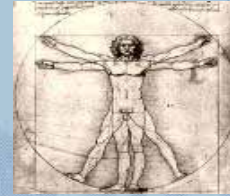




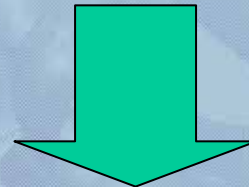
Biclustering of Expression Data Using Simulated Annealing

Kenneth Bryan, Pàdraig Cunningham, Nadia Bolshakova

Il genoma di molti organismi è stato sequenziato:



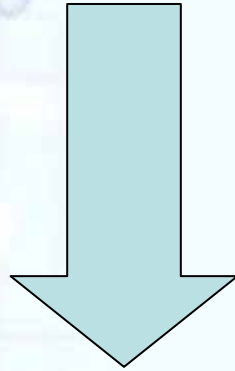
La maggior parte dei geni (> nell'uomo) ha funzione sconosciuta.



- **Conoscere la funzione di tutti i geni (studiando l'espressione genica)**
- **Raggruppare i geni in base alle funzioni**

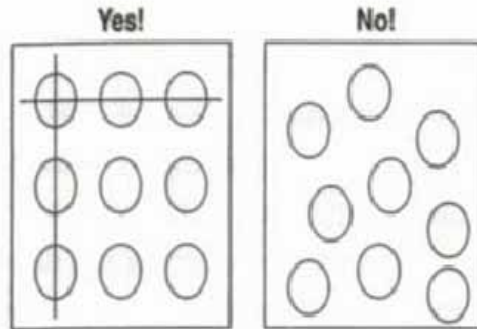
Tecnologia microarray

MICROARRAY

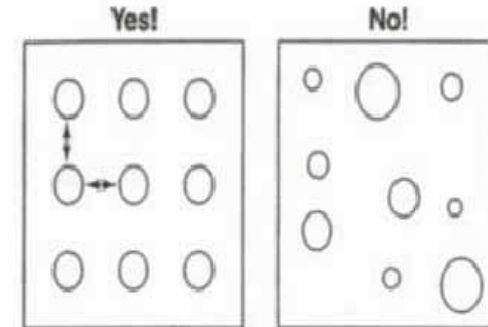


matrice ordinata di elementi microscopici su un substrato planare, che consente il legame di specifici geni o di prodotti di geni.

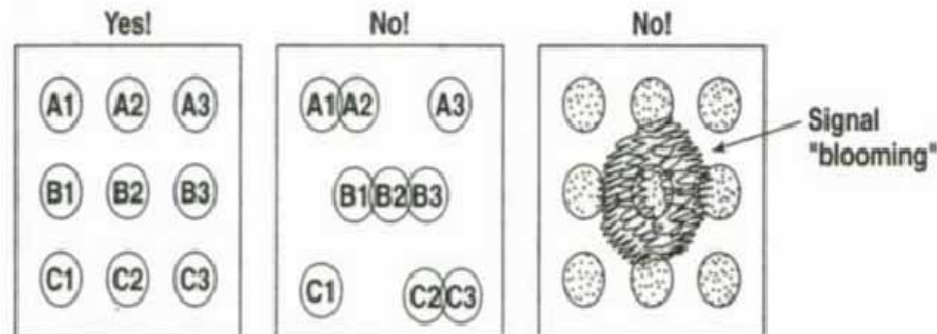
Rows & columns



Uniform size and spacing



Unique address

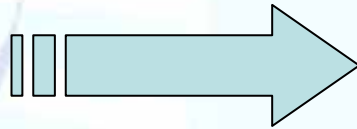


Rapida deposizione, individuazione e quantificazione degli spot

Substrari:

- **Vetro (favorisce il legame con la sonda)**
- **Materiali plastici**
- **Silicio**
- **Filtri di nylon**
- **Nitrocellulosa**

PLANARE

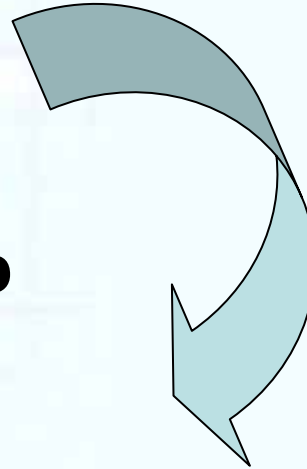


- Ottima automatizzazione della procedura di stampa tramite pin e ugelli ik-jet
- Accurato scattering grazie alla precisa individuazione della distanza tra gli elementi ottici dello scanner e la superficie

Vantaggio:

Monitorare l'attività di migliaia di geni contemporaneamente e in condizioni diverse:

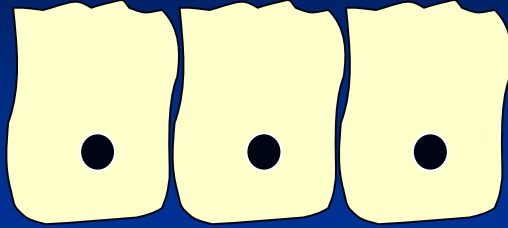
- **Tipo di tessuto**
- **Paziente**
- **Ambienti di sviluppo**



- 1. Studio di tessuti**
- 2. Scoprire le basi molecolari di una malattia(espressione genica, co-regolazione)**
- 3. Analisi del meccanismo di azione dei farmaci**

