# 利用蛋白質-配體交互作用與化合物結構為基礎之虛擬藥物篩選群集分析

學生：葛振寧　　　　　　　　　　　　　　　指導教授：楊進木 博士

國立交通大學生物資訊研究所碩士班

## 摘　　要

　　我們發展了一個針對虛擬藥物篩選後處理(post analysis)的兩階段階層式分群分析法。此方法利用蛋白質-配體交互作用與化合物結構做為兩階段分析的主要原則。在第一階段，篩選出的候選化合物與目標蛋白質之蛋白質-配體三維結構與交互作用資訊將轉換成一維的實數表示，並採用階層式分群法針對候選化合物做第一階段的分群。在第二階段中，我們以 atom-pair 一維結構分析轉換法，淬取第一階段之分群的分子拓樸結構資訊。每一個經過交互作用分群後的群集將再進一步根據結構相似度做細分。兩階段（交互作用與藥物結構）階層式分群分析用在虛擬藥物篩選結果之組織化與視覺化分析，可以提升分析的速率與命中率，節省時間與經費，並且有助於未來實驗測試藥物的挑選與進一步分析。本方法以一組具有五種不同分子藥物目標的資料做驗證，包含胸腺嘧啶激酶(thymidine kinase)抑制劑，二氫葉酸還原酶(dihydrofolate reductase)抑制劑，雌激素受體(estrogen receptor)促進劑，雌激素受體抑制劑與神經胺酸酶(neuraminidase)抑制劑。經過在這些重要的分子藥物目標的分群分析測試後，本方法可以提供訂定分群界線之可能參考值，並能幫助研究人員有效的從虛擬篩選後產生的大量資料中找出具代表性的測試候選藥物，減少時間與金錢的花費。除了上述五個重要藥物目標之外，我們的方法也實際應用到幽門螺旋桿菌之莽草酸激酶(*Helicobacter pylori* shikimate kinase, HpSK)的抑制劑篩選分析。在對 CMC 藥物資料庫的虛擬藥物篩選後，我們由前 300 名的可能藥物分子中，經兩階段階層式分群分析後選出 23 種具代表性的藥物結構。經過合作實驗室的酵素抑制性測試後發現五個實際測試的結構中有一個具有莽草酸激酶之抑制性。此結果證明我們的方法不僅對虛擬藥物篩選與分析有效，並且確實有助於提升先導藥物開發流程的篩選速度與命中率。

Cluster analysis of Structure-based Virtual Screening by Using Protein-ligand Interactions

and Compound Structures

Student : Cheng-Neng Ko                    Adviser : Dr. Jinn-Moon Yang

Institute of Bioinformatics
National Chiao Tung University

**ABSTRACT**

We developed a cluster analysis method for post analysis of structure-based virtual screening. The analysis was composed of two stages based on protein-ligand interactions and compound structures, respectively. The first stage was to generate a protein-ligand interaction cluster by translating 3D structural binding information from a protein-ligand complex into a 1D real number representation, and using hierarchical clustering method to preliminarily cluster our screening results. In the second stage, we extracted molecular topology by atom-pair representation of each compound to re-grouping the clusters derived from the first stage. Each interaction cluster could be further divided into sub-clusters according to their topological similarities. The two-staged cluster analysis could be used to organize, analyze, and visualize the data of virtual screening and mining the representative candidates for future biological test. We validated this method on data sets having five classes: thymidine kinase inhibitors, dihydrofolate reductase inhibitors, estrogen receptor agonist, estrogen receptor antagonists and neuraminidase inhibitors. Our method on these pharmaceutical interest targets provided a suggestion of cluster threshold and helped to mining diversely representative structures from large number of virtual screening data. Our method also has been applied on the practical inhibitor screening analysis for *Helicobacter pylori* shikimate kinase (HpSK). After virtual screening in CMC database, we selected compounds from top 300 and selected 23 representative candidates. Five of 23 representative candidates were tested *in vivo*, and one of the five candidates, furosemide, was identified being able to inhibit HpSK by cooperated laboratory of Dr. Wen-Ching Wang.

# Acknowledgements

# CONTENTS

# List of Tables

# List of Figures

the target protein TK (PDB id: 1kim). (b) Hierarchical clustering of protein-ligand interaction of 53 docked poses on TK (PDB id: 1kim). Each docked pose is represented as one line in the heat map in the middle of the figure, and the red being the lowest protein-ligand interaction energy and the green being the highest energy. The left side of the heat map shows the hierarchical clustering results on the TK, including the dendrogram. Docked poses in the heat map are rearranged according to the order given by hierarchical clustering marked by the black bar 'c' in the right side of the heat map .The hot spots identified from overlapping known active compounds were also shown in the top side of the heat map. (c) Overlay of docked poses of the cluster with most number of known active compounds, and shown the important hydrogen bonds between protein and ligand. (d) Overlay of docked poses of the cluster with most number of unknown compounds, and shown the important hydrogen bonds between protein and ligand. The blue frames in the heat map were the major interaction that

**Figure 24.** (d) The dendrogram of Hierarchical clustering of 61 known compound structures. The descriptor was calculated by atom-pair representation, using tanimoto coefficient for measuring distance between two molecules. Under the reference threshold (tanimoto coefficient = 0.55), There were three major clusters, a, b, and c. (a) In the cluster a, all 10 ER $\alpha$ agonists were grouped within the cluster. (b) In the cluster b, all 11 ER $\alpha$ antagonists were also grouped within the cluster. (c) In the cluster c, all 10 TK inhibitors and 14 NA were