

Compare

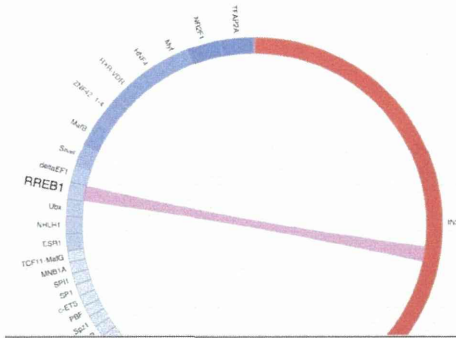
positive #1 mTOR (H.sapiens) | control #2 mTOR (H.sapiens) | compare

Results

Job	Date	Gene List	Repeat Mask	UpStream	Mode	DownStream	Scoring	Pareto	Delete
#1 mTOR (H.sapiens)	2015/12/16 16:34:03 ↓ 2015/12/16 16:34:40	INS (458)	yes	2000	mode2	200	jaspar	30	<a href="#">delete</a>
#2 mTOR (H.sapiens)	2015/12/16 16:36:58 ↓ 2015/12/16 17:01:38	INS (1) IGF1 (1) IRS1 (116) PIK3R2 (33) PIK3R3 (53) PTEN (3) PDPK1 (135) AKT3 (3) RRAGA (16) HIF1A (24) RPS6KB1 (27) RPS6KB2 (10) EIF4B (48) EIF4EBP1 (54) EIF4E (1) MAPK1 (5) MAPK3 (8)	yes	2000	mode2	200	transfac32	30	<a href="#">delete</a>

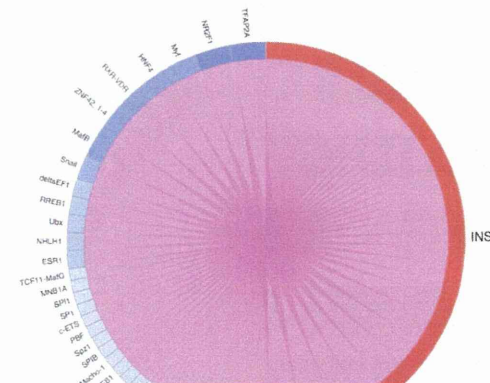
The screenshot shows the SHOE web application interface. At the top, there are navigation buttons for 'Job Input' and 'Queue List'. The main content area displays sequence alignments for three different regions, with a 'Pareto front' plot on the right showing the relationship between MA Score and PSSM Score. Below the alignments is a detailed results table with columns for Gene, NM, TF, Region, strand, NA score, PSSM score, ID, MOTIF, CONSENSUS, Hit, and Pareto. A red arrow points to the 'Insulin' entry in the table, which is highlighted in blue. The table also includes a 'Code Means' column with color-coded icons for each entry.

Gene	NM	TF	Region	strand	NA score	PSSM score	ID	MOTIF	CONSENSUS	Hit	Pareto
INS	NM_000207	Snail1	62-70	+	2	-0.281	SNAIL	CGTGG	CGTGG	1	0
INS	NM_000207	Usk	68-92	+	1.67	-0.241	USK	CGTGG	CGTGG	1	0
INS	NM_000207	NR2A	353-357	+	1.67	-0.581	NR2A	CGTGG	CGTGG	1	1
INS	NM_000207	Usk	299-302	+	1.45	-0.241	USK	CGTGG	CGTGG	1	2
INS	NM_000207	Usk	142-145	+	1.45	-0.241	USK	CGTGG	CGTGG	1	2
INS	NM_000207	YY1	247-252	+	3.15	-0.081	YY1	CGTGG	CGTGG	0.92	0
INS	NM_000207	TCF11-MaxC	29-34	+	2.02	-0.921	TCF11	CGTGG	CGTGG	0.92	1
INS	NM_000207	YY1	411-424	+	1.92	-1.081	YY1	CGTGG	CGTGG	0.92	7
INS	NM_000207	Deaf2	353-358	+	1.78	-1.1	DEAF2	CGTGG	CGTGG	0.92	5
INS	NM_000207	ENF42-1-4	20-25	+	1.72	-0.814	ENF42	CGTGG	CGTGG	0.92	2
INS	NM_000207	ENF42-1-4	32-37	+	1.02	-0.814	ENF42	CGTGG	CGTGG	0.92	10
INS	NM_000207	ENF42-1-4	21-24	+	6.8	-0.851	ENF42	CGTGG	CGTGG	0.92	15
INS	NM_000207	TFE3	352-359	+	1.87	-0.581	TFE3	CGTGG	CGTGG	0.9	2
INS	NM_000207	TFE3	54-60	+	1.02	-0.811	TFE3	CGTGG	CGTGG	0.9	9
INS	NM_000207	E74A	240-244	+	2.49	-1.221	E74A	CGTGG	CGTGG	0.86	1
INS	NM_000207	Ena1	146-151	+	3.18	-1.581	ENA1	CGTGG	CGTGG	0.83	0
INS	NM_000207	Ena1	37-42	+	2.79	-1.871	ENA1	CGTGG	CGTGG	0.83	7
INS	NM_000207	Ena1	37-42	+	2.76	-1.871	ENA1	CGTGG	CGTGG	0.83	8



