Master's thesis

# EMCIP

An Ensemble Model For Cdr1 Inhibitor Prediction

**Student:** The-Chuong Trinh
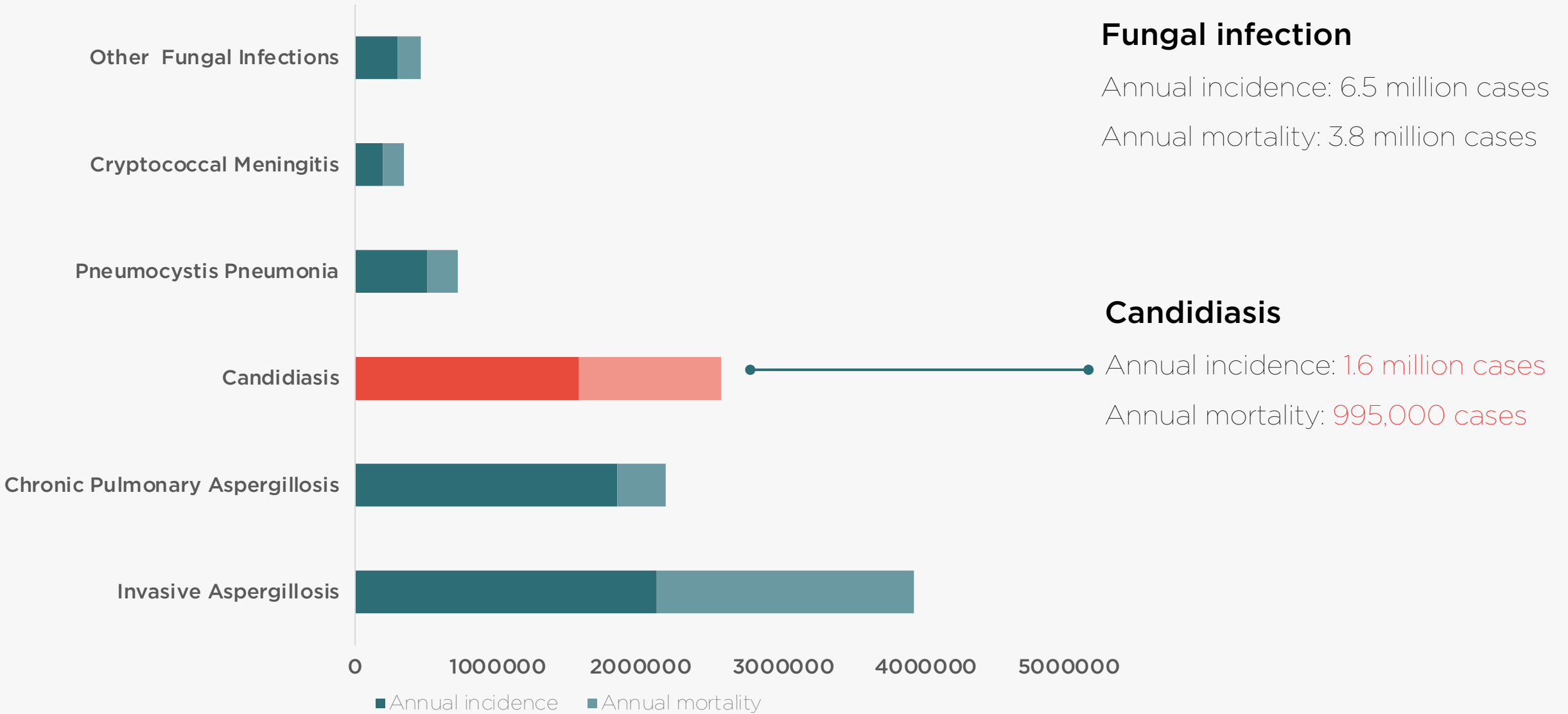
**Supervisor:** Prof. Ahcène Boumendjel

Laboratoire Radiopharmaceutiques Biocliniques

UGA Université Grenoble Alpes

Inserm La science pour la santé From science to health

# 1. INTRODUCTION

# FUNGAL INFECTION

**Fungal infection**

Annual incidence: 6.5 million cases

Annual mortality: 3.8 million cases

**Candidiasis**

Annual incidence: 1.6 million cases

Annual mortality: 995,000 cases

Chart categories (top to bottom):
- Other Fungal Infections
- Cryptococcal Meningitis
- Pneumocystis Pneumonia
- Candidiasis
- Chronic Pulmonary Aspergillosis
- Invasive Aspergillosis

