
Análisis de secuencias. Patrones, perfiles y dominios.

Curso de verano de
Bioinformática y Biología Computacional
de la UCM, Madrid 2005

Federico Abascal
Museo Nacional de Ciencias Naturales

Recordatorio

Queremos comparar secuencias porque creemos que nos pueden hablar de la historia evolutiva de las proteínas, donde quizás podamos encontrar huellas de sus características funcionales y estructurales.

Para poder hacer esta comparación lo mejor posible: debemos encontrar el **alineamiento** que con mayor probabilidad (nunca sabremos si es el real) refleje qué cambios se han producido.

Limitación del alineamiento entre pares de secuencias

Problema: las mismas proteínas alinean de forma distinta según la matriz de sustitución y las penalizaciones por gaps utilizadas.

¿Cómo podemos saber cuál es el mejor alineamiento?

Observación: cuantas más secuencias, mayor cantidad de información, menor incertidumbre.

¿Cómo utilizar la información de muchas secuencias?

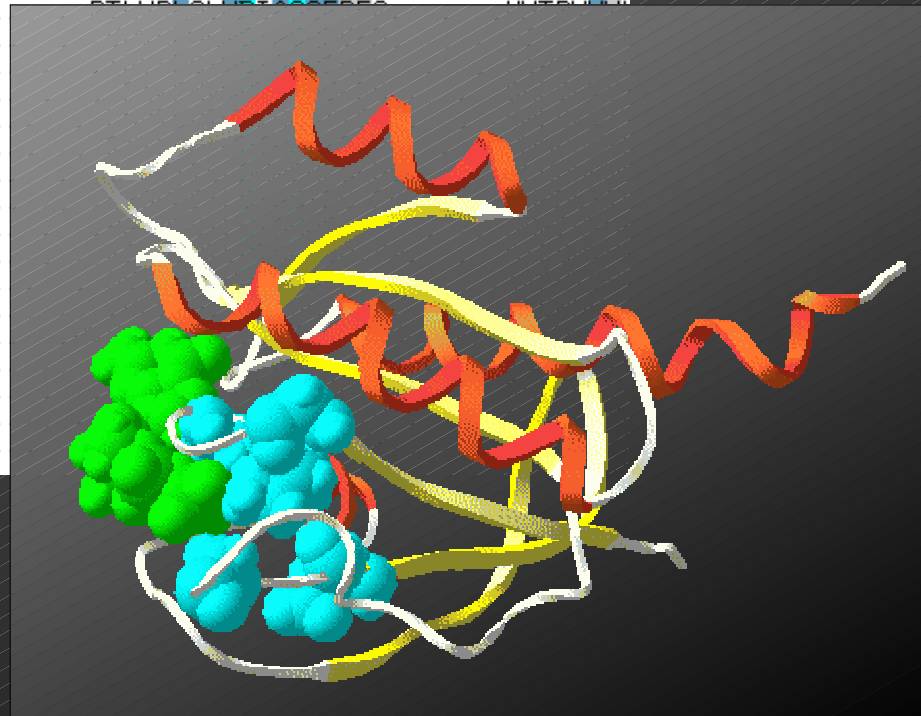
Construyendo un alineamiento múltiple.

```
# Matrix: BLOSUM62
# Gap_penalty: 10.0
# Extend_penalty: 0.5
DGHFVPNITLGQP (prot 1)
| |||.|.:::.
D-HFVDNTVFGQE (prot 2)
# Score: 296.0

# Matrix: BLOSUM45
# Gap_penalty: 10.0
# Extend_penalty: 0.5
DGHFVPN-ITLGQP (prot 1)
| |||.| :..|:.
D-HFVDNTVFGQGEH (prot 2)
# Score: 130.5
```

Alineamiento múltiple

```
NILCVGETGLGKSTLMDTLFNTKFEGEPATHTQPGVQLQSN.TYDLQES.....NVRCLKLTIVSTVGFQD.QI.....NKEDSYK  
KLLLIQDSGVQKTCVLFRFSEDAFNSTFIS..TIGIDFKIR.TIELDG.....KRIKQIWDTAGQERFR.....TITTAYY  
KLLIIQDSGVQKSSLLRFADNTFSGSYIT..TIGVDFKIR.TVEING.....EKVKLQIWDTAGQERFR.....TITSTYY  
KILIIQNSSVQKTSFLFRYADDSFTPAFVS..TVGIDFKVK.TIYRND.....KRIKQIWDTAGQERYR.....TITTAYY  
KILIIQESGVQKSSLLRFDDTDFPELAA..TIGVDFKVK.TISVDG.....NKAKLAIWDTAGQERFR.....TLTPSY  
KVVLIQDSGVQKSNLLSRFTRNEFNLESKS..TIGVEFATR.SIQVDG.....KTIKAQIWDTAGQERYR.....AITSAYY  
KFLVIGNAGTGKSCLLHQFIEKKFKDDSNH..TIGVEFGSK.IINVGG.....KYVKLQIWDTAGQERFR.....SVTRSYY  
KIIVIQDSNVGKTCLTFRFCGGTFDPKTEA..TIGVDFREK.TVEIEG.....EKIKVQVWDTAGQERFR.....SMVEHY  
KIWLIGNAGVGKTCVRRFTQGLFPPGQGA..TIGVGFMIK.TVEING.....EKVKLQIWDTAGQERFR.....SITQSY  
..MLVQDSGVQKTCVLRFRKDGAFLAGTFIS.TVGIDFRNK.VLDVDG.....VKVKLQMWDTAGQERFR.....SVTHAY  
KLWLLGSGSVQKSSALRYVKNDFKSILP...TVGCAFFTK.VVDVGA.....TSLKLEIWDTAGQEKYH.....SVCHLY  
KVCLLQDGTGVQKSSIVWRFVEDSFDPNINP..TIGASFMTK.TVQYQN.....ELHKFLIWDTAGQERFR.....ALAPMY  
KLWLLGESAVQKSSVLRFRVKGQFHEFQES..TIGAAFLTQ.TVCLDD.....TTVKFEIWDTAGQERYH.....SLAPMY  
KVVLLGEGCVGKTSVLRVYCNKFNPKHIT..TLQASFLT.KLNIGG.....KRVNLAIWDTAGQERFH.....ALGPIY  
KLWFLGQSVGKTSLITRFMYDSFDNTYQA..TIGIDFLSK.TMYLED.....RTVRLQLWDTAGQERFR.....SLIPSY  
KLLALQDSGVQKTTFLYRYTDNKFNPKFIT..TVGIDFREKRVVYNAQGPNGSSGKAFKVLQLWDTAGQERFR.....SLTTAF  
KVILLQDGGVQKSSLMNRYVTNKFDTQLFH..TIGVEFLNK.DLEVDG.....HFVT.MQIWDTAGQERFR.....SLRTPFY  
KVLVIGELGVGKTSIIKRYVHQLFSQHYRA..TIGVDFALK.VLNWDS...  
KMVVVQNGAVQKSSMIQRYCKGIFTKDYKK..TIGVDFLER.QIQVND...  
KVVVVQDLYVGKTSLIHRFCKNVFDRDYKA..TIGVDFEIE.RFEIAG...  
KLVLVQDGGTGKTTFVKRHLTGEFEKYYVA..TLGVEVHPLVFHTNRG...  
KIWVLQDGTSGKTSLTTCFAQETFGKQYKQ..TIGLDFLRRITLPGN...  
KIICLQDSAVQKSKLMEFLMDGFQPPQQLS..TYALTLYKH.TATVDG...  
RVVLIQEQGVQKSTLANIFAGVHDSMDSDC..EVLGEDTYERTLMVDG...  
KVVVLQSGGVQKSALTVQFVTGTFIEKY...DPTIEDFYRKEIEVDS...  
RLVVVQGGGVQKSALTIQFIQSYFVTDY...DPTIEDSYTKQCVIDD...  
KVIMVQSGGVQKSALTLQFMYDEFVEDY...EPTKADSYRKKVVLDG...  
KIAILQYRSVQKSSLTIQFVEGQFVDSY...DPTIENTFTKLITVNG...  
RVVVVGTAGVGKSTLLHKWASGNFRHEYLP..TIENTYCQLLGCSHG...  
RVAVLGAPGVGKTAIRQFLFGDYPERHR...PTDGPRLYRPAVLLDG...  
KCVVVQDGAVGKTCLLISYTTNKFPSYV...TVFDNYAVT..VMIGG...  
KVVLVQDGGCGKTSLLMVAFADGAFPEYTP..TVFERYMVN..LQVKG...  
KIWVVQDSQCGKTAALLHVFAKDCFPENYV...TVFENYAS..FEIDT...  
KCVLVQDSAVGKTSVLRFTSETFPEAYKP..TVYENTGVD..VFMDD...  
RTLMVQLDAACKTTIYKIKLGETVTTIP..TIGENWETVEY
```



Otra limitación de las comparaciones entre pares

Problema: si dos homólogos han divergido mucho (parecido $< 20-25\%$), BLAST no es capaz de distinguir ese parecido del azar.

