

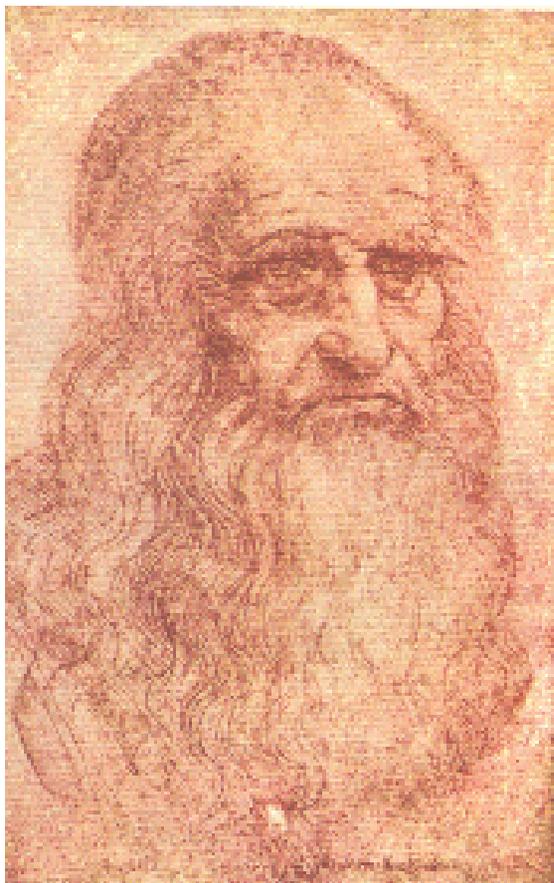
1. Similarità e omologia.
2. Allineamenti di sequenze.
3. Sostituzioni e gap.
4. Algoritmo di allineamento

Similarita' e omologia

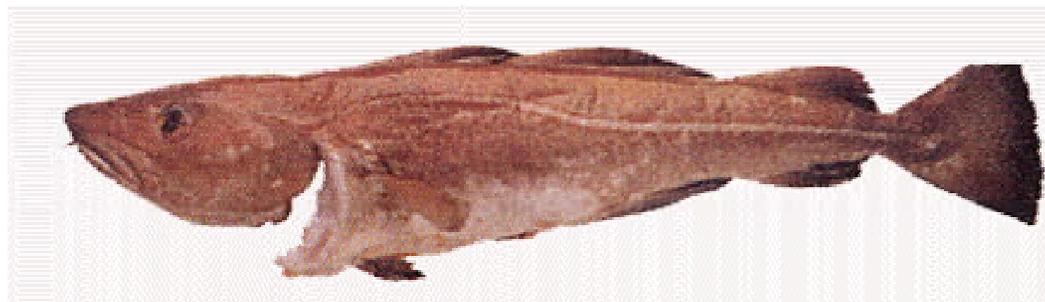
- Due sequenze sono **simili** se possono essere allineate in modo che molti ammino acidi corrispondenti sono identici o simili
- Tecnicamente due o piu' sequenze possono essere definite **omologhe** se derivano da un progenitore comune
- L'omologia tra due sequenze si deduce dalla loro similarità in sequenza o funzione.
- Il concetto di similarita' si puo' estendere anche alle strutture 3D

L'omologia tra due proteine/geni si deduce dalla loro similarita' in sequenza (struttura) o funzione.

Specie = uomo



Specie = merluzzo



→ La differenza
e' evidente!

Allineamento di sequenze

Qual è la corrispondenza tra gli amminoacidi (nucleotidi) che più probabilmente rispecchia l'evoluzione delle due proteine (geni)?

Allineamento di sequenza

Scopo: minimizzare la distanza evolutiva tra le sequenze da allineare, quindi minimizzare le differenze (che equivale a massimizzare le similarità) tra le componenti (nucleotidi o amminoacidi) delle sequenze stesse

L'*ipotesi* di allineamento *piu' ragionevole* e' quella che coinvolge il numero minore di eventi di mutazione per passare da una sequenza all'altra

allineamenti di sequenze

Applicazioni

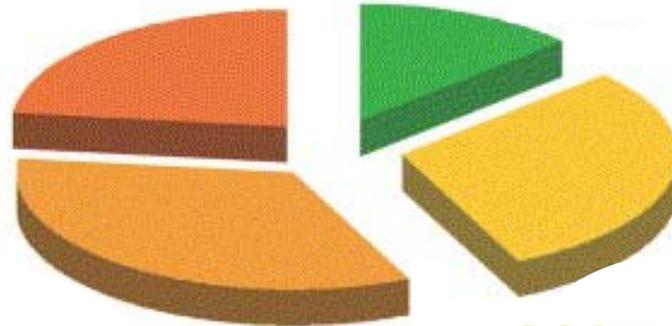
- Riconoscimento di funzione
 - basta stabilire che due sequenze sono simili
- Filogenia
 - occorre misurare la similarita' quantitativamente
- Model building
 - occorre costruire il miglior allineamento possibile

La distribuzione delle sequenze

Delle sequenze proteiche note ad oggi
meno del 15% hanno una caratterizzazione
funzionale sperimentale soddisfacente
Il 25% delle proteine non ha omologhe note

**sequenze con
nessuna omologa**

**sequenze
annotate**

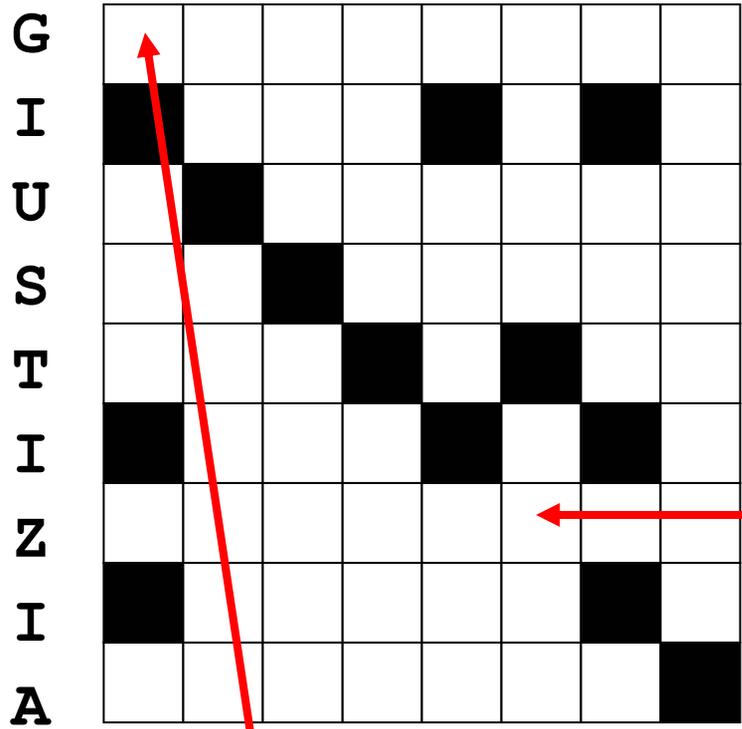


**sequenze con
omologhe distanti**

**sequenze con
omologhe vicine**

Il **dotplot** (disegno a punti) fornisce una visione immediata della similarita' tra due stringhe di caratteri

I U S T I T I A



latino

- IUSTITIA

italiano

GIUST**Z**IA

inserzione

mutazione

R E S P U B L I C A

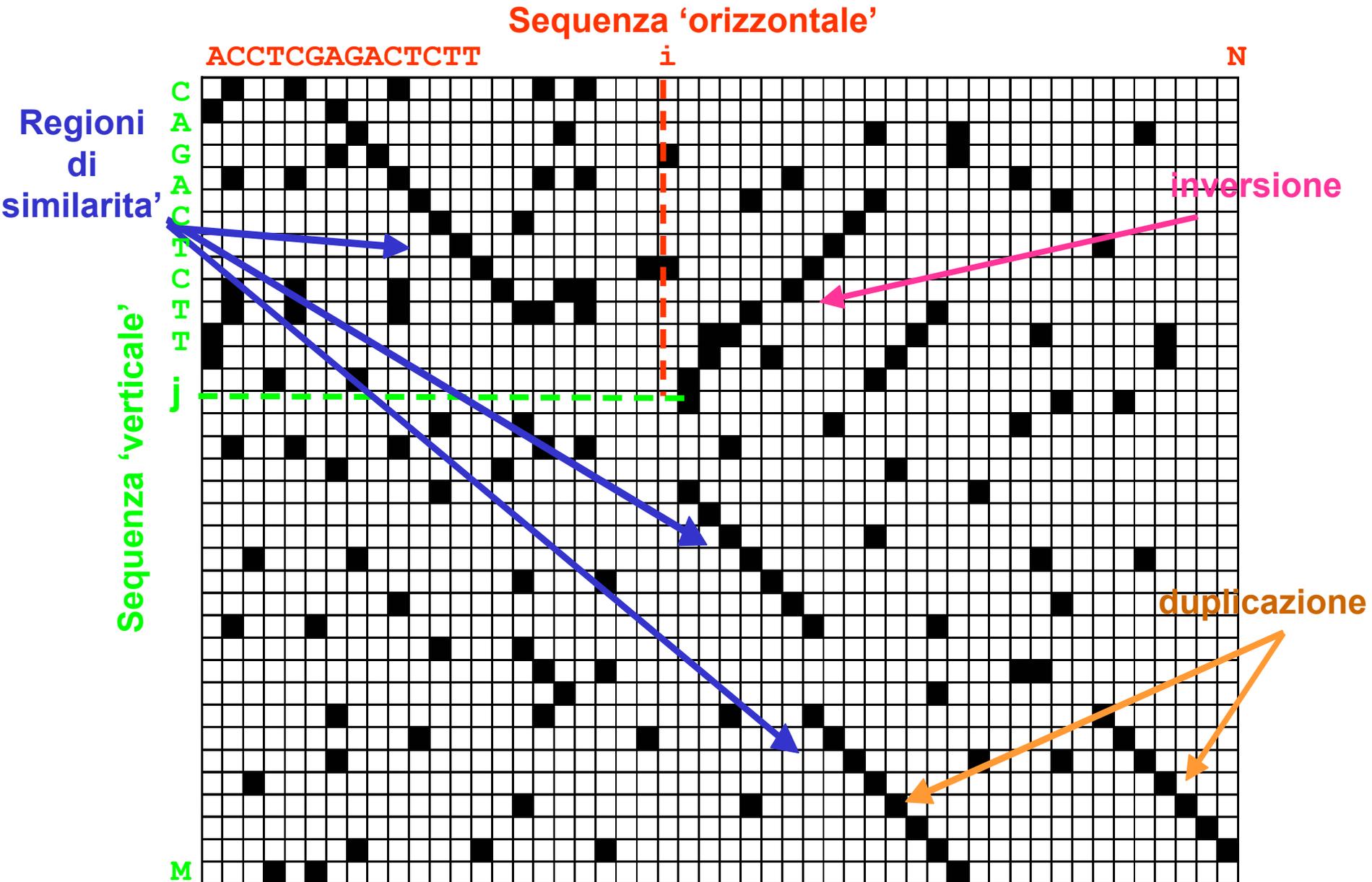
R	■								
E		■							
P			■						
U				■					
B					■				
L						■			
I							■		
C								■	