Sim. z 0.5 Pareto s 30 filter Download (CSV, Pareto Plot, Selected Alignment) Chart Alignments Reset Sorting

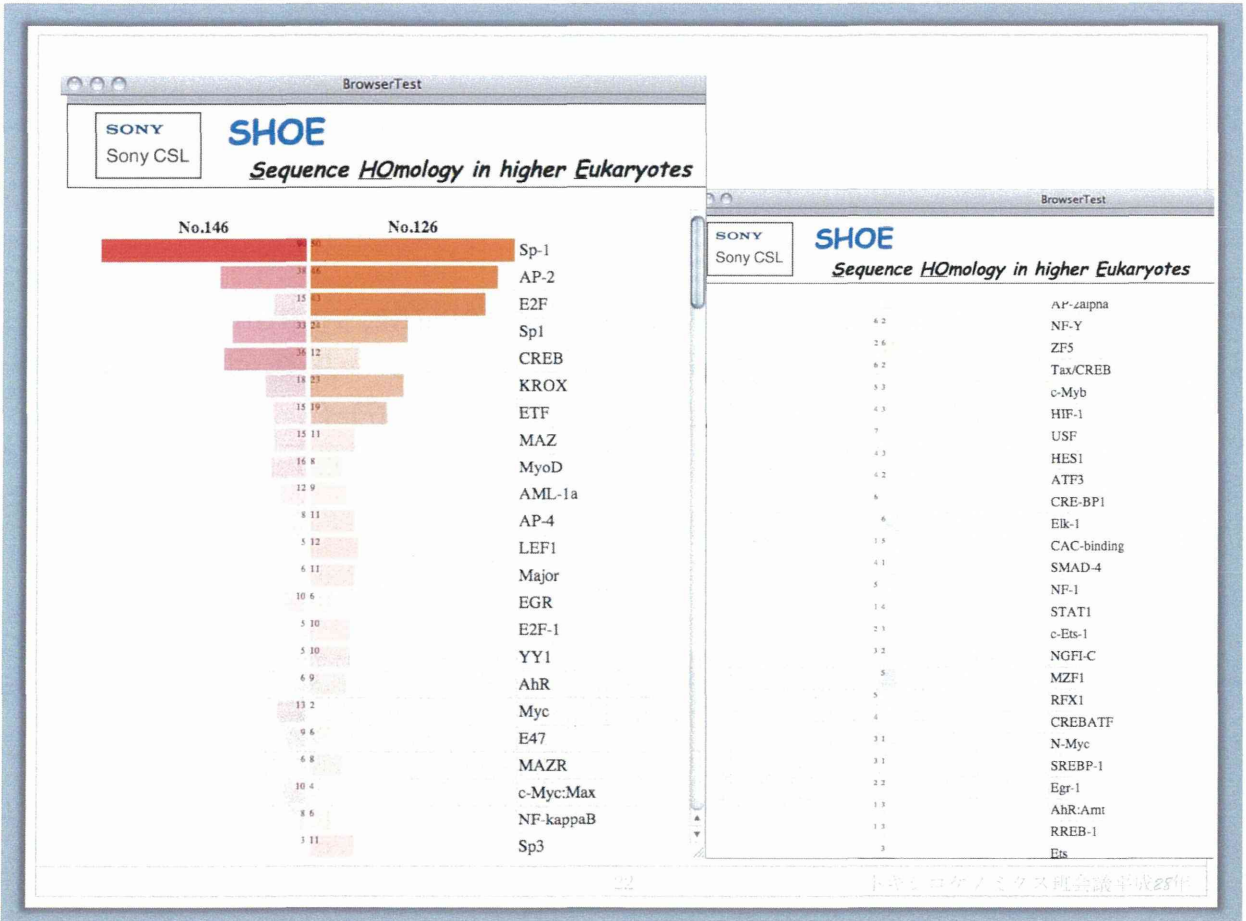
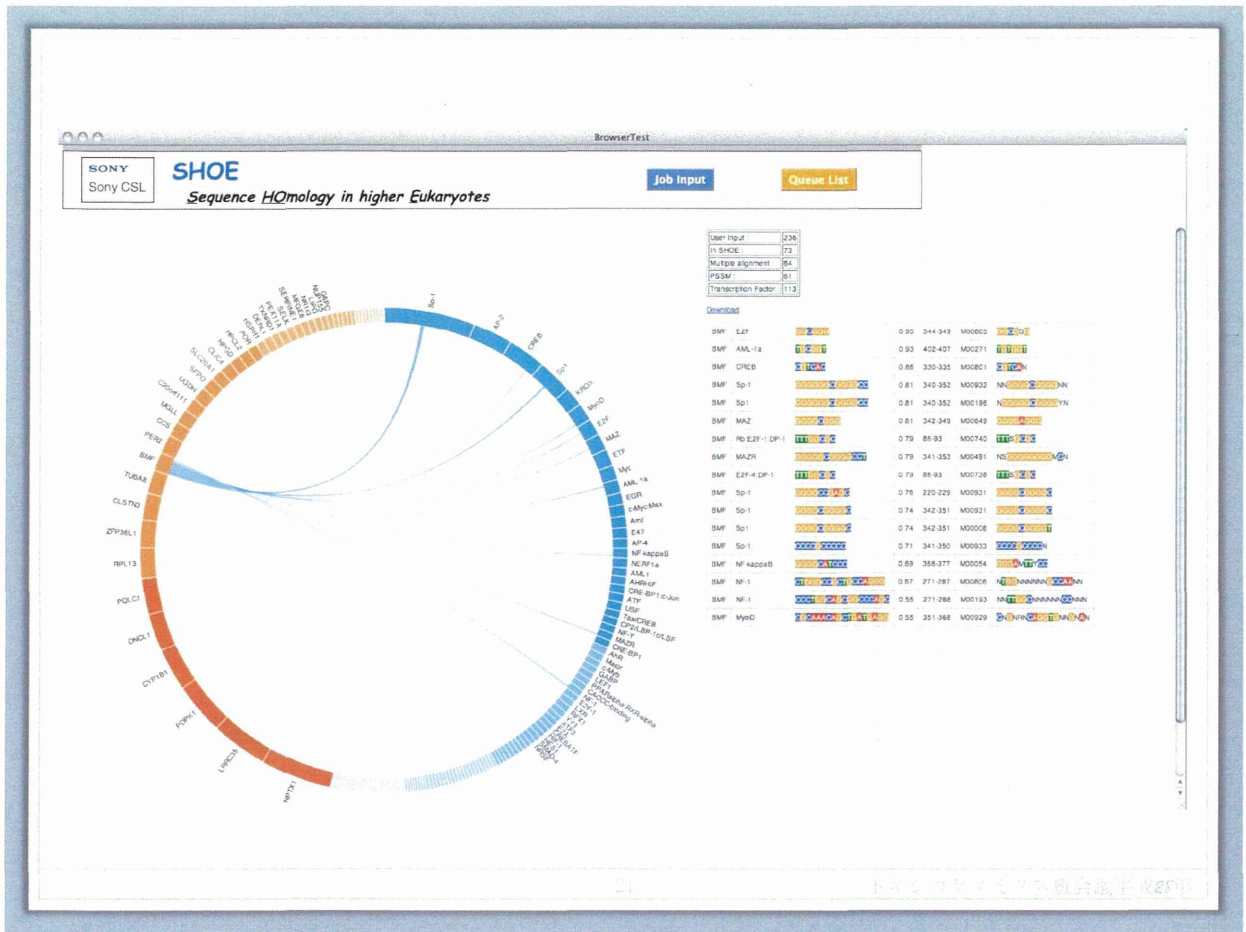
Gene	NM	TF	Region	Strand	MA Score	PSM Score	ID	MOTIF	CONSENSUS	Sim.	Pareto
INS	NM_000207	del-t:EP1	63-70	-	2.43	-0.74	1	CACTG	CACTG	0.56	9
INS	NM_000207	Smal1	180-191	-	3.18	-1.58	1	CACTG	CACTG	0.83	9
INS	NM_000207	Ubx	48-53	-	1.67	-0.26	1	CACTG	CACTG	1	9
INS	NM_000207	TFAP2A	106-114	-	4.06	-2.64	1	CACTG	CACTG	0.63	9
INS	NM_000207	Myf	266-273	-	5.96	-5.13	1	CACTG	CACTG	0.58	9
INS	NM_000207	EBF2	51-57	-	3.29	-2.11	1	CACTG	CACTG	0.71	9
INS	NM_000207	Small	65-70	-	2	-0.28	1	CACTG	CACTG	1	9
