



# 医薬品の標的分子や副作用 の予測における機械学習

山西芳裕

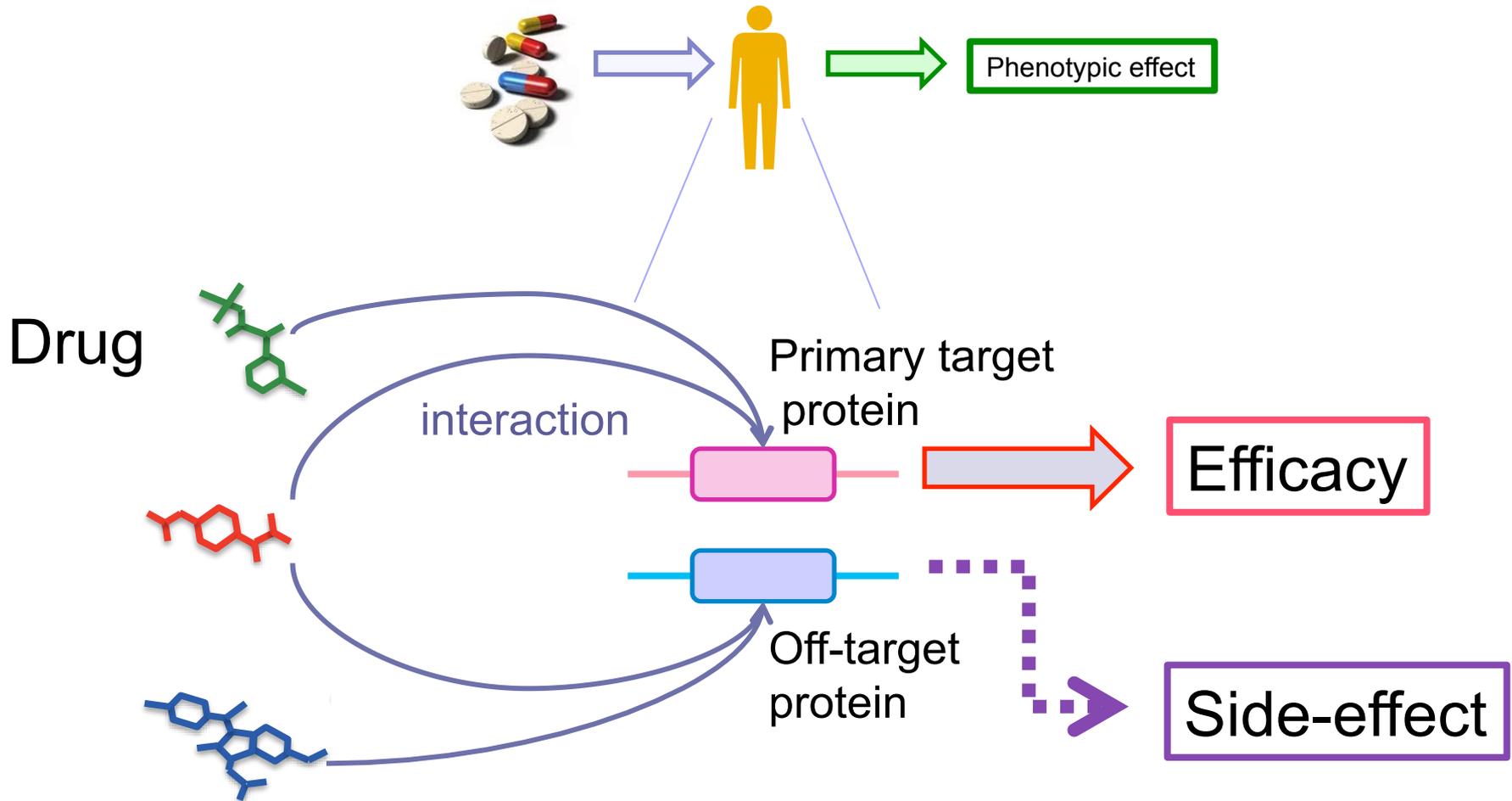
九州大学生体防御医学研究所



# Outline

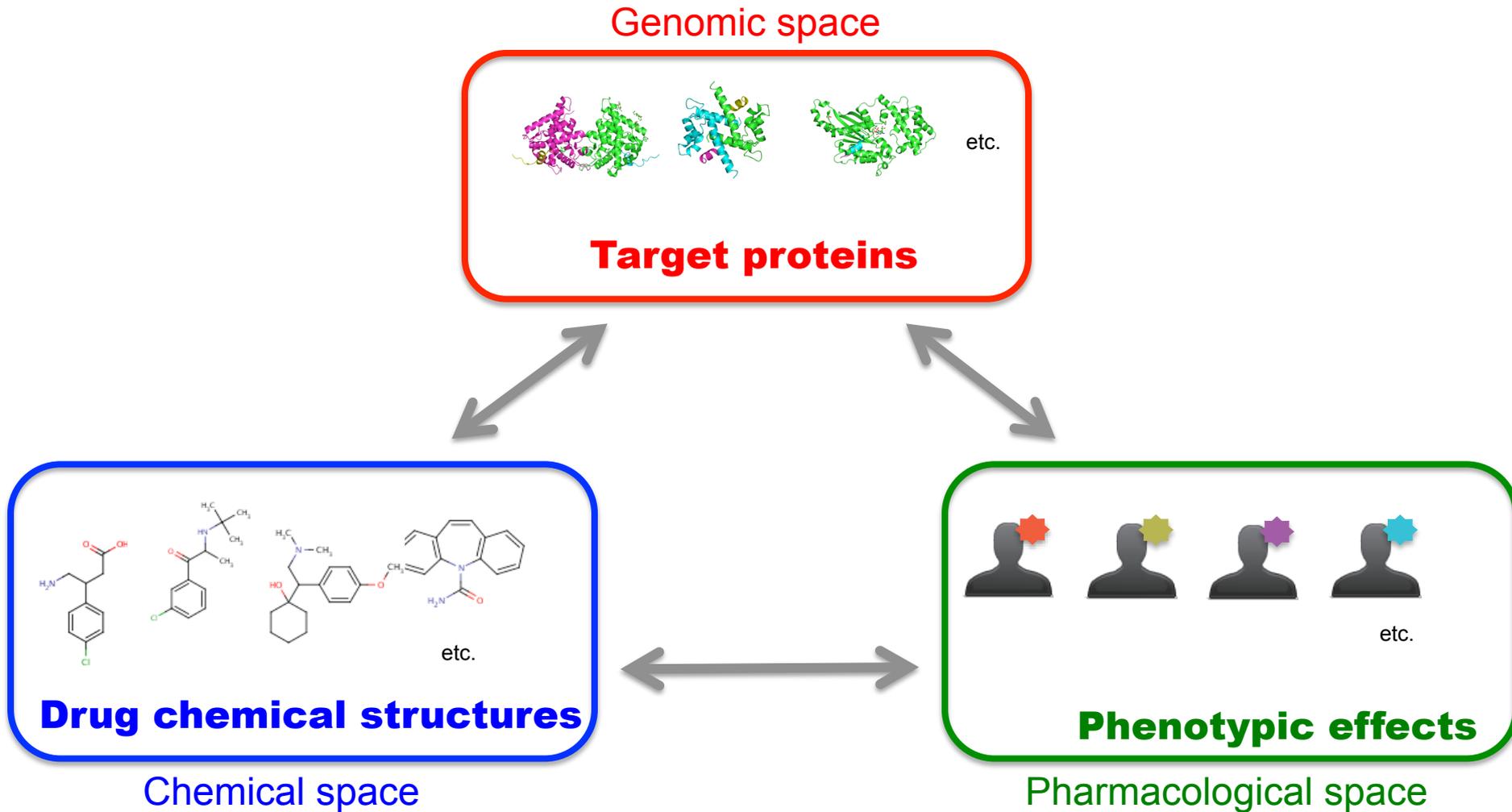
- Background
  - Chemical, genomic and pharmacological spaces
- Methods for pharmaceutical applications
  - Drug target prediction from chemical and genomic data
  - Side-effect prediction from biological data
- Results
- Concluding remarks

# Drug-target interaction



Identification of interactions between drugs and target proteins is crucial in the drug development