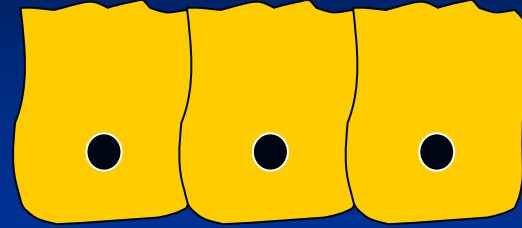
Gene A

Cellula normale



Basso livello di espressione

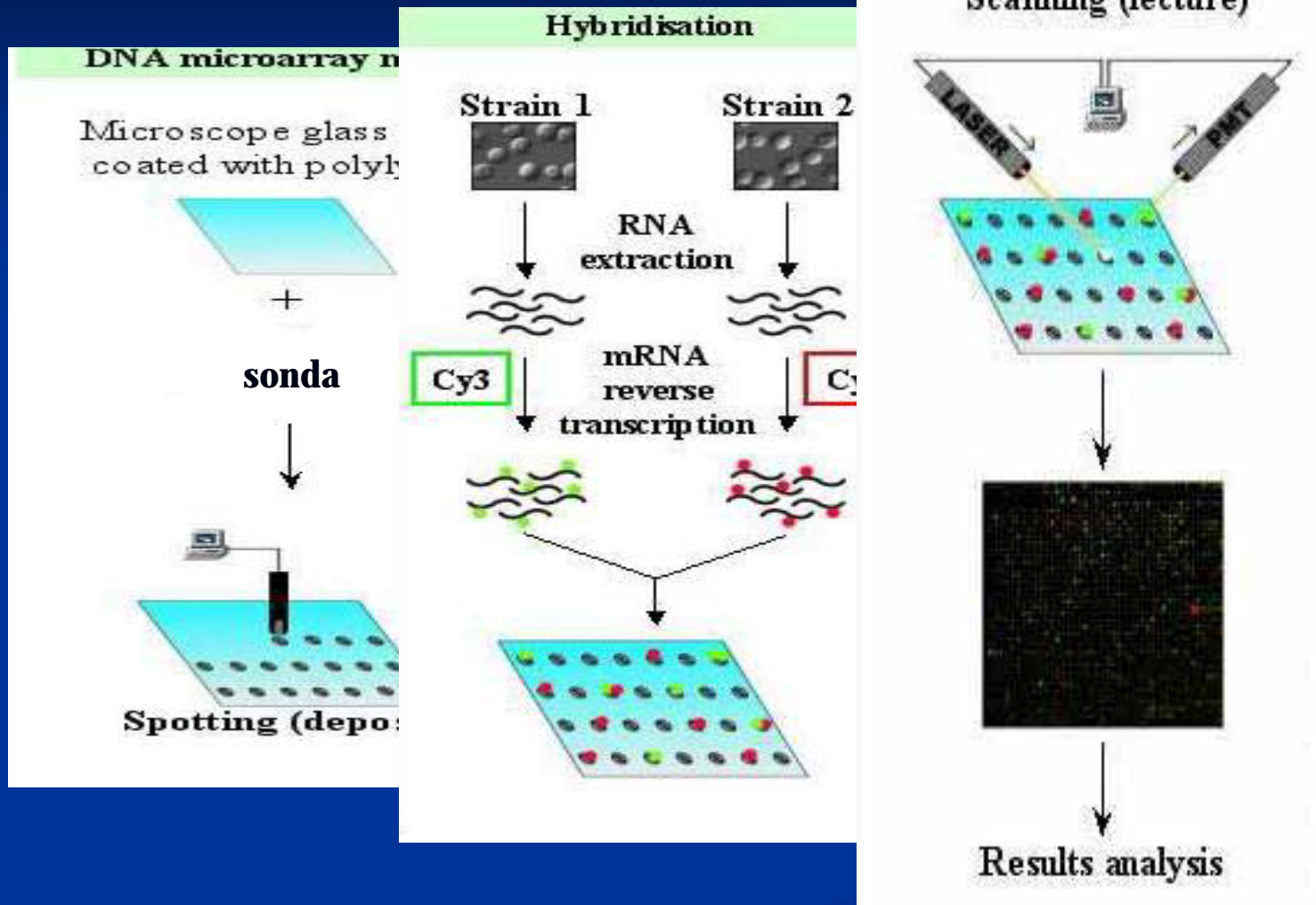
Cellula malata cancro



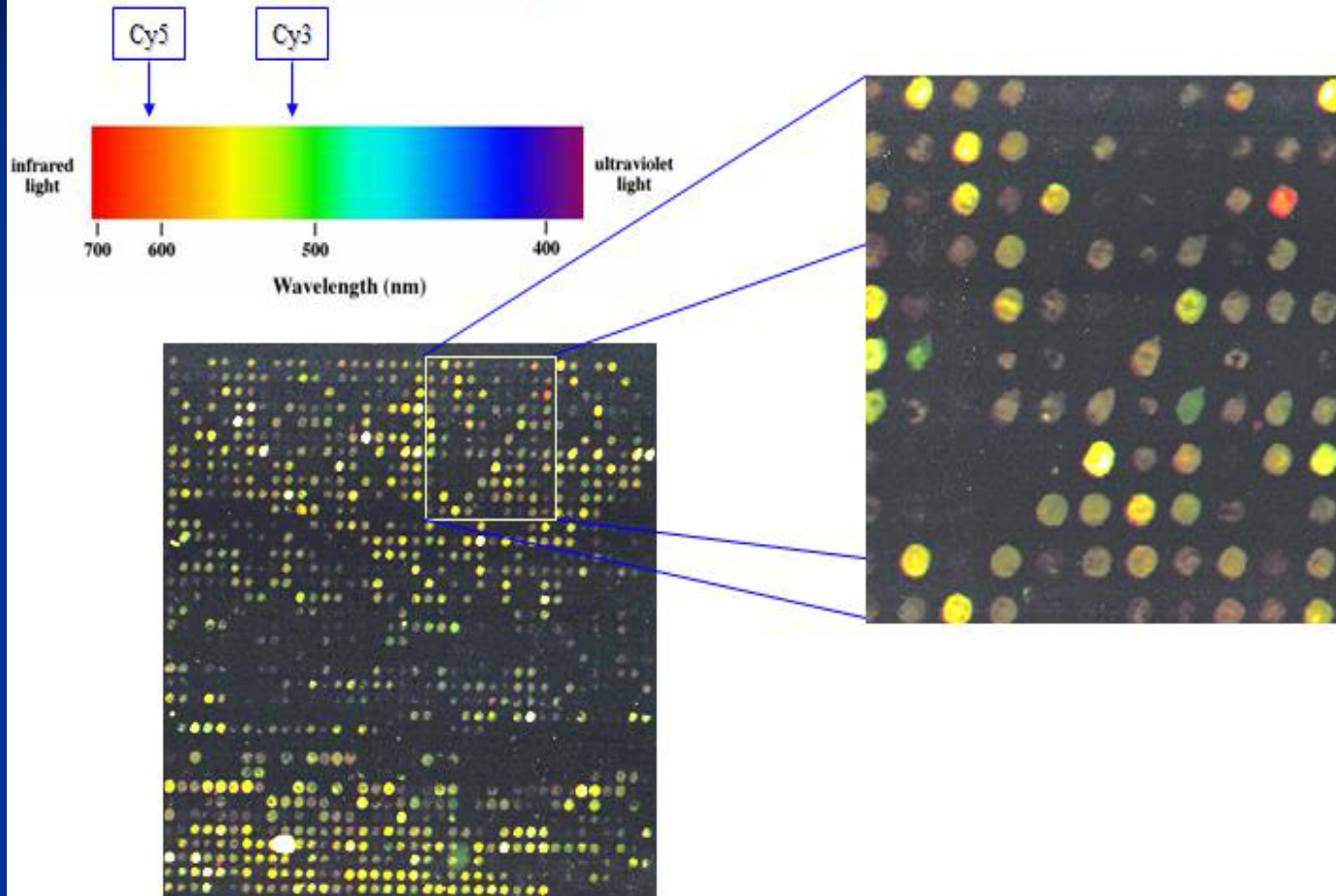
Alto livello di espressione

Il gene A potrebbe avere un ruolo importante nell'insorgenza del cancro!

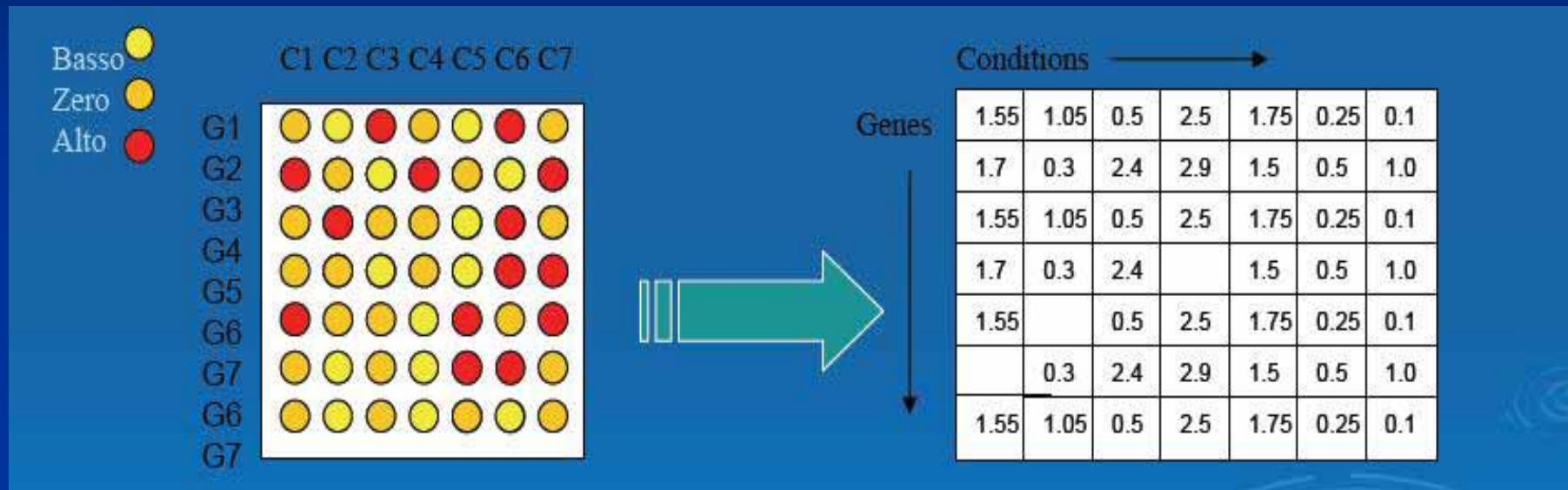
Preparazione standard : tecnica spotting



MicroArray con Fluorocromi



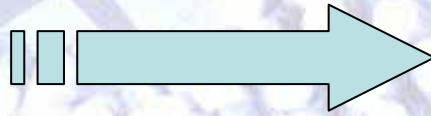
Gli esperimenti vengono eseguiti tramite una matrice simile a questa:



I diversi colori indicano diversi livelli di espressione.

Per facilitare l'analisi computazionale la matrice viene convertita in una matrice numerica.

A (n x m)