INS	NM_000207	RELA	242-251	-	4.79	-3.12	1	CACTG	CACTG	0.65	9
INS	NM_000207	YY1	267-272	-	3.1	-1.06	1	CACTG	CACTG	0.92	9
INS	NM_000207	MHL1	108-119	-	5.73	-3.81	1	CACTG	CACTG	0.61	9
INS	NM_000207	RREB1	210-229	-	8.72	-9.89	1	CACTG	CACTG	2.0	9
INS	NM_000207	E74A	240-246	-	2.49	-1.23	1	CACTG	CACTG	0.5	9
INS	NM_000207	IRF2	242-259	-	8.21	-9.13	1	CACTG	CACTG	0.5	9
INS	NM_000207	MNB1A	353-359	-	1.67	-0.58	1	CACTG	CACTG	0.5	9
INS	NM_000207	SRF	131-150	-	6.76	-8.21	1	CACTG	CACTG	0.5	9
INS	NM_000207	TCF11-MaFG	29-34	-	2.01	-0.92	2	CACTG	CACTG	0.5	9
INS	NM_000207	REL	242-251	-	4.79	-3.79	1	CACTG	CACTG	0.5	9
INS	NM_000207	Smal1	221-230	-	3.17	-1.58	1	CACTG	CACTG	0.5	9



