

## Future works on AGCT

- \* Visualization
  - \* 3D graphs, PCA balls, validation graphs, etc.
- \* Methodology : Cluster inference
- \* Tests on performance & memory

**Finding Gene Network regulated by the toxicity equivalent factor (TEF) of TCDD and TCDF chemicals .**

# Algorithms for clustering

## \* Affinity propagation, Bergman K-means, EM, Hierarchical clustering

Computer global optimum with small and sensitive clusters.

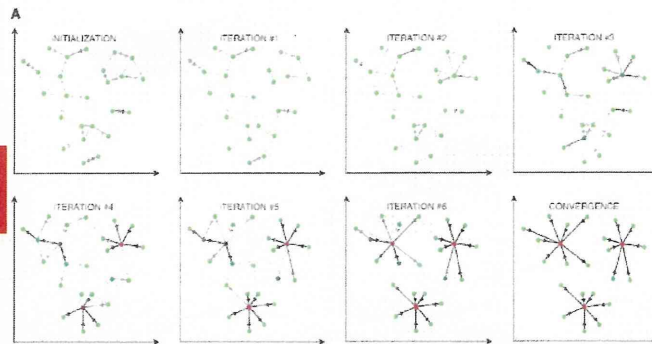


Fig. 1. How affinity propagation works. (A) Affinity propagation is illustrated for two-dimensional data points, where negative Euclidean distance (squared error) was used to measure similarity. Each point is colored according to the current evidence that it is a cluster center (exemplar). The darkness of the arrow directed from point  $i$  to point  $k$  corresponds to the strength of the transmitted message that point  $i$  belongs to exemplar point  $k$ . (B) "Responsibilities"  $r(i,k)$  are sent from data points to candidate exemplars and indicate how strongly each data point favors the candidate exemplar over other candidate exemplars. (C) "Availabilities"  $a(i,k)$  are sent from candidate exemplars to data points and indicate to what degree each candidate exemplar is available as a cluster center for the data point. (D) The effect of the value of the input preference (common for all data points) on the number of identified exemplars (number of clusters) is shown. The value that was used in (A) is also shown, which was computed from the median of the pairwise similarities.

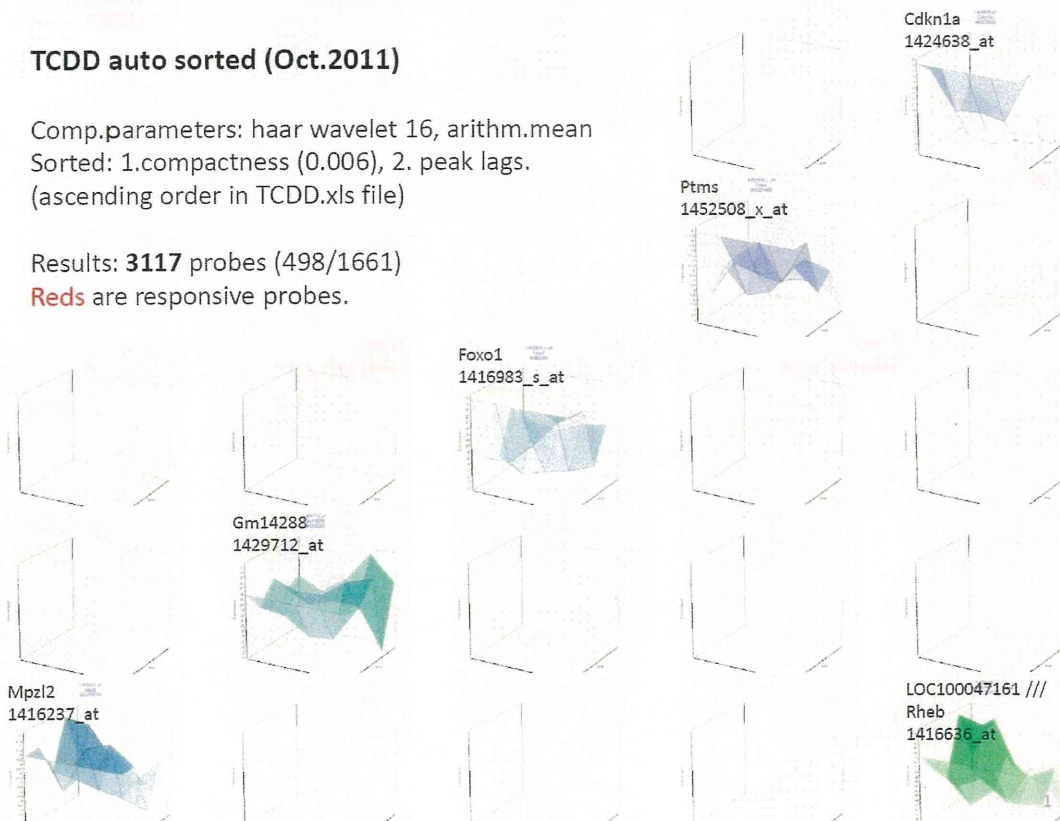
27 www.sciencemag.org SCIENCE VOL 315 16 FEBRUARY 2007

トキシコゲノミクス第3班会議 973

### TCDD auto sorted (Oct.2011)

Comp.parameters: haar wavelet 16, arithm.mean  
Sorted: 1.compactness (0.006), 2. peak lags.  
(ascending order in TCDD.xls file)

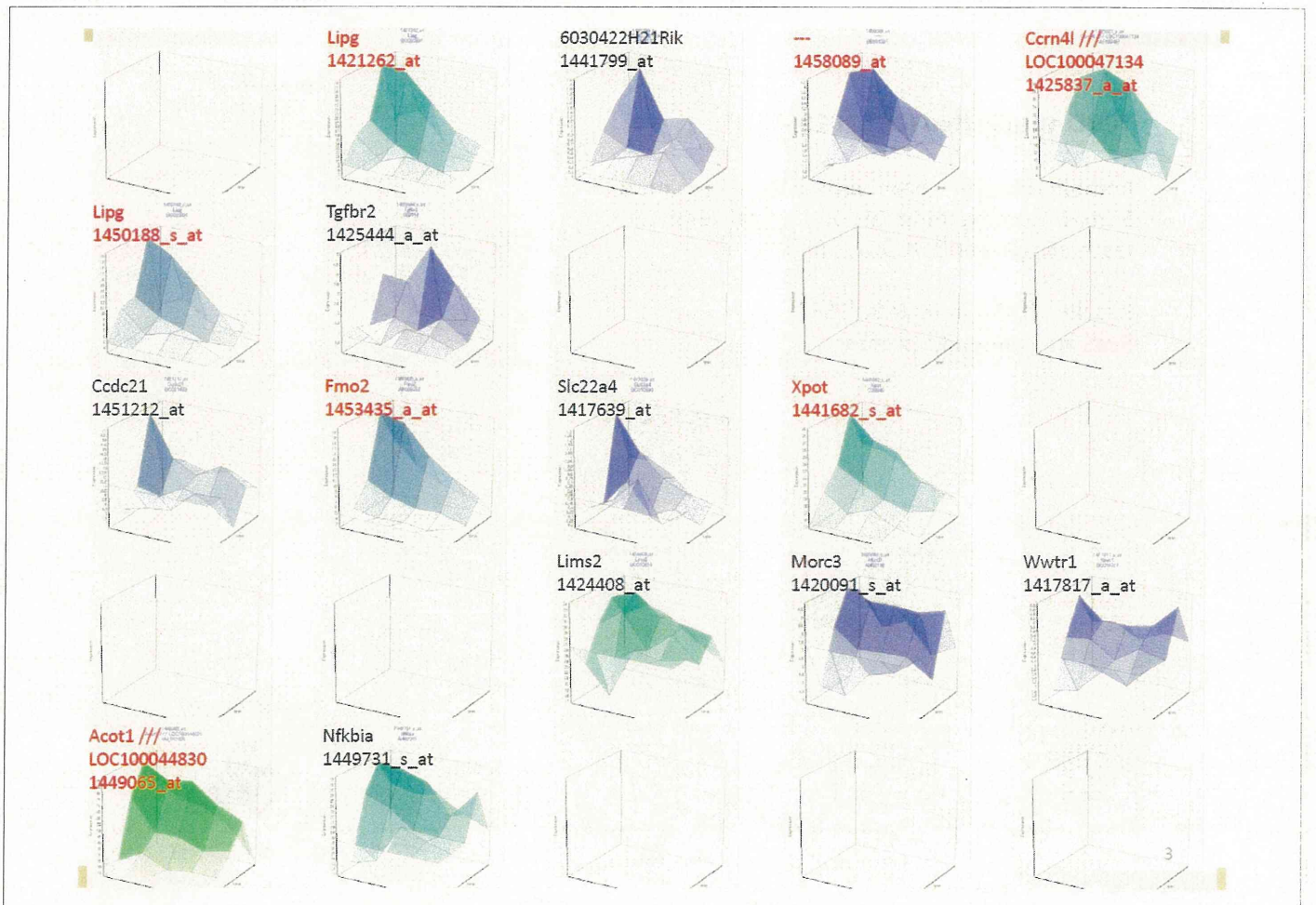
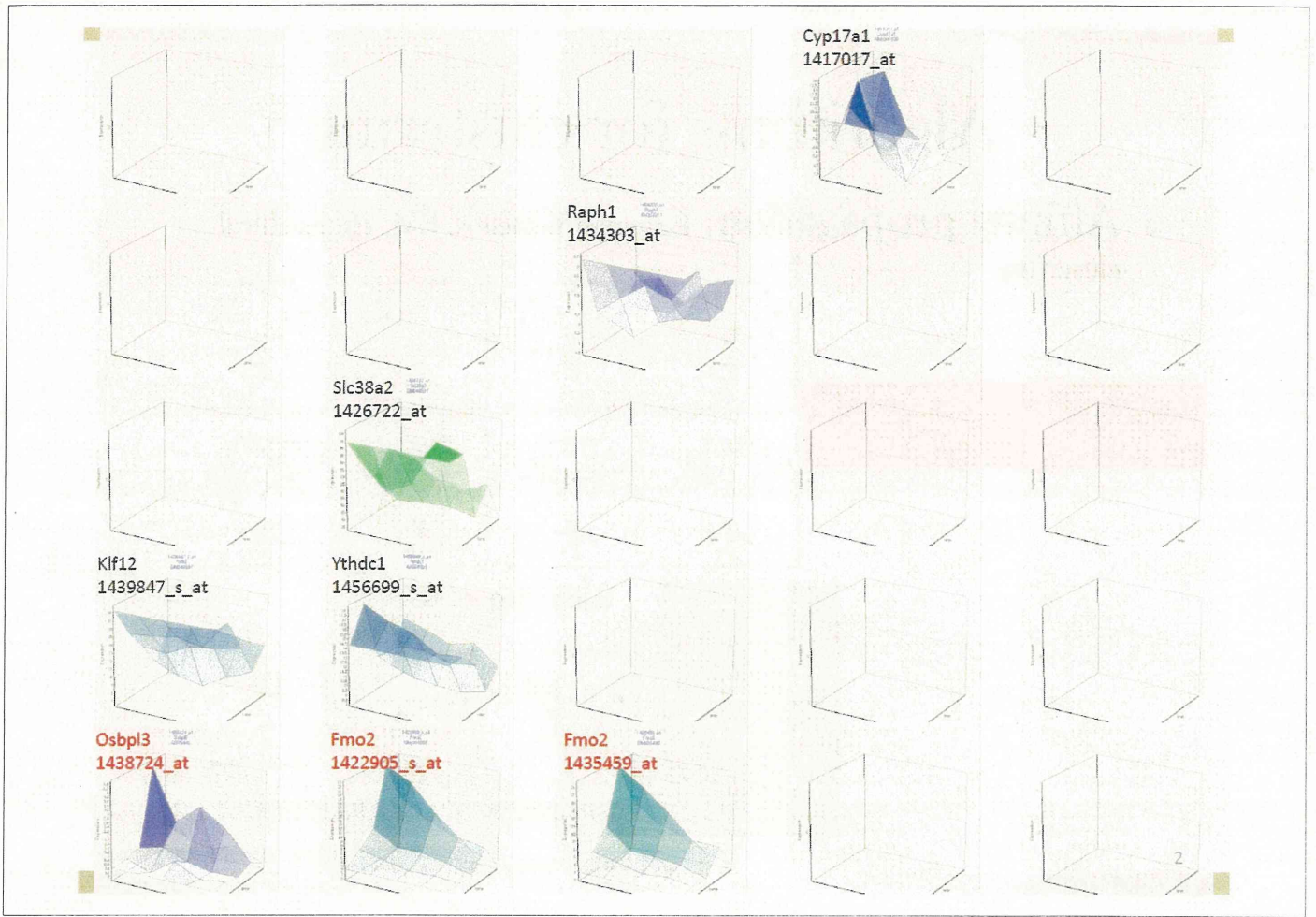
Results: 3117 probes (498/1661)  
Reds are responsive probes.



