

Schedule Week 3

4 – 6 September

- ***Wednesday***
 - Quiz
 - ***Alignments /Scoring Systems***
- ***Friday –***
 - Searching for homologs

Scoring Systems

PAM Matrices

- ***PAM means "percent accepted mutations"***
 - accepted means fixed in the population and is therefore a more complex process than simply mutation
- ***PAM-1 therefore is a scoring system for sequences in which 1% of the residues have undergone mutation***
- ***PAM-250 represents 250% mutation, i.e., an average of 2.5 accepted mutation per residue - a very distant relationship***
- ***PAM tries to model what happens at long evolutionary distances based on a simple Markov model derived from closely related sequences.***

Scoring Systems

PAM Matrices

- ***Accepted point mutations - tabulate actual mutations by looking at proteins that are sufficiently closely related that there is no ambiguity in alignment***
 - Relying on actual observed mutations is why we call it empirical
 - Sequences no more than 15% different so that changes can be thought of as a single evolutionary step – no multiple mutations or back mutations
 - 1572 changes in 71 families
 - Consider a tree to correctly count changes

Scoring Systems

PAM Matrices

- **Counting mutations. How many changes in these sequences?**

A**E**IKC

A**E**IKD

AD**L**KD

G**D**L**R**D

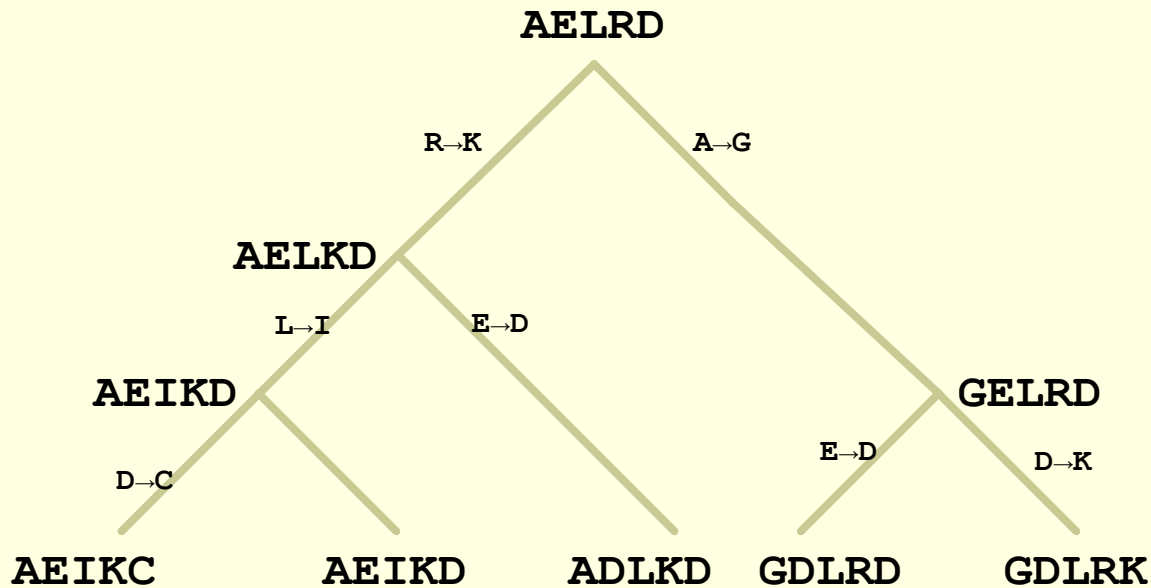
G**D**L**R**K

~~A ↔ G ?
E ↔ 6 ?
I ↔ L 6 ?
K ↔ R 6 ?
C ↔ 3 ?
C ↔ 1 ?
D ↔ K ?~~

Scoring Systems

PAM Matrices

- Counting mutations (correctly)



A → G	1
D → C	1
D → K	1
E → D	2
L → I	1
R → K	1

From Previous page

A ↔ G	6 ?
D ↔ C	3 ?
D ↔ K	3 ?
E ↔ D	6 ?
L ↔ I	6 ?
K ↔ R	6 ?
C ↔ K	1 ?

See also fig 5.1 in text

Scoring Systems

PAM matrices

- **The goal is to make a matrix that reflects 1% change**

- data were collected for amounts of change up to 15% so they have to be scaled
- first calculate *mutability* – how likely is it that each kind of residue undergoes a change
- then use eq. 5.6 to find a constant Λ so that the total change is 1%

$$0.01 = \sum_b f_b (1 - \Lambda m_b)$$

where f_b is the frequency of residue b and m_b is its mutability

- A mutation probability matrix can then be calculated by multiplying this 1 PAM matrix by itself some number of times. This is a *Markov process*, we assume that the only thing that affects the observed mutation is the current residue

Scoring Systems

PAM Matrices

- **Mutation probability matrix** - probability that residue in column b will be replaced by residue in row a after some amount of evolution

$$M_{bb} = 1 - \Lambda m_b \qquad M_{ab} = \frac{\Lambda m_b A_{ab}}{\sum_a A_{ab}}$$

m_b = mutability of residue b (probability of mutating)

A_{ab} = number of accepted point mutations

Λ = proportionality constant

Scoring Systems

PAM Matrices

- ***Mutation Probability Matrix***

- Probability that the residue represented by the column will mutate into the residue represented by the row after a specified amount of mutation, for instance, 1 PAM
- Not symmetric. The probability of $A \rightarrow S \neq S \rightarrow A$

- ***Relatedness odds matrix***

- Log-odds form of the mutation probability matrix
 - Log-odds matrix compares what is expected of homologous sequences to what is expected of unrelated (random) sequences.
- Remember that a log-odds matrix compares two models, in this case a model of relationship by homology (the mutation probability matrix) and relationship by chance

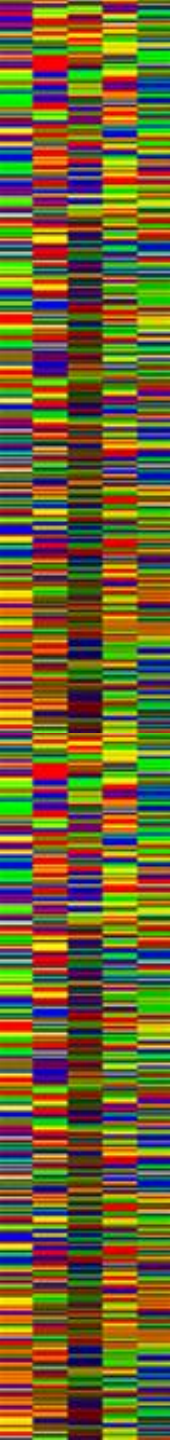
$$R_{ij} = \log(M_{ij} / f_i)$$

- Where f_i is the frequency of amino acid residue i at random

Scoring Systems

PAM Matrices

- *1 PAM matrix can be extended to simulate longer distances by multiplying it by itself.*
- *Dayhoff recommended PAM 250*
- *More modern recommendations suggest PAM 125 is a better general choice*
- *Which is best depends on the sequences!*



Scoring Systems

PAM Matrices

- ***Problems with the PAM approach:***
 - Not all positions are the same, e.g., internal vs external
 - Evolutionary rates vary greatly within a sequence
 - Each position has a unique three dimensional environment
 - Environment changes over evolutionary time as surrounding residues change
 - The most mutable positions were inadvertently selected as the basis for the calculation
 - proteins change more rapidly at the least constrained positions and most slowly at buried “core” positions

