

Alineamientos Múltiples de Secuencias

Luis E. Garreta U

Curso de Bioinformática

Pontificia Universidad Javeriana – Cali
Facultad de Ingeniería - Carrera de Biología

17 de septiembre de 2018

Introducción

- Hemos visto cómo comparar una secuencia con otra (alineamiento de pares)
- Hemos visto cómo comparar una secuencia con muchas otras en una BD (muchos alineamientos de pares - BLAST)
- Ahora veremos cómo comparar múltiples secuencias simultáneamente, no de dos en dos.

- Las secuencias biológicas a menudo se agrupan en familias
 - ✓ Genes relacionados de un organismo (parálogos)
 - ✓ Genes relacionados de distintas especies (ortólogos)
 - ✓ Secuencias dentro de una población (variantes polimórficas)
- Dos secuencias pueden tener un alineamiento no muy bueno entre ellas, pero pueden alinearse vía una tercera
 - ✓ Identificación de familias y regiones conservadas

Alineamientos Múltiples de Secuencias

- Un alineamiento múltiple de secuencias es un alineamiento de más de dos secuencias.
- Estas secuencias, como en el caso de los alinamientos por parejas pueden ser ADN, ARN o proteína.

```
Hsa_TM666 ALTLHYDRYTTSRRLDPIPLKCVGGTAGCDSYTPKVIQCQNKGDGVDVQWECKTDLDI 103
Ptr_TM666 ALTLHYDRYTTSRRLDPIPLKCVGGTAGCDSYTPKVIQCQNKGDGVDVQWECKTDLDI 103
Ppy_TM666 ALTLHYDRYTTSRRLDPIPLKCVGGTAGCDSYTPKVIQCQNKGDGVDVQWECKTDLDI 103
Mml_TM666 ALTLHYDRYTTSRRLDPIPLKCVGGTAGCDSYTPKVIQCQNKGDGVDVQWECKTDLDI 103
Mfa_TM666 ALTLHYDRYTTSRRLDPIPLKCVGGTAGCDSYTPKVIQCQNKGDGVDVQWECKTDLDI 103
Mne_TM666 ALTLHYDRYTTSRRLDPIPLKCVGGTAGCDSYTPKVIQCQNKGDGVDVQWECKTDLDI 103
Ssc_TM666 ALTLHYDRYTTSRRLDPIPLKCVGGTAGCDSYTPKVIQCQNKGDGVDVQWECKTDLDI 103
Bta_TM666 ALTLHYDRYTTSRRLDPIPLKCVGGTAGCDSYTPKVIQCQNKGDGVDVQWECKTDLDI 103
Cfa_TM666 ALTLHYDRYTTSRRLDPIPLKCVGGTAGCDSYTPKVIQCQNKGDGVDVQWECKTDLDI 103
Mmu_TM666 ALTLHYDRYTTSRRLDPIPLKCVGGTAGCDSYTPKVIQCQNKGDGVDVQWECKTDLDI 104
Rno_TM666 ALTLHYDRYTTSRRLDPIPLKCVGGTAGCDSYTPKVIQCQNKGDGVDVQWECKTDLDI 104
Ocu_TM666 ALTLHYDRYTTSRRLDPIPLKCVGGTAGCDSYTPKVIQCQNKGDGVDVQWECKTDLDI 97
Laf_TM666 ALTLHYDRYTTSRRLDPIPLKCVGGTAGCDSYTPKVIQCQNKGDGVDVQWECKTDLDI 89
Mdo_TM666 ALTLHYDRYTTSRRLDPIPLKCVGGTAGCDSYTPKVIQCQNKGDGVDVQWECKTDLDI 119
Gga_TM666 VLTLLHGRYTTARRTAAPVQLQICGGSAGCS-DIPEVVQCYNRGWDGVDVQWECKADLEN 94
Xla_TM666 TITLYADRYTTARRSAPVPQLKICGGSAGCHAMPVQVQVQCHNRGWDGVDVQWECKADMDN 93
Xtr_TM666 AITLYADRYTTARRSAPVPQLKICGGSAGCHAMPVQVQVQCHNRGWDGVDVQWECKADMDN 93
Dre_TM666 VLTLYGRYTTARRSSPVPQLQICGGSAGCSFTPEVVQCYNRGSDGIDQWECKADMDN 93
Ssa_TM666 VLTLYGRYTTARRSSAPVPQLQICGGSAGCSFTPEVVQCYNRGSDGIDQWECKADMDN 93
Tru_TM666 VLTLYGRYTTARRSSPVPQLQICGGSAGCSFTPEVVQCYNRGSDGIDQWECKADMDN 99
Tni_TM666 TLTLYGRYTTARRSSPVPQLKICGGSAGCSFTPEVVQCYNRGSDGIDQWECKADMDN 89
Gac_TM666 ALTLYKRYTTARRSAPVPQLQICGGSAGCSFTPEVVQCYNRGSDGIDQWECKADMDN 92
Ppr_TM666 VLTLYKRYTTARRSSPVLQQLCAGGTAGCGSFVPEVVQCYNRGSDGIDQWECKADMDN 93
Cel_TM666 AITLHGKMTTGRVSPFTQLKCVGG-SAKGAFTPKVVQCANQGFDSGVDVQWCRDADLPH 96
Cre_TM666 AITLHGKMTTGRVSPFTQLKCVGG-SAKGAFTPKVVQCANQGFDSGVDVQWCRDADLPH 96
Cbr_TM666 AITLHGKMTTGRVSPFTQLKCVGG-SAKGQFSPKVVQCANQGFDSGVDVQWCRDADLPH 96
```

```
. : * * . * * * . * : * * : . * : : * * : * * * * * . : :
```

