トランスクリプトーム解析におけるバイオインフォマティクス要素技術~私の相場観~

東京大学大学院農学生命科学研究科 アグリバイオインフォマティクス教育研究ユニット 門田 幸二(かどた こうじ)

http://www.iu.a.u-tokyo.ac.jp/~kadota/kadota@iu.a.u-tokyo.ac.jp



門田の部分は

■ バイオインフォマティクス人材育成講座

- □ 7月:マイクロアレイ受託実験や遺伝子発現データ解析
- □ 8月:次世代シーケンサー見学
- □ 8月:統計解析言語Rや遺伝子発現データのクラスタリング

このあたりの補講



■ 第一部のねらい

□ クラスタリング(やエントロピー)などのバイオインフォマティクスの**基本的なスキルを身につけるだけで様々な局面に応用可能である**という二つの事例を紹介するとともに、これらの具体的な計算例を示すことで数式アレルギー緩和に貢献



■ 第二部のねらい

- □ 次世代シーケンサーのデータだって、統計解析言語Rでお手軽に解析できる
- □ 遺伝子発現行列にしたら後は同じ
- □ 高度なプログラミング能力がなくてもバイオインフォマティクスの世界で生存可能
- □ 必要な情報はインターネットのみで十分(@沖縄)





自己紹介



- 1995年3月
 - □ 高知工業高等専門学校・工業化学科 卒業
- 1997年3月
 - □ 東京農工大学・工学部・物質生物工学科 卒業
- 1999年3月
 - □ 東京農工大学・大学院工学研究科・物質生物工学専攻 修士課程修了
- 2002年3月
 - □ 東京大学・大学院農学生命科学研究科・応用生命工学専攻 博士課程修了
 - □ 学位論文:「cDNAマイクロアレイを用いた遺伝子発現解析手法の開発」(指導教官:清水謙多郎教授)
- **2002/4/1~**
 - □ 産総研・生命情報科学研究センター 産総研特別研究員
- **2003/11/1~**
 - □ 放医研・先端遺伝子発現研究センター 研究員
- **2005/2/16~**
 - □ 東京大学・大学院農学生命科学研究科 特任助手
- 2007/4/1~現在

□ 東京大学·大学院農学生命科学研究科 特任助数

高専時代の成績もたいしたことない門田が、かれこれ10年以上 バイオインフォマティクスの分野で楽しくやってます。

アグリバイオインフォマティクス プログラム

トランスクリプトームとは

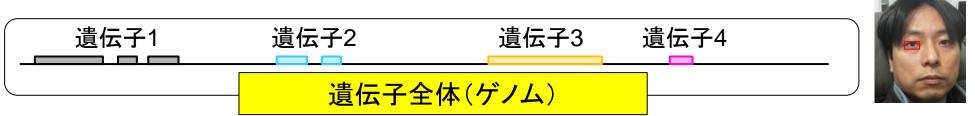
- ある特定の状態の組織や細胞中に存在する全RNA(転写物、transcripts)の総体
- 様々なトランスクリプト―ム解析技術
 - □マイクロアレイ
 - cDNAマイクロアレイ、Affymetrix GeneChip、タイリングアレイなど
 - □配列決定に基づく方法
 - ■EST、SAGEなど、次世代シーケンサー
 - □電気泳動に基づく方法
 - Differential Display、AFLPなど

調べたい組織でどの遺伝子がどの程度発現しているのかを一度に観察

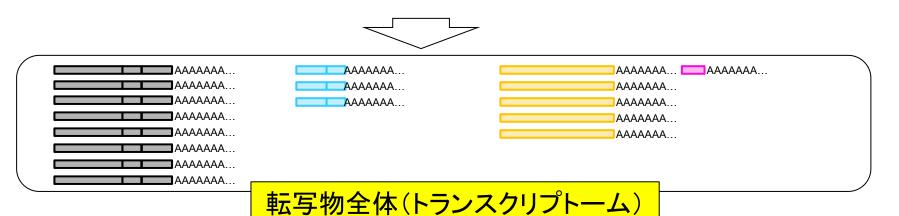
Sep 17 2010

トランスクリプトームとは

■ ある状態のあるサンプル(例:目)のあるゲノムの領域



・どの染色体上のどの領域にどの遺伝子があるかは調べる個体(例:ヒト)が同じなら不変(目だろうが心臓だろうが...)



- ・遺伝子1は沢山転写されている(発現している)
- ・遺伝子4はごくわずかしか転写されてない

• . . .

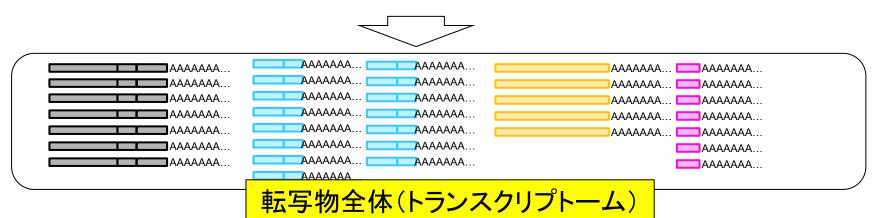
Sep 17 2010 ...

トランスクリプトームとは

■ ある状態のあるサンプル(例:目)のあるゲノムの領域

遺伝子1 遺伝子2 遺伝子3 遺伝子4 遺伝子全体(ゲノム)

・どの染色体上のどの領域にどの遺伝子があるかは調べる個体(例:ヒト)が同じなら不変(目だろうが心臓だろうが...)



- ・遺伝子2は光刺激に応答して発現亢進
- ・遺伝子4も光刺激に応答して発現亢進

光刺激

トランスクリプトーム情報を得る手段

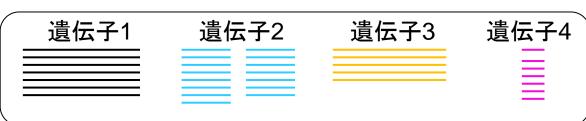
■ 光刺激前(T1)の目のトランスクリプトーム



これがいわゆる 「遺伝子発現行列」

T1T2遺伝子187遺伝子2315遺伝子355遺伝子417

■ 光刺激後(T2)の目のトランスクリプトーム



- マイクロアレイ
- 電気泳動に基づく方法
- 配列決定に基づく方法

トランスクリプトーム取得(マイクロアレイ)

よく研究されている生き物は多数の遺伝子 (の配列情報)がわかっている

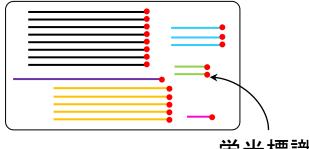
遺伝子1 遺伝子2 遺伝子3 遺伝子4

Mage
Courtes Val
Afrymetrix

わかっている遺伝子(の配列の相補鎖)を搭載した"チップ"

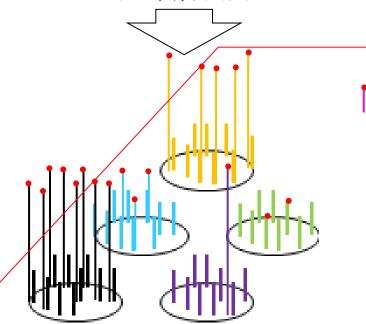
- ・メーカーによって搭載されている遺伝子の 種類が異なる
- →搭載されていない遺伝子(未知遺伝子含む、例:遺伝子4)の発現情報は測定不可... Sep 17 2010

光刺激前(T1)の目の トランスクリプトーム



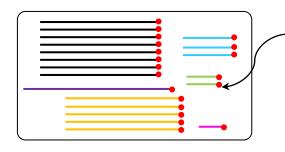
蛍光標識

ハイブリダイゼーション (二本鎖形成)



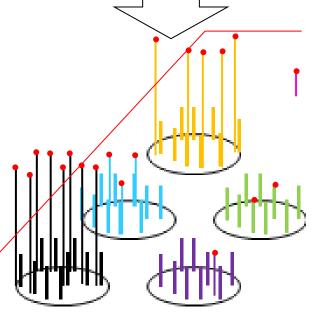
トランスクリプトーム取得(マイクロアレイ)