BLAST no es capaz de encontrar homólogos remotos

Observación: cuando hacemos un alineam. múltiple vemos qué posiciones son más importantes.

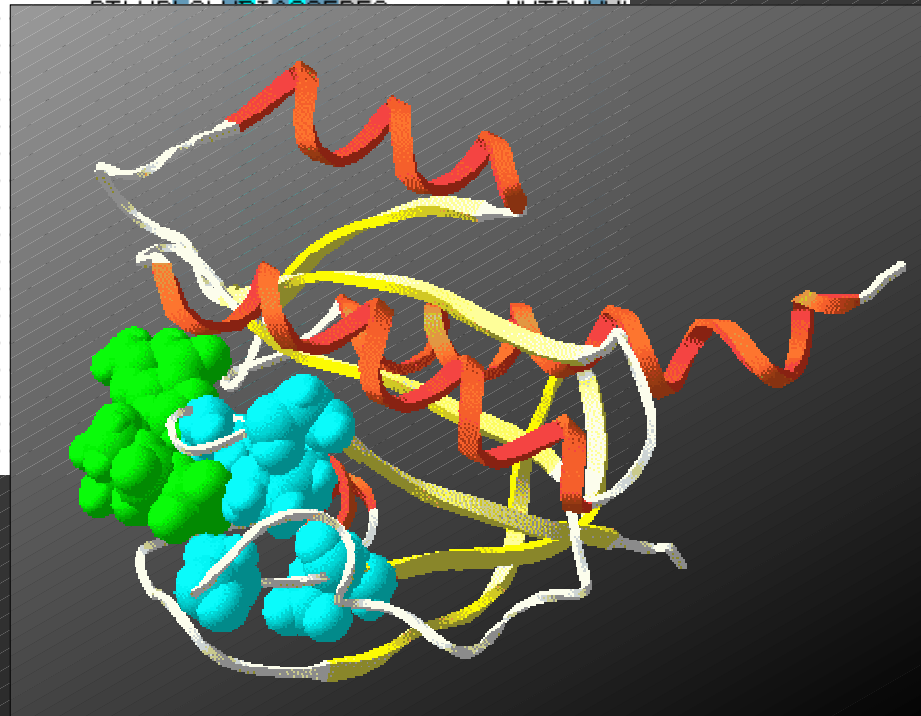
Idea: si las coincidencias en el alineamiento entre dos secuencias se producen en los sitios más importantes, la confianza en que sean homólogas ha de aumentar

Objetivo: utilizar la información de los alineam. múltiples para hacer búsquedas de homólogos más sensibles.

¿Cómo aprovechar la información del alineamiento múltiple?

Alineamiento múltiple

```
NILCVGETGLGKSTLMDTLFNTKFEGEPATHTQPGVQLQSN.TYDLQES.....NVRCLKLTIVSTVGFQD.QI.....NKEDSYK  
KLLLIQDSGVQKTCVLFRFSEDAFNSTFIS..TIGIDFKIR.TIELDG.....KRIKQIWDTAGQERFR.....TITTAYYF  
KLLIIQDSGVQKSSLLRFADNTFSGSYIT..TIGVDFKIR.TVEING.....EKVKLQIWDTAGQERFR.....TITSTYYF  
KILIIQNSSVQKTSFLFRYADDSFTPAFVS..TVGIDFKVK.TIYRND.....KRIKQIWDTAGQERYR.....TITTAYYF  
KILIIQESGVQKSSLLRFTDDTFDPELAA..TIGVDFKVK.TISVDG.....NKAKLAIWDTAGQERFR.....TLTPSYYP  
KVVLIQDSGVQKSNLLSRFTRNEFNLESKS..TIGVEFATR.SIQVDG.....KTIKAQIWDTAGQERYR.....AITSAYYF  
KFLVIGNAGTGKSCLLHQFIEKKFKDDSNH..TIGVEFGSK.IINVGG.....KYVKLQIWDTAGQERFR.....SVTRSYYF  
KIIVIQDSNVGKTCLTFRFCGGTFDPKTEA..TIGVDFREK.TVEIEG.....EKIKVQVWDTAGQERFR.....SMVEHYYP  
KIWLIGNAGVGKTCVRRFTQGLFPPGQGA..TIGVGFMIK.TVEING.....EKVKLQIWDTAGQERFR.....SITQSYYP  
..MLVQDSGVQKTCLLVRFKDGAFLAGTFIS.TVGIDFRNK.VLDVDG.....VKVKLQMWDTAGQERFR.....SVTHAYYP  
KLWLLGSGSVQKSSALRYVKNDFKSILP...TVGCAFFTK.VVDVGA.....TSLKLEIWDTAGQEKYH.....SVCHLYYP  
KVCLLQDGTGVQKSSIVWRFVEDSFDPNINP..TIGASFMTK.TVQYQN.....ELHKFLIWDTAGQERFR.....ALAPMYYP  
KLWLLGESAVQKSSVLRVFKGQFHEFQES..TIGAAFLTQ.TVCLDD.....TTVKFEIWDTAGQERYH.....SLAPMYYP  
KVVLLGEGCVGKTSVLRVYCNKFNPKHIT..TLQASFLT.KLNIGG.....KRVNLAIWDTAGQERFH.....ALGPIYYF  
KLWFLGQSVGKTSLITRFMYDSFDNTYQA..TIGIDFLSK.TMYLED.....RTVRLQLWDTAGQERFR.....SLIPSYYP  
KLLALQDSGVQKTTFLYRYTDNKFNPKFIT..TVGIDFREKRVVYNAQGPNGSSGKAFKVHLQLWDTAGQERFR.....SLTTAFFY  
KVILLQDGGVQKSSLMNRYVTNKFDTQLFH..TIGVEFLNK.DLEVDG.....HFVT.MQIWDTAGQERFR.....SLRTPFFY  
KVLVIGELGVGKTSIIKRYVHQLFSQHYRA..TIGVDFALK.VLNWDS...  
KMVVVQNGAVQKSSMIQRYCKGIFTKDYKK..TIGVDFLER.QIQVND...  
KVVVVQDLYVGKTSLIHRFCKNVFDRDYKA..TIGVDFEIE.RFEIAG...  
KLVLVQDGGTGKTTFVKRHLTGEFEKYYVA..TLGVEVHPLVFHTNRG...  
KIWVLQDGTSGKTSLTTCFAQETFGKQYKQ..TIGLDFLRRITLPGN...  
KIICLQDSAVQKSKLMERFLMDGFQPPQLS..TYALTLYKH.TATVDG...  
RVVLIQEQGVQKSTLANIFAGVHDSMDSDC..EVLGEDTYERTLMVDG...  
KVVVLQSGGVQKSALTVQFVTGTFIEKY...DPTIEDFYRKEIEVDS...  
RLVVVQGGGVQKSALTIQFIQSYFVTDY...DPTIEDSYTKQCVIDD...  
KVIMVQSGGVQKSALTLQFMYDEFVEDY...EPTKADSYRKKVVLDG...  
KIAILQYRSVQKSSLTIQFVEGQFVDSY...DPTIENTFTKLITVNG...  
RVVVVGTAGVGKSTLLHKWASGNFRHEYLP..TIENTYCYLLGCSHG...  
RVAVLGAPGVGKTAIRQFLFGDYPERHR...PTDGPRLYRPAVLLDG...  
KCVVVQDGAVGKTCLLISYTTNKFPSYVVP..TVFDNYAVT..VMIGG...  
KVVLVQDGGCGKTSLLMVFADGAFPEYTP..TVFERYMVN..LQVKG...  
KIWVVQDSQCGKTAALLHVFAKDCFPENYVP..TVFENYAS..FEIDT...  
KCVLVQDSAVGKTSLLVRFSETFPEAYKP..TVYENTGVD..VFMDD...  
RTLMVQLDAACKTTTLYKLYKLGSETVTTIP..TIGENWETVEY
```



Guión de la charla. Patrones, perfiles y dominios.