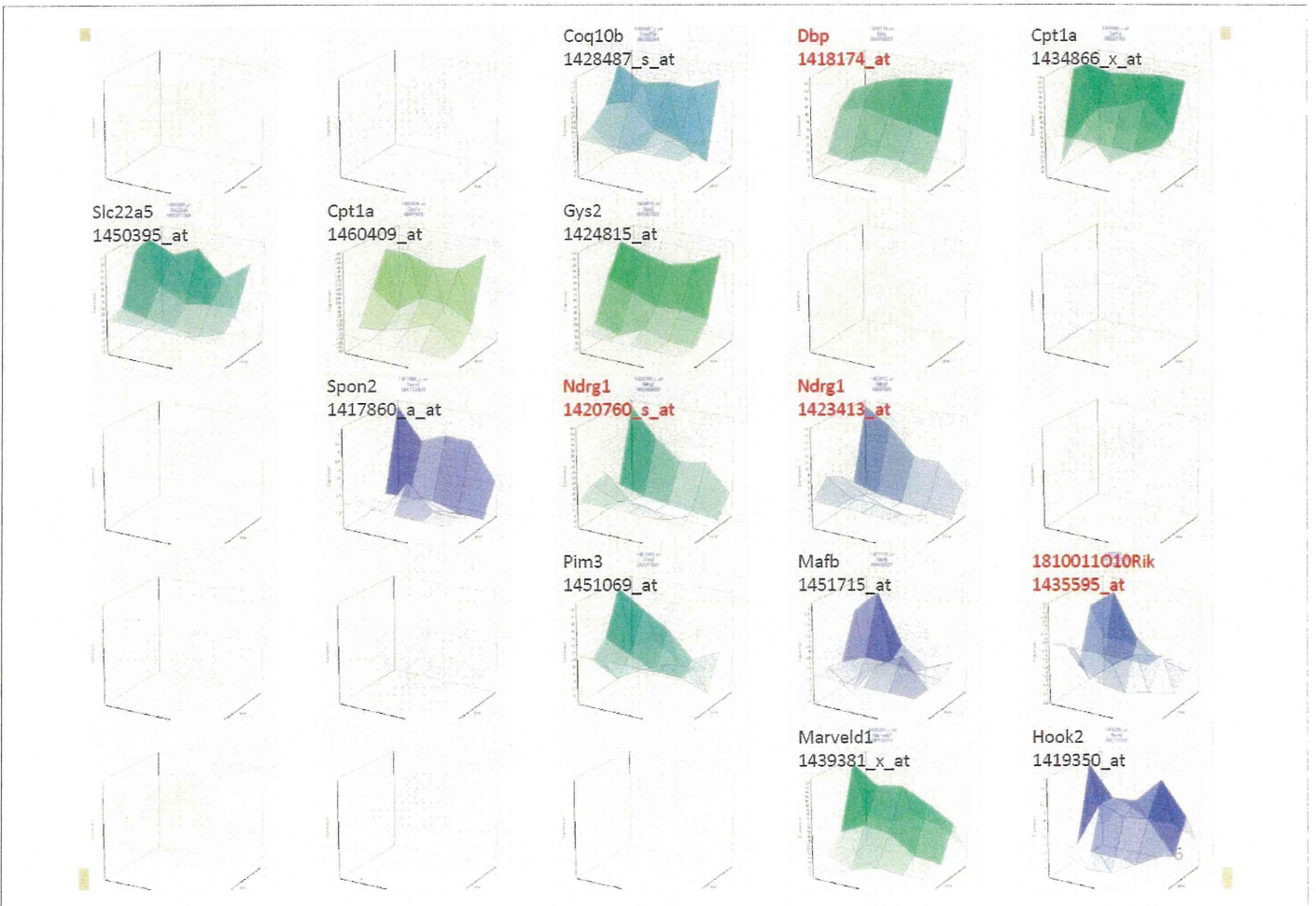
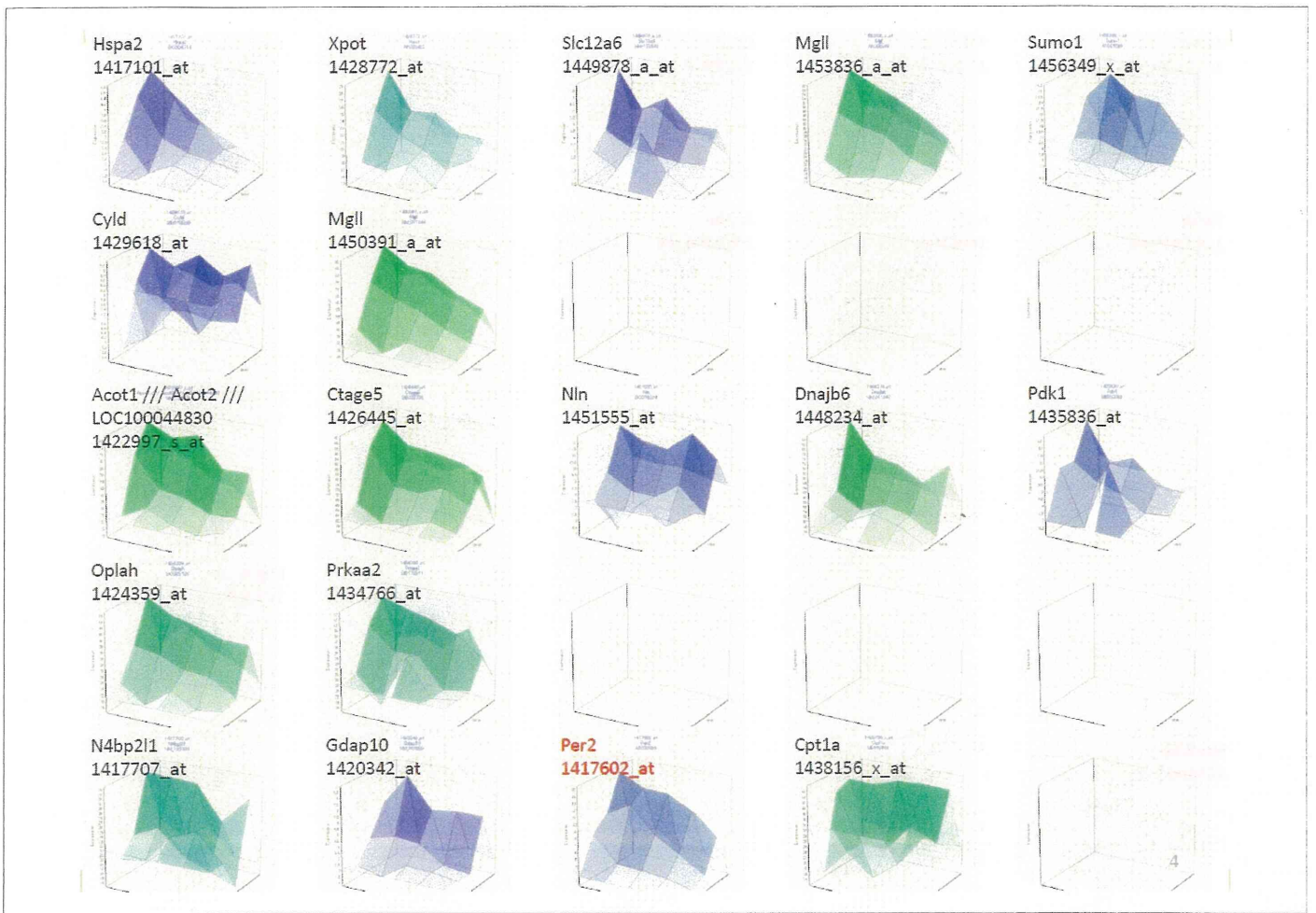
28