Definiciones

■ Residuo:

- ✓ Se refiere a secciones homólogas de las secuencias
- ✓ A veces se llama residuo a cada columna del alineamiento

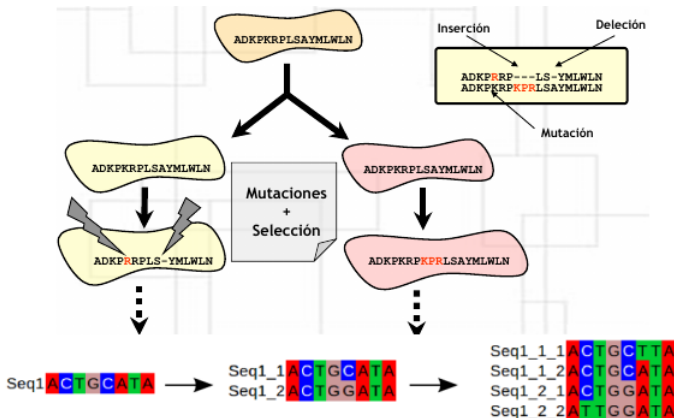
■ Los residuos se pueden originarse ya sea

- ✓ **Evolutivamente:** presumiblemente provenientes de un ancestro común
- ✓ **Estructuralmente:** suelen ocupar lugares relevantes en la estructura 3D

beta globin	NFRLLGNVLVCVLAHHF-GKEFTPPVQAAAYQKVVGAVANALAHKYH-----] Secuencias alineadas
myoglobin	YLEFISECIIQVLQSKH-PGDFGADAQGAMNKALELFRKDMASNYKELGFQG	
neuroglobin	SFSTVGESLLYMLEKCL-GPAFTPATRAANSQLYGAVVQAMSRGWDGE---	
soybean	QFVVVKEALLKTIKAAV-GDKWSDLSRAWEVAYDELAALAIKKA-----	
rice	HFEVVKFALLDTIKKEVPADMWS PAMKSAWSEAYDHLVAAIKQEMKPAE---	
	: : : : : * . . :	Residuo

Un Alineamiento nos cuenta una historia

En todos los casos los algoritmos de alineamiento múltiple asumen que las secuencias que se están alineando descienden de un antepasado común y lo que se intenta hacer es alinear las posiciones homólogas.



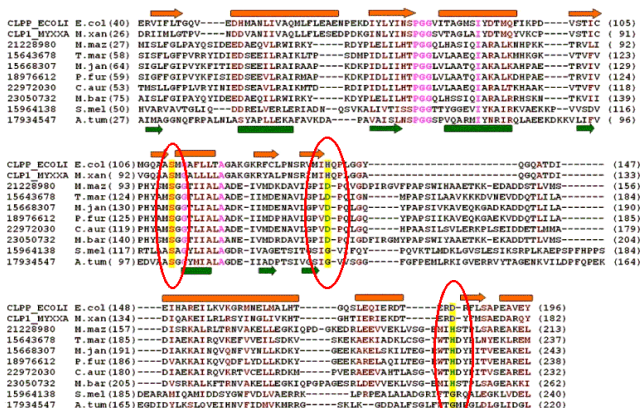
- Para hacer un alineamiento, generalmente necesitamos seleccionar:
 - ✓ 1. Las secuencias homólogas a alinear
 - ✓ 2. El software que utilice una función de puntuación óptima
 - ✓ 3. Los parámetros adecuados (fundamentalmente huecos)
- No hay un alineamiento perfecto
 - ✓ Las secuencias evolucionan más rápido que las estructuras o funcionalidades
 - ✓ La secuencia puede variar y la estructura o función seguir invariante

Aplicaciones de los Alineamientos Múltiples

- Dar información acerca de la función, estructura y evolución de una secuencia
 - ✓ Al conocer cómo se alinea respecto a un grupo de secuencias
 - ✓ Válido para análisis de genes, proteínas o poblaciones
- Encontrar miembros distantes de una familia de proteínas
 - ✓ Es muy frecuente que estén distantes, y el alineamiento de pares no suele ser lo suficientemente preciso para encontrarlos
- Primer paso (y el más importante) en la generación de árboles filogenéticos

Aplicaciones:

Determinar elemento común que corresponde al Sitio Activo de la Sestrina



Aplicaciones:

Extrapolación: determinar la función de una proteína desconocida con base en proteína conocidas

```
chite  ---ADKPKRPLSAYMLWLNSARES IKRENPDFK-VTEVAKKGGELWRGLKD
wheat  --DPNPKRAPSAFFVFMGEFEEFKQKNPKNKSVAAVGKAAGERWKSISE
trybr  KKDSNAPKRAMTSFMFFSSDFRS----KHS DLS-IVEMSKAAGA AWKELGP
unknown ----KPKRPR SAYNIYVSESFQ----EAKDDS-AQGK LKLVNEAWKNLSP
          ***. :.:. :. . . . : . . . * . *: *

chite  AATAKQNYIRALQEYERNGG-
wheat  ANKLKGEYNKAIAAYNKGESA
trybr  AEKDKERYKREM-----
unknown AKDDRIRYDNEMKSWEEQMAE
          * : . * . :
```

Aplicaciones:

Bloques o patrones comunes que corresponden a Motivos o Dominios de Proteínas

```
chite  ---ADKPKRPLSAYMLWLNLSARESISKRENPDFK-VTEVAKKGGELWRGLKD
wheat  --DPNKPKRAPSAFFVFMGEFREFEFQKNPKNKSVAAVGKAAGERWKSLSSE
trybr  KKDSNAPKRAMTSFMFFSSDFRS----KHSDFS-IVEMSKAAGAAWKELGP
mouse  -----KPKRPRSAYNIVVSESFQ----EAKDDS-AQGKCLKVNEAWKNLSP
```

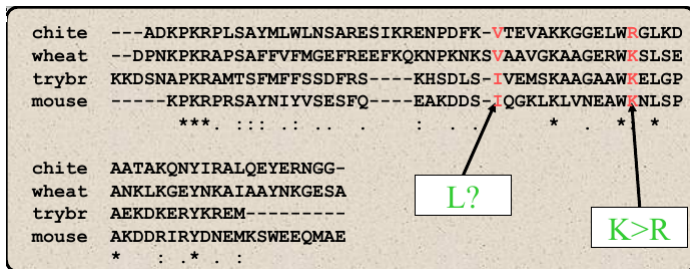
```
      ***. ::: .: .. .      : . .      * . *: *
```

```
chite  AATAKQNYIRALQEYERNGG-
wheat  ANKCLKGEYNKAIAAYNKGESA
trybr  AEKDKERYKREM-----
mouse  AKDDRIRYDNEMKSWEEQMAE
```

```
*      : . * . :
```

Aplicaciones

Elementos representativos comunes que definen el perfil de una familia de proteínas



Aplicaciones

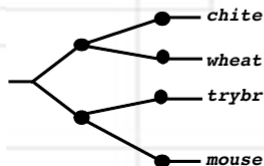
El alineamiento me guía para la construcción árboles filogenéticos

```
chite  ---ADKPKRPLSAYMLWLNLSARESISIKRENPDFK-VTEVAKKGGELWRGLKD
wheat  --DPNPKKRAPSFAFFVFMGEFFREEFKQKNPKNKSVAAVGKAAGERWKSLSSE
trybr  KKDSNAPKRAMTSFMFFSSDFRS----KHS DLS-IVEMSKAAGA AWKELGP
mouse  -----KPKRPRSA YNIYVSESFQ-----EAKDDS-AQGK LKL VNEAWKNLSP
```

```
***. : : : . : . . * . *: *
```

```
chite  AATAKQNYIRALQEYERNGG-
wheat  ANKLKGEYNKAIAAYNKGESA
trybr  AEKDKERYKREM-----
mouse  AKDDRI RYDNEMKSWEEQMAE
```

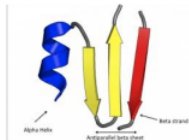
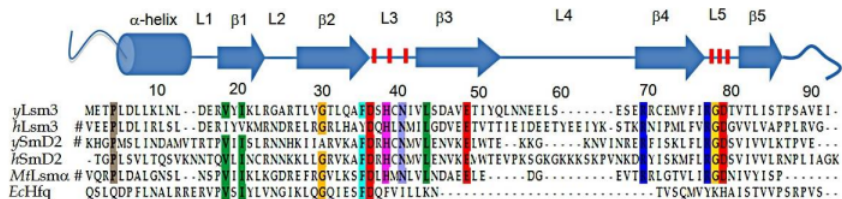
```
* : . * . :
```



-Evolución
-Paralogía/Ortología

Aplicaciones de los AMS:

Determinación de la estructura 3D de una proteína



Selección de Secuencias para AMS

- Es su trabajo el de seleccionar las secuencias
- Dos factores a tener en cuenta:
 - ✓ El número de secuencias
 - ✓ La naturaleza de las secuencias
- Un número razonable de secuencias: 20 a 50
 - ✓ Inicia con 10 o 15 secuencias
 - ✓ Alineamientos pequeños son fáciles de mostrar y analizar
- Tipos de secuencias:
 - ✓ ADN, ARN, Proteínas (AA)
- Evite secuencias muy similares o muy distantes