Latino **RESPUBLICA**

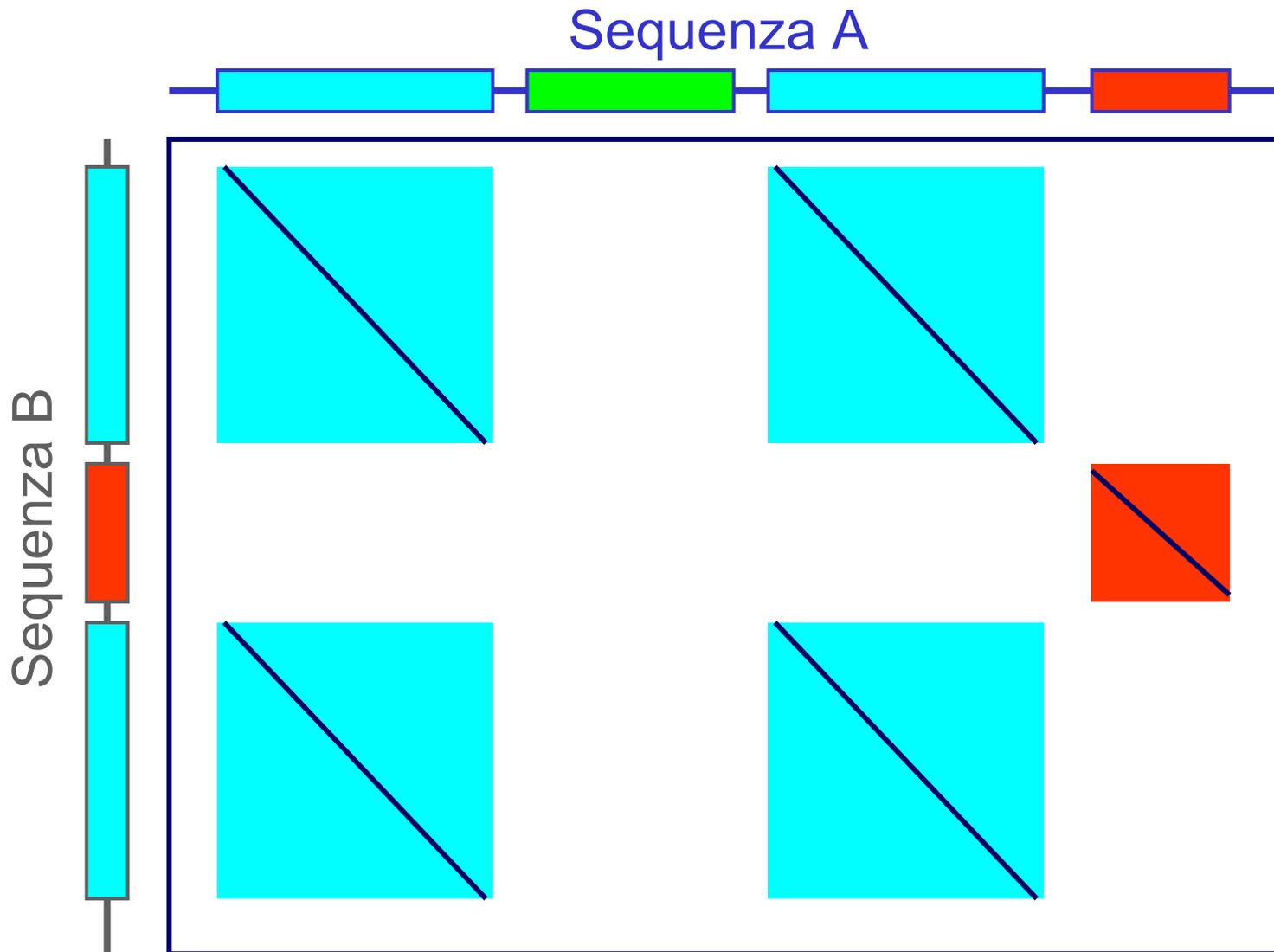
Inglese **RE-PUBLIC-**

delezioni

Dot plot



Dotplot di proteine domain-shuffled



Sequenza 1 (19 caratteri)

LAMIAPRIMASEQCREATA

Sequenza 2 (22 caratteri)

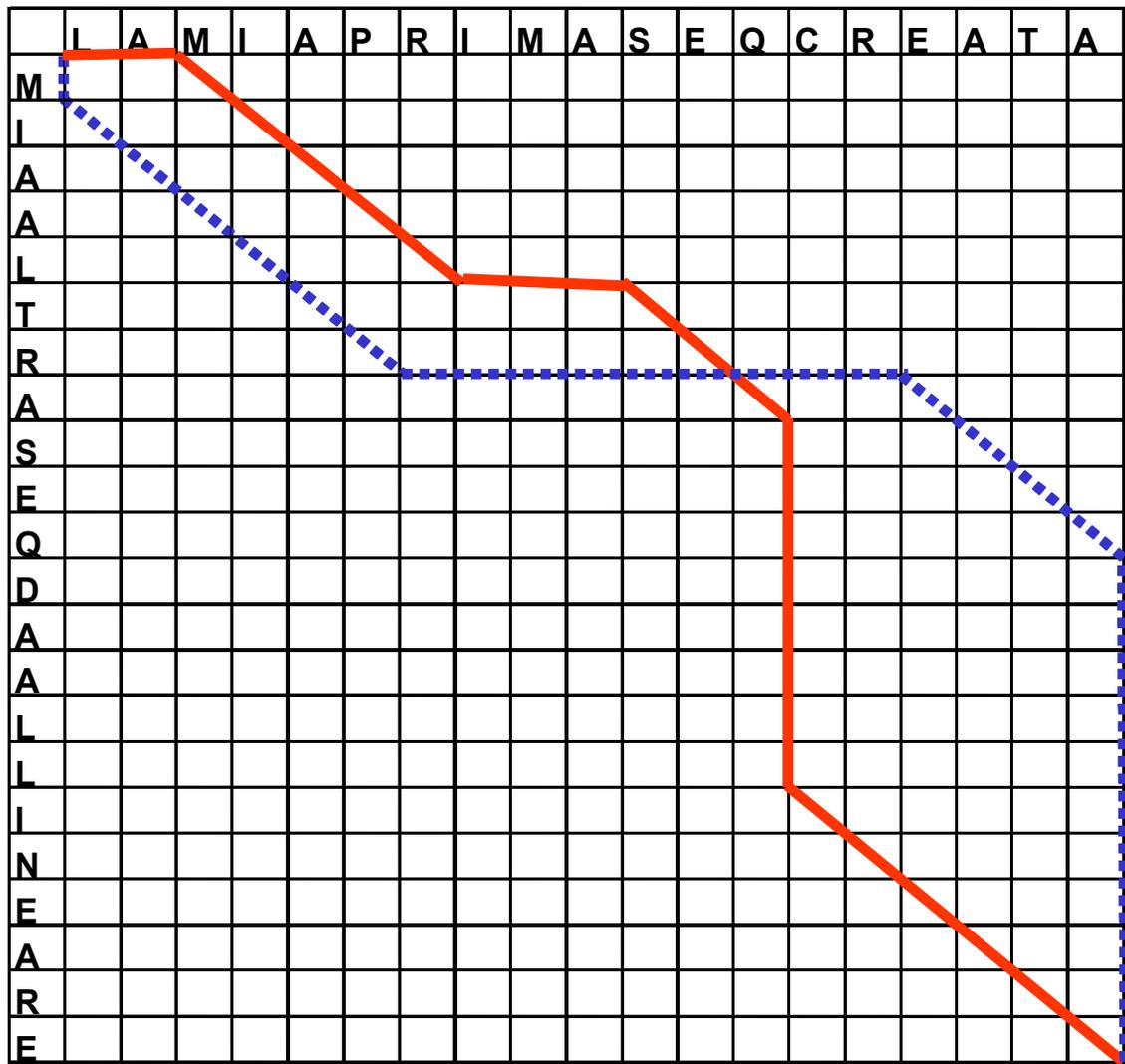
MIAALTRASEQDALLINEARE

	L	A	M	I	A	P	R	I	M	A	S	E	Q	C	R	E	A	T	A
M	0	0	1	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
I	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
A	0	1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	1	0	1
A	0	1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	1	0	1
L	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
T	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0
R	0	0	0	0	0	0	1	0	0	0	0	0	0	0	1	0	0	0	0
A	0	1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	1	0	1
S	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0
E	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	1	0	0	0
Q	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
D	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	0	1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	1	0	1
A	0	1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	1	0	1
L	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
L	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
I	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
N	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
E	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	1	0	0	0
A	0	1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	1	0	1
R	0	0	0	0	0	0	1	0	0	0	0	0	0	0	1	0	0	0	0
E	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0

	L	A	M	I	A	P	R	I	M	A	S	E	Q	C	R	E	A	T	A
M	0	0	1	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
I	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
A	0	1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	1	0	1
A	0	1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	1	0	1
L	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
T	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0
R	0	0	0	0	0	0	1	0	0	0	0	0	0	0	1	0	0	0	0
A	0	1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	1	0	1
S	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0
E	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	1	0	0	0
Q	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
D	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	0	1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	1	0	1
A	0	1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	1	0	1
L	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
L	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
I	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
N	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
E	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	1	0	0	0
A	0	1	0	0	1	0	0	0	0	1	0	0	0	0	0	0	1	0	1
R	0	0	0	0	0	0	1	0	0	0	0	0	0	0	1	0	0	0	0
E	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0

10 punti

LAMIAP---RIMASEQ-----CREATA
 --MIA-ALTR--ASEQDAALLIN--EARE



LAMIAPRIMASEQCREATA-----

-MIAALTR-----ASEQDAALLINEARE

LAMIAPRIMASEQ-----CREATA

--MIAAL---TRASEQDAALLINEARE

Allineamento

- 1) Punteggio per la corrispondenza di aa/basi
- 2) Penalizzazione di inserzioni / delezioni
- 3) Algoritmo che effettui l'allineamento
- 4) Misura della significativita' dell'allineamento

Allineamento

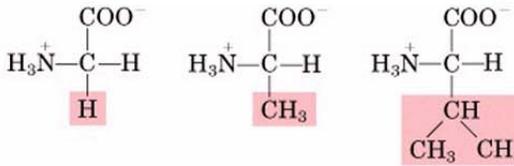
- 1) Punteggio per la corrispondenza di aa/basi
- 2) Penalizzazione di inserzioni / delezioni
- 3) Algoritmo che effettui l'allineamento
- 4) Misura della significativita' dell'allineamento

1. Punteggio (Score)

- Identita' tra basi o amminoacidi
- Proprieta' chimico-fisiche degli amminoacidi
- Numero minimo di basi che e' necessario mutare per ottenere la mutazione osservata
- Frequenze di rimpiazzo osservate in famiglie di proteine

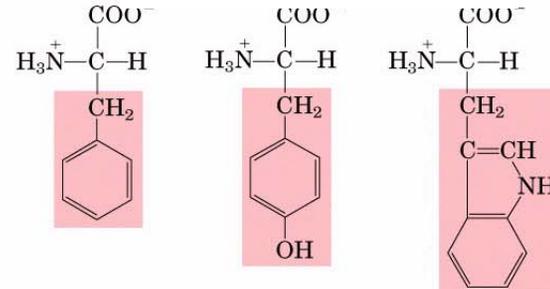
Twenty standard Amino Acids

non polari

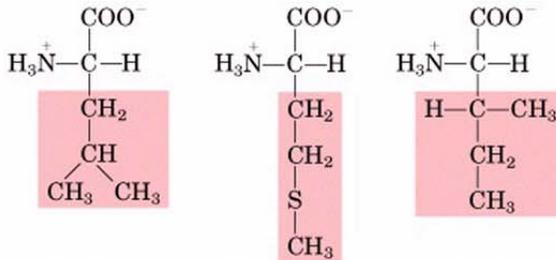


Glicina(**G**) Alanina(**A**) Valina(**V**)

aromatici

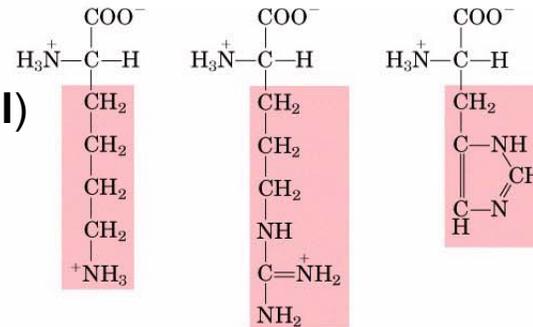


Phenilalanina(**F**) Tiroisina(**Y**) Triptofano(**W**)



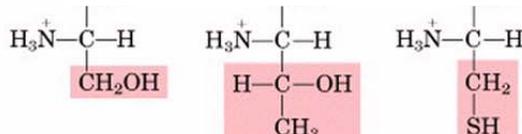
Leucina(**L**) Metionina(**M**) Isoleucina(**I**)

basici



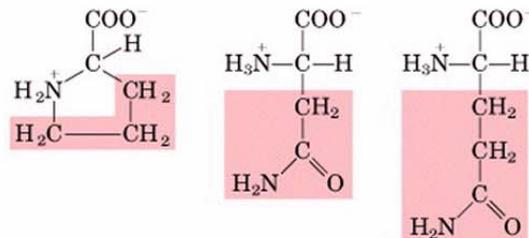
Lisina(**K**) Arginina(**R**) Histidina(**H**)

polari, non carichi

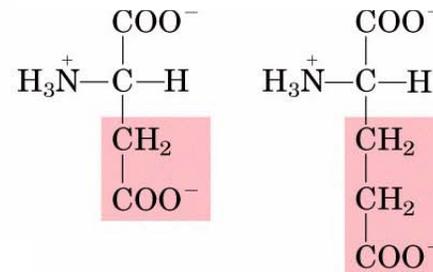


Serina(**S**) Treonina(**T**) Cisteina(**C**)

acidi



Prolina(**P**) Asparagina(**N**) Glutammina(**Q**)