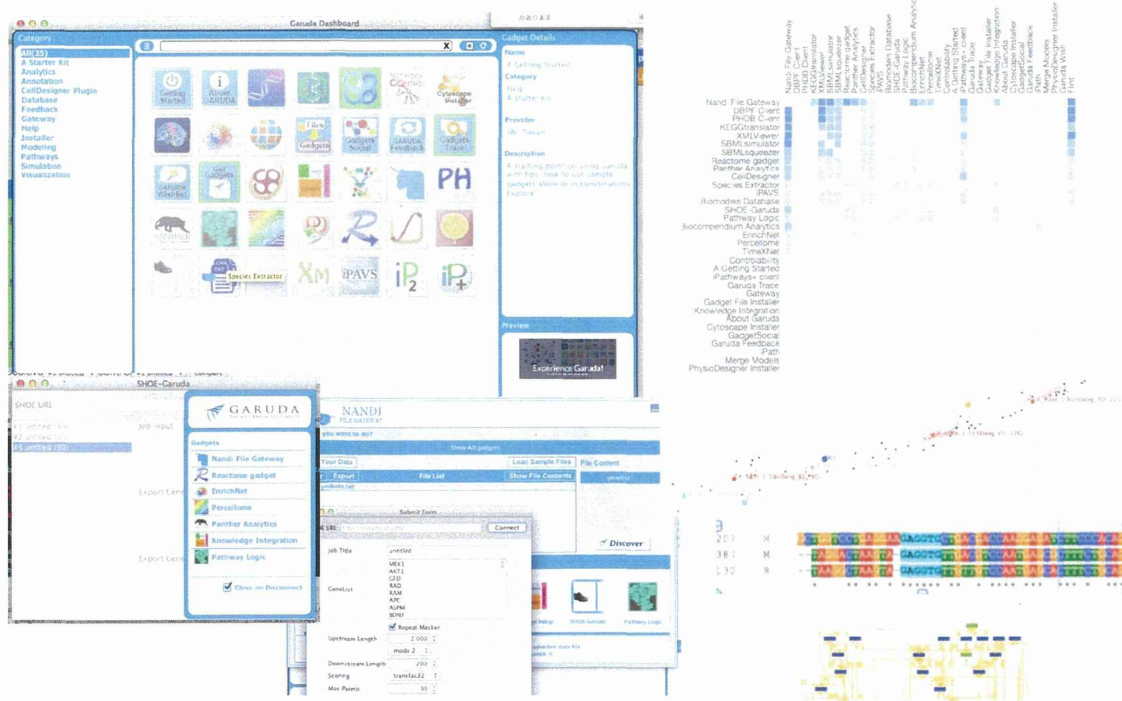
Sim. z 0.5 Pareto s 30 filter Download (CSV, Pareto Plot, Selected Alignment) Chart Alignments Reset Sorting

