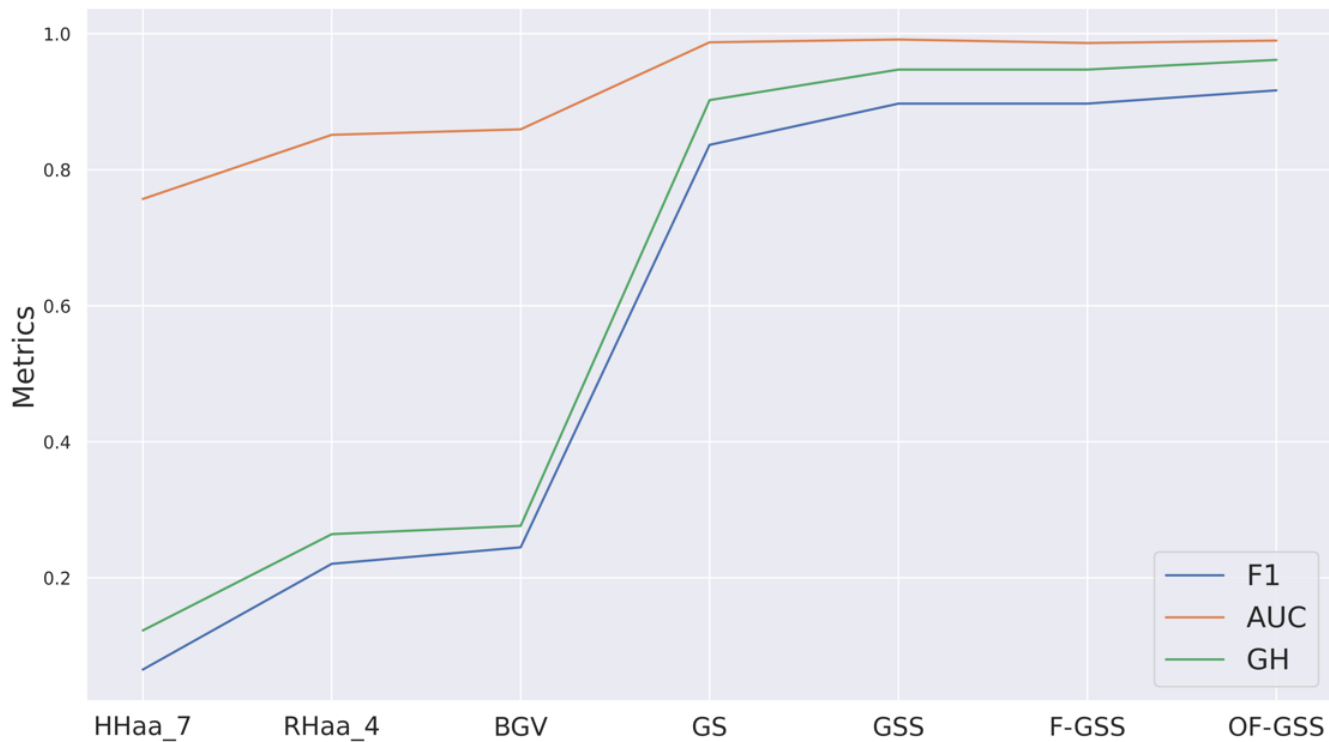


# DISCUSSION

## Pharmacophore Performance comparison



## DISCUSSION



## Pharmacophore External Validation

### OF-GSS

- F1 = 0,917
- AUC = 0,990
- GH = 0,962



RESULT

2

QSAR





## 2.1. Mô hình học máy phân loại

---

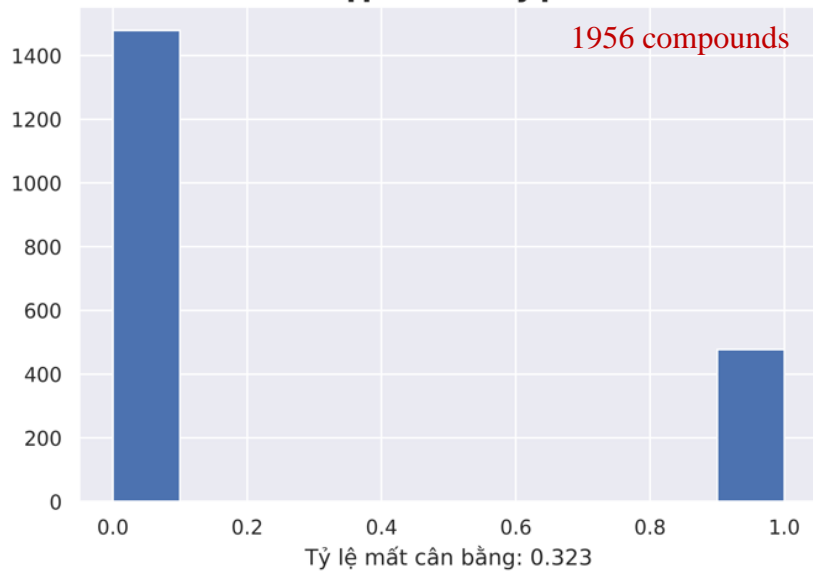


## QSAR – Classification

### Data preparation

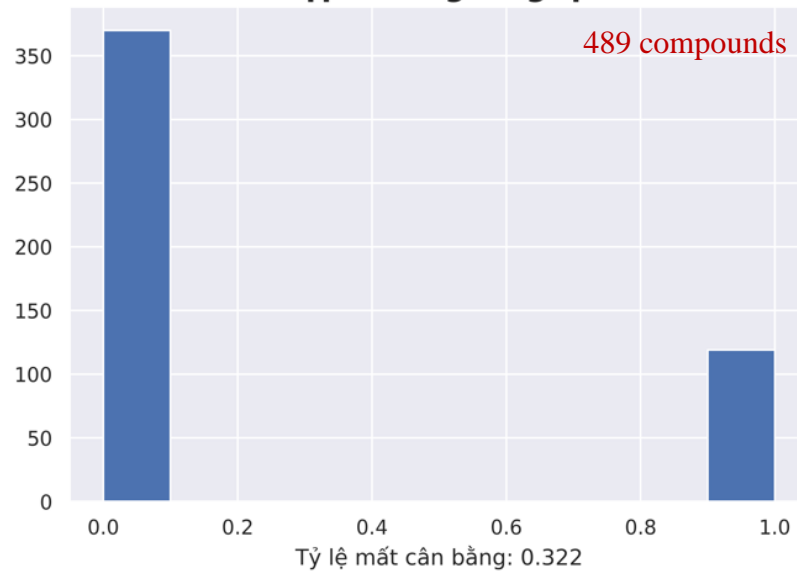
Tập huấn luyện

1956 compounds



Tập đánh giá ngoại

489 compounds

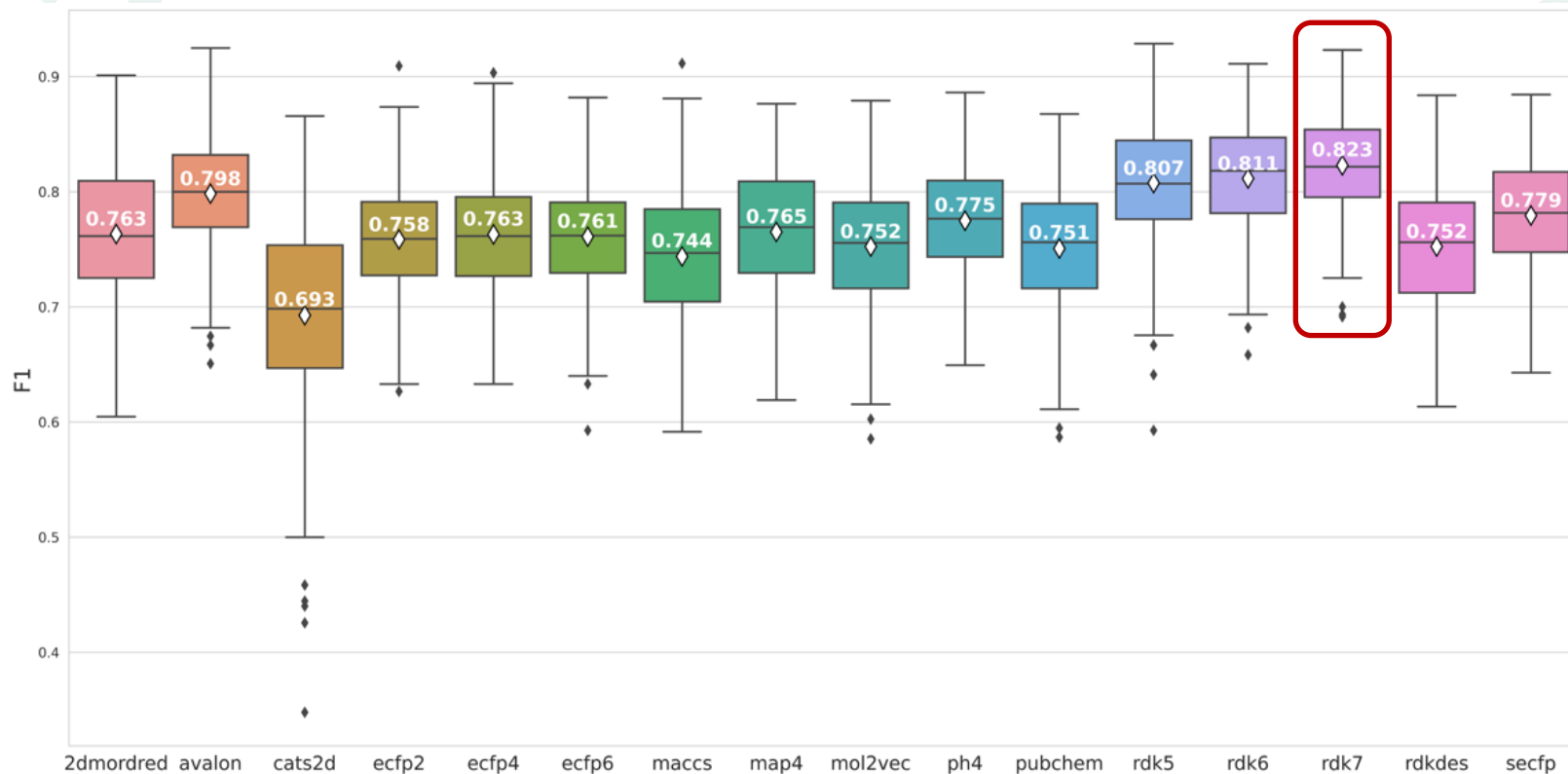