Scoring Systems

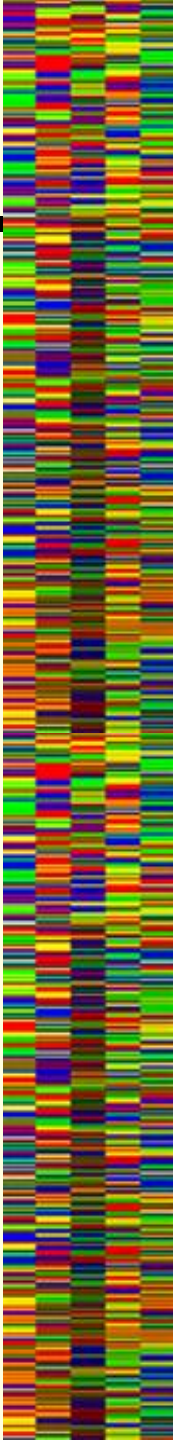
BLOSUM (BLOOcks SUbstitution Matrix)

- ***Based on PROSITE signatures***
 - Signatures are short expressions like C-X-X-C-X-X-X-S-T
- ***Locally align sequences to each signature to get "blocks"***
- ***Blocks are locally conserved regions, i.e., more constrained regions likely to be related to structure/function***
- ***Blocks contain sequences at all different evolutionary distances and may be highly biased (e.g., many identical sequences)***

Scoring Systems

BLOSUM Matrices (Ch 5.1)

- ***Dealing with bias and distance***
 - Cluster all sequences with less than X% identities
 - Clustered sequences count as 1 sequence
 - if X is 100% it simply removes identical sequences
 - if $X < 100\%$ it reduces the weight on closely related sequences
 - Calculate substitution frequencies and log-odds matrix
- ***This gives a BLOSUM X table***
 - BLOSUM 62 - sequences greater than 62% identical are clustered
 - BLOSUM 80 - sequences greater than 80% identical are clustered



Scoring Systems

BLOSUM Matrices

- METHYLTRANSFERASE BI**

TCMN_STRGA (331)	IADLGGGDGWFLAQILRRHPHATGLLMDLPRVA	74
TCMO_STRGA (173)	FVDLGGARGNLA AHLHRAHPHLRATCFDLPEME	81
ZRP4_MAIZE (204)	LVDVGGGIGAAAQAISKAFPHVKCSVLDLAHVV	68
CHMT_POPTM (204)	LVDVGGGTGAVVNTIVSKYPSIKGINFDLPHVI	41
COMT_EUCGU (205)	VVDVGGGTGAVLSMIVAKYPSMKGINFDLPHVI	42
COMT_MEDSA (204)	LVDVGGGTGAVINTIVSKYPTIKGINFDLPHVI	47
CRTF_RHOSH (205)	LMDVGGGTGAFLAAVGRAYPLMELMLFDLPVVA	59
OMTA_ASPPA (250)	VVDVGGGRGHL SRRV SQKHPHLRFIVQDLPAVI	47

Unweighted (BLOSUM 100) count of transitions for column 1, total $(n^2 - n)/2$ transitions

$$\begin{array}{llll} c_{FF} = 0 & c_{FI} = 1 & c_{FL} = 4 & c_{FV} = 2 \\ & c_{II} = 0 & c_{IL} = 4 & c_{IV} = 2 \\ & & c_{LL} = 6 & c_{LV} = 8 \\ & & & c_{VV} = 1 \end{array}$$

Scoring Systems

BLOSUM Matrices

- **Unweighted (BLOSUM 100) count (c_{ij}) of transitions for column 1**

N=28	F	I	L	V
F	0	1	4	2
I		0	4	2
L			6	8
V				1

- **$N = 28$ transitions, $f_{ij} = c_{ij} / N$ (foreground frequencies)**

	F	I	L	V
F	0.000	0.036	0.143	0.071
I		0.000	0.143	0.071
L			0.214	0.286
V				0.036

- **Log-Odds $s_{ij} = \log_2(f_{ij} / p_i p_j)$ - Background frequencies, p_i , from database**

◦ $p_F = 0.0397$ $p_I = 0.0529$

$p_L = 0.0917$

$p_V = 0.0649$

	F	I	L	V
F	undef	3.1	4.3	3.8
I		undef	3.9	3.4
L			4.7	4.6
V				3.1

Scoring Systems

BLOSUM Matrices

- **Background frequencies are calculated from the overall frequencies of the letters – a random background model**
 - in the database (as on previous slide)
 - or in the alignment (if the alignment is large enough)

$$p_{ii} = p_i^2 \quad p_{ij} = 2 p_i p_j$$

- **if the probabilities of the letters, p_i , are**

$$\circ p_F = 0.0397 \quad p_I = 0.0529 \quad p_L = 0.0917 \quad p_V = 0.0649$$

	F	I	L	V
F	0.0016	0.0042	0.0072	0.0052
I		0.0027	0.0097	0.0069
L			0.0084	0.0119
V				0.0042

Scoring Systems

BLOSUM Matrices

- ***Avoiding undefined logs***
- ***Undefined logs occur when either the foreground or background frequencies are zero***
- ***The problem arises because the sample is too small – rare items (residues) or events (mutations) are so uncommon they do not occur in the sample***
- ***Solution: use a pseudocount***
 - add one to every count (AKA plus one prior)
 - add a total of one distributed according to background

Scoring Systems

BLOSUM Matrices

- **Unweighted (BLOSUM 100) with plus one prior**

N=38	F	I	L	V
F	1	2	5	3
I		1	5	3
L			7	9
V				2

- **$N = 38$ transitions, $f_{ij} = c_{ij} / N$ (foreground frequencies)**

	F	I	L	V
F	0.026	0.053	0.132	0.079
I		0.026	0.132	0.079
L			0.184	0.237
V				0.053

	F	I	L	V
F	0.000	0.036	0.143	0.071
I		0.000	0.143	0.071
L			0.214	0.286
V				0.036

without pseudocount

- **Log-Odds $s_{ij} = \log_2(f_{ij} / p_i p_j)$ - Background frequencies, p_i , from database**

	F	I	L	V
F	4.0	3.6	4.2	3.9
I		3.3	3.8	3.5
L			4.5	4.3
V				3.7

	F	I	L	V
F	undef	3.1	4.3	3.8
I		undef	3.9	3.4
L			4.7	4.6
V				3.1

without pseudocount

Scoring Systems

BLOSUM Matrices – BLOSUM 80

- *Looking for sequences > 80% identical*

TCMN_STRGA (331)	IADLGGGDGWFLAQILRRHPHATGLLMDLPRVA	74
TCMO_STRGA (173)	FVDLGGARGNLA AHLHRAHPHLRATCFDLPEME	81
ZRP4_MAIZE (204)	LVDVGGGIGAAAQAISKAFPHVKCSVLDLAHVV	68
CHMT_POPTM (204)	LVDVGGGTGAVVNTIVSKYPSIKGINFDLPHVI	41
COMT_EUCGU (205)	VVDVGGGTGAVLSMIVAKYPSMKGINFDLPHVI	42
COMT_MEDSA (204)	LVDVGGGTGAVINTIVSKYPTIKGINFDLPHVI	47
CRTF_RHOSH (205)	LMDVGGGTGAFLAAVGRAYPLMELMLFDLPVVA	59
OMTA_ASPPA (250)	VVDVGGGRGHL SRRVSQKHPHLRFIVQDLP AVI	47