# Heterogeneous omics data

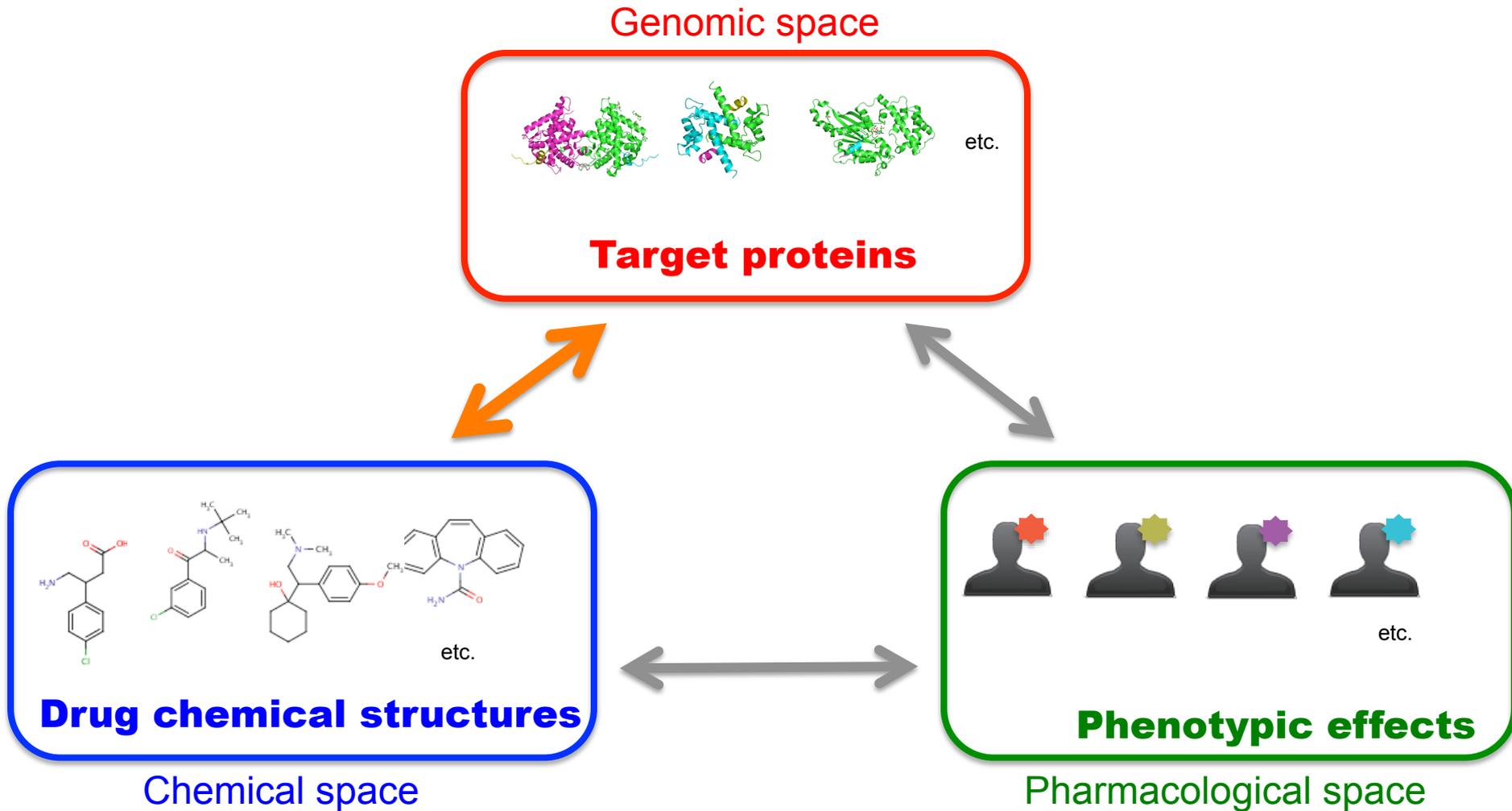




# Outline

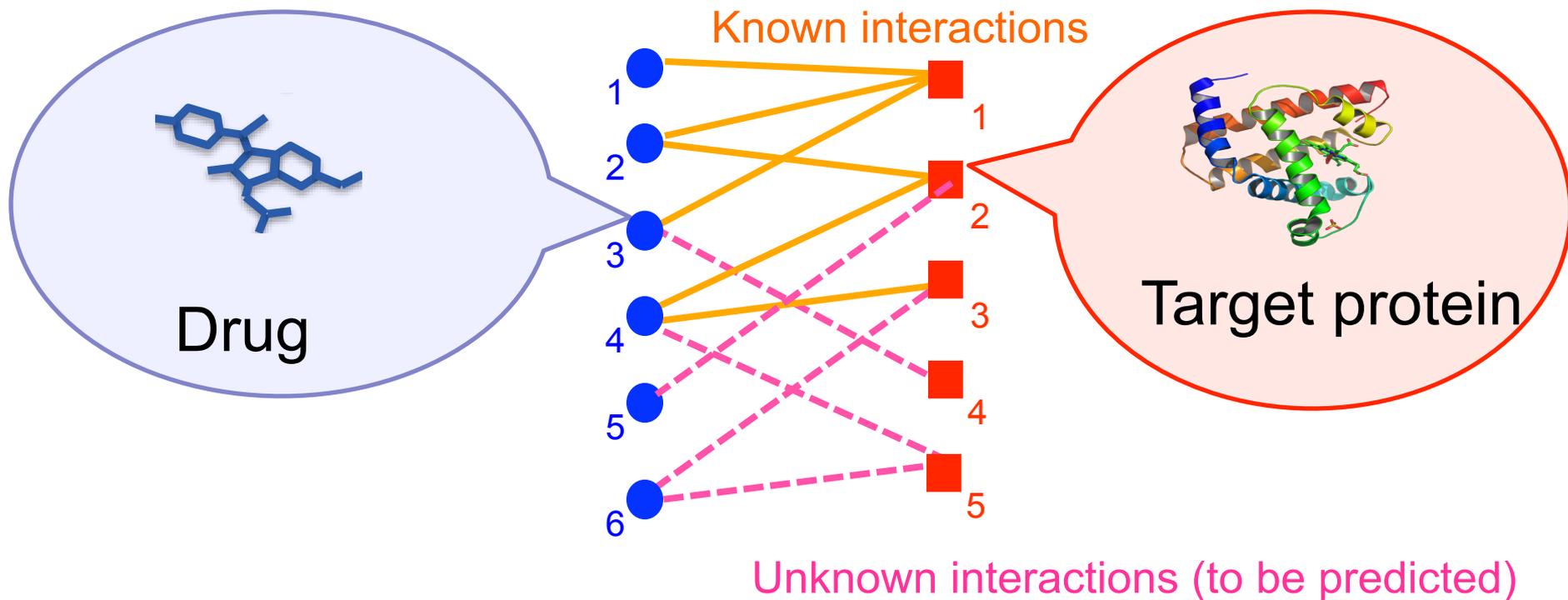
- Background
  - Chemical, genomic and pharmacological spaces
- Methods for pharmaceutical applications
  - Drug target prediction from chemical and genomic data
  - Side-effect prediction from biological data
- Results
- Concluding remarks

# Heterogeneous omics data



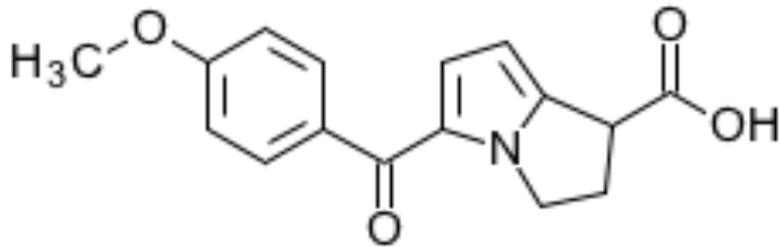
# Objective

- Prediction of unknown drug-target interactions on a large scale from chemical and genomic data



# Examples of the data structures

Drug



D02944

Target protein

```
MAHAAQVGLQDATSPIMEELITFHDHALMIIFLICFLVLYA  
LFLTLLTTKLTNTNISDAQEMETVWTILPAIILVLIALPSLRIL  
YMTDEVNDPSLTIKSIGHQWYWTYEYTDYGGLIFNSYML  
PPLFLEPGDLRLLDVDNRVVLPIEAPIRMMITSQDVLHSW  
AVPTLGLKTD AIPGRLNQTTFTATRPGVYYGQCSEICGAN  
HSFMPIVLELIPLKIFEMGPVFTL
```

Chemical graph structure

Amino acid sequence

# Chemogenomic approach

Strategy: Chemically similar drugs are predicted to interact with similar target proteins

(Yamanishi et al, *Bioinformatics*, 2008; Faulon et al., *Bioinformatics*, 2008; Jacob et al, *Bioinformatics*, 2008, Yabuuchi et al, *Mol Sys Bio*, 2011)

Chemical structure similarity for drugs

is evaluated by a graph kernel: (Mahe et al, *J Chem Inf Model*, 2005)

$$(K_x)_{ij} = k_x(\mathbf{x}_i, \mathbf{x}_j) \text{ for } i, j = 1, 2, \dots, n_x$$

Genomic sequence similarity for target proteins

is evaluated by a string kernel: (Saigo et al, *Bioinformatics*, 2004)

$$(K_z)_{ij} = k_z(\mathbf{z}_i, \mathbf{z}_j) \text{ for } i, j = 1, 2, \dots, n_z$$



# Binary classification approach

- Classification of drug-target pairs into the interaction class or non-interaction class

Support Vector Machine (SVM) with pairwise kernels

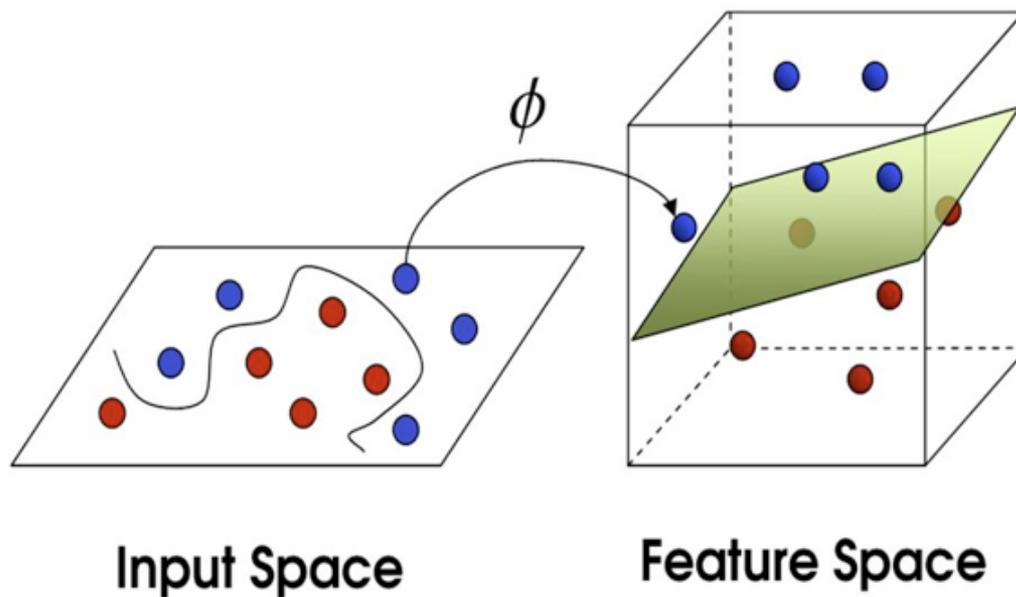
- Faulon et al., *Bioinformatics*, 24:225-233, 2008
- Jacob and Vert, *Bioinformatics*, 24:2149-2156, 2008
- Yabuuchi et al, *Mol Sys Bio*, 2011

# SVM with pairwise kernels (pairwise SVM)

ordinary SVM :

$$f(x) = \sum_{i=1}^n a_i k(x_i, x') + b$$

where  $x$  : an object



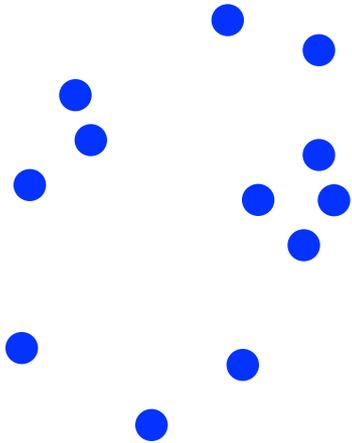
pairwise SVM :