# RESULT

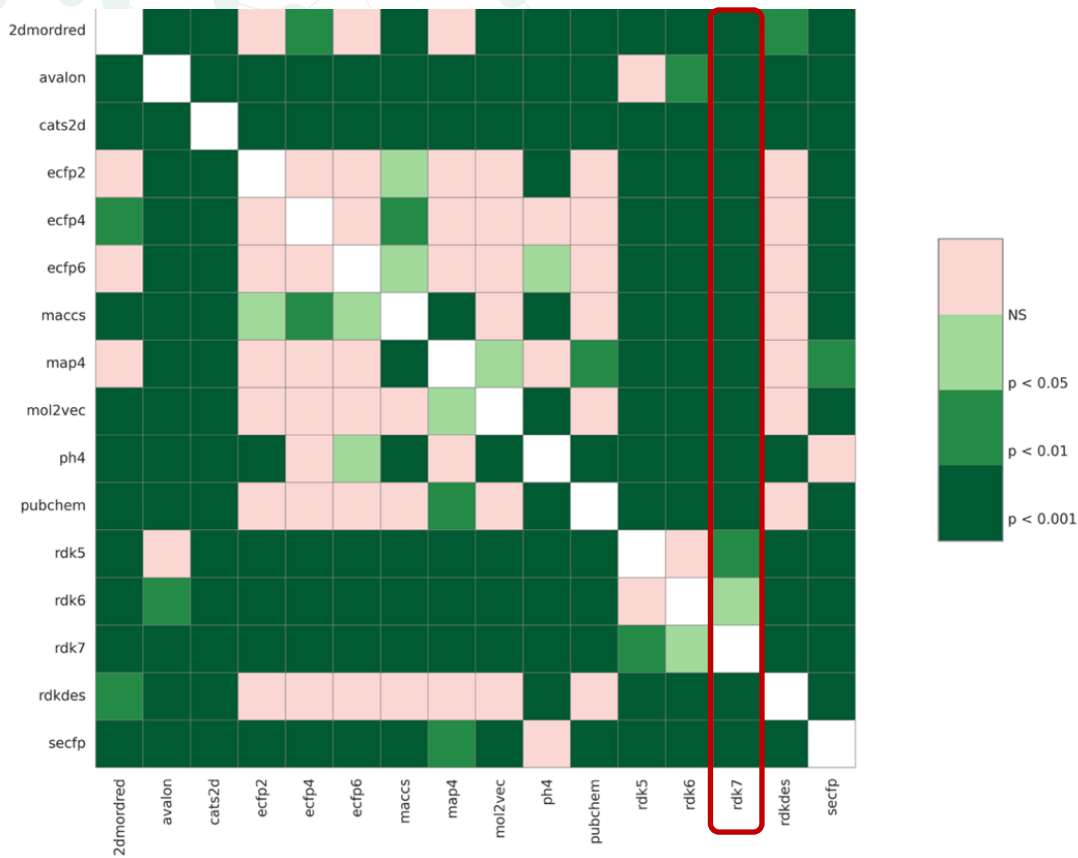
## QSAR – Classification

### Feature selection





## RESULT



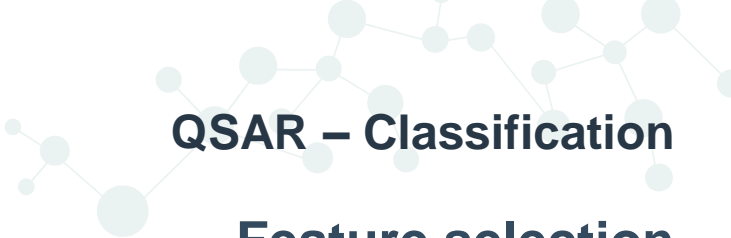
## QSAR – Classification

### Data selection



**Primary endpoint**

Highest F1 score



**QSAR – Classification**  
**Feature selection**

**Endpoint**



**01**

**Higher F1 score**

**02**

**Least features**



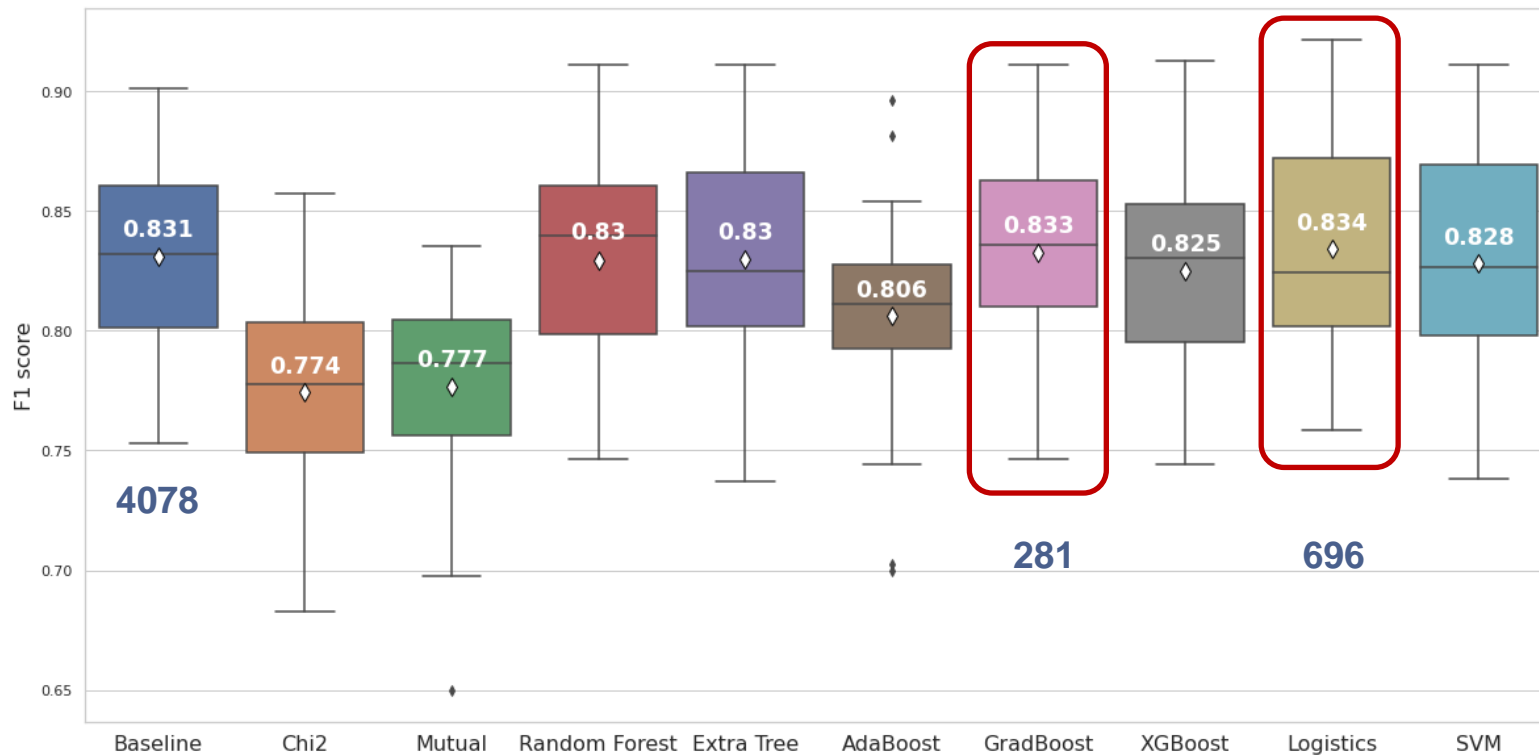




# RESULT

## QSAR – Classification

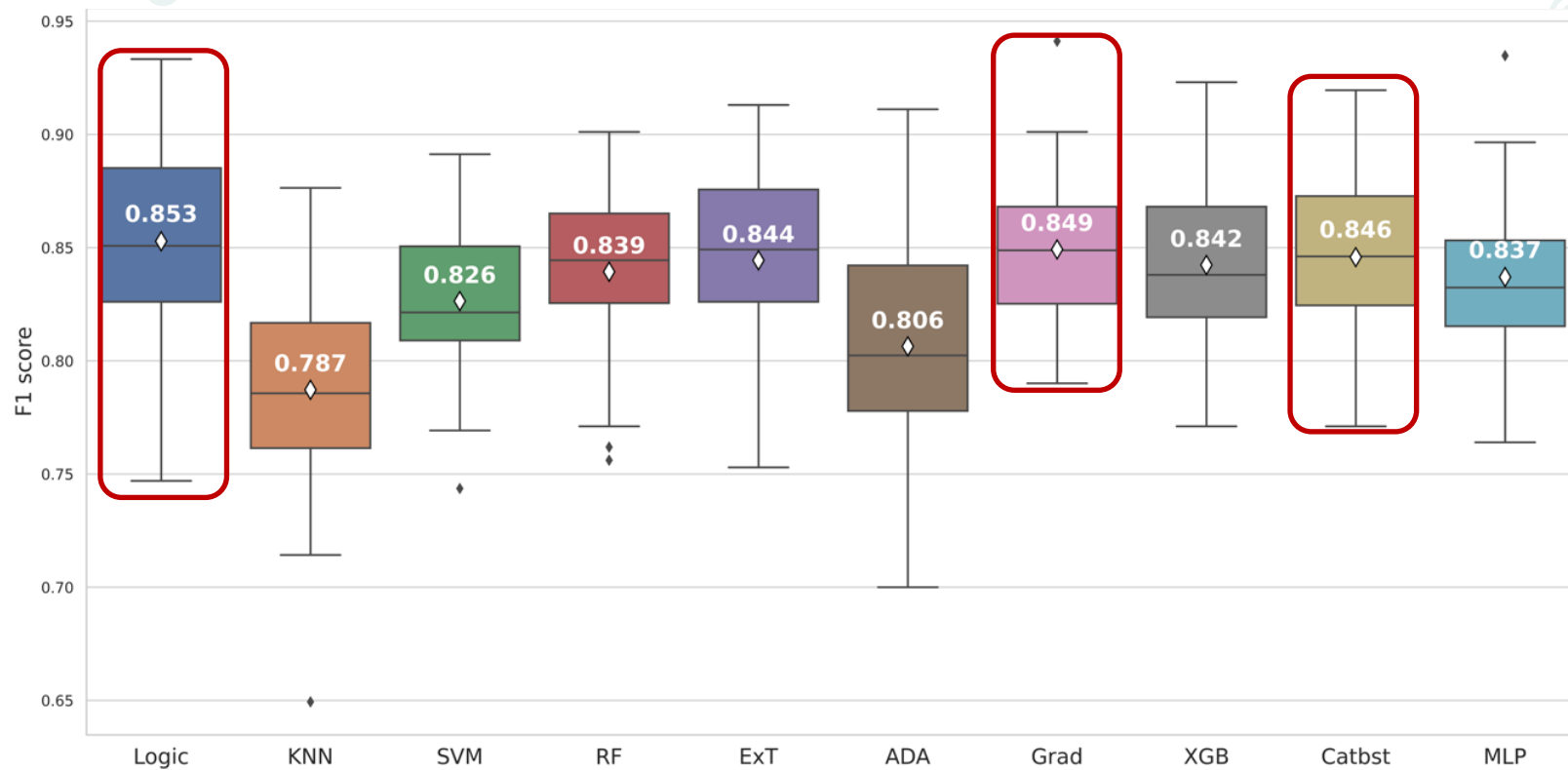