Gene	NM	TF	Region	Strand	MA Score	PSM Score	ID	MOTIF	CONSENSUS	Sim.	Pareto
INS	NM_000207	MNB1A	353-357	+	1.67	-0.58	1	CACTG	CACTG	1	1
INS	NM_000207	Ubx	299-302	+	1.45	-0.26	1	CACTG	CACTG	1	2
INS	NM_000207	Smal1	65-70	+	2	-0.28	1	CACTG	CACTG	1	2
INS	NM_000207	Ubx	162-165	+	1.45	-0.26	1	CACTG	CACTG	1	2
INS	NM_000207	Ubx	48-51	+	1.67	-0.26	1	CACTG	CACTG	1	0
INS	NM_000207	YY1	411-416	+	1.92	-1.08	1	CACTG	CACTG	0.92	7
INS	NM_000207	ZNFX2_1-4	21-26	+	0.8	-0.85	5	CACTG	CACTG	0.92	15
INS	NM_000207	ZNFX2_1-4	32-37	+	1.02	-0.85	4	CACTG	CACTG	0.92	10
INS	NM_000207	YY1	267-272	+	3.1	-1.06	1	CACTG	CACTG	0.92	0
INS	NM_000207	DoF2	353-358	-	1.78	-1	1	CACTG	CACTG	0.92	5
INS	NM_000207	TCF11-MaFG	29-34	-	2.01	-0.92	2	CACTG	CACTG	0.92	1
INS	NM_000207	ZNFX2_1-4	20-25	-	1.72	-0.85	3	CACTG	CACTG	0.92	2
INS	NM_000207	EBF	353-357	-	1.67	-0.58	1	CACTG	CACTG	0.9	2
INS	NM_000207	EBF	56-60	-	1.02	-0.81	1	CACTG	CACTG	0.9	9
INS	NM_000207	EBF	326-327	-	3.29	-2.11	1	CACTG	CACTG	0.92	1





