

How where contigs assembled?

Sequence alignment

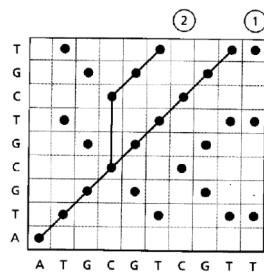
BLAST scores

Demo of NCBI and UCSC

How do we score these alignments?

Which alignment do we choose?

How do we score gaps?



ATGCGTCGTT
|||||||
ATGCGTCGT

ATG - - CGTCGTT
||| |||
ATGCGTCGT

Basic sequence alignment and
database searching:
BLAST nucleotide searches

Sophisticated programming

- Dynamic programming
- Look up tables
 - Database broken up into short segments which are screened first for exact matches (word size). If a “hit” occurs, then try to extend.
- Needle-Wunsch alignments
 - Optimal alignment

Basic Local Alignment Search Tool (BLAST)

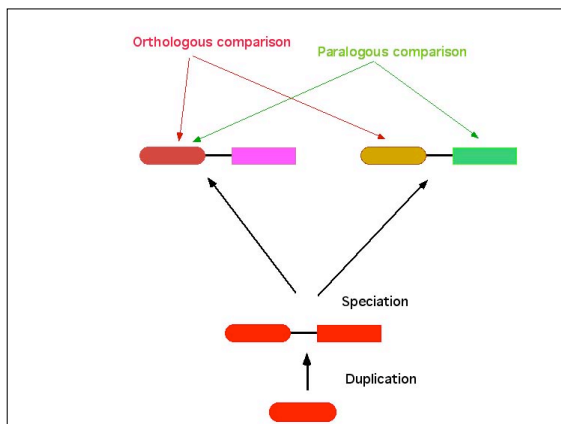
- Compares a QUERY (your sequence) to a database (i.e. GenBank).
- Different programs depending on sequence
 - BLASTN (DNA vs DNA)
 - BLASTP (AA vs AA)
 - BLASTX (Translated DNA vs AA)
 - TBLAST: Translated database searches

Main questions for today

- I have a DNA/protein sequence, is it similar to anything else in the database?
- My sequence is similar to something in the database, is it significant?

What does statistically significant similarity mean?

- Common ancestry
 - Homologous (all or none)
 - Orthologous: separated by speciation event: same locus
 - Paralogous: Separated by gene duplication: different locus
- May provide clues to similar function



Three scores from BLAST

- Raw score
 - No information unless know scoring matrix
- Bit Score
 - Takes into account scoring matrix
- E-value
 - Easy to interpret statistic, takes into account database search size.

Raw score

- Sum up score for matches
- Subtract penalties for miss-match
- Subtract penalties for gaps (open and extend)

C	C	A	A	G	A	C
C	C	A	T	G	A	C
1	1	1	-2	1	1	1

Raw Score = 1+1+1-2+1+1+1 = 4

Bit score: take into account the scoring system

- Normalize raw scores
- Two values describe scoring system
 - K and λ (λ) (both values determined by program)

$$S' = \frac{\lambda S - \ln K}{\ln 2}$$

- Takes into account the search space
 - m and n (database and query size)

- One equation from score to E-value

The Expect value (E) is a parameter that describes the number of hits one can "expect" to see just by chance when searching a database of a particular size. Typical cut-offs biologists use are $\sim E = 10^{-5}$

- 3

Expect score: take into account search space

Score = 64.9 bits (33) Expect = 0.001
 Identical = 27/29 (93%)
 Strand = Plus / Plus

Query: 1 aggggggattagggggggggggggggggggggg
 Subject: 1 aggggggattagggggggggggggggggggggg

CF9 time: 0.04 user secs 0.02 sys

Lambda: 1.22

Bapped: 1.33

Matrix: blastn matrix 1 -2

Gap Penalties: Gapopen: 5, Extension: 2

Number of Bits to MB: 1

Number of Sequences: 1

Number of extensions: 1

Number of sequences better than 10.0: 1

Length of query: 29

Length of database: 5,162,042,099

Effective length of database: 4,564,059,483

Effective search space: 4178333347

Effective search space used: 4178333347

T: 0

X1: 6 (11.5 bits)