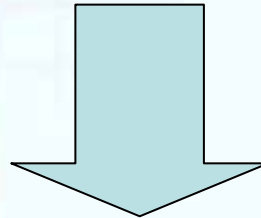


n = numero di geni

**m = numero
condizioni (array)**

**a_{ij} = livello di espressione del gene *i*
nella condizione *j***

Analisi molto difficile



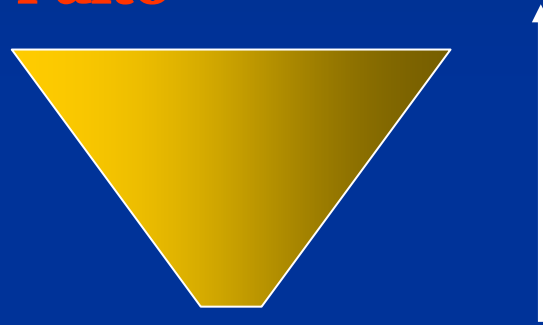
DATA MINING

**(insieme di strumenti per
estrapolare pattern)**

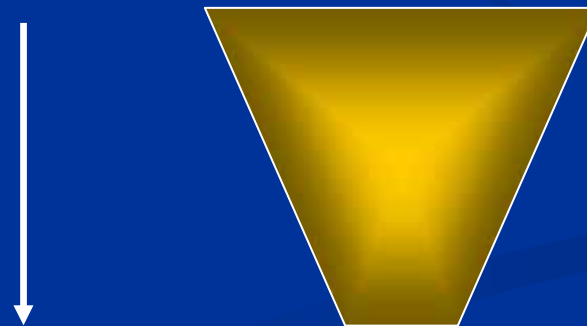
CLUSTERING

Tecnica di analisi di dati per selezionare e raggruppare elementi omogenei, basato sul concetto di distanza tra due elementi.

a) Basso verso l'alto



b) Alto verso il basso



Lo scopo, nell'analizzare la matrice, è quello di poter identificare gruppi di geni con profili di espressione simili. Spesso i geni appartenenti allo stesso cluster sono detti coespressi.