-cómo utilizar la información de los alineamientos múltiples

- secuencias consenso y expresiones regulares
- perfiles y perfiles-hmm

-algunas bases de datos de patrones y perfiles:

- Prosite
- Pfam

-búsquedas en bases de datos:

- PSI-BLAST
- HMMer
- búsqueda con secuencias intermedias

Guión de la charla. Patrones, perfiles y dominios.

-cómo utilizar la información de los alineamientos múltiples

-secuencias consenso y expresiones regulares

-perfiles y perfiles-hmm

-algunas bases de datos de patrones y perfiles:

-Prosite

-Pfam

-búsquedas en bases de datos:

-PSI-BLAST

-HMMer

-búsqueda con secuencias intermedias

Definición de motivo

```
NILCVHETGLGK...TMDTLFNTK...FEGEPATHTQP...VQLQSN, TYDLQES,.....NVRLKLTIVSTVGF...D, QI.....NKEDSKF  
KLLLI...GDSGVGK...T...VLFRRFSEDAFNSTFIS...TIGIDFKIR, TIELDG,.....KRIK...IWIWDTAGQ...RFR,.....TITTAYYF  
KLLI...GDSGVGK...S...LLRFADNTFSGSYIT...TIGVDFKIR, TVEING,.....EKVK...IWIWDTAGQ...RFR,.....TITSTYYF  
KILII...GNSVVGK...T...FLFRYADDSFTPAFVS...TVGIDFKVK, TIYRND,.....KRIK...IWIWDTAGQ...EYR,.....TITTAYYF  
KILII...GESGVGK...S...LLRFTDDT...DPPELAA...TIGVDFKVK, TISVDA,.....NKAK...IWIWDTAGQ...EYR,.....TLTPSYF  
KVV...IGDSGVGK...S...LSRFRTRNEFNLESKS...TIGVEFATR, SIQVDA,.....KTIK...IWIWDTAGQ...EYR,.....AITSAYYF  
KFL...IGNAGTGK...S...LHQFIEKKFKDDSNH...TIGVEFGSK, IINVGG,.....KYVK...IWIWDTAGQ...EYR,.....SVTRSYF  
KII...IGDSNVGK...T...C...TFRFCCGGT...FPDKTEA...TIGVDFREK, TVEIEG,.....EKIK...VQVWDTAGQ...EYR,.....SMVEHYF  
KIV...IGNAGVGK...T...C...VRRFTQGL...PPGQGA...TIGVGFMIK, TVEING,.....EKV...IWIWDTAGQ...EYR,.....SITQSYF  
...NLVGD...SGV...GK...T...C...L...VRFKDGAF...LAGTFIS...TVGIDFRNK, VLDVDA,.....VKV...LQMWDTAGQ...EYR,.....SVTHAYF  
KVL...LLG...SGV...GK...S...L...R...RYVKNDF...KSILP...TVGCAFFTK, VVDVGA,.....TSL...IWIWDTAGQ...EYH,.....SVCHLYF  
KVL...LLG...D...T...G...V...GK...S...S...I...W...RFV...EDS...F...DPNINP...TIGASFMTK, TVQYQN,.....ELH...FLIWDTAGQ...EYR,.....ALAPMYF  
KVL...LLG...S...A...V...GK...S...S...L...R...F...V...K...Q...F...H...E...F...Q...E...S...TIGAAFLTQ, TVCLDD,.....TTW...FEIWDTAGQ...EYH,.....SLAPMYF  
KVL...LLG...E...G...C...V...GK...T...S...L...R...Y...C...E...N...K...F...N...D...K...H...I...T...TLQASFLT, KLNIGG,.....KRW...IWIWDTAGQ...EYH,.....ALGPIYF  
KVL...FLG...E...Q...S...V...GK...T...S...L...T...R...F...M...Y...D...S...F...D...N...T...Y...Q...A...TIGIDFLSK, TMYLED,.....RTV...LQLWDTAGQ...EYR,.....SLIPSYF  
KVL...ALG...D...S...G...V...GK...T...T...F...Y...R...Y...T...D...N...K...F...N...P...K...F...I...T...TVGIDFREKRVVYNAQGPNGSSGKAFK...V...H...L...Q...L...W...D...T...A...G...Q...E...Y...R,.....SLTTAFFY  
KVL...LLG...D...G...G...V...GK...S...S...L...N...R...Y...V...T...N...K...F...D...T...Q...L...F...H...TIGVEFLNK, DLEVDG,.....HFVT...MQIWDTAGQ...EYR,.....SLRTPFYF  
KVL...VIGELGVGK...T...S...I...K...R...Y...V...H...Q...L...F...S...Q...H...Y...R...A...TIGVDFALK, VLNWDS,.....RTL...V...L...Q...L...W...D...I...A...G...Q...E...Y...R,.....NMTRVYF  
KVL...V...V...G...N...G...A...V...GK...S...S...M...Q...R...Y...C...K...G...I...F...T...K...D...Y...K...K...TIGVDFLER, QIQVND,.....EDW...L...M...L...W...D...T...A...G...Q...E...E...F...D,.....AITKAYF  
KVL...V...V...G...D...L...Y...V...GK...T...S...L...H...R...F...C...K...N...V...F...D...R...D...Y...K...A...TIGVDFEIE, RFEIAG,.....IPY...L...Q...I...W...D...T...A...G...Q...E...K...F...K,.....CIASAYF  
KVL...LVG...D...G...G...T...G...K...T...T...F...K...R...H...L...T...G...E...F...E...K...K...Y...V...A...TIGVEVHPLVFHTNRG,.....PI...F...N...V...W...D...T...A...G...Q...E...K...F...K,.....GLRDGYF  
KVL...V...L...G...D...G...T...S...G...K...T...S...L...T...C...F...A...Q...E...T...F...G...K...Q...Y...K...Q...TIGLDFLRRITLPGN,.....LNV...L...Q...I...W...D...I...G...G...Q...T...I...G,.....KMLDKYF  
KVL...I...C...L...G...D...S...A...V...GK...S...K...L...H...E...R...F...L...M...D...G...F...Q...P...Q...L...S...TYALTLYKH, TATVDA,.....RTI...V...D...F...W...D...T...A...G...Q...E...R...F...Q,.....SMHASYF  
RVL...L...I...G...E...Q...G...V...GK...S...T...L...A...N...I...F...A...G...V...H...D...S...M...D...S...D...C...EVLGEDTYERTLMVDG,.....ESA...I...L...L...D...M...W...E...N...K...G...E...N...E,.....WLHDHCF  
RVL...V...L...G...S...G...G...V...GK...S...A...L...T...V...Q...F...V...T...G...T...F...I...E...K...Y...DPTIEDFYRKEIEVDS,.....SPS...L...E...I...L...D...T...A...G...T...E...Q...F...A,.....SMRDLYF  
RVL...V...V...G...G...G...V...GK...S...A...L...T...I...Q...F...I...Q...S...Y...F...V...T...D...Y...DPTIEDSYTKQCVIDD,.....RAA...D...I...L...D...T...A...G...Q...E...E...F...G,.....AMREQMYF  
KVL...V...V...G...S...G...G...V...GK...S...A...L...T...L...Q...F...M...Y...D...E...F...V...E...D...Y...EPTKADSYRKKVLDG,.....EEV...D...I...L...D...T...A...G...Q...E...D...Y...A,.....AIRDNYF  
KVL...A...L...G...Y...R...S...V...GK...S...S...T...I...Q...F...V...E...G...Q...F...V...D...S...Y...DPTIENTFTKLITVNG,.....QEY...H...L...V...D...T...A...G...Q...E...Y...S,.....IFPQYSI  
RVL...V...V...G...T...A...G...V...GK...S...T...L...L...H...K...W...A...S...G...N...F...R...H...E...Y...L...P...TIENTYQQLLGCSHG,.....VLS...H...I...T...D...S...K...S...G...D...N...R,.....ALQRHWI  
RVL...V...A...P...G...V...GK...T...F...I...R...Q...F...L...F...G...D...Y...P...E...R...H...R...PTDGPRLYRPVALLDG,.....AVY...D...L...S...R...D...G...V...A...G...P...G...S...P...G...G...P...E...E...W...P...D...A...K...D...W...S...L...C  
KVL...V...G...D...G...A...V...GK...T...L...L...I...S...Y...T...T...N...K...F...P...S...E...Y...V...P...TVFDNYAVT, VMIGG,.....EPY...T...L...G...L...F...D...T...A...G...Q...E...Y...D,.....RLRPLSYF  
KVL...V...L...G...D...G...C...G...K...T...L...L...M...V...F...A...D...G...A...F...P...E...S...Y...T...P...TVFERYMVN, LQVKG,.....KPV...H...L...I...W...D...T...A...G...Q...D...Y...D,.....RLRPLFYF  
KVL...V...V...D...S...Q...C...G...K...T...A...L...L...H...V...F...A...K...D...C...F...P...E...N...Y...V...P...TVFENYAS, FEIDT,.....QRI...E...L...L...W...D...T...S...S...Y...Y...D,.....NVRPLSYF  
KVL...V...D...S...A...V...GK...S...L...L...V...R...F...T...S...E...T...P...E...A...Y...K...P...TVYENTGVD, VFMDG,.....IQI...S...L...L...W...D...T...A...G...N...A...F...R,.....SIRPLSYF  
RVL...M...V...D...A...A...G...K...T...T...I...Y...K...I...K...L...G...E...T...V...T...T...P...TIGENVETVEY,.....KNT...S...E...T...N...W...A...G...T...D...K...T...R,.....PLWRHYF
```

