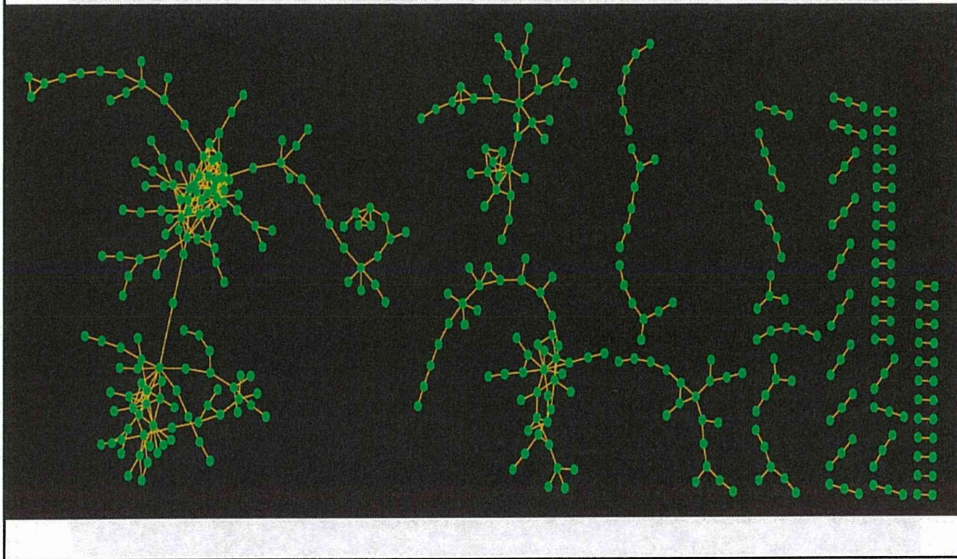
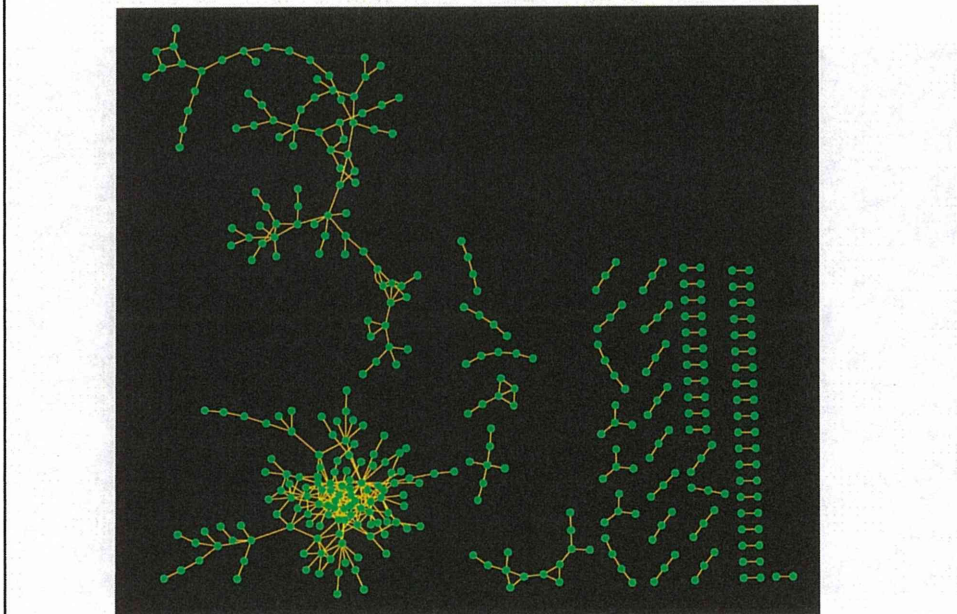


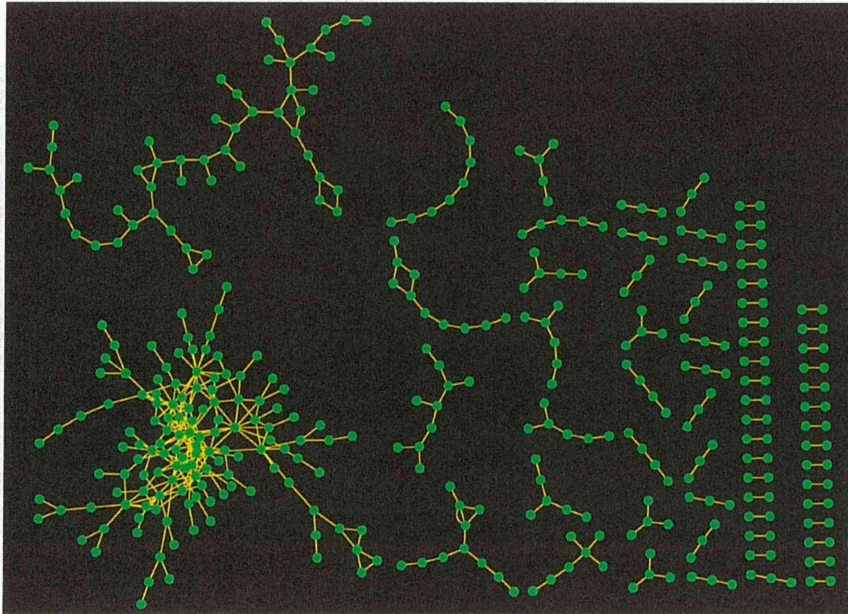
T1 specific gene regulatory network  
(時刻T1のみで現れるリンク)