Aspartato(**D**) Glutammato(**E**)

es. 1 Matrice di Sostituzione

	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y
A	2																			
C	-2	12																		
D	0	-5	4																	
E	0	-5	3	4																
F	-4	-4	-6	-5	9															
G	1	-3	1	0	-5	5														
H	-1	-3	1	1	-2	-2	6													
I	-1	-2	-2	-2	1	-3	-2	5												
K	-1	-5	0	0	-5	-2	0	-2	5											
L	-2	-6	-4	-3	2	-4	-2	2	-3	6										
M	-1	-5	-3	-2	0	-3	-2	2	0	4	6									
N	0	-4	2	1	-4	0	2	-2	1	-3	-2	2								
P	1	-3	-1	-1	-5	-1	0	-2	-1	-3	-2	-1	6							
Q	0	-5	2	2	-5	-1	3	-2	1	-2	-1	1	0	4						
R	-2	-4	-1	-1	-4	-3	2	-2	3	-3	0	0	0	1	6					
S	1	0	0	0	-3	1	-1	-1	0	-3	-2	1	1	-1	0	2				
T	1	-2	0	0	-3	0	-1	0	0	-2	-1	0	0	-1	-1	1	3			
V	0	-2	-2	-2	-1	-1	-2	4	-2	2	2	-2	-1	-2	-2	-1	0	4		
W	-6	-8	-7	-7	0	-7	-3	-5	-3	-2	-4	-4	-6	-5	2	-2	-5	-6	17	
Y	-3	0	-4	-4	7	-5	0	-1	-4	-1	-2	-2	-5	-4	-4	-3	-3	-2	0	10

PAM 250

Matrici 20x20

Point Accepted Mutations

es. 2 Matrice di Sostituzione

	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y
A	4																			
C	0	9	-3																	
D	-2	-3	6																	
E	-1	-4	2	5																
F	-2	-2	-3	-3	6															
G	0	-3	-1	-2	-3	6														
H	-2	-3	-1	0	-1	-2	8													
I	-1	-1	-3	-3	0	-4	-3	4												
K	-1	-3	-1	1	-3	-2	-1	-3	5											
L	-1	-1	-4	-3	0	-4	-3	2	-2	4										
M	-1	-1	-3	-2	0	-3	-2	1	-1	2	5									
N	-2	-3	1	0	-3	0	1	-3	0	-3	-2	6								
P	-1	-3	-1	-1	-4	-2	-2	-3	-1	-3	-2	-2	7							
Q	-1	-3	0	2	-3	-2	0	-3	1	-2	0	0	-1	5						
R	-1	-3	-2	0	-3	-2	0	-3	2	-2	-1	0	-2	1	5					
S	1	-1	0	0	-2	0	-1	-2	0	-2	-1	1	-1	0	-1	4				
T	0	-1	-1	-1	-2	-2	-2	-1	-1	-1	-1	0	-1	-1	-1	1	5			
V	0	-1	-3	-2	-1	-3	-3	3	-2	1	1	-3	-2	-2	-3	-2	0	4		
W	-3	-2	-4	-3	1	-2	-2	-3	-3	-2	-1	-4	-4	-2	-3	-3	-2	-3	11	
Y	-2	-2	-3	-2	3	-3	2	-1	-2	-1	-1	-2	-3	-1	-2	-2	-2	-1	2	7

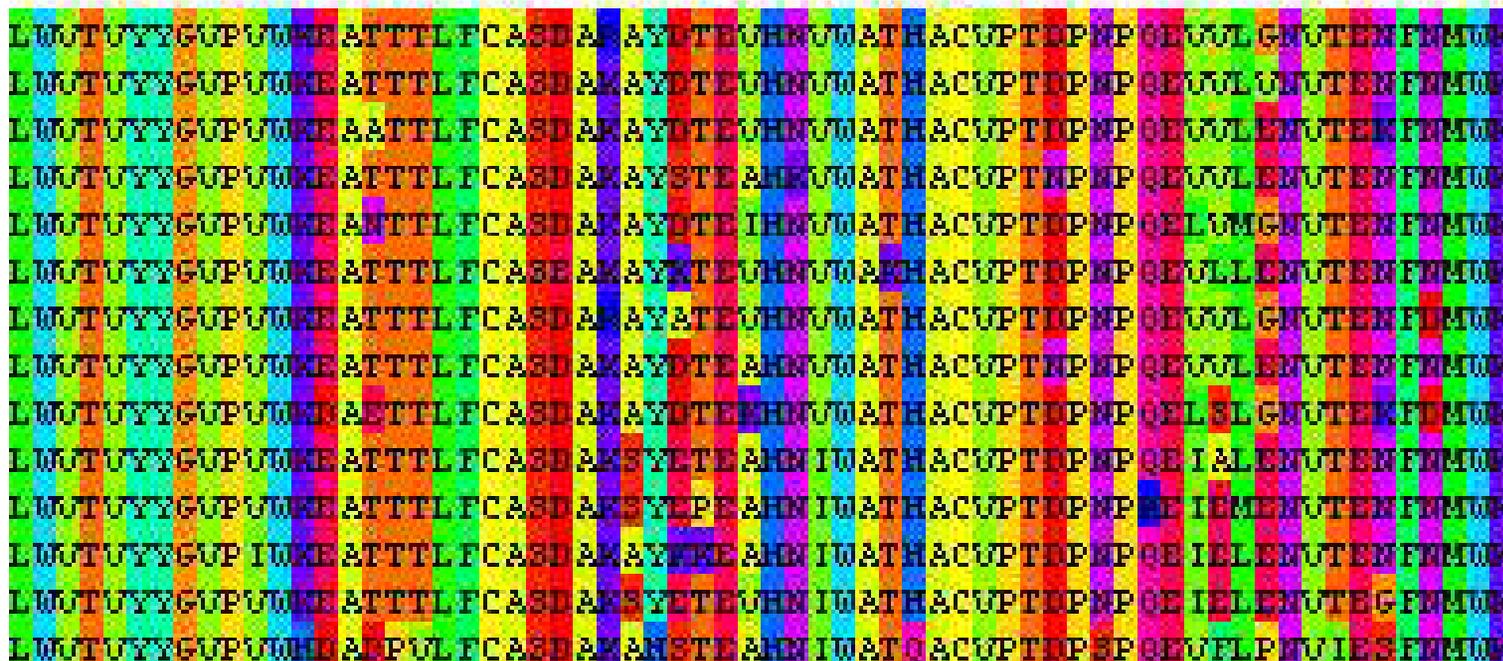
BLOSUM 62

Blocks Substitution Matrix

(dalla banca dati derivata BLOCKS)

Matrici 20x20

Matrici sito specifiche, basate su regole empiriche



```
LWUTUYYGUPUWKE ATTLFCASDARAYDTEUHNWATHACUPTIDPMPGEUOLGNUTENFNMWQ
LWUTUYYGUPUWKE ATTLFCASDARAYDTEUHNWATHACUPTIDPMPGEUOLGNUTENFNMWQ
LWUTUYYGUPUWKE AATTLFCASDARAYDTEUHNWATHACUPTIDPMPGEUOLENUTENFNMWQ
LWUTUYYGUPUWKE ATTLFCASDARAYSTEAMENWATHACUPTIDPMPGEUOLENUTENFNMWQ
LWUTUYYGUPUWKE ANTTLFCASDARAYDTEIHNWATHACUPTIDPMPGELUWGNUTENFNMWQ
LWUTUYYGUPUWKE ATTLFCASEAMAYKTEUHNWATHACUPTIDPMPGEVLLLENUTENFNMWQ
LWUTUYYGUPUWKE ATTLFCASDARAYSTEUHNWATHACUPTIDPMPGEUOLGNUTENFNMWQ
LWUTUYYGUPUWKE ATTLFCASDARAYDTEAMNWATHACUPTIDPMPGEUOLENUTENFNMWQ
LWUTUYYGUPUWKE AATTLFCASDARAYDTEHWNWATHACUPTIDPMPGELSLGNUTENFNMWQ
LWUTUYYGUPUWKE ATTLFCASDARASYETEAMNIWATHACUPTIDPMPQEIALENUTENFNMWQ
LWUTUYYGUPUWKE ATTLFCASDARASYEPEAMNIWATHACUPTIDPMPQEIEENUTENFNMWQ
LWUTUYYGUPUWKE ATTLFCASDARASYETEAMNIWATHACUPTIDPMPQEIELENUTENFNMWQ
LWUTUYYGUPUWKE ATTLFCASDARANSTEAMNIWATHACUPTIDPEPGEUFLPNVIESFNMWQ
```

PAM N : Percent/Point Accepted Mutations (dove N è il numero di mutazioni accettate per 100 aa)

BLOSUM N : BLOcks SUBstitution Matrix (dove N è la % max di identità di sequenza tra le omologhe allineate)

Le matrici PAM e BLOSUM riportano i \log_2 di:

$$\frac{f_{ij}}{f_i \times f_j}$$

dove

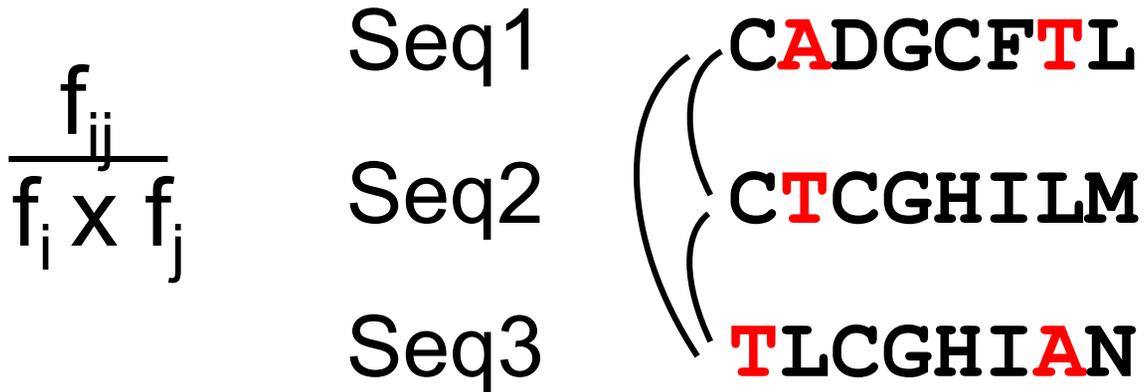
f_{ij} rappresenta la frequenza con cui troviamo due amminoacidi allineati tra loro

es. f_{AT} è il numero di volte in cui osserviamo una Ala allineata con una Thr diviso il numero totale di coppie in un allineamento

f_i ed f_j rappresentano le frequenze con cui compaiono due amminoacidi in un allineamento

es. f_A ed f_T sono il numero di volte che abbiamo una Ala o una Thr diviso il numero totale di amminoacidi in un allineamento