X-axis: 0, 1000000, 2000000, 3000000, 4000000, 5000000

Legend: Annual incidence, Annual mortality

David W Denning, *The Lancet Infectious Disease*, 2024

3

# CDR1 INHIBITORS

## Azoles

The most widely used for treating and preventing Candida infections



**Clotrimazole**



**Ketoconazole**



**Fluconazole**

## Cdr1 inhibitors



**Enniatin B**



**Beauvericin**



**Fk506**

## Cdr1 efflux pump

One key resistance mechanism of Candida



Extracellular

Cdr1  Cdr1

Plasma membrane

Intracellular

Azole drugs

# VIRTUAL SCREENING

## Virtual screening

Assess whether a compound is a good drug using computation models (Walter et al., 1998)



Compound → Virtual screening model → Prediction: Good! → Experiments

## Pros

- Much faster than experimental screening in wet labs

- Test $10^8$ compounds within a day

- Much cheaper than experimental screening

# WORK PACKAGES

Collect and curate relevant compounds for Cdr1 transporters.

## DATA CURATION

Discover optimals machine learning algorithms and train a deep learning model to predict Cdr1 inhibitors

## AI MODEL DEVELOPEMENT

1

2

3

## MOLECULAR REPRESENTATION

Discover and design optimal molecular representation schemes for our models.

Synthesize/Buy prioritized compounds
and test their potency by bioassays.

**PROSPECTIVE SCREENING**

**VIRTUAL SCREENING**

Use model to screen
potential compounds
from a large library

# 2. METHODS

# QSAR
## Quantitative structure-activity relationship

**QSAR** is a mathematical model showing relationship

between **biological activity** and **molecular properties.**

$$Bio\_activity = f(D_1, D_2, D_3,..., D_n)$$

1 — Bemis Murcko Scaffold Split
2 — Molecular Featurize
3 — Train Model
4 — Evaluate Model
5 — Make Prediction

Original dataset

External test dataset

Training dataset

QSAR model

New dataset

Good model

# MOLECULAR REPRESENTATION

CC(=O)NC1=CC=C(C=C1)O

SMILES string

2D Molecule

Molecular graph

## Molecular fingerprints (Avalon, RDKit)

| | | | | | | |
|---|---|---|---|---|---|---|
| 1 | 0 | 1 | ... | ... | ... | 1 |

2D Molecule

Substructure 1

Substructure 2

Embedding vector

Machine learning algorithms

Prediction

## Molecular descriptors (Mordred)

2D Molecule

**Molecular weights**
**Atom types**
**H-bond donnors**
**....**

Molecular descriptors

151.2
12
1
...
...
...
...
...
2

Embedding vector

Machine learning algorithms

Prediction

# MESSAGE PASSING

**Halicin molecular graph**

**Halicin structures**



**- Message passing mechanism**

$$m_v^{t+1} = \sum_{w \epsilon N(v)} M_t(h_v^t, h_w^t, e_{vw}) \quad (1)$$

$$h_v^{t+1} = U_t(h_v^t, m_v^{t+1}) \quad (2)$$

$$h_G = \mathrm{R}(\{h_v^T \mid v \epsilon G\}) \quad (3)$$

- **Node features:** atomic number, number of bonds for each atom, formal charge, chirality, number of bonded hydrogens, hybridization, aromaticity, atomic mass.

- **Edge features:** bond type (single/double/triple/aromatic), conjugation, ring membership, stereochemistry.