### Feature selection



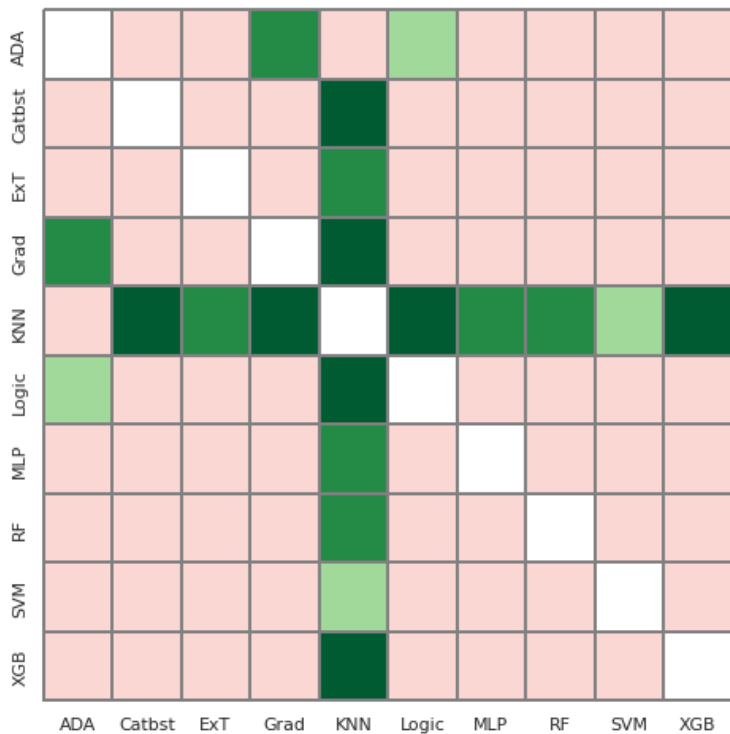
# RESULT

## QSAR – Classification

### Model selection



# RESULT



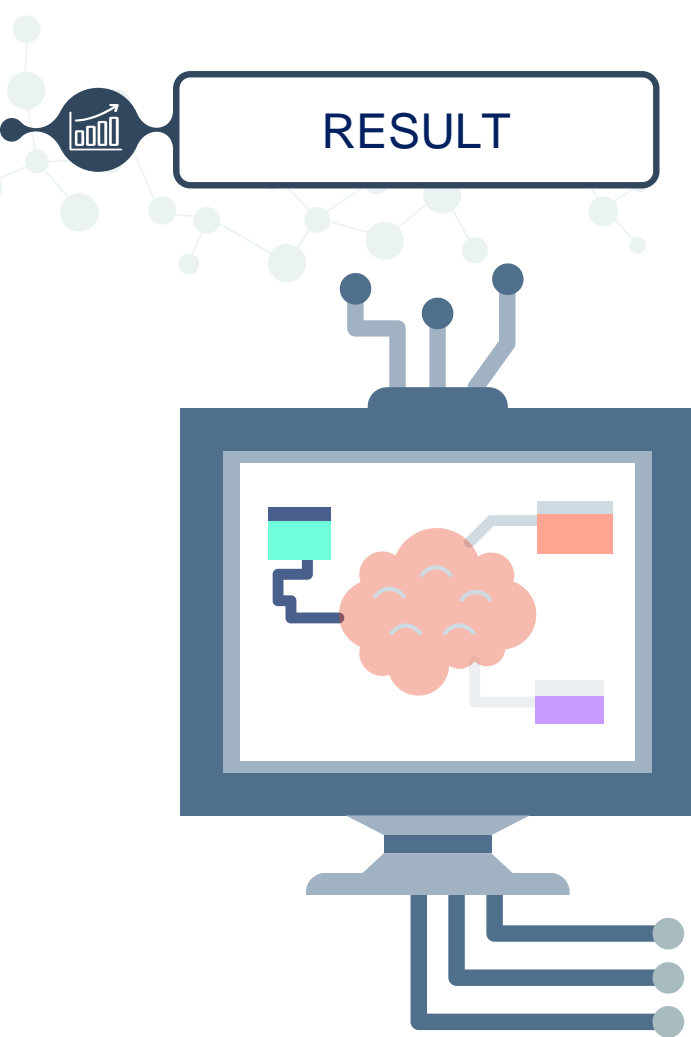
## QSAR – Classification

### Model selection

0

#### Optimized Model

- Logistic regression
- **Gradient Boosting**
- CatBoost



## QSAR – Classification

### Model Optimization

#### Data sampling

Over-sampling, under-sampling



#### Hyperparameter Tuning

Grid search,

#### Probability

Sigmoid (Platt Scaling) or isotonic (Isotonic Regression).