Esempio di calcolo di punteggio per la sostituzione di una Ala con una Thr



$$\text{Tot aa} = 3 \times 8 = 24$$

$$\text{Tot coppie aa allineati} = 3 \times 8 = 24$$

$$f_{AT} = 2 / 24 = 0.083$$

$$f_A = 2 / 24 = 0.083$$

$$f_T = 3 / 24 = 0.12$$

$$f_A \times f_T = 0.010$$

$$\frac{f_{AT}}{f_A \times f_T} = \frac{0.083}{0.010} = 8$$

$$\ln_2 \left[\frac{f_{AT}}{f_A \times f_T} \right] = 3$$

Esempio di calcolo di punteggio per la sostituzione di una Ala con una Thr

$$\frac{f_{ij}}{f_i \times f_j}$$

Seq1 CADGCF**TL**
Seq2 C**T**CGHILM
Seq3 **T**LCGHIA**N**

$$\text{Tot aa} = 3 \times 8 = 24$$

$$\text{Tot coppie aa allineati} = 3 \times 8 = 24$$

Nelle 3 sequenze considerate l'allineamento **A-T** occorre 8 volte più che per caso

$$\frac{f_{AT}}{f_A \times f_T} = \frac{0.083}{0.010} = 8$$

$$\ln_2 \left(\frac{f_{AT}}{f_A \times f_T} \right) = 3$$

Esempio di calcolo di punteggio per la sostituzione di una Ala con una Thr

$$\frac{f_{ij}}{f_i \times f_j}$$

Seq1 C**A**DGCF**T**L
Seq2 C**T**CGHILM
Seq3 **T**LCGHIA**N**

$$\text{Tot aa} = 3 \times 8 = 24$$

$$\text{Tot coppie aa allineati} = 3 \times 8 = 24$$

In una matrice di sostituzione scriveremmo **3** all'incrocio tra Ala(**A**) e Thr(**T**)

$$\frac{f_{AT}}{f_A \times f_T} = \frac{0.083}{0.010} = 8$$

$$\ln_2 \left[\frac{f_{AT}}{f_A \times f_T} \right] = 3$$

Matrici sito specifiche, basate su regole empiriche

```
LWUTUYYGUPVUWE ATTLFCASD AAYDTEUHNWATHACUPTDPPNPGEUOLGNUTENFNMW
LWUTUYYGUPVUWE ATTLFCASD AAYDTEUHNWATHACUPTDPPNPGEUOLGNUTENFNMW
LWUTUYYGUPVUWE AATTLFCASD AAYDTEUHNWATHACUPTDPPNPGEUOLENUTENFNMW
LWUTUYYGUPVUWE ATTLFCASD AAYSTEAMHWATHACUPTDPPNPGEUOLENUTENFNMW
LWUTUYYGUPVUWE ANTTLFCASD AAYDTEIHNWATHACUPTDPPNPGEUOLGNUTENFNMW
LWUTUYYGUPVUWE ATTLFCASE AAYKTEUHNWATHACUPTDPPNPGEULLENUTENFNMW
LWUTUYYGUPVUWE ATTLFCASD AAYATEUHNWATHACUPTDPPNPGEUOLGNUTENFNMW
LWUTUYYGUPVUWE ATTLFCASD AAYDTEAMHWATHACUPTDPPNPGEUOLENUTENFNMW
LWUTUYYGUPVUWE AATTLFCASD AAYDTEIHNWATHACUPTDPPNPGEULSLGNUTENFNMW
LWUTUYYGUPVUWE ATTLFCASD AASYETEAMNIWATHACUPTDPPNPGEIALENUTENFNMW
LWUTUYYGUPVUWE ATTLFCASD AASYEPEAMNIWATHACUPTDPPNPGEIEMENUTENFNMW
LWUTUYYGUPVUWE ATTLFCASD AAYEKEAMNIWATHACUPTDPPNPGEIELENUTENFNMW
LWUTUYYGUPVUWE ATTLFCASD AASYETEAMNIWATHACUPTDPPNPGEIELENUTEGFNMW
LWUTUYYGUPVUWE AALPULFCASD AANSTEAMNIWATHACUPTDPPNPGEUFLPNVIESFNMW
```

PAM N : Percent/Point Accepted Mutations (dove N è il numero di mutazioni accettate per 100 aa)

BLOSUM N : BLOcks SUBstitution Matrix (dove N è la % max di identità di sequenza tra le omologhe allineate)

- **PAM 1** (ottenuta da allineamento di proteine con sequenze differenti per 1 aa su 100) può essere usata per generare matrici per distanze evoluzionistiche maggiori:

moltiplicandola ripetutamente per se stessa.

$$\mathbf{PAM2} = \mathbf{PAM1} * \mathbf{PAM1}$$

etc etc

- **PAM250:**
 - 2,5 mutazioni per residuo
 - equivalente al 20% matches rimanenti tra due sequenze, cioè l'80% delle posizioni amminoacidiche sono cambiate.
 - È la matrice usata di default in molti programmi di analisi.

- Le matrici **BLOSUM** sono state sviluppate per allineare sequenze scarsamente correlate. Hanno largamente soppiantato le PAM.
- Sono ottenute dalla banca dati derivata **BLOCKS** contenente allineamenti di regioni di proteine strettamente correlate, allineabili senza gap.
- **BLOSUM62**: ottenuta da allineamenti di proteine con un massimo di 62 % di identità di sequenza. Molto usata. (Corrisponde approssimativamente ad una PAM110).

es. 1 Matrice di Sostituzione

	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y
A	2																			
C	-2	12																		
D	0	-5	4																	
E	0	-5	3	4																
F	-4	-4	-6	-5	9															
G	1	-3	1	0	-5	5														
H	-1	-3	1	1	-2	-2	6													
I	-1	-2	-2	-2	1	-3	-2	5												
K	-1	-5	0	0	-5	-2	0	-2	5											
L	-2	-6	-4	-3	2	-4	-2	2	-3	6										
M	-1	-5	-3	-2	0	-3	-2	2	0	4	6									
N	0	-4	2	1	-4	0	2	-2	1	-3	-2	2								
P	1	-3	-1	-1	-5	-1	0	-2	-1	-3	-2	-1	6							
Q	0	-5	2	2	-5	-1	3	-2	1	-2	-1	1	0	4						
R	-2	-4	-1	-1	-4	-3	2	-2	3	-3	0	0	0	1	6					
S	1	0	0	0	-3	1	-1	-1	0	-3	-2	1	1	-1	0	2				
T	1	-2	0	0	-3	0	-1	0	0	-2	-1	0	0	-1	-1	1	3			
V	0	-2	-2	-2	-1	-1	-2	4	-2	2	2	-2	-1	-2	-2	-1	0	4		
W	-6	-8	-7	-7	0	-7	-3	-5	-3	-2	-4	-4	-6	-5	2	-2	-5	-6	17	
Y	-3	0	-4	-4	7	-5	0	-1	-4	-1	-2	-2	-5	-4	-4	-3	-3	-2	0	10

PAM 250

Matrici 20x20

Point Accepted Mutations

es. 1 Matrice di Sostituzione

Ala



	^	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y
A	2																			
C	~	12																		
D	0	-5	4																	
E	0	-5	3	4																
F	-4	-4	-6	-5	9															
G	1	-3	1	0	-5	5														
H	-1	-3	1	1	-2	-2	6													
I	-1	-2	-2	-2	1	-3	-2	5												
K	-1	-5	0	0	-5	-2	0	-2	5											
L	-2	-6	-4	-3	2	-4	-2	2	-3	6										
M	-1	-5	-3	-2	0	-3	-2	2	0	4	6									
N	0	-4	2	1	-4	0	2	-2	1	-3	-2	2								
P	1	-3	-1	-1	-5	-1	0	-2	-1	-3	-2	-1	6							
Q	0	-5	2	2	-5	-1	3	-2	1	-2	-1	1	0	4						
R	-2	-4	-1	-1	-4	-3	2	-2	3	-3	0	0	0	1	6					
S	1	0	0	0	-3	1	-1	-1	0	-3	-2	1	1	-1	0	2				
T	1	-2	0	0	-3	0	-1	0	0	-2	-1	0	0	-1	-1	1	3			
V	0	-2	-2	-2	-1	-1	-2	4	-2	2	2	-2	-1	-2	-2	-1	0	4		
W	-6	-8	-7	-7	0	-7	-3	-5	-3	-2	-4	-4	-6	-5	2	-2	-5	-6	17	
Y	-3	0	-4	-4	7	-5	0	-1	-4	-1	-2	-2	-5	-4	-4	-3	-3	-2	0	10

Le Ala sono facilmente rimpiazzabili

es. 1 Matrice di Sostituzione

Cys



	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y
A	2																			
C	-2	12																		
D	0	-3	4																	
E	0	-5	3	4																
F	-4	-4	-6	-5	9															
G	1	-3	1	0	-5	5														
H	-1	-3	1	1	-2	-2	6													
I	-1	-2	-2	-2	1	-3	-2	5												
K	-1	-5	0	0	-5	-2	0	-2	5											
L	-2	-6	-4	-3	2	-4	-2	2	-3	6										
M	-1	-5	-3	-2	0	-3	-2	2	0	4	6									
N	0	-4	2	1	-4	0	2	-2	1	-3	-2	2								
P	1	-3	-1	-1	-5	-1	0	-2	-1	-3	-2	-1	6							
Q	0	-5	2	2	-5	-1	3	-2	1	-2	-1	1	0	4						
R	-2	-4	-1	-1	-4	-3	2	-2	3	-3	0	0	0	1	6					
S	1	0	0	0	-3	1	-1	-1	0	-3	-2	1	1	-1	0	2				
T	1	-2	0	0	-3	0	-1	0	0	-2	-1	0	0	-1	-1	1	3			
V	0	-2	-2	-2	-1	-1	-2	4	-2	2	2	-2	-1	-2	-2	-1	0	4		
W	-6	-8	-7	-7	0	-7	-3	-5	-3	-2	-4	-4	-6	-5	2	-2	-5	-6	17	
Y	-3	0	-4	-4	7	-5	0	-1	-4	-1	-2	-2	-5	-4	-4	-3	-3	-2	0	10

La Cys non è facilmente rimpiazzabile

es. 1 Matrice di Sostituzione

	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y
A	2																			
C	-2	12																		
D	0	-5	4																	
E	0	-5	3	4																
F	-4	-4	-6	-5	9															
G	1	-3	1	0	-5	5														
H	-1	-3	1	1	-2	-2	6													
I	-1	-2	-2	-2	1	-3	-2	5												
K	-1	-5	0	0	-5	-2	0	-2	5											
L	-2	-6	-4	-3	2	-4	-2	2	3	6										
M	-1	-5	-3	-2	0	-3	-2	2	0	4	6									
N	0	-4	2	1	-4	0	2	-2	1	-3	-2	2								
P	1	-3	-1	-1	-5	-1	0	-2	1	-3	-2	-1	6							
Q	0	-5	2	2	-5	-1	3	-2	1	-2	-1	1	0	4						
R	-2	-4	-1	-1	-4	-3	-2	3	3	-3	0	0	0	1	6					
S	1	0	0	0	-3	1	-1	-1	0	-3	-2	1	1	-1	0	2				
T	1	-2	0	0	-3	0	-1	0	0	-2	-1	0	0	-1	-1	1	3			
V	0	-2	-2	-2	-1	-1	-2	4	-2	2	2	-2	-1	-2	-2	-1	0	4		
W	-6	-8	-7	-7	0	-7	-3	-5	-3	-2	-4	-4	-6	-5	2	-2	-5	-6	17	
Y	-3	0	-4	-4	7	-5	0	-1	-4	-1	-2	-2	-5	-4	-4	-3	-3	-2	0	10

Arg e Lys tendono a sostituirsi

es. 1 Matrice di Sostituzione

	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y
A	2																			
C	-2	12																		
D	0	-5	4																	
E	0	-5	3	4																
F	-4	-4	-6	-5	9															
G	1	-3	1	0	-5	5														
H	-1	-3	1	1	-2	-2	6													
I	-1	-2	-2	-2	1	-3	-2	5												
K	-1	-5	0	0	-5	-2	0	-2	5											
L	-2	-6	-4	-3	2	-4	-2	2	3	6										
M	-1	-5	-3	-2	0	-3	-2	2	0	4	6									
N	0	-4	2	1	-4	0	2	-2	-1	-3	-2	2								
P	1	-3	-1	-1	-5	-1	0	-2	-1	-3	-2	-1	6							
Q	0	-5	2	2	-5	-1	3	-2	-1	-2	-1	1	0	4						
R	-2	-4	-1	-1	-4	-3	2	-2	3	-3	0	0	0	1	6					
S	1	0	0	0	-3	1	-1	-1	2	-3	-2	1	1	-1	0	2				
T	1	-2	0	0	-3	0	-1	0	4	-2	-1	0	0	-1	-1	1	3			
V	0	-2	-2	-2	-1	-1	-2	-2	-2	2	2	-2	-1	-2	-2	-1	0	4		
W	-6	-8	-7	-7	0	-7	-3	-5	3	-2	-4	-4	-6	-5	2	-2	-5	-6	17	
Y	-3	0	-4	-4	7	-5	0	-1	-4	-1	-2	-2	-5	-4	-4	-3	-3	-2	0	10

Polari e non polari tendono a non sostituirsi

Allineamento

- 1) Punteggio per la corrispondenza di aa/basi
- 2) Penalizzazione di inserzioni / delezioni**
- 3) Algoritmo che effettui l'allineamento
- 4) Misura della significativita' dell'allineamento

2. Penalizzazioni per inserzioni/delezioni

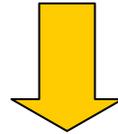
Le matrici di sostituzione sono state derivate da allineamenti che non presentavano inserzioni/delezioni. Queste vanno dunque trattate separatamente, su base empirica.