Qué tipo de secuencias seleccionar: ADN o Proteínas

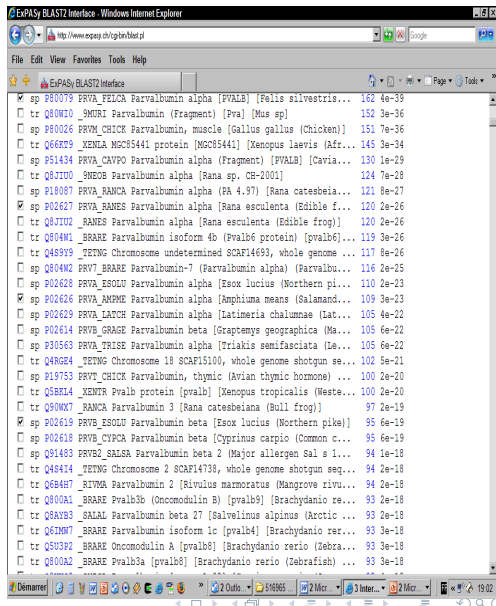
- Secuencias de ADN son más complicadas de alinear que proteínas ya que se nucleótidos se repiten demasiado
- Mayoría métodos trabajan con proteínas
- Si su ADN es codificante, trasládalo a proteínas
- Si las secuencias son homologas:
 - ✓ **En toda su extensión:** Use alineamientos progresivos
 - ✓ **Solo en partes:** Use descubrimiento de motivos.

Obtención de Secuencias con BLAST

- El método más conveniente para seleccionar sus secuencias es a través de un servidor BLAST
- Servidores BLAST tienen integrados múltiples métodos de alineamiento:

Obtención de Secuencias con BLAST

- Seleccione algunas del tope
- Seleccione uniformemente algunas del medio que no sean tan idénticas
- La idea es tener secuencias:
 - ✓ No tan similares
 - ✓ No tan distantes



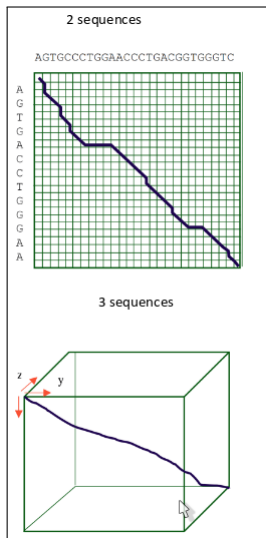
Muchos Algoritmos y Herramientas

- Existen muchos algoritmos
- La mayoría son métodos aproximativos donde no se garantiza el alineamiento perfecto:
- Todos los métodos tienen pros y contras:
 - ✓ CLUSTALW: alineamiento progresivo
 - ✓ MUSCLE: alineamiento iterativo
 - ✓ T-COFFEE: alineamiento basado en consistencia
 - ✓ EXPRESSO: alineamiento basada en la estructuras 3D

Algoritmos para Alineamiento Múltiple

- Existen cinco tipos de algoritmos:
 - ✓ 1. Métodos exactos
 - ✓ 2. Alineamiento progresivo
 - ✓ 3. Aproximaciones iterativas
 - ✓ 4. Métodos basados en la consistencia
 - ✓ 5. Métodos basados en la estructura
- Las aproximaciones no son excluyentes
 - ✓ Las tres últimas, por ejemplo, utilizan alineamiento progresivo

Métodos exactos



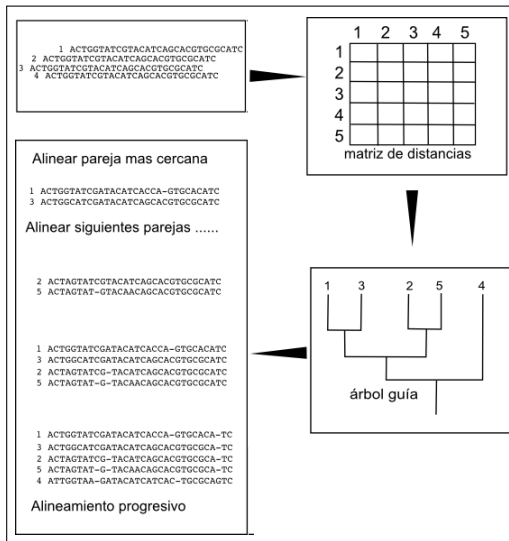
- Similar a alineamiento global de pares pero ahora con más de 3 secuencias:
 - ✓ Se basan en programación dinámica
- Aseguran un alineamiento óptimo (perfecto), pero son lentos
 - ✓ No son factibles ni en espacio ni en tiempo si tenemos más de unas pocas secuencias
- Se prefieren los métodos inexactos, mucho más rápidos
 - ✓ ClustalW
 - ✓ MUSCLE
 - ✓ T-COFFEE

Alineamiento progresivo

■ “Progresivo”

- ✓ Calcula alineamientos de pares entre las secuencias consideradas
- ✓ Elige el mejor alineamiento de entre ellos
- ✓ Añade progresivamente más secuencias al alineamiento escogido