■ 光刺激前(T1)の目のトランスクリプトーム

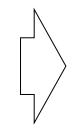


蛍光標識

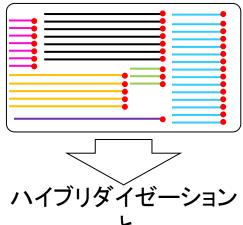
ハイブリダイゼーション (二本鎖形成)



専用の検出器で各 遺伝子に対応する 領域の蛍光シグナ ル強度を測定



光刺激後(T2)の目の トランスクリプトーム



ラッダイビーン: と シグナル検出

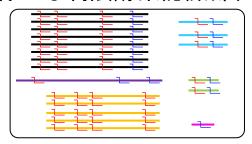
	T1	
遺伝子1	8	
遺伝子2	3	
遺伝子3	5	
遺伝子4	?	
遺伝子5		

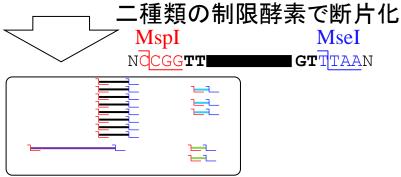
T2	
7	
15	
5	
?	

トランスクリプトーム取得(電気泳動)

■ cDNA-AFLP(HiCEPの場合)

様々な制限酵素認識部位をもつトランスクリプトーム



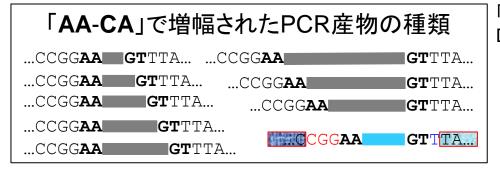




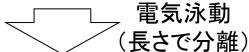


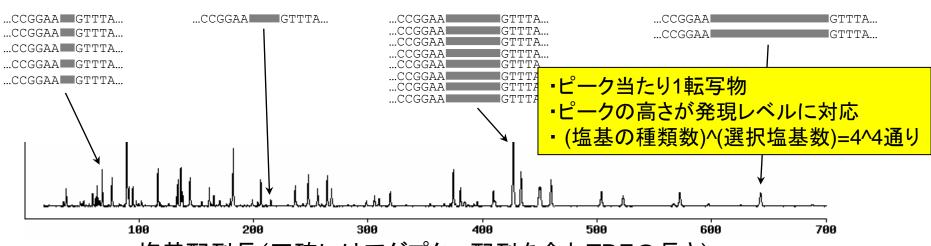
トランスクリプトーム取得(電気泳動)

■ cDNA-AFLP(HiCEPの場合)



「転写物由来配列断片(Transcripts Derived Fragments; TDFs)」という表現が業界標準(たぶん)

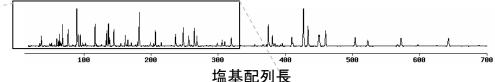




塩基配列長(正確にはアダプター配列を含むTDFの長さ)

トランスクリプトーム取得(電気泳動)

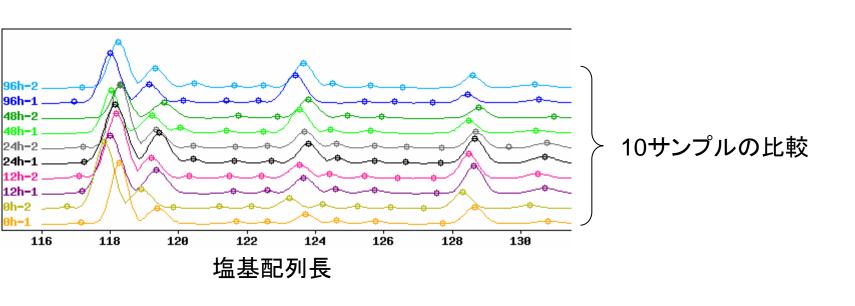
■ cDNA-AFLP(HiCEPの場合)



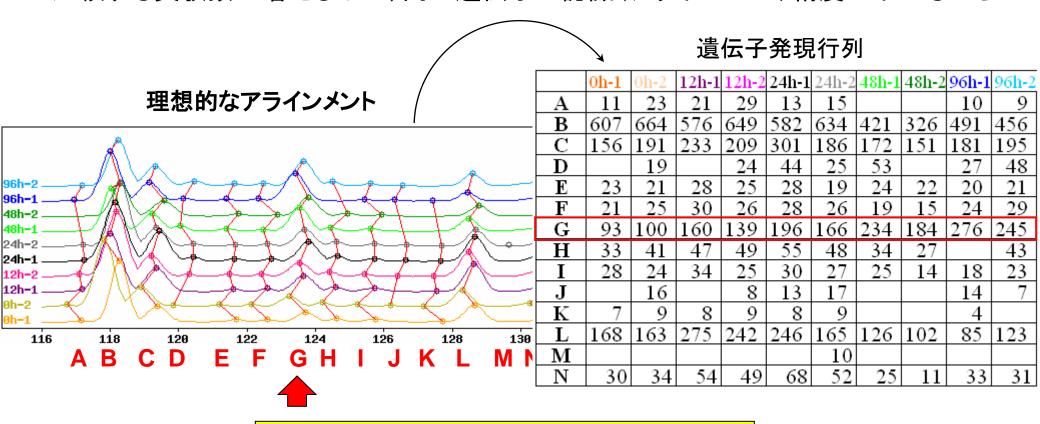
AA-CA (T3)(T2)光刺激前 (T1)遺伝子2は光刺激に応答して発現亢進!

のような解析を目でやるのは労力がかかり大変です → **バイオインフォマティクス**

- マイクロアレイ(や塩基配列データ)では遺伝子発現行列が出発点
- 電気泳動データは遺伝子発現行列の作成が困難 比較する実験数が増えるほど、同一遺伝子の認識(アラインメント)精度が下がるから



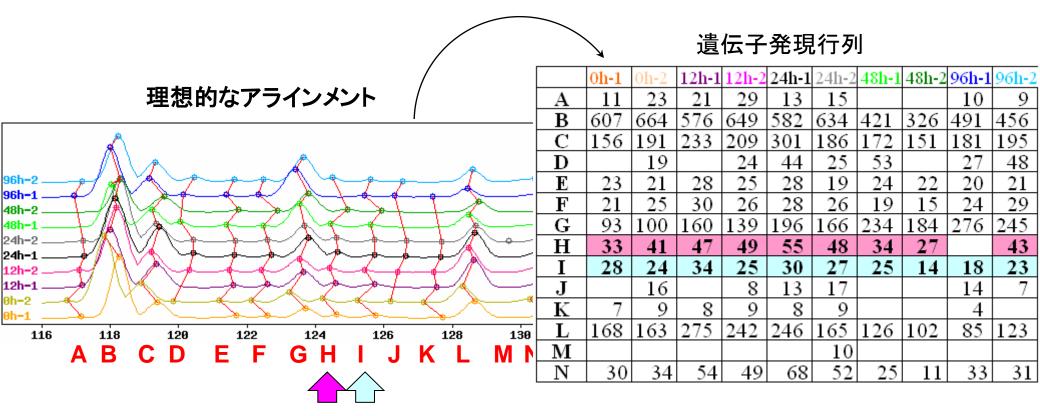
- マイクロアレイ(や塩基配列データ)では遺伝子発現行列が出発点
- 電気泳動データは**遺伝子発現行列の**作成が困難 比較する実験数が増えるほど、同一遺伝子の認識(アラインメント)精度が下がるから