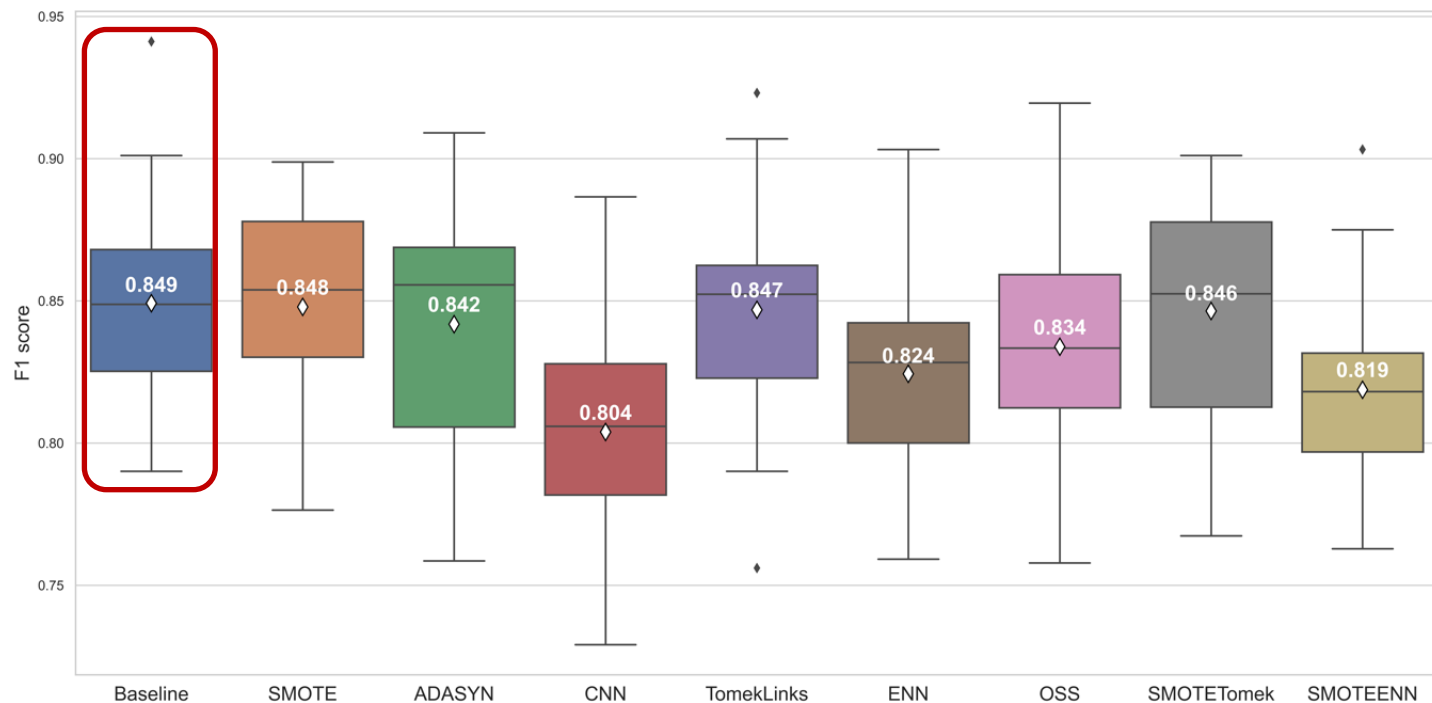


# RESULT

Data sampling

## QSAR – Classification

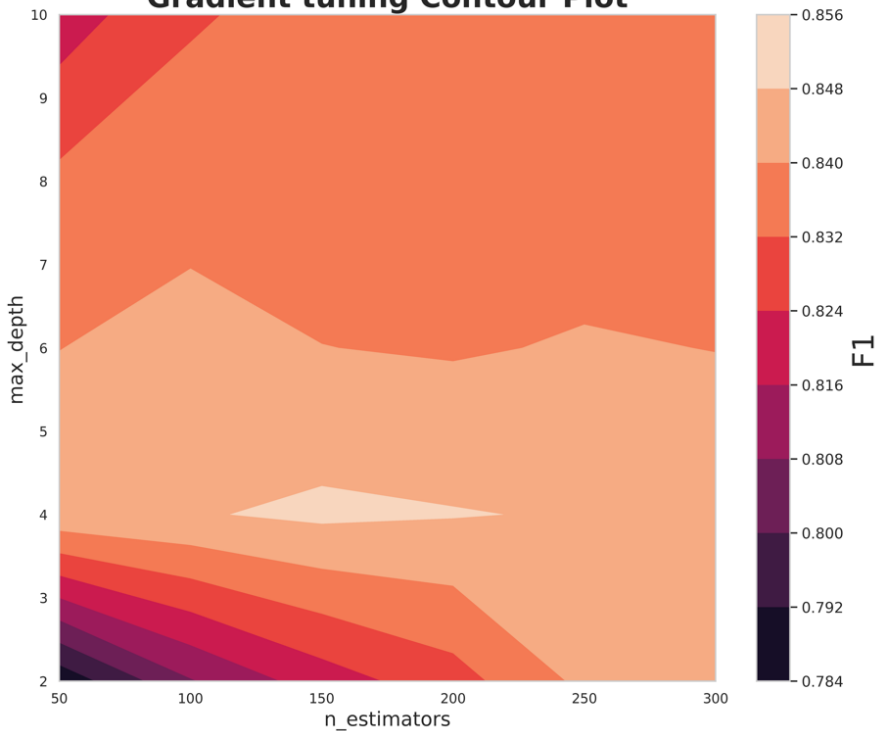
## Model Optimization



# RESULT

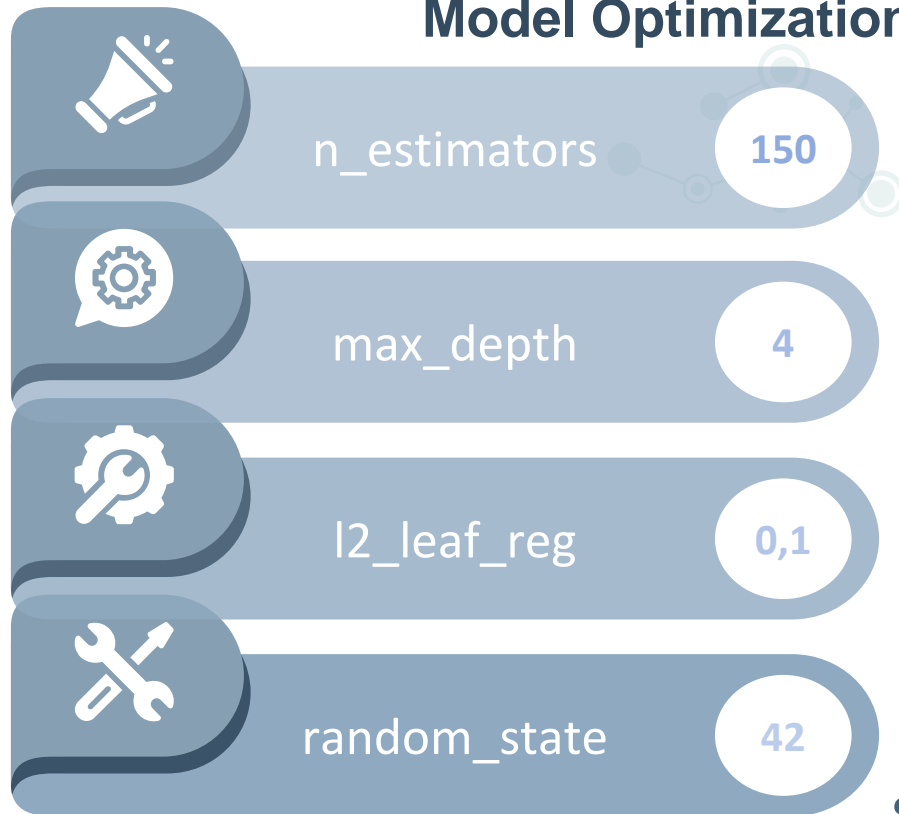
## Tuning

Gradient tuning Contour Plot



# QSAR – Classification

## Model Optimization

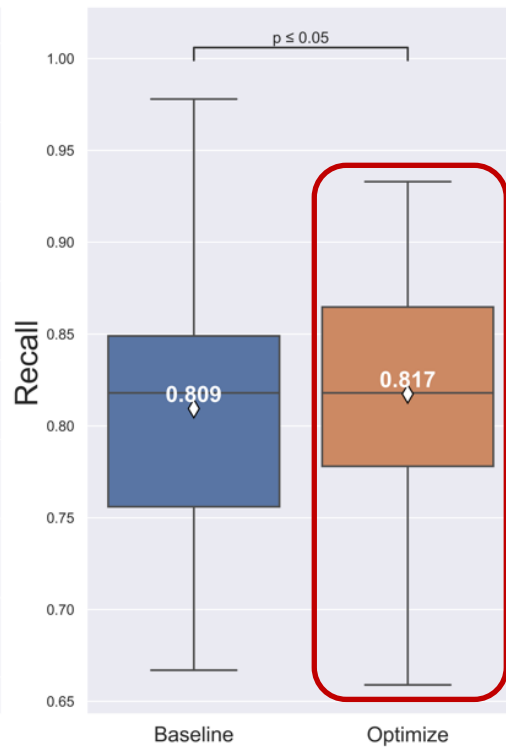
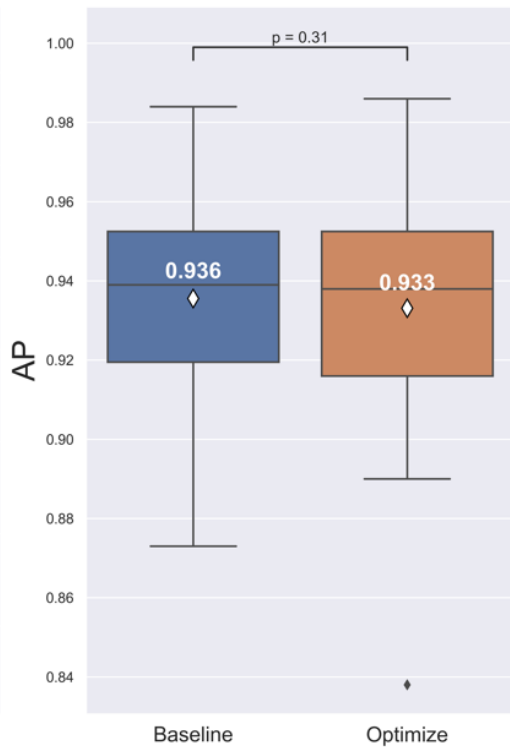
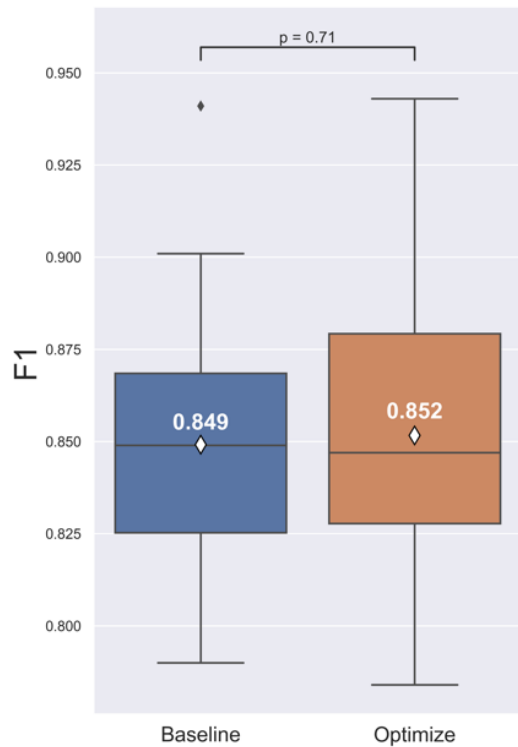