Son pequeñas zonas conservadas.

Se suelen corresponder con características funcionales de las proteínas:

- centros activos
- sitios de unión de ligandos
- etc

Motivos

Secuencias consenso y patrones

¿Cómo aprovechar la información del alineamiento múltiple?

-Secuencias consenso:

```
AGTVATVSC
AGTSATHAC
IGRCARGSC
IGEMARLAC
IGDYARWSC
.....
```

IGTVARVSC <= Ejemplo de secuencia consenso

-Patrones o expresiones regulares:

(para caracterizar motivos)

```
ALRDFATHDDF
SMTAEATHDSI
ECDQAATHEAS
```



A-T-H-[DE]

Patrones (expresiones regulares)

¿Cómo expresarse regularmente?

- Cualquier aminoácido: x
- Ambigüedad:
 - $[A,B]$ A, o B...
 - $\{A,B..\}$ cualquiera menos A y B.
- Repetición: $A(2,4)$ significa A-A o A-A-A o A-A-A-A
- N terminal: $<$, C-terminal: $>$

Ejemplo: $[AC]-x-V-x(4)-\{E,D\}$.

[Ala or Cys]-any-Val-any-any-any-
any-{any but Glu or Asp}

Patrones: un ejemplo

Ejemplo:

AGTVATVSC

AGTSATHAC

IGRCARGSC

IGEMARLAC

IGDYARWSC

.....

IGTVARVSC



Ejemplo de secuencia consenso

[AI]-G-X-X-A-[RT]-[SA]-C



Ejemplo de patrón

Construcción de un patrón:

- Más o menos subjetivo. Ensayo y error.
- Objetivo: alta sensibilidad, alta especificidad



Debemos construirlos en torno a motivos conservados.

Suficientemente cortos (sensibilidad),
suficientemente largos (especificidad)

Patrones (III)

Ventajas y desventajas de los patrones

-Su construcción es bastante laboriosa, pero...

existen algunos métodos automáticos (PRATT: <http://www.ebi.ac.uk/pratt/>)

existen bases de datos donde expertos hacen ese trabajo por nosotros

(Prosite: <http://www.expasy.org/prosite>)

-Muy estrictos.

Básicamente distingue posiciones importantes y no importantes (con 'X'), pero en la Naturaleza hay una mayor gradación.

Si una proteína nueva se sale de la regla general, no será detectada con el patrón.

Guión de la charla. Patrones, perfiles y dominios.

-cómo utilizar la información de los alineamientos múltiples

-secuencias consenso y expresiones regulares

-perfiles y perfiles-hmm

-algunas bases de datos de patrones y perfiles:

-Prosite

-Pfam

-búsquedas en bases de datos:

-PSI-BLAST

-HMMer

-búsqueda con secuencias intermedias

Dominios

“Unidad estructural independiente”, en otras áreas se le da un sentido diferente (en estudios genéticos de delección a veces se utiliza como sinónimo de la parte mínima de la secuencia capaz de realizar la función).

Se muestra el antígeno de Histocompatibilidad de clase I: dominios $\alpha_1, 2$ y 3 y proteína beta-2-microglobulina.

¿dos dominios o uno?

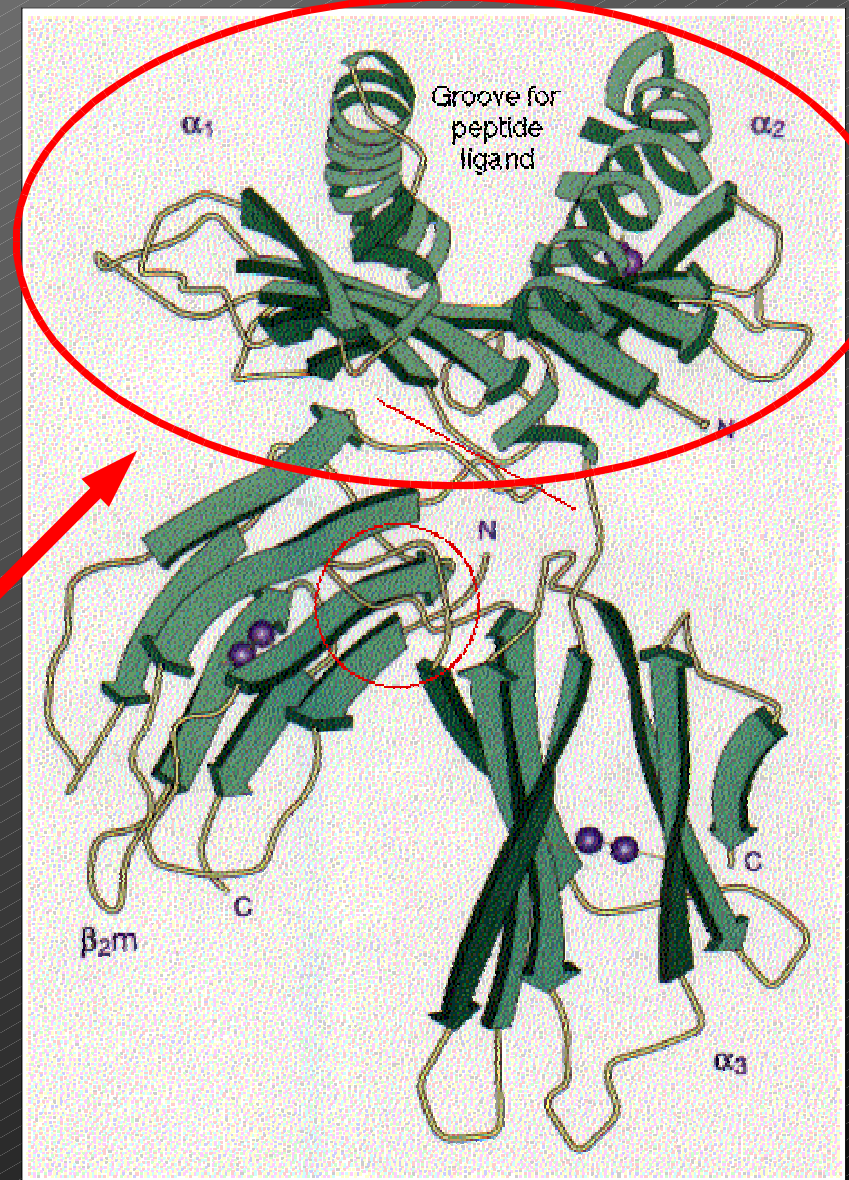


Imagen tomada de:
<http://users.rcn.com/jkimball.ma.ultranet/BiologyPages/T/TertiaryStructure.html>, que a su vez la tomó de P. J. Bjorkman from Nature 329:506, 1987