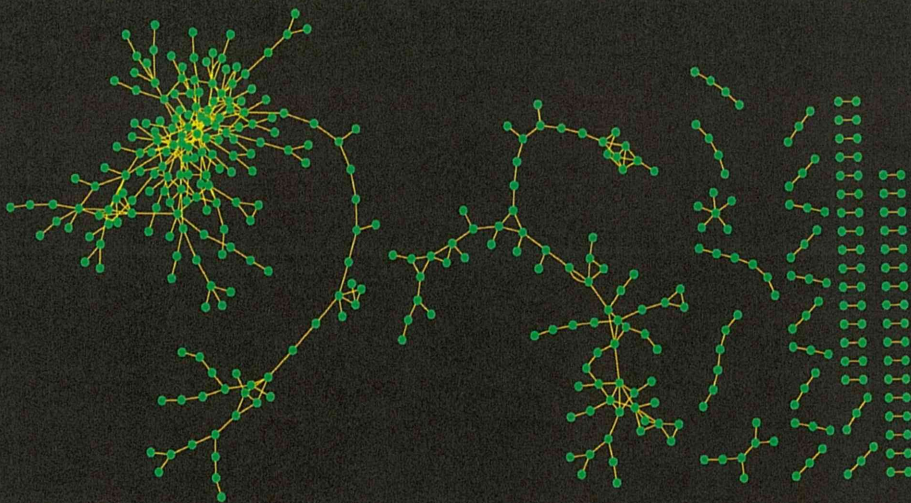
T2 specific gene regulatory network



### T3 specific gene regulatory network

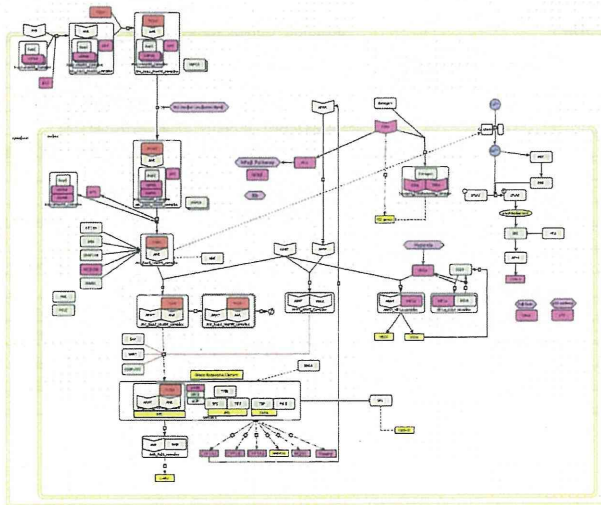


### T4 specific gene regulatory network

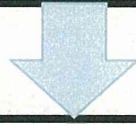


# Results3

パスウェイ情報を用いて重要な遺伝子セットを同定する。



パスイェ情報を用いて重要な遺伝子セツトを同定する。

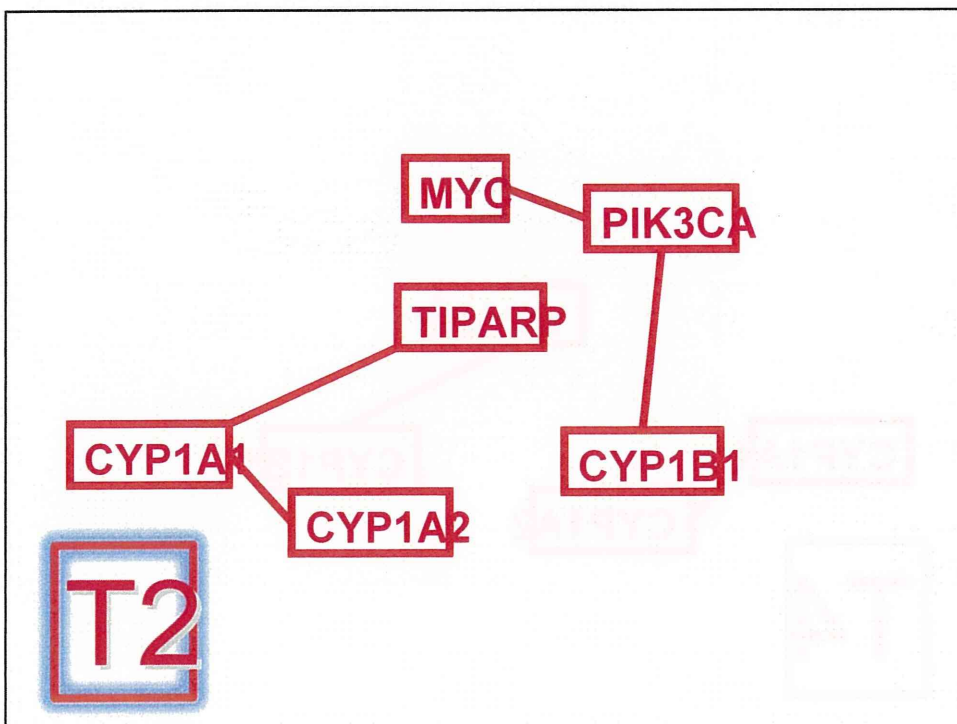
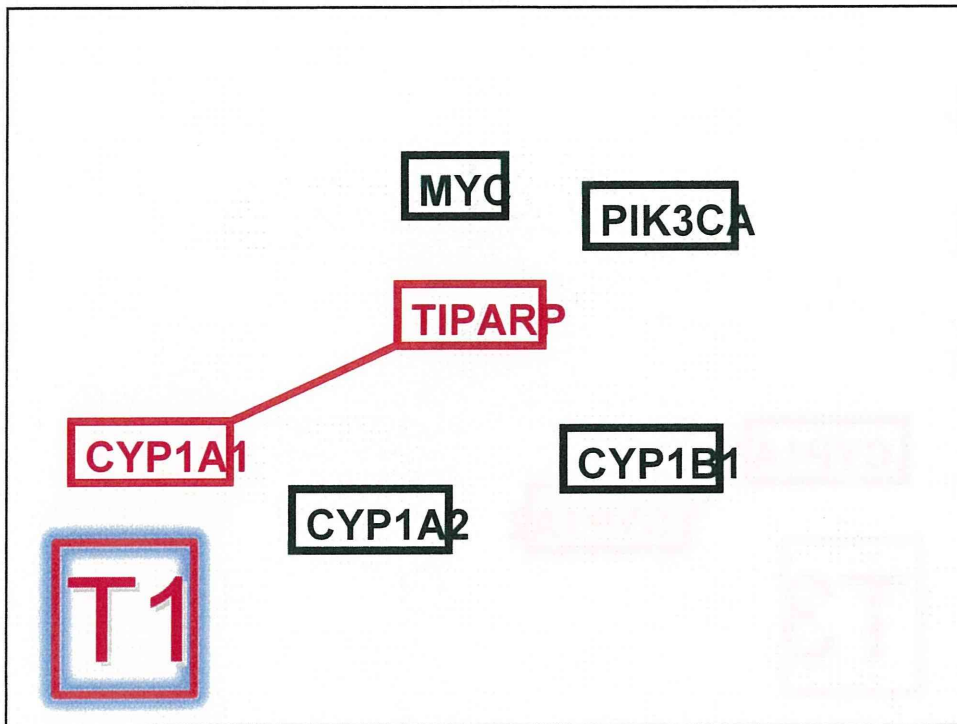