$$f(x, z) = \sum_{i=1}^{n_x \times n_z} a_i k((x, z)_i, (x', z')) + b$$

where  $(x, z)$  : a compound – protein pair

# Supervised bipartite graph inference

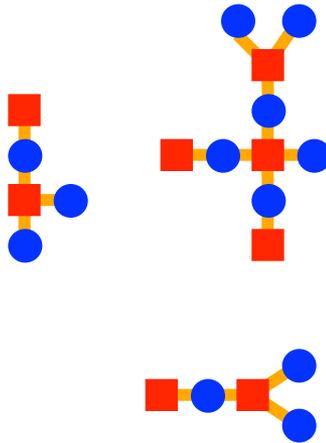
Chemical space



● known drug

Compounds with similar structures are close to each other

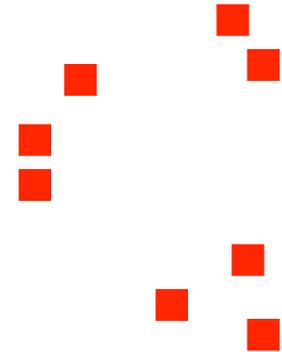
Interaction space



— known interaction

Interacting drugs and targets are connected on the graph

Genomic space

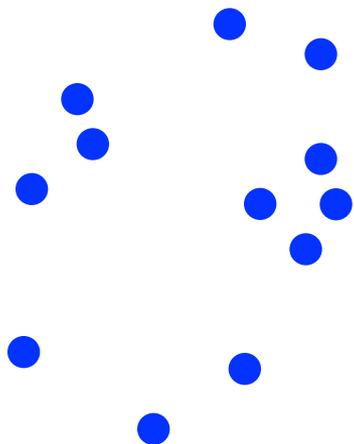


■ known target

Proteins with similar sequences are close to each other

# Step 1: embedding drugs and targets on the known graph into a unified feature space $\mathbf{R}^d$

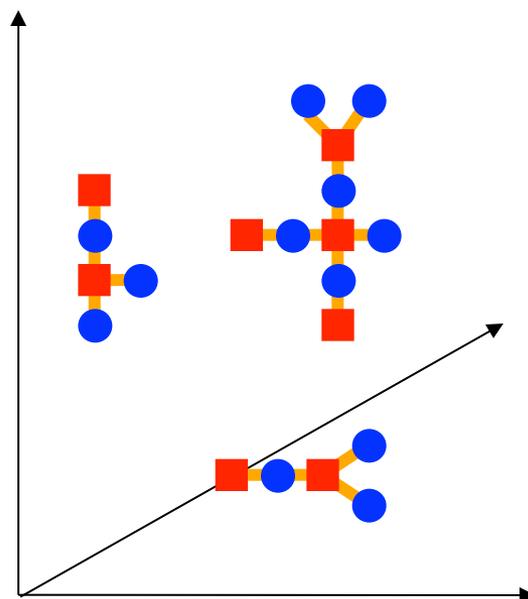
Chemical space



● known drug

Compounds with similar structures are close to each other

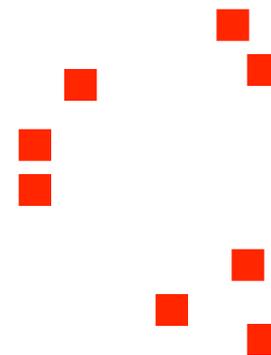
Interaction space



— known interaction

Interacting drugs and targets are close to each other

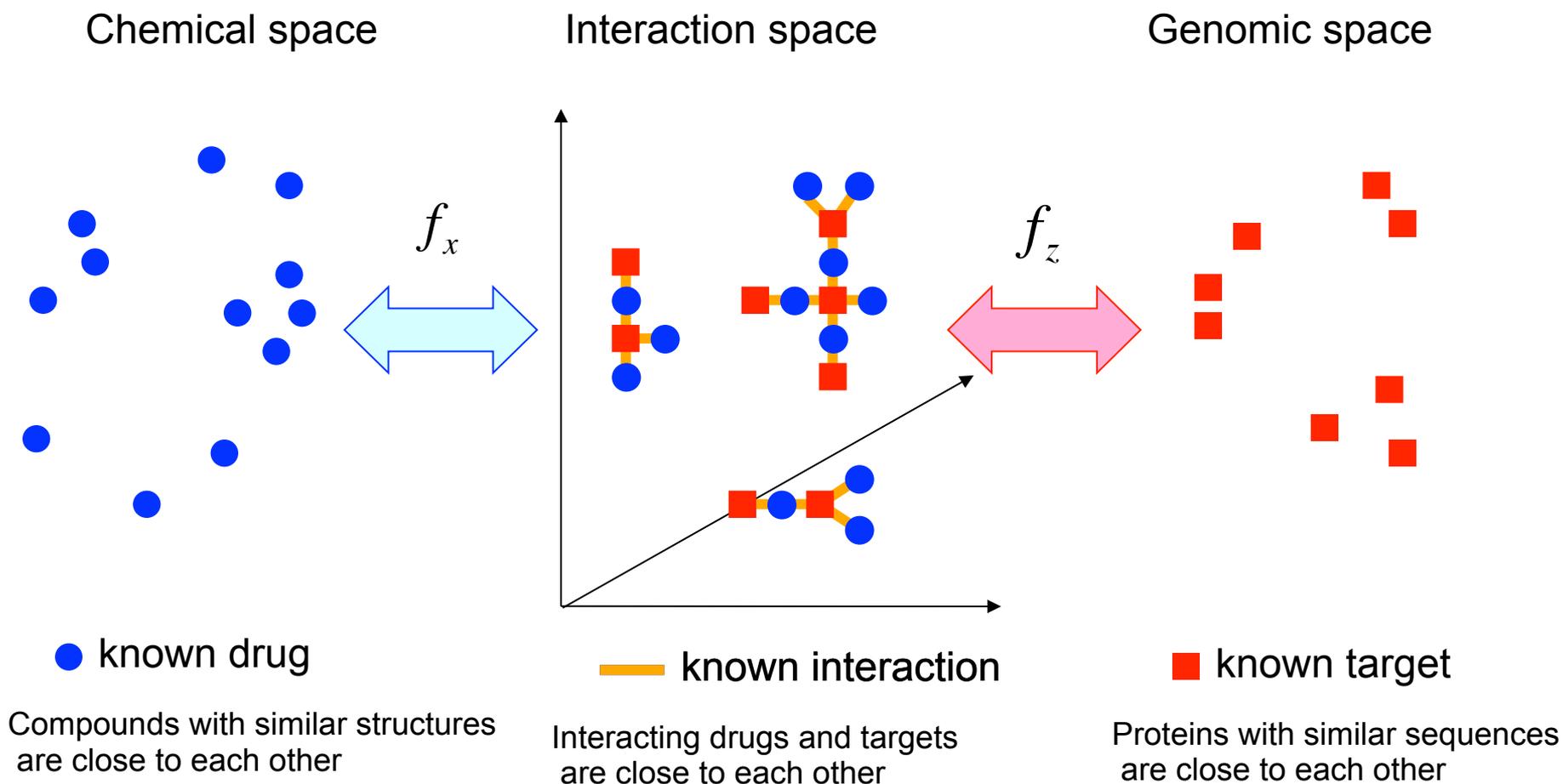
Genomic space



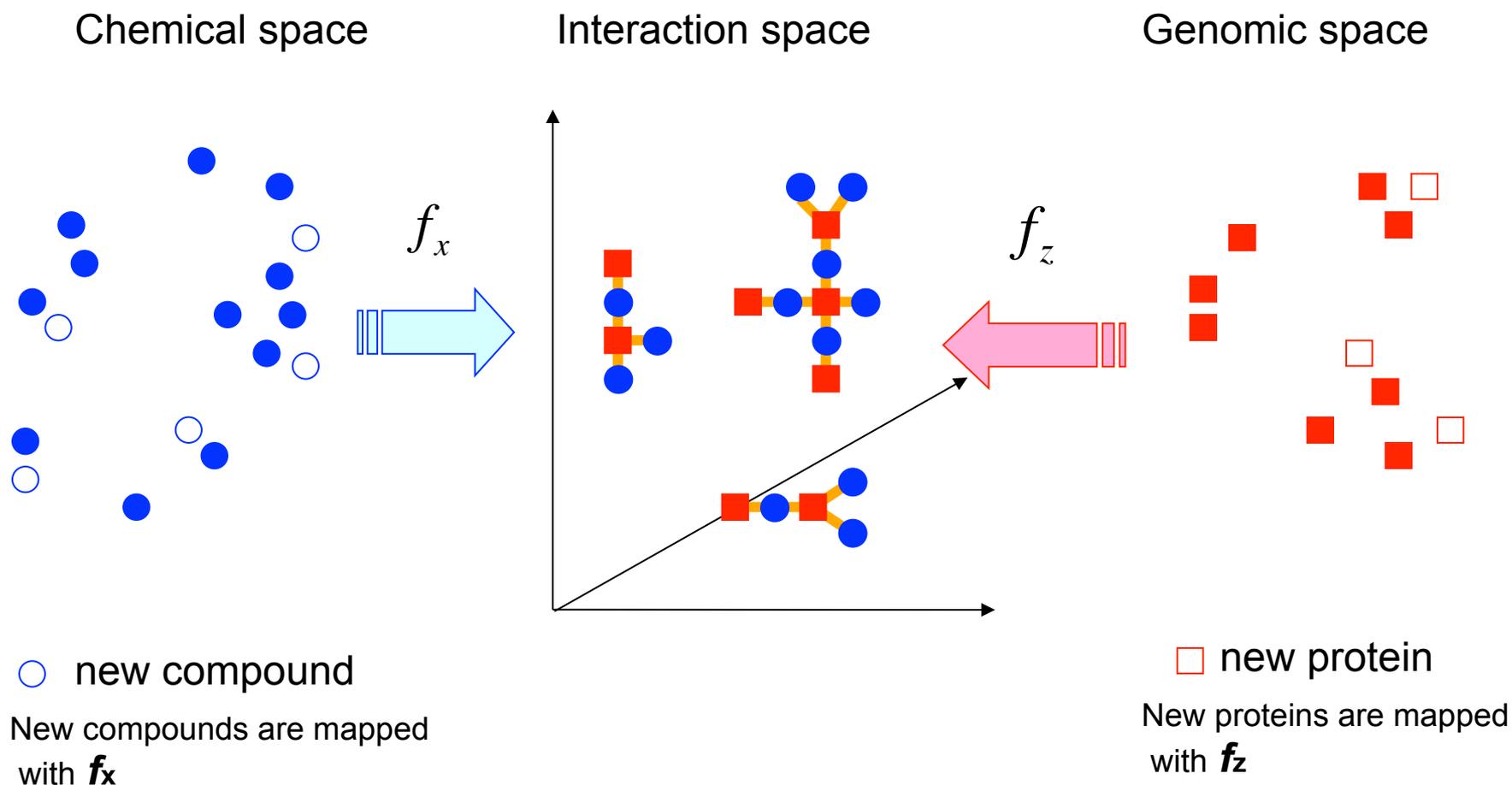
■ known target

Proteins with similar sequences are close to each other

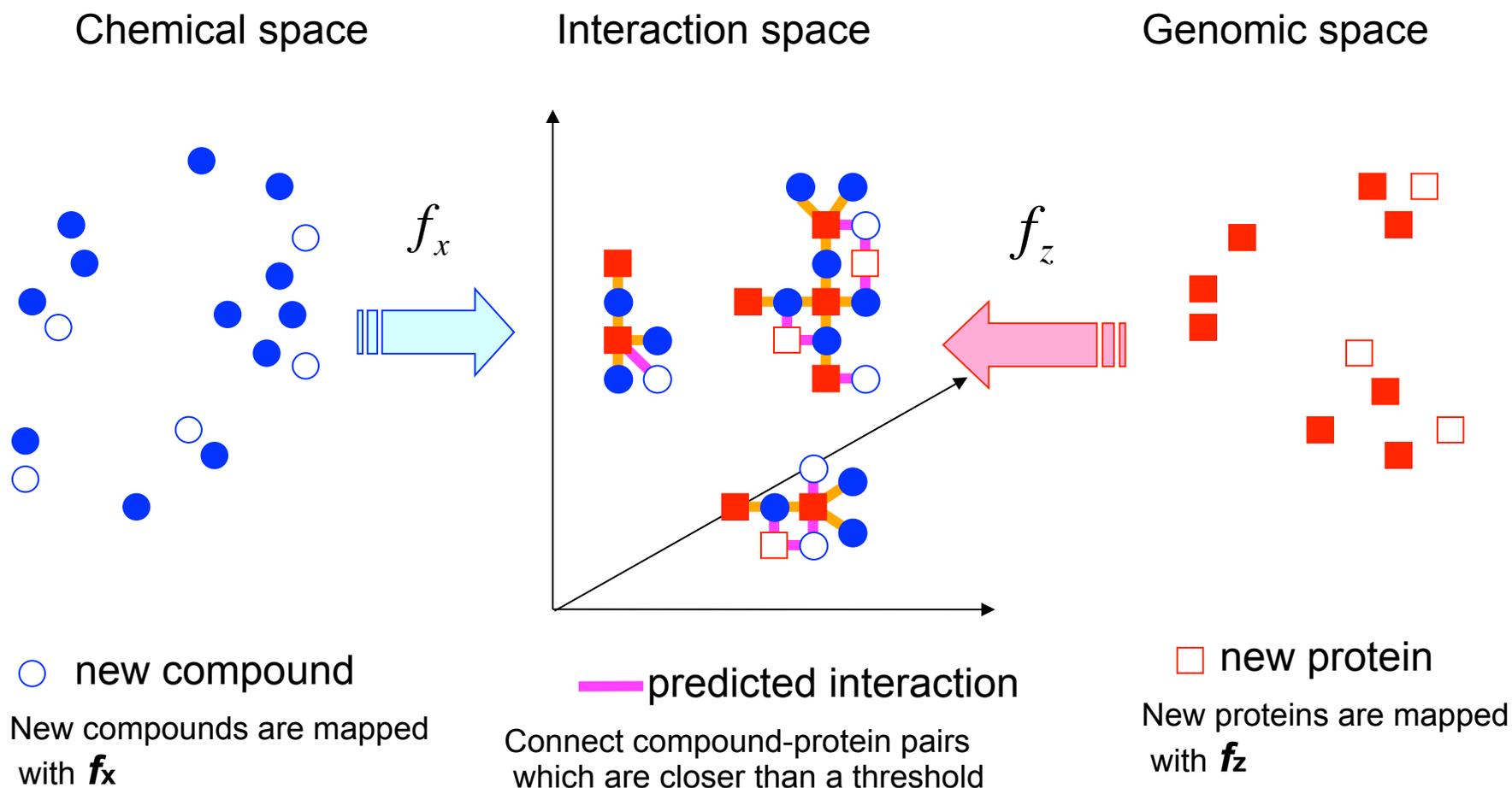
## Step 2. Learning a model between the chemical/genomic space and the interaction space