Nell'allineare due sequenze un algoritmo tenderebbe a massimizzare il punteggio (corrispondenza amminoacidi identici o simili) inserendo un gran numero di gap.

Questo modo di procedere rispecchia l'evoluzione?

2. Penalizzazioni per inserzioni/delezioni

In natura inserzioni e delezioni sono spesso letali



Dobbiamo quindi penalizzare inserzioni e delezioni. Cioè associare ad esse un punteggio negativo che si sottrae al punteggio totale per un allineamento

2. Penalizzazioni per inserzioni/delezioni

In natura la delezione di una serie di basi/amminoacidi contigui è un evento più probabile rispetto alla delezione indipendente dello stesso numero di basi/amminoacidi in posizioni non contigue

Distinguiamo l'**inizio** di un gap:

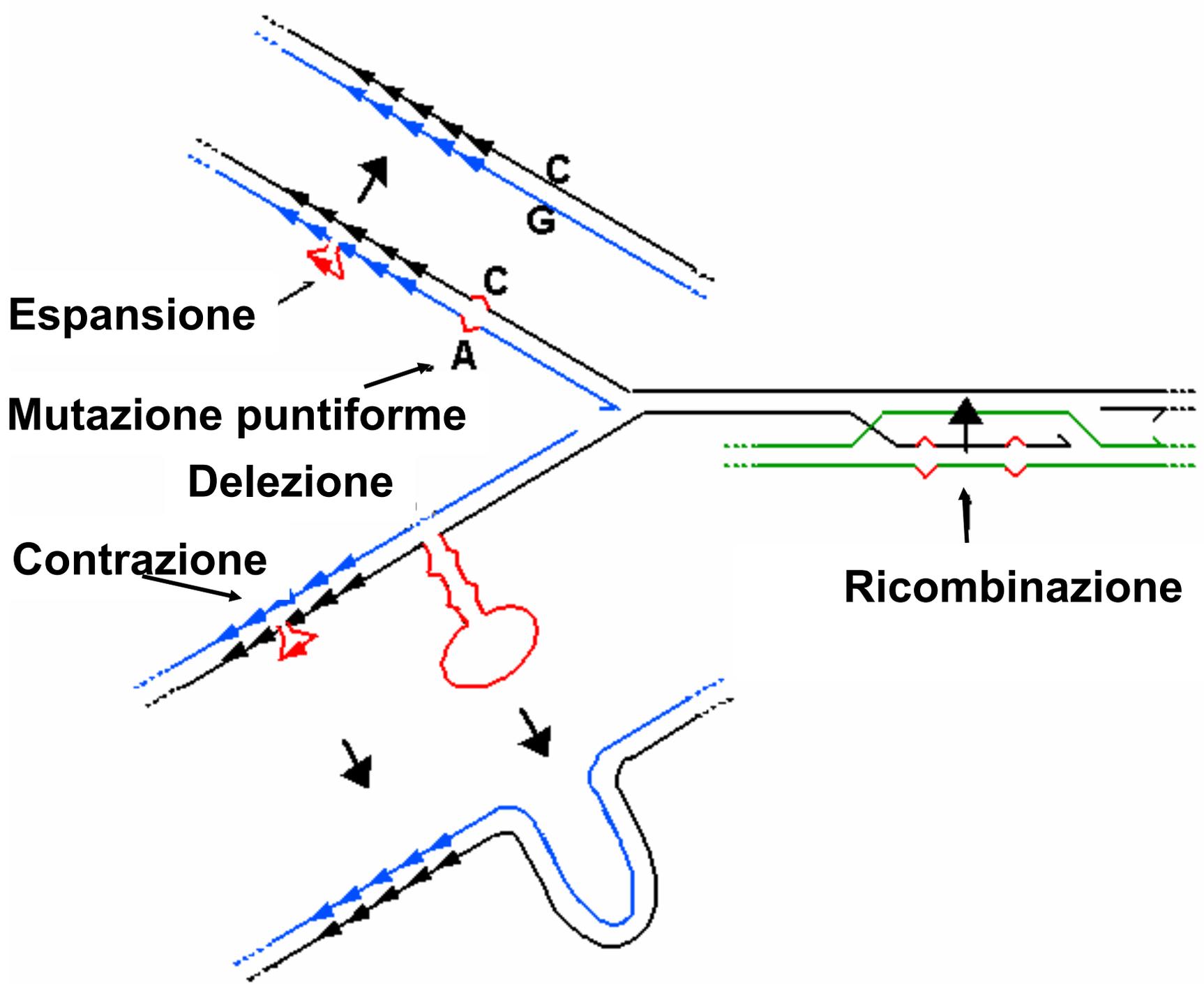
EGQTCA

AG-TCL

dall'**estensione** di un gap:

EGQQQTCA

AG---TCL



2. Penalizzazioni per inserzioni/delezioni

In natura la delezione di una serie di basi/amminoacidi contigui è un evento più probabile rispetto alla delezione indipendente dello stesso numero di basi/amminoacidi in posizioni non contigue

Distinguiamo l'**inizio** di un gap:

EGQTCA

AG-TCL

dall'**estensione** di un gap:

EGQQQTCA

AG---TCL

2. Penalizzazioni per inserzioni/delezioni

In natura la delezione di una serie di basi/amminoacidi contigui è un evento più probabile rispetto alla delezione indipendente dello stesso numero di basi/amminoacidi in posizioni non contigue

Distinguiamo l'**inizio** di un gap:

EGQTCA

AG-TCL

dall'**estensione** di un gap:

EGQQQTCA

AG---TCL

Penalizziamo maggiormente l'inizio di un gap rispetto alla sua estensione

es.:

-11 inizio

-1 estensione

Allineamento

- 1) Punteggio per la corrispondenza di aa/basi
- 2) Penalizzazione di inserzioni / delezioni
- 3) Algoritmo che effettui l'allineamento**
- 4) Misura della significativita' dell'allineamento

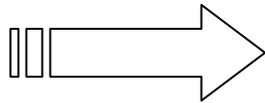
Qualche definizione...

dal nome di un matematico persiano al-Kharezmi (al-Khawarizmi) vissuto nel IX secolo

Algoritmo: Sequenza di istruzioni semplici (univocamente interpretabili) che permette di giungere ad un risultato in un numero determinato di passi.

Un esempio di algoritmo

1. *Prendere una padella antiaderente*
2. *versare un cucchiaino di olio extra vergine d'oliva*
3. *farlo scaldare per qualche secondo.*
4. *rompere l'uovo e adagiarlo al centro della padella*
5. *dare una spolverata di sale solo sul tuorlo*
6. *spegnere quando attorno all'albume si è formata una crosticina*



Qualche definizione...

dal nome di un matematico persiano al-Kharezmi (al-Khawarizmi) vissuto nel IX secolo

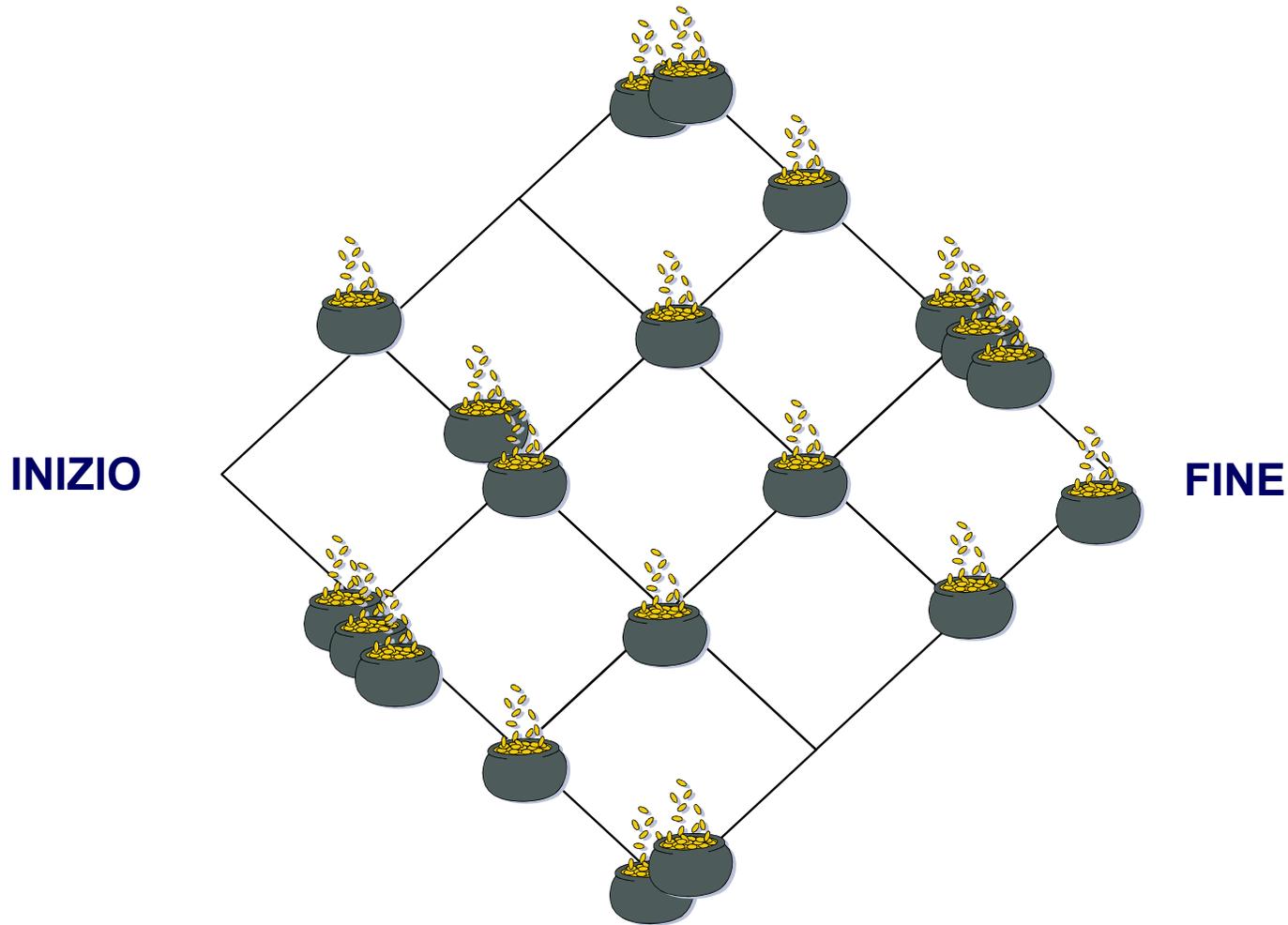
Algoritmo: Sequenza di istruzioni semplici (univocamente interpretabili) che permette di giungere ad un risultato in un numero determinato di passi.

Programma: descrizione di un algoritmo in uno specifico linguaggio di programmazione (C, fortran, perl).