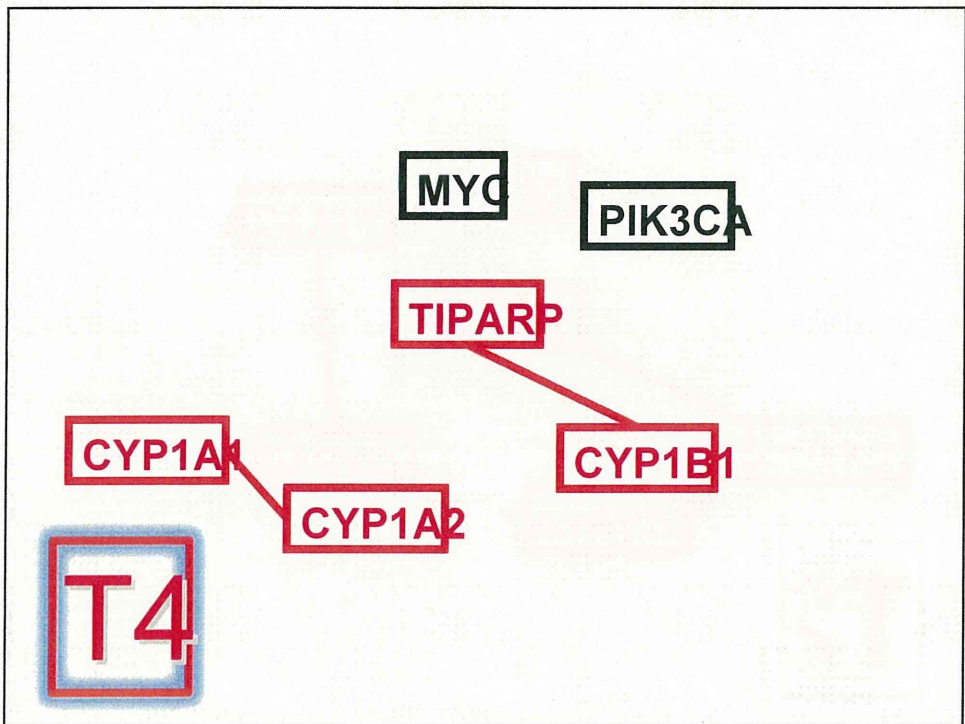
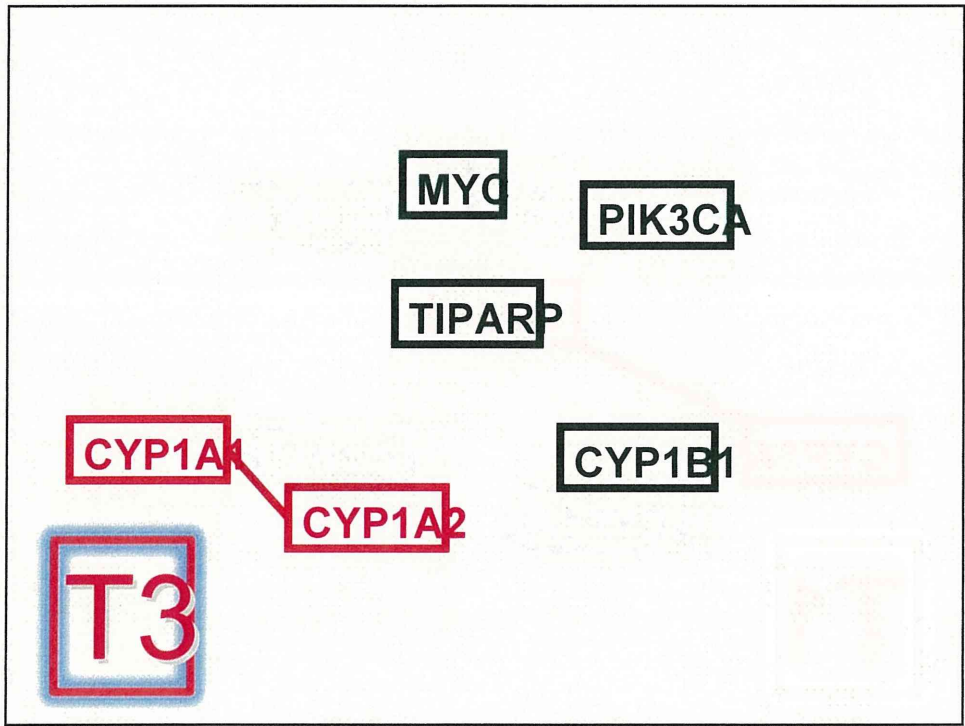
これらの遺伝子間のリンクが、TCDD投与後の時間経過に応じて、どのように変化するかを探索する。

以下、時間変化に応じたネットワークの変化の例を示す。

## 17 important genes

Symbol	Synonym	Description
UBE2J1	UBE2M=NEDD8	ubiquitin-conjugating enzyme E2 J1
TRIAP1	p53	TP53 regulated inhibitor of apoptosis 1
TIPARP		TCDD-inducible poly(ADP-ribose) polymerase
RRAS2	p23	related RAS viral (r-ras) oncogene homolog 2
PTGS2	COX2	prostaglandin-endoperoxide synthase 2
PIK3CA		phosphatidylinositol 3-kinase catalytic alpha polypeptide
nqo1		NAD(P)H dehydrogenasequinone 1
notch1	p300	Notch gene homolog 1 (Drosophila)
NFKBIA		nuclear factor of kappa light polypeptide gene enhancer in B-cells inhibitor alpha
MYC		myelocytomatosis oncogene
HSP90AA1	HSP90	heat shock protein 90 alpha (cytosolic) class A member 1
esrrg		estrogen-related receptor gamma
Epas1		HIF-1 alpha-like factor endothelial PAS domain protein 1
CYP1B1		cytochrome P450 family 1 subfamily b polypeptide 1
Cyp1a2		cytochrome P450 family 1 subfamily a polypeptide 2
CYP1A1	AHRR	cytochrome P450 family 1 subfamily a polypeptide 1
Cdk2		cyclin-dependent kinase 2





## Application of network reconstruction technique to Percellome

- 投与量に依存した、遺伝子間制御ネットワークの構造変化の探索。
- 組織特異的な遺伝子間制御ネットワークの構築
- ある投与量と組織において、遺伝子間制御ネットワークが時間とともにどのように変化するのか？



# 班会議平成24年2月9日

## クラスタリングとプロモータ解析によりネットワーク検出のアプローチ

株式会社ソニーコンピュータサイエンス研究所  
ポリリチャーフ・ナターリア

1

トキシコゲノミクス第3班会議

## Agenda

- \* **AGCT** A Geometric Clustering Tool (from 2008)
  - \* Clustering Percellome data based on similarity of gene expression profile. Application on TCDD and TCDF (2,3,7,8-Tetrachlorodibenzo-p-dioxin and 2,3,7,8-Tetrafuran) chemicals.

- \* Data normalization
- \* Unsupervised gene clustering
- \* Sorting clusters upon validity
- \* Issues remained



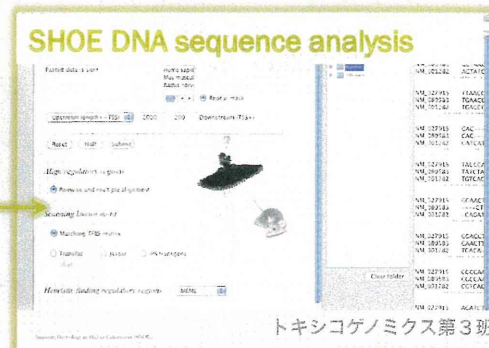
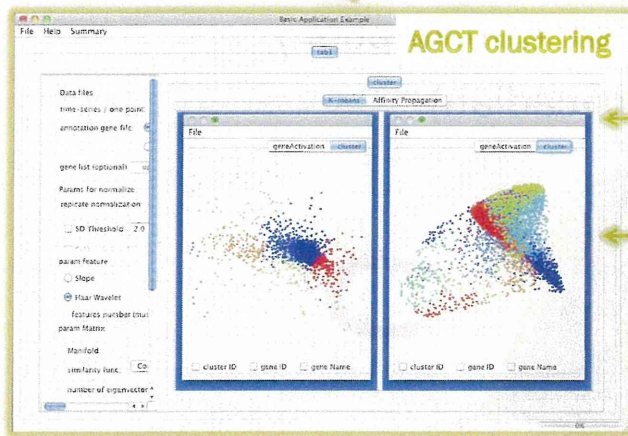
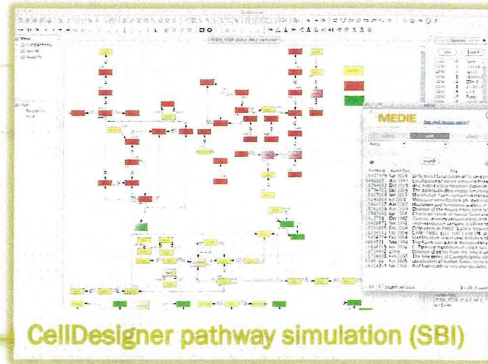
- \* **SHOE** Sequence Homology in Higher Eukaryote
  - \* Phylogenetic footprinting for discovery of transcription regulation network (from 2011)