X2: 26 (130.0 bits)

S1: 22 (122.0 bits)

S2: 17 (133.4 bits)

$$E = mn2^{-S'}$$

$$E = (9)(4,564,059,483)2^{-44.9}$$

$$E = 0.001$$

What if database was the *Drosophila* genome?

$$E = (9)(1,800,000)2^{-44.9}$$

$$E = .00000049 (4.9e^{-7})$$

Summary of statistics

• Raw score

- First pass - no information because scoring system unknown

• Bit Score

- Take into account the scoring system

• E-value

- Takes into account the search space. Easy to interpret. ONLY USE E-VALUES!

BLAST against database: Important input parameters

- Database
- Filter
- (Expect)
- (Word size)
- Matrix (proteins)
- Gap open and extend

Affine Gap Penalty

- Usually do not need to change
- Penalty to open Gap
 - Default 5 for nucleotide, 11 for protein
- Penalty to extend gap
 - Default 2 for nucleotide, 1 for protein

Databases

- NR (non-redundant)
 - GenBank+EMBL+DDBJ+PDB
- ESTs
 - Expressed Sequence tags, organism subsets
- HTGS (High-Throughput Genomic Sequence)
 - Low quality, but may find new matches
- PAT (Patent)
- Vector
- Organism specific genomic sequences
- Month
 - Monthly updates to search for new hits

Filter

- Mask out low complexity regions (ATATATATATATATATA)
- Only masks Query sequence
- Can mask out repeats like Lines & Sines
- Mask out “lookup table”
- Masked regions appear as Xs or Ns in Query sequence

What types of searches to perform on protein coding regions?

Protein searches are more sensitive

- More characters to compare.
 - 20 amino acids versus 4 bases
- Substitution matrices.
- No “wobble” codon positions.
- Potentially smaller database size.

- # Protein searches are more sensitive
- More characters to compare.
 - 20 amino acids versus 4 bases
 - Substitution matrices.
 - No “wobble” codon positions.
 - Potentially smaller database size.

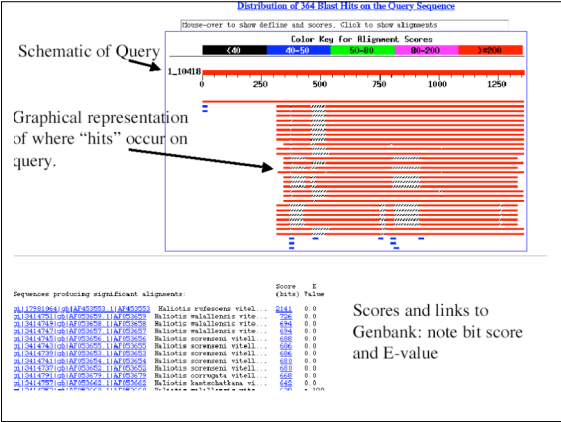
Protein Scoring Matrices

- Multiple matrices of each type optimized for particular divergence. Higher number for greater divergence.
- PAM (PAM250)
 - Align all sequences of varying divergence and calculate probability of substitution between different amino acids
- BLOSSUM (BLOSSUM62)
 - Align blocks of sequences with varying divergence and calculate probability of substitution between different amino acids

- # Protein Scoring Matrices
- Multiple matrices of each type optimized for particular divergence. Higher number for greater divergence.
 - PAM (PAM250)
 - Align all sequences of varying divergence and calculate probability of substitution between different amino acids
 - BLOSSUM (BLOSSUM62)
 - Align blocks of sequences with varying divergence and calculate probability of substitution between different amino acids

[illegible][illegible]

BLAST example



[illegible] $9.2e^{-2}$

Alignments

Score = 2141 bits (1080), Expect = 0.0
Identities = 1224/1272 (96%)
Strand = Plus / Plus