Perfiles

Si queremos alinear la secuencia FKTLGCCLLV al perfil, el mejor alineamiento será:


```

F K L L S H C L L V
F K A F G Q T M F Q
Y P I V G Q E L L G
F P V V K E A I L K
F K V L A A V I A D
L E F I S E C I I Q
F K L L G N V L V C
F K T L G C C L L V
  
```

Y la puntuación:

60 25 -6 27 28 -26 22 33 26 -16

Lo cual suma en total: **173.**

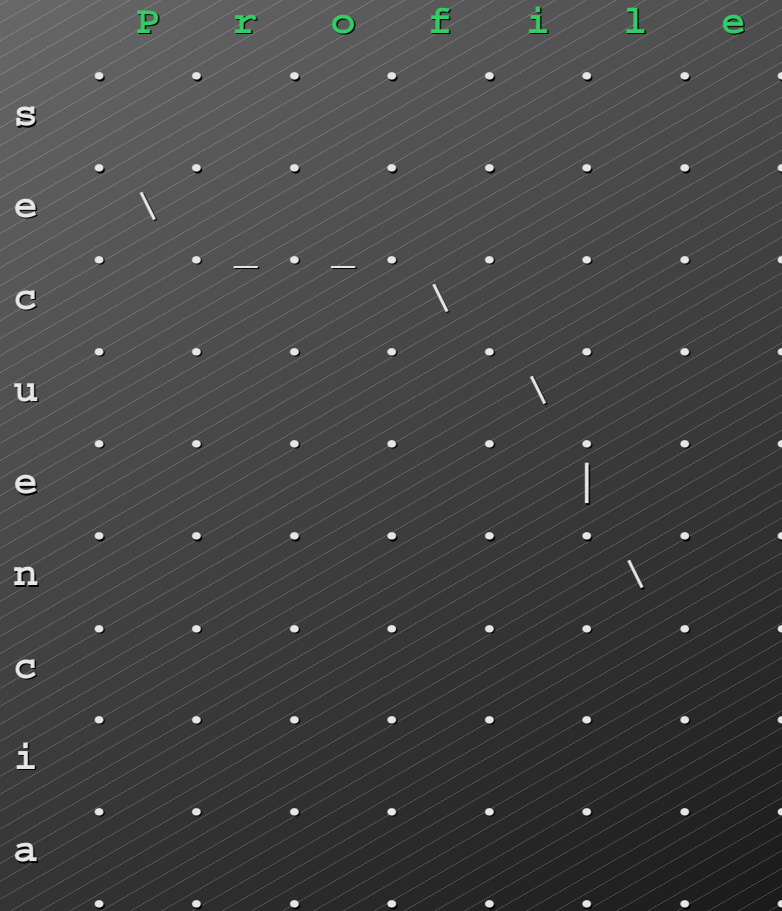


	F	K	L	L	S	H	C	L	L	V
F	60	-30	12	14	-26	-29	-15	4	12	-29
K	-26	25	-25	-27	-6	4	-15	-27	-26	0
L	14	-28	19	27	-27	-20	-9	33	26	-21
L	14	-28	19	27	-27	-20	-9	33	26	-21
S	-22	-8	-16	-21	11	2	-1	-24	-19	-4
H	-13	-12	-25	-25	-16	14	-22	-22	-23	-10
C	-22	-33	-18	-18	-22	-26	22	-24	-19	-7
L	3	-15	10	14	-17	-10	-9	25	12	-11
L	3	-15	10	14	-17	-10	-9	25	12	-11
V	0	-25	22	25	-19	-26	6	19	16	-16
A	-18	-10	-1	-8	8	-3	3	-10	-2	-8
C	-22	-33	-18	-18	-22	-26	22	-24	-19	-7
D	-35	0	-32	-33	-7	6	-17	-34	-31	0
E	-27	15	-25	-26	-9	23	-9	-24	-23	-1
F	60	-30	12	14	-26	-29	-15	4	12	-29
G	-30	-20	-28	-32	28	-14	-23	-33	-27	-5
I	3	-27	21	25	-29	-23	-8	33	19	-23
K	-26	25	-25	-27	-6	4	-15	-27	-26	0
L	14	-28	19	27	-27	-20	-9	33	26	-21
M	3	-15	10	14	-17	-10	-9	25	12	-11
N	-22	-6	-24	-27	1	8	-15	-24	-24	-4
P	-30	24	-26	-28	-14	-10	-22	-24	-26	-18
Q	-32	5	-25	-26	-9	24	-16	-17	-23	7
R	-18	9	-22	-22	-10	0	-18	-23	-22	-4
S	-22	-8	-16	-21	11	2	-1	-24	-19	-4
T	-10	-10	-6	-7	-5	-8	2	-10	-7	-11
V	0	-25	22	25	-19	-26	6	19	16	-16
W	9	-25	-18	-19	-25	-27	-34	-20	-17	-28
Y	34	-18	-1	1	-23	-12	-19	0	0	-18

Perfiles (II)

¿Cómo utilizar un perfil para buscar homólogos?

El mismo algoritmo usado para alinear dos secuencias (Smith & Waterman) sirve para alinear una secuencia y un perfil.

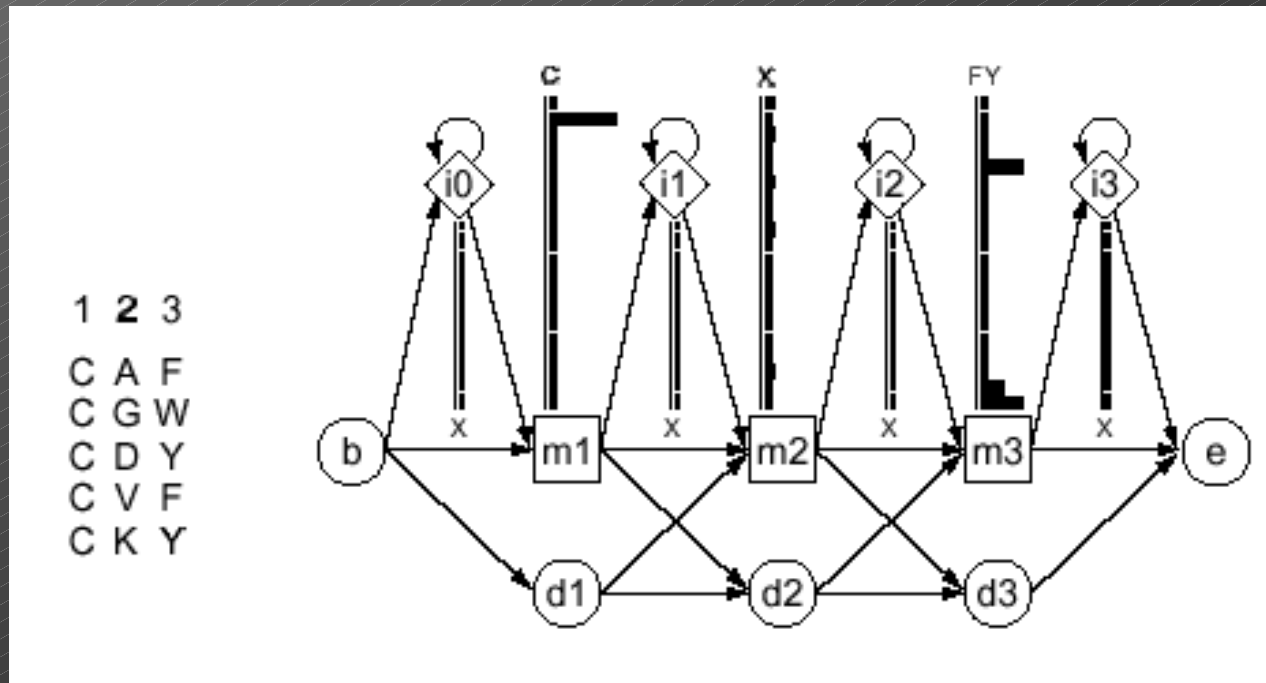


Perfiles de tipo HMM

Perfiles de tipo HMM (hidden markov model)

La base probabilística de los perfiles simples es pobre, especialmente en cuanto a la penalización de gaps.