Gの発現が徐々に上昇している

- マイクロアレイ(や塩基配列データ)では遺伝子発現行列が出発点
- 電気泳動データは**遺伝子発現行列**の作成が困難

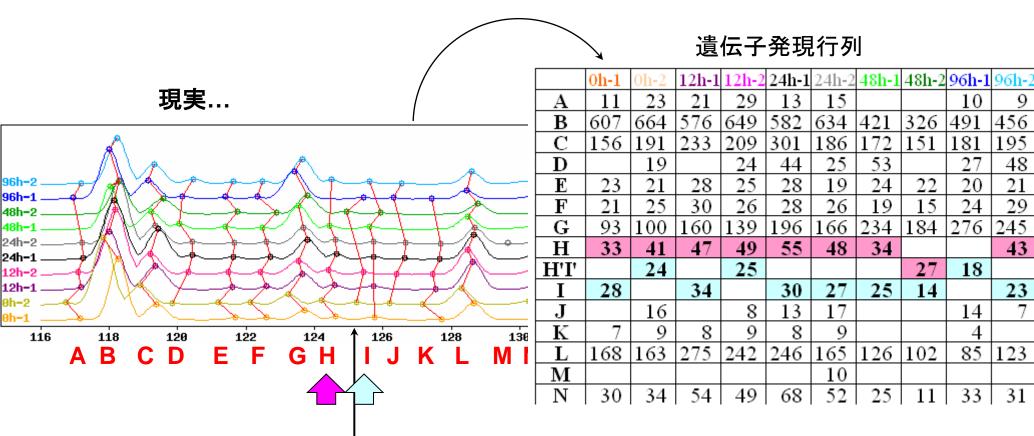
比較する実験数が増えるほど、同一遺伝子の認識(アラインメント)精度が下がるから



- マイクロアレイ(や塩基配列データ)では遺伝子発現行列が出発点
- 電気泳動データは**遺伝子発現行列**の作成が困難

Sep 17 2010

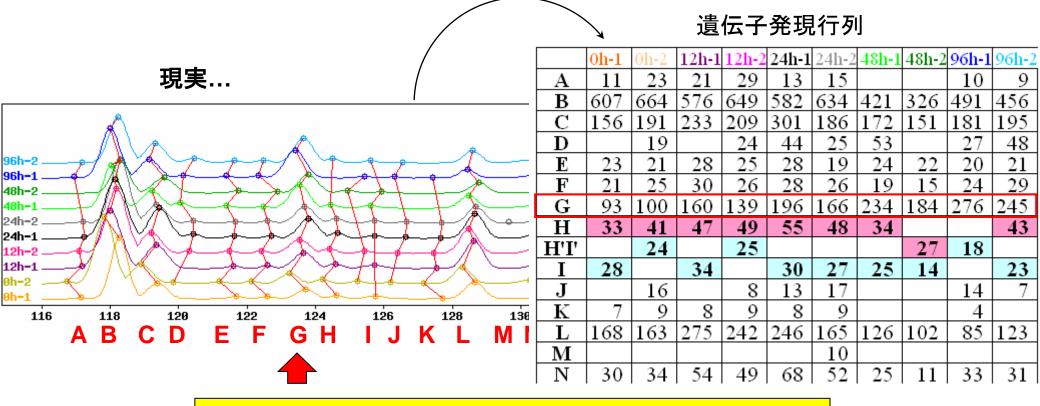
比較する実験数が増えるほど、同一遺伝子の認識(アラインメント)精度が下がるから



16

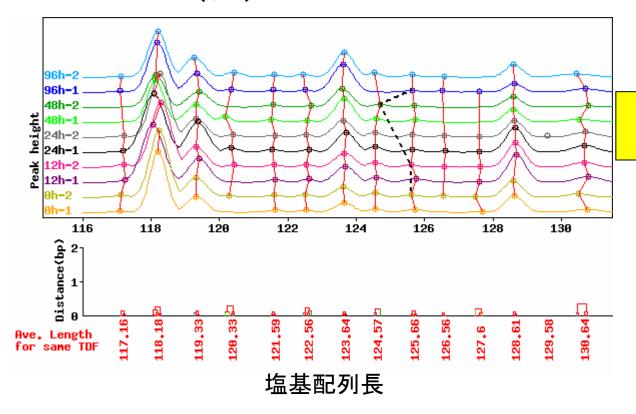
- マイクロアレイ(や塩基配列データ)では遺伝子発現行列が出発点
- 電気泳動データは**遺伝子発現行列**の作成が困難

比較する実験数が増えるほど、同一遺伝子の認識(アラインメント)精度が下がるから

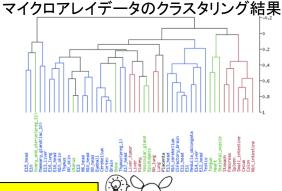


Gの発現パターンは本当に全部G由来?!

GOGOT法



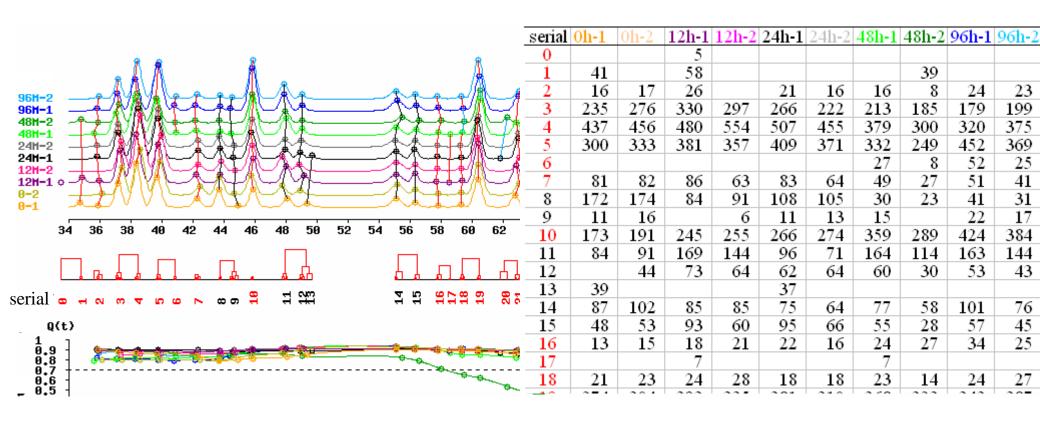
「フラグメント長補正」 + 「Complete-linkage clustering」

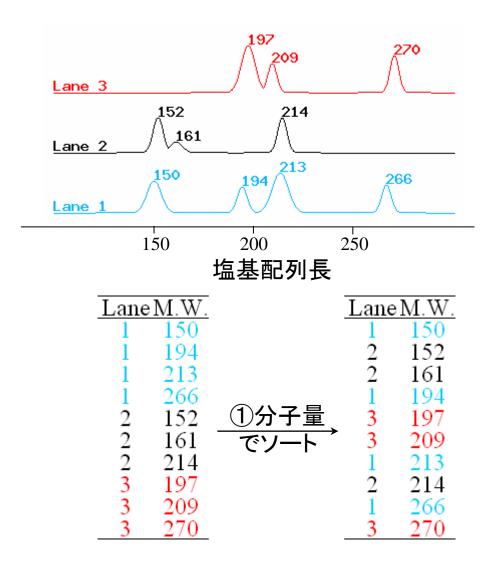


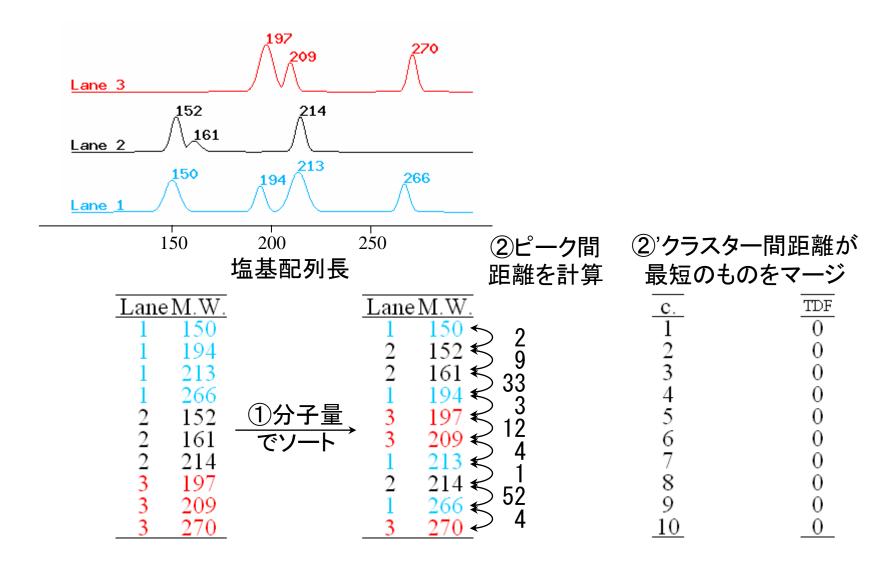
赤線を引く(同一ピークの認識)手段として、マイクロアレイ解析などでよく用いられている「クラスタリング」を利用可能