Query: 1 tgcagcagcagacaagagcaaacagcaacccatcatcagctaccccggtacattcaagggc 60
 |||
 Sbjct: 127 tgcagcagcagacaagagcaaacagcaacccatcatcagctaccccggtacattcaagggc 186

Query: 61 catgttatcaaaagcaatgcaaatottotgcaagaatggatggatgcaaatgaaactgga 120

Objet: 187 catgttatcaaaagcaaatgcaaatottotgcaagaatggatggatgcaaatgcaaatgcaaatgga 246

Query: 121 cgtggaattaacatgattgattcattatccacaaacttacacttcagtagttccaggg 180
 |||
 Sbjct: 247 agaggcattaacatgattgattcattatccacaaacttacacttcagtagttccaggg 306

Query: 181 gcatgtgtattccggggaacataatcttatccgactaatgactotattgagaagtataac 240
 |
 Subject: 267

Query: 241 gtaagtgttgcaactgctatggagcgatgggaacccaacttatgactctcttggaatgtaat 300

Subjet: 367 gtaagtgttgaactgctatggagogatgggaaccccaacctatgagtctcttgaaatgtaatt 426

1: [AF453553](#). Reports *Haliotis rufescen...*[gi:17981064]

```

DEFINITION      AF435353.1 11166 bp mRNA linear INV 30-MAY-2002
ACCESSION       Malattia refenscens vitelline envelope sperm linear receptor (VERL)
KEYWORDS        mRNA, partial cds
VERSION         AF435353.1 GI:17981054
SOURCE          Malattia refenscens (California red abalone)
ORGANISM        Malattia refenscens (Mollusca; Gastropoda; Orthogastropoda;
                Vetigastropoda; Bivalvia; Haliotidae; Haliotis.
REFERENCE       1 11166
AUTHOR          Galindo-B., Moy-C., Swanson,W.J. and Vacquier,V.D.
TITLE           The VERL, the egg vitelline envelope receptor
JOURNAL         J. Exp. Zool. 292: 117-127 (2002)
MEDLINE         200215453
FORMED         12810535
NOTE            (based on 11166)
AUTHOR          Galindo,B.T., Moy-C., Swanson,W.J. and Vacquier,V.D.
TITLE           Submitted (29-MAY-2002) MBRD, VCMD-Scripts Institution of
JOURNAL         Oceanographic Sciences, La Jolla, CA 92093-0322, USA
                Location/Qualifiers
                /organism="Malattia refenscens"
                /mol_type="mRNA"
                /db_xref="taxon:6454"
                /number="11166"
                /db_xref="VERL"
                /number="11166"
                /db_xref="VERL"
                /analysis="start"
                /product="vitelline envelope sperm linear receptor"
                /protein_id="AF435353.1"
                /db_xref="GI:17981054"
                /translation="MSGNGNGGQCFPLFVETVQMDAPFQVFPVSLSLTSLT

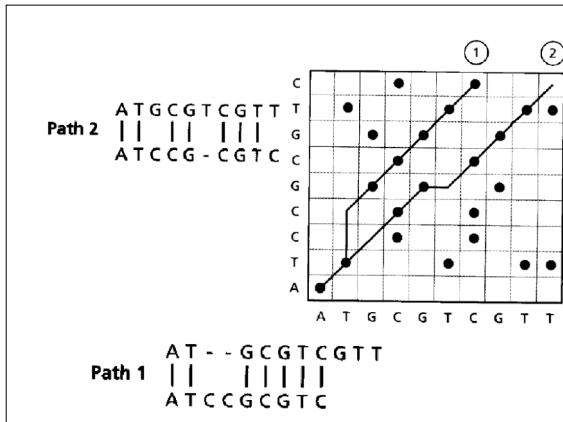
```