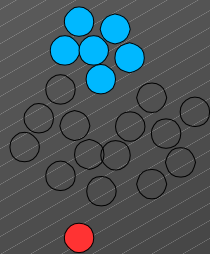
Los HMM son más sólidos (y complejos)



Perfiles (III)

Posible problema de los perfiles: sesgo en la representación de las secuencias.

Solución: asignar distintos “pesos” a las secuencias en función de cuánta información extra aporten al perfil.



ATCILYC
ATCILYC
MTCIRYC
ATCDLYC
ATCILYC
ATCILYC
SRCDRMC

Guión de la charla. Patrones, perfiles y dominios.

-cómo utilizar la información de los alineamientos múltiples

- secuencias consenso y expresiones regulares
- perfiles y perfiles-hmm

-algunas bases de datos de patrones y perfiles:

- Prosite
- Pfam

-búsquedas en bases de datos:

- PSI-BLAST
- HMMer
- búsqueda con secuencias intermedias

Prosite (I)

PROSITE:

<http://us.expasy.org/prosite/>

-caracterizan motivos
conocidos con
expresiones regulares y/o
perfiles.

-gran cantidad de
información para cada
familia de proteínas.

-baja cobertura: sólo 1.245
familias

```
ID MOLYBDOPTERIN_EUK; PATTERN.
AC PS00559;
DT DEC-1991 (CREATED); NOV-1995 (DATA UPDATE); JUL-1998 (INFO UPDATE).
DE Eukaryotic molybdopterin oxidoreductases signature.
PA [GA]-x(3)-[KRNQHT]-x(11,14)-[LIVMFYWS]-x(8)-[LIVMF]-x-C-x(2)-[DEN]-R-
PA x(2)-[DE].
NR /RELEASE=38,80000;
NR /TOTAL=50(50); /POSITIVE=45(45); /UNKNOWN=0(0); /FALSE_POS=5(5);
NR /FALSE_NEG=2; /PARTIAL=5;
CC /TAXO-RANGE=??E??; /MAX-REPEAT=1;
DR P48034, ADO_BOVIN , T; Q06278, ADO_HUMAN , T; P11832, NIA1_ARATH , T;
DR P39867, NIA1_BRANA , T; P27967, NIA1_HORVU , T; P16081, NIA1_ORYSA , T;
DR P39865, NIA1_PHAVU , T; P54233, NIA1_SOYBN , T; P11605, NIA1_TOBAC , T;
DR P11035, NIA2_ARATH , T; P39868, NIA2_BRANA , T; P27969, NIA2_HORVU , T;
DR P39866, NIA2_PHAVU , T; P39870, NIA2_SOYBN , T; P08509, NIA2_TOBAC , T;
DR P49102, NIA3_MAIZE , T; P27968, NIA7_HORVU , T; P36858, NIA_ASPNG , T;
DR P43100, NIA_BEABA , T; P27783, NIA_BETVE , T; P43101, NIA_CICIN , T;
DR P17569, NIA_CUCMA , T; P22945, NIA_EMENI , T; P39863, NIA_FUSOX , T;
DR P36842, NIA_LEPMC , T; P39869, NIA_LOTJA , T; P17570, NIA_LYCES , T;
DR P08619, NIA_NEUCR , T; P36859, NIA_PETHY , T; P49050, NIA_PICAN , T;
DR P23312, NIA_SPIOL , T; Q05531, NIA_USTMA , T; P36841, NIA_VOLCA , T;
DR P07850, SUOX_CHICK , T; P51687, SUOX_HUMAN , T; Q07116, SUOX_RAT , T;
DR P80457, XDH_BOVIN , T; P08793, XDH_CALVI , T; P47990, XDH_CHICK , T;
DR P10351, XDH_DROME , T; P22811, XDH_DROPS , T; P91711, XDH_DROSU , T;
DR P47989, XDH_HUMAN , T; Q00519, XDH_MOUSE , T; P22985, XDH_RAT , T;
DR P80456, ADO_RABIT , P; P17571, NIA1_MAIZE , P; P39871, NIA2_MAIZE , P;
DR Q01170, NIA_CHLVU , P; P39882, NIA_LOTTE , P;
DR P39864, NIA_PHYIN , N; Q12553, XDH_EMENI , N;
DR P27034, BGLS_AGRU , F; P03598, COAT_TOBSV , F; P19235, EPOR_HUMAN , F;
DR P20054, PYR1_DICDI , F; Q23316, YHC6_CAEEL , F;
3D 1SOX;
DO PDOC00484;
//
```

Prosite (II)

Lo que podemos hacer con Prosite:

- buscar con una secuencia para ver si se parece a alguno de los patrones o perfiles descritos en la base de datos.
- encontrar información de una familia determinada y ver qué proteínas pertenecen a ella.
- buscar con uno de los patrones o perfiles descritos contra una base de datos de secuencias.
- etcétera

Pfam (I)

Pfam: <http://www.sanger.ac.uk/Pfam/>

- caracterizan dominios de proteínas con perfiles HMM.
- gran cantidad de información.
- alta cobertura (7.316 familias, 73% swiss-prot y TrEMBL)



Rick:



Caspasa 9:

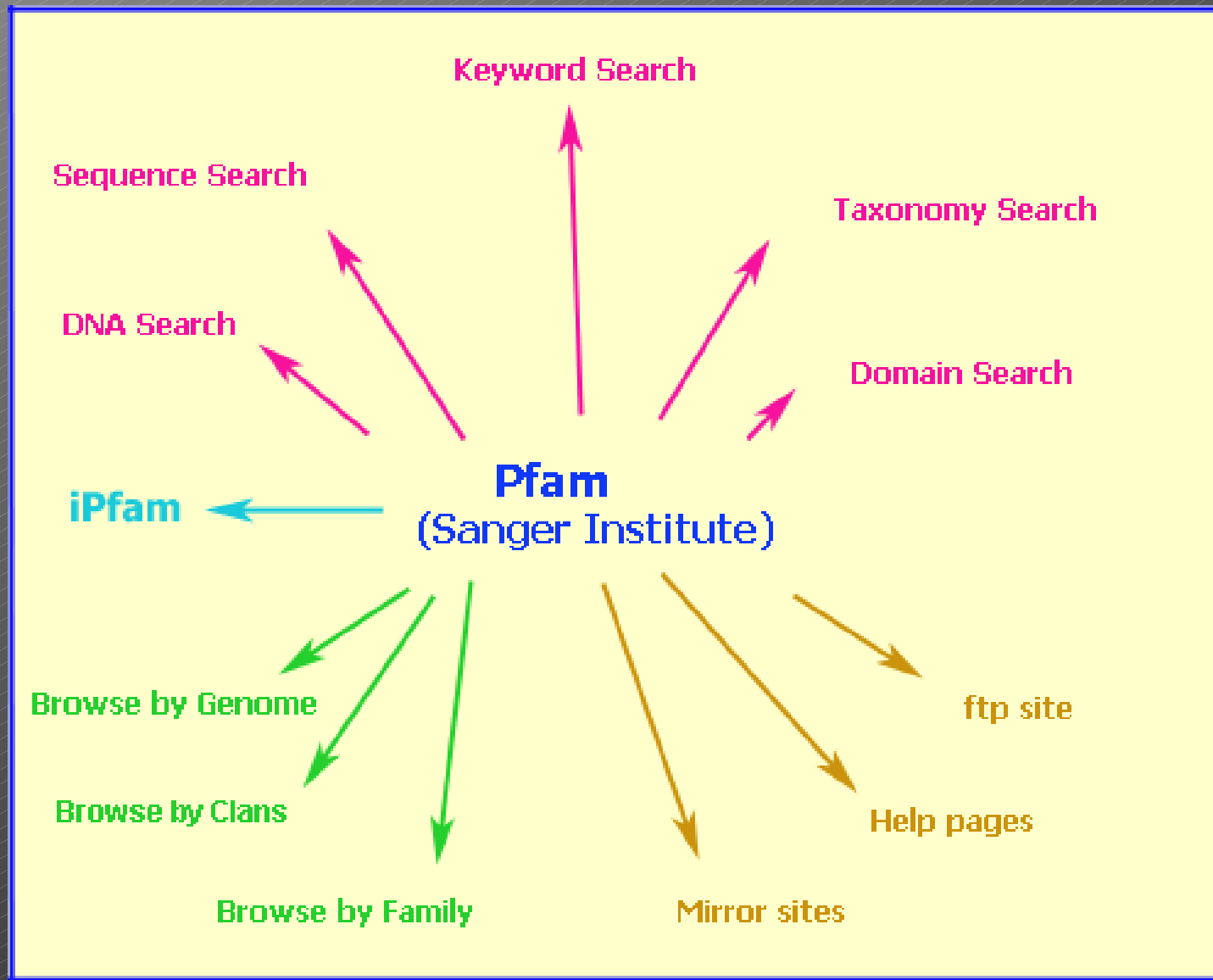


-Clasifican dominios y no proteínas completas (el dominio es la unidad evolutiva básica)

-Interfaz web muy útil:

- alineamientos
- distribución filogenética
- organización de dominios
- búsqueda usando perfiles-hmm
- etc.