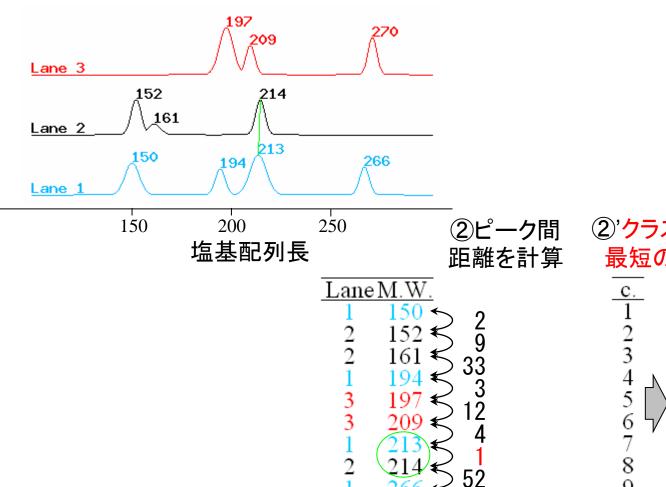
であることを講義中にふと思いついた(2006年ごろ)

■ ピークのアラインメントがとれている = 遺伝子発現行列を作れている

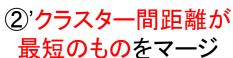




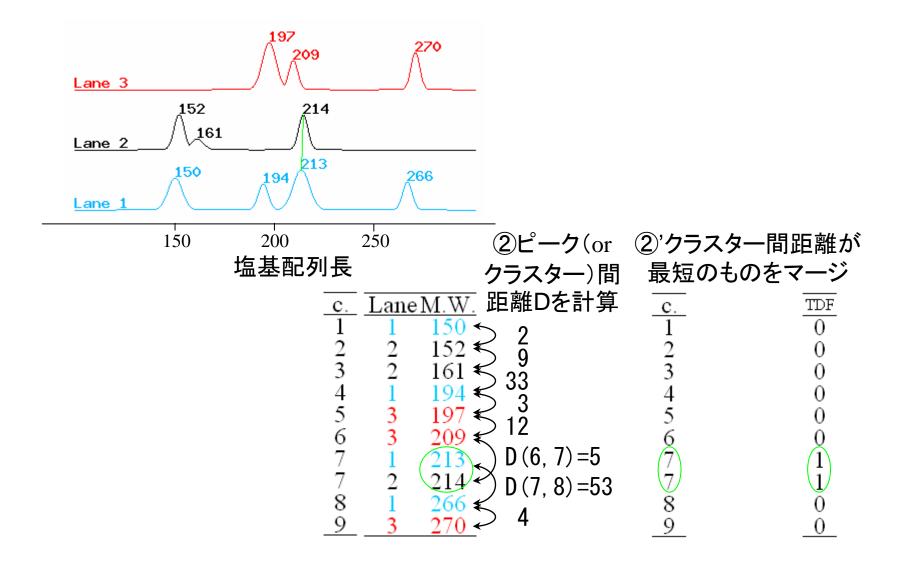


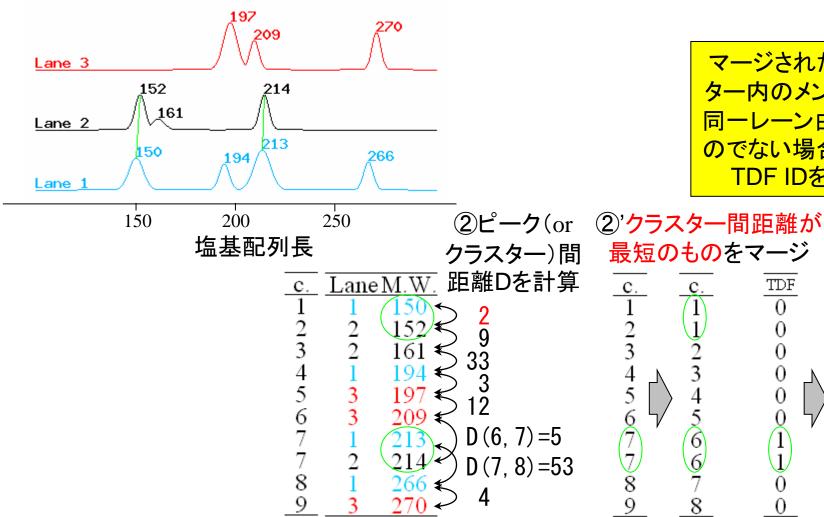


マージされたクラス ター内のメンバーが 同一レーン由来のも のでない場合に同じ TDF IDを付与



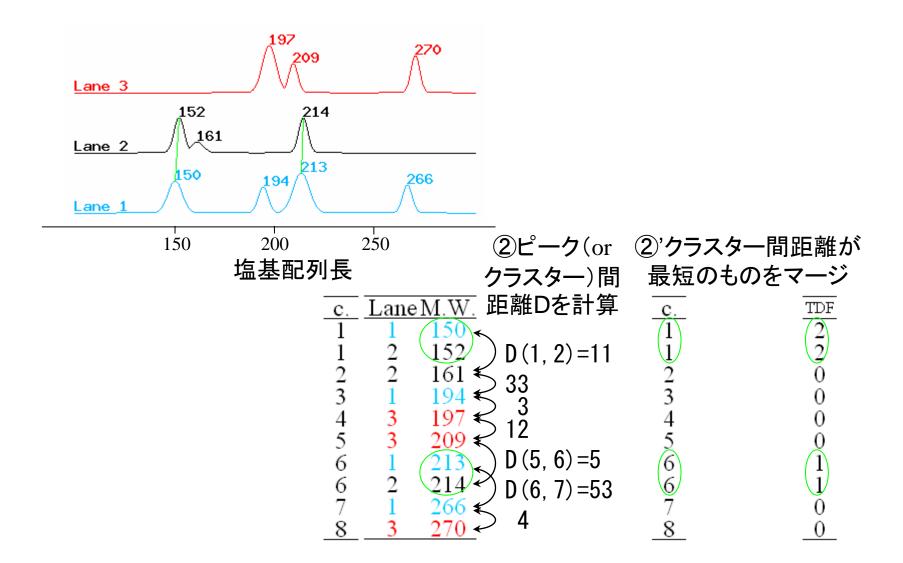
7/1	• • •		
<u>c.</u>	<u>c.</u>	TDF	TDF
1	1	0	0
2	2	0	0
2 3	2	0	0
4	4	0 _/	0 /
$\begin{array}{c} 4\\5\\6 \end{array}$	> 5 6	0 [0
65/	6	0 5/	0/
7	(7)	0	$(1)^{\star}$
8 9	$\sqrt{7}$	0	\bigcup
9	8	0	0
10	_9_	_0_	0