トキシコゲノミクス第3班会議

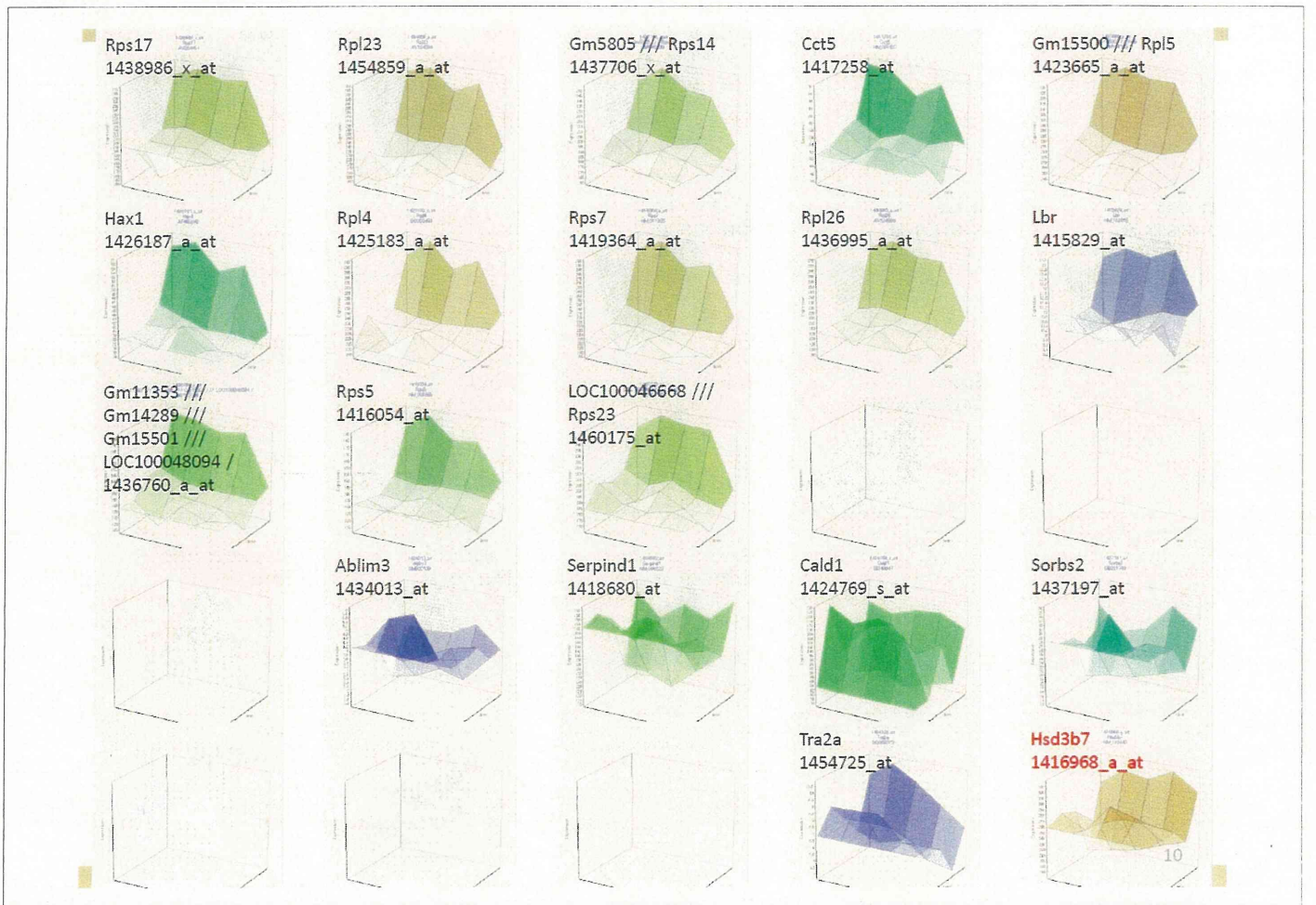
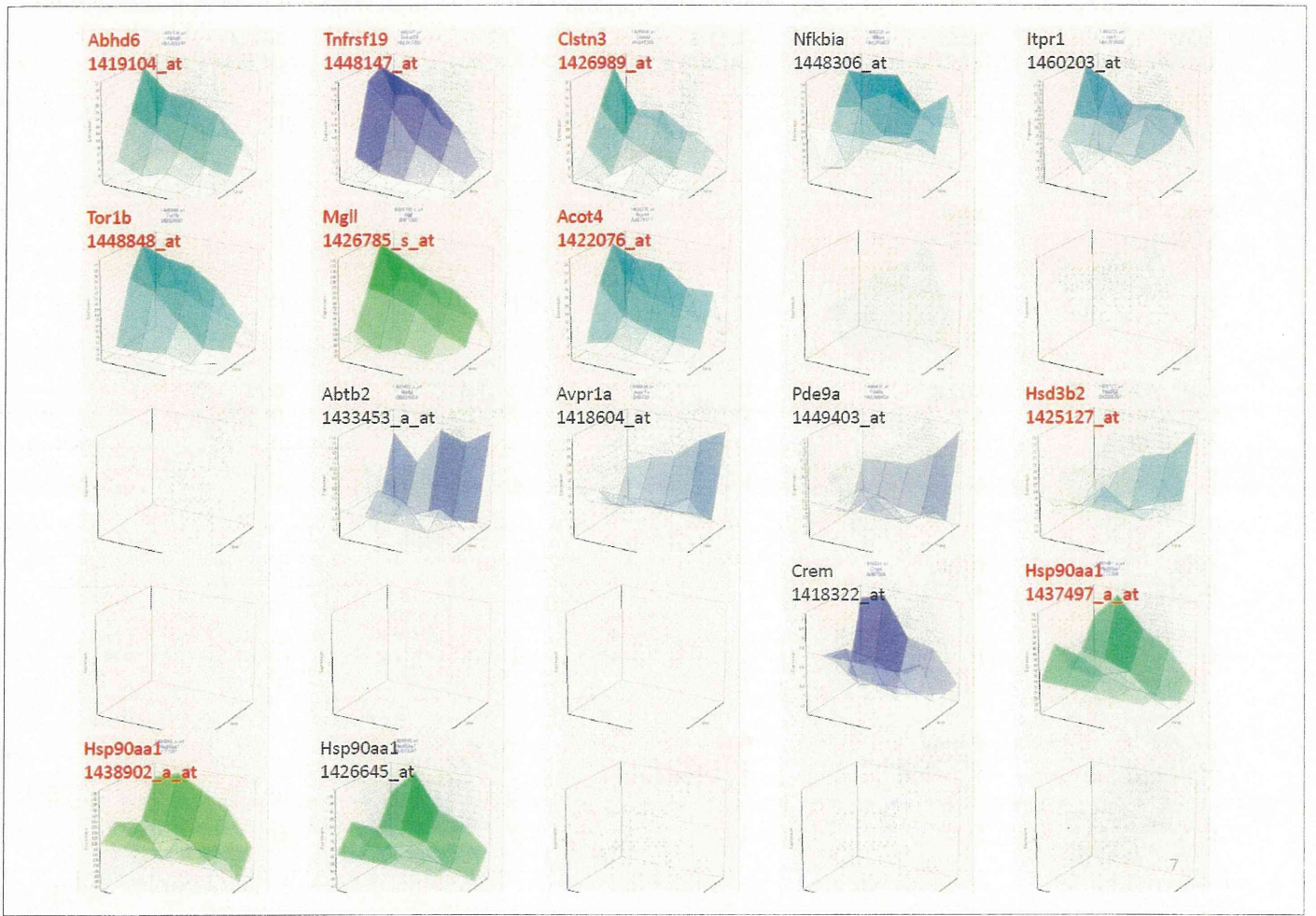














**TCDF auto sorted (Oct.2011)**

Parameters : haar wavelet 16, arithmmean

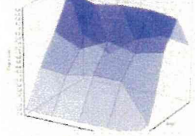
Sorted: compactness, peak lag (ascending order)

Results: 3771 probes (369/1600 clusters)

4 spaces separate clusters, **red** are responsive probes.

**Blues** are new and nice, brown are noise.

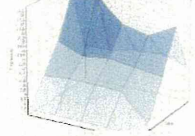
Ppm1l  
1438012\_at



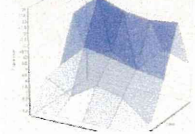
Plekha8  
1454819\_at



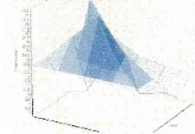
Homer2  
1424367\_a\_at



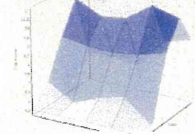
Arhgap21  
1428368\_at



Cpeb2  
1458518\_at



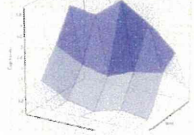
Esrra  
1460652\_at



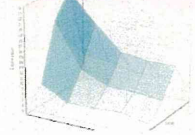
BC016495  
1447503\_at



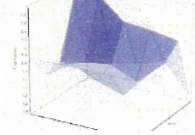
Por  
1416933\_at



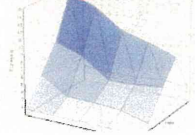
Lmo7  
1455056\_at



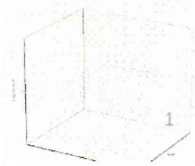
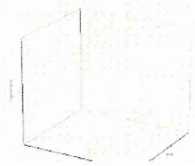
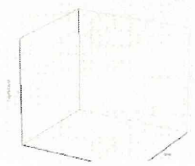
Fgf1  
1441042\_at