2

トキシコゲノミクス第3班会議

# Pecellome analysis pipeline



Plugged-in

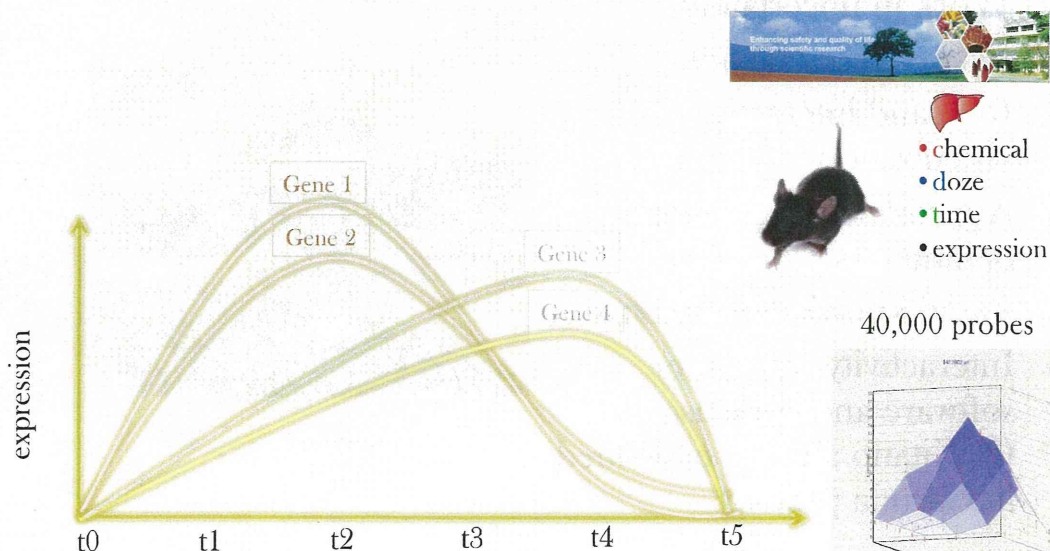
3

トキシコゲノミクス第3班会議

## AGCT A Geometric Clustering Tool

# Atlas of Cell Life by AGCT

- \* AGCT reconstructs gene network basing on the similarities of the expression profiles of genes



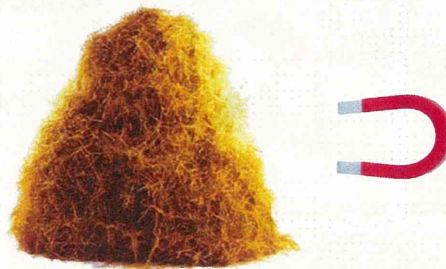
40,000 probes

トキシコゲノミクス第3班会議

4

## データの現状

- \* 全発現プロファイルが4万プローブを含む
  - \* 初期反応遺伝子 (responsive transcription)
  - \* Circadian 遺伝子 (basal transcription)
  - \* 無反応遺伝子 (not circadian)
  - \* 実験誤差



5

トキシコゲノミクス第3班会議



## Software: robust (independent but interactive )

- \* Fast
- \* Work in uncertainty
- \* Macro/micro flexible
- \* Rich and interactive visualization
- \* Approximate/generate parameters
- \* Rapid knowledge annotation
- \* Interactivity with other software and database (H.Kitano already explained in Garuda project)



6

トキシコゲノミクス第3班会議

# AGCT

## A Gene Geometric Clustering Tool

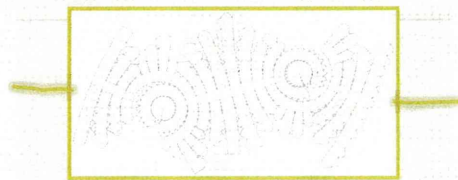
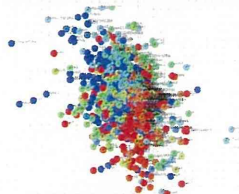
7

トキシコゲノミクス第3班会議

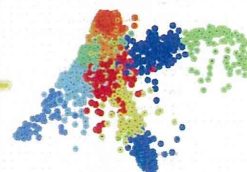
### Processing data on AGCT

1. 時系列データ前処理:線形回帰/ウェーブレット変換
2. 遺伝子間の類似度マトリックス
3. 低次元に落とすためにSpectral clusteringを行う。通常の主成分分析も行う。
4. 発見的なClustering法を使って構造上でデータの分割を行う。
5. 結果のinteractive visualizationやscenario 記録を行う。

PCA :  $M \times N$  matrix



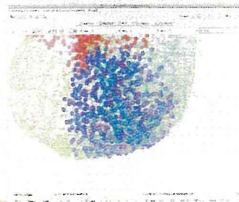
Spectral clustering:  
 $M \times M$  matrix



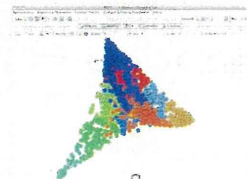
orthogonal matrix to compute  
one dimension per cluster/gene

### Examples of different network topologies

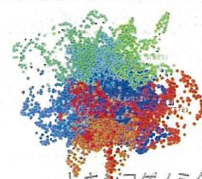
Mouse Stem cell



TCDD affected mouse liver cell



Influenza affected mouse bronchi cell



8

トキシコゲノミクス第3班会議

# AGCT interface 1

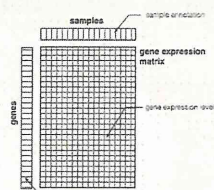
## Input file format (tab-delimited .txt)

$5 \times 1 \times 4 \times 3 = 60$  columns

Species	Mouse430.2.txt	Groups to specify						/* Constants												
E (entity)	5	Doze0	1	Doze0.1	1	Doze0.3	1	Doze10	1	Doze30	1	/* Groups								
EG (entity group)	1																			
T (time)	4		2		4		8		24											
R (replicate)	3																			
DATA	7881																			
1426408_at	13.88638	13.74065	10.62874	10.29534	12.14933	13.41148	16.33447	16.8473	17.75223	14.50626										
1430123_a_at	329.4958	304.0648	317.0229	320.1172	314.5138	331.6112	273.2262	351.2784	337.1702	303.972										
1427909_at	79.69938	65.37696	76.33909	74.55177	67.03818	75.51203	79.76383	76.7676	66.0889	65.07177										
1449231_at	1.998361	0.6321552	1.147684	0.8781132	1.167874	2.397003	2.210826	0.4704294	2.583703	0.3943034										
1419477_at	23.11418	13.63653	16.88693	7.034207	10.33751	11.61575	10.32749	9.681456	14.54885	8.335682										
1430115_at	4.780659	4.343996	3.392105	2.823855	3.745805	4.923362	2.997585	3.623773	5.434154	3.351972										
1433236_at	0.6044968	2.002586	1.370757	0.2526424	0.247167	0.2072143	0.4506208	1.558007	0.9185823	0.3346251										
1448996_at	1.737852	1.323068	0.7732058	1.240441	1.236961	0.4244971	0.6674667	0.6260773	0.9274884	0.4314171										
1458025_at	3.499104	2.880943	3.601012	4.745797	4.247523	4.485948	4.44686	3.311373	4.278707	5.641602										