# RESULT

## QSAR – Classification

## Model Optimization





## QSAR – Classification

### External Validation

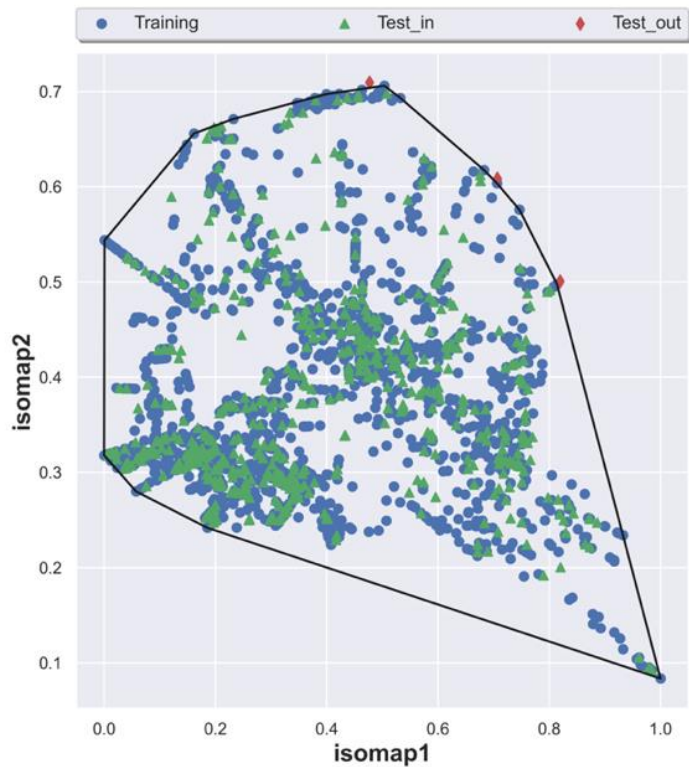
	Internal Validation			External Validation		
	AP	F1	Độ nhạy	AP	F1	Độ nhạy
<b>Baseline</b>	0,936	0,849	0,809	0,928	0,873	0,864
<b>Optimize</b>	0,933	0,852	<b>0,817</b>	<b>0,938</b>	<b>0,874</b>	<b>0,865</b>