- ## ■ El programa de alineamiento progresivo más usado, pero ya no el mejor es ClustalW



Algoritmo de ClustalW

Clustal implementa el algoritmo de Feng y Doolittle, que consta de 3 fases:

- 1 Alineamiento global 2 a 2 mediante el algoritmo de NW
- 2 Se crea un árbol guía
- 3 Se crea el alineamiento múltiple paso a paso

Fase 1: Alineamiento global de pares

- Ejemplo: cinco globinas muy conocidas, bastante distantes:
✓ NP_000509, NP_005359, NP_067080, 1FSL, 1D8U
- Para 5 secuencias tendremos 10 alineamientos
- Para n secuencias tendremos $n!/(2(n-2)!)$ alineamientos

SeqA	Name	Length	SeqB	Name	Length	Score
1	beta_globin	147	2	myoglobin	154	25.17
1	beta_globin	147	3	neuroglobin	151	15.65
1	beta_globin	147	4	soybean_globin	144	13.19
1	beta_globin	147	5	rice_globin	166	21.09
2	myoglobin	154	3	neuroglobin	151	16.56
2	myoglobin	154	4	soybean_globin	144	8.33
2	myoglobin	154	5	rice_globin	166	12.99
3	neuroglobin	151	4	soybean_globin	144	17.36
3	neuroglobin	151	5	rice_globin	166	18.54
4	soybean_globin	144	5	rice_globin	166	43.06

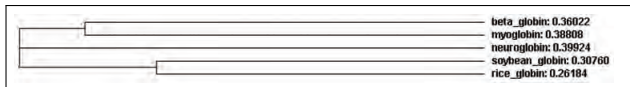
Fase 2: Creación del árbol guía

- La longitud de las ramas depende de las distancias
- Se unen las ramas de las secuencias con distancias más cortas
- Árbol:

Formato de
Newid (.nwk)

```
{
(
  beta_globin:0.36022,
  myoglobin:0.38808)
:0.06560,
  neuroglobin:0.39924,
  (
    soybean_globin:0.30760,
    rice_globin:0.26184)
  :0.13652);
```

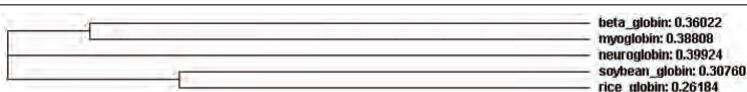
Representación
Gráfica:



Fase 3: Creación del alineamiento múltiple

- Se seleccionan las dos secuencias más cercanas según el árbol guía
- Se seleccionan las dos secuencias más cercanas siguientes
 - ✓ Si alguna coincide, se añade al alineamiento de pares, dando lugar a un alineamiento de 3+ secuencias, o perfil
 - ✓ Si ninguna coincide con las anteriores, se realiza su alineamiento de pares, para después realizar alineamiento de perfiles.
- El alineamiento continúa hasta llegar a la raíz del árbol

Fase 3: Creación del alineamiento múltiple



Notación:

Punto (.): Coincidencia

Dos puntos (:): Coincidencia Alta.

Asterisco (*): Coincidencia Exacta

beta_globin -----MVHLT**PEEK**SAVTALWGKVN--VDEVGGEALGRLLVVYPWTQR**FES**FG- 47
myoglobin -----MGLSD**GEW**QLVLNVWGKVEADIPGHGQEVLI**RLFKGH**PETLEK**FD**KPK- 48
neuroglobin -----MERPE**PELI**RQSWRAVSRS**PLE**HGTVLFA**RLFA**LEPDLL**PLFQ**YNCR 47
soybean_globin -----MVAFT**EKQ**DALVSS**FEAF**KANIPQYSVVFT**YSILE**KAPAK**DLFS**FLA- 49
rice_globin MALVEDNNNAVAVSF**SE**EQBALVLKSWAILKKD**SANIAL**RRFFLK**IFEV**APSASQ**MF**SFLR- 59

beta_globin DLSTPD**AVM**GNPKVKA**HG**KKVLGAFSDGLAHLN**DLK**GT**FAT**-----LS**ELH**CDKLHVD**P** 101
myoglobin HLKSEDEMKAS**EDLKK**HGATVLTALGGILKKKG**HHEA**EIKP-----LAQ**SH**ATKHKIPV 102
neuroglobin QFSSPEDCLSS**PEFLD**HIRK**VML**VIDAAVTINVEDLSS**LEB**Y-----LASLG**RKH**RAVG**VKL**S 104
soybean_globin NGVDPT--N**PKLTG**HAEKLFALVRDSAGQLKASGTVVAD----AALGS**VH**AQKAV**TD**P 101
rice_globin --NSDVPLEKN**PKLKT**HAMSVFVMTCEAAQLRKAGKVTVR**DTTL**KRLGATH**LKY**GVGDA 117

beta_globin ENFRLLGNVLVCVLA**HHF**GKEFTPPVQAA**YQ**KVVAGVANALAHKYH----- 147
myoglobin KYLEFISECIIQVLQ**SKH**PGDFGADAQGAMNKAL**ELFR**KDMASNYKELGFQ**G** 154
neuroglobin SFSTVGESLLYMLEKCLG-PAFT**PATRA**AWSQ**LYG**AVVQAMSRGWDGE---- 151
soybean_globin QFVVVKEALL**TKI**KA**AVG**-DKWSDELSRAWEVAYDELA**AAIK**KA----- 144
rice_globin HPEVV**KF**ALLDTIKBEVPADMWS**PAMK**SAWS**EAYD**HLVAAIKQEMK**PAE**--- 166