Percent Identity

	1	2	3	4	5	6	7
0 TCMN_STRGA	39.4	33.3	36.4	39.4	36.4	48.5	36.4
1 TCMO_STRGA		33.3	30.3	30.3	30.3	36.4	45.5
2 ZRP4_MAIZE			48.5	45.5	48.5	42.4	42.4
3 CHMT_POPTM				81.8	93.9	48.5	45.5
4 COMT_EUCGU					78.8	51.5	48.5
5 COMT_MEDSA						48.5	45.5
6 CRTF_RHOSH							36.4
7 OMTA_ASPPA							

Scoring Systems

BLOSUM Matrices – BLOSUM 80

```
TCMN_STRGA ( 331) IADLGGGDGWFLAQILRRRHPHATGLLMDLPRVA 74
TCMO_STRGA ( 173) FVDLGGARGNLA AHLHRAHPHLRATCFDLPEME 81
ZRP4_MAIZE ( 204) LV DVGGGIGAAAQAISKAFPHVKCSVLDLAHVV 68
CRTE_RHOSH ( 205) LMDVGGGTGAF LA AVGRAYPLMELMLFDLPVVA 59
OMTA_ASPPA ( 250) V DVGGGRGHL SRRV SQKHPHLRFIVQDLP AVI 47

CHMT_POPTM ( 204) LV DVGGGTGAVVNTIVSKYPSIKGINFDLPHVI 41
COMT_EUCGU ( 205) V DVGGGTGAVLSMIVAKYPSMKGINFDLPHVI 42
COMT_MEDSA ( 204) LV DVGGGTGAVINTIVSKYPTIKGINFDLPHVI 47
```



each sequence counts as one



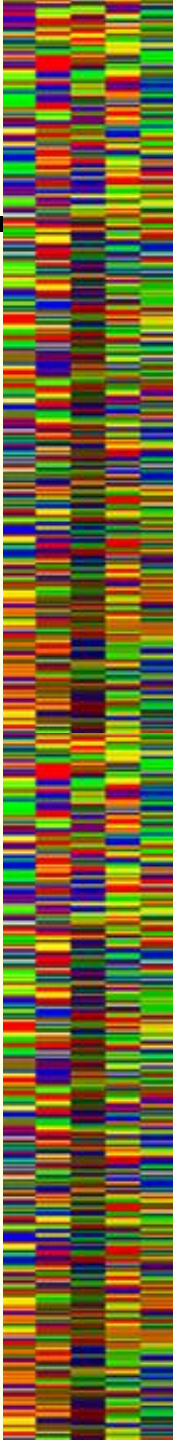
three sequences together count as one

N=15	F	I	L	V
F	0	1	2.67	1.33
I		0	2.67	1.33
L			2.33	3.33
V				0.33

N=25	F	I	L	V
F	1	2	3.67	2.33
I		1	3.67	2.33
L			3.33	4.33
V				1.33

with pseudocount

- **Equivalent to 6 sequences so number of transitions should equal 15 (6*5/2)**



Scoring Systems

BLOSUM Matrices – BLOSUM80

- ***foreground pseudocounts from previous slide***

N=25	F	I	L	V
F	1	2	3.67	2.33
I		1	3.67	2.33
L			3.33	4.33
V				1.33

N=25	F	I	L	V
F	0.040	0.080	0.147	0.093
I		0.040	0.147	0.093
L			0.133	0.173
V				0.053

- ***Log-odds (same background as before)***

	F	I	L	V
F	4.6	4.3	4.4	4.2
I		3.9	3.9	3.8
L			4.0	3.9
V				3.7

	F	I	L	V
F	4.0	3.6	4.2	3.9
I		3.3	3.8	3.5
L			4.5	4.3
V				3.7

BLOSUM100 (ungrouped)

Scoring Systems

BLOSUM Matrices

- ***Derived from a very large set of conserved sequence motifs.***
- ***Represents more core of protein and less surface than PAM***
- ***BLOSUM is always a blend of evolutionary distances***

PAM Matrices

- ***Derived from a relatively small set of closely related, small, and mostly globular proteins***
 - Note that JTT update (discussed in text) does not consider trees
- ***Biased towards surface***
- ***Cover entire evolutionary range from identical to random***

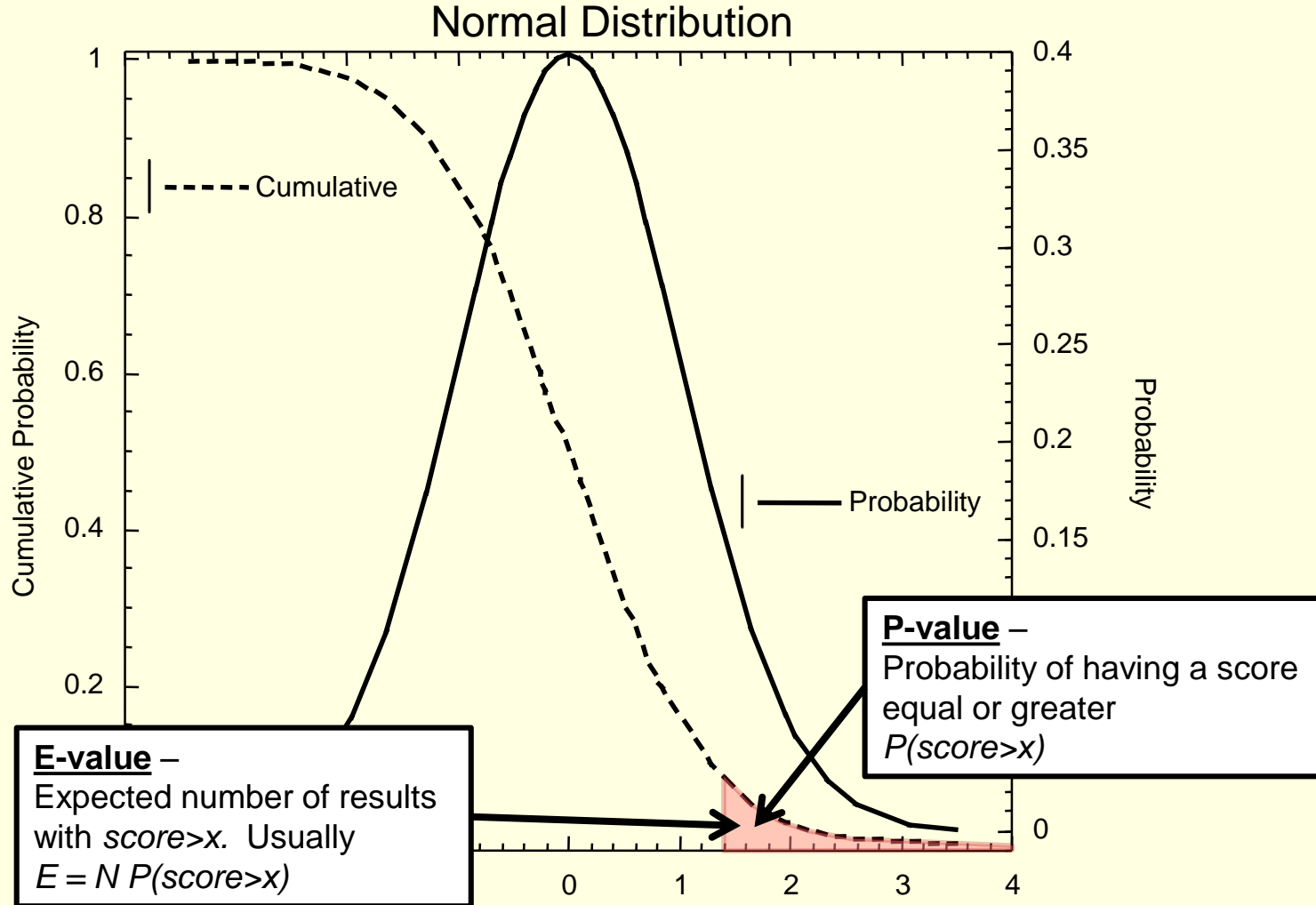
Scoring Systems

Target Frequencies

- *Alignment programs “try” to produce alignments in which the letter-pairs occur such that the highest possible score is produced*
- *EVERY scoring matrix implies a set of target frequencies*
- *The best alignments are produced when the target frequencies implied by the scoring matrix are the same as the correct alignment*
- *Disagreement between the scoring matrix target frequencies, and the biologically correct pair frequencies is one of the main reasons that the mathematically optimal alignment is not necessarily the biologically correct one (approximation made for gaps is the other)*