Sequences producing significant alignments:	Score	E
	(bits)	Value

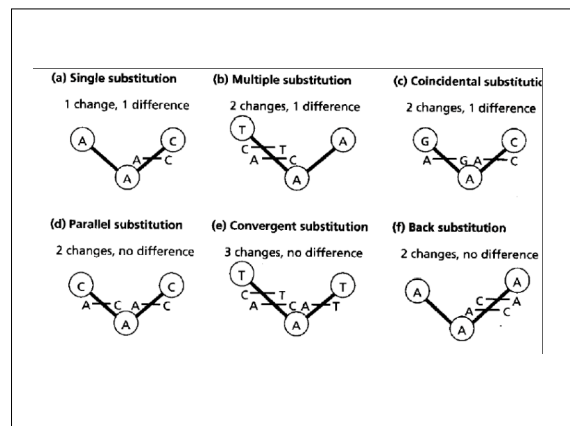
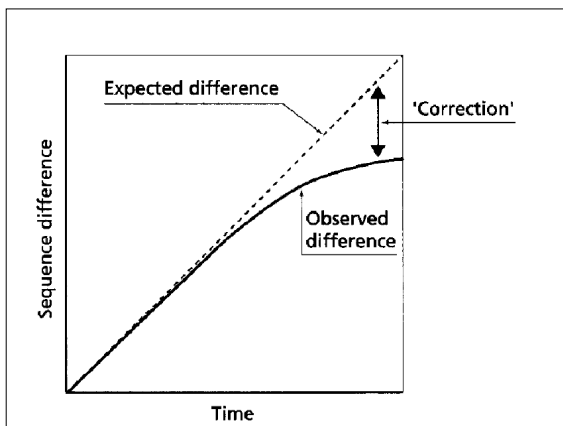
Sequences producing significant alignments:		Score E	Value
		(bits)	
gi 25314548 ref NM_019715.3 	Homo sapiens zona pellucida glyco-	823	0.0
gi 140493 ref U00609.1 	Homo sapiens ZP3 protein (ZP3);	821	0.0
gi 292790 emb J55677.1 	H.sapiens mRNA for ZP3 gene	820	0.0
gi 1475567 emb X82639.1 	M. radiata mRNA for zona pellu-	776	0.0
gi 30133887 ref Z22644.1 	Macaca fascicularis ZPC protein	763	0.0
gi 63833 ref J57825.1 	ZP3=perm receptor protein [ma]	753	0.0
gi 16552292 db X650478.1 	Homo sapiens cDNA FLJ32226 fis,	634	e-174
gi 267342 ref J22469.1 	Pig zona pellucida glycoprote-	598	e-167
gi 83309 ref U004565.1 	Pig mRNA for zona pellucida-	577	e-163
gi 133353 ref U004565.1 	Cat mRNA for zona pellucida	575	e-162
gi 633051 ref U045070.1 	Dog mRNA for zona pellucida-	575	e-162
gi 158474 ref U027678.1 	Felis catus zona pellucida C-	570	e-161
gi 133353 ref U004565.1 	Cat mRNA for zona pellucida	569	e-161
gi 13333338 ref NM_137974.2 	Ms taurus zona pellucida glyco-	563	e-159
gi 1458926 ref U05775.1 	Ms taurus zona pellucida C-	563	e-159
gi 1458926 ref U05782.1 	Oryetolagus cuniculus zona p-	558	e-157
gi 9212413 ref U078395.1 	Pseudomys australis zona pellicu-	555	e-156
gi 133353 ref U004565.1 	Cat mRNA for zona pellucida	552	e-155
gi 17556082 ref NM_011776.1 	Ms musculus zona pellucida gly-	549	e-154
gi 1663731 ref U020076.1 	Mouse zona pellucida sperm-bi-	549	e-154
gi 21435514 ref AF151621.1 	Lepus lagurus zona pellucida g-	547	e-154

Sequences producing significant alignments:	Score	E
	(bits)	Value
gi 338852221 tRNA-Glu-211 Homo sapiens	315	6e-54