# Multi-instance 3D Graph neural network



**GNN encoder**

**DNN**

Conformation 1

$\mathcal{G}_1 = (\mathcal{V}, \mathcal{E}, \mathcal{R})$

Conformation 2

$\mathcal{G}_2 = (\mathcal{V}, \mathcal{E}, \mathcal{R})$

$x_i \in \mathbb{R}^{1 \times m}$

$e_{i,i} \in \mathbb{R}^{1 \times e}$

Conformation 3

$\mathcal{G}_3 = (\mathcal{V}, \mathcal{E}, \mathcal{R})$

**Message passing**

$\sum_i$

Prediction 1

Prediction 2

Prediction 3

$\sum_j$

**Final Prediction**

**GNN:** Graph neural network

$\sum_i$ = Sum || Max || Mean operation

**DNN:** Deep neural network

$\sum_j$ = Mean operation

13

**GRAPH CONVOLUTION BLOCK**

**FULLY CONNECTED BLOCK**

2D features

3D features

node features

edge features

Sum operation

Sum, Max, and Mean operation

Transformer Conv

Linear

Normalization

Self-Attention Graph Pooling

y

# EVALUATION METRICS

## Confusion matrix

**Prediction**

|  | Active | Inactive |
|---|---|---|
| **Active** | True positive (TP) | False negative (FN) |
| **Inactive** | False positive (FP) | True negative (TN) |

**Ground truth**

**Precision**
$$\frac{TP}{TP + FP}$$

**Recall**
$$\frac{TP}{TP + FN}$$

**Specificity**
$$\frac{TN}{FP + FN}$$

**False positive rate (FPR)**
$$\frac{FP}{FP + FN}$$

## 1. Average precision

The area under Precision Recall curve

## 2. F1-score

The harmonic mean of Precision and Recall

## 3. ROC-AUC

The area under the

Receiver operating characteristic (ROC) curve

## 4. Balanced accuracy

The average between Recall and Specificity

# 3. RESULTS

# DATASET

**Sources:** Public repository (PubChem, ChEMBL), Literatures, Chemical patents (US11174267B2)

**Partition method:** Bemis Murcko scaffold and Local outlier factor

**Active/Inactive ratio** in each subset: 1/4.5

**Active compounds**
130 compounds
**7.2%**

**Inactive compounds**
1684 compounds
**92.8%**

**Validation set (DL)**
106 compounds

**Training set**
423 compounds

**External test set**
177 compounds

**Hard test set**
1108 compounds

**Composition in
the original dataset**

**Composition in
the each subset**

**Inactive compounds outside
the chemical space (Outliers)**

# MOLECULAR REPRESENTATION

## Meta-analysis

**16** types of molecular representations



Boxplots comparing the BM 10-fold cross validation results based on average precision



Heat map illustrating the results of Wilcoxon signed-rank tests based on average precision

Boxplots comparing the BM 10-fold cross validation results
based on F1-score

Heat map illustrating the results of
Wilcoxon signed-rank tests based on F1-score