Gli algoritmi di allineamenti esatti sono esempi di Programmazione Dinamica

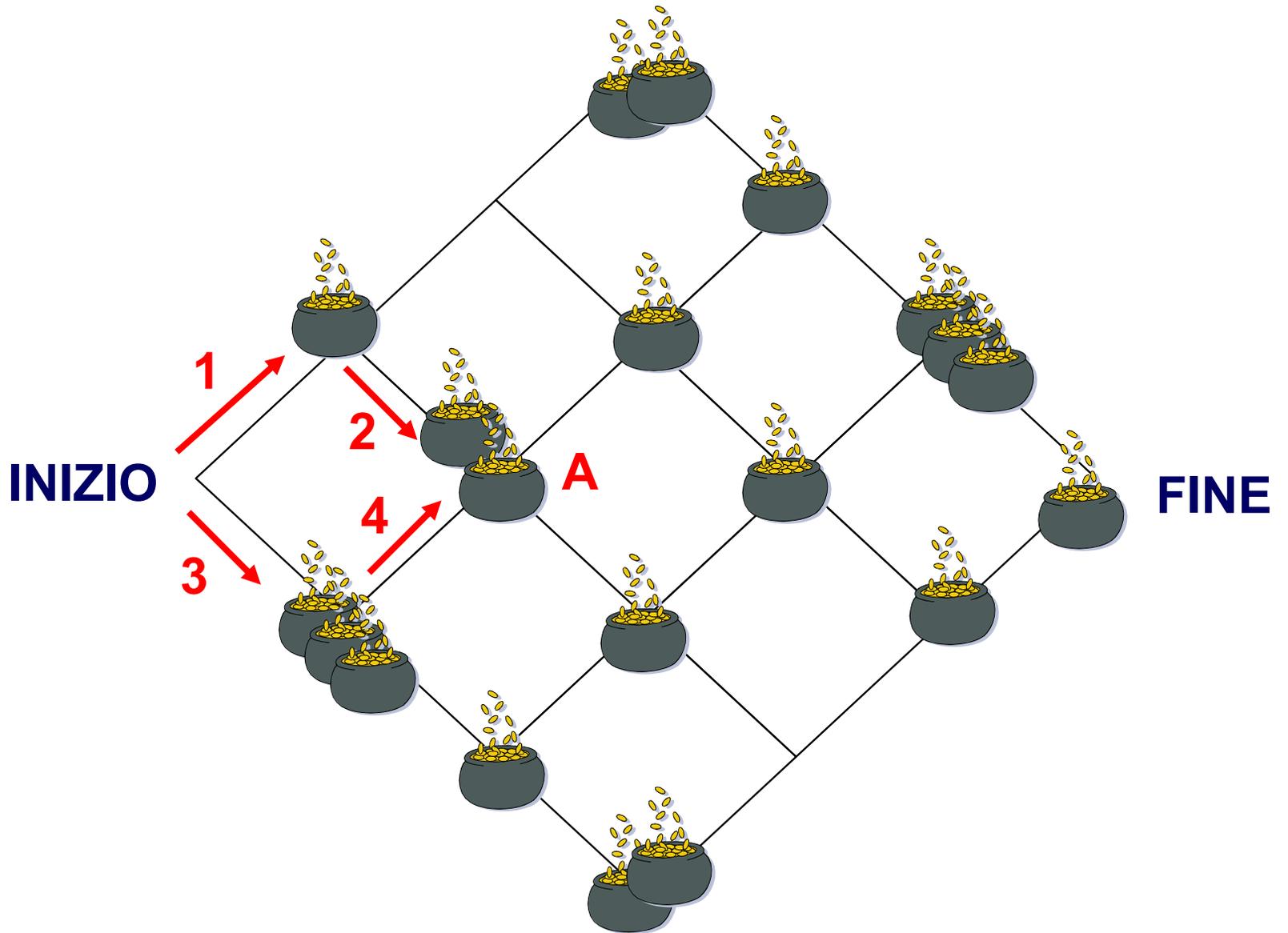
Dividere il problema in tanti sotto-problemi.
Utilizzare la soluzione di ciascun sottoproblema per risolvere i successivi.

Il gioco delle pentole d'oro: regole

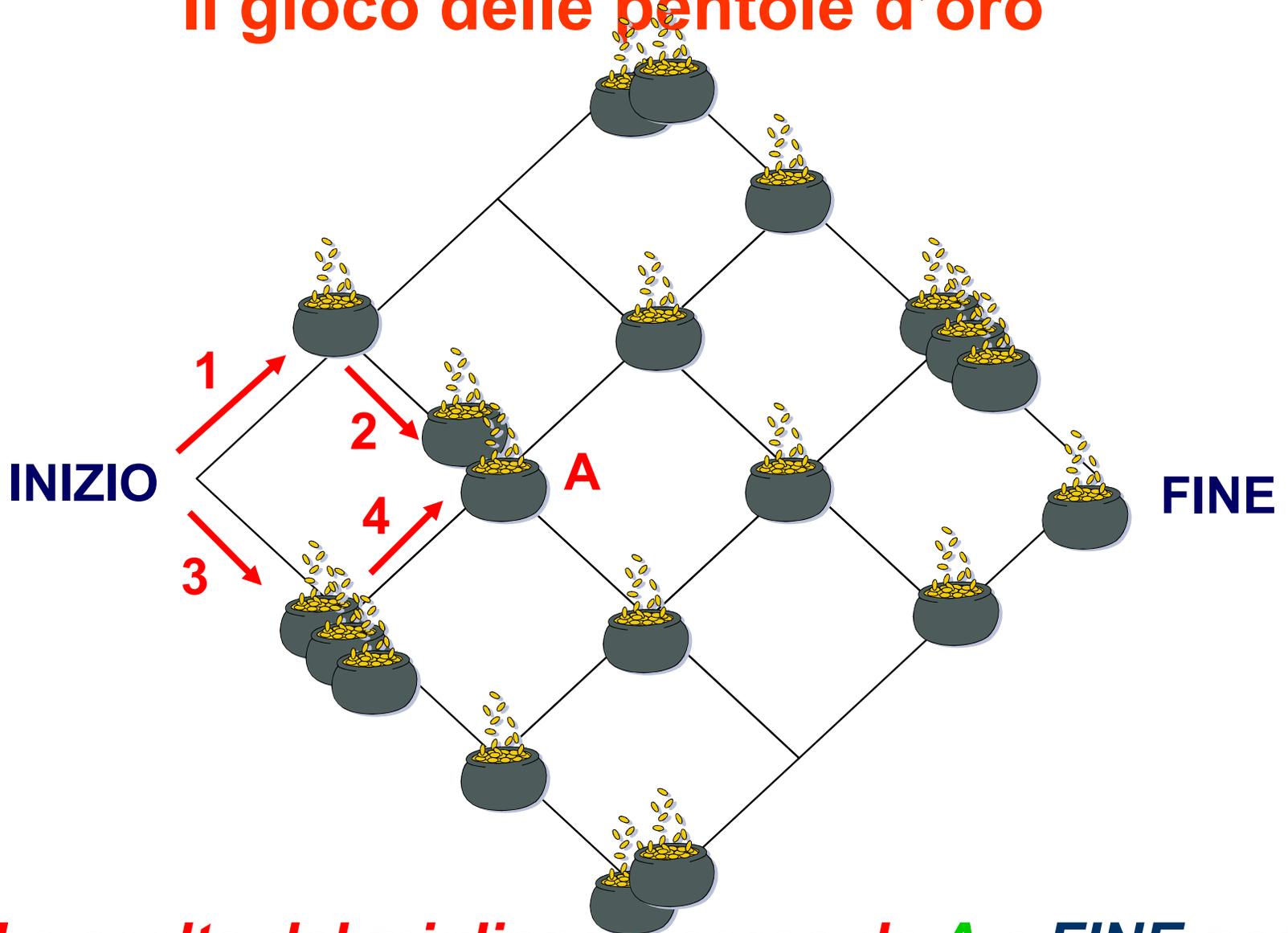


Andare da INIZIO a FINE senza mai passare per lo stesso punto e senza tornare indietro, raccogliendo il max numero di pentole d'oro