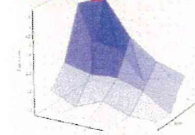
Per2  
1417602\_at



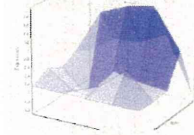
Dnajb6  
1448234\_at



Osbp13  
1428484\_at



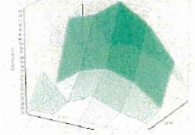
Zbtb20  
1443471\_at



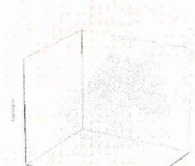
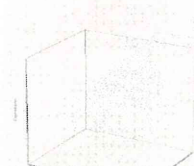
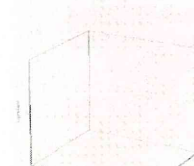
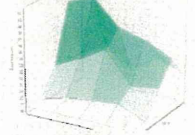
LOC631639:///Lonrf1  
1455665\_at



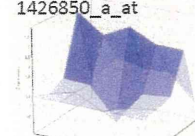
Prnp  
1448233\_at



Xpot  
1428949\_at



Map2k6  
1426850\_a\_at



Alas1  
1424126\_at



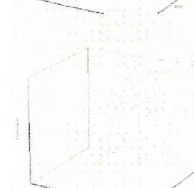
Nampt  
1455320\_at



Nampt  
1417190\_at



DOH4S114  
1436736\_x\_at



Xpot  
1428772\_at



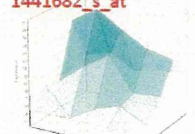
Tmtc4  
1428113\_at



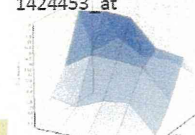
Slc25a46  
1418134\_at



Xpot  
1441682\_s\_at



Pcyt1a  
1424453\_at



Xpo6  
1422759\_a\_at



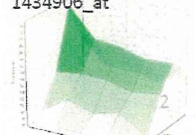
Lin7a  
1435805\_at



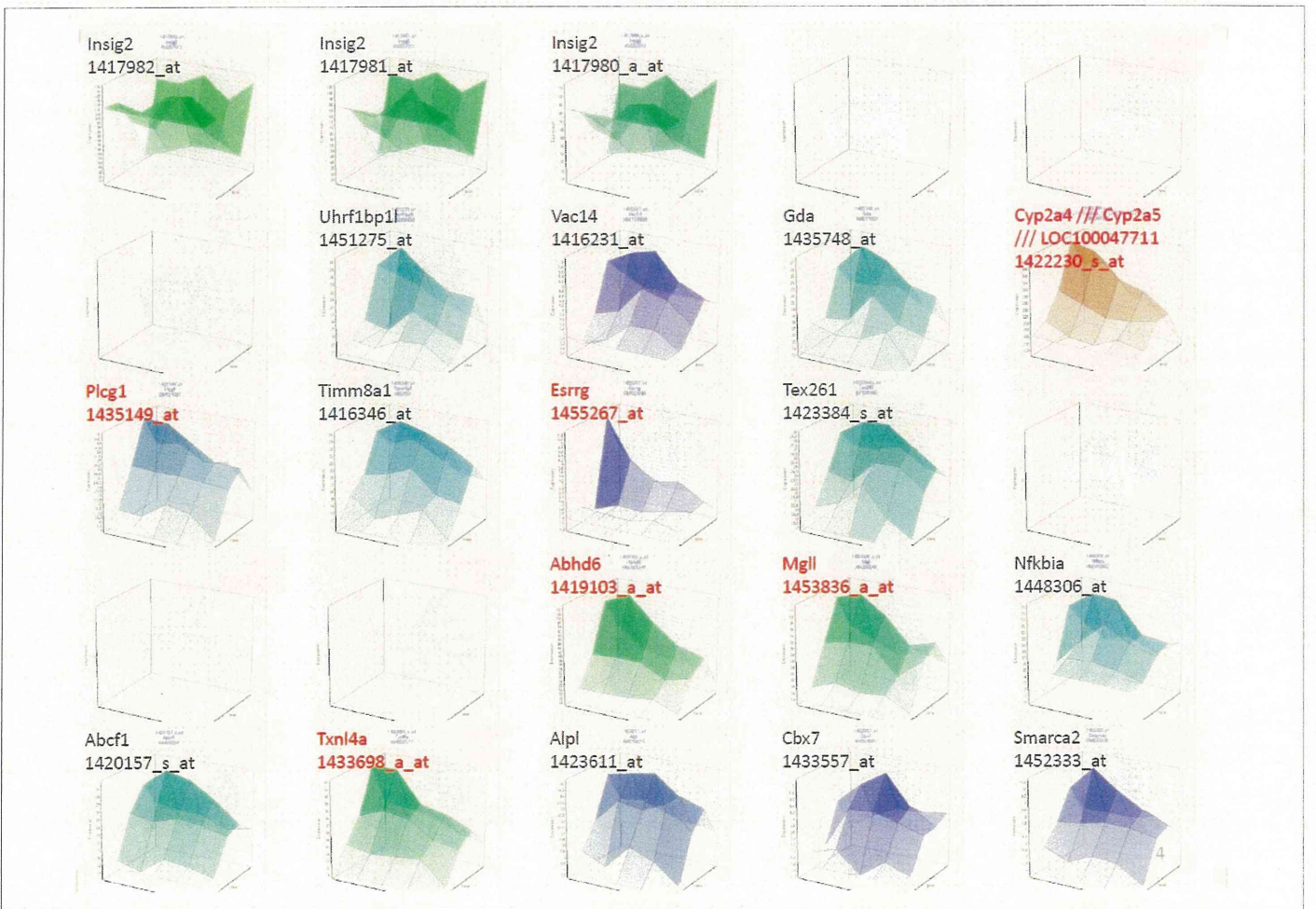
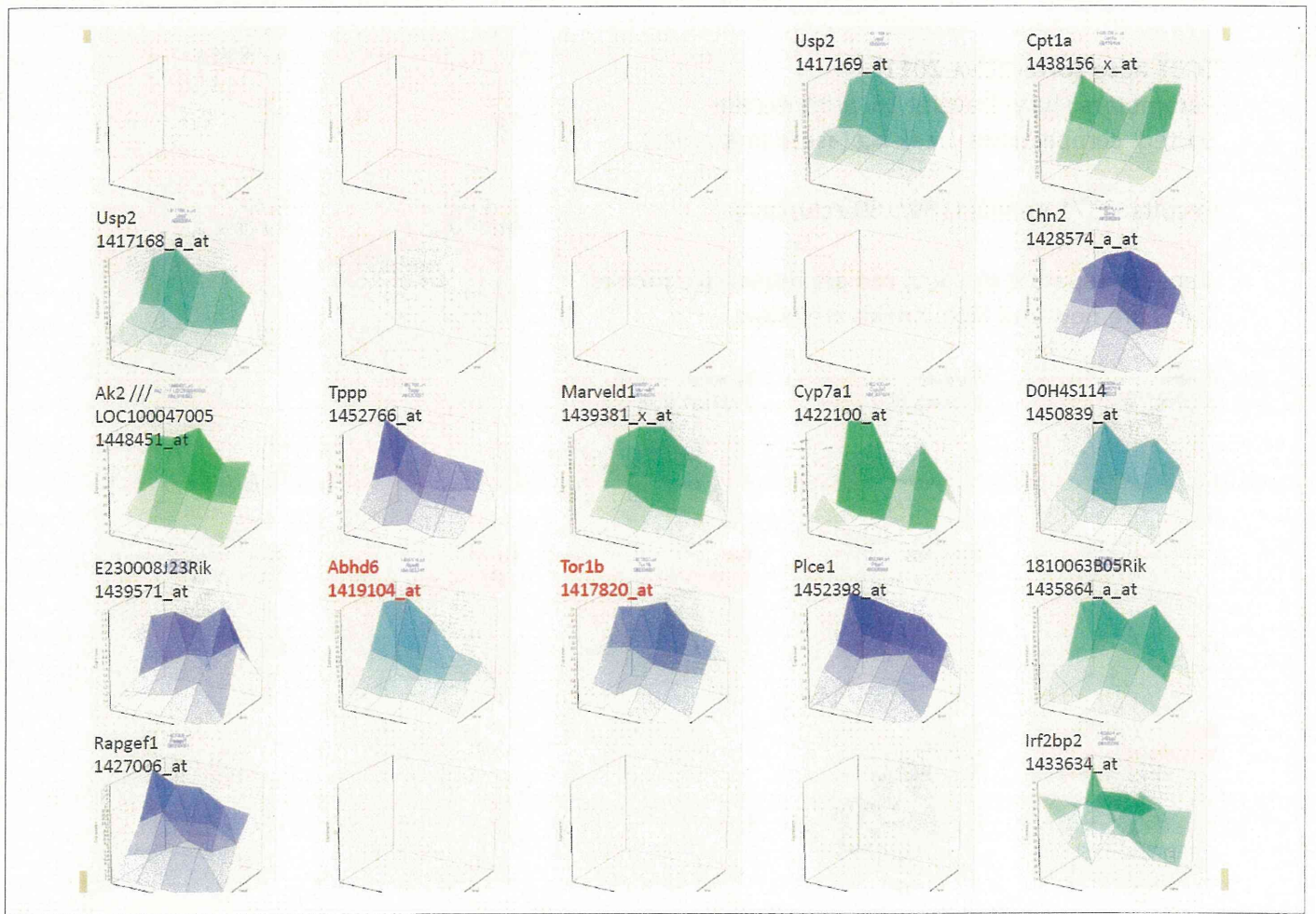
Ptplad1  
1452427\_s\_at