**Molecular fingerprints:** RDK5, RDK6, RDK7, Avalon, Gobbi Pharmacophore fingerprints

**Molecular descriptors:** Mordred descriptors

# MOLECULAR REPRESENTATION
## Model selection

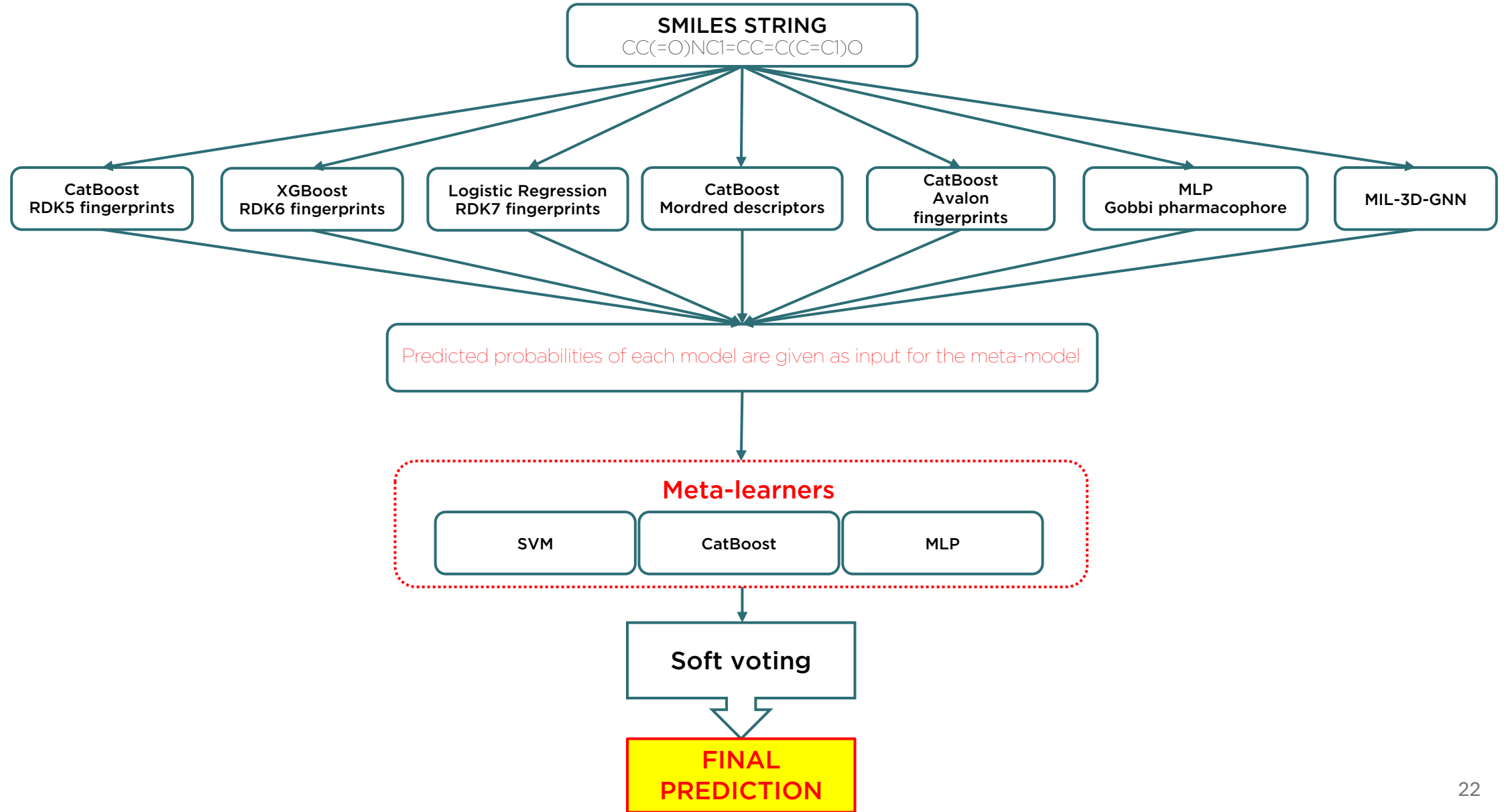**10** machine learning algorithms: Logistic regression, K-nearest neighbors, Support vector machine, Random forest, Extra tree, AdaBoost, Gradient boosting, XGBoost, CatBoost, and Multi-layer perceptron
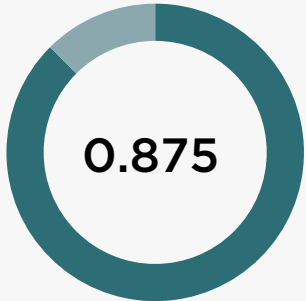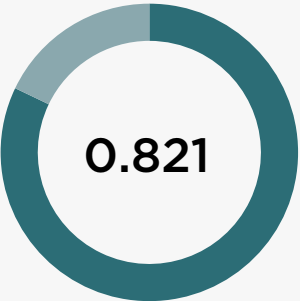


Yandex CatBoost

RDK6 fingerprints

Avalon fingerprints

Mordred descriptors

RDK5 fingerprints

RDK7 fingerprints

Ph4_Gobbi fingerprints

dmlc XGBoost

Yandex CatBoost

Yandex CatBoost

# MODEL PERFORMANCE

**Validation set**



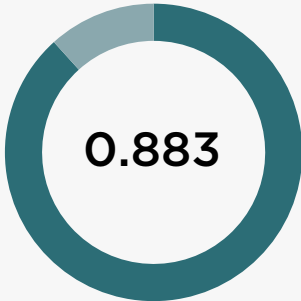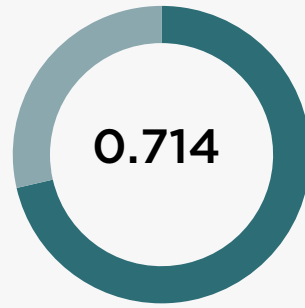| 0.875 | 0.821 | 0.906 | 0.883 |
|:---:|:---:|:---:|:---:|
| Average precision | F1-score | ROC-AUC | Balanced accuracy |