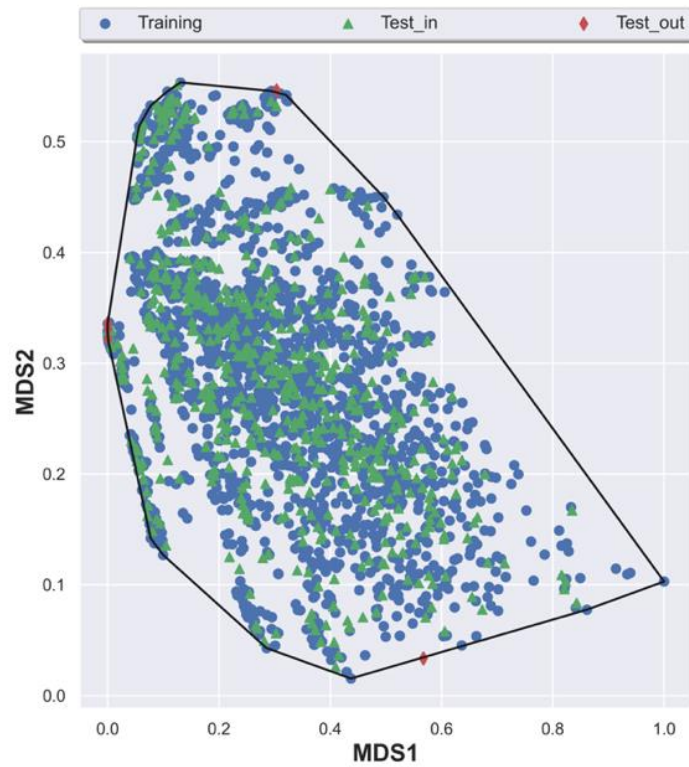
# RESULT

Convex hull



# QSAR – Classification

Applicability domain





RESULT

QSAR – Classification

Applicability domain

Isomap	MDS
CHEMBL394650	CHEMBL3310415
CHEMBL252831	CHEMBL2236598
CHEMBL19043	CHEMBL3310409
	CHEMBL1914564
	CHEMBL252757



RESULT

3

# Docking



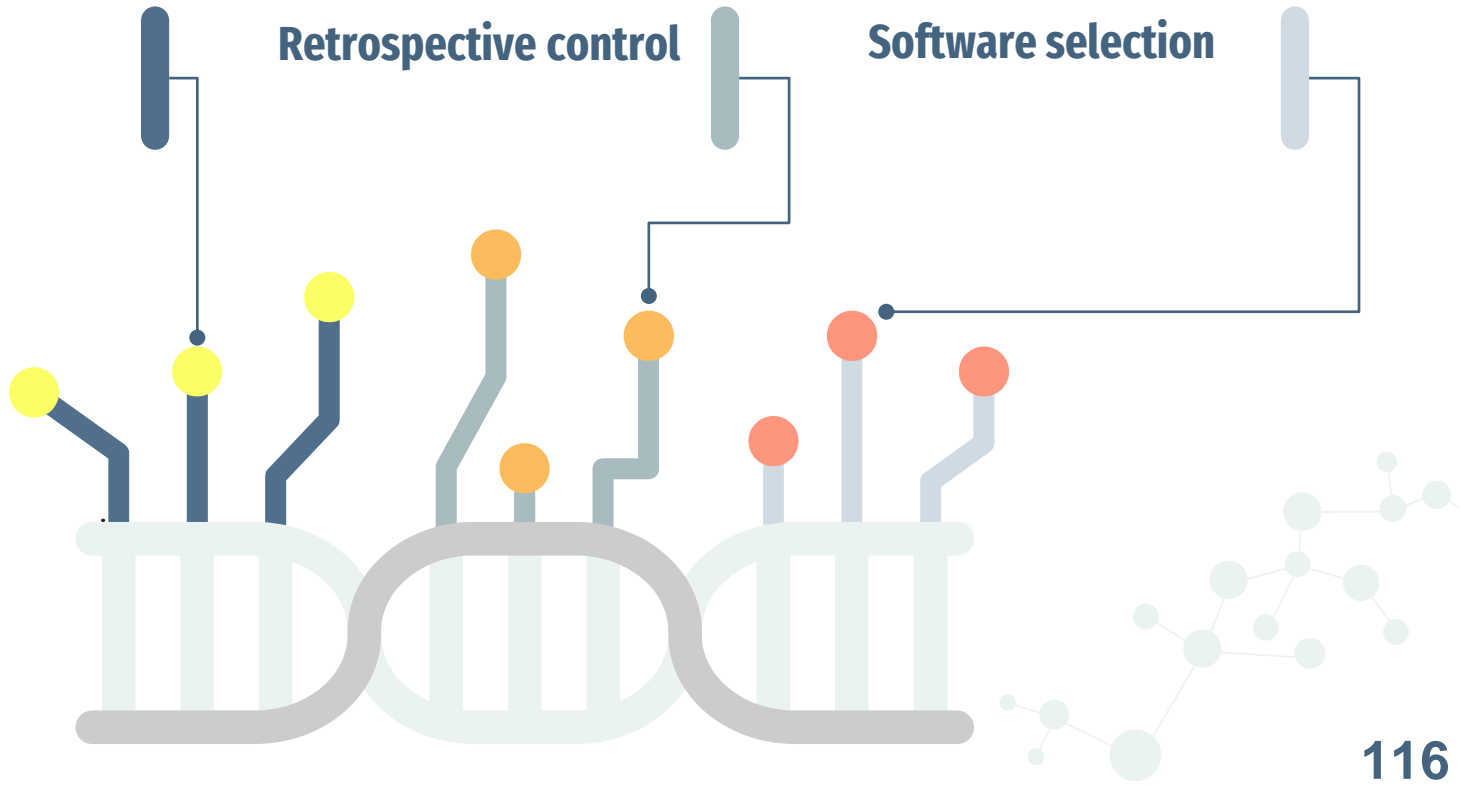


# MOLECULAR DOCKING

Re-docking

Retrospective control

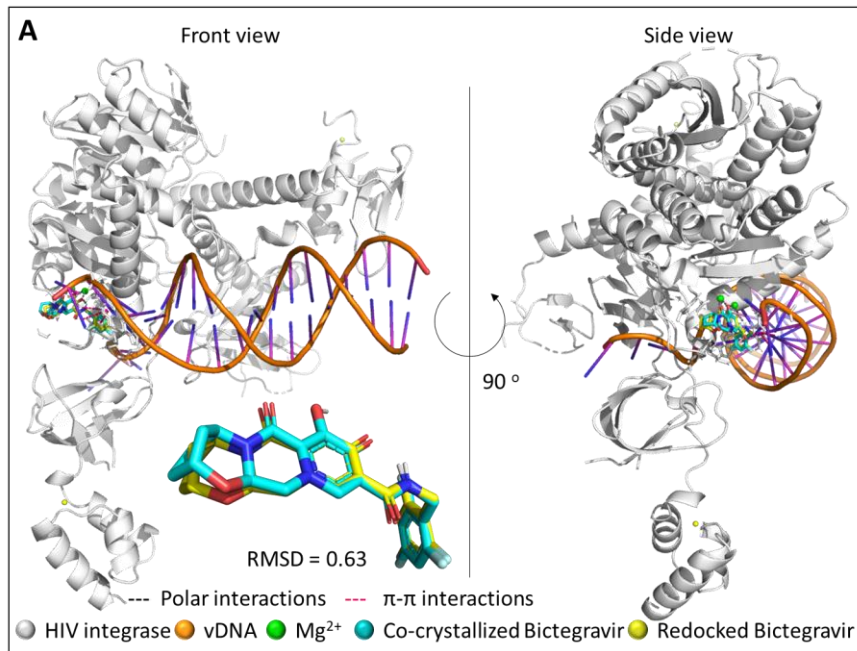
Software selection





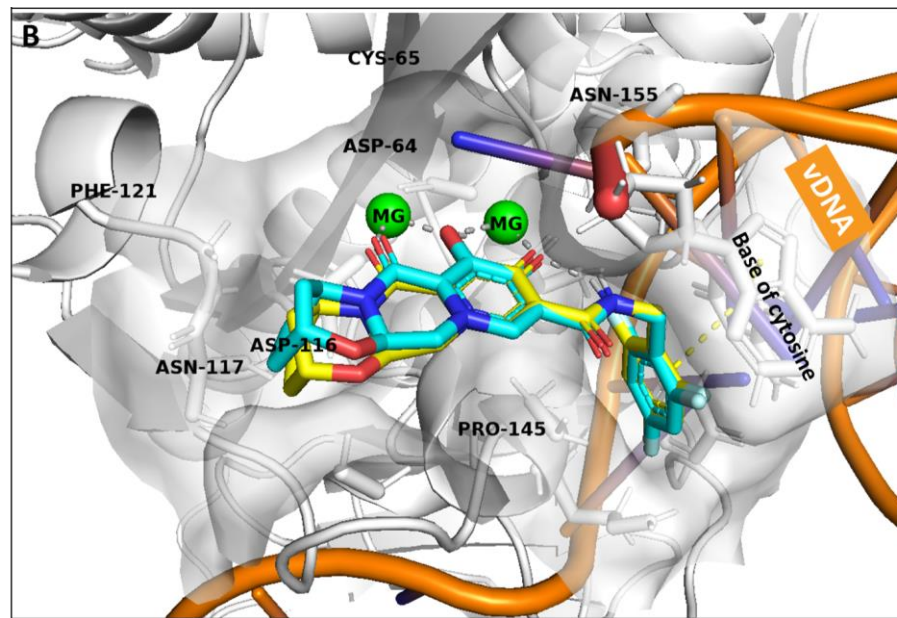
# RESULT

## Autodock-GPU



# MOLECULAR DOCKING

## Re-docking

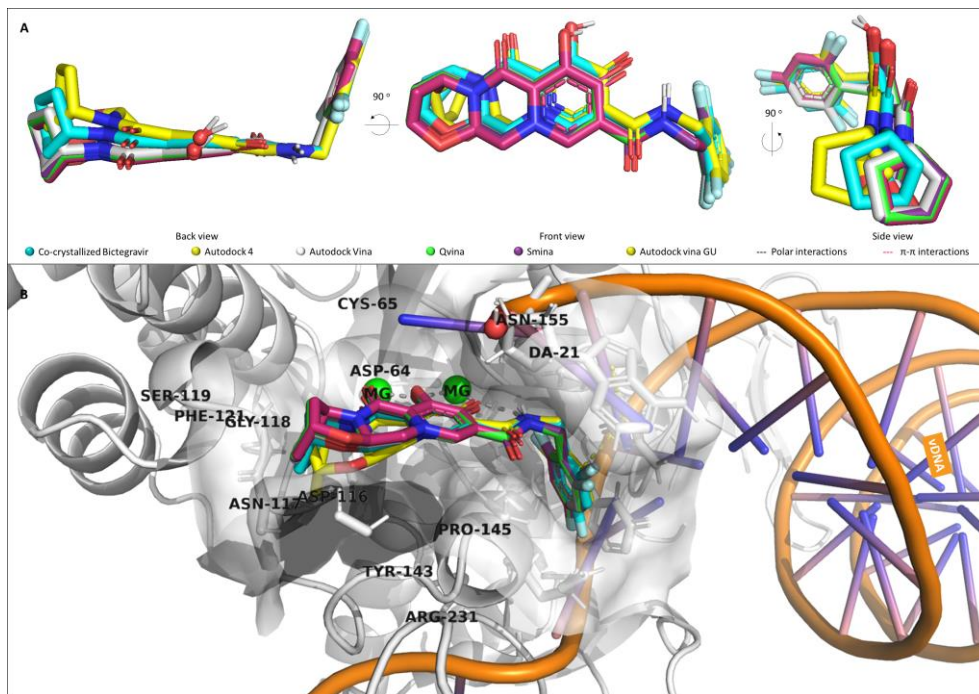




# RESULT

# MOLECULAR DOCKING

## Re-docking



Softwares	RMSD (Å)	Best docking score(kcal/mol)
<b>Smina</b>	0,677	-10,10
<b>Qvina2</b>	0,716	-9,30
<b>Autodock Vina 1.2.3</b>	0,636	-9,86
<b>Vina-GPU</b>	0,792	-9,80
<b>Autodock-GPU</b>	0,630	-11,11

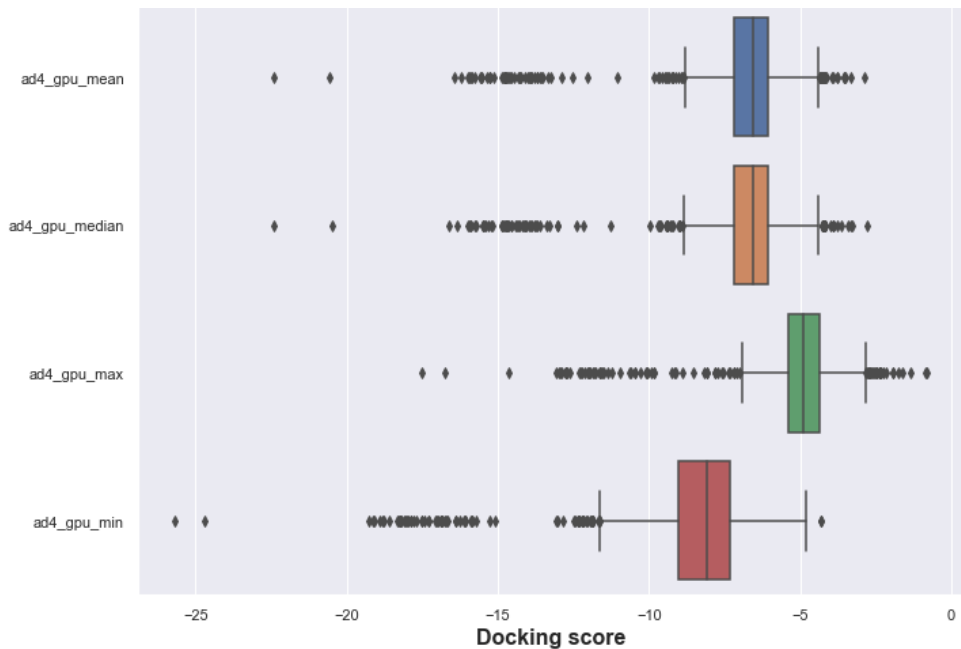
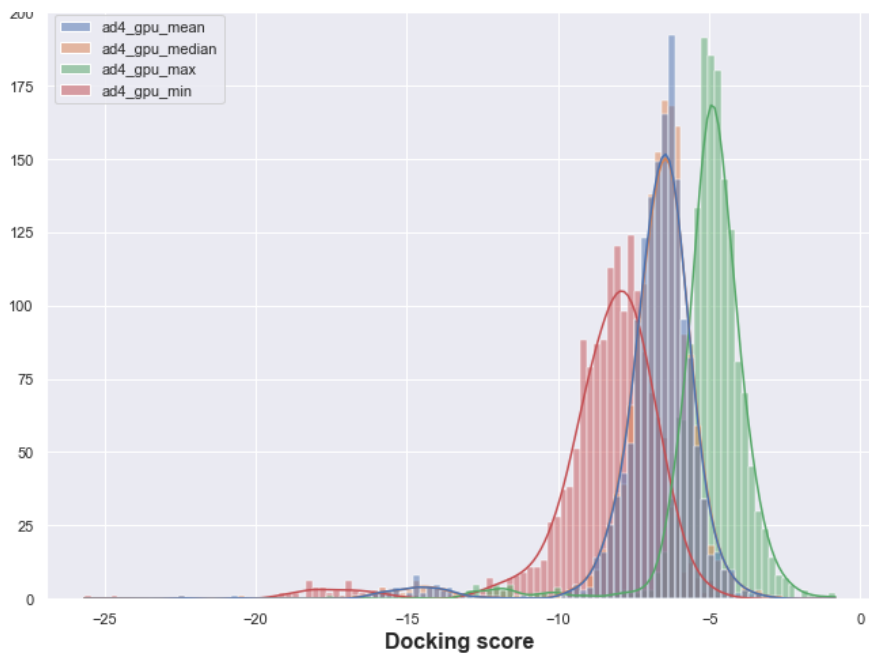