Sequences producing significant alignments:		Score	E Value
		(bits)	
gq_123497582.gem	ACQ935454.1	Fan troglodytes clone UNK R4-84M6...	215 6e-54
gq_123497586.gem	ACI16566.1	Hom sapiens chromosome UNK clon...	215 6e-54
gq_123497587.gem	ACI16566.1	Homo sapiens chromosome UNK clon...	215 6e-54
gq_123497588.gem	ACQ93130.1	Papio anubis clone RP41-48D0E, W...	202 4e-50
gq_123497474.gem	ACQ93142.1	Casus familiaris clone RP81-92D9...	191 9e-50
gq_123497589.gem	ACQ93142.1	Felis tatus clone RP82-88D19, W...	189 1e-50
gq_123497472.gem	ACQ9321.3	Felis cat clone RP86-21M22, WO...	125 4e-35
gq_123497479.gem	ACR87420.3	Mus musculus chromosome 5 clone ...	122 3e-33
gq_123497543.gem	ACQ93728.2	Rattus norvegicus clone CH230-1...	122 3e-35
gq_123578851.gem	ACQ127408.1	Rattus norvegicus clone CH230-1...	122 5e-26
gq_130581181.gem	ACI130196.3	Rattus norvegicus clone CH230-20...	122 5e-26
gq_130581182.gem	ACI130196.3	Rattus norvegicus clone CH230-20...	122 5e-26
gq_13175511.gem	ACQ106885.1	Hom sapiens chromosome UNK clon...	99 6e-19
gq_123607851.gem	ACI16561.1	Hom sapiens chromosome UNK clon...	99 6e-19
gq_9931702.gem	ACI161642.3	Hom sapiens chromosome 1 clone ...	89 7e-16



One thing we did not consider was multiple substitutions



More sensitive searches

- Profiles
 - Position Specific Scoring Matrix (PSSM)
- PSI-Blast
 - -automated PSSM
- Hidden Markov Models (HMMs)

PAM250 substitution matrix
Matrices can be optimized for different distances
Newer matrices are BLOSUM

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

- Made from alignment from one gene family.
- Substitution matrix for each position in alignment.
- Incorporates the types of changes that occur for that protein family at each position.
- Still uses general substitution matrix since the alignment would not cover all possible changes
- Changes gap penalty based upon alignment (remember ClustalW).

[illegible][illegible]

PAM 250 (Seq. 1 vs 2): 0

```
Seq 1: Lys Ser Thr Val Ser Lys
Seq 2: Asp Val Val Val Ser Arg
```



```
Seq 1: Lys Ser Thr Val Ser Lys
Seq 2: Asp Val Val Val Ser Arg
```

PAM 250 (Seq. 1 vs 2): 0 -1

Seq 1: Lys Ser Thr Val Ser Lys
Seq 2: Asp Val Val Val Ser Arg

PAM 250 (Seq. 1 vs 2): 0 -1 0

Seq 1: Lys Ser Thr Val Ser Lys
Seq 2: Asp Val Val Val Ser Arg

PAM 250 (Seq. 1 vs 2): 0 -1 0 4 2 3 = 8

Alignment of sp18.

```
red18      MRSVLVLLCVLLMAICAADK-----KSTVSKENAAMAKVMKIFLDSRTDRFK-KRIE 51
white18    MRSVLVLLCVLLMAICAADK-----KTTYSKENAAAMIAKIMKIFLDAGGPKF-KRFE 51
asin18     MRSVLVLLCVLLMAICAADK-----KTSYSEKENAAMVKMMPKLDMKAGPK-EYIE 51
pink18     MRFLLLLCLVLGAVSQVQCKRPNVWGKIVVKEKNKAAMKIGPMEYLDALVKLPKRHLWL 60
green18    MRSVLVLLCVLMGVGCAVPD-----DVTVSRQEQSQTSGRMVNFLDEEMHKLV-KRFR 51
**.*.:*****:.*. * . . :. :. :. :. :. :. :.
```

```

red18      KIGYPITPPQYTLLYYNRRRLMDWCHNYVEVSKKIIILGGKNLKNKPFMRGRIIGWKN 111
white18    NMGVYPTPPQWTLLYYNRRRLMECHTSPFVSKKIIILGGKNLKNKPFMRGRIIGWKN 111
asim18     DMGYPTTPQWTLLYYNRRRLMECHTSPFVSKKIIILGGKNLKNKPFMRGRIIGWKS 111
pink18     GANWKLQKFPBDEMRYLAIKRLIKVCHGYTIWSQRILMLKYRPLNEKYFKKVGRYLAWN 120
green18    DMRNWLGPGFVFLKLVNRRRMYRCMDYARSKSLQLQKHLFVNKKTLTKMGFVGYRN 111