**Hard test set**

0.124

False positve rate

**External test set**

**0.755**

Average precision

**0.714**
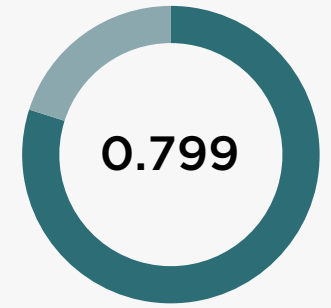
F1-score

**0.884**

ROC-AUC

**0.799**

Balanced accuracy

👍 The **generalizability** of the EMCIP model and its effectiveness on **unseen** data

# DEPLOYMENT

scikit learn   PyTorch   PyTorch geometric   Streamlit

# EMCIP GUI

## EMCIP: an Ensemble Model for Cdr1 Inhibitor Prediction

### Main Menu

▷ **Predict a batch**    ▷ Predict a molecule    ▷ About    ← ........ Batch prediction

## 1. Upload CSV File

Upload your file

☁ Drag and drop file here
Limit 200MB per file • CSV                                    [ Browse files ]  ← ........ Upload a csv file

📄 data_deploy.csv  0.8KB                                              ✕

Your file has 10 molecules

| | ID | Standardize_smile |
|---|---|---|
| 0 | 24818973 | Cc1ccc(OCc2cc(C(=O)NCC(C)(C)N3CCOCC3)no2)cc1C |
| 1 | spiroindolinone_24 | Cc1c(Cl)ccc2c1NC(=O)[C@@]21c2c(c(O)n(-c3ccccc3)c2O)[C@H]2CN(C(=O)NCc3ccccc3 |
| 2 | 44601952 | COc1cc(OC)cc(-c2sc3ccc(OC)cc3c2-c2ccc(OCCN3CCCCC3)cc2)c1 |
| 3 | 44601949 | COc1cccc(-c2sc3ccc(OC)cc3c2-c2ccc(OCCN3CCCCC3)cc2)c1 | ← ........ SMILES and ID table |
| 4 | N_ethylmaleimid | CCN1C(=O)C=CC1=O |
| 5 | spiroindolinone_13 | CC(C)(C)OC(=O)N1CCN2[C@H](C1)c1c(c(O)n(Cc3ccccc3)c1O)[C@@]21C(=O)N(Cc2ccc |
| 6 | spiroindolinone_9 | Cc1c(Cl)ccc2c1NC(=O)[C@@]21c2c(c(O)n(-c3ccccc3)c2O)[C@H]2CN(C(=O)OC(C)(C)C) |
| 7 | spiroindolinone_10 | Cc1c(Cl)ccc2c1NC(=O)[C@@]21c2c(c(O)n(-c3ccc(F)cc3)c2O)[C@H]2CN(C(=O)OC(C)(C) |
| 8 | spiroindolinone_11 | Cc1c(Cl)ccc2c1NC(=O)[C@@]21c2c(c(O)n(-c3cc(F)cc(F)c3)c2O)[C@H]2CN(C(=O)OC(C) |
| 9 | spiroindolinone_18 | Cc1cccc2c1NC(=O)[C@@]21c2c(c(O)n(-c3ccccc3)c2O)[C@H]2CN(C(=O)OC(C)(C)C)CCN |