# RESULT

## Autodock-GPU

# MOLECULAR DOCKING

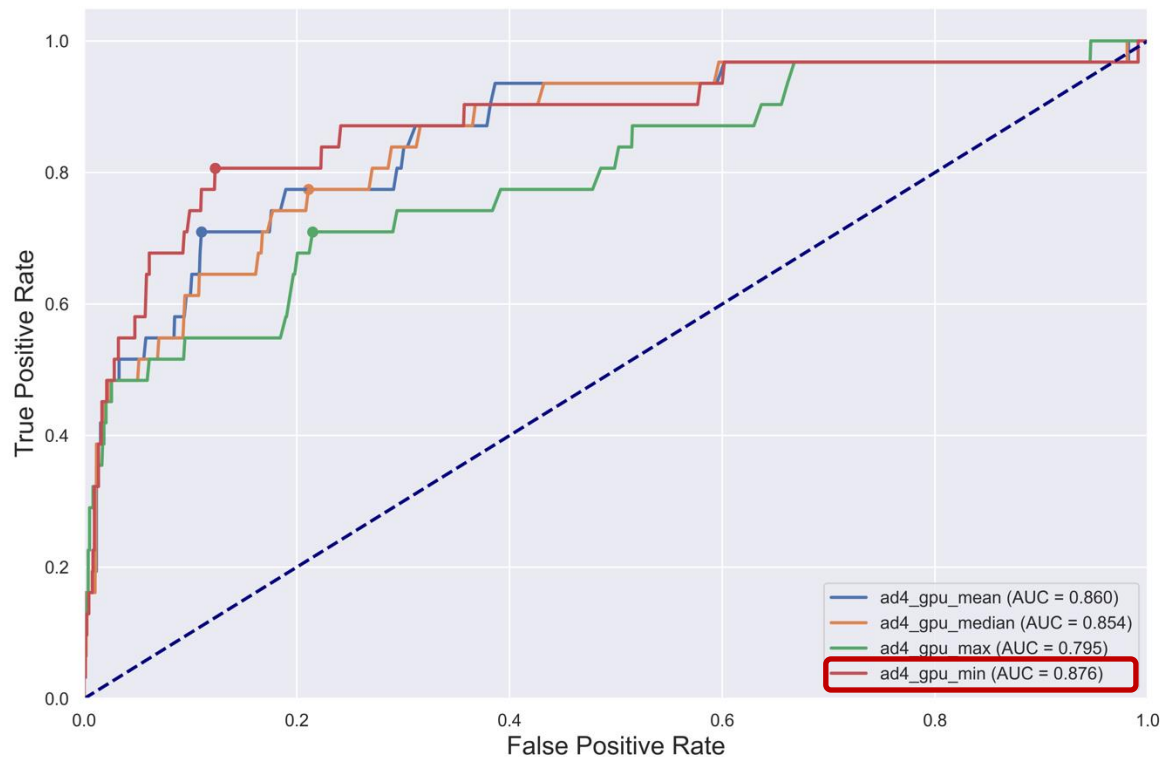
## Retrospective control





## RESULT

### Autodock-GPU



## MOLECULAR DOCKING

### Retrospective control



#### Autodock-GPU\_min:

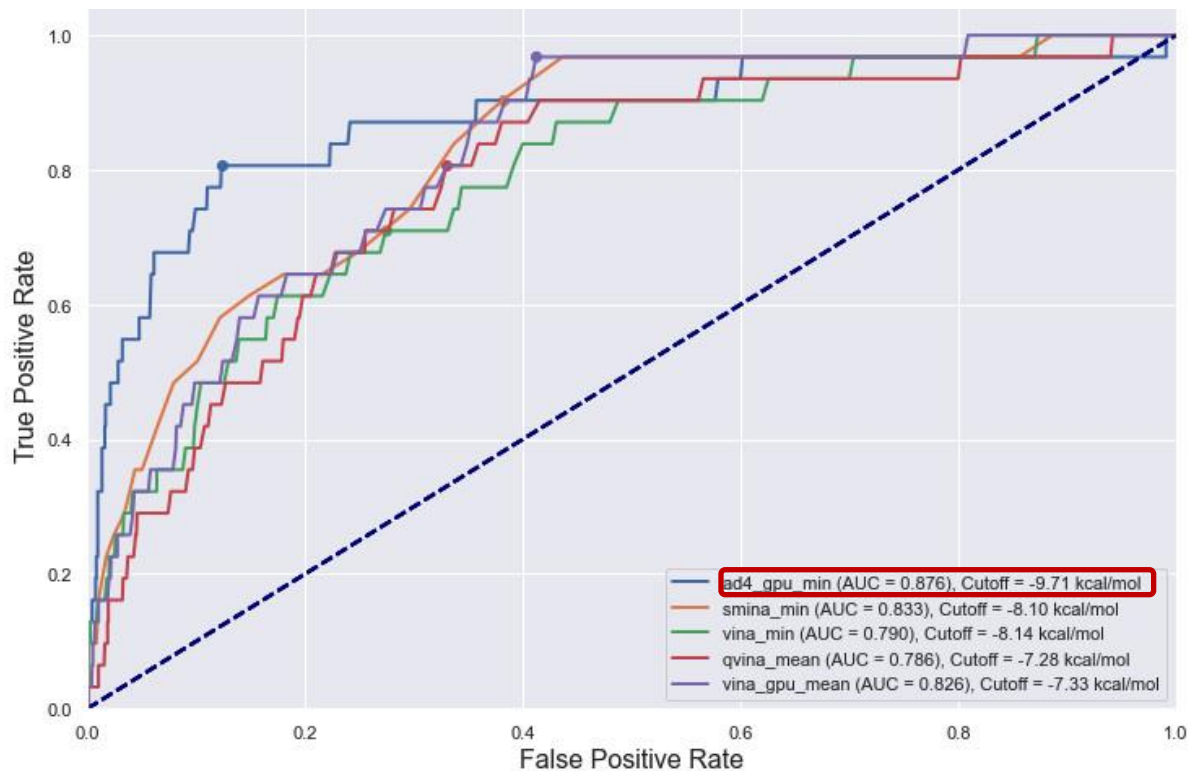
- AUC = 0,876
- G-mean max = 0,841
- TPR = 0,806
- FPR = 0,123
- cutoff = -9,71 kcal/mol