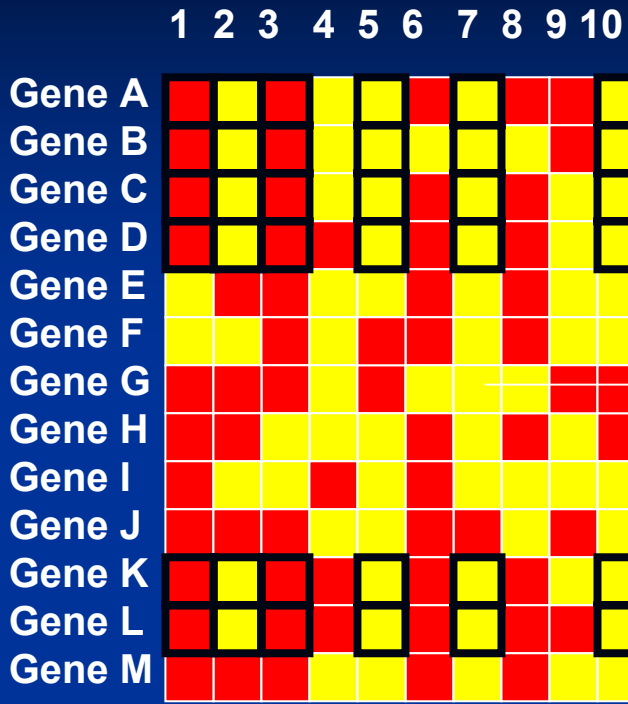
La ragione per cui si cercano geni coespressi sono:

- 1) si è evidenziato che geni funzionalmente correlati sono coespressi**
- 2) geni coespressi possono dare informazioni sui meccanismi regolatori**

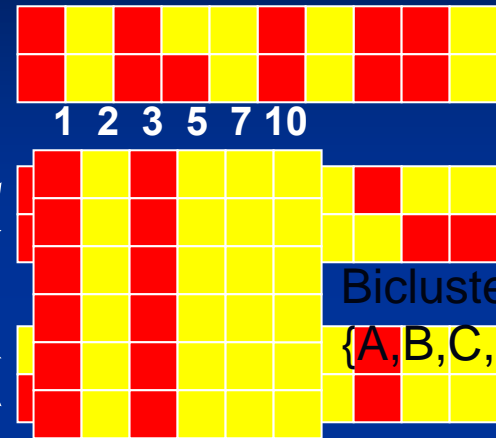
	Conditions							
	A	B	C	D	E	F	G	H
Gene 1	Red	Green	White	Red	Green	Green	Red	Red
Gene 2	White	White	White	White	White	White	White	White
Gene 3	White	White	White	White	White	White	White	White
Gene 4	Red	Green	White	Red	Green	Green	Red	Red
Gene 5	White	White	White	White	White	White	White	White
Gene 6	White	Green	White	White	Green	Green	White	White
Gene 7	White	Green	White	White	Green	Green	White	White
Gene 8	White	White	White	White	White	White	White	White
Gene 9	Red	Green	White	Red	Green	Green	Red	Red

Clustering

	A	B	C	D	E	F	G	H
Gene 1	Red	Green	White	Red	Green	Green	Red	Red
Gene 4	Red	Green	White	Red	Green	Green	Red	Red
Gene 9	Red	Green	White	Red	Green	Green	Red	Red



Clustering...



Bicluster {1,2,3,5,7,10}
{A,B,C,D,E,F}

Non esistono livelli di espressione simili considerando tutte le condizioni...

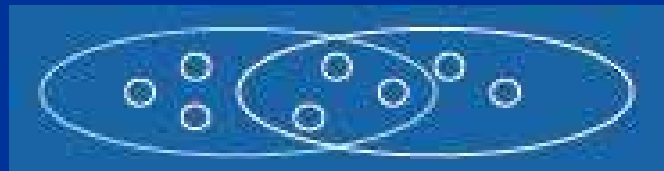
Soluzione: Cluster simultaneo di righe e colonne (Biclustering)

BICLUSTERING

Tecnica di data mining che raggruppa i geni rispetto a sottoinsiemi di condizioni.

Ciò ci permette di rappresentare meglio:

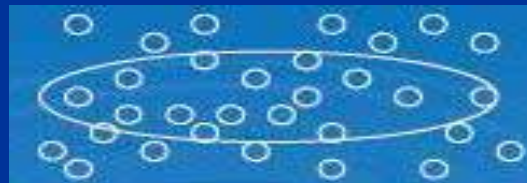
- **Cluster correlati**



- **Segnali locali**



- **Rumore**



Bicluster = sottomatrice della matrice di espressione genica che soddisfa alcune caratteristiche di omogeneità

Bicluster(I,J) = sottomatrice $k \times s$

$I = \{i_1, \dots, i_k\}$ $J = \{j_1, \dots, j_s\}$

Può essere di vari modi:

- **valori costanti**
- **valori costanti sulle righe o sulle colonne**
- **valori coerenti**

Algoritmo di Cheng e Church (node deletion)

Similarità = coerenza di geni e condizioni nel subset
mean squared residue score (H)

$$R(a_{ij}) = a_{ij} - a_{iJ} - a_{Ij} + a_{IJ}$$

a_{ij} **media dei valori nella subm. lungo la linea i**

a_{Ij} **media dei valori nella subm. lungo la colonna j**

a_{IJ} **media di tutti i valori nella submatrice**

$$H_{IJ} = \frac{1}{|I| |J|} \sum_{i \in I, j \in J} R^2$$

- **Matrice costante: 0**
- **Matrice con singolo elemento: 0**
- **Matrice con elementi scelti random in un intervallo [a,b] , $H = (b-a)^2/12$**

H elevato  **i dati non sono correlati**

H basso  **i dati sono correlati**

Media-matrice = 6.5

1	2	3	2	
4	5	6	5	media r.
7	8	9	8	
10	11	12	11	

$$R(1) = 1 - 2 - 5.4 + 6.5 = 0.1$$

$$R(2) = 2 - 2 - 6.4 + 6.5 = 0.1$$

.....