# Step 3. Predicting unknown interactions involving new compounds/proteins after the projection



# Step 3. Predicting unknown interactions involving new compounds/proteins after the projection



# Mapping to a unified space

Let us consider two functions to map each compound  $\mathbf{x}$  and each protein  $\mathbf{z}$  onto a unified Euclidian space

$$\begin{cases} f_x(\mathbf{x}) = (f_x^{(1)}(\mathbf{x}), \dots, f_x^{(d)}(\mathbf{x}))^T \in \mathbf{R}^d \\ f_z(\mathbf{z}) = (f_z^{(1)}(\mathbf{z}), \dots, f_z^{(d)}(\mathbf{z}))^T \in \mathbf{R}^d \end{cases}$$

We find  $f_x$  and  $f_z$  which minimize

$$R(f_x, f_z) = \frac{\sum_{(\mathbf{x}_i, \mathbf{z}_j) \in E} (f_x(\mathbf{x}_i) - f_z(\mathbf{z}_j))^2 + \lambda_1 \|f_x\|^2 + \lambda_2 \|f_z\|^2}{\sum_{(\mathbf{x}_i, \mathbf{z}_j) \in V_x \times V_z} (f_x(\mathbf{x}_i) - f_z(\mathbf{z}_j))^2}$$

where  $V_x$  (resp.  $V_z$ ) is a set of drugs (resp. target proteins),

$E$  is a set of interactions, and  $\lambda_1$  and  $\lambda_2$  are regularization parameters

# Extraction of multiple features

Successive features  $f_x^{(q)}$  and  $f_z^{(q)}$  ( $q = 1, 2, \dots, d$ ) are obtained by

$$\begin{pmatrix} f_x^{(q)} \\ f_z^{(q)} \end{pmatrix} = \arg \min \frac{\sum_{(\mathbf{x}_i, \mathbf{z}_j) \in E} (f_x(\mathbf{x}_i) - f_z(\mathbf{z}_j))^2 + \lambda_1 \|f_x\|^2 + \lambda_2 \|f_z\|^2}{\sum_{(\mathbf{x}_i, \mathbf{z}_j) \in V_x \times V_z} (f_x(\mathbf{x}_i) - f_z(\mathbf{z}_j))^2}$$

under the following orthogonality constraints:

$$f_x \perp f_x^{(1)}, \dots, f_x^{(q-1)}, \quad f_z \perp f_z^{(1)}, \dots, f_z^{(q-1)}$$

Prediction score for a given pair of compound  $\mathbf{x}'$  and protein  $\mathbf{z}'$ :

$$g(\mathbf{x}', \mathbf{z}') = \sum_{q=1}^d f_x^{(q)}(\mathbf{x}') \cdot f_z^{(q)}(\mathbf{z}')$$

# Algorithm

By the representer theorem, features can be expanded as

$$f_x(\mathbf{x}) = \sum_{j=1}^{n_x} \alpha_j k_x(\mathbf{x}_j, \mathbf{x}), \quad f_z(\mathbf{z}) = \sum_{j=1}^{n_z} \beta_j k_z(\mathbf{z}_j, \mathbf{z})$$

Kernel Gram matrices:

$$(K_x)_{ij} = k_x(\mathbf{x}_i, \mathbf{x}_j), \quad i, j = 1, 2, \dots, n_x$$

$$(K_z)_{ij} = k_z(\mathbf{z}_i, \mathbf{z}_j), \quad i, j = 1, 2, \dots, n_z$$

The solution can be obtained by finding  $\alpha_q$  and  $\beta_q$  which minimizes

$$R(\alpha, \beta) = \frac{\begin{pmatrix} \alpha \\ \beta \end{pmatrix}^T \begin{pmatrix} K_x D_x K_x & -K_x A K_z \\ -K_z A^T K_x & K_z D_z K_z \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} + \lambda_1 \alpha^T K_x \alpha + \lambda_2 \beta^T K_z \beta}{\begin{pmatrix} \alpha \\ \beta \end{pmatrix}^T \begin{pmatrix} K_x^2 & 0 \\ 0 & K_z^2 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix}}$$

under the following constraints:

$$\alpha^T K_x \alpha_1 = \dots = \alpha^T K_x \alpha_{q-1} = 0, \quad \beta^T K_z \beta_1 = \dots = \beta^T K_z \beta_{q-1} = 0$$

where

$D_x$  (resp.  $D_z$ ): degree matrix of drugs (resp. target proteins),

$A$ : adjacency matrix of drug-target interactions

It is reduced to the generalized eigenvalue problem:

$$\begin{pmatrix} K_x D_x K_x + \lambda_1 K_x & -K_x A K_z \\ -K_x A^T K_x & K_z D_z K_z + \lambda_2 K_z \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \rho \begin{pmatrix} K_x^2 & 0 \\ 0 & K_z^2 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

# Drug-target interaction data for human: Gold standard data

	Statistics
Number of drugs	1874
Number of target proteins (Total in human genome)	436 (23196)
Number of drug-target interactions	6769

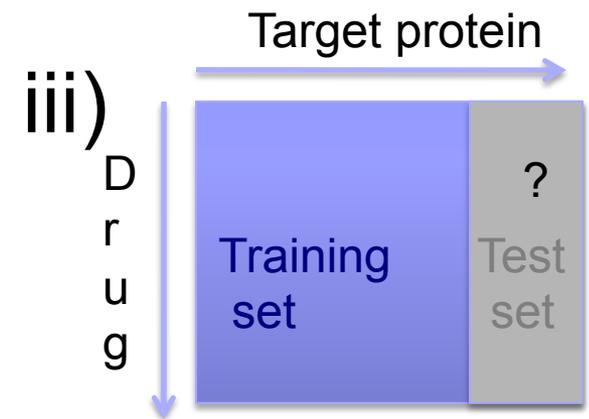
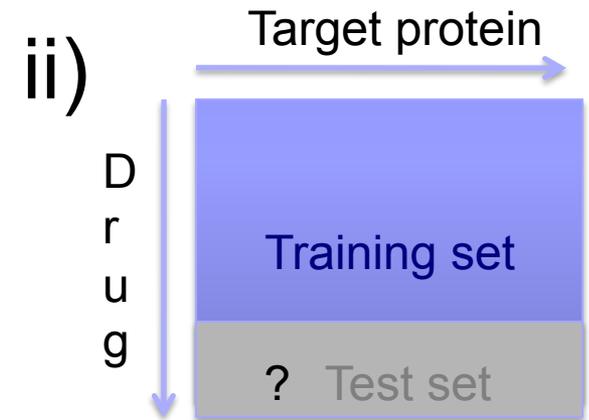
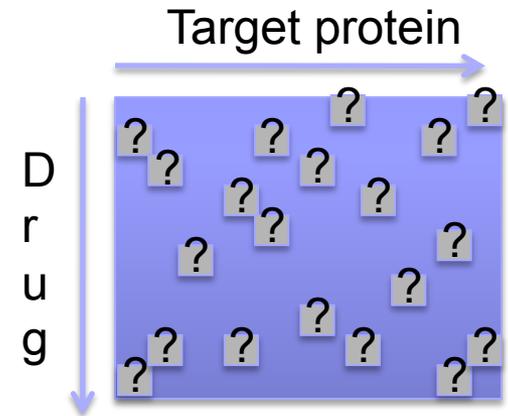
KEGG DRUG (December, 2011)

# Cross-validation (CV) i)

i) Pairwise CV  
(Missing interaction detection)

ii) Blockwise CV I  
(new drug identification)

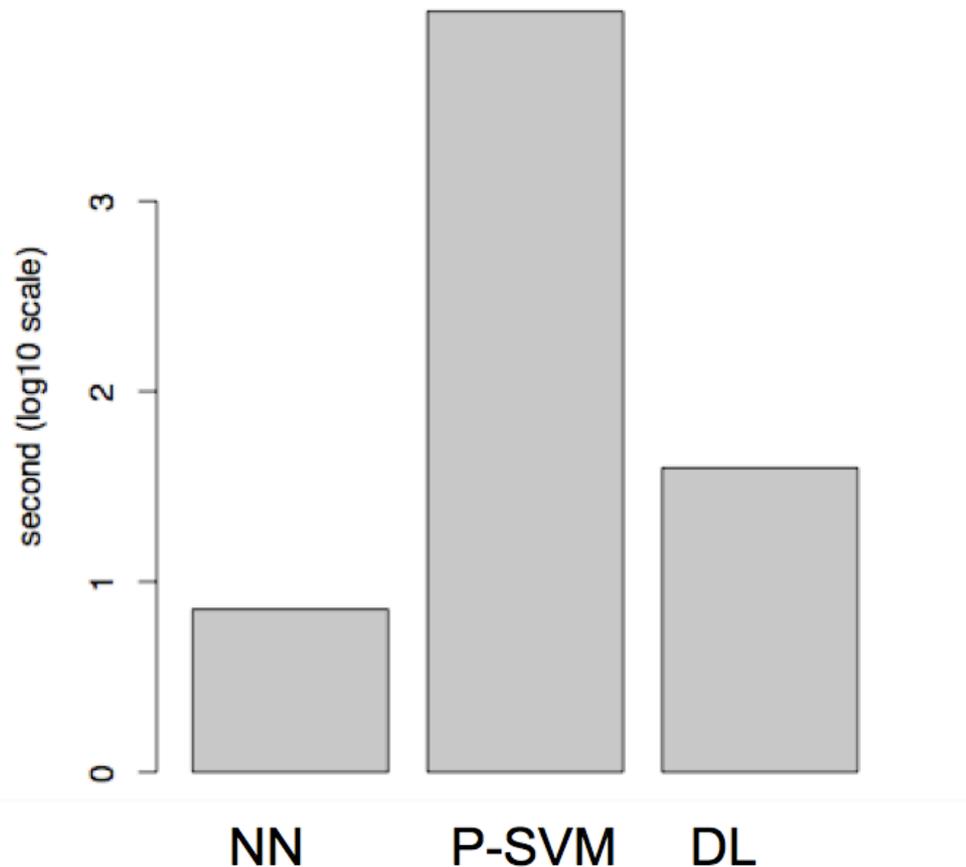
iii) Blockwise CV II  
(new target identification)