25

# EMCIP: an Ensemble Model for Cdr1 Inhibitor Prediction

**Main Menu**

▷ **Predict a batch**　　▷ Predict a molecule　　▷ About

## 3. Cdr1 Inhibitors Prediction

Predict ◁········································· **Predict molecules**

Completed

Processing 3D Graph neural network: 10/10

| | Standardized SMILES | Probability | Prediction |
|---|---|---|---|
| 5 | CC(C)(C)OC(=O)N1CCN2[C@H](C1)c1c(c(O)n(Cc3ccccc3)c1O)[C@@]21C(=O)N(Cc2cccc | 0.9958 | 1 |
| 9 | Cc1cccc2c1NC(=O)[C@@]21c2c(c(O)n(-c3ccccc3)c2O)[C@H]2CN(C(=O)OC(C)(C)C)CCN | 0.9955 | 1 |
| 1 | Cc1c(Cl)ccc2c1NC(=O)[C@@]21c2c(c(O)n(-c3ccccc3)c2O)[C@H]2CN(C(=O)NCc3ccccc3 | 0.9955 | 1 |
| 6 | Cc1c(Cl)ccc2c1NC(=O)[C@@]21c2c(c(O)n(-c3ccccc3)c2O)[C@H]2CN(C(=O)OC(C)(C)C)( | 0.9954 | 1 |
| 7 | Cc1c(Cl)ccc2c1NC(=O)[C@@]21c2c(c(O)n(-c3ccc(F)cc3)c2O)[C@H]2CN(C(=O)OC(C)(C) | 0.9949 | 1 |
| 8 | Cc1c(Cl)ccc2c1NC(=O)[C@@]21c2c(c(O)n(-c3cc(F)cc(F)c3)c2O)[C@H]2CN(C(=O)OC(C) | 0.9948 | 1 |
| 2 | COc1cc(OC)cc(-c2sc3ccc(OC)cc3c2-c2ccc(OCCN3CCCCC3)cc2)c1 | 0.9864 | 1 |
| 3 | COc1cccc(-c2sc3ccc(OC)cc3c2-c2ccc(OCCN3CCCCC3)cc2)c1 | 0.979 | 1 |
| 4 | CCN1C(=O)C=CC1=O | 0.0206 | 0 |
| 0 | Cc1ccc(OCc2cc(C(=O)NCC(C)(C)N3CCOCC3)no2)cc1C | 0.0035 | 0 |

**Probability of being a Cdr1 inhibitor** ·········· (Probability / Prediction columns)

**Result table** ·········· (points to result table)

Successfully predicted your data

# EMCIP GUI

## EMCIP: an Ensemble Model for Cdr1 Inhibitor Prediction

**Main Menu**

▷ Predict a batch     ▷ **Predict a molecule**     ▷ About ·········· **Predict a molecules**

### 1. Input SMILES

Please, input your SMILES

COc1cc(C=CC(=O)CC(=O)C=Cc2ccc(O)c(OC)c2)ccc1O ◁········· **Input a SMILES string**

### 2. Cdr1 Inhibitor Prediction

Predict ◁·········· **Predict**

Standardization completed.

Graph dataset created...

Completed

Processing 3D Graph neural network: 1/1

28

**EMCIP**

**GUI**

# EMCIP: an Ensemble Model for Cdr1 Inhibitor Prediction

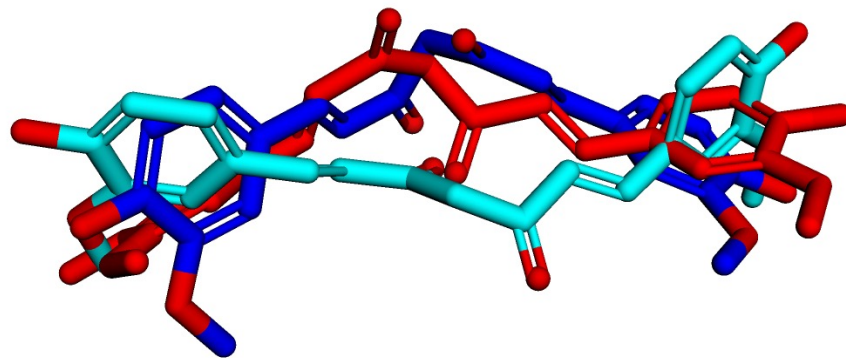## Main Menu

▷ Predict a batch    ▷ **Predict a molecule**    ▷ About

Number of displayed conformations

3

1    50

These are the first 3 conformations of your molecule

**Generated conformations**

The probability of your molecule to be a CDR1 inhibitor is: 0.9937756190299815

**Probability of being a Cdr1 inhibitor**

Restart

# 4. DISCUSSION

## ❖ Conclusion

- **EMCIP:** The first ensemble machine learning model specific for Cdr1 inhibitor prediction.

- **MIL-3D-GNN:** A novel 3D graph neural network for multi-instance learning.

- **Promising results on the external test set and validation set,** demonstrating the generalizability of the EMCIP model and its effectiveness on unseen data.

- **A practical GUI for EMCIP**, making it accessible and user-friendly even to AI non-experts.

- **A practical workflow**, conducting **ligand-based predictive AI models for other targets.**

## ❖ Limitations

- **Data scarcity**: Test more compounds to augment the dataset.

- **Lack of experimental structure of Cdr1 protein**: Prevents integration of protein information and implementation of structure-based drug design.

# THANK YOU FOR YOUR ATTENTION

**Email:** the-chuong.trinh@etu.univ-grenoble-alpes.fr

EMCIP Online Version