$$R(12) = 12 - 11 - 7.4 + 6.5 = 0.1$$

5.4 6.4 7.4
media c.



$$H = (0.01 \times 12) / 12 = 0.01$$

Se i valori fossero stati scelti in modo casuale tra [1,12]

$$H = (12-1)^2 / 12 = 10.08$$

OBIETTIVO = data una matrice di espressione **A** e un valore δ , cercare un bicluster con massimo un $H=\delta$

• Input: **A** e $\delta \geq 0$ minimo **H** accettabile

• Output: **A(I,J)**, detto δ bicluster ($H \geq \delta$)

• Inizializzazione: **A(I,J)=A**

• Iterazione: si calcola a_{iJ} per ogni $i \in I$, a_{Ij} per ogni j di **J**, a_{IJ} e **H(I,J)**.

se $H \leq \delta$ si ritorna **A(I,J)**

si cerca la riga con valore più alto $d(i) = \frac{1}{|J|} \sum_{j \in J} (a_{ij} - a_{iJ} - a_{Ij} + a_{IJ})^2$

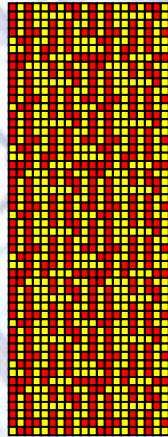
e si cerca la colonna con valore più alto $d(j) = \frac{1}{|I|} \sum_{i \in I} (a_{ij} - a_{iJ} - a_{Ij} + a_{IJ})^2$

si elimina la colonna o la riga aggiornando **I** o **J**

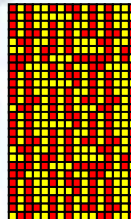
Input:

**Matrice di
dati**

$\delta = 300$



Score: 1,052



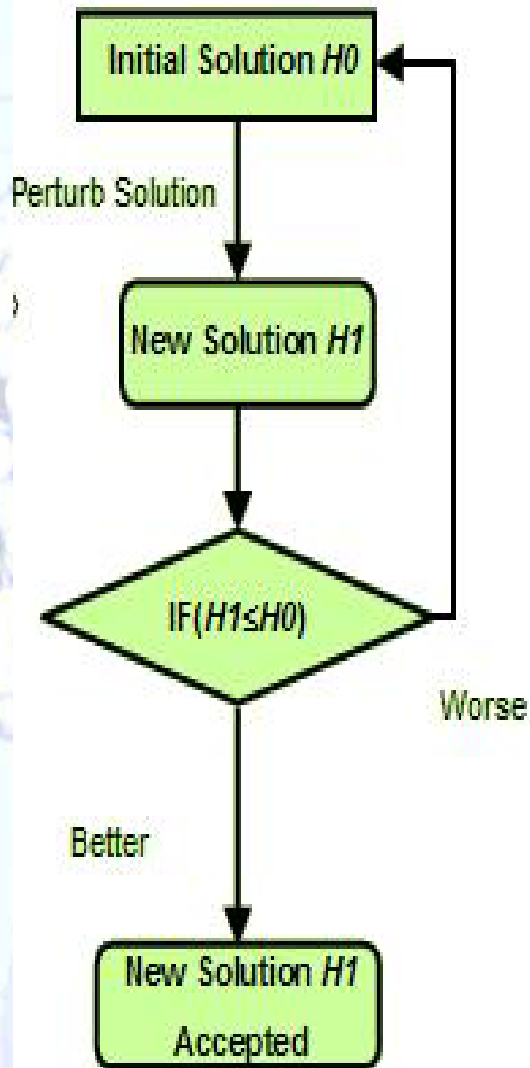
Score: 543



Score: 423



Score: 300



Ricerca Locale(s);

$s^* = s$; fine=false

repeat

$s :=$ migliore soluzione I(s)

if $f(s) < f(s^*)$ then

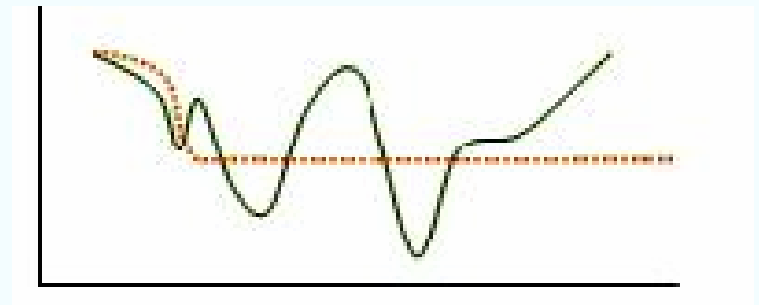
$s^* := s$;

else fine=true;

until not fine;

return s^* ;

end



SIMULATED ANNEALING

Nato per simulare l'annealing dei solidi

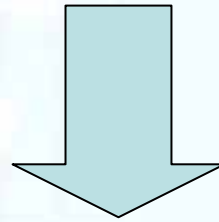


Processo termico atto ad ottenere stati di bassa energia in un solido immerso in un bagno di calore.

- 1) Si aumenta la temperatura ad un valore tale che il solido si scioglie.**
- 2) Si decresce lentamente la temperatura fino a quando la particelle si sistemano in uno stato base in cui l'energia è minima.**



Il sistema si dice in equilibrio termico alla temperatura T se la probabilità $p(E_i)$ di uno stato avente energia E_i è governata dalla distribuzione di Boltzmann.



$$p(E_i) = g_i \frac{\exp\left(-\frac{E_i}{K_B T}\right)}{\sum_J \exp\left(-\frac{E_j}{K_B T}\right)}$$

ALGORITMO DI METROPOLIS

Algoritmo che simula l'evoluzione del solido verso l'equilibrio termico.