## SHOE is connected to Garuda platform



## 自動で行う環境構築

- Makefile, makeに対応する形でVagrantfile, vagrant upなどのコマンドを使う。設定ファイルに必要なツールを自動ダウンロードができるようになっている。所定OSのインストールまで含めてすべて自動でやってくれる
  - OSイメージのインストールすらユーザ操作の必要がない
  - 開発環境を数行のコマンドを打つだけで終わらせる

```
$ git clone https://github.com/10up/varying-vagrant-vagrants.git
$ cd varying-vagrant-vagrants
$ vagrant up
(please wait for about twenty minutes)
$ open http://192.168.50.4/
```

# Future plans

- Create Plugin to CellDesigner to make a workflow
  - Percellome ->SHOE Transcription Regulation analysis ->CellDesigner transcription regulation map (in progress)
  - Percellome -> AGCT Cluster Analysis->SHOE Transcription Regulation analysis -> Transcription regulation map on CellDesigner
- Establish connections with other Garuda gadgets (pathway, literature mining)

25

トキシログenomixスミダシ学会 平成25年

# Collaborators

- Richard Nock (NICTA & the Australian National University, Canberra, Australia )
- Frank Nielsen (Sony Computer Science Laboratories Inc.)
- Kazuhiro Shibanaï (Tokyo Institute of Technology)
- Kodai Takata (Tokyo Institute of Technology)
- SBI Garuda Team

26

トキシログenomixスミダシ学会 平成25年



# **Xsight: an ensemble machine learning based platform for big data analysis in toxicology studies**

The Systems Biology Institute  
SBX

**We developed, “Xsight”, a computational learning system that integrates a large number of machine-learning models for big data analysis**



