

## Deep Transfer Learning-based Hit Compound Classification

## for Therapeutic Targets in Heart Failure Drug Discovery on Small Datasets

#### Parajaree Ungudonpakdee, Thanasan Kumdee

11th grade, Kamnoetvidya Science Academy, Rayong, Thailand Scientific advisor 1: Thanasan Nilsu, Department of Biology and Environmental Science, Kamnoetvidya Science Academy (KVIS), Academic Teacher, Ph.D. (Applied Biological Sciences: Environmental Health) Scientific advisor 2: Bundit Boonyarit, School of Information Science and Technology, Vidyasirimedhi Institute of Science and Technology (VISTEC), Ph.D. student, Master of Science (Biochemistry)

## INTRODUCTION

#### HEART FAILURE





# **108 BILLION**

US dollars

One of the biggest death contributor & huge economic burden

## INTRODUCTION



Abnormal cardiac muscle function

Not enough O<sub>2</sub> for body Compensatory Mechanism



Worsen Heart Failure

## INTRODUCTION





The Process of Drug Discovery



## The Process of Drug Discovery



Machine Learning Model



#### OBJECTIVE



"To construct a deep transfer learning model which classifies active compound for specific receptor targets "





give rise to

**Compound Type** 

Chemical Structure

Classification based on Bioactivity Value



#### **Chemical Structure**

**Compound Type** 









#### Trained Machine Learning Model





#### Active or Inactive Compound

#### **Trained** Machine Learning Model

#### OBJECTIVE









ASSITS PHARMACEUTICAL INDUSTRY **REDUCE THE NUMBER OF HEART FAILURE PATIENTS** 

#### I. DATA PREPARATION

#### **II. FEATURE EXTRACTION**

## METHODOLOGY

#### III. MODEL TRAINING

#### **IV. PERFORMANCE EVALUATION**

## METHODOLOGY

#### I. DATA PREPARATION

FOCUSED RECEPTOR TARGETS	PRESENT NEGATIVE INOTROPES
Beta 1 Adrenergic receptor(ADBR1)	Beta blockers
Antiogensin-Converting enzyme(ACE)	ACE inhibitors
Mineralcorticoid receptor(MCR)	Diuretics

Table 1: Focused receptor targets and present negative inotropes

### CN1CCC[C@H]1c2cccnc2

SMILES Simplified Molecular – Input Line Entry System

Compound Dataset

#### **pX** value

pX Bioactivity Datapoints

**INPUT** 

## CN1CCC[C@H]1c2cccnc2

SMILES Simplified Molecular – Input Line Entry System





COMPOUND CLASSIFICATION

## METHODOLOGY

#### **II. FEATURE EXTRACTION**

#### METHODOLOGY II. FEATURE EXTRACTION



Diagram 1: FP2VEC technique

## METHODOLOGY

#### III. MODEL TRAINING

#### METHODOLOGY III. MODEL TRAINIG



Diagram 2: Model architecture

#### METHODOLOGY III. MODEL TRAINIG

#### N=1000nM dataset

N=500nM dataset

N=300nM dataset

N=100nM dataset

## Hyperparameter Tuning with hyperparameters:

Embedding vector (k): 50, 100, 150 Window size of filter (h): 3, 4, 5 Feature map size (n): 256, 512, 1024

Training

#### METHODOLOGY III. MODEL TRAINIG

81/100	9/100	1/100
Train	Validation	Test
dataset	dataset	dataset
Train	Calculate	Evaluate
model	loss	model

## METHODOLOGY

#### **IV. PERFORMANCE EVALUATION**

METHODOLOGY IV.PERFORMANCE EVALUATION

## ROC-AUC scores (Area under the receiver operating characteristic curve scores)



X (nM)	Target Protein Family	AUC	Average AUC
	GPCR Subfamily A3	0.940	
1000	Nuclear Receptor Subfamily 3	0.853	0.858
	GPCR Subfamily A17	0.780	
500	GPCR Subfamily A3	0.899	
	Nuclear Receptor Subfamily 3	0.859	0.838
	GPCR Subfamily A17	0.756	
300	GPCR Subfamily A3	0.890	
	Nuclear Receptor Subfamily 3	0.840	0.825
	GPCR Subfamily A17	0.745	
100	GPCR Subfamily A3	0.911	
	Nuclear Receptor Subfamily 3	0.859	0.847
	GPCR Subfamily A17	0.771	

Efficiency of multi-task model at k= 50, h = 5, n = 1024

X (nM)	Target Protein Family	Multi-Task		Single Task	
		AUC	Average AUC	AUC	Average AUC
1000	GPCR Subfamily A3	0.940		0.923	
	Nuclear Receptor Subfamily 3	0.853	0.858	0.824	0.849
	GPCR Subfamily A17	0.780		0.802	
500	GPCR Subfamily A3	0.899		0.878	
	Nuclear Receptor Subfamily 3	0.859	0.838	0.829	0.819
	GPCR Subfamily A17	0.756		0.751	
300	GPCR Subfamily A3	0.890		0.846	
	Nuclear Receptor Subfamily 3	0.840	0.825	0.832	0.815
	GPCR Subfamily A17	0.745		0.767	
100	GPCR Subfamily A3	0.911		0.884	
	Nuclear Receptor Subfamily 3	0.859	0.847	0.803	0.818
	GPCR Subfamily A17	0.771		0.767	

Efficiency of multi-task & single task model at k= 50, h = 5, n = 1024



Comparison of Logarithm Loss for Each Task of

Multi-task Model (Above) & Single Task Model (Below)

#### **Best hyperparameters**

- Embedding vector (k): 50
- Window size of filter (h): 5
- Feature map size (n): 1024

## Best bioactivity criteria - X = 1000 nM

#### **Result AUC score**

- GPCR Subfamily A3 : 0.940
- Nuclear Receptor Subfamily 3 : 0.853
- GPCR Subfamily A17 : 0.780

## **FUTURE PLANS**



# **Thank You!**