- Dato lo stato corrente i del solido con energia E_i
- Perturbandolo si genera un nuovo stato j con E_j
- Se $E_i - E_j \geq 0$ lo stato j diviene il nuovo stato corrente
- Altrimenti lo stato j diviene il nuovo stato corrente con probabilità

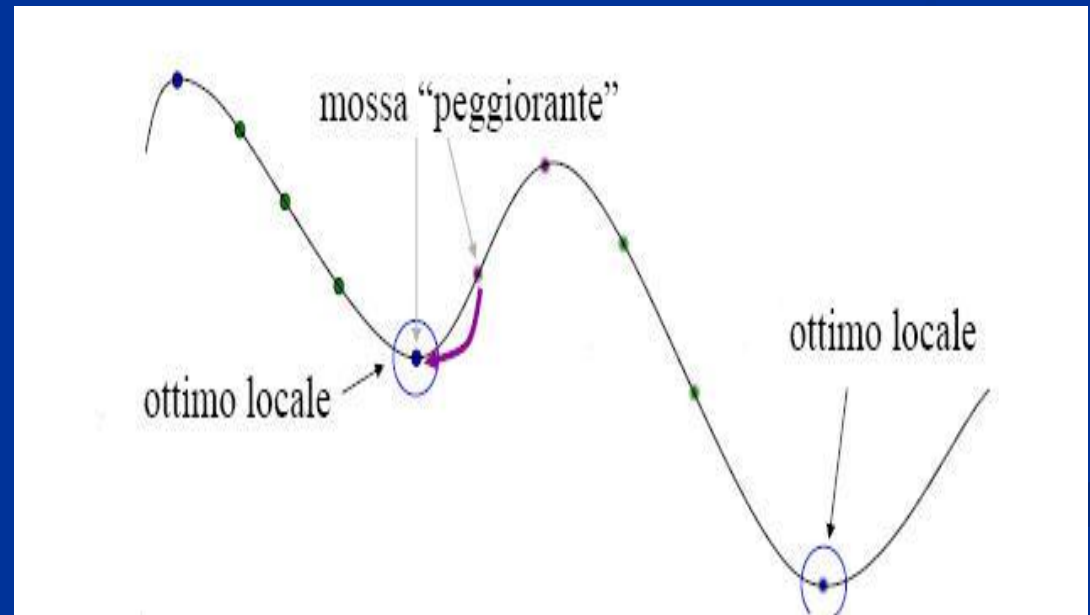
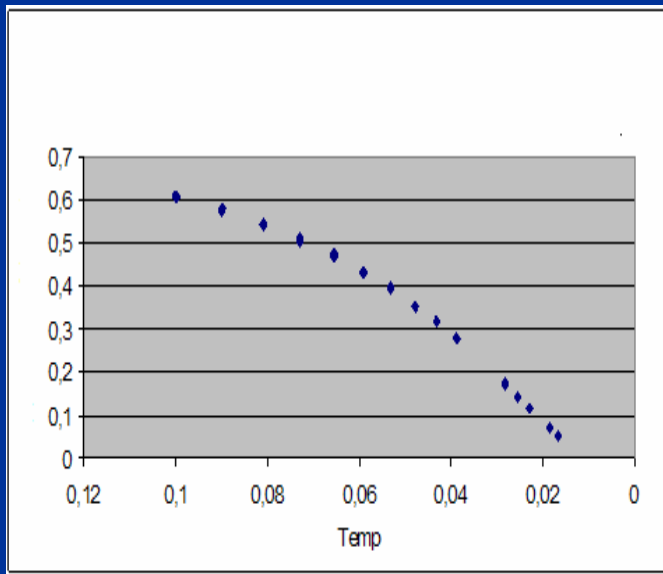
$$e^{\left(\frac{E_i - E_j}{kT}\right)}$$

T = temperatura del bagno di calore

K = costante di Boltzmann

All'inizio la ricerca salta da un punto all'altro individuando le direzioni o le aree in cui è più probabile trovare l'ottimo globale.

A basse temperature le soluzioni vengono localizzate in un dominio promettente.

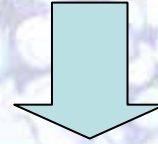


Simulated annealing problemi di ottimizzazione

- **soluzione di un problema = stati di un sistema fisico**
- **costo di una soluzione = energia di uno stato**
 - funzione fitness che definisce se una perturbazione determina un miglioramento**
- **parametro di controllo = temperatura**

To ?

velocità di raffreddamento



Elevata = blocco in un minimo

Bassa = aumentano gli stati, quindi il costo

**La temperatura viene ridotta ad ogni stato
moltiplicandola per una costante chiamata**

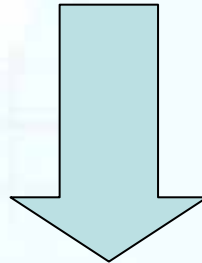
$0 < \text{cooling ratio} > 1$

**Metodo
logaritmico**

•La procedura che genera nuove soluzioni deve essere il più veloce possibile.

Criterio di interruzione

- **numero di soluzioni trovate è $<$ di un valore fissato**
- **valore finale di temperatura**



Stabiliti in base al problema che si sta affrontando

Parametri iniziali:

- **Funzione fitness = mean squared residue score**
- **T_0 = definita tramite esperimenti**
- **$T(i) = T(i-1)/(1+\sigma)$**
 $\sigma = 0.1$ cooling ratio

criteri di arresto:

- **$\delta = 300, 200, 100$ \longrightarrow node deletion (CCND/2)**
- **bicluster minimo 10x10 \longrightarrow node deletion modificato (CCND2)**

SAB($f, x_0, t_0, r, a, s, M, l, \delta$)

$x \leftarrow x_0$

$t \leftarrow t_0$

While($t > t_{Min}$)

 While($a_{count} < a$ AND $s_{count} < a$)