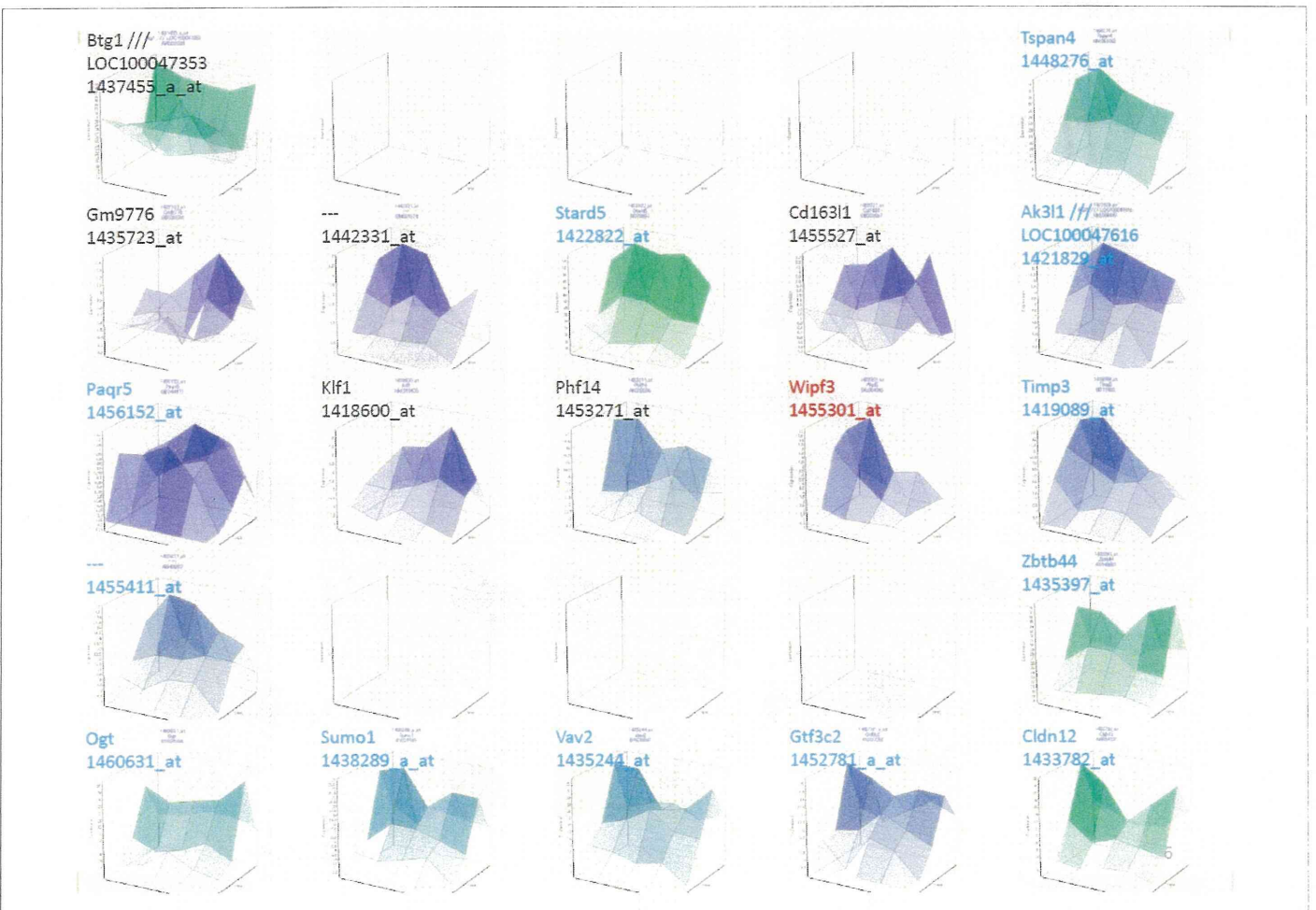
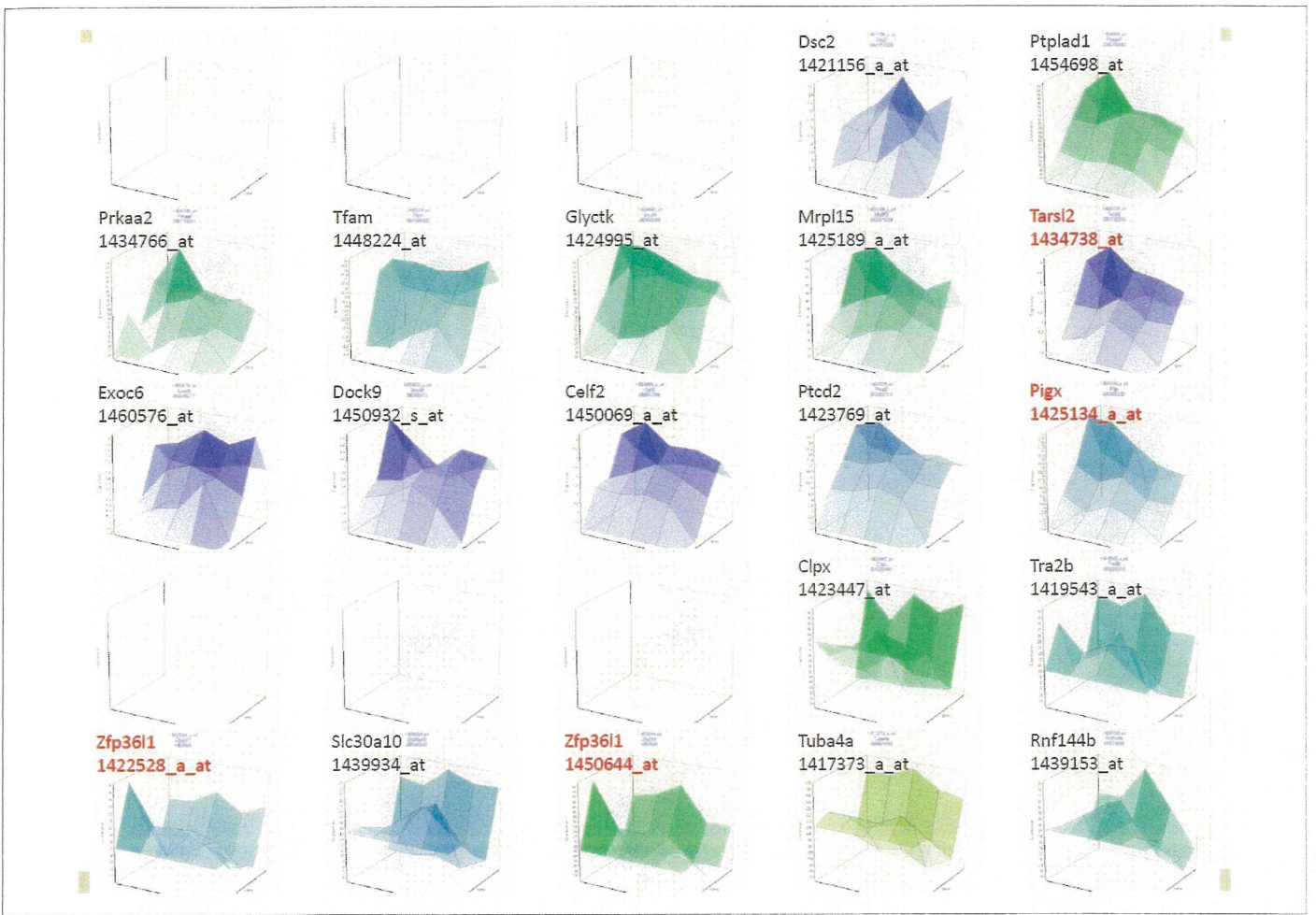
O610005C13Rik  
1434906\_at