Lo que podemos hacer con Pfam



Guión de la charla. Patrones, perfiles y dominios.

-cómo utilizar la información de los alineamientos múltiples

- secuencias consenso y expresiones regulares
- perfiles y perfiles-hmm

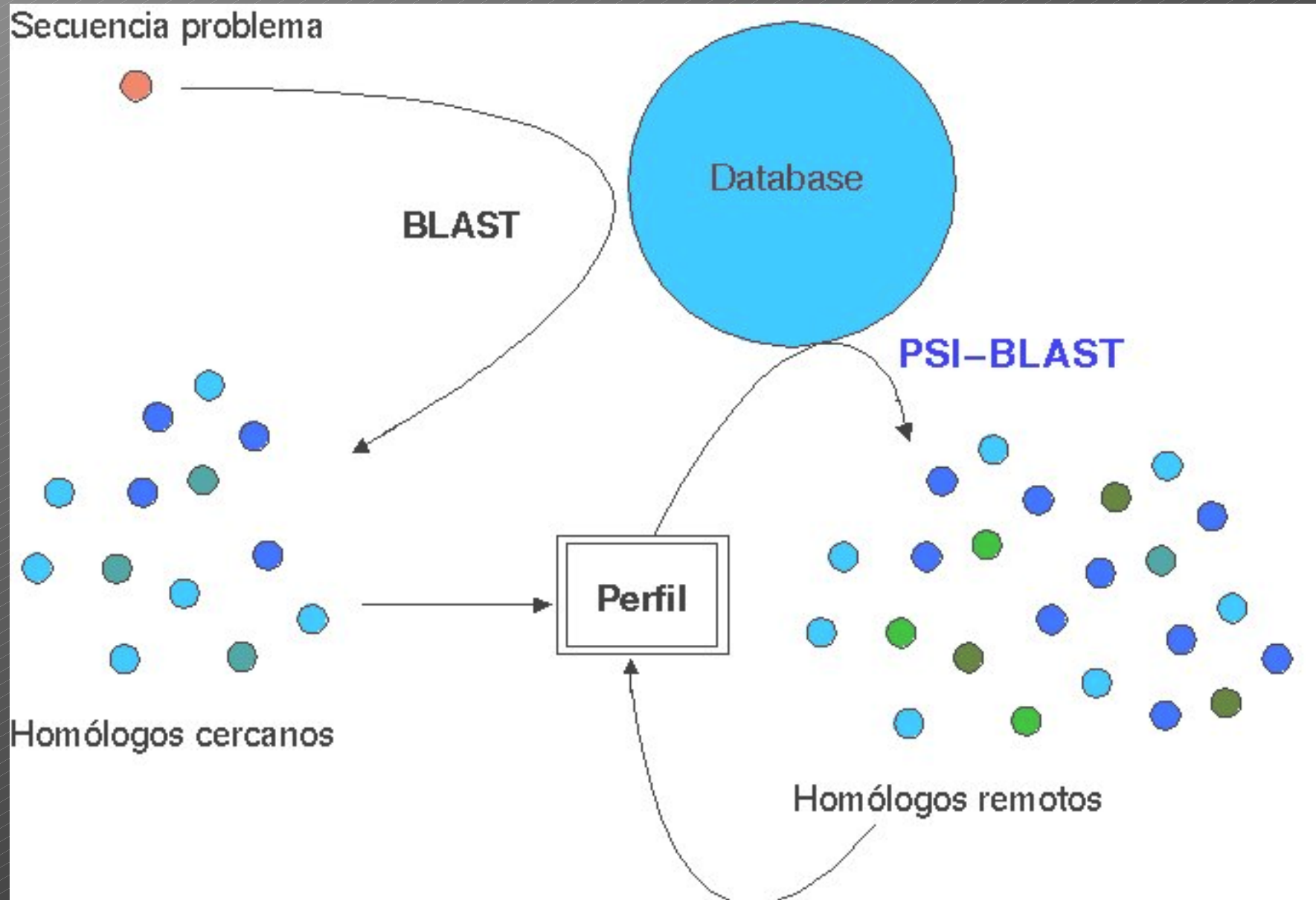
-algunas bases de datos de patrones y perfiles:

- Prosite
- Pfam

-búsquedas en bases de datos:

- PSI-BLAST
- HMMer
- búsqueda con secuencias intermedias

Búsqueda de homólogos con PSI-BLAST



Búsqueda de homólogos con PSI-BLAST

Demostración del funcionamiento de PSI-BLAST.

Página de PSI-BLAST:

<http://www.ncbi.nlm.nih.gov/BLAST/>

Secuencia de:

>gi|2501594|sp|Q57997|Y577_METJA PROTEIN MJ0577

MSVMYKKILYPTDFSETAEIALKHVKAFKTLKAEVILLHVIDEREIKKRDIFSLLLGVAGLNKSVEEFE
NELKNKLTEEAKNKMENIKKELEDVGFKVKDIIVVGIPHEEIVKIAEDEGVDIIMGSHGKTNLKEILLG
SVTENVIKKSINKPVLVVKRNS

(es el ejemplo que se sigue en el tutorial del NCBI:
<http://www.ncbi.nlm.nih.gov/Education/BLASTinfo/psi1.html>)

Búsqueda de homólogos con HMMer

<http://hmmer.wustl.edu/>

Es el método más sensible, pero es muy lento.

Requiere demasiados recursos, por lo que no se puede utilizar a través de la web.



Strategia básica:

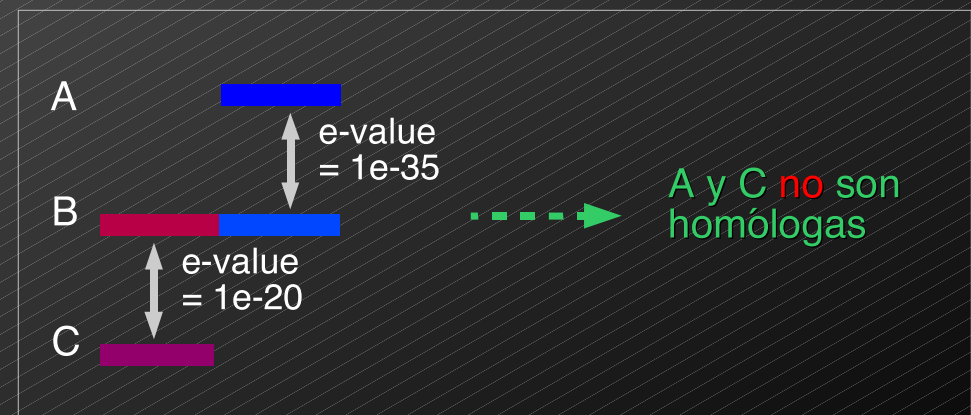
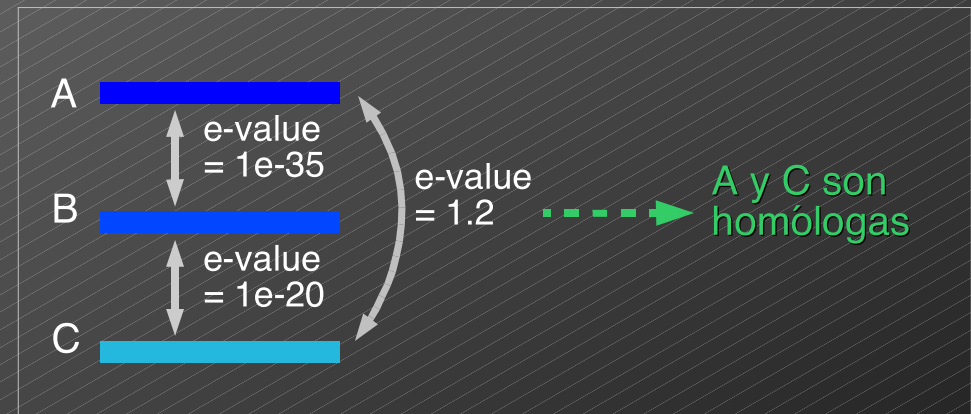
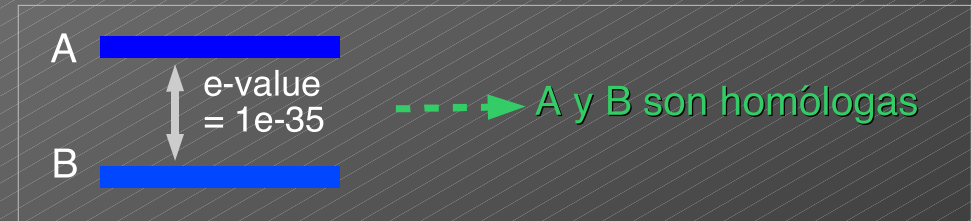
- 1.- obtener homólogos (p.e. con **BLAST**) y construir un alineamiento múltiple (p.e. con **Clustalw**).
- 2.- transformar el alineam. múltiple en un perfil HMM:
1º: **hmmbuild**, 2º **hmmcalibrate**.
- 3.- búsqueda con el perfil HMM en una base de datos de secuencias: **hmmsearch**.
- 4.- con los nuevos homólogos que encontremos podemos volver al paso "2".