$x_{new} \leftarrow \text{GenerateNewSolution}(M, x)$

 if $f(x_{new}) < f(x)$

 then $x \leftarrow x_{new}$

 else if $\exp(-\frac{\Delta E}{T}) > \text{random}(0,1)$

 then $x \leftarrow x_{new}$

 if $(f(x_{new}) \leq \delta)$ AND $\text{SizeOf}(x_{new}) > \text{SizeOf}(x)$

 then $l \leftarrow \text{SizeOf}(x_{new})$

$t \leftarrow \text{cool}(t, \text{rate})$

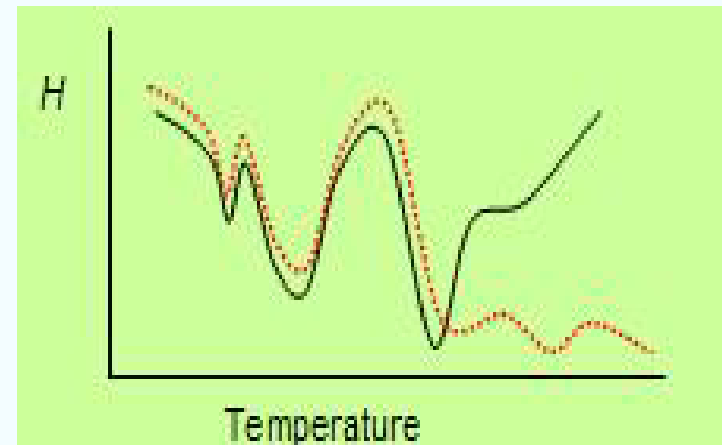
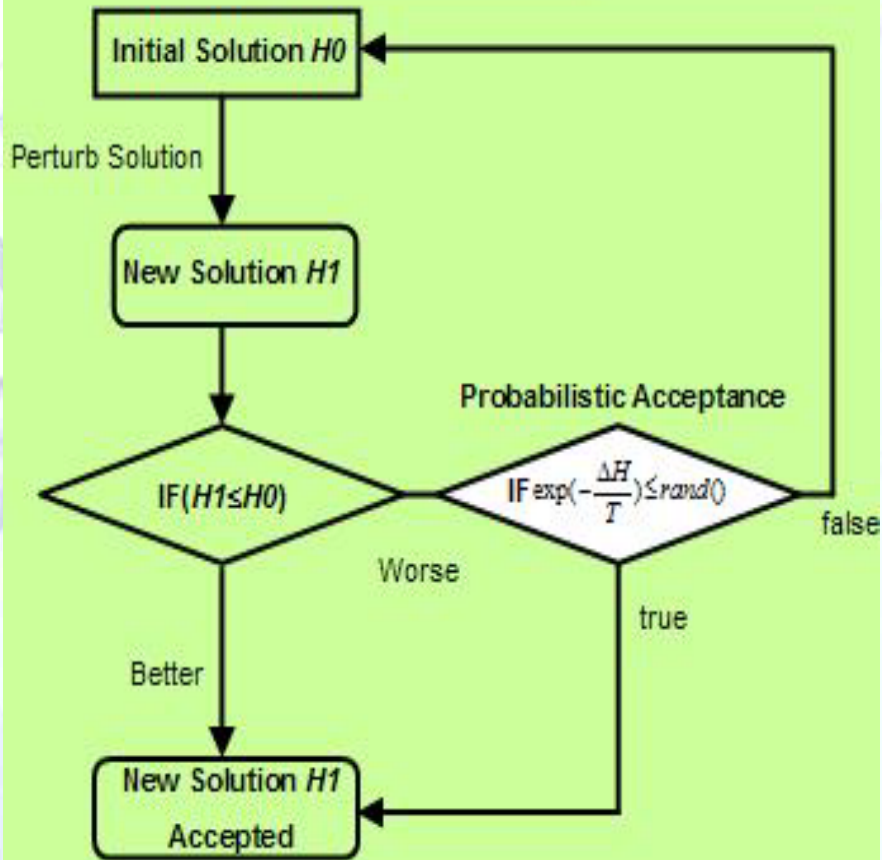
$x \leftarrow \text{NodeAdditon}(x)$

return x

Funzione genera(M,x)

```
rowx ← number of rows in x
colx ← number of rows in x
if(rowx ≥ colx)
  then generate random  $r$ , range[0,  $\frac{row_x}{col_x}$ ] ∈ ℝ
  if ( $r = 0$ )
    then flip random column in  $M$  to generate  $X_{new}$ 
  else
    then flip random row in  $M$  to generate  $X_{new}$ 
else if(rowx < colx)
  then generate random  $r$ , range[0,  $\frac{col_x}{row_x}$ ] ∈ ℝ
  if ( $r = 0$ )
    then flip random row in  $M$  to generate  $X_{new}$ 
  else
    then flip random column in  $M$  to generate  $X_{new}$ 
return  $x_{new}$ 
```


Simulated Annealing

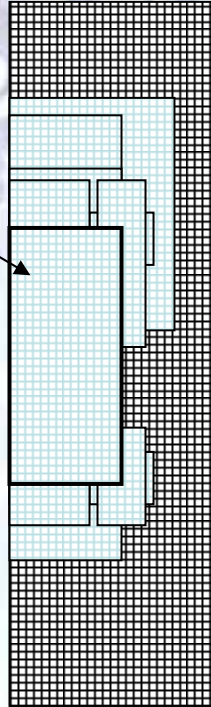


Simulated Annealing

Per evitare di cercare lo stesso bicluster questo viene sostituito con numeri random

Righe
(Geni)

Colonne
(Condizioni)