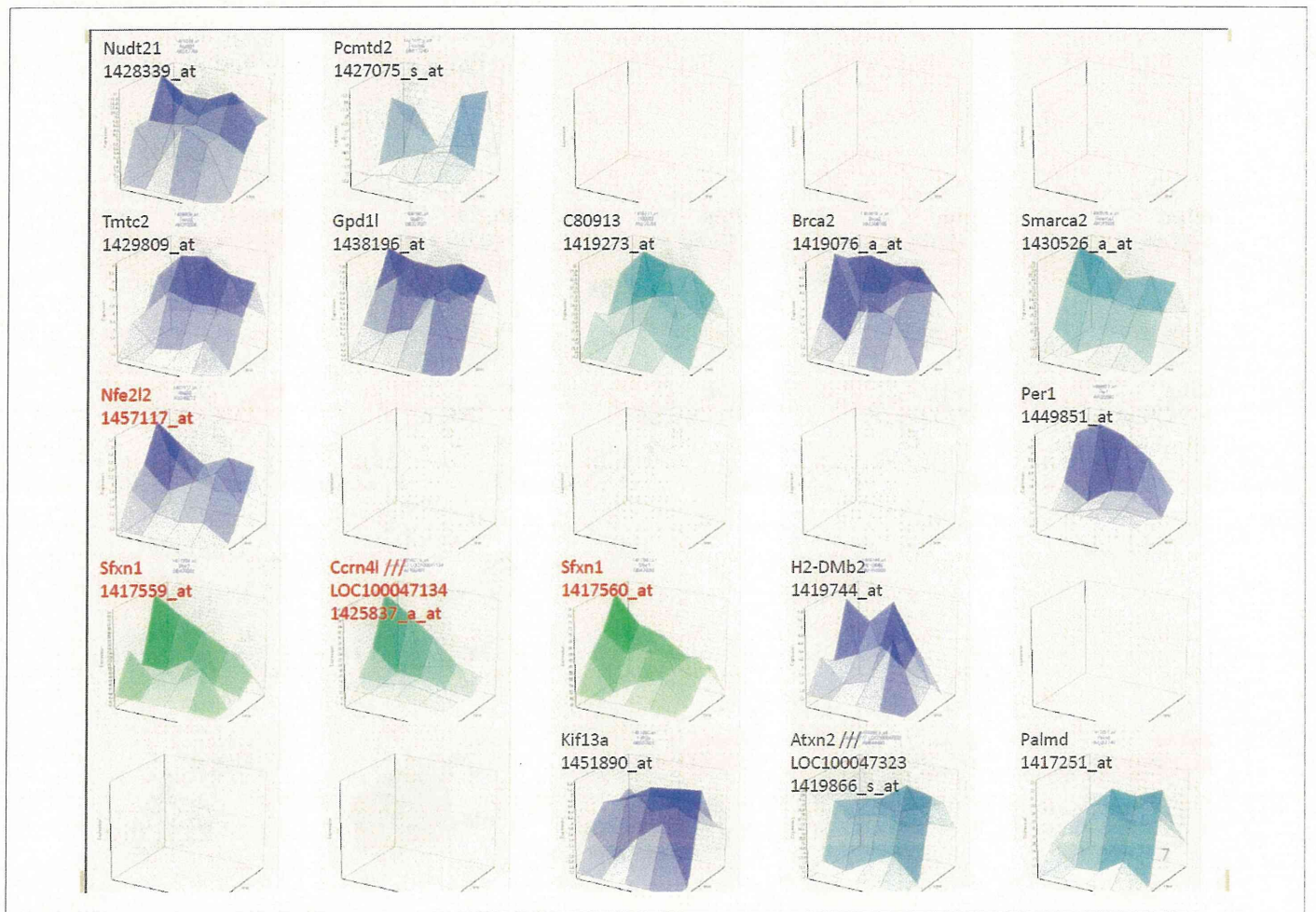




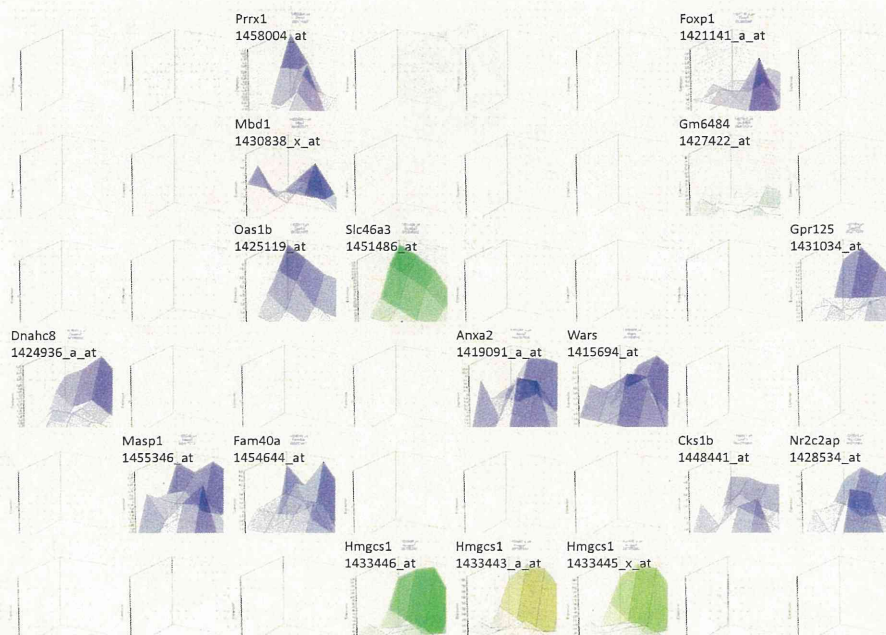


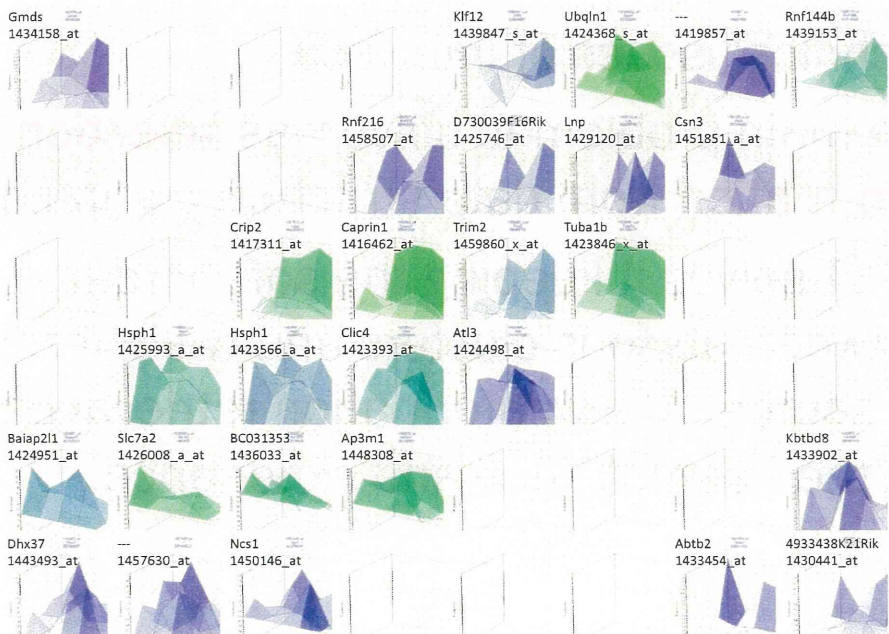
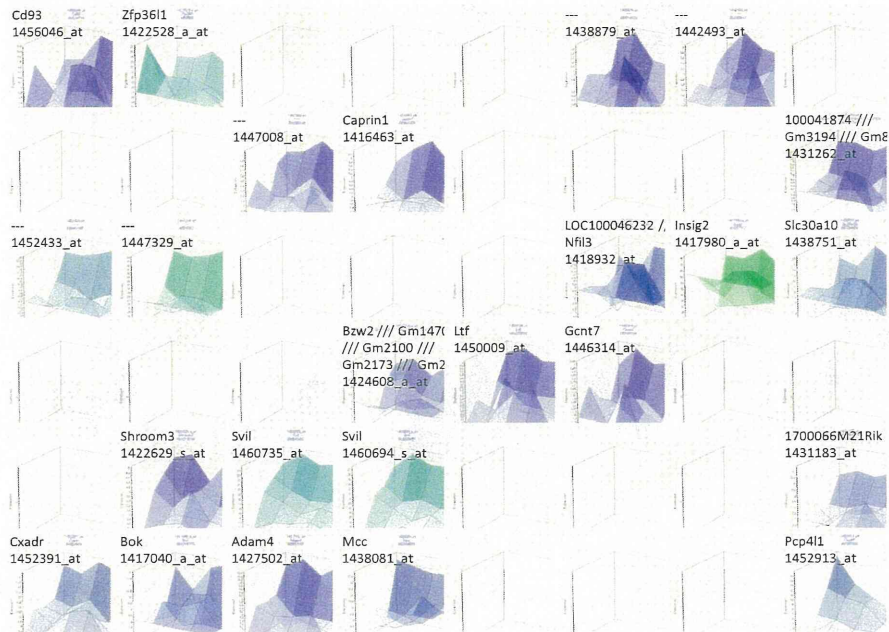






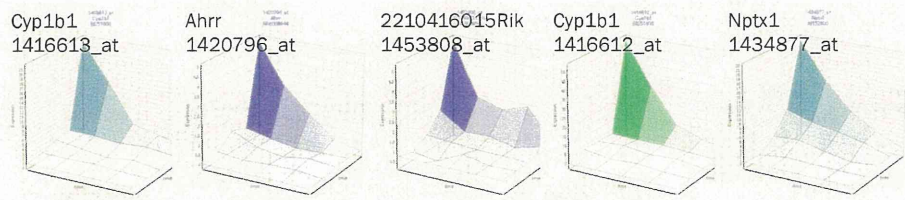
## TCDF 20,000 probes unsorted



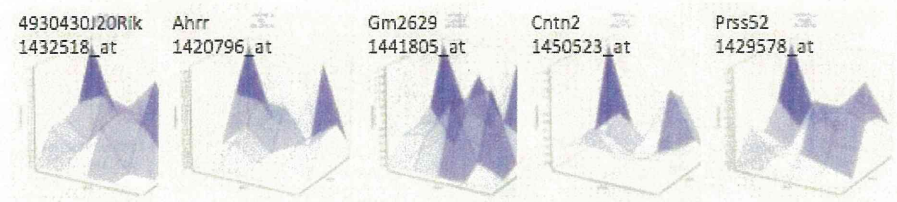




## TCDD



## TCDF



## What was the problem?

- \* technical issues of interface
- \* the main problem to solve was **selection of information** from data structure computed.
- \* AGCT solves it keeping the soft border between obvious clusters and twilight zone, as source of new information retrieval basing on the User intuition and experience.

# SHOE

## Sequence Homology in Higher Eukaryotes

47

トキシコゲノミクス第3班会議

# SHOE Phylogenetic footprinting tool

Basic Application Example

File Help

Upload use example

Pasted data is sent.

Homo sapie  
Mus muscul  
Rattus norv

Repeat mask

Upstream length (-TSS) 2000 200 Downstream (TSS+)

Reset Halt Submit

Align regulatory regions

Pairwise and multiple alignment

Scanning known motif

Matching TFBS matrix

Transfac Jaspar iPS transgene

Heuristic finding regulatory regions MEME

'Users/NataP/Desktop/SHC

Results

Promoter

Regulatory

TFBS match

Clear folder

CLUSTAL W (1.83) multiple sequence alignment

```
NM_027915 -----C
NM_080583 GCTAAC
NM_001282 ACTATC
.....
NM_027915 TTAACC
NM_080583 TGAACC
NM_001282 TGAGCC
.....
NM_027915 CAC---
NM_080583 CAC---
NM_001282 CATCAT
..
NM_027915 TACCCA
NM_080583 TATCTA
NM_001282 TGTCAC
.....
NM_027915 GCAACT
NM_080583 ----CT
NM_001282 -CAGAT
.....
NM_027915 CGACCT
NM_080583 CAACTT
NM_001282 TCACA-
.....
NM_027915 CGCCAA
NM_080583 CGCCAA
NM_001282 CGTCAC
.....
NM_027915 AGATCT
```

Sequence Homology in Higher Eukaryotes (SHOE)

トキシコゲノミクス第3班会議



# AGCT-Shoe Plug-in

ACCT cluster result file

Load ACCT file

upload /TCDD\_sortCompactness0.6PeaksLag.txt

Input AGCT output file as is.