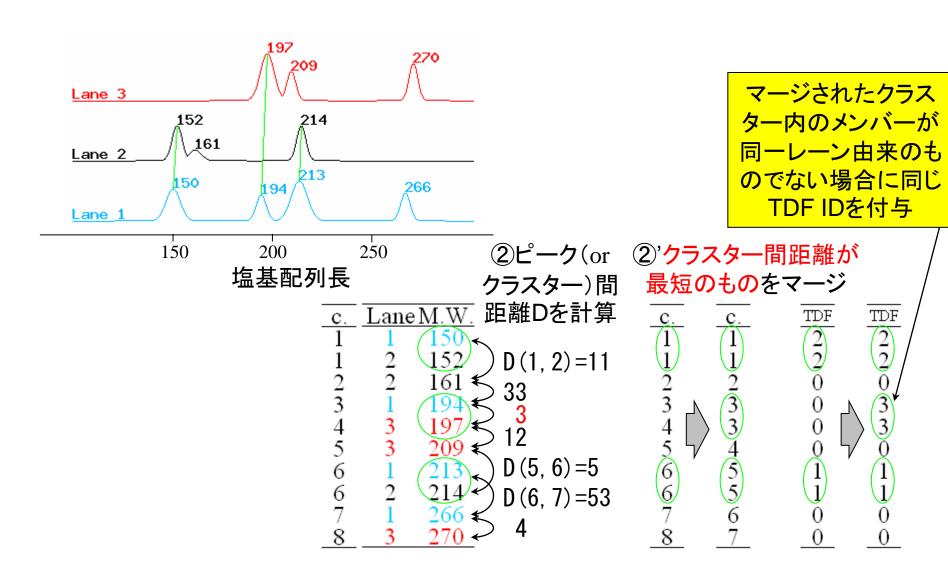


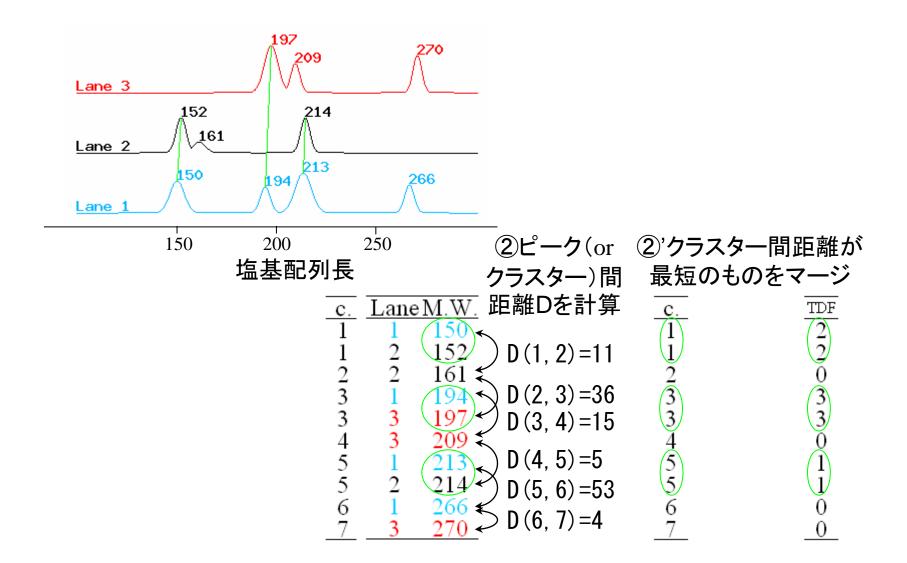


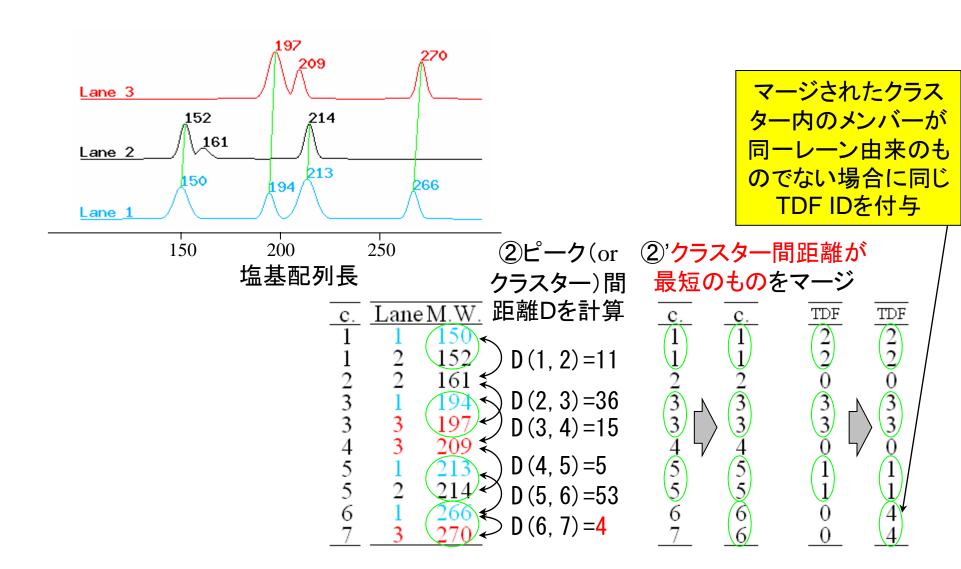
マージされたクラス ター内のメンバーが 同一レーン由来のも のでない場合に同じ TDF IDを付与

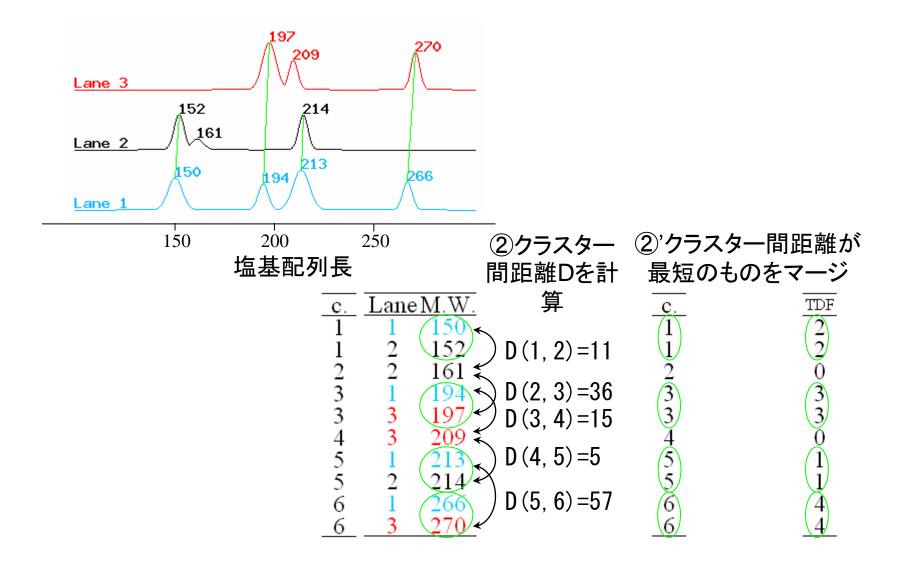
- 17 4 7		•	/
<u>c.</u>	$\frac{\overline{c}}{1}$	$\frac{\overline{\text{TDF}}}{0}$	$\frac{\overline{\text{TDF}}}{2}$
2 3	$\left(\underline{1}\right)$	Ŏ	$\overline{2}$
3	2	0	0
4 🙏	3	ŏ 7	0
5 6	\ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \	0	$\rangle \stackrel{0}{\circ}$
$\begin{pmatrix} 6 & 7 \\ 7 & 7 \end{pmatrix}$	6	$\begin{pmatrix} 0 & 1 \\ 1 & 1 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 1 \end{pmatrix}$
<i>(1)</i>	6	Ú)	Ú
8	7	0	0
9	_8_	_0_	_0_