Sequence Comparison

Alignment Significance

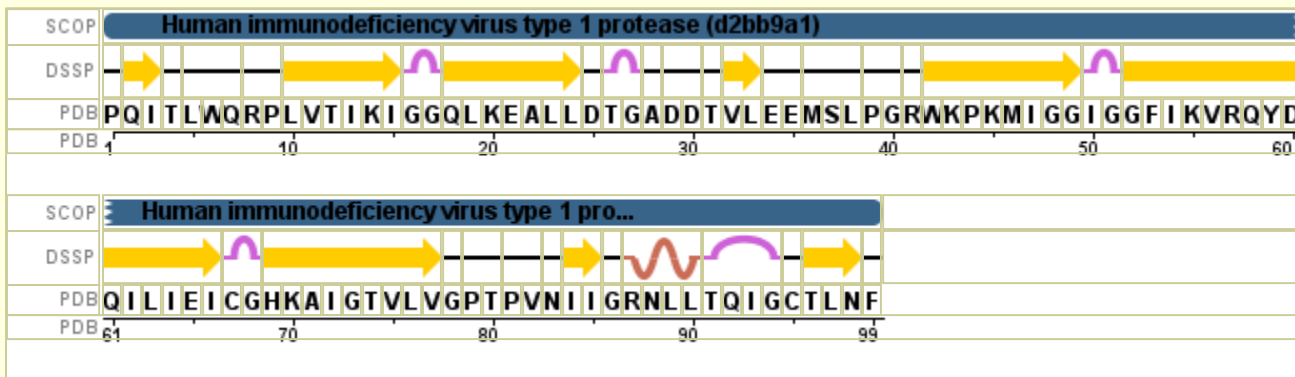
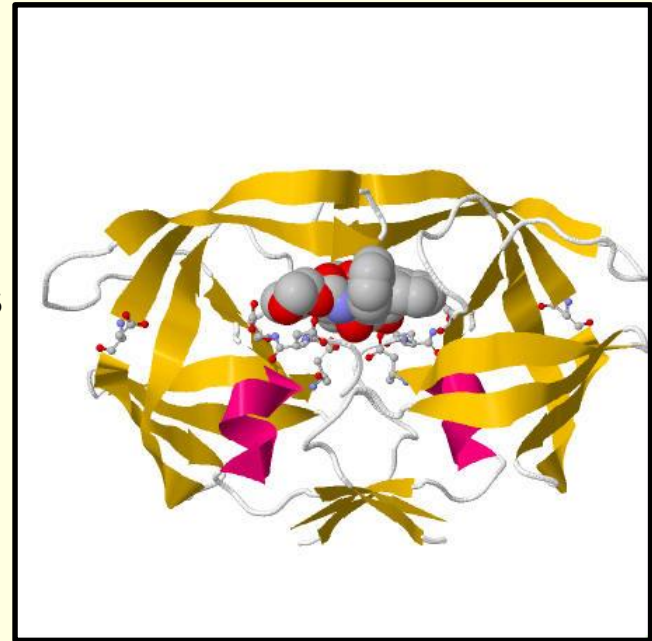


Sequence Comparison

Practical Alignments

- **HIV-1 protease**

- Cleaves viral polyprotein into active proteins
- 99 residues
- Aspartic protease
- Active site residues 25-27
- Active site flap: 46..50,52..56



Sequence Comparison

Practical Alignments

- **GALV – Gibbon Ape Leukemia Virus**

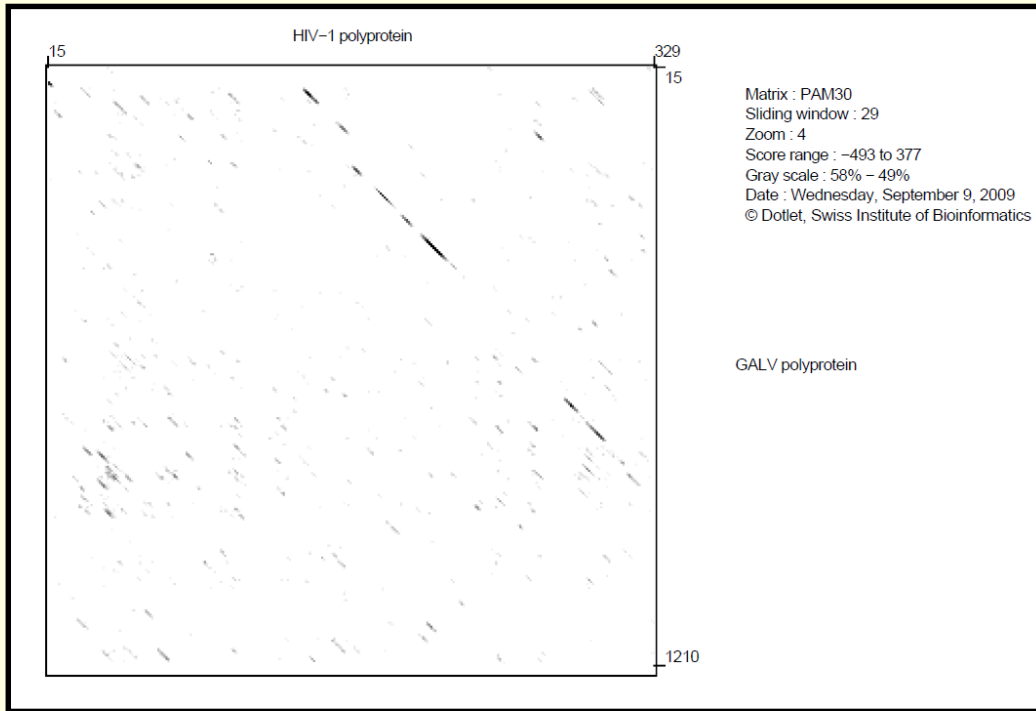
- Does it have the same mechanism (viral protease cleaves polyprotein)?
- Where is protease active site?
- Does it recognize the same cleavage site?

```
>P21414 Gibbon Ape Leukemia virus (GALV) entire polyprotein
gsqgsdplpe prvtltvegt pieflvdtga ehsvltqpmg kvgsrrtve gatgskvypw
ttkrllkigh kqvthsflvi pecpapllgr dlltklkaqi qfsaegpqt wgerptmclv
lnleeyrlh ekvpvssidp swlqlfptvw aeragmglan qvppvvelr sgaspvavrq
ypmskeareg irphiqkfld lgvlvpcrsp wntpllpvkk pgtdnyrpvq dlreinkrvq
dihptvpnpy nllsslppsy twysvldldk affclrlhpn sqplfafewk dpekngtgql
twtrlpqgfk nsptlfdeal hrdlapfral npqvvllqyv ddllvapty edckkgtqkl
lqelsklgyr vsakkaqlcq revtylgyll kegkrwltpa rkatvmkipv pttprqvref
lgtagfcrlw ipgfaslaap lypltkesip fiwteehqqa fdhikkalls apalalpdlt
kpftlyider agvargvltq tlgpwrrpva ylskkldpva sgwptclkav aavalllkda
dkltlgqnv t viashslesi vrqppdrwmt narmthyqsl llnervsfap pavlnpatll
pvseatpvh rcseilaeet gtrrdledqp lpgvptwytd gssfitegkr ragapivdgk
rtvwasslpe gtsaqkaelv altqalrlae gkniniytds ryafatah gaiyqrgll
tsagkdiknk eeilalleai hlprrvaiih cpghqrgsnp vatgnrrade aakqaalstr
vlagttkpqe piepaqektr preltpdrgk efikrlhqlt hlgpekllql vnrtsl lipn
lqsavrevts qcqacamtna vttyretgkr qrgdrpgvyw evdfteikpg rygnkyllvf
idtfsgwvea fptktetali vckkileeil prfgipkvlg sdngpafvaq vsqglatqlg
inwklhcayr pqssgqverm nrtiketltk laletggkdw vtllplallr arntpgrfgl
tpyeilyggp ppilesgetl gpddrflpvl fthlkaleiv rtqiwdqike vykpgvtvip
hpfqvgdqvl vrrhrpssle prwkgpylvl lttptavkvd giaawvhash lkpappsapd
eswelekt dh plklrirrrr desak
```