Title 1	Title 2	Title 3	Title 4	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	
1	SPEC	Affinit	Probe																	
2	ID	Num	Comp	Davie	Silhou	num i	% in i	num	peak	peak	max									
3	319	1	0	1.46	1	0	0	5	0	0	18.2									
4	ID	GB_ACC	SPT	Speci	Annot	Seque	Seque	Targe	Repre	Gene	Gene	ENTR	RefSeq	Transcript	ID	Gene	Gene	Gene	labeled	maxA
5	1424	AK00	Not_d	Mus	11-M	Conse	GenB	qb AK	AK00	cyclin	Cdkn1a	12575	NM_001111099	/	0000	0000	0004			18.2
6																				
7	ID	Num	Comp	Davie	Silhou	num i	% in i	num	peak	peak	max									
8	487	1	0	1	1	0	0	5	0	0	14.9									
9	ID	GB_ACC	SPT	Speci	Annot	Seque	Seque	Targe	Repre	Gene	Gene	ENTR	RefSeq	Transcript	ID	Gene	Gene	Gene	labeled	maxA
10	1452	BC02	Not_d	Mus	11-M	Conse	GenB	qb BC	BC02	Mus	Ptns	69202	NM_026988		0006	0000	0003			14.9
11																				
12	ID	Num	Comp	Davie	Silhou	num i	% in i	num	peak	peak	max									
13	892	1	0	1	1	0	0	5	0	0	24.7									
14	ID	GB_ACC	SPT	Speci	Annot	Seque	Seque	Targe	Repre	Gene	Gene	ENTR	RefSeq	Transcript	ID	Gene	Gene	Gene	labeled	maxA
15	1416	AI462	Not_d	Mus	11-M	Conse	GenB	qb AI	AI462	forkh	Foxo1	56458	NM_019739		0001	0005	0003			24.7
16																				
17	ID	Num	Comp	Davie	Silhou	num i	% in i	num	peak	peak	max									
18	978	1	0	1	1	0	0	10	0	0	31.0									
19	ID	GB_ACC	SPT	Speci	Annot	Seque	Seque	Targe	Repre	Gene	Gene	ENTR	RefSeq	Transcript	ID	Gene	Gene	Gene	labeled	maxA
20	1429	AK00	Not_d	Mus	11-M	Conse	GenB	qb AK	AK00	RIKEN	6230	1001	NM_001033123	/	0006	0005	0003			31.0
21																				
22	ID	Num	Comp	Davie	Silhou	num i	% in i	num	peak	peak	max									
23	1021	1	0	1	1	0	0	5	0	0	21.9									
24	ID	GB_ACC	SPT	Speci	Annot	Seque	Seque	Targe	Repre	Gene	Gene	ENTR	RefSeq	Transcript	ID	Gene	Gene	Gene	labeled	maxA
25	1416	BC01	Not_d	Mus	11-M	Conse	GenB	qb N	BC01	myel1	Mozl2	14012	NM_007962		0007	0016	0005			21.9
26																				
27	ID	Num	Comp	Davie	Silhou	num i	% in i	num	peak	peak	max									
28	1023	1	0	1	1	0	0	5	0	0	78.0									
29	ID	GB_ACC	SPT	Speci	Annot	Seque	Seque	Targe	Repre	Gene	Gene	ENTR	RefSeq	Transcript	ID	Gene	Gene	Gene	labeled	maxA
30	1416	NM_0	Not_d	Mus	11-M	Conse	GenB	qb N	NM_0	simila	LOC1	1000	NM_053075	/// X	0007	0005	0000			78.0
31																				
32	ID	Num	Comp	Davie	Silhou	num i	% in i	num	peak	peak	max									
33	1452	1	0	1	1	0	0	5	0	0	13.7									

Close Enrich column

49

トキシコゲノミクス第3班会議

## SHOE RESULTS on HUMAN-MOUSE-RAT promoter analysis

pronetbeansdesktopapp.ProNetBeansDesktopApp

Relseq ID of genes in clusters

10 Top motifs

71 ID 574

72 NM\_144862

73 NM\_001045529

74 NM\_133784

75 NM\_012006

76 NM\_010907

77

78 ID 667

79 NM\_001002012

80 NM\_008301

81 NM\_001081056

82 NM\_133648

83 NM\_133649

84 NM\_011844

85 NM\_009460

86 NM\_001128169

87 NM\_001128170

88 NM\_001128171

89 NM\_173369

90 NM\_011844

91

92

93 ID 587

94 NM\_012006

95 NM\_134188

96 NM\_146034

97 NM\_029447

98 NM\_001037940

99 NM\_001037941

100 NM\_001127367

101 NM\_011847

102 NM\_172665

103 NM\_153122

104

105 ID 411

106 NM\_133896

107 XM\_895068

108 NM\_011066

109 NM\_013495

110

111 ID 10

112 NM\_021449

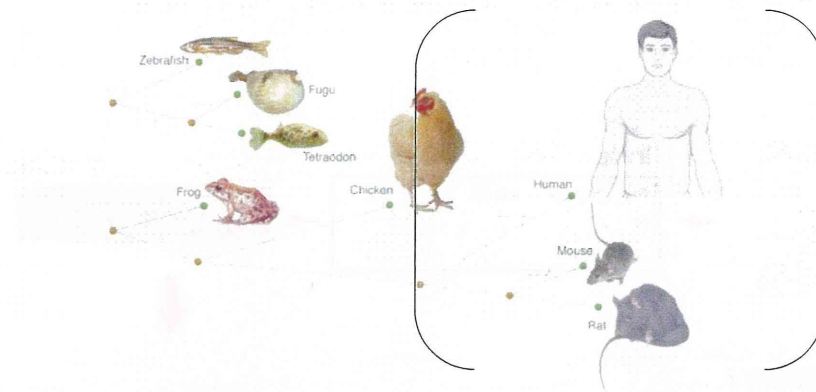
113

114

115

50

# Phylogenetic footprinting finds evidence of functionality



- TFs
  - Brain-specific
- (Suzuki et. al, 2004)

## Regulation conserved!

Transcription Factor (TF)

```

Human  CTCCCACCTCCCAATTCGCGCAGGGCCGCGAG--ACTATAGGCGCTGC
Mouse  CTCTCGACCCCTCCAAATTC--CACACAGGGCTCTCTGACTATATGAGCGCTG
Rat     CTCTCGACCCCTCCAAATTC--CACACAGGGCTCTCTGACTATATGAGCGCTCTT
    * * * * *
Human  GGCTGAGCCCTCTGCGCTGGGACAGCTAGAGAGAGGGCCGCGCTGGGAGATCGCTCTC
Mouse  AGCTGA-CAGGCGGACCGCAGATGCGGAGATGGCCGCGCTGGGAGATCGCTCTCC
Rat     AGCTGA-CAGGCGGACCGCAGATGCGGAGATGGCCGCGCTGGGAGATCGCTCTCC
    * * * * *
Human  AAGCCCTGCTGTCGCGCTCCCGCTGAGCGAGGTTGTTGGG
Mouse  AAGCCCTGCTGTCGCGCTCCCGCTGAGCGAGGTTGTTGGG
Rat     AAGCCCTGCTGTCGCGCTCCCGCTGAGCGAGGTTGTTGGG
    * * * * *
    
```

ear

## Regulation not conserved

```

Human  ----TTCCTTCCCA-TGTTATTTATTTGAAAATGATGGCTGGAGATCATGAGTAGAG--
Mouse  ----GTGTTTCCCA-TGCTTTGTTTTCTC--TGACTCCCTGGGG-CATGGGACACATTT
Rat     AGCTGTATATTAATAAATACTTTGTTTT-----TACAGGAAAG-GAAGATAAAGACT-
    * * * * *
Human  --AAAATGAAATTTGGGCTGAGGGAT-TAAA--CCCTGCTGCTGCTGCTATATGACAT
Mouse  ACAACGGAAATACTGGGTGGAAGAA-TGAAAGCCCATGCTG-----
Rat     ACGCTTTAAAACA---AATAGAGGAACTGGAAATTTCTATGCTG-----
    * * * * *
    
```

brain

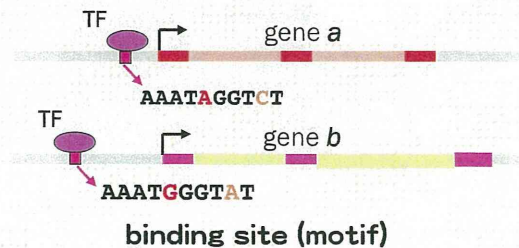
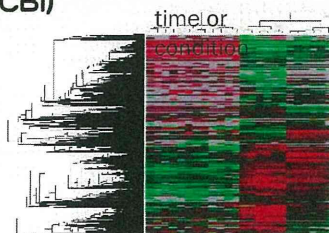
51

トキシコゲノミクス第3 班会議

# Co-regulated genes

Similar way of expression suggests **regulation** by the same Transcription Factor (TF)

GE (NCBI) Gene Expression Omnibus



## Motif discovery tools

&

## Comparative genomics tools

- MEME (Bailey et.al.)
- Consensus (Stormo et.al.)
- Gibbs sampler (Lawrence et.al.)
- Yebis (Yada et.al.)....

- MONKEY (Moses et.al)
- FootPrinter (Tompa et.al)
- PhyMe (Sinha et.al)
- PhyloGibbs (Siddharthan et.al) ....

52

トキシコゲノミクス第3 班会議



鼠 S

### STEP 1

## DATASET of ORTHOLOGOUS PROMOTERS



refseq ID

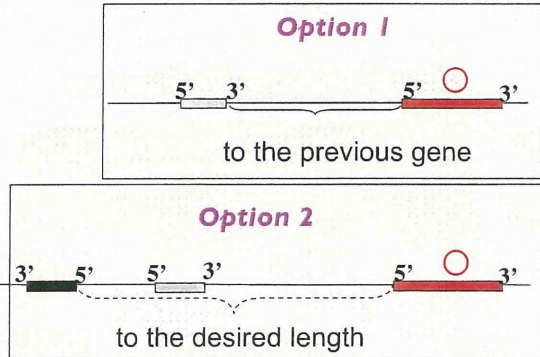
Human gene A  
Human gene B  
Mouse gene C  
Rat gene D

Index file: 7,000genes

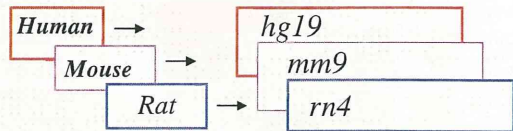
Gene	Human	Mouse	Rat
Acly	NM_198830	NM_134037	NM_016987
Alcoa	NM_184043	NM_007438	NM_012495
Fdft1	NM_004462	NM_010191	NM_019238
Gpd2	NM_000408	NM_010274	NM_012336