Búsqueda con secuencias intermedias

No utiliza información del alineamiento múltiple, pero puede superar las limitaciones de métodos sencillos como BLAST.

Propiedad transitiva de la homología:
si dos proteínas A y B son homólogas, y a su vez B y C son homólogas, entonces A y C también son homólogas, aunque no se parezcan entre sí.

La propiedad transitiva sólo es aplicable cuando los dominios de las proteínas se corresponden unos con otros (“la unidad evolutiva son los dominios”).



Guión de la charla. Patrones, perfiles y dominios.

- cómo utilizar la información de los alineamientos múltiples

 - secuencias consenso y expresiones regulares

 - perfiles y perfiles-hmm

- algunas bases de datos de patrones y perfiles:

 - Prosite

 - Pfam

- búsquedas en bases de datos:

 - PSI-BLAST

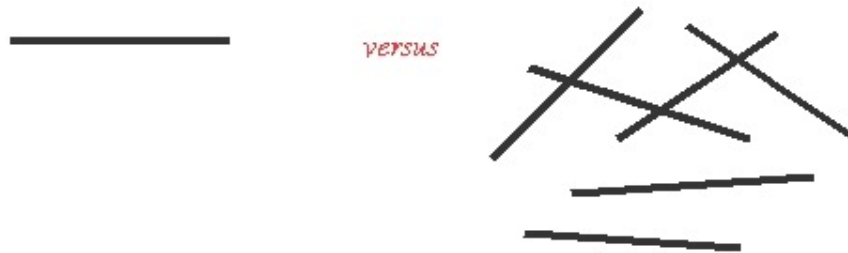
 - HMMer

 - búsqueda con secuencias intermedias

Formas de comparar secuencias (I)

1 secuencia contra una base de datos de secuencias.

BLAST(web/local), FASTA(web/local), Smit & Waterman(web/local)



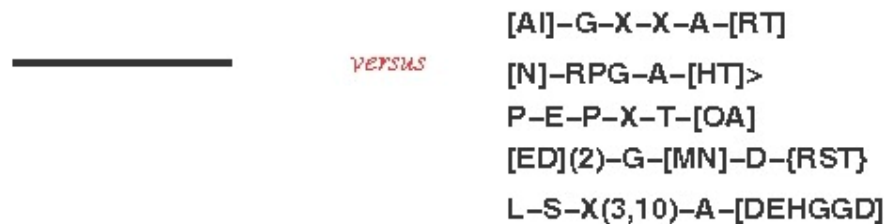
1 patrón contra una base de datos de secuencias.

ScanProsite (web), ps_scan(local)



1 secuencia contra una base de datos de patrones.

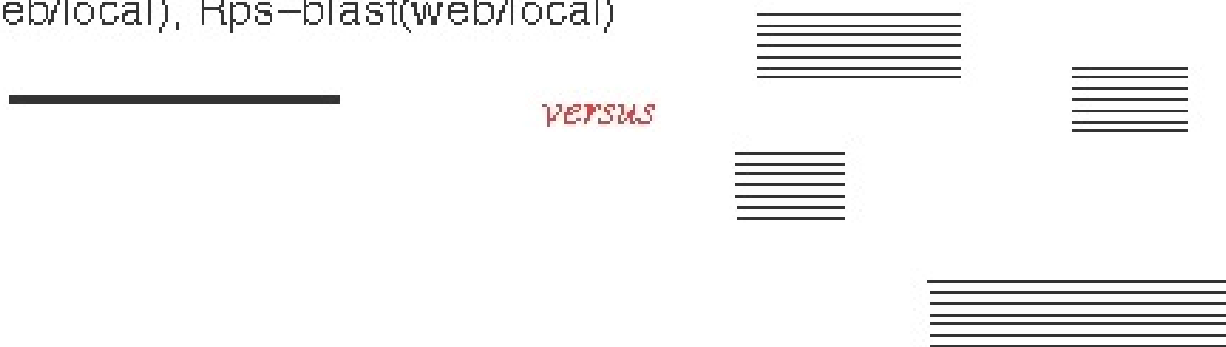
ScanProsite(web), ps_scan(local), MotifScan(web)



Formas de comparar secuencias (y II)

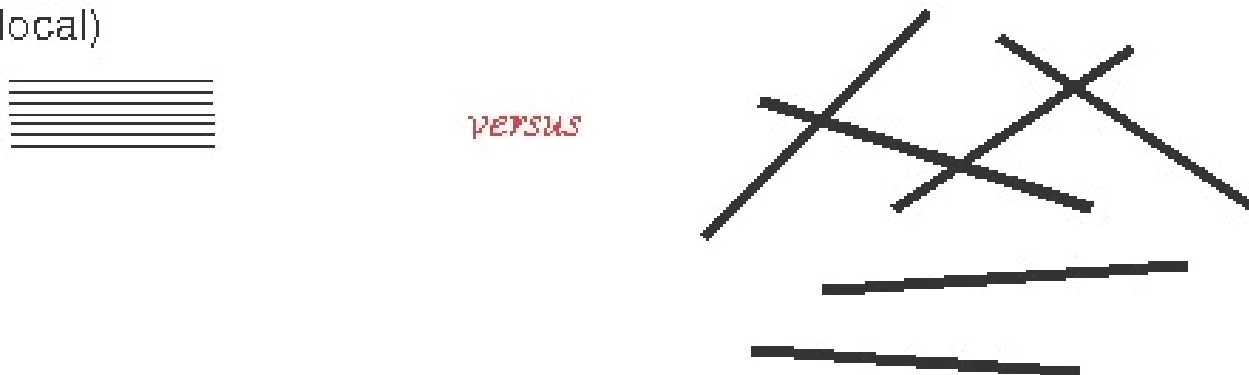
1 secuencia contra una base de datos de perfiles o HMMs.

ScanProsite(web), MotifScan(web), Pfam (hmmpfam)(web/local),
Impala(web/local), Rps-blast(web/local)



1 perfil o un perfil-HMM contra una base de datos de secuencias

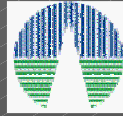
Bioccelerator (profileSearch)(web), hmmsearch(local), PSI-
BLAST(web/local)



Agradecimientos

Algunas figuras han sido tomadas de...

-Paulino Gómez Puertas



Centro de Astrobiología

-Oswaldo Trelles



Arquitectura de Computadores
Universidad de Málaga

-Joaquín Dopazo



Bioinformatics Unit
CNIO