Ejercicio 1

- Realizar el anterior alineamiento múltiple mediante el ClustalW del EBI: <http://www.ebi.ac.uk/clustalW> y compare con el alineamiento anterior:
- Cinco beta globinas bastante distantes:

```
>beta_globin 2hhbB NP_000509.1 [Homo sapiens]
MVHLTPEEKSAVTALWGKVNVDDEVGGEALGRLLVYPWTQRRFFESFGDLSTPDVAVMGNPVKVKAHGKKVLG
AFSDGLAHLNLDLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVAN
ALAHKYH
>myoglobin 2MM1 NP_005359.1 [Homo sapiens]
MGLSDGEWQLVLNVWGKVEADIPGHGQEVLRFLFKGHPETLEKFDKFKHLKSEDEMKASEDLKKHGATVL
TALGGILKKKGHHEAIKPLAQSHATKHKIPVKYLEFISECIIQVLQSKHPGDFGADAQGAMNKALELFR
KDMASNYKELGFQG
>neuroglobin 1OJ6A NP_067080.1 [Homo sapiens]
MERPEPELIRQSWRAVSRSPLEHGTVLFARLFALEPDLLPLFQYNCRQFSSPEDCLSSPEFLDHIRKVML
VIDAAVTNVEDLSSLEEYLASLGRKHRAVGVLKSSFSTVGESLLYMLEKCLGPAFTPATRAAWSQLYGAV
VQAMSRGWGDE
>soybean_globin 1FSL leghemoglobin P02238 LGBA_SOYBN [Glycine max]
MVAFTEKQDALVSSSFSAFANIPQYSVVFYTSILEKAPAAKDLFSFLANGVDPTNPFLTGHAEKLFALV
RDSAGQLKASGTVVADAALGSVHAQKAVTDPQFVVVKEALLKTIKAAVGDKWSDELSRAWEVAYDELA
IKKA
>rice_globin 1D8U rice Non-Symbiotic Plant Hemoglobin NP_001049476.1
MALVEDNNAVAVSFSEEQALVLKSWAILKKDSANIALRFFLKIFEVAPSASQMFSLRNSDVPLEKNPK
LKTTHAMSVFVMTCEAAALQRLKAGKVTVRDITLRLGATHLKYGVGDAHFVVKFALLDTIKEEVPADMWS
PAMKSAWSEAYDHLVAIAIKQEMKPAE
```

Ejercicio 2

- Realizar el siguiente alineamiento múltiple mediante el ClustalW del EBI: <http://www.ebi.ac.uk/clustalW>
- Cinco globinas bastante cercanas:

>human_NP_000509

MVHLTPEEKSAVTALWGKVNVDDEVGGEALGRLLVYPWTQRRFFESFGDLSTPDVAMGNPKVKAHGKKVLG
AFSDGLAHLNLDNLKGTFFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVAN
ALAHKYH

>Pan_troglodytes_XP_508242

MVHLTPEEKSAVTALWGKVNVDDEVGGEALGRLLVYPWTQRRFFESFGDLSTPDVAMGNPKVKAHGKKVLG
AFSDGLAHLNLDNLKGTFFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVAN
ALAHKYH

>Canis_familiaris_XP_537902

MVHLTAEKSLVSGLWGKVNVDDEVGGEALGRLLIVYPWTQRRFDSFGDLSTPDVAMSNKVKAHGKKVLN
SFSDGLKNLDNLKGTFAKLSELHCDKLHVDPENFKLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVAN
ALAHKYH

>Mus_musculus_NP_058652

MVHLTDAEKSAVSLWAKVNPDEVGGEALGRLLVYPWTQRYFDSFGDLSSASAIMGNPKVKAHGKKVIT
AFNEGLKNLDNLKGTFFASLSELHCDKLHVDPENFRLLGNAIVIVLGHHLGKDFTPAQAQAFQKVVAGVAT
ALAHKYH

>Gallus_gallus_XP_444648

MVHWTAEEKQLITGLWGKVNVAECGAELARLLIVYPWTQRRFFASFGNLSSPTAILGNPMVRAHGKKVLT
SFGDAVKNLNDNIKNTFSQLSELHCDKLHVDPENFRLLGDILIVLAHFSKDFTPECQAQAWQKLVVRVVAH
ALARKYH

Alineamiento Beta Globinas Cercanas

- The coloring scheme (ClustalW program) includes groups such as:
 - ✓ acidic amino acids (blue),
 - ✓ basic amino acids (magenta), and
 - ✓ hydrophobic residues (red).