Sequence Comparison

Practical Alignments

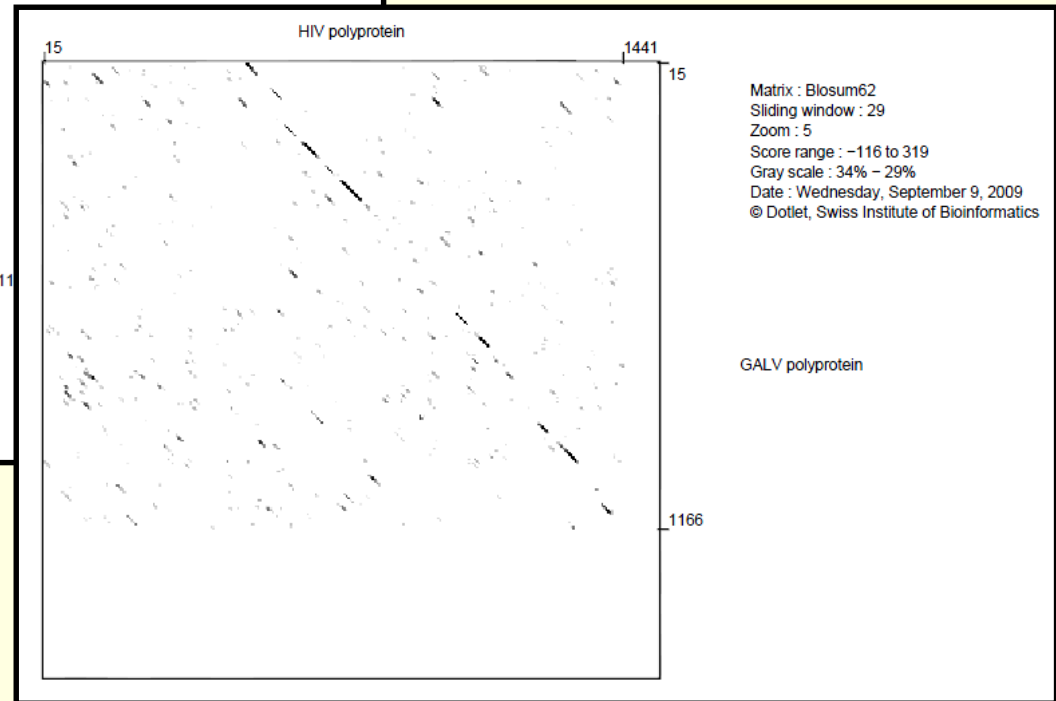
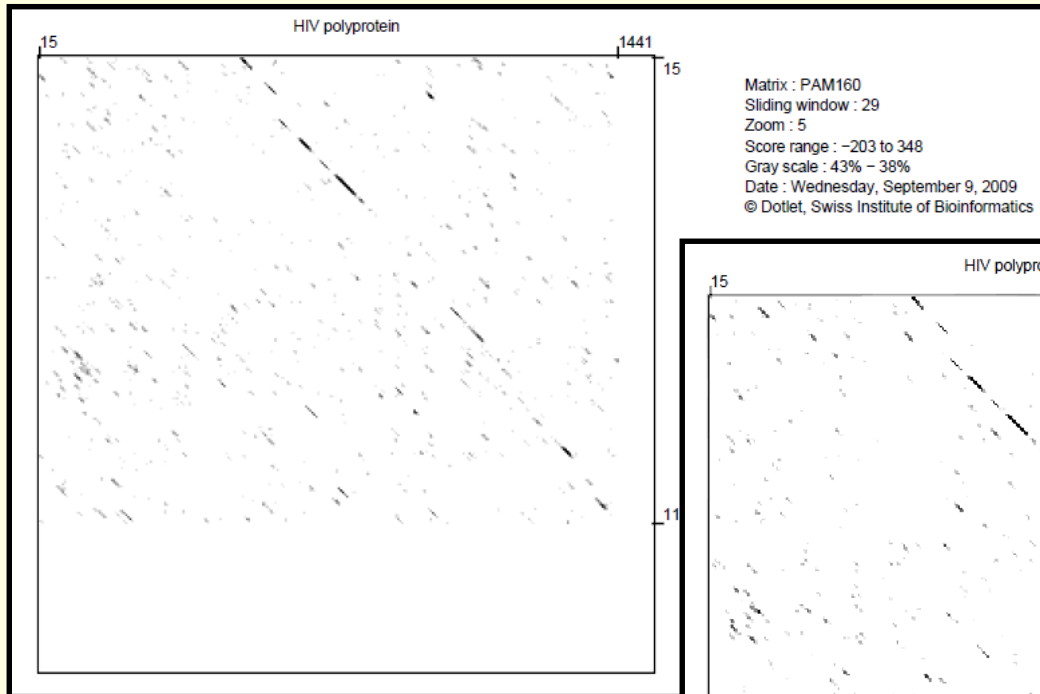
- *HIV-1 polyprotein vs GALV polyprotein*