# Performance evaluation by 3-fold cross-validation: AUC (Area under the ROC curve)

CV	Method	AUC	(S.D)
Pairwise CV (Missing interaction detection)	Nearest neighbor	0.894	(0.009)
	Pairwise SVM	0.968	(0.006)
	Proposed method	0.972	(0.007)
Blockwise CV I (New drug identification)	Nearest neighbor	0.808	(0.010)
	Pairwise SVM	0.855	(0.013)
	Proposed method	0.869	(0.006)
Blockwise CV II (New target identification)	Nearest neighbor	0.711	(0.009)
	Pairwise SVM	0.805	(0.003)
	Proposed method	0.811	(0.009)

# Computational cost



NN: nearest neighbor

P-SVM: pairwise SVM

DL: distance learning  
(proposed method)



# Comprehensive prediction of unknown drug-target interactions

- Test drugs: all compounds in KEGG LIGAND and all drugs in KEGG DRUG
- Test target proteins: all human proteins in KEGG GENES
- All gold standard interaction data are used in the training



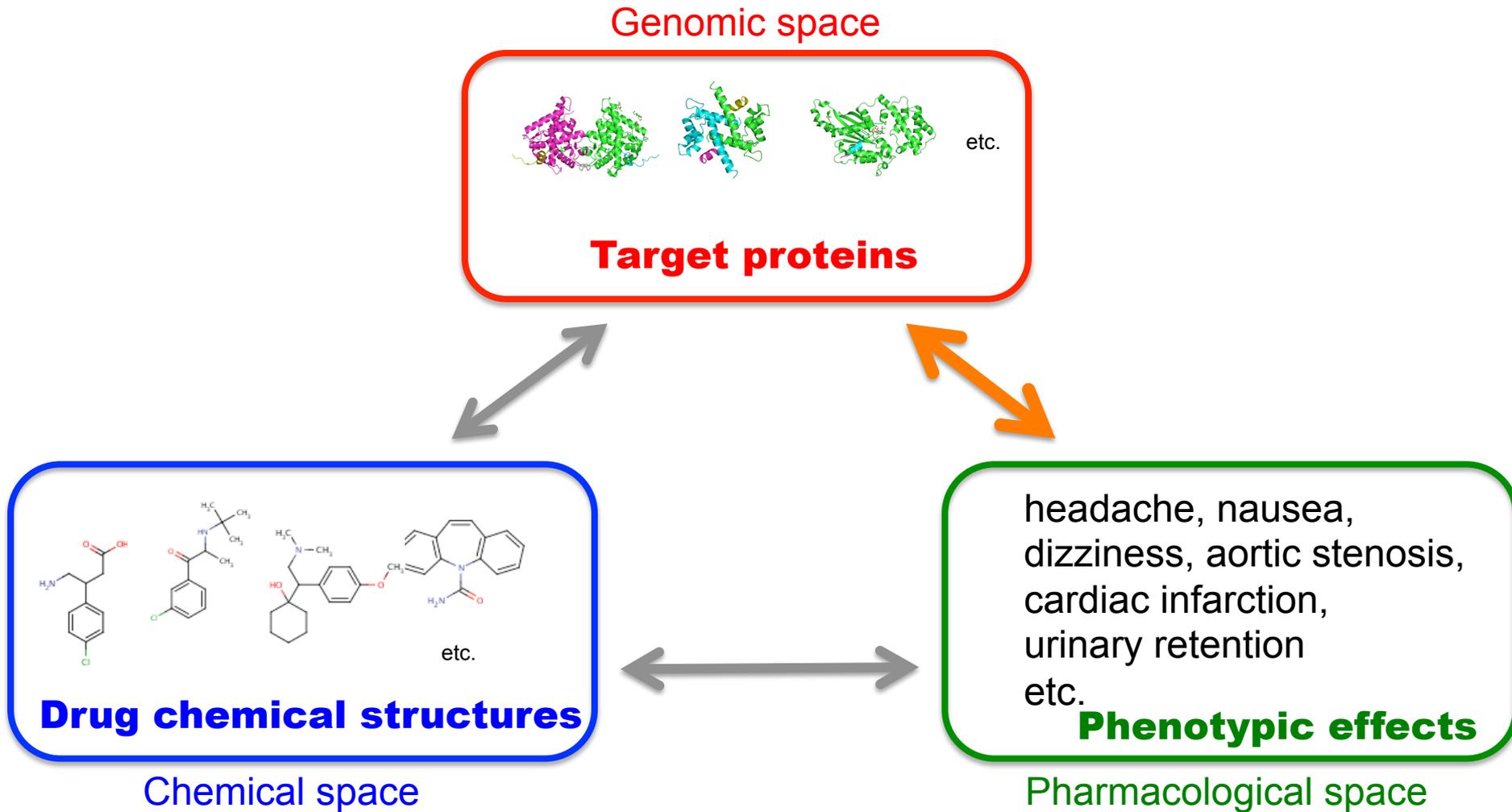
# Summary

- The proposed method can predict unknown drug-target interactions on a large scale.
- The prediction is performed based on the integration of chemical and genomic data.
- It does not need 3D structure information of the target proteins.

# Outline

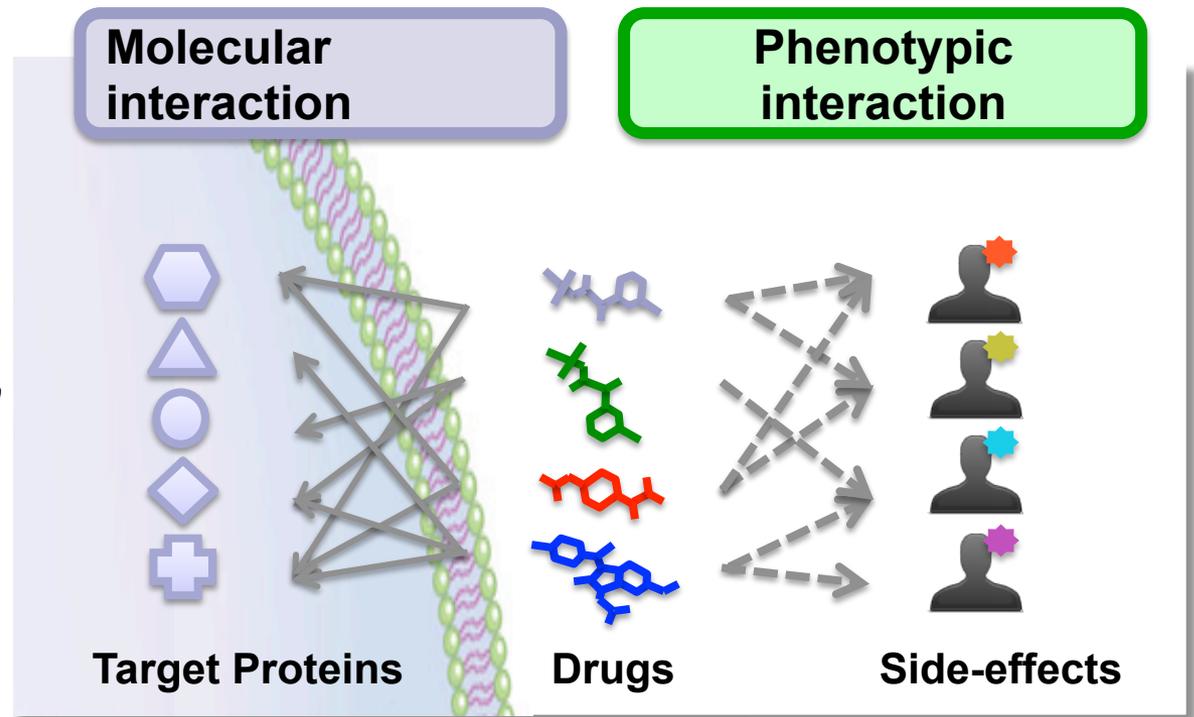
- Background
  - Chemical, genomic and pharmacological spaces
- Methods for pharmaceutical applications
  - Drug target prediction from chemical and genomic data
  - Side-effect prediction from biological data
- Results
- Concluding remarks

# Heterogeneous omics data

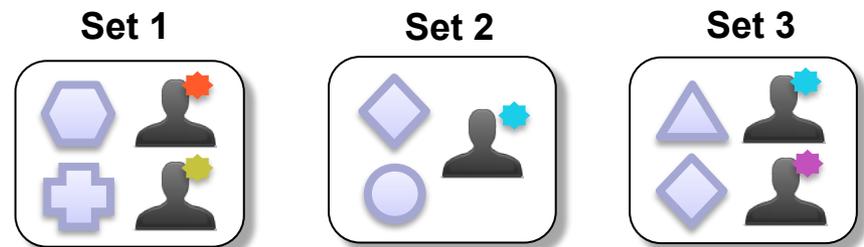


# Objective

*Two interaction data*



*Extraction of correlated sets  
of **target proteins** and **side-effects***



# Correlation analysis of target proteins and side-effects

Target protein profile for each drug:  $\mathbf{x}=(x_1, x_2, \dots, x_p)^T$

indicating presence/absence of 1368 target proteins (DrugBank, Matador)

Side-effect profile for each drug:  $\mathbf{y}=(y_1, y_2, \dots, y_q)^T$

indicating presence/absence of 1339 side-effects (SIDER)

Goal: extraction of target proteins and side-effects which share similar sets of drugs in terms of occurrence.

# Ordinary canonical correlation analysis (OCCA)

- 1 For each drug, set two binary vectors indicating presence/absence of interactions;

Target protein profile:  $\mathbf{x} = [x_1, \dots, x_p]^T$ , Side-effect profile:  $\mathbf{y} = [y_1, \dots, y_q]^T$

- 2 Consider two linear combinations for drug  $i$  with weight vectors  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$ ;

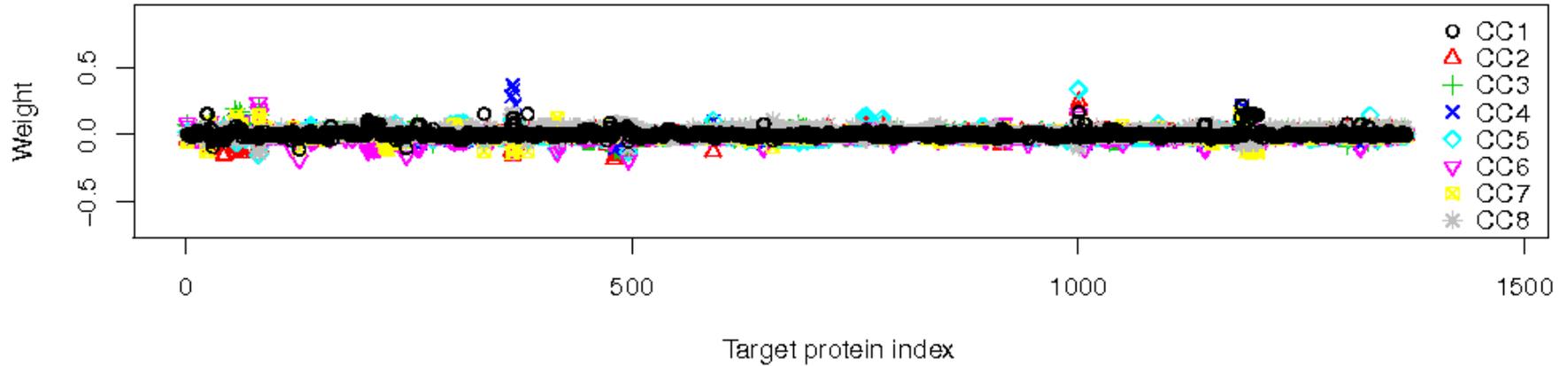
$$u_i = \boldsymbol{\alpha}^T \mathbf{x}_i \quad \text{and} \quad v_i = \boldsymbol{\beta}^T \mathbf{y}_i \quad (u \text{ and } v \text{ are centered})$$