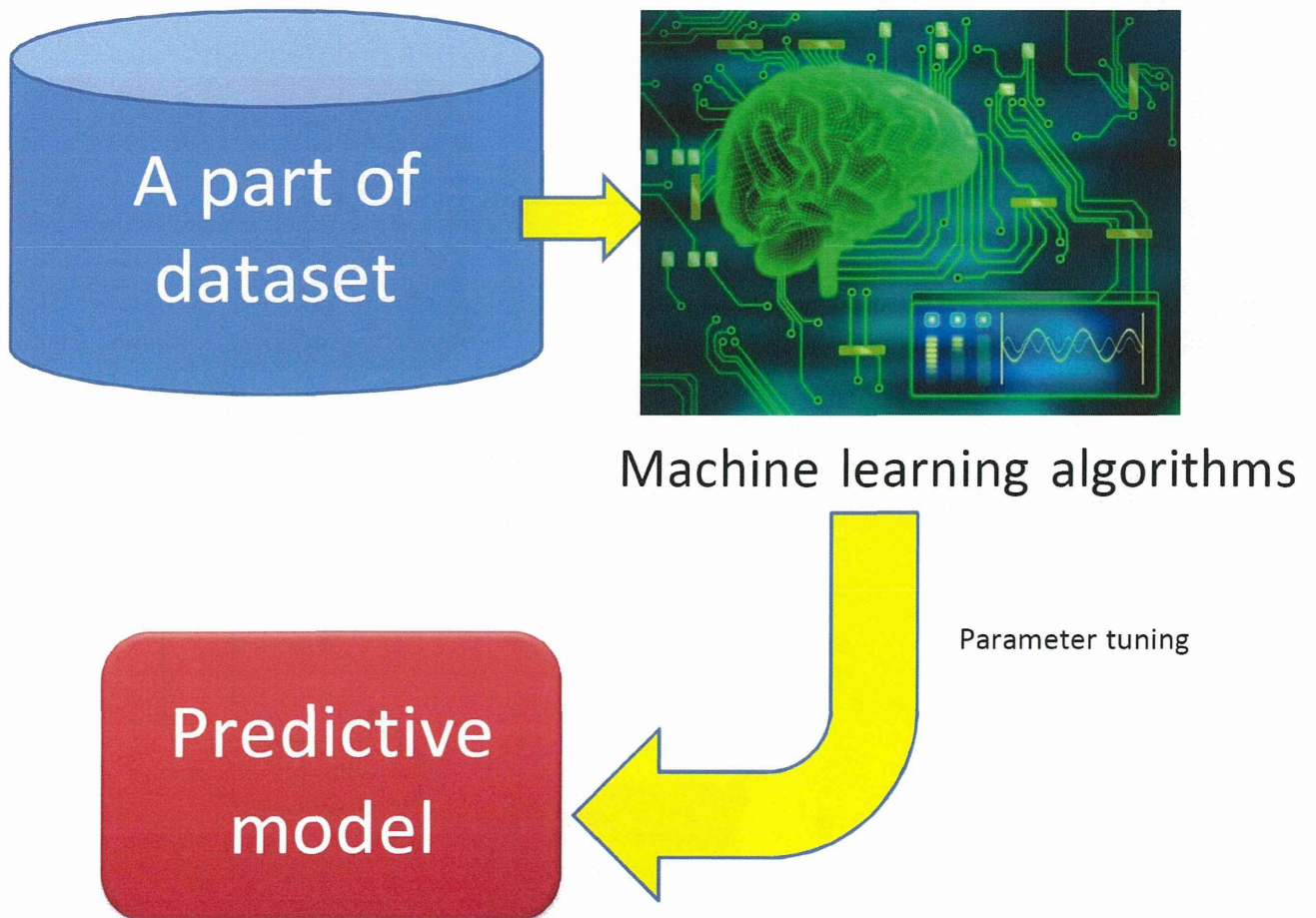
# Machine learning

- Procedure of machine learning is composed of two steps.

**Step 1: Training; Use a part of data to build predictive model**

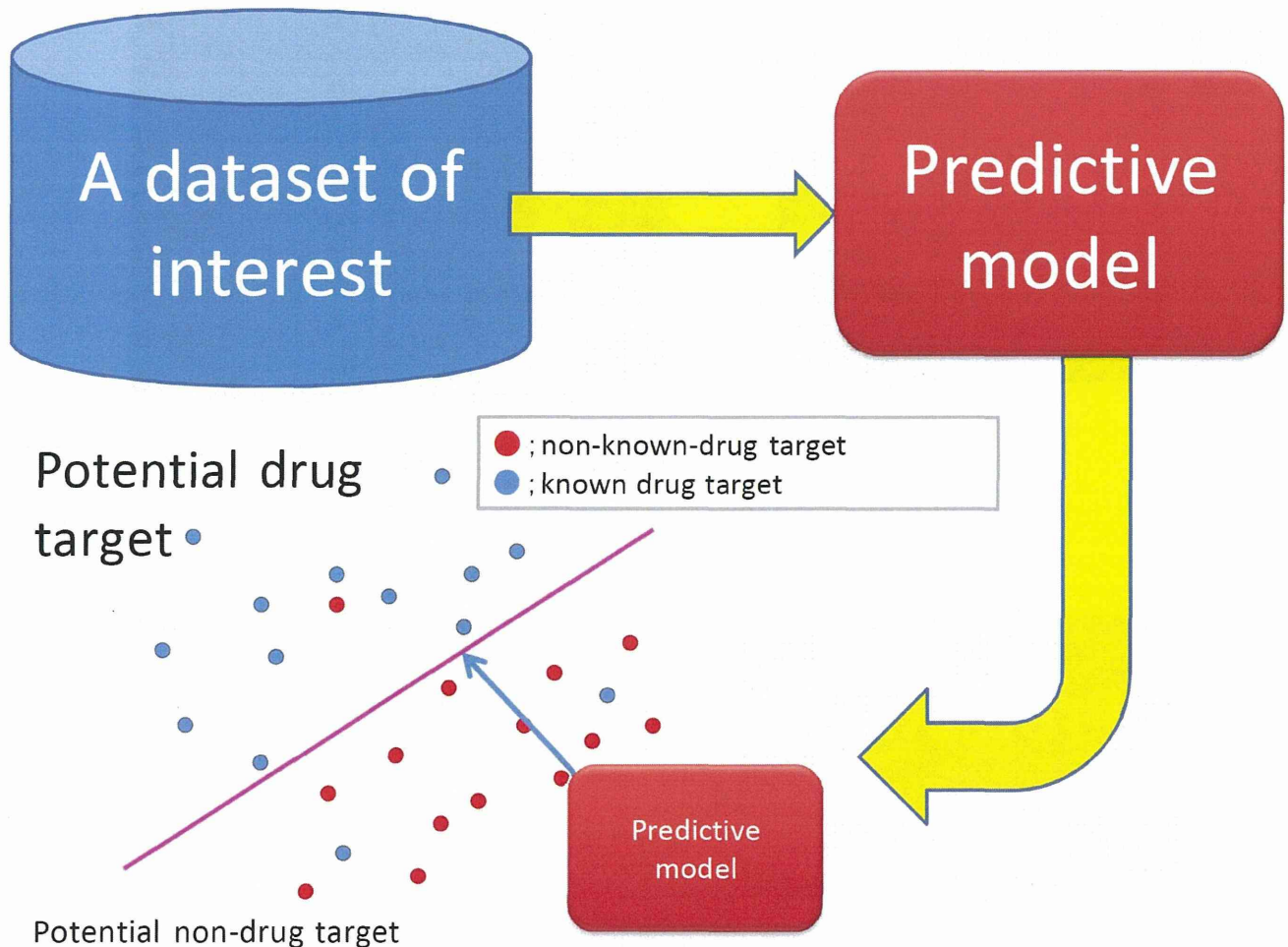
**Step 2: Test and Prediction; Using remaining dataset to test their predictive performance. Then, apply predictive model to a dataset of interest to make predictions, i.e., inference of potential side effect of a compound as well as potential candidate drug targets**