Sequence Comparison

Practical Alignments

- *HIV-1 polyprotein vs GALV polyprotein*



Sequence Comparison

Practical Alignments

• Global alignment of HIV and GALV polyprotein

ALIGN calculates a global alignment of two sequences
version 2.2u
Please cite: Myers and Miller, CABIOS (1989) 4:11-17
hiv-1 protease 99 bp 99 aa vs.
GALV polyprotein 1165 bp 1165 aa
using matrix file: BLOSUM50, gap open/ext: -14/-4
3.7% identity in 1165 aa overlap; Global score: -4163

```
hiv-1 -----PQITL-----  
          :...:  
GALV  GSQGSDDLPEPRVTLTVEGTPIEFLVDTGAEHSVLTQPMGKVGSRRTVVEGATGSKVYPW  
  
hiv-1 -----  
GALV  TTRRLKIKGHQVTHSFLVIPECPAPLLGRDLLTKLKAQIQFSAEGPQVTWGERPTMCLV  
  
hiv-1 -----  
GALV  LNLEEEYRLHEKFPVSSIDPSWLQLFPTVWAERAGMGLANQVPPVVVELRSGSPVAVRQ  
  
hiv-1 -----WQRPVLTVKIGG-----  
          :. :. :. :  
GALV  YPMSKEAREGIRPHIQKFLDLGLVLPGRSPWNTPLLPVKKPGTNDYRPVQDLREINKRVQ  
  
hiv-1 -----  
GALV  DIHPTVPNPNYLLSSLPPSYTWYSVLDLKDFAFFCLRHLHNSQPLFAFEWKDPEKNGTGQL  
  
hiv-1 -----  
GALV  TWTRLPGQFKNSPTLFDEALHRDLAPFRALNPQVLLQYVDDLLVAAPTYEDCKKGTQKL  
  
hiv-1 -----  
GALV  LQELSKLGYRVSAKKAQLCQREVITYLGYLLKEGRWLT PARKATVMKIPVPTTPRQVREF  
20  
hiv-1 -----QLKEALDT-----  
          :. :. :. :. :  
GALV  LGTAGFCRLWIPGFASLAAPLYPLTKESIPFIWTEEHQQAQFDHIKALLSAPALAPDLT  
  
hiv-1 -----  
GALV  KPFTLYIDERAGVARGVLTQTLPWRRFVAYLSKKLDPVASGWPTCLKAVAVALLLKDA  
  
hiv-1 -----GADDTV-----  
          :. :.  
GALV  DKLTILGNVTVIASHSLESIVRQPDRWMTNARMTHYQSLLLNERNVSFAPPVLPATLL
```

```
hiv-1 -----LEDINLPG-----  
          :. :. :.  
GALV  PVSEATPVHRCSEILAEETGTRRDLEDQPLPGVPTWYTDGSSFITEGKRRAGAPIVDGK  
  
hiv-1 -----  
GALV  RTVWASSLPEGTSQAQKAEVALTQALRLAEGKNNIYTDSTRYAFATAHIGAIYKQRGLL  
  
hiv-1 -----  
GALV  TSAGKDIKNKEEILALLEAIHLPRRVAIIHCPGHQSGNSPVATGNRRADEAAKQALSTR  
  
hiv-1 -----KWPKMIGGIGG--FIK-----  
          :. :. :. :. :. :.  
GALV  VLAGTTKQPEPIEPAQEKTRPRELTPDRGKEFIKRLHLQTLHGPEKLLQLVNRTSLIPN  
  
hiv-1 ----VRQYD-----  
          :. :.  
GALV  LQSAVREVTSQCQACAMTNAVTTYRETGKRQRGDRPGVYWEVDFTEIKPGRYGNKYLLVF  
  
hiv-1 -----QILIEICGKKAIGTVL---VGPTPVNIIGRNMLTQIG  
          :. :. :. :. :. :. :. :.  
GALV  IDTFSGWVEAFPPTKTETALIVCKKILEEILPRFGIPKVLGSDNGPAFVAQVQGLATQLG  
  
hiv-1 CTLNF-----  
          :. :.  
GALV  INWKLHCAYRPQSSGQVERMRTIKETLTKLALETGGKDWVTLPLALLRARNTPGRFGL  
  
hiv-1 -----  
GALV  TPYEILYGGPPPILESGETLGPDDRFLVPLVFTHLKALEIVRTQIWDQIKEYKPGTVTIP  
  
hiv-1 -----  
GALV  HPFQVGDQVLVRRHRPSSLEPRWKGPIYLVLLTTPAVKVDGIAAWVHASHLKAPPSPAD  
  
hiv-1 -----  
GALV  ESWELEKTDHPLKLRIRRRRDESAK
```

Sequence Comparison

Practical Alignments

• HIV polyprotein vs GALV polyprotein

using matrix file: BL50 (15/-5), gap-open/ext: -14/-4 E(limit) 0.05

29.0% identity in 93 aa overlap (501-592:10-102); score: 107 E(10000): 0.029

```
      510      520      530      540      550      560
hiv-1  RPVVTVRVGGQKLEALLDTGADTVELEINLPKWKPKMIGGIGGFIKVRQYDQILIEIC
      . . . . . : : : : : . . . . . : : : : : . . . . . : : : : :
GALV   EPRVTLTVEGTPIEFVLDVTGAEHSVLTQPMGKVGSRRTVVEGATGSKVYPWTKRLKLG
      10      20      30      40      50      60

      570      580      590
hiv-1  GKKAIGTILVGPT-PVNIIGRNMLTQIGCTLNF
      . . . . . : : : : : . . . . . : : : : :
GALV   HKQVTHSFLVIPECPAPILGRDLLTKLKAQIQF
      70      80      90      100
```

A

40.0% identity in 80 aa overlap (1020-1095:625-704); score: 120 E(10000): 0.0014

```
      1020      1030      1040      1050      1060      1070
hiv-1  QLETEPIVGAETFFYVDGAANRETKKKGAGYVDRGRQKV--VSLTE-TTNQKTELQAIHL
      . . . . . : : : : : . . . . . : : : : : . . . . . : : : : :
GALV   DLEDQPLPGVPTWYTDGSSFITTEGKRRAQPIVDGKRTVWASSLEPGTSAQKAEALVALTQ
      630      640      650      660      670      680

      1080      1090
hiv-1  ALQDS-GSEVNIVTDSQYAL
      . . . . . : : : : : . . . . . : : : : :
GALV   ALRLAEGKNINIYTDSRYAF
      690      700
```

C

27.9% identity in 247 aa overlap (595-835:160-399); score: 296 E(10000): 2.6e-21

```
      600      610      620      630      640      650
hiv-1  SPIETVPVKLPGMDGPRVKQWPLTEEEKIKALTEICKDMEKEGKILKIGPENPYNTPVFA
      . . . . . : : : : : . . . . . : : : : : . . . . . : : : : :
GALV   NQVPPVVVELRSGASPVAVRQYPMSEKEAREGIRPHIQKFLDLGLVPC--RSPFWNTPLLP
      160      170      180      190      200      210

      660      670      680      690      700
hiv-1  IKKDKSTKWRKLVNFRELNKRQDFWEVQLGIPHPAGLKKKKS-----VTVLDVGDAYFS
      . . . . . : : : : : . . . . . : : : : : . . . . . : : : : :
GALV   VKKPGTNDYRPVQDLREINKRVQD---IHPTVPNPNYLLSSLPSTWYSVLDLKDFFC
      220      230      240      250      260      270

      710      720      730      740      750      760
hiv-1  VPLDEDFRKYTAFTIPSINNETPGIRYQYNVLPQGWKGSIPAIFQSSMTKILEPFRKTNPE
      . . . . . : : : : : . . . . . : : : : : . . . . . : : : : :
GALV   LRLHPNSQPLFAFEWKDPEKGNTE-QLTWTRLPQGFKNSTPLFDEALHRDLAPFRALNPFQ
      280      290      300      310      320      330

      770      780      790      800      810      820
hiv-1  IVIYQYMDLTVGSDLEIGQHRTKIEELREHLLKWGFTTDPKHKHQK-EPPFLWMGYELHP
      . . . . . : : : : : . . . . . : : : : : . . . . . : : : : :
GALV   VLLQYVDDLVAAPT-YEDCKKGTQKLLQELSKLGYRVSAKKAQLCQREVTVLGYLLKE
      340      350      360      370      380      390

      830
hiv-1  DKWTVQP
      . . .
GALV   GKRWLTP
```