Il gioco delle pentole d'oro

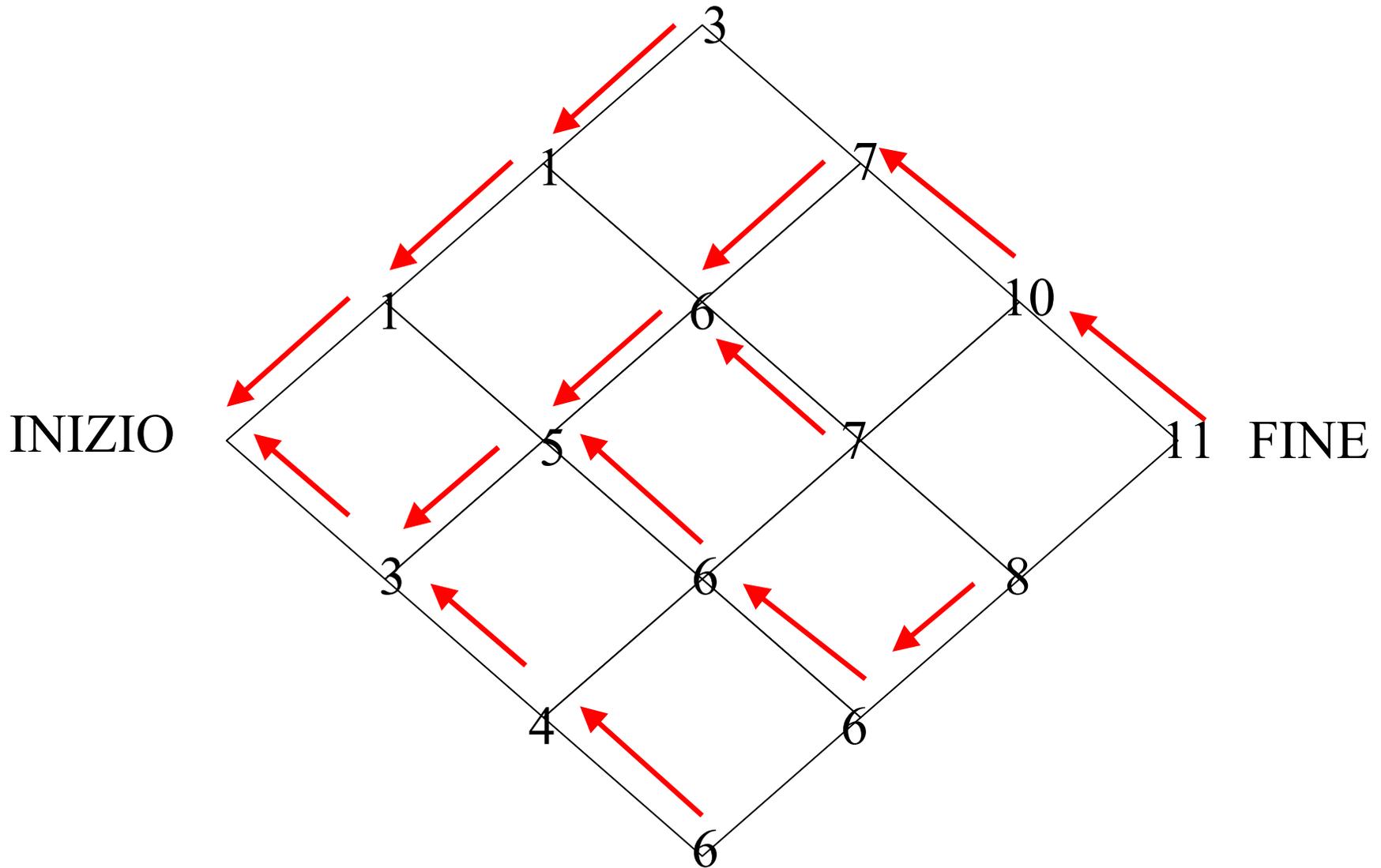


Il gioco delle pentole d'oro

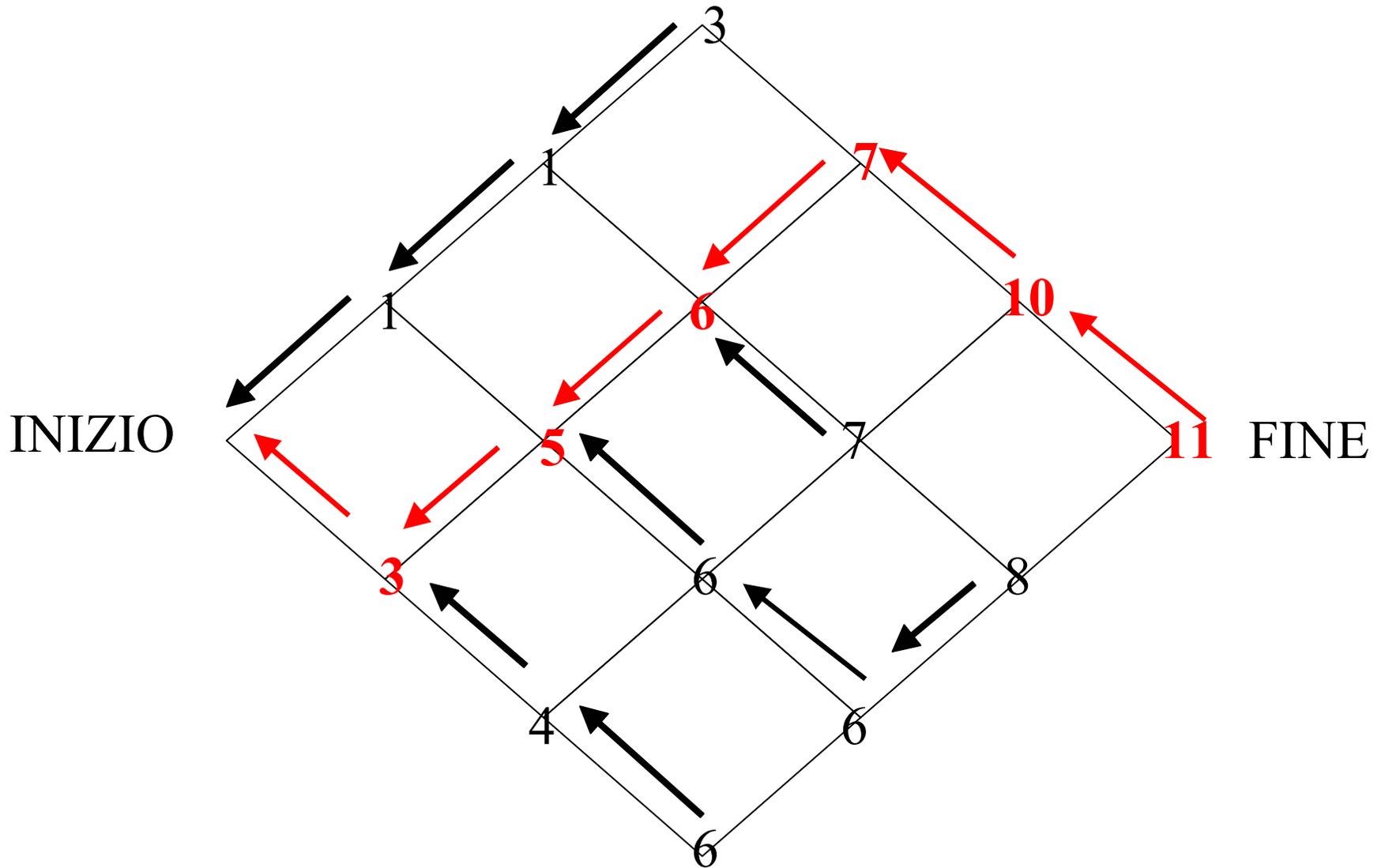


*La scelta del miglior percorso da **A** a **FINE** non dipende da come sono arrivata da **INIZIO** ad **A**.*

Il gioco delle pentole d'oro: soluzione



Il gioco delle pentole d'oro: soluzione



Rispetto alle pentole d'oro, le matrici di similarita' contengono anche **valori negativi**



Un algoritmo di allineamento tendera' a saltarli e ad inserire un numero eccessivo di **inserzioni e delezioni (GAP)**



Sono necessarie le **penalizzazioni** di inserimento dei GAP (sia *di apertura*, che *di continuazione*)

Es. Come allineiamo le due sequenze:

HEAGAWGHEE

PAWHEAE

?

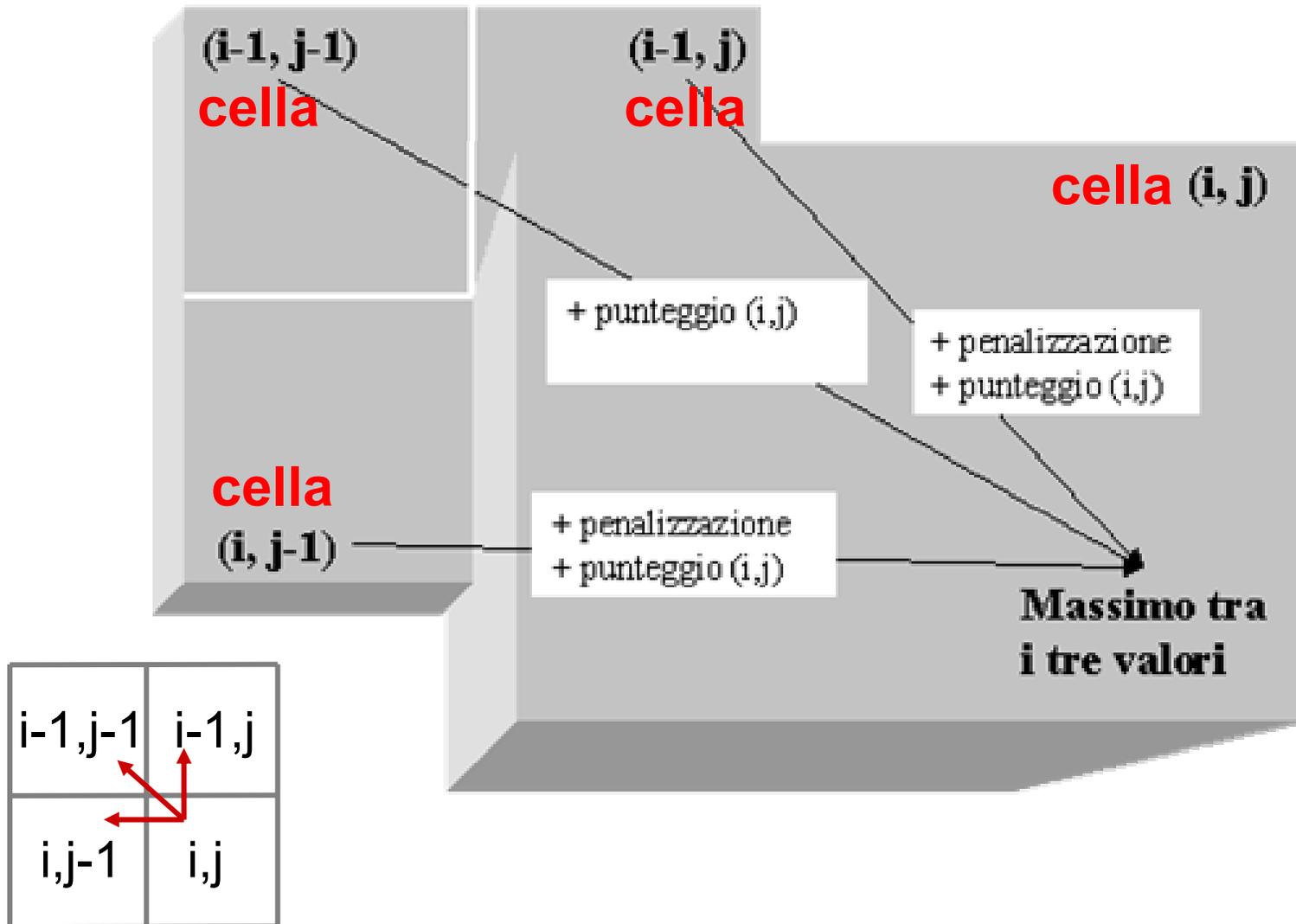
Un esempio di algoritmi di allineamento: matrici cumulative

HEAGAWGHEE vs PAWHEAE

Step1: costruzione della matrice sito-specifica
(valori da BLOSUM45)

	H	E	A	G	A	W	G	H	E	E
P	-2	-1	-1	-2	-1	-4	-2	-2	-1	-1
A	-2	-1	5	0	5	-3	0	-2	-1	-1
W	-3	-3	-3	-3	-3	15	-3	-3	-3	-3
H	10	0	-2	-2	-2	-3	-2	10	0	0
E	0	6	-1	-3	-1	-3	-3	0	6	6
A	-2	-1	5	0	5	-3	0	-2	-1	-1
E	0	6	-1	-3	-1	-3	-3	0	6	6

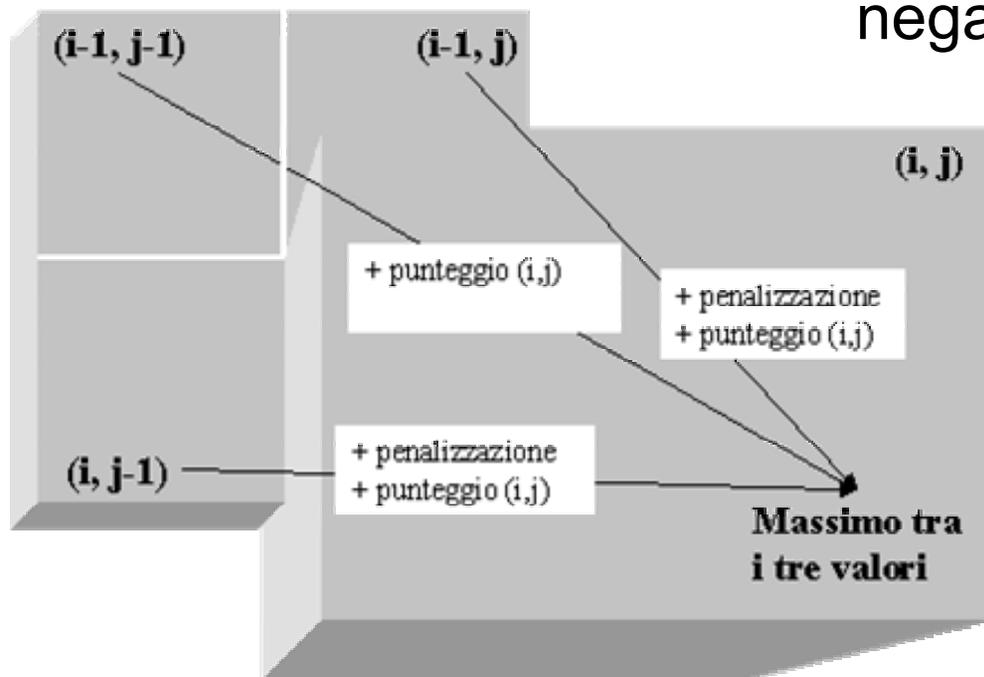
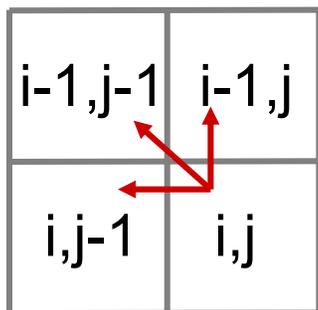
Algoritmo per l'allineamento: matrici cumulative



Step2: costruzione di una matrice *cumulativa*, in cui ogni elemento rappresenta il punteggio massimo ottenibile per arrivare dall'inizio fino a quel punto

$$S_{i,j} = S_{i,j} \text{ (sito-specifica)} + \max \begin{bmatrix} S_{i-1,j-1} \\ S_{i-1,j} + g \\ S_{i,j-1} + g \end{bmatrix}$$

dove g è la penalizzazione per in./del. e assume valori negativi



Matrice cumulativa

delezioni



H E A G A W G H E E

inserzioni

P
A
W
H
E
A
E

0	-8	-16	-24	-32	-40	-48	-56	-64	-72	-80
-8										
-16										
24										
-32										
-40										
-48										
-56										