# Training: Model building





# Prediction: prediction based on models obtained from step 1



A large number of machine learning techniques based on various statistical techniques.

**One machine learning method will generate one predictive model.**

Bagged CART	Multivariate Adaptive Regression Splines
Bagged Flexible Discriminant Analysis	Naive Bayes
Boosted Generalized Additive Model	Nearest Shrunken Centroids
Boosted Generalized Linear Model	Neural Network
Boosted Linear Model	Neural Networks with Feature Extraction
Boosted Logistic Regression	Oblique Trees
Boosted Tree	partDSA
CART	Partial Least Squares
Conditional Inference Random Forest	Penalized Discriminant Analysis
Conditional Inference Tree	Penalized Linear Discriminant Analysis
Cost-Sensitive CART	Penalized Logistic Regression
Extreme Learning Machine	Penalized Multinomial Regression
Flexible Discriminant Analysis	Quadratic Discriminant Analysis
Gaussian Process	Quadratic Discriminant Analysis with Stepwise Feature Selection
Generalized Additive Model using Splines	Radial Basis Function Network
Generalized Linear Model	Random Forest
Generalized Linear Model with Stepwise Feature Selection	Regularized Random Forest
Generalized Partial Least Squares	Robust Regularized Linear Discriminant Analysis
Glmnet	Robust SIMCA
Greedy Prototype Selection	RCC-Based Classifier
High Dimensional Discriminant Analysis	Self-Organizing Maps
k-Nearest Neighbors	Shrinkage Discriminant Analysis
Learning Vector Quantization	SIMCA
Least Squares Support Vector Machine with Radial Basis Function Kernel	Sparse Partial Least Squares
Linear Discriminant Analysis	Stabilized Linear Discriminant Analysis
Linear Discriminant Analysis with Stepwise Feature Selection	Stochastic Gradient Boosting
Maximum Uncertainty Linear Discriminant Analysis	Support Vector Machines with Class Weights
Mixture Discriminant Analysis	Support Vector Machines with Linear Kernel
Model Averaged Neural Network	Support Vector Machines with Radial Basis Function Kernel
Multi-Layer Perceptron	Gaussian Process with Radial Basis Function Kernel
Multivariate Adaptive Regression Spline	Sparse Linear Discriminant Analysis

## Ensemble learning

Integrating a larger number of predictive model from a large number of machine learning methods to make predictions