CLUSTAL 2.1 multiple sequence alignment

```
human_NP_000509          MVHLTPEEKSAVTALNGKVINDEVGGEALGRLLVVYPWTQRFFESFGDLS 50
Pan_troglodytes_XP_508242 MVHLTPEEKSAVTALNGKVINDEVGGEALGRLLVVYPWTQRFFESFGDLS 50
Canis_familiaris_XP_537902 MVHLTAEKSLVSGLNGKVINDEVGGEALGRLLVYPWTQRFFDSFGDLS 50
Mus_musculus_NP_058652    MVHLTDAEKSAVSLWAKVNPDEVGGEALGRLLVVYPWTQRIFYDSFGDLS 50
Gallus_gallus_XP_444648   MVHWTAEKQLITGLNGKVINVAECGAELARLLVYPWTQRFFASFGNLS 50
                          *** *  **  :: ** **  *  *  ***  *****  *  *** **

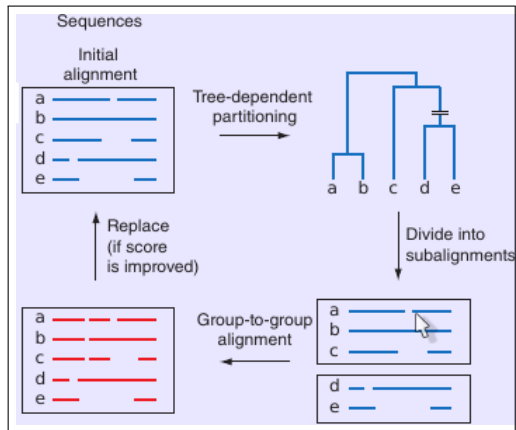
human_NP_000509          IPDAVMGNPKVKAHGKKVLGAFSDGLAHLDNLKGITFATLSSELHCDKLEVD 100
Pan_troglodytes_XP_508242 IPDAVMGNPKVKAHGKKVLGAFSDGLAHLDNLKGITFATLSSELHCDKLEVD 100
Canis_familiaris_XP_537902 IPDAVMSNAKVAHGKKVLNSFSDGLKLNLDNLKGITFAKLSSELHCDKLEVD 100
Mus_musculus_NP_058652    SASAIMGNPVKVAHGKKVITAFNEGLKLNLDNLKGITFASLSSELHCDKLEVD 100
Gallus_gallus_XP_444648   SPTAILGNPMVRAHGKKVLTSFGDAVKLNLDNLKNTFSQISELHCDKLEVD 100
                          . .  *::*  *::*****  ::::  :  *****  **  *****

human_NP_000509          PENFRLLGNVLVCVLAHHFGKEFTPFVQAAYQKVVAGVANALAHKYH 147
Pan_troglodytes_XP_508242 PENFRLLGNVLVCVLAHHFGKEFTPFVQAAYQKVVAGVANALAHKYH 147
Canis_familiaris_XP_537902 PENFRLLGNVLVCVLAHHFGKEFTPFVQAAYQKVVAGVANALAHKYH 147
Mus_musculus_NP_058652    PENFRLLGNIAIVTVLGHHLGKDFTPAAQAAYQKVVAGVATALAHKYH 147
Gallus_gallus_XP_444648   PENFRLLGDILITVLAHHFSKDFTECCQAANQKLVRVVAHALARKYH 147
                          *****  ::  **  *::***  *****  **  *****
```


ClustalW y huecos

- ClustalW sigue la política: “una vez se encuentra un hueco, siempre hay un hueco”
 - ✓ Cuando hay un hueco en un alineamiento, se fomenta que se conserve en alineamientos posteriores
 - ✓ Da al alineamiento múltiple una estructura de “bloques”
- Generalmente, alineamientos con más huecos (menos compactos) coinciden mejor con la filogenia y la estructura de proteínas conocidas como la globina

Algoritmos Iterativos



- Calculan una solución subóptima mediante un alineamiento progresivo
- Modifican el alineamiento y repiten el proceso el alineamiento como inicio hasta que la solución converge
- Se corrige el error de los alineamiento progresivo, donde una vez que cometemos un error, no lo podemos corregir
- Algoritmos conocidos: MAFFT, MUSCLE

Aproximaciones basadas en la consistencia

- Este enfoque combina alineamientos progresivos e iterativos y métodos probabilísticos
- Esta estrategia suele generar alineamientos de secuencia mucho más precisos
- ProbCons y T-Coffee son los dos algoritmos más conocidos
- M-COFFEE: usa muchos algoritmos (T-COFFEE, ClustalW, MAFFT, MUSCLE, and ProbCons)