# RESULT

## Autodock-GPU



# MOLECULAR DOCKING

## Retrospective control



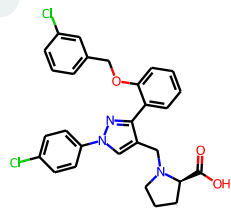
4

# Screening

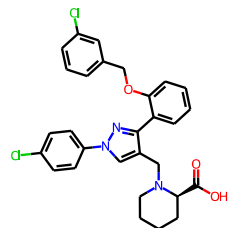




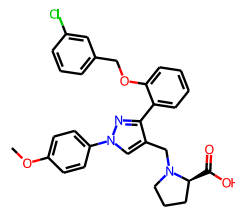
# RESULT



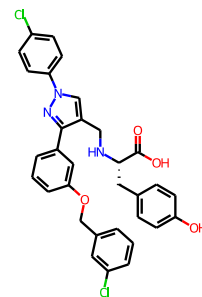
4.3.Pro  
pChEMBL = 7,00



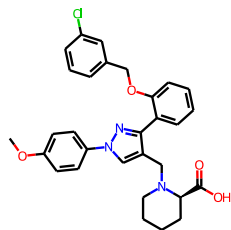
4.3.Pipe  
pChEMBL = 6,92



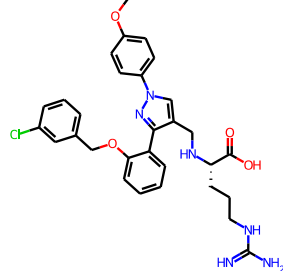
6.3.Pro  
pChEMBL = 7,13



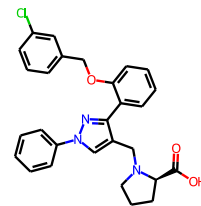
4.3.Tyr  
pChEMBL = 7,06



6.3.Pipe  
pChEMBL = 7,08



6.3.Arg  
pChEMBL = 7,40



1.3.Pro  
pChEMBL = 7,16

# VIRTUAL SCREENING

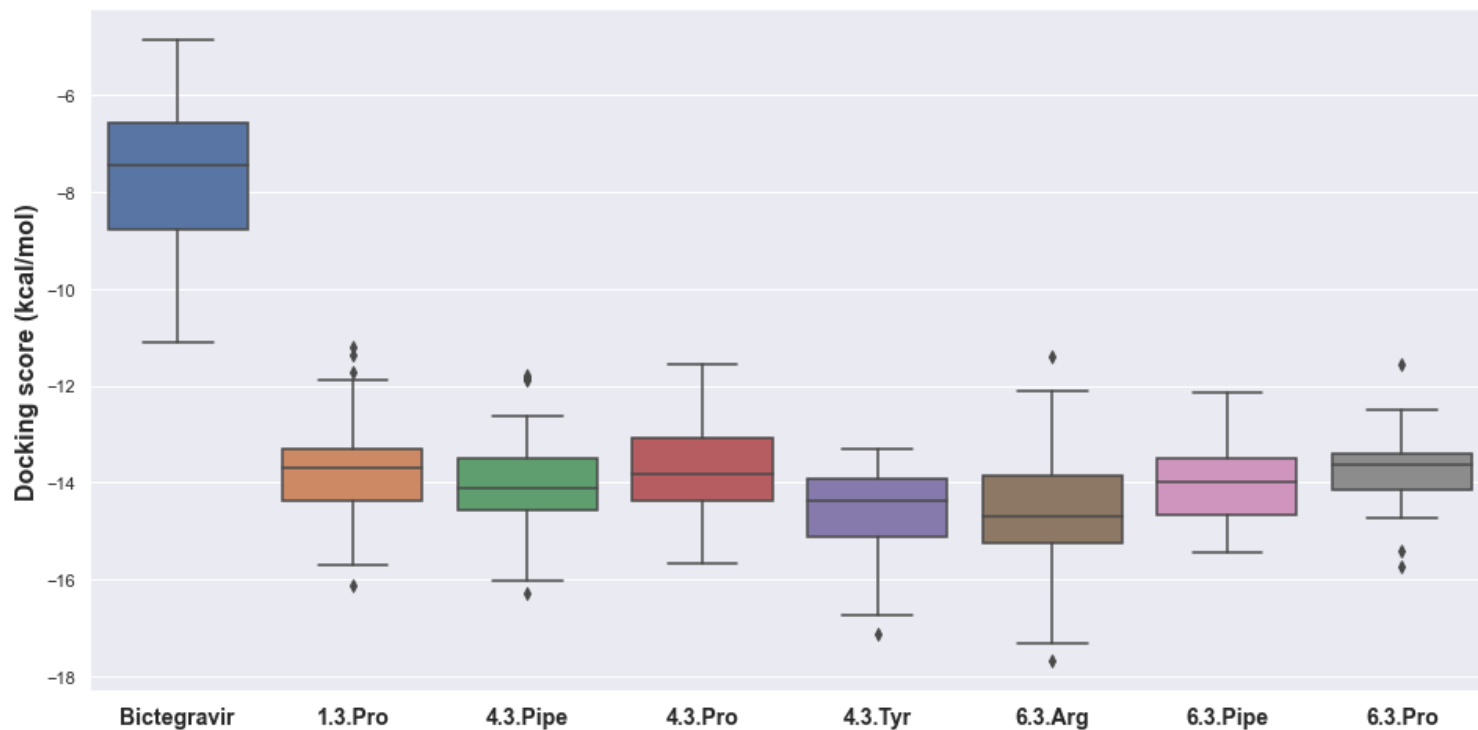




# RESULT

## VIRTUAL SCREENING

### Molecular Docking





## RESULT

## VIRTUAL SCREENING

### Molecular Docking

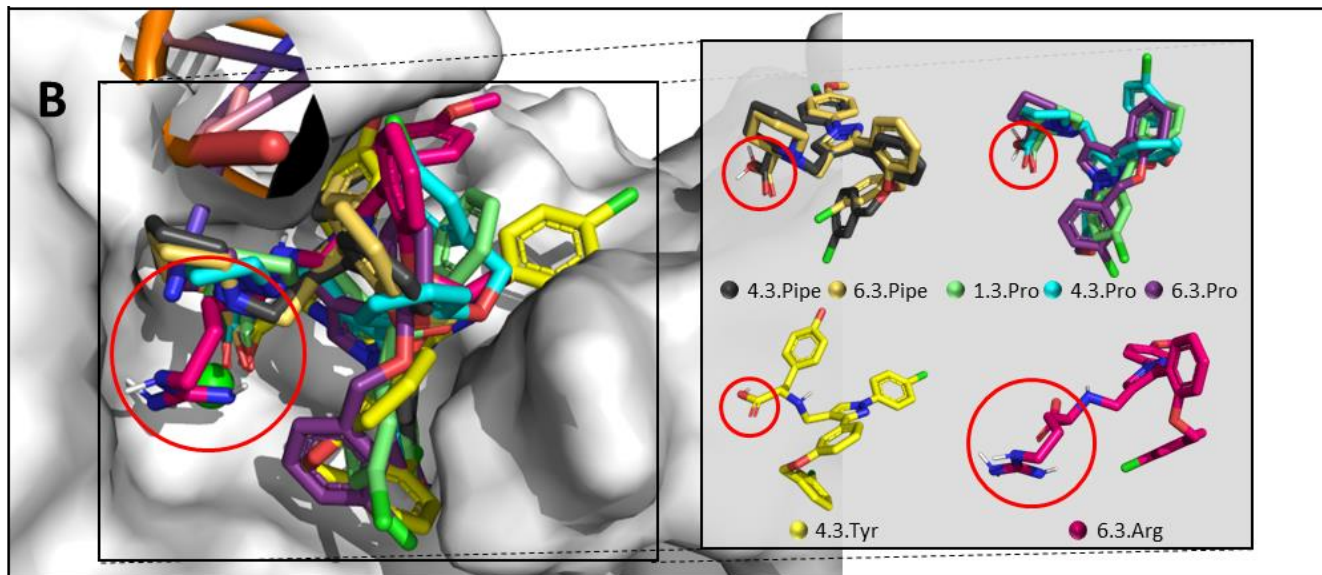


**Group 1:** 4.3.Pipe, 6.3.Pipe

**Group 2:** 1.3.Pro, 4.3.Pro, 6.3.Pro

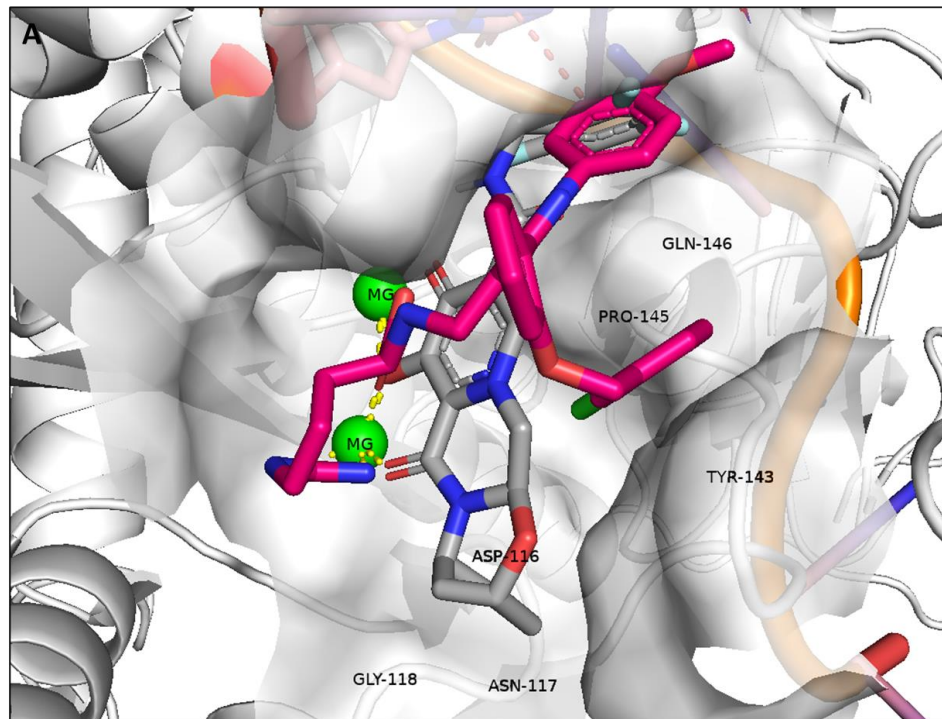
**Group 3:** 4.3.Tyr

**Group 4:** 6.3.Arg



# DISCUSSION

## Group 4



# VIRTUAL SCREENING

## Molecular Docking