© 2001 Nature Publishing Group <http://genetics.nature.com> commentary

Minimum information about a microarray experiment (MIAME)—toward standards for microarray data



# Gene ontology file

1429433.st	B1083827	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	gb B1083827 D B1083827	BAT2 domain cc Bat2c	226022	NM_001001680	Not defined	Not defined	Not defined		
1429434.st	B2E47269	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	gb AK0191841 B B2E47269	phosphatidylinositol 3-kinase	167068	NM_0038039	0006006	glo 0005942	glo 0004428	protein or phosphatase	
1429435.st	B2E47269	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	gb AK0191841 B B2E47269	phosphatidylinositol 3-kinase	167068	NM_0038039	0006006	glo 0005942	glo 0004428	protein or phosphatase	
1429436.st	B0826803	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	gb AK0191841 B B2E47269	PRP40 pre-mRNP Prpf40a	56194	NM_018785	0006337	mR 0006834	nuc 0005515	protein binding // inf	
1429437.st	B0826803	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	gb AK0191841 B B2E47269	PRP40 pre-mRNP Prpf40a	56194	NM_018785	0006337	mR 0006834	nuc 0005515	protein binding // inf	
1429438.st	AV318005	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	gb B0826803 C B0826803	BCL6 interactin Dicer	71458	NM_028510	0006350	trn 0006834	nuc 0003714	transcript on corpse	
1429439.st	BF020557	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	gb B0300720 C BF020557	excision repair Erc6b	71991	NM_028042	0002059	prc 0000109	nuc 0004842	ubiquitin protein lga	
1429440.st	BI734299	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	gb BI734299 D BI734299	Riken cDNA 181810041.15Rik	72301	XM_128189	Not defined	Not defined	Not defined	Not defined	
1429441.st	AK006369	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	gb BI686465 D AK006369	F-box protein 31 Fbxo30	71065	NM_027969	0005511	ubs 000	Not defined	0003270	zinc ion binding // h
1429442.st	AA560280	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	gb B1125490 C AA560280	RNA pseudourid Ribad2	271642	NM_173450	0001922	pas 000	Not defined	0003723	RNA binding // mfn
1429443.st	AK014396	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	74020	NM_028719	Not defined	Not defined	Not defined	Not defined	
1429444.st	AK004371	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	86995	NM_028804	0007264	sma 000	0003622	mtm 0000198	nucleotide binding //
1429445.st	AK014681	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	70994	NM_001081046	Not defined	Not defined	0005509	calcium ion binding //	
1429446.st	B0229980	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	62844	NM_025441	Not defined	0005834	Not defined	Not defined	
1429447.st	AK003753	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	62525	NM_145820	Not defined	0016820	mer 000	Not defined	
1429448.st	BI240219	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	52453	XM_125673	Not defined	Not defined	Not defined	Not defined	
1429449.st	BF136085	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	74400	NM_001037821	0043727	pos 001	0019717	syn 0000371	translation repressor
1429450.st	AK006308	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	71884	Not defined	Not defined	Not defined	Not defined	Not defined	
1429451.st	AK011950	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	61707	NM_026005	Not defined	Not defined	Not defined	Not defined	
1429452.st	AK007420	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	60776	NR_003516	0006654	pin 000	Not defined	0004909	phosphatidylserine d
1429453.st	BI202930	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	67212	NM_028035	0006412	trn 000	0005739	mtx 0003735	structural constituent
1429454.st	BI123170	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	86691	NM_026709	0006987	ens 000	0006822	mtx 0002085	guanylate nucleotid ex
1429455.st	BI123170	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	86691	NM_026709	0006987	ens 000	0006822	mtx 0002085	guanylate nucleotid ex
1429456.st	AV251549	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	20938	NM_025290	0006310	trn 000	0005624	nuc 0003999	DNA-directed RNA p
1429457.st	BD001825	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	76943	NM_175249	0006829	lad 000	0005578	satv 000	Not defined
1429458.st	AV083806	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	70249	XM_001476141	0007165	mul 000	0005976	etr 0004872	receptor activity // i
1429459.st	BB499147	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	70249	XM_001476141	0007165	mul 000	0005976	etr 0004872	receptor activity // i
1429460.st	AK019500	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	70249	XM_001476141	0007165	mul 000	0005976	etr 0004872	receptor activity // i
1429461.st	AK013101	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	70249	XM_001476141	0007165	mul 000	0005976	etr 0004872	receptor activity // i
1429462.st	AK011759	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	70249	XM_001476141	0007165	mul 000	0005976	etr 0004872	receptor activity // i
1429463.st	BB612385	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	69908	NM_172402	0006810	trn 000	0005739	mtx 0002919	transporter activity
1429464.st	AK005847	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	100079	NM_178143	0006468	prc 000	0006834	nuc 0000168	nucleotide binding //
1429465.st	AK002310	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	100079	NM_178143	0006468	prc 000	0006834	nuc 0000168	nucleotide binding //
1429466.st	BB118542	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	66328	NM_026400	0002152	met 000	Not defined	0005488	stereocamer activity
1429467.st	AK003775	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	13487	NM_021353	0006919	trn 000	0005987	mtx 0002515	transporter activity
1429468.st	BB498931	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	66328	NM_026400	0002152	met 000	Not defined	0005488	stereocamer activity
1429469.st	AK014910	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	66328	NM_026400	0002152	met 000	Not defined	0005488	stereocamer activity
1429470.st	AK003742	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	66328	NM_026400	0002152	met 000	Not defined	0005488	stereocamer activity
1429471.st	AK014969	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	66328	NM_026400	0002152	met 000	Not defined	0005488	stereocamer activity
1429472.st	AV292561	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	66328	NM_026400	0002152	met 000	Not defined	0005488	stereocamer activity
1429473.st	AV292561	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	66328	NM_026400	0002152	met 000	Not defined	0005488	stereocamer activity
1429474.st	BE283373	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	66328	NM_026400	0002152	met 000	Not defined	0005488	stereocamer activity
1429475.st	AK013361	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	66328	NM_026400	0002152	met 000	Not defined	0005488	stereocamer activity
1429476.st	CG063818	Not defined	Mus musculus	11-Mar-09	Consensus seq, GenBank	Not defined	Not defined	66328	NM_026400	0002152	met 000	Not defined	0005488	stereocamer activity

Use NCBI GEO annotation file as is, but substitute blank fields with 'Not\_defined'