Find  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$  that maximize the canonical correlation coefficient;

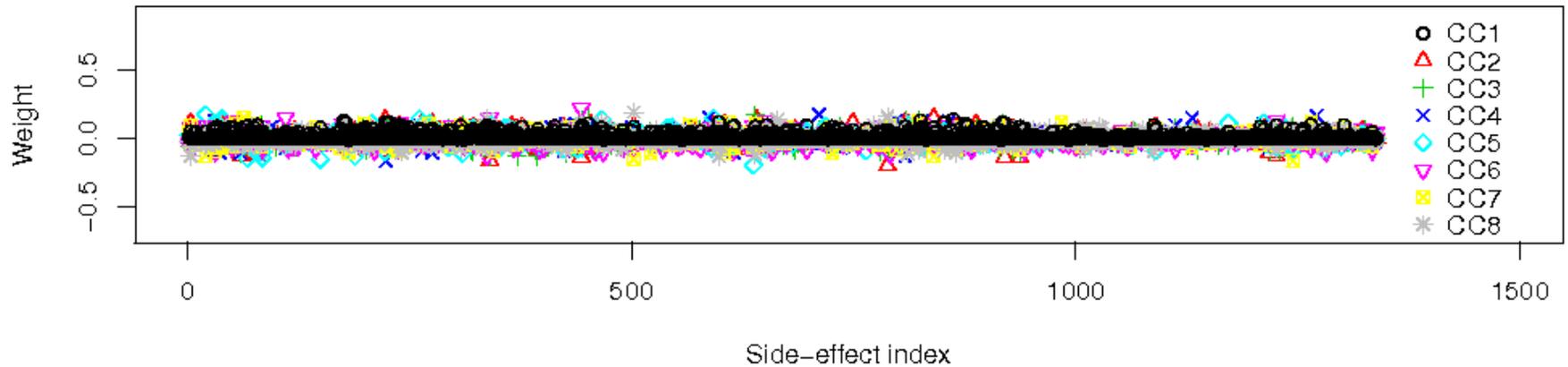
$$\rho = \text{corr}(u, v) = \frac{\sum_{i=1}^n \boldsymbol{\alpha}^T \mathbf{x}_i \cdot \boldsymbol{\beta}^T \mathbf{y}_i}{\sqrt{\sum_{i=1}^n (\boldsymbol{\alpha}^T \mathbf{x}_i)^2} \sqrt{\sum_{i=1}^n (\boldsymbol{\beta}^T \mathbf{y}_i)^2}}$$

# Weight vectors in OCCA

OCCA weight for target proteins



OCCA weight for side-effects



# Sparse canonical correlation analysis (SCCA)

2 Consider two linear combinations for drug  $i$  with weight vectors  $\alpha$  and  $\beta$ ;

$$u_i = \alpha^T \mathbf{x}_i \quad \text{and} \quad v_i = \beta^T \mathbf{y}_i \quad (u \text{ and } v \text{ are centered})$$

Find  $\alpha$  and  $\beta$  that maximize the canonical correlation coefficient;

$$\rho = \text{corr}(u, v) = \frac{\sum_{i=1}^n \alpha^T \mathbf{x}_i \cdot \beta^T \mathbf{y}_i}{\sqrt{\sum_{i=1}^n (\alpha^T \mathbf{x}_i)^2} \sqrt{\sum_{i=1}^n (\beta^T \mathbf{y}_i)^2}}$$

## SCCA

3 In order to impose sparsity on  $\alpha$  and  $\beta$  we give an  $L_1$  penalty;

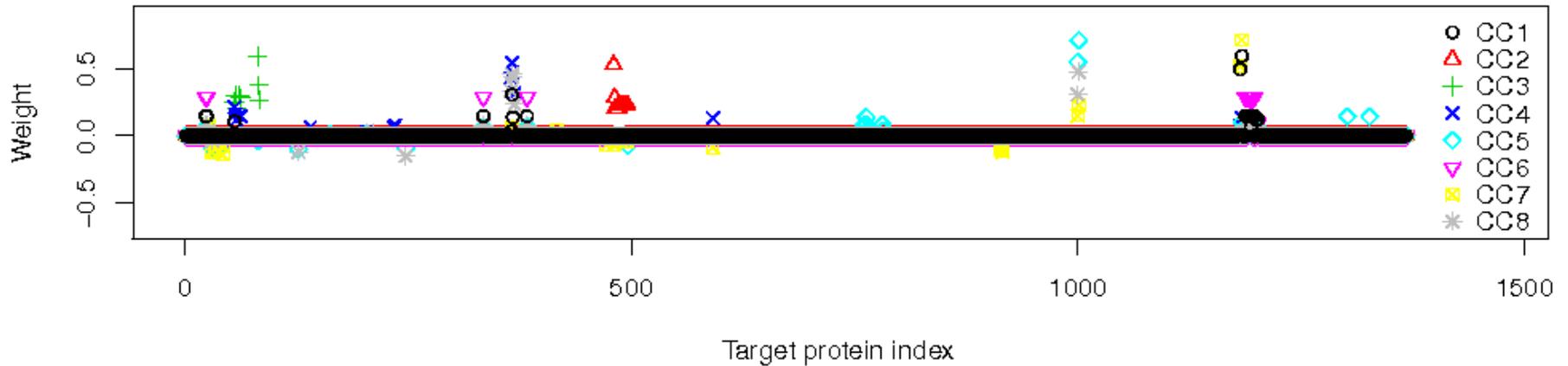
$$\|\alpha\|_1 \leq c_1 \sqrt{p}, \quad \|\beta\|_1 \leq c_2 \sqrt{q},$$

$0 < c_1 \leq 1$  and  $0 < c_2 \leq 1$  are sparsity parameters.

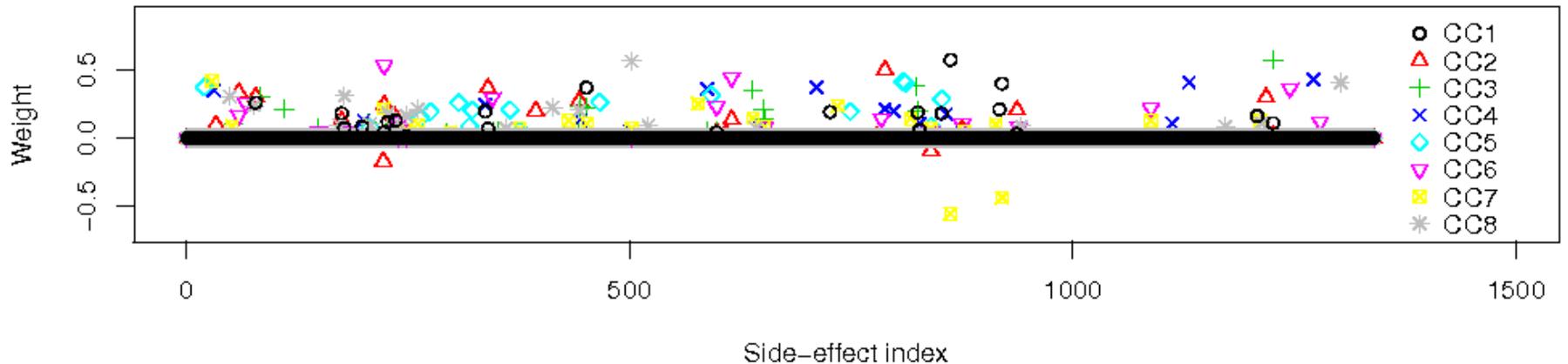
4 For  $n$  drugs, given  $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T$ ,  $Y = [\mathbf{y}_1, \dots, \mathbf{y}_n]^T$ , the weight vectors  $\alpha$  and  $\beta$  can be optimized by solving penalized matrix decomposition of the matrix  $X^T Y$  (Witten *et al.*, *Biostatistics* 2009).