Taller1: Comparaciones de diferentes Algoritmos de MSA

- Realizar el alineamiento múltiple de las siguientes 9 globinas usando:
(a) CLUSTALW, (b) MAFFT, (b) MUSCLE, y (d) T-COFFEE.
 - ✓ (1) hbb_human (human NP_000509.1, 1HBB);
 - ✓ (2) hbb_chimp (Pan_troglodytes XP_508242.1, no structure);
 - ✓ (3) hbb_dog (Canis lupus familiaris NP_001257813.1, 2QLS|B);
 - ✓ (4) hbb_mouse (Mus_musculus NP_058652.1, 3HRW|B);
 - ✓ (5) hbb_chicken (Gallus_gallus NP_990820.1, 1HBR|B);
 - ✓ (6) myoglobin (human NP_005359.1, 3RGK);
 - ✓ (7) neuroglobin (human NP_067080.1, 1OJ6|A);
 - ✓ (8) globin_soybean (Glycine max leghemoglobin A, NP_001235928.1, 1FSL); and
 - ✓ (9) globin_rice (Oryza sativa (japonica cultivar-group) NonSymbiotic Plant Hemoglobin NP_001049476.1, 1D8U).

Taller1: Secuencias

Las secuencias anteriores deben parecerse a las siguientes:

```
>hbb_human NP_000509.1
MVHLTPEEKSAVTALWGKVVNDEVGGEALGRLLVYPWTQRRFFESFGDLSTPDVAVMGNPVKVKAHGKKVLGAFSDGLA
HLDNLKGTGFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHKYH
>hbb_chimp Pan_troglodytes XP_508242
MVHLTPEEKSAVTALWGKVVNDEVGGEALGRLLVYPWTQRRFFESFGDLSTPDVAVMGNPVKVKAHGKKVLGAFSDGLA
HLDNLKGTGFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHKYH
>hbb_dog Canis_familiaris XP_537902
MVHLTAEKSLVSGLWGKVVNDEVGGEALGRLLIYPWTQRRFFDSFGDLSTPDVAVMSNAKVKAHGKKVLNSFSGLK
NLDNLKGTFAKSELHCDKLHVDPENFKLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHKYH
>hbb_mouse Mus_musculus NP_058652
MVHLTDAEKSAVSLWAKVNPDEVGGEALGRLLVYPWTQRYFDSFGDLSSASAIMGNPKVKVKAHGKKVITAFNEGLK
NLDNLKGTFAKSELHCDKLHVDPENFRLLGNAIVLGHHLGKDFTPAAQAAFQKVVAGVATALAHKYH
>hbb_chicken Gallus_gallus XP_444648
MVHWTAEKQLITGLWGKVVNVAECGAEALARLLIYPWTQRRFFASFGNLSSTAILGNPMVRAHGKKVLTSFGDAVK
NLDNIKNTFSQSELHCDKLHVDPENFRLLGDILIVLAAHFSKDFTPECQAAWQKLVVRVVAHALARKYH
>myoglobin_human 2MM1 Homo sapiens NP_005359.1
MGLSDGEWQLVLNVWVGKVEADIPGHGQEVLIIRLFKGGHPETLEKFDKFKHLKSEDEMKAEDLKKHGATVLTALGGIL
KKKGHHEAEIKPLAQSHATKHKIPVKYLEFISECIIQVLQSKHPGDFGADAQGAMNKALELFRKDMASNYKELGFQG
>neuroglobin_human 1OJ6A Homo sapiens NP_067080.1
MERPEPELIRQSWRAVRSRPLEHGTVLFARLFALPDLLPLFQYNCRQFSSPEDCLSSPEFLDHIRKVMLVIDAAVT
NVEDLSSLEEYLASLGRKHRAVGKLSFSSTVGESLLYMLEKCLGPAFTPATRAAWSQLYGAVVQAMSRGWGDE
>globin_soybean 1FSL leghemoglobin P02238 LGBA_SOYBN [Glycine max]
MVAFTKQDALVSSFEAFKANIPQYSVVFYSILEKAPAAKDLFSFLANGVDPTNPKLTGHAEKLFALVRDSAGQL
KASGTVVADAALGSVHAQKAVTDPQFVVVKEALLKTIKAAVGDKWSDELSRAWEVAYDELAIAIKKA
>globin_rice 1D8U rice Non-Symbiotic Plant Hemoglobin NP_001049476.1
MALVEDNNAVAVSFSEEQALVLKSWAILKKDSANIALRFFLKIFEVAPSASQMFSFLRNSDVPLEKNPKLKTHAMS
VFVMTCEAAQRLRKAGKVTVRDITLTKRLGATHLKYGVGDAHFEVVKFALLDTIKEEVPADMMWSPAMKSAWSEAYDHL
VAAIKQEMKPAE
```

Taller1: Trabajo Casa

- Realizar y comparar los alineamientos con los cuatro algoritmos ClustalW, MAFF, Muscle, T-Coffee
- Investigar alineamientos con estructuras 3D y realizar el anterior alineamiento utilizando la herramienta EXPRESSO
- Respondes:
 - ✓ Para usted, cuál algoritmo presenta el mejor alineamiento? Porqué?
 - ✓ Qué elementos comunes describe este grupo de globinas?
 - ✓ Existe algún bloque común?
 - ✓ Qué tipo de algoritmo realiza el mejor alineamiento y en general como lo hace?