User's orthologs

Pur2 NM\_004462 NM\_010191 NM\_019238  
Gpd2 NM\_000408 NM\_010274 NM\_012336



Genes locations INFO + Genomic Data

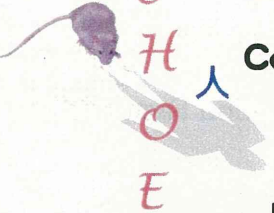


Updated!  
トキシゲノミクス第3班会議

鼠 S

### STEP 2

## Construction of multiple alignment



human-mouse (Ssearch)

```

230 240 250 260 270 280
AACTAAAGGCCCGCAGGGGAGAGTAAATTAAAGCGCTAATTAGGAGTAAATGGAGGGGAG
: : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
AATTAATGCTCACTGGGAA--GAATTAAAGGCTCAACTGGAAGATAATGAAGGAGGC
650 660 670 680 690 700

290 300 310 320 330 340
ACGCAGAAAGCCCTTACTCTTGGCCCTCAGGGAAAAGGAGTTTCTCTCACTGCCTG
: : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
CAAAAGA-GCCTGCTTATTCTTGGCCCTTAGGGGAGAGAGTTTGTCTCTTTGCCAG
710 720 730 740 750 760

```

human-rat

```

270 280 290 300 310 320
GTAATGAAGGGGAGCAGCAGAAAGCCCTTACTCTTGGCCCTCAGGGAAAAGGAGTT
: : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
AGAGGTAAGTGGCAGATAATACGCTGCCTTAAATCTTGGCCCTTAGGGGAGAGGAGTT
11510 11520 11530 11540 11550 11560

330 340 350 360 370 380
TCCTCTCACTGCCTGAGAATAG-GAAGGTGGCTGGCAGAAAAGTCCAAAAGGGAAGAGA
: : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
TGCTCTTTGCCAGTGGGAGTGCAGTGGCTGGCAGAAAAGTCTGAGCGCCATGG
11570 11580 11590 11600 11610 11620

```

human-mouse-rat (ClustalW)

```

Human  AACTAAAG-GCCCCGAGGGGAGAGTAAATTAAAGCGCTAATTAGGAGTAAATGGAGGGGGA
Mouse  AATTAATGCTCACTGGGAAAGATTAA---AGGGTCAACTGGAAGATAATGAAGGAGGC
Rat    -----AGAGGTAAGT----GGCA
          * * * * *

Human  GAGCAGAAAGCCCTTACTCTTGGCCCTCAGGGAAAAGGAGTTTCTCTCACTGCCT
Mouse  G-CAAAAGAGCCTGCTTATTCTTGGCCCTTAGGGGAGAGAGTTTGTCTCTTTGCCA
Rat    GATAATACGCCCTGCTTAAATCTTGGCCCTTAGGGGAGAGGAGTTTGTCTCTTTGCCA
          * * * * * * * * * * * * * * * * * * * * * * * * * * * * * * * * * *

```

### STEP3

How good is the alignment?

Pattern frequencies tables for “good” and “random” alignments

74 patterns

A, T, G, C, - X AA, AT, AG, AC, A-, TT, TG, TC, T-, GG, GC, G-, CC, C-, --

835 orthologous alignments  
(238,800bp)

1260 random alignments  
(239,600bp)

$$MAscore = \log_{10} \frac{\Pr g1 * \Pr g2 * \Pr g3 \dots \Pr gn}{\Pr r1 * \Pr r2 * \Pr r3 * \Pr rn} = \log_{10} \frac{\prod \Pr(c|good\_alignment)}{m}{\prod \Pr(c|random\_alignment)}{m}$$

c - probability of pattern in each column in *good alignment* and *random alignment* tables,  
m - motif length. 55 トキシコゲノミクス第3班会議

BIOBASE  
BIOLOGICAL DATABASES



### STEP4

Calculation of PSSM score (PMscore)

498 human-rodent matrices

	①	②	③	④	⑤	⑥	⑦	⑧	⑨	⑩	⑪
A	4	5	3	0	4	3	3	2	1	1	1
C	1	2	0	0	0	0	0	1	3	4	6
G	2	2	7	2	3	7	0	4	3	1	1
T	3	1	0	8	3	0	7	3	3	4	2

$$PMscore = \sum_{i=1}^m \log_{10} \frac{\text{count}_{x_i} + \text{pseudocount}_{x_i}}{\sum_{x=A,T,G,C} \text{count}_{x_i} + \sum_{x=A,T,G,C} \text{pseudocount}_{x_i}}$$

where pseudocount = 1  
m is a motif length

Probability of motif at each position in human sequence

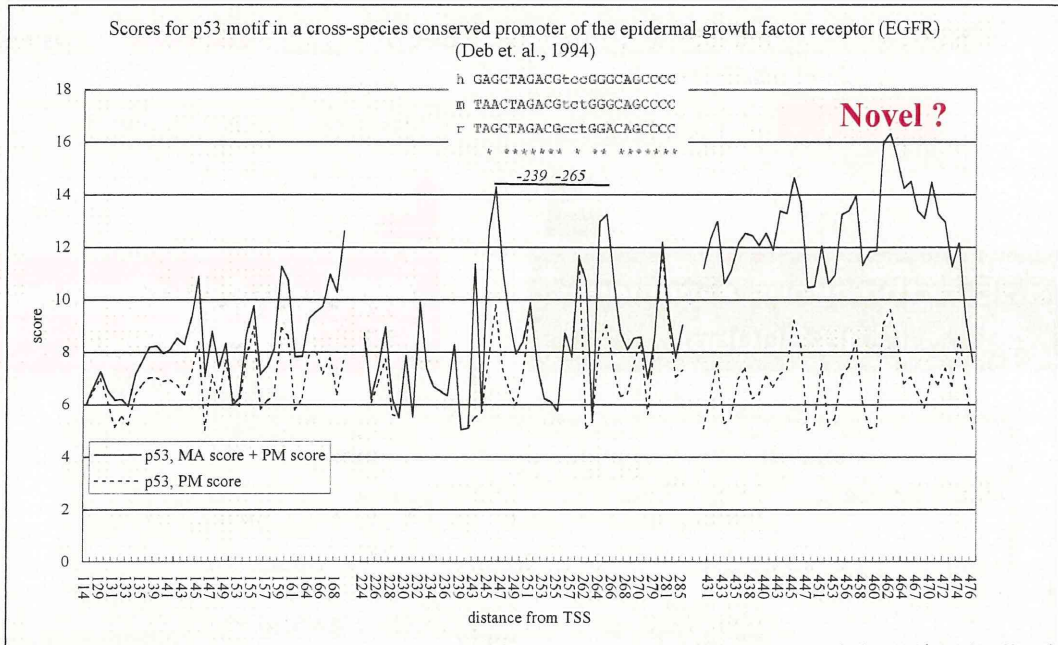
human	-0.48	-0.39	-0.27	-0.22	-0.57	-0.27	-0.27	-0.57	-0.57	-0.47	-0.33	-4.46
A	A	A	G	T	G	G	T	T	C	C	C	
mouse	A	T	G	T	G	G	T	C	A	C	C	
rat	A	T	C	T	G	G	T	A	A	C	C	



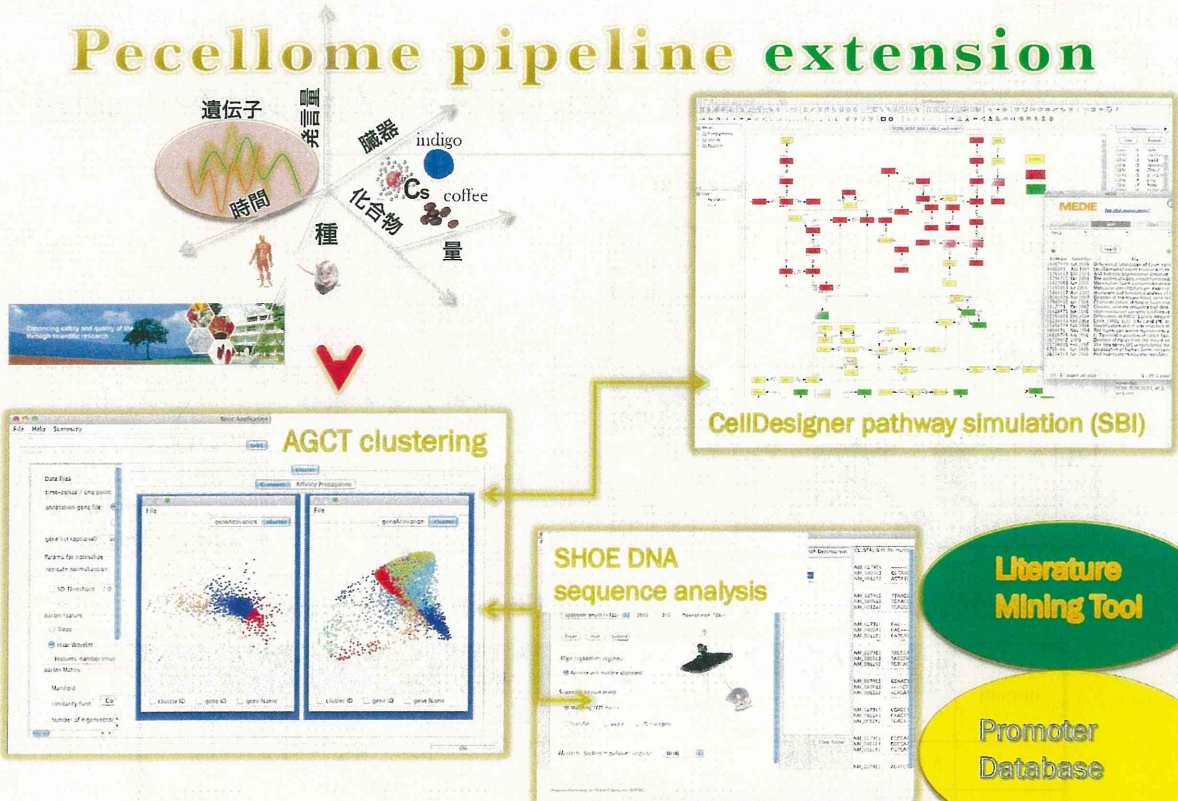


# Finding motif candidates

Multiple alignment score (MA) + Transfac PSSM score (PM)



# Pecellome pipeline extension



# Acknowledgements

- \* Hiroaki Kitano (Sony CSL)
- \* Keigo Oka (Tokyo University)
- \* Frank Nielsen (Sony CSL)
- \* Richard Nock (University of Martinique)



## TCDD / TCDF

2012 Feb  
ym

## TCDD

- 2,3,7,8-Tetrachlorodibenzo-p-dioxin (TCDD)