B

25.7% identity in 148 aa overlap (1184-1317:844-991); score: 120 E(10000): 0.0014

```
      1190      1200      1210      1220      1230
hiv-1  VAKEIVASCDKQCQLKG-----EAMHGQVDCSPGI-WQLDCTHLE----GKIIIVAVHV
      . . . . . : : : : : . . . . . : : : : : . . . . . : : : : :
GALV   AVREVTSQCQACAMTNAVTTYRETGKRQRGDRPGVYWEVDFTEIKPGRYGNKYLIVFIDT
      850      860      870      880      890      900

      1240      1250      1260      1270      1280
hiv-1  ASGYIEAEVIPAETGQETAYFILK-LAGRWPV-KVVHTDNGSNFTSAAVKAACWANIKQ
      . . . . . : : : : : . . . . . : : : : : . . . . . : : : : :
GALV   FSGWVEAFPTKTETALIVCKKILEEILPRFGIPKVLGSDNGPAFVAQVSQGLATQLGINW
      910      920      930      940      950      960

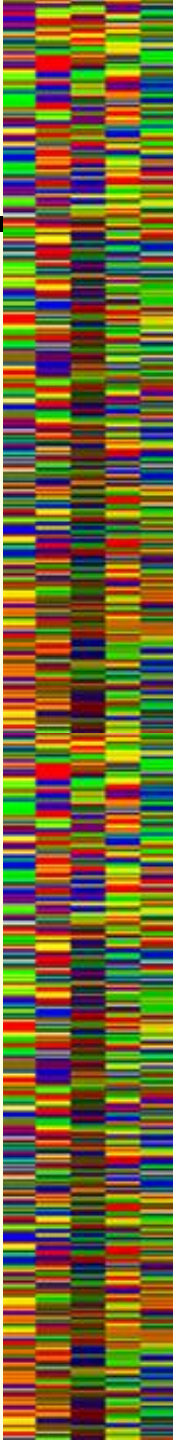
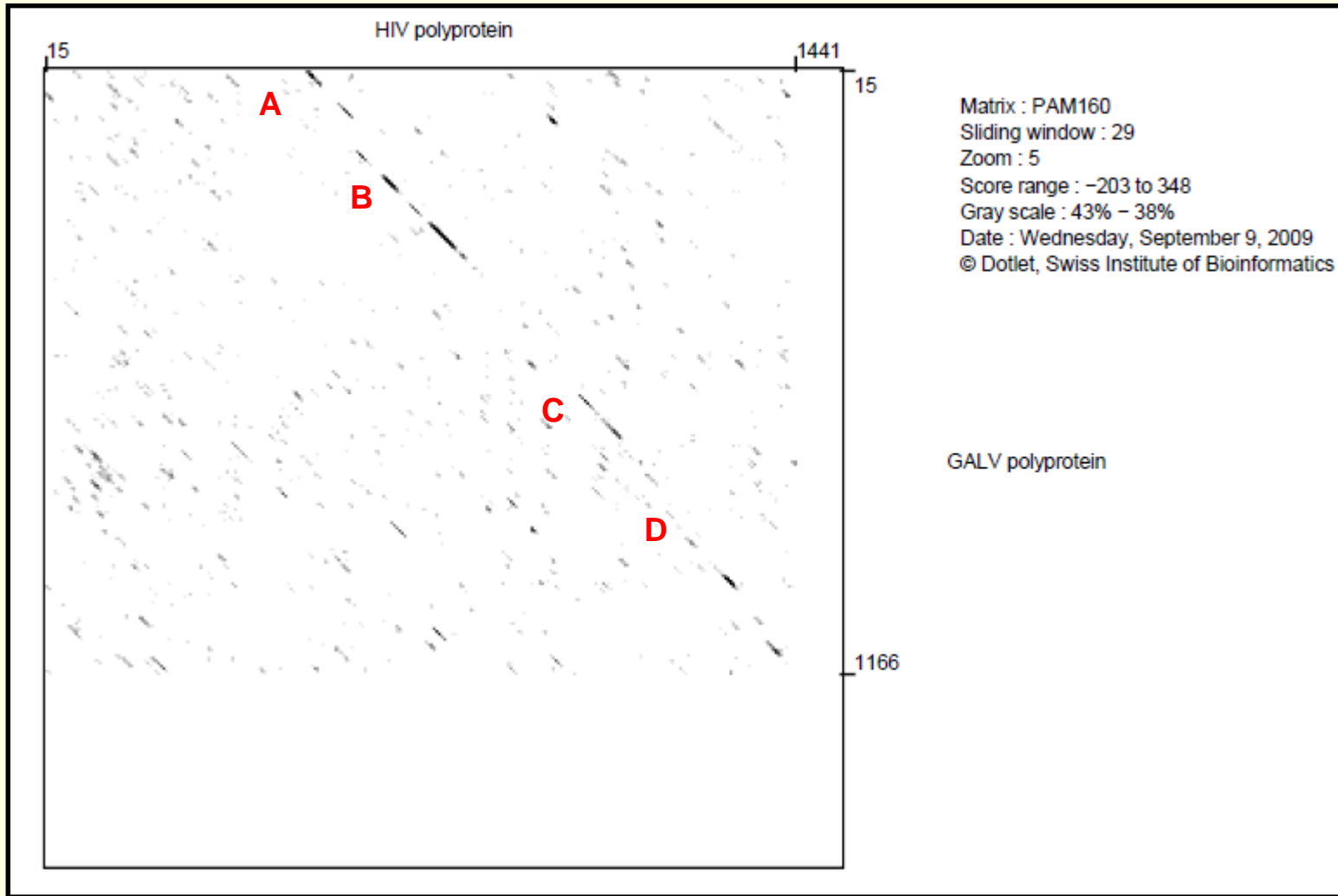
      1290      1300      1310
hiv-1  EFGIPYNPQSQGVSMNKLKIIQGV
      . . . . . : : : : : . . . . . : : : : :
GALV   KLHCAYRPQSSGQVERMNRITIKETLTKL
      970      980      990
```

D

Sequence Comparison

Practical Alignments

- *HIV-1 polyprotein vs GALV polyprotein*



Sequence Comparison

Practical Alignments

Statistics: (shuffled [500]) MLE statistics: Lambda= 0.1032; K=0.0432
Threshold: E() < 1 score: 63
Algorithm: Smith-Waterman (SSE2, Michael Farrar 2006) (7.0 Oct 2008)
Parameters: BL50 matrix (15:-5), open/ext: -5/-2
Scan time: 0.090

>>>AAA21284, 99 aa vs TMP.q2 library

>>P63120 HERV-K endogenous retrovirus type K, protease (156 aa)

Waterman-Eggert score: 166; 29.3 bits; E(1) < 2.4e-05

35.3% identity (57.8% similar) in 102 aa overlap (2-94:4-101)

Entrez Lookup Re-search database General re-search

```

                10      20      30      40      50      60      70
AAA212 QITLWQRPLVTVKIGGQLKEALLDTGADDTVLEDINL-PGKW-KPKMI-G--GIGGFIKVRQYDQILIEIC-GKKAI-GT
      ..  :..  . : : . : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
P63120 QVSE-NRPVCKAIIQGKQFEGLVDTGADVSIIA-LNQPKNWPKQKAVTGLVGIGTASEVYQSTEILH--CLGPDNQEST
                10      20      30      40      50      60      70

                80      90
AAA212 V--LVGPTPVNIIGRNMLTQIG
      : ..  :..  : : : : : :
P63120 VQPMITSIPLNLWGRDLLQQWG
                80      90      100
```