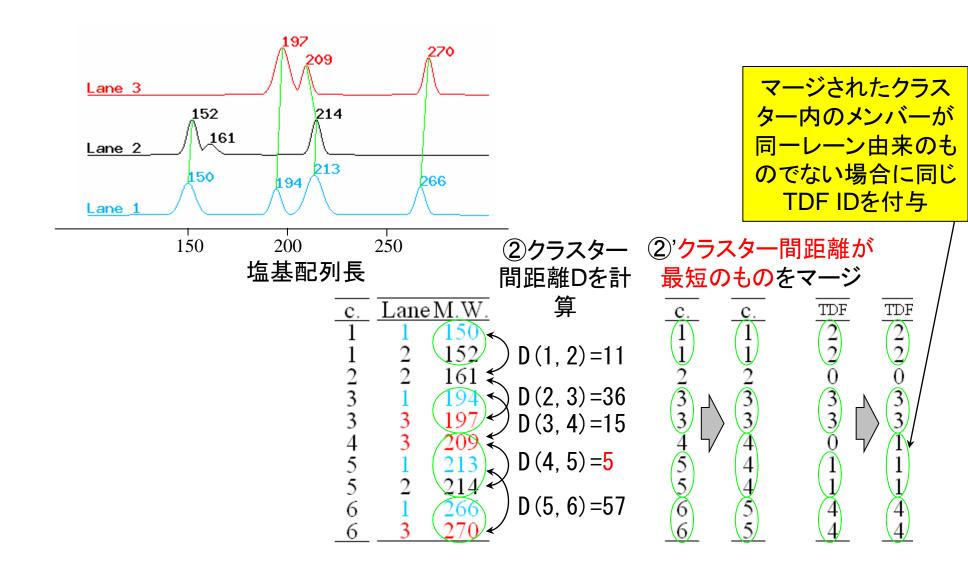


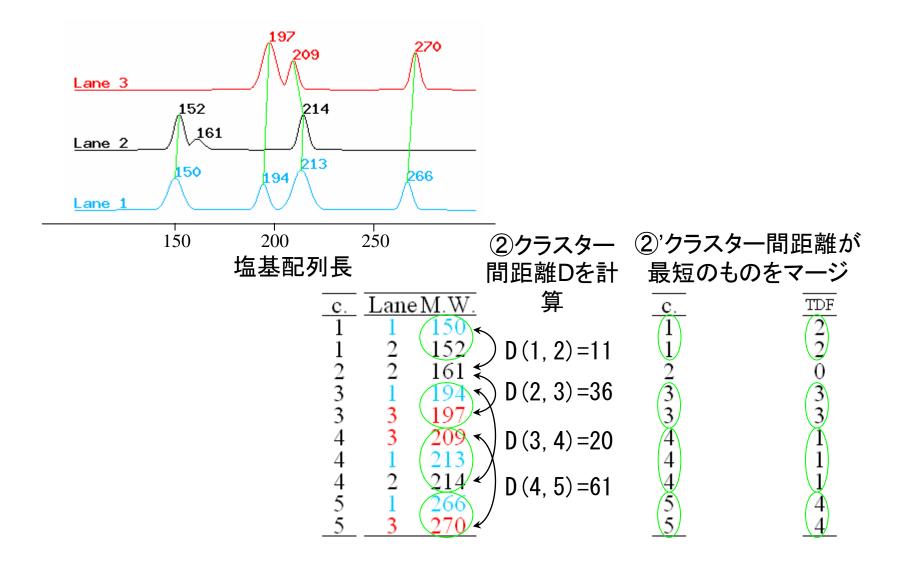


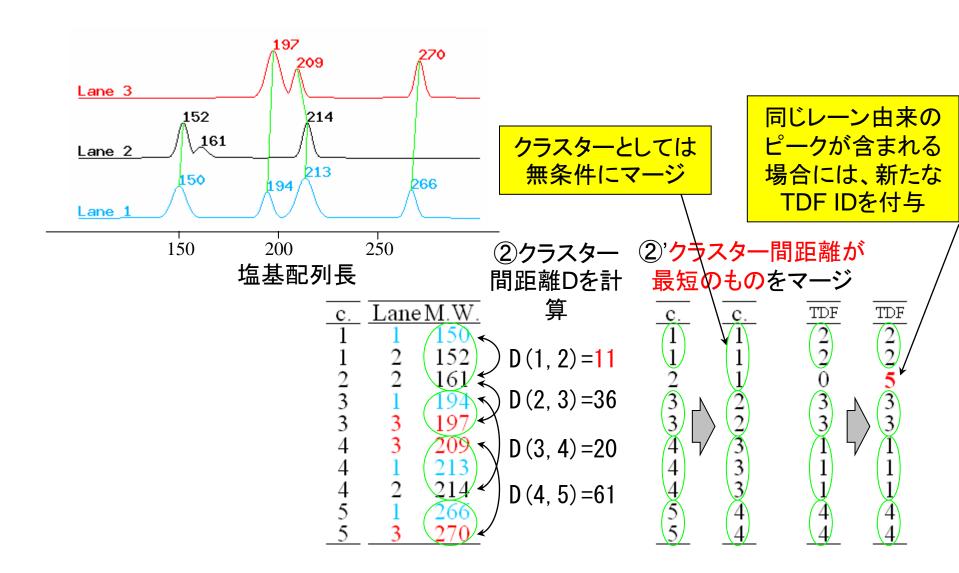


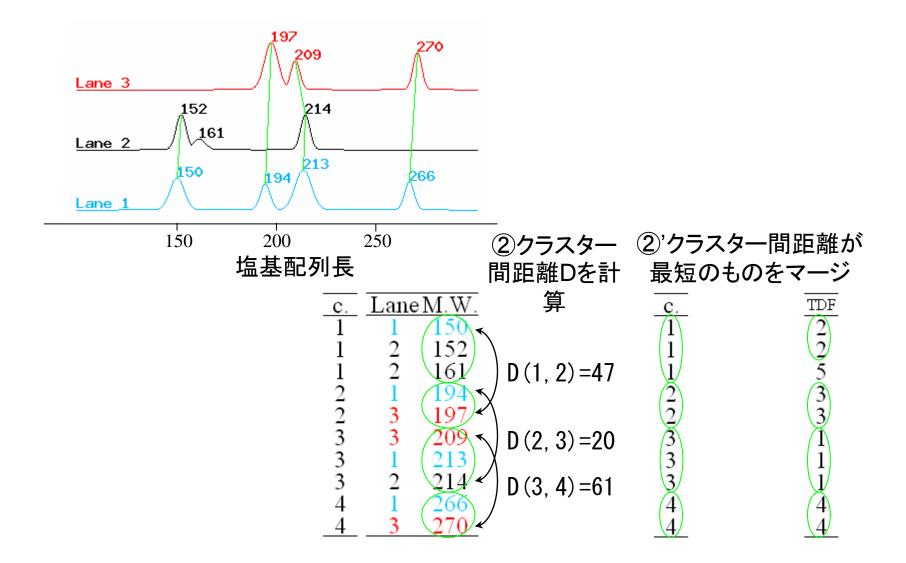


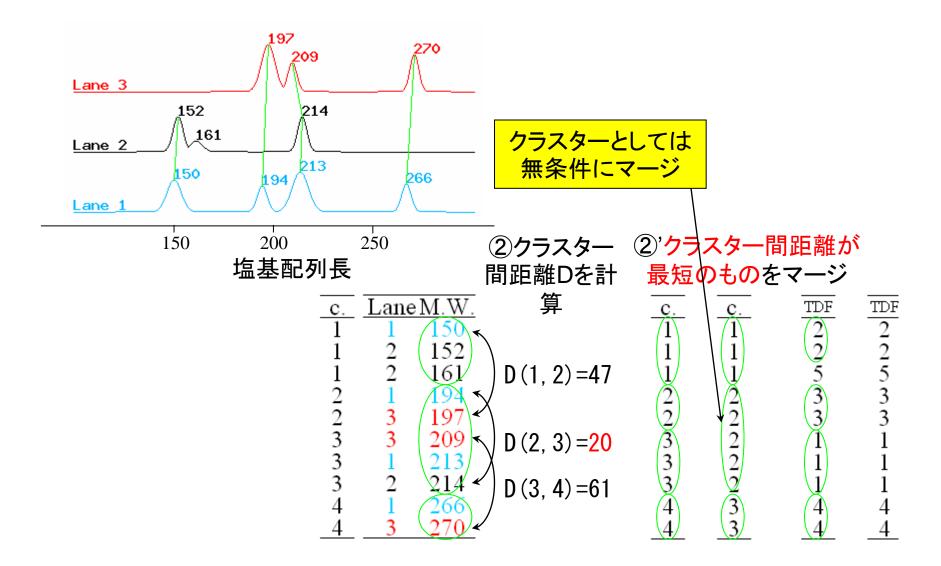


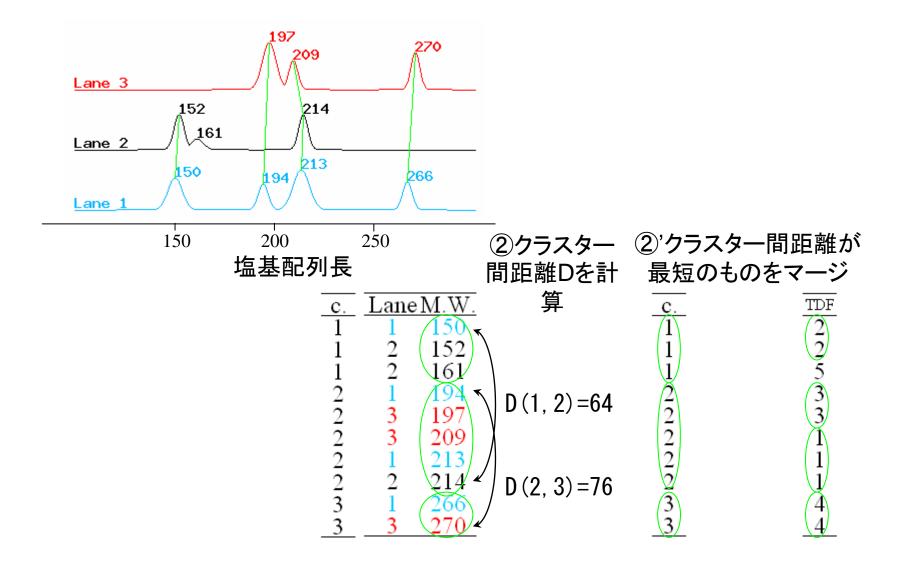


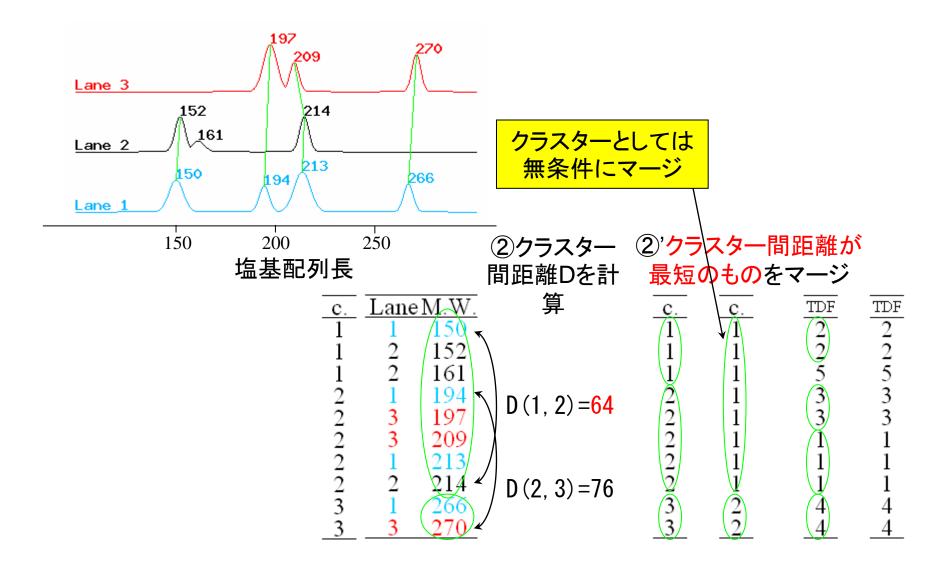


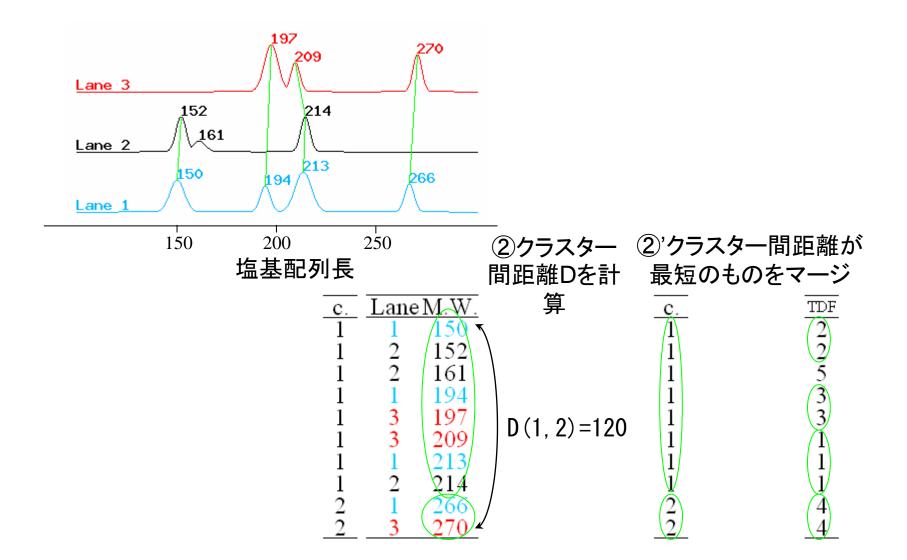






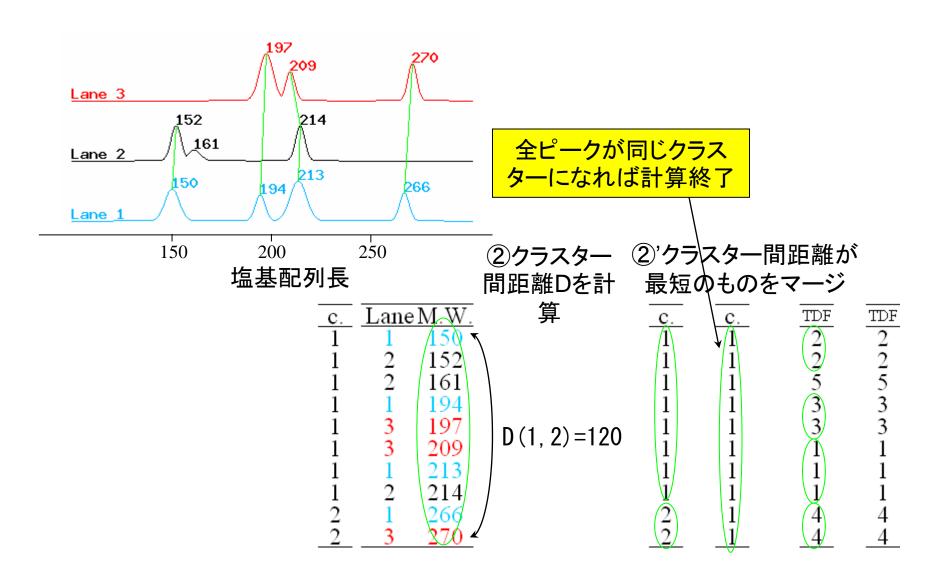






٠,

解析例(Peak alignment via clustering)



様々な遺伝子発現行列

横に関係して

1. 二群間比較

	A群		B群		
A1	A2		В1	B2	
$x_{1,1}^A$	$x_{1,2}^A$		$x_{1,2}^B$	$x_{1,2}^B$	
$x_{2,I}^A$	$x_{2,2}^A$		$x_{2,2}^B$	$x_{2,2}^B$	
$x_{i,I}^A$	$x_{i,2}^A$		$x_{i,2}^{B}$	$x_{i,2}^B$	
$x_{n,1}^A$	$\chi_{n,2}^A$		$x_{n,2}^B$	$x_{n,2}^B$	
	$\begin{array}{c} x_{1,1}^A \\ x_{2,1}^A \\ \dots \\ x_{i,1}^A \end{array}$	$ \begin{array}{c cccc} & A & A & A & A & A & A & A & A & A & $	$x_{1,1}^{A}$ $x_{1,2}^{A}$ $x_{2,1}^{A}$ $x_{2,2}^{A}$ $x_{i,1}^{A}$ $x_{i,2}^{A}$	$egin{array}{cccccccccccccccccccccccccccccccccccc$	$egin{array}{cccccccccccccccccccccccccccccccccccc$