# AGCT interface 1

The screenshot shows the AGCT3 application interface. The 'Data files' section has 'time-series / one point' set to 'upload' with the path '/s/NataP/Desktop/AGCT1.0\_demo\_2/TCDD\_sel\_788'. The 'annotation gene file' is set to 'load from library' with 'Human' selected. The 'gene list (optional)' is set to 'upload' with the path '/Users/NataP/Desktop/AGCT1.0\_demo\_2/Labels'. The 'Params for normalize' section has 'replicate normalization' set to 'Geometric Mean' and 'SD Threshold' set to '2.0'. The 'param feature' section has 'Slope' selected. The 'param Matrix' section has 'number of eigenvectors' set to '20'. The 'Manifold' section has 'similarity func' set to 'Cosine Similarity'. The 'cluster' section shows a table with columns 'File', 'Filter', and 'Search'. The table contains the following data:

File	Filter	Search
Cl	N	C
0	1	0
1	462	0
2	553	0
3	214	0
4	981	0

java -jar -Xmx1G(40G) -Xms1G(40G) AGCT3.jar

# Pre-processing: Slopes

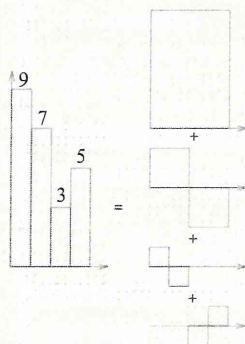
A **time series** for gene **i** and ligand **j** is mapped to a slope **s<sub>ij</sub>** using a conventional linear regression fit:

$$S_{ij} = \frac{\text{cov}(t, x_{ij})}{\text{var}(t)} = \frac{\sum_k (t_k - \text{avg}(t))(x_{ijk} - \text{avg}(x_{ij}))}{\sum_k (t_k - \text{avg}(t))^2}$$

Here, “avg” denotes the average, “var” is the variance and “cov” is the covariance; furthermore, *t* denotes the set of time stamps and *x<sub>ij</sub>* denotes the set of measurements for gene *i* and ligand *j*. *k* spans {2, 4, 8, 24}.

# Pre-processing: Haar wavelet transform

The Haar wavelet is the first known **wavelet** and was proposed in **1909** by **Alfred Haar**.



One-dimensional Haar wavelet transform

Resolution	Averages	Detail coefficients
4	[ 9 7 3 5 ]	
2	[ 8 4 ]	[ 1 -1 ]
1	[ 6 ]	[ 2 ]

wavelet transform (filter bank)

$$\begin{bmatrix} 6 & 2 & 1 & -1 \end{bmatrix}$$

basis coefficients

Figure 1: Decomposition of a small histogram (left) into a set of four Haar wavelets components (right). Taking a smooth curve instead of an histogram would necessitate much more coefficients to get an appropriate fitting.

## Using the vector space from linear algebra

We can think of histogram as piecewise-constant functions on the half-open intervals  $[0, 1]$ .

We'll let  $V^0$  be the vector space of all functions, i.e., **block function**:

$$\phi(t) = \begin{cases} 1 & \text{if } t \in [0, 1) \\ 0 & \text{otherwise} \end{cases}$$

and **mother wavelet** :

$$\psi(t) = \begin{cases} 1 & \text{if } t \in [0, 1/2) \\ -1 & \text{if } t \in [1/2, 1) \\ 0 & \text{otherwise} \end{cases}$$

The space  $V^j$  will include all piecewise constant functions defined on the interval  $[0, 1]$  with constant pieces over each of  $2^j$  equal subintervals. Thus, the spaces  $V^j$  are nested as:

$$V^0 \subset V^1 \subset V^2 \subset \dots$$

Any curve can be obtained (eventually up to any finite precision) by the sum of translated/dilated curves that are born from (2), (3):

$$\phi_i^j = \sqrt{2^j} \phi(2^j t - i), \quad j \in \mathbb{N}, i = 0, 1, \dots, 2^j - 1 \quad (4)$$

$$\psi_i^j = \sqrt{2^j} \psi(2^j t - i), \quad j \in \mathbb{N}, i = 0, 1, \dots, 2^j - 1. \quad (5)$$

What we want to do is to compute each time series as a sum of such coefficients, and more precisely, as a sum of the form  $\alpha \phi_0^0 + \sum_{i,j} \beta_i^j \psi_i^j$ ,  $\alpha \in \mathbb{R}, \beta_i^j \in \mathbb{R}$ . We have normalized the functions ( $\int_0^1 f^2(t) dt = 1$ ) for orthonormality reasons, but if we had chosen to remove the square-roots above,  $\phi_0^0$  would just be the global average of the curve. Since the mother wavelet is the curve which gives an indication of the slopes at each scale-dilation factors, this first "average" coefficient is removed from the analysis.

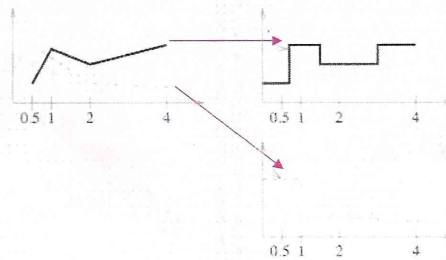


Figure 2: Curves are best fit to histograms on a time step tailored to have exact Haar approximation with few coefficients.



# Similarity matrix

param Normalize  
 replicate normalization: Arithmetic Mean

SD Threshold 2.0  
 Subtract circadian

param feature  
 Slope  
 Haar Wavelet

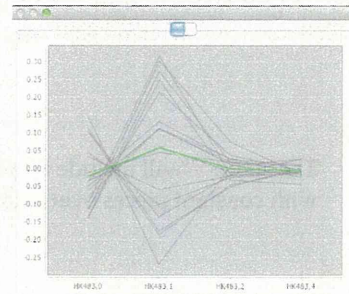
features number (must be power of 2) 8

param Matrix  
 Manifold:  
 similarity fun: Heat Kernel  
 Cosine Similarity  
 Absolute CS  
 number of eigenvectors: 20

clustering

math. artifact serves Biology

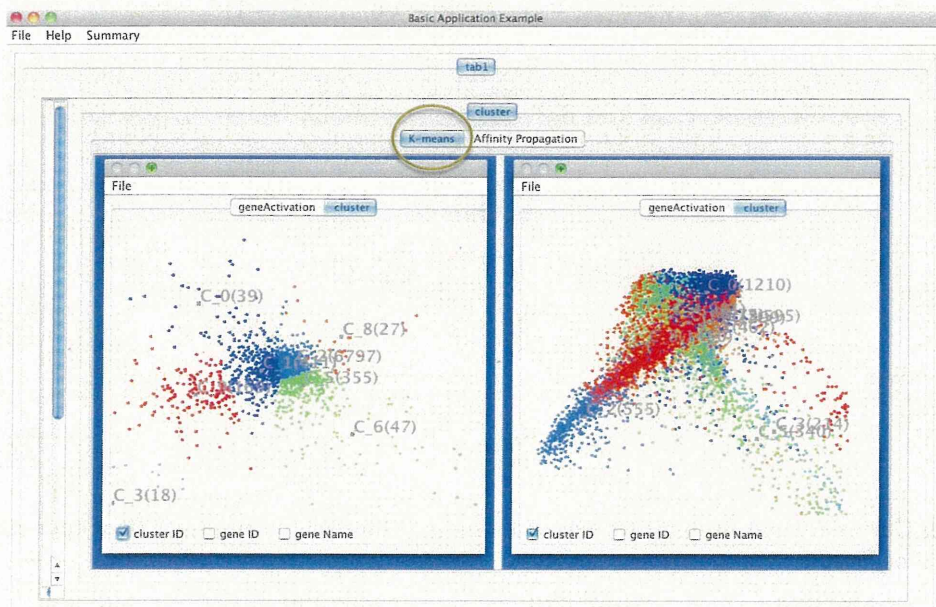
$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n (A_i)^2} \times \sqrt{\sum_{i=1}^n (B_i)^2}}$$