Mean Squared
Residue Score

1000

**Bicluster
massimo**

Delta si blocca e la
dimensione aumenta

Dati

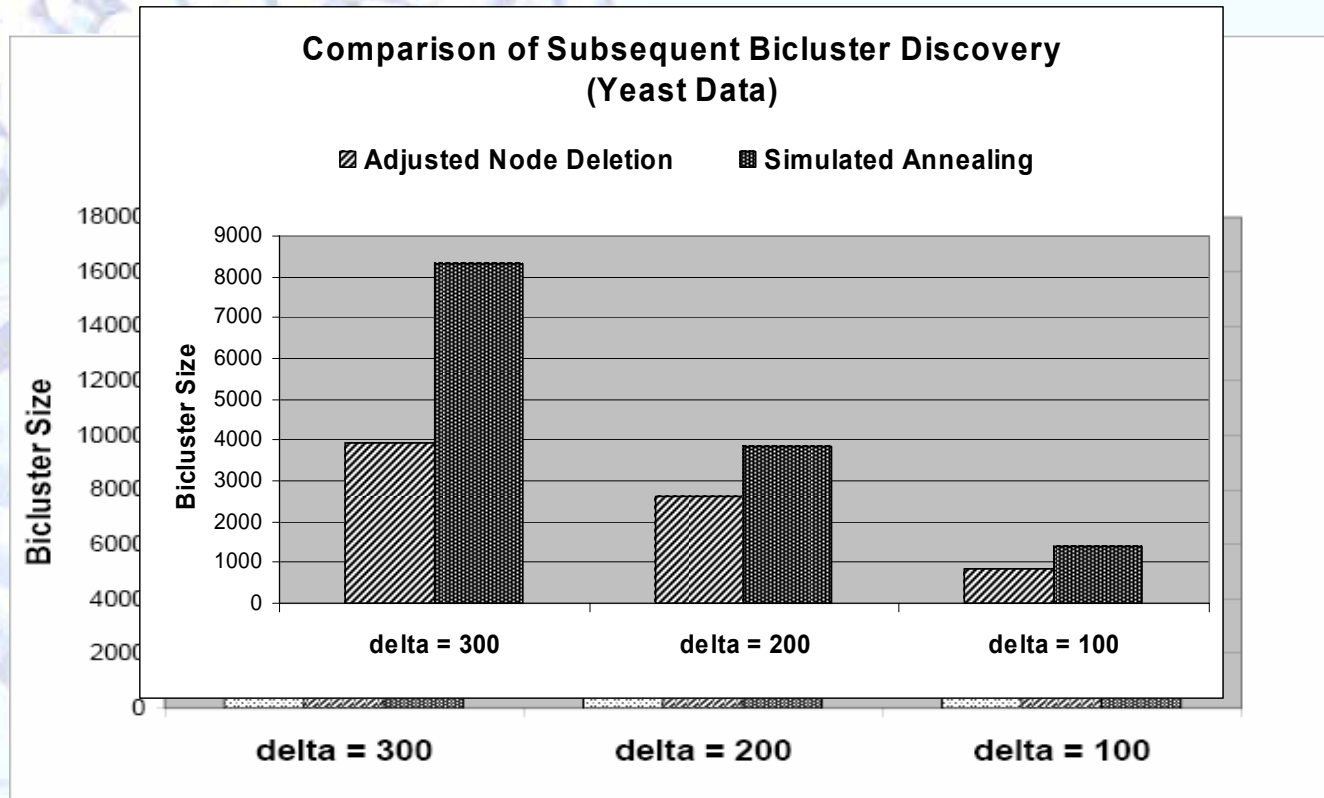
geni

condizioni

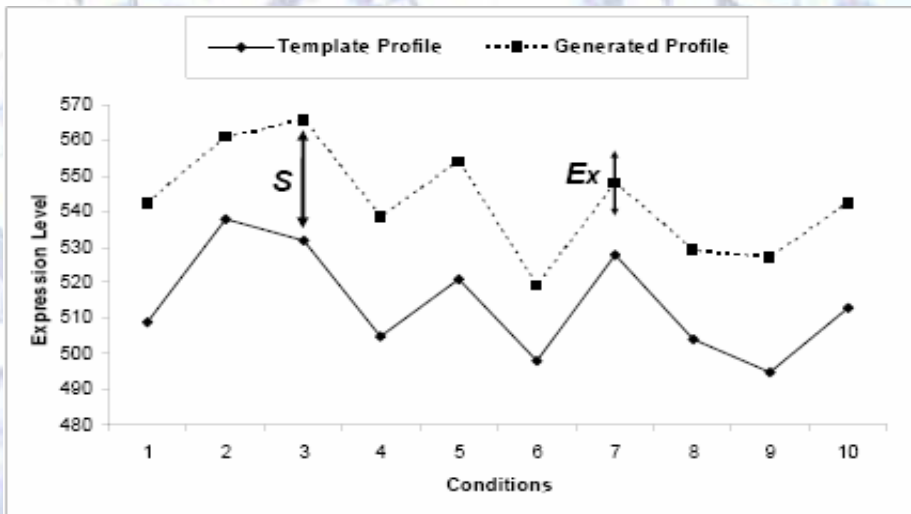
2884

17

lievito



Partendo dal profilo di espressione iniziale si effettua uno shift random



$$E_x = \sigma(x) \cdot \epsilon \cdot r_x$$

$\sigma(x)$ = deviazione standard del nuovo profilo genico

ϵ = costante [0,1] (livello errore e qualità del bicluster)

r_x = variabile random [-1,1] (il livello di espressione sia più alto o iù basso)

$$X = \{x_1, \dots, x_n\}$$

$$Y = \{x_1 + (S + E_1), \dots, x_n + (S + E_n)\}$$

2774 geni 27 condizioni sclerodermia

3051 geni 38 condizioni linfoma

	Node Deletion	Adjusted Node Deletion	Simulated Annealing	Node Deletion	Adjusted Node Deletion	Simulated Annealing
Yeast	First Bicluster			Second Bicluster		
$\delta = 300$	15165	15750	16460	<i>9012</i>	3930	8320
200	8463	9540	10360	<i>4972</i>	2630	3860
100	2520	2700	2940	<i>1260</i>	830	1390
Scleroderma						
$\delta = 300$	13590	18260	18230	<i>4320</i>	6780	6310
200	7296	12920	13210	<i>7876</i>	3290	4030
100	2730	5170	5140	<i>1570</i>	830	850
Lymphoma						
$\delta = 300$	1344	3320	3220	<i>518</i>	1740	1810
200	1032	2510	2460	<i>300</i>	1370	1200
100	851	1780	1790	<i>136</i>	1050	810

Cercare con SAB biclusters che abbiano un autentico significato biologico.

550 geni su 2884 sono annotati nella banca dati KEGG.

Bic.	# Genes	Dominant FM	Genes in FM	P-value
1	81	Ribosomal Proteins(96)	61	2.15×10^{-38}
		Glycolysis/Gluco-genesis(26)	5	0.16
2	59	Basal Transcription Factors(10)	6	1.6×10^{-4}
		Nucleotide Metabolism(81)	16	4.02×10^{-3}

P-value = probabilità che (x geni/ totale dei geni) appartenga alla funzione y

Più è basso e più è significativo

Bicluster 1 contiene 75% di geni che codificano per proteine ribosomali

Bicluster 2 contiene 40% di geni che codificano per fattori di trascrizione.

Il SAB trova biclusters significativi ma, in questi non domina totalmente un gruppo funzionale

- **dati incompleti**
- **la ricerca top down trova biclusters che non rispecchiano la situazione *'in vivo'***

FUTURO: **SAB** → **bottom up**