# Weight vectors in SCCA

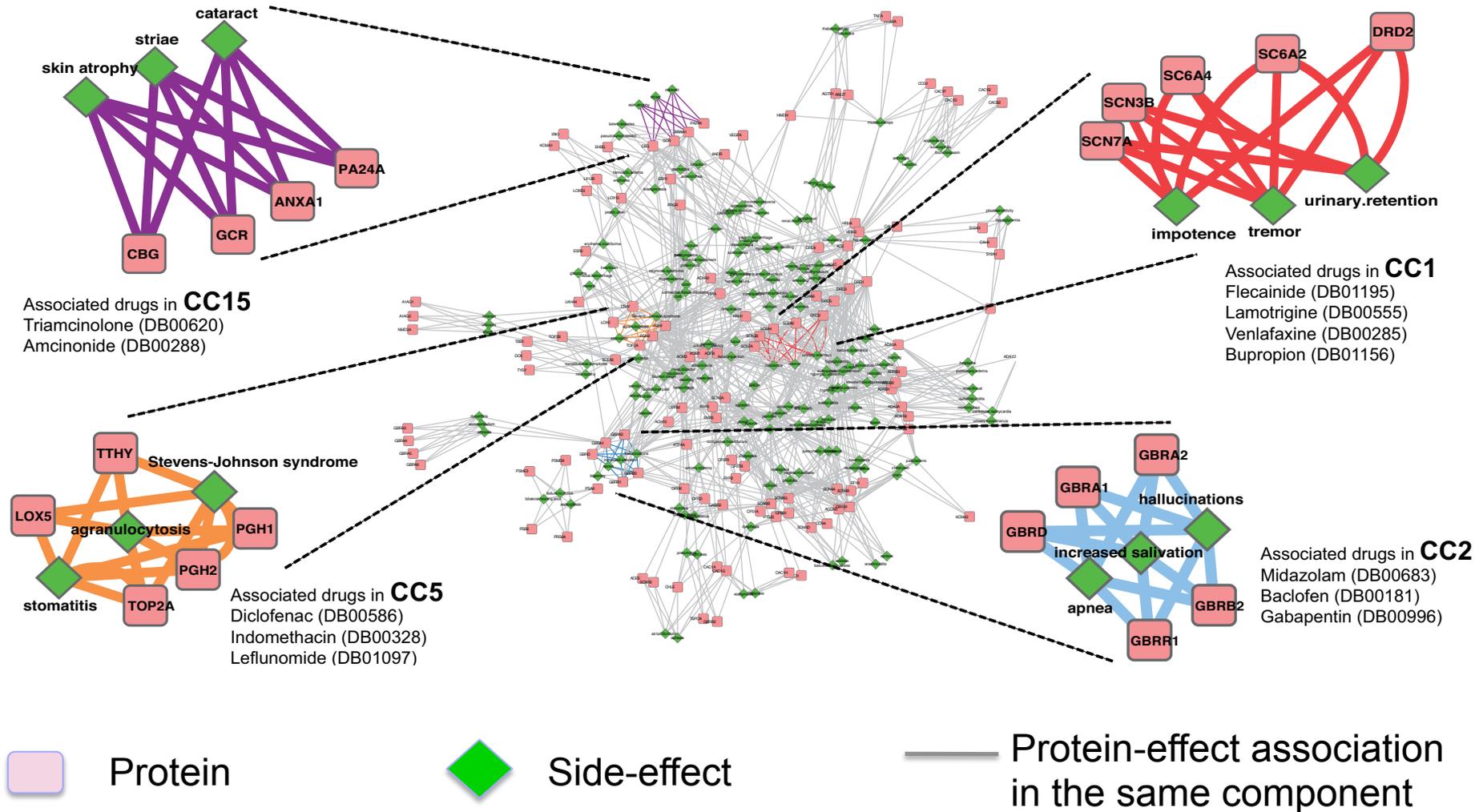
SCCA weight for target proteins



SCCA weight for side-effects



# Extracted protein-effect association network across 80 components in SCCA

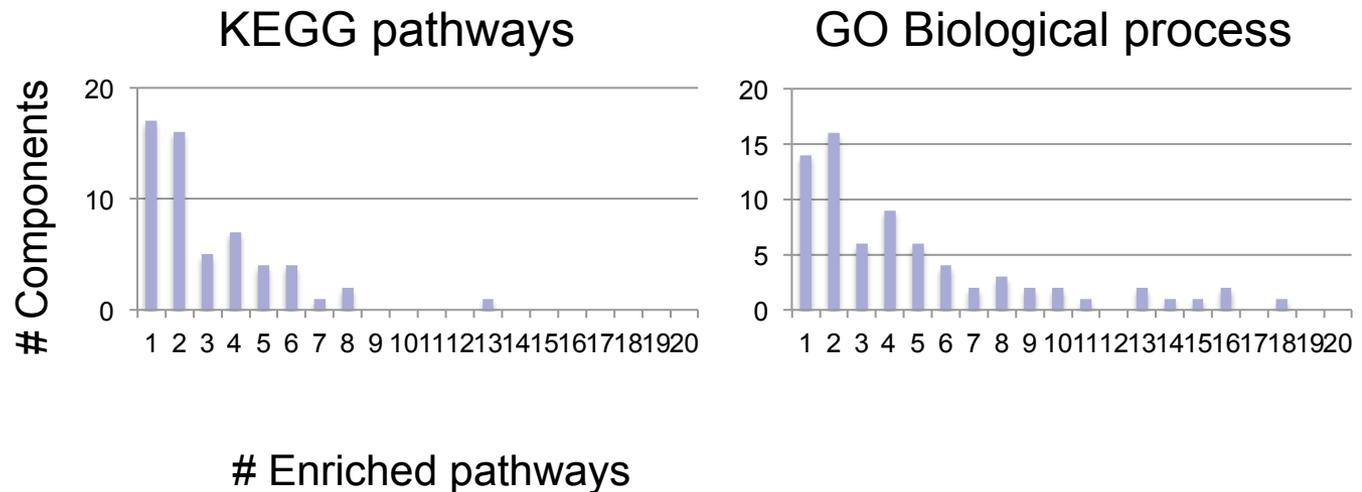


# Biological evaluation of extracted components

- Evaluation of the components by pathway enrichment analysis
  - How each of the correlated sets is enriched with proteins that are involved in a specific pathway?
- Pathway databases
  - KEGG : pathway maps
  - GO (Gene Ontology) : biological process terms

# Pathway enrichment of proteins in components

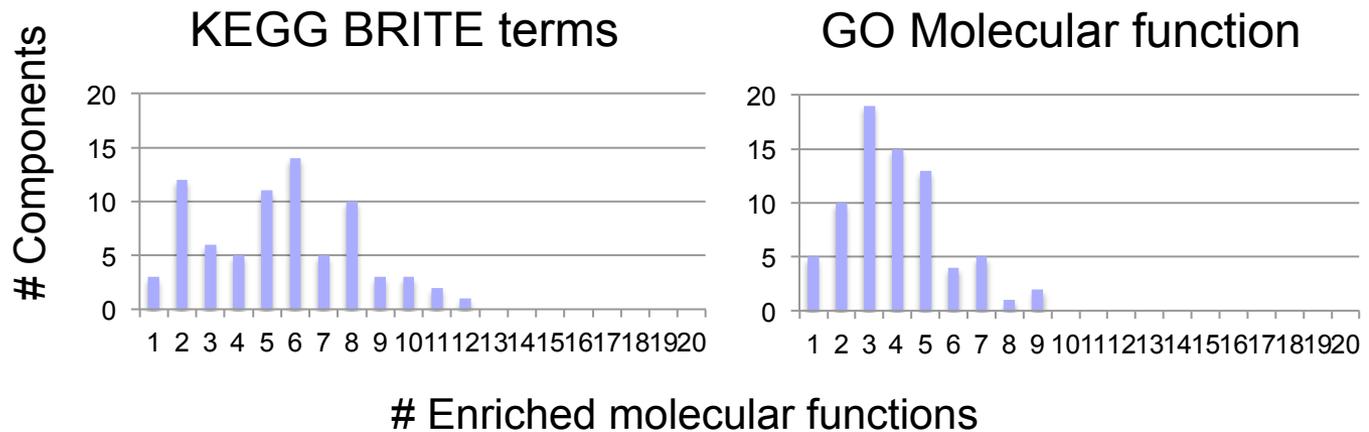
## Frequency of components v.s. Number of enriched pathways



The components were enriched with a small number of pathways.

# Molecular function enrichment of proteins in components

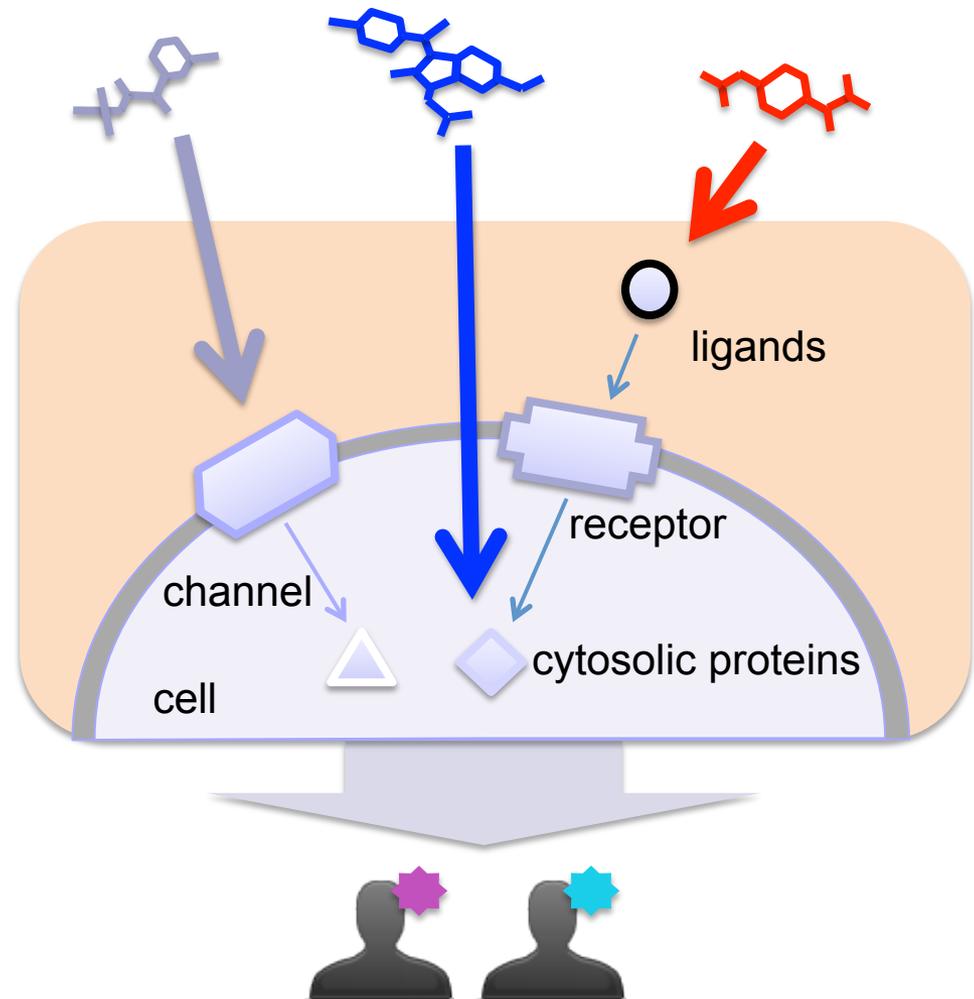
## Frequency of components v.s. molecular functions



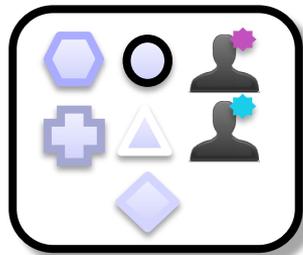
The components were enriched with a small number of pathways, but with various number of different molecular functions.

# Biologically relevant interpretation of correlated sets

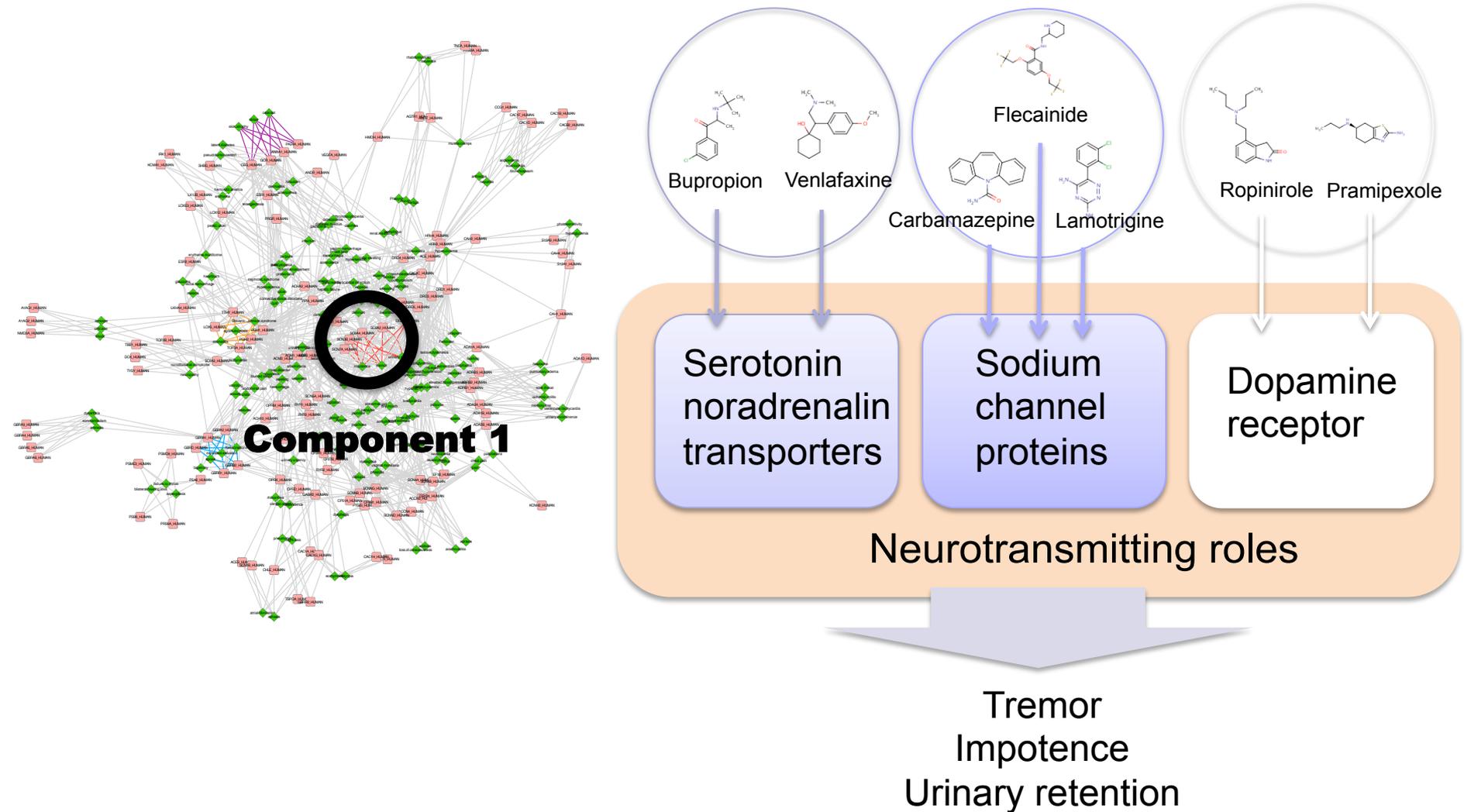
Interactions between drugs with proteins of different molecular functions tend to induce the regulation of the whole protein pathway, which lead to certain side-effects.