2. 様々な組織(条件)

	S1	S2	S3	S4	
gene 1	$X_{1,1}$	$x_{1,2}$	$x_{1,3}$	$x_{1,4}$	
gene 2	$x_{2,1}$	$x_{2,2}$	$x_{2,3}$	$x_{2,4}$	
gene i	$X_{i,1}$	$X_{i,2}$	$X_{i,3}$	$X_{i,4}$	
gene n	$X_{n,1}$	$X_{n,2}$	$X_{n,3}$	$X_{n,4}$	

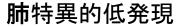
光刺激

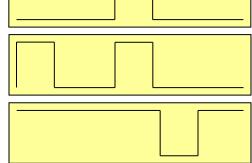
3. 時系列データ

	T1	T2	T3	T4	
gene 1	$x_{1,1}$	$x_{1,2}$	$x_{1,3}$	$x_{1,4}$	
gene 2	$x_{2,1}$	$x_{2,2}$	$x_{2,3}$	$x_{2,4}$	
gene i	$X_{i,1}$	$X_{i,2}$	$X_{i,3}$	$X_{i,4}$	
gene n	$X_{n,1}$	$X_{n,2}$	$X_{n,3}$	$X_{n,4}$	

脳特異的高発現

心臓と脳特異的高発現





組織特異的遺伝子検出にエントロピーを利用

■ 遺伝子iのエントロピー $H(x_i) = -\sum_{j=1}^{N} p_{ij} \log_2(p_{ij})$, where $p_{ij} = x_{ij} / \sum_{j=1}^{N} x_{ij}$

N:組織数(jの数) = 8

Hの取りうる範囲: $0 \le H \le \log_2 N \to 0 \le H \le 3$

				ı				
	x_{ij}	遺伝子1	遺伝子2	遺伝子3	事任子4	遺伝子5		
	組織1	1	5	6	4	10		
	組織2	0	2	6	4	10		
	組織3	0	1	6	4	10	:	
j	組織4	9	2	6	4	10	:	
	組織5	0	4	6	10	4		'
	組織6	0	6	6	4	10		
	組織7	0	3	6	4	10		
	組織8	0	5	6	4	10		
	$\sum_{j} x_{ij}$	10	28	48	38	74		

				_	_	1
p_{ij}	遺伝子』	遺伝子2	遺伝子3	遺伝子4	遺伝子5	
1	0.10	0.18	0.13	0.11	0.14	
2	0.00	0.07	0.13	0.11	0.14	
3	0.00	0.04	0.13	0.11	0.14	
4	0.90	0.07	0.13	0.11	0.14	
5	0.00	0.14	0.13	0.26	0.05	7/
6	0.00	0.21	0.13	0.11	0.14	
7	0.00	0.11	0.13	0.11	0.14	
8	0.00	0.18	0.13	0.11	0.14	
Σ_j	1.00	1.00	1.00	1.00	1.00	

			i			
- <i>p</i> ij log2(<i>p</i> ij)	遺伝子』	遺伝子2	遺伝子3	遺伝子4	遺伝子5	
1	0.33	0.44	0.38	0.34	0.39	
2	0.00	0.27	0.38	0.34	0.39	
3	0.00	0.17	0.38	0.34	0.39	
4	0.14	0.27	0.38	0.34	0.39	
5	0.00	0.40	0.38	0.51	0.23	
6	0.00	0.48	0.38	0.34	0.39	
7	0.00	0.35	0.38	0.34	0.39	
8	0.00	0.44	0.38	0.34	0.39	
\sum_{j}	0.47	2.83	3.00	2.90	2.96	

組織特異的遺伝子は低いエントロピー

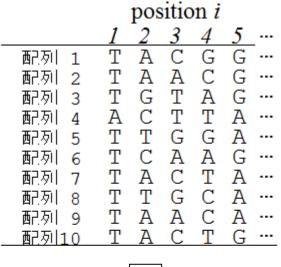
そうでないものは高い値

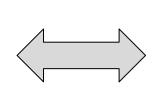
配列モチーフなどの表現にエントロピーを利用

■ position iの情報量 $IC_i = \frac{\log_2(N)}{2} - H(x_i)$

N:塩基**の**種類数=4

Hの取りうる範囲: $0 \le H \le \log_2 N$





2 7	Sequence logoは
1.5	エントロピーを計算 してるだけです
27 1	$ p_{1,4} = 90\% $ $ p_{5,3} = 50\% $
0.5	$\begin{cases} p_{5,1} = 50\% \end{cases}$
0	
	<u> </u>
IC	1.53 0.24 0.03 0.03 1.00

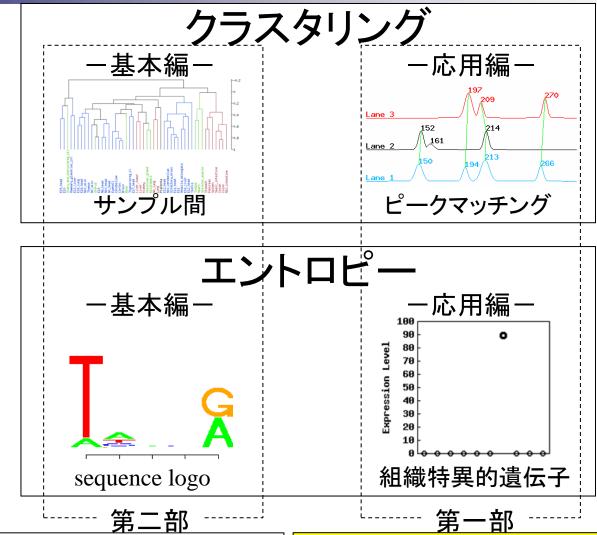


	x_{ij}	1	2	3	4	5	
	Aの数(j=1)	1	5	3	2	5	
j	Cの数 <i>(j=2)</i>	0	2	3	3	0	 $ \mathcal{Y} $
	Gの数 (j=3)	0	1	2	2	5	 L,
	Tの数 (j=4)	9	2	2	3	0	\
	$\sum_{j} x_{ij}$	10	10	10	10	10	

	p_{ij}	1	2	3	4	5	
\rangle	1	0.1	0.5	0.3	0.2	0.5	
	2	0.0	0.2	0.3	0.3	0.0	
	3	0.0	0.1	0.2	0.2	0.5	
	4	0.9	0.2	0.2	0.3	0.0	
	Σ_j	1.0	1.0	1.0	1.0	1.0	

	$-p_{ij}\log_2(p_{ij})$	1	2	3	4	5				
\rangle	1	0.33	0.50	0.52	0.46	0.50				
	2	0.00	0.46	0.52	0.52	0.00				
	3	0.00	0.33	0.46	0.46	0.50				
	4	0.14	0.46	0.46	0.52	0.00				
	$H = \Sigma_j$	0.47	1.76	1.97	1.97	1.00				





次世代シーケンサデータもRのコピペで解析可能
→ 頭脳労働

バイオインフォ要素技術の習得は大事だが、それだけでも様々な種類の実験データに対応可能



謝辞



東京大学 大学院農学生命科学研究科

清水謙多郎 教授

中井雄治 特任准教授

放射線医学総合研究所 先端遺伝子発現研究G

安倍真澄 グループリーダー

荒木良子 チームリーダー

グラント

若手研究(B) (H19-20年度):「全生物種のトランスクリプトーム解析を加速するHiCEPデータ高速解析手法の開発」(代表)

挿絵担当:門田雅世