# AGCT Results Panel 2

Principal Component structure

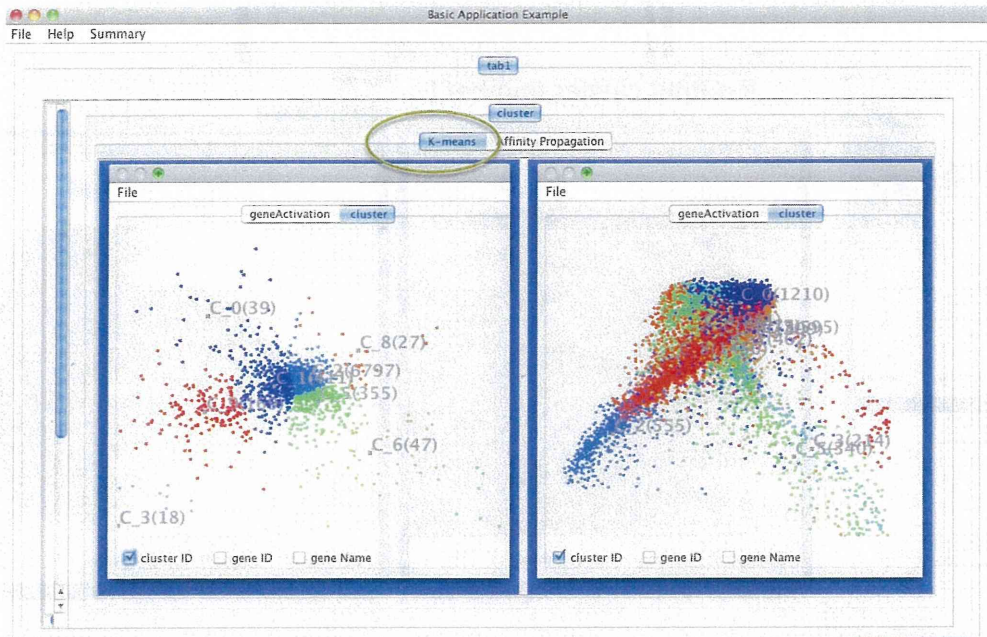
Spectral clustering structure (MANIFOLD)



# AGCT K-means Results 1

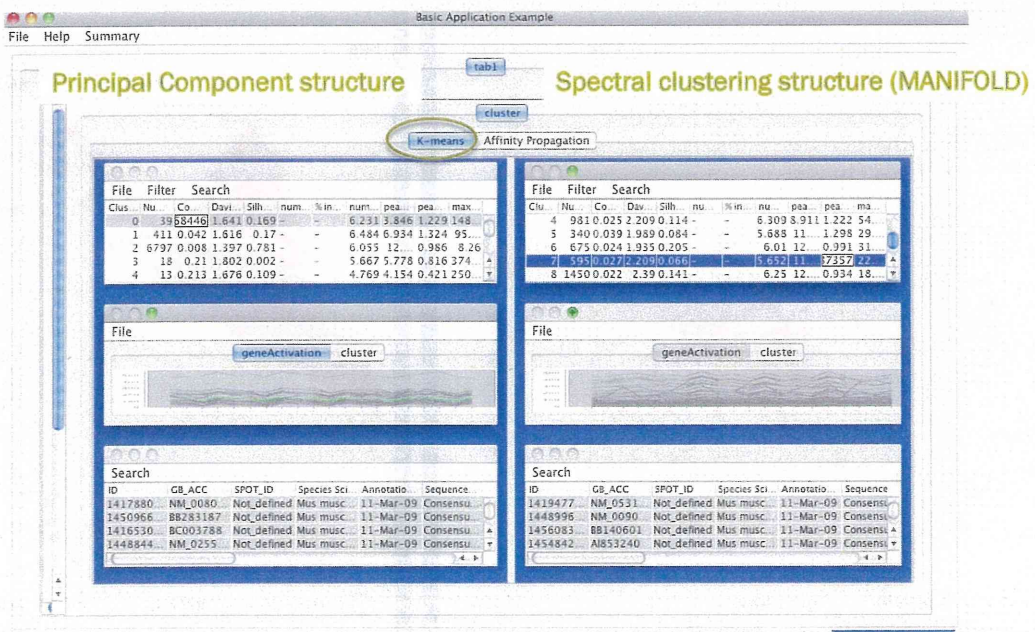
Principal Component structure

Spectral clustering structure (MANIFOLD)