Sequence Comparison

Practical Alignments

Statistics: (shuffled [500]) MLE statistics: Lambda= 0.1745; K=0.05741
Threshold: E() < 1 score: 39
Algorithm: Smith-Waterman (SSE2, Michael Farrar 2006) (7.0 Oct 2008)
Parameters: BL50 matrix (15:-5), open/ext: -10/-2
Scan time: 0.070

>>>AAA21284, 99 aa vs TMP.q2 library

>>P63120 HERV-K endogenous retrovirus type K, protease (156 aa)

Waterman-Eggert score: 137; 38.6 bits; E(1) < 3.7e-08

33.0% identity (57.4% similar) in 94 aa overlap (7-94:8-101)

Entrez Lookup Re-search database General re-search

```
      10      20      30      40      50      60      70      80
AAA212 QRPLVTVKIGGQLKEALLDTGADDTVLEDINLPGKW-KPKMIGG---IGGFIKVRQYDQILIEICGKKAIGTV--LVGPT
      . . . . . : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
P63120 NRPVCKAIIQ GKQFEGLVDTGADVSIIALNQWPKNWPQKAVTGLVGIGTASEVYQSTEILHCLGPDNQESTVQPMITSI
      10      20      30      40      50      60      70      80
```

90

AAA212 PVNIIGRNMLTQIG

: : : : : : : :

P63120 PLNLWGRDLLQQWG

90 100

Sequence Comparison

Practical Alignments

Statistics: (shuffled [500]) MLE statistics: Lambda= 0.1904; K=0.06576
Threshold: E() < 1 score: 36
Algorithm: Smith-Waterman (SSE2, Michael Farrar 2006) (7.0 Oct 2008)
Parameters: BL50 matrix (15:-5), open/ext: -12/-2

>>>AAA21284, 99 aa vs TMP.q2 library

>>P63120 HERV-K endogenous retrovirus type K, protease (156 aa)

Waterman-Eggert score: 131; 39.9 bits; E(1) < 1.5e-08

33.0% identity (57.4% similar) in 94 aa overlap (7-94:8-101)

Entrez Lookup Re-search database General re-search

```
      10      20      30      40      50      60      70      80
AAA212 QRPLVTVKIGGQLKEALLDTGADDTVLEDINLPGKW-KPKMIGG---IGGFIKVRQYDQILIEICGKKAIGTV--LVGPT
      . . . . . : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
P63120 NRPVCKAIIQ GKQFEGLVDTGADVSI IALNQWPKNWPQKAVTGLVGIGTASEVYQSTEILHCLGPDNQESTVQPMITSI
      10      20      30      40      50      60      70      80
```

90

AAA212 PVNIIGRNMLTQIG

. : : : : :

P63120 PLNLWGRDLLQQWG

90 100

Sequence Comparison

Practical Alignments

Statistics: (shuffled [500]) MLE statistics: Lambda= 0.2295; K=0.1712
Threshold: E() < 1 score: 34
Algorithm: Smith-Waterman (SSE2, Michael Farrar 2006) (7.0 Oct 2008)
Parameters: BL50 matrix (15:-5), open/ext: -50/-1
Scan time: 0.050

>>>AAA21284, 99 aa vs TMP.q2 library

>>P63120 HERV-K endogenous retrovirus type K, protease (156 aa)

Waterman-Eggert score: 83; 30.0 bits; E(1) < 1.4e-05
36.1% identity (66.7% similar) in 36 aa overlap (7-42:8-43)
Entrez Lookup Re-search database General re-search

	10	20	30	40
AAA212	QRPLVTVKIGGQ	LKEALLDTGAD	DTVLEDINLPG	KW
 : . :	: : : : : : : . :
P63120	NRPVCKAIIQG	KQFEGLVDTG	ADVSIIALNQ	WPKNW
	10	20	30	40

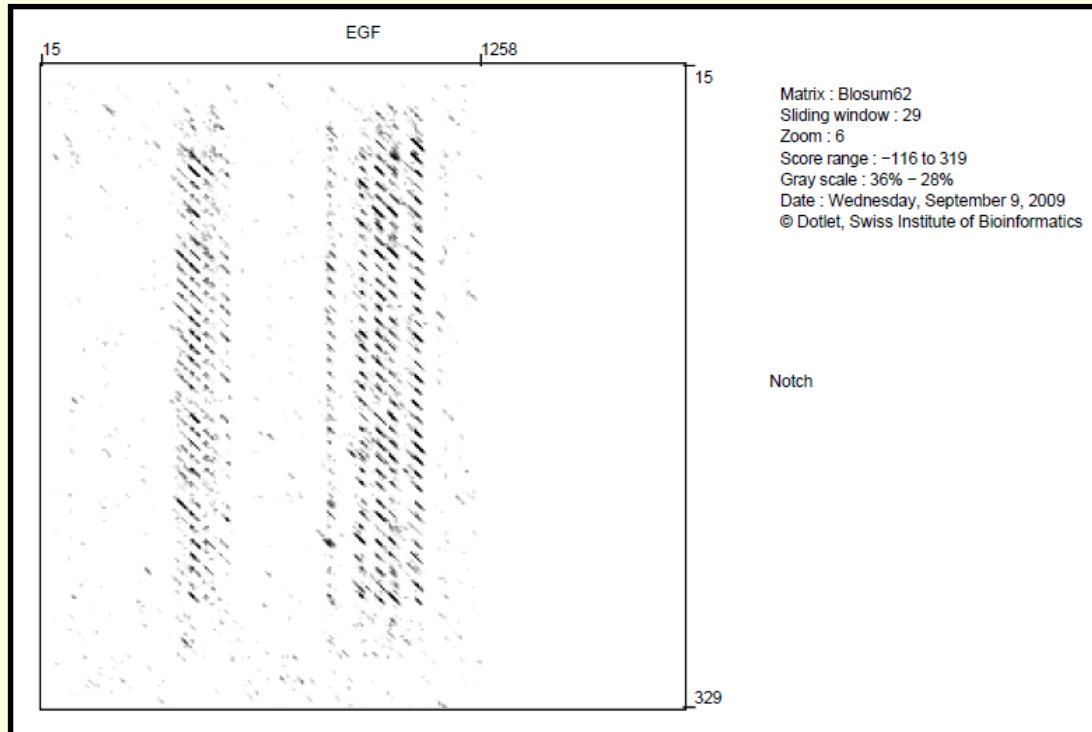
Waterman-Eggert score: 56; 21.1 bits; E(1) < 0.0069
36.8% identity (68.4% similar) in 19 aa overlap (76-94:83-101)
Entrez Lookup Re-search database General re-search

	80	90
AAA212	LVGPTPVNIIG	RNMLTQIG
	..	: : : : : :
P63120	MITSIPLN	LWGRDLLQ
	90	100

Sequence Comparison

Practical Alignments

- **Proteins with repeats**
 - Mouse Epidermal Growth Factor (EGF)
 - Drosophila Notch



Schedule Week 3

4 – 6 September

- ***Wednesday - Alignments /Scoring Systems***
- ***Friday – Searching for homologs***

- ***9 Sep Monday***
 - Searching for homologs (Ch 5.3-5.5, Ch 4.6-4.7)
- ***11 Sep Wednesday***
 - Searching for homologs
- ***13 Sep Friday***
 - Sequence motifs (Ch 4.8-4.10)

Database Searching

Sequence database searching

- *Null hypothesis: a pair sequences are unrelated*

If the pair of sequences have a more significant match than could reasonably be expected from unrelated (not homologous) sequences

Then they must be related (homologous)