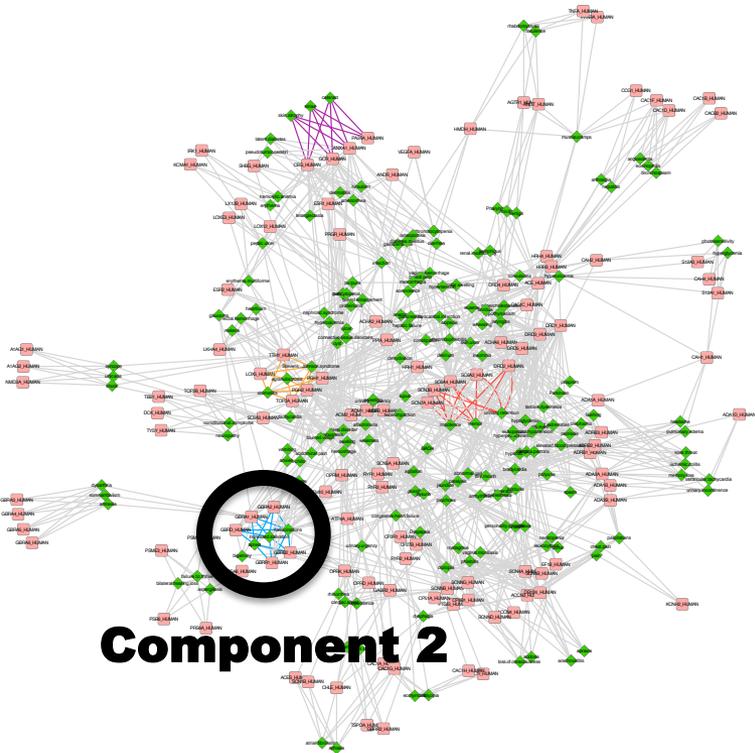
Component



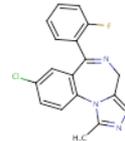
# An example of extracted components



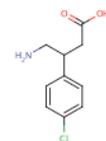
# An example of extracted components



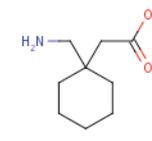
anti-anxiety, anesthesia adjuvants, ...



Midazolam



Baclofen



Gabapentin

GABA  
receptor  
subunits

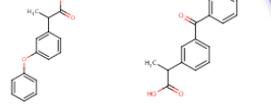
GABA receptor

Increased salivation  
Hallucinations  
apnea

# An example of extracted components

## Component 5

Anti-inflammatory



Fenoprofen Ketoprofen

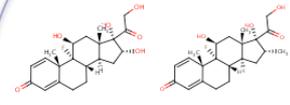
Prostaglandin G/H synthase 1 & 2  
Arachidonate 5-lipoxygenase

Arachidonic acid  
metabolism

Stevens-Johnson syndrome  
Stomatitis  
Agranulocytosis

## Component 15

Anti-inflammatory



Triamcinolone Dexamethazone

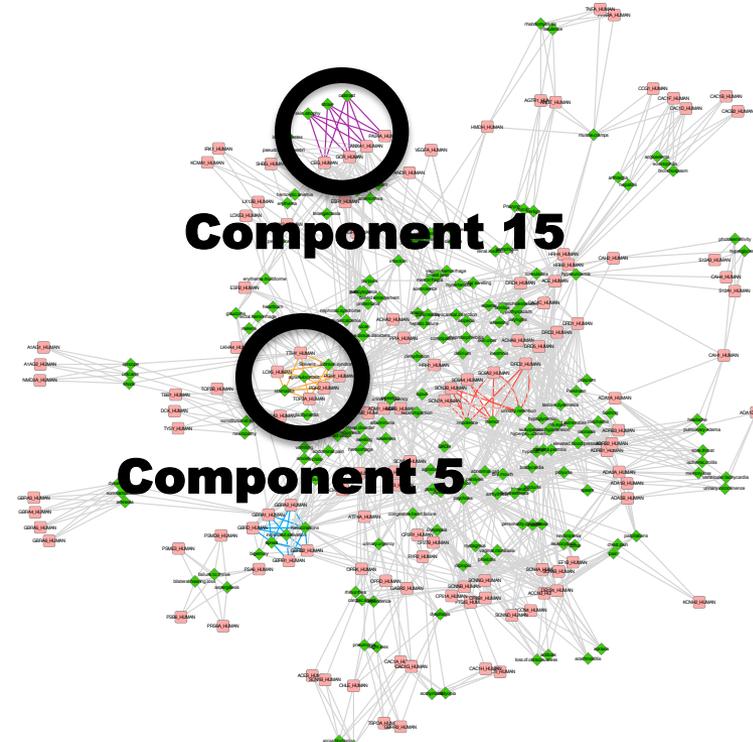
Annexin 1A  
Glucocorticoid receptor  
corticosteroid-binding  
globulin  
phospholipase A2

Glucocorticoid effects

Striae  
Skin atrophy  
Cataract

Component 15

Component 5



# Performance evaluation for side-effect prediction by 5-fold cross-validation

Prediction score :

$$\mathbf{y}_{new} = B^{-T} A^T \mathbf{x}_{new}$$

where  $A = [\alpha_1, \dots, \alpha_m]$ ,  $B = [\beta_1, \dots, \beta_m]$ ,  $m$  : number of components

Method	AUC $\pm$ S.D.	AUPR $\pm$ S.D.
Chemical structure based approach:		
Random	0.5000 $\pm$ 0.0000	0.0556 $\pm$ 0.0000
OCCA	0.8355 $\pm$ 0.0010	0.3753 $\pm$ 0.0016
SCCA	0.8708 $\pm$ 0.0007	0.3766 $\pm$ 0.0030
Target protein based approach:		
Random	0.5000 $\pm$ 0.0000	0.0556 $\pm$ 0.0000
OCCA	0.8850 $\pm$ 0.0007	0.4067 $\pm$ 0.0006
SCCA	<b>0.8895 <math>\pm</math> 0.0002</b>	<b>0.4103 <math>\pm</math> 0.0018</b>

# Examples of predicted side-effects by SCCA

- For 730 uncharacterized drugs in DrugBank

## Top 20 predicted side-effects

	Drug	Side-effect	Score
1	Cinnarizine	<i>tremor</i>	1.176339
2	Benzocaine	<i>diplopia</i>	1.163882
3	Cinnarizine	<i>constipation</i>	1.143190
4	Bepridil	<i>tremor</i>	1.046472
5	Cinnarizine	<i>somnolence</i>	0.996260
6	Cinnarizine	<i>dry.mouth</i>	0.961833
7	Cinnarizine	<i>angioedema</i>	0.955471
8	Cinnarizine	<i>insomnia</i>	0.950757
9	Benzocaine	<i>nausea</i>	0.947105
10	Alprenolol	<i>dizziness</i>	0.943918

	Drug	Side-effect	Score
11	Benzocaine	<i>diarrhea</i>	0.937409
12	Nisoldipine	<i>tremor</i>	0.933258
13	Nitrendipine	<i>tremor</i>	0.933258
14	Lercanidipine	<i>tremor</i>	0.933258
15	Benzocaine	<i>syncope</i>	0.928905
16	Promazine	<i>tachycardia</i>	0.926875
17	Promazine	<i>somnolence</i>	0.924103
18	Benzocaine	<i>vomiting</i>	0.922091
19	Bepridil	<i>nervousness</i>	0.918828
20	Cinnarizine	<i>nervousness</i>	0.918615

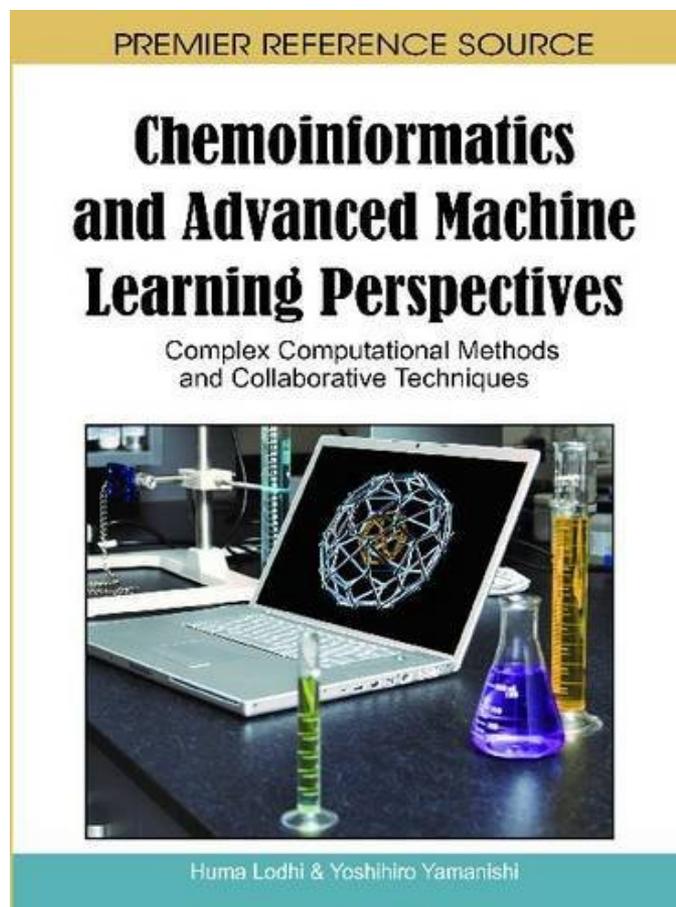
- Side-effects reported in the literature are colored in red.



# Summary

- The proposed method enables us to relate target proteins with side-effects of drugs.
- It may be useful to predict unknown side-effects in the drug development.

# Machine learning in chemoinformatics



Lodhi, H. and Yamanishi, Y.,  
Chemoinformatics and Advanced Machine Learning Perspectives,  
IGI Global, 2010.

# Acknowledgements

- Kyushu University
  - Mikita Suyama, Tetsuya Sato
- Kyushu Institute of Technology
  - Hiroto Saigo
- Kyoto University
  - Minoru Kanehisa, Susumu Goto, Masaaki Kotera, Masataka Takarabe, Sayaka Mizutani, Yosuke Nishimura
- JST
  - Yasuo Tabei
- Curie Institute - Inserm U900 - Mines ParisTech
  - Jean-Philippe Vert, Véronique Stoven, Kevin Bleakley, Edouard Pauwels