Inserire una riga ed una colonna di inserzioni e delezioni

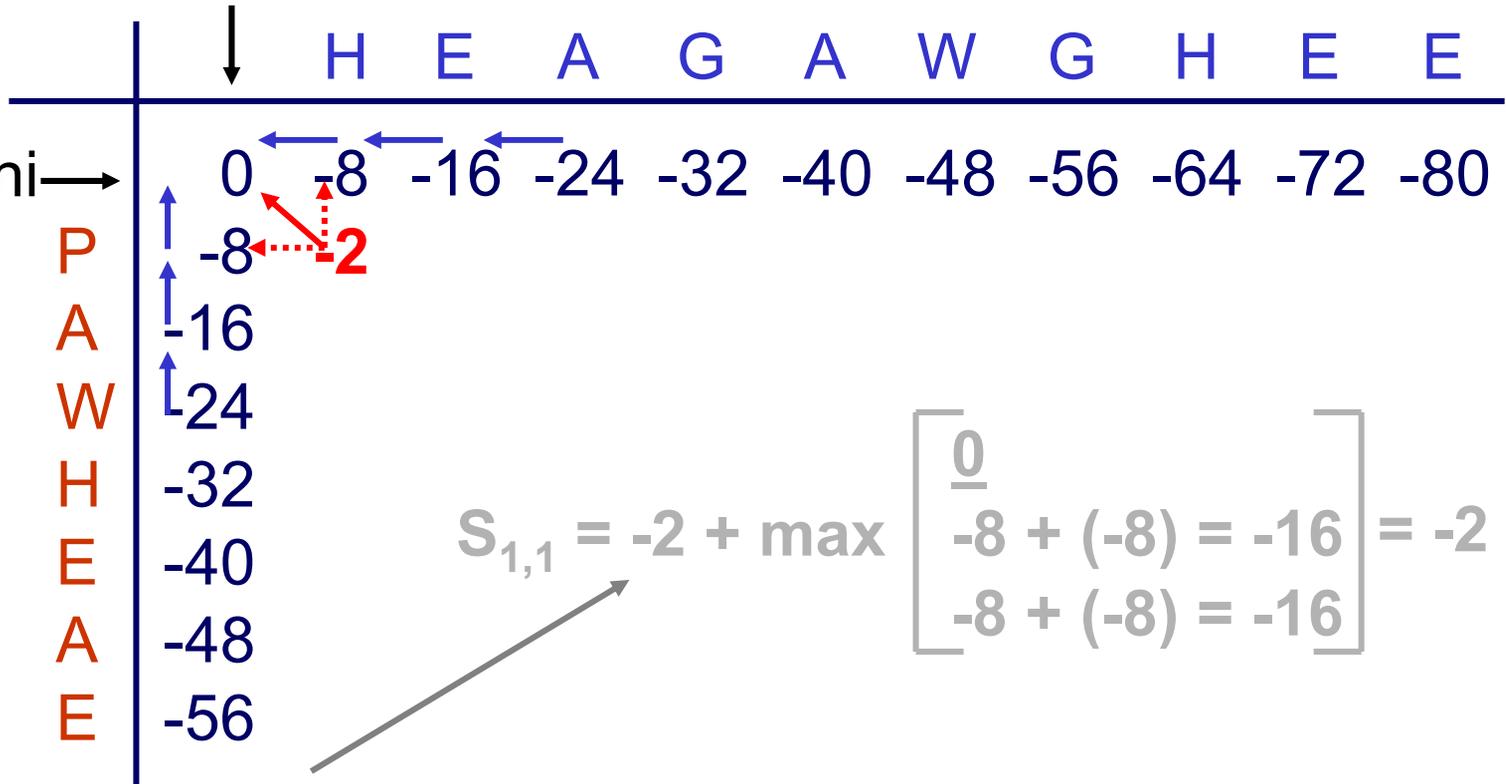
	H	E	A	G	A	W	G	H	E	E
P	-2	-1	-1	-2	-1	-4	-2	-2	-1	-1
A	-2	-1	5	0	5	-3	0	-2	-1	-1
W	-3	-3	-3	-3	-3	15	-3	-3	-3	-3
H	10	0	-2	-2	-2	-3	-2	10	0	0
E	0	6	-1	-3	-1	-3	-3	0	6	6
A	-2	-1	5	0	5	-3	0	-2	-1	-1
E	0	6	-1	-3	-1	-3	-3	0	6	6

Penalizzazione GAP = -8

Matrice sito-specifica originaria

Matrice cumulativa

delezioni



$$S_{1,1} = -2 + \max \begin{bmatrix} 0 \\ -8 + (-8) = -16 \\ -8 + (-8) = -16 \end{bmatrix} = -2$$

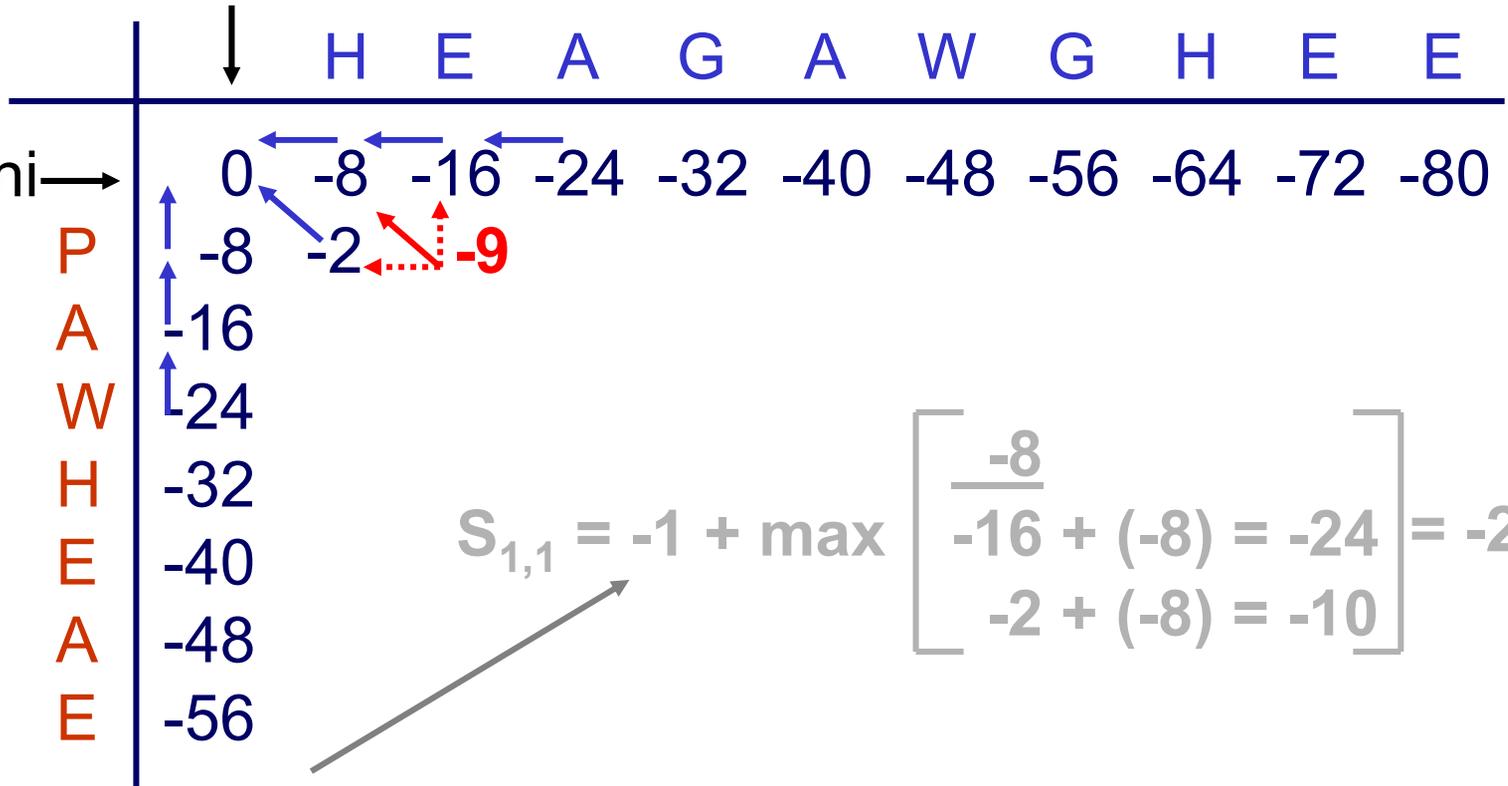
	H	E	A	G	A	W	G	H	E	E
P	<u>-2</u>	-1	-1	-2	-1	-4	-2	-2	-1	-1
A	-2	-1	5	0	5	-3	0	-2	-1	-1
W	-3	-3	-3	-3	-3	15	-3	-3	-3	-3
H	10	0	-2	-2	-2	-3	-2	10	0	0
E	0	6	-1	-3	-1	-3	-3	0	6	6
A	-2	-1	5	0	5	-3	0	-2	-1	-1
E	0	6	-1	-3	-1	-3	-3	0	6	6

Penalizzazione GAP = -8

Matrice sito-specifica originaria

Matrice cumulativa

delezioni



$$S_{1,1} = -1 + \max \begin{bmatrix} -8 \\ -16 + (-8) = -24 \\ -2 + (-8) = -10 \end{bmatrix} = -2$$

	H	E	A	G	A	W	G	H	E	E
P	<u>-2</u>	-1	-1	-2	-1	-4	-2	-2	-1	-1
A	-2	-1	5	0	5	-3	0	-2	-1	-1
W	-3	-3	-3	-3	-3	15	-3	-3	-3	-3
H	10	0	-2	-2	-2	-3	-2	10	0	0
E	0	6	-1	-3	-1	-3	-3	0	6	6
A	-2	-1	5	0	5	-3	0	-2	-1	-1
E	0	6	-1	-3	-1	-3	-3	0	6	6

Penalizzazione GAP = -8

Matrice sito-specifica originaria

Matrice cumulativa

delezioni



H E A G A W G H E E

inserzioni

	0	-8	-16	-24	-32	-40	-48	-56	-64	-72	-80
P	-8	-2	-9								
A	-16	-10	-3								
W	24										
H	-32										
E	-40										
A	-48										
E	-56										

$$S_{2,2} = -1 + \max$$

$$\begin{bmatrix} \underline{-2} \\ -9 + (-8) = -17 \\ -10 + (-8) = -18 \end{bmatrix} = -3$$

	H	E	A	G	A	W	G	H	E	E
P	-2	-1	-1	-2	-1	-4	-2	-2	-1	-1
A	-2	<u>-1</u>	5	0	5	-3	0	-2	-1	-1
W	-3	-3	-3	-3	-3	15	-3	-3	-3	-3
H	10	0	-2	-2	-2	-3	-2	10	0	0
E	0	6	-1	-3	-1	-3	-3	0	6	6
A	-2	-1	5	0	5	-3	0	-2	-1	-1
E	0	6	-1	-3	-1	-3	-3	0	6	6

Penalizzazione GAP = -8

*Matrice sito-specifica
originaria*

Matrice cumulativa

	H	E	A	G	A	W	G	H	E	E	
P	0	-8	-16	-24	-32	-40	-48	-56	-64	-72	-80
A	-8	-2	-9	-17	-26	-33	-44	-50	-58	-65	-73
W	-16	-10	-3	-4	-12	-15	-26	-34	-44	-53	-62
H	-24	-19	-13	-6	-7	-15	0	-11	-22	-33	-44
E	-32	-14	-19	-15	-8	-9	-11	-2	0	-8	-16
A	-40	-22	-8	-17	-18	-9	-12	-13	-2	6	4
E	-48	-32	-17	-3	-11	-12	-12	-12	-12	-3	5
E	-56	-40	-19	-12	-6	-12	-15	-15	-12	-5	3

	H	E	A	G	A	W	G	H	E	E
P	-2	-1	-1	-2	-1	-4	-2	-2	-1	-1
A	-2	-1	5	0	5	-3	0	-2	-1	-1
W	-3	-3	-3	-3	-3	15	-3	-3	-3	-3
H	10	0	-2	-2	-2	-3	-2	10	0	0
E	0	6	-1	-3	-1	-3	-3	0	6	6
A	-2	-1	5	0	5	-3	0	-2	-1	-1
E	0	6	-1	-3	-1	-3	-3	0	6	6

*Matrice sito-specifica
originaria*

	H	E	A	G	A	W	G	H	E	E	
P	0	-8	-16	-24	-32	-40	-48	-56	-64	-72	-80
A	-8	-2	-9	-17	-26	-33	-44	-50	-58	-65	-73
W	-16	-10	-3	-4	-12	-15	-26	-34	-44	-53	-62
H	-24	-19	-13	-6	-7	-15	0	-11	-22	-33	-44
E	-32	-14	-19	-15	-8	-9	-11	-2	0	-8	-16
A	-40	-22	-8	-17	-18	-9	-12	-13	-2	6	4
E	-48	-32	-17	-3	-11	-12	-12	-12	-12	-3	5
E	-56	-40	-19	-12	-6	-12	-15	-15	-12	-5	3

HEAGAWGH-EE

-P--AW-HEAE

massimo
punteggio
ottenibile

Step3: percorso a ritroso attraverso le celle che hanno permesso di ottenere i punteggi migliori → allineamento

		H	E	A	G	A	W	G	H	E	E
	0	-8	-16	-24	-32	-40	-48	-56	-64	-72	-80
P	-8	-2	-9	-17	-26	-33	-44	-50	-58	-65	-73
A	-16	-10	-3	-4	-12	-15	-26	-34	-44	-53	-62
W	-24	-19	-13	-6	-7	-15	0	-11	-22	-33	-44
H	-32	-14	-19	-15	-8	-9	-11	-2	0	-8	-16
E	-40	-22	-8	-17	-18	-9	-12	-13	-2	6	4
A	-48	-32	-17	-3	-11	-12	-12	-12	-12	-3	5
E	-56	-40	-19	-12	-6	-12	-15	-15	-12	-5	3

HEAGAWGH-EE

delezione

-P-AW-HEAE

inserzione



1. Similarità e omologia.
2. Allineamenti di sequenze.
3. Sostituzioni e gap.
4. Algoritmo di allineamento