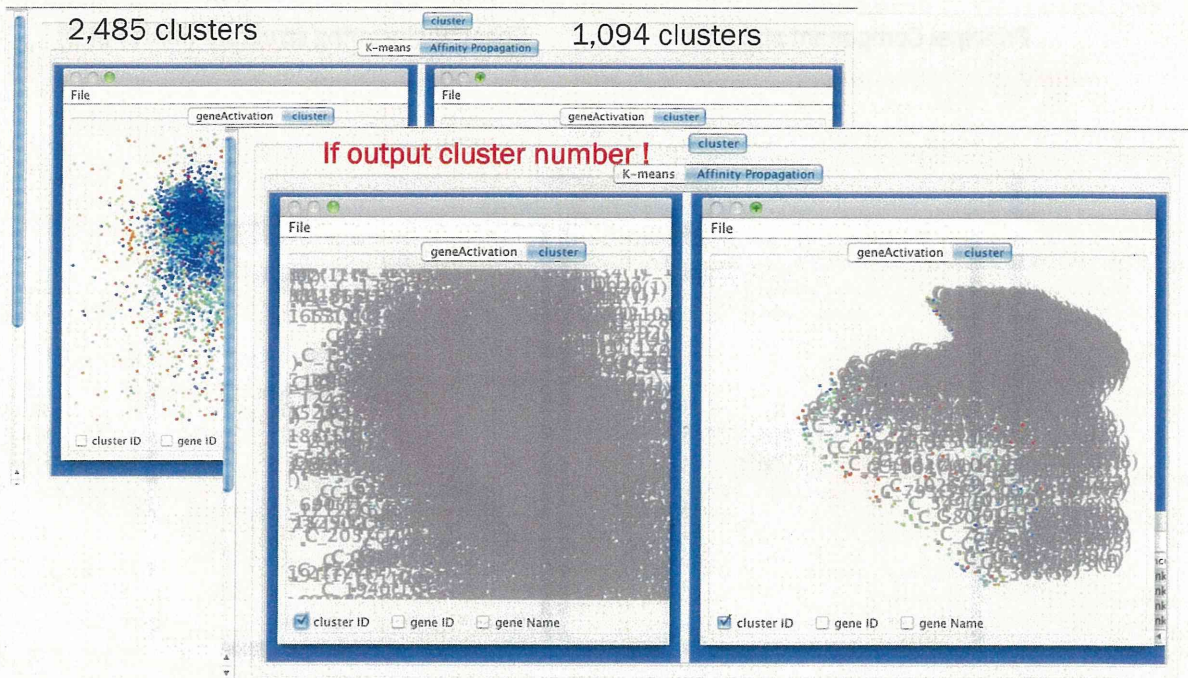


3 班会議

# AGCT K-means Results 2



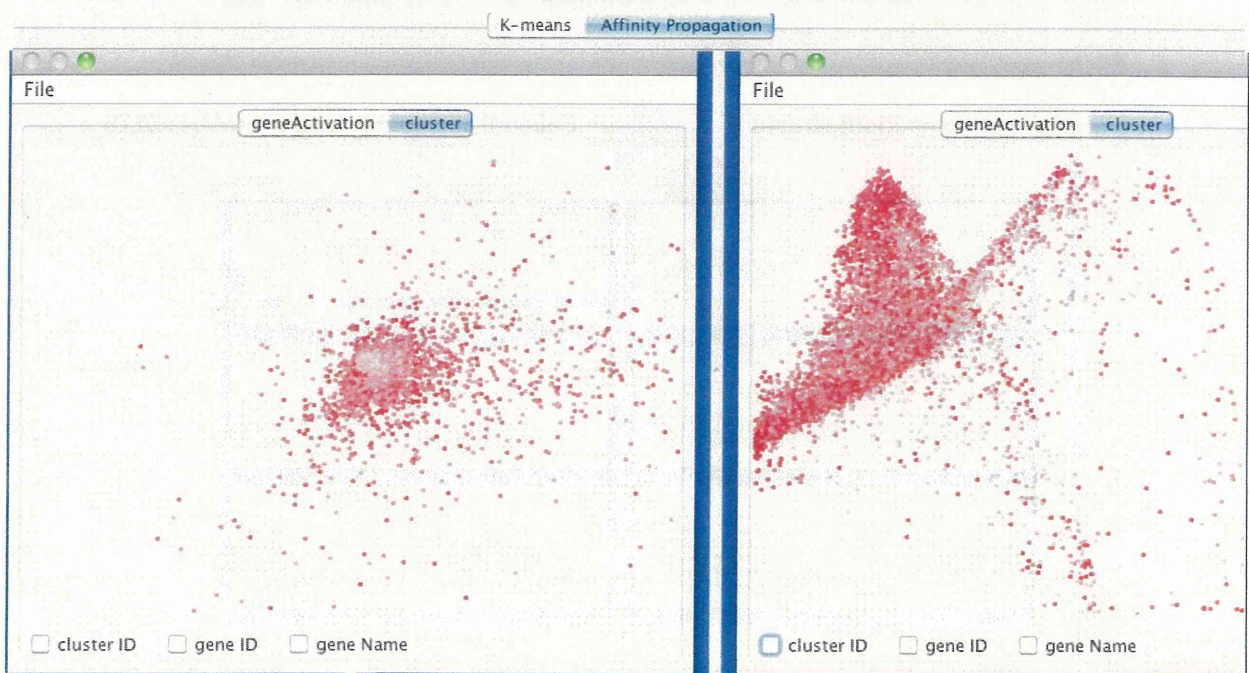
# AGCT Affinity Propagation (AP) Results 1



21

トキシコゲノミクス第3班会議

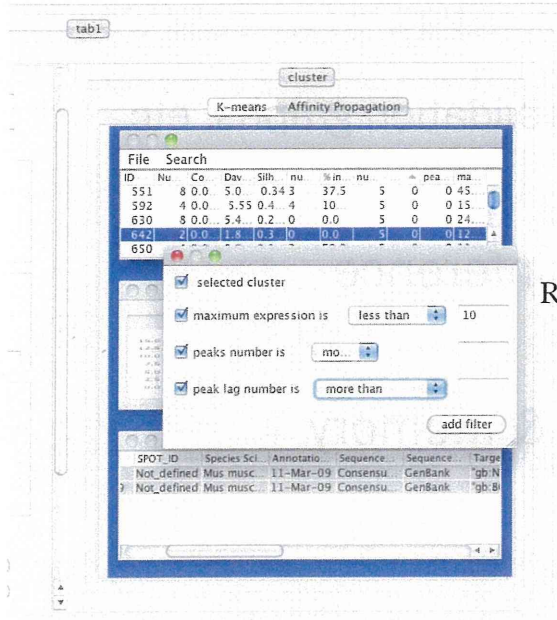
# AGCT Cluster Validation Panel



22

トキシコゲノミクス第3班会議

# Filtering data



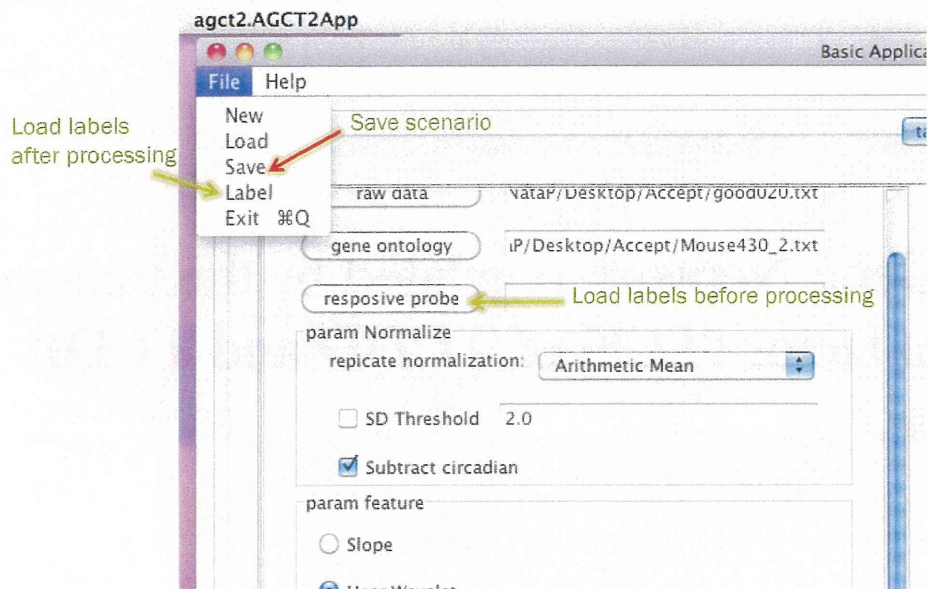
Right click on mouse  
and filtering menu shows up !

- Data below the thresholds will disappear. You keep the rest of data in clusters.
- Filtering can be done on selected/all clusters.

23

トキシコゲノミクス第3 班会議

# Scenario



24

トキシコゲノミクス第3 班会議