```

```
red18      QWILKRRQWHMVR--VMRRYKASIAKKIVAMKVADLPCN----- 149
white18    QWVLKRRQWEMVR--VMRRYKSTAIKKIVAMKVADLPCN----- 149
asimi18    QWAVRQGWGMVR--VSRHSTSTAIKKIVAMKVADLPCN----- 149
pink18     YLIVFRMWIGVLKK--NLKRSEITKPMQKLLDTKDGLPCPVRIKH 165
green18    YGVIRLYADVPRDVGQFPGPKMTAAMRKYSKDPGTFPCCKNEKRRG 158
```

Age	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100	101	102	103	104	105	106	107	108	109	110	111	112	113	114	115	116	117	118	119	120	121	122	123	124	125	126	127	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239	240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255	256	257	258	259	260	261	262	263	264	265	266	267	268	269	270	271	272	273	274	275	276	277	278	279	280	281	282	283	284	285	286	287	288	289	290	291	292	293	294	295	296	297	298	299	300	301	302	303	304	305	306	307	308	309	310	311	312	313	314	315	316	317	318	319	320	321	322	323	324	325	326	327	328	329	330	331	332	333	334	335	336	337	338	339	340	341	342	343	344	345	346	347	348	349	350	351	352	353	354	355	356	357	358	359	360	361	362	363	364	365	366	367	368	369	370	371	372	373	374	375	376	377	378	379	380	381	382	383	384	385	386	387	388	389	390	391	392	393	394	395	396	397	398	399	400	401	402	403	404	405	406	407	408	409	410	411	412	413	414	415	416	417	418	419	420	421	422	423	424	425	426	427	428	429	430	431	432	433	434	435	436	437	438	439	440	441	442	443	444	445	446	447	448	449	450	451	452	453	454	455	456	457	458	459	460	461	462	463	464	465	466	467	468	469	470	471	472	473	474	475	476	477	478	479	480	481	482	483	484	485	486	487	488	489	490	491	492	493	494	495	496	497	498	499	500	501	502	503	504	505	506	50
-----	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	----

```
Pos. : 1   2   3   4   5   6
Seq 1:Lys Ser Thr Val Ser Lys
Seq 2:Asp Val Val Val Ser Arg
```

Profile vs. Sequence 1: 24
Profile vs. Sequence 2: 17

Compute each sequence comparison to Profile

Pos. : 1 2 3 4 5 6
 Seq 1:Lys Ser Thr Val Ser Lys
 Seq 2:Asp Val Val Val Ser Arg

Profile vs. Sequence 1:24 0
 Profile vs. Sequence 2:17 17

Pos. : 1 2 3 4 5 6
 Seq 1:Lys Ser Thr Val Ser Lys
 Seq 2:Asp Val Val Val Ser Arg

PAM 250 (Seq. 1 vs 2) : 0 -1 0 4 2 3 = 8
 Profile vs. Sequence 1:24 0 8 34 16 24 = 106
 Profile vs. Sequence 2:17 17 18 34 16 26 = 128

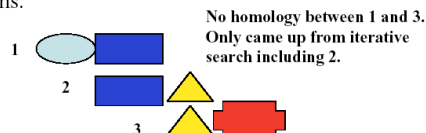
Look at alignment and see why scores are different for the 3 valines

Psi-Blast

- Use BlastP to search protein database.
- Blast makes alignments with all hits with an E-value lower than given value (default e^{-3}).
- Calculates position specific scoring matrix and compares to the database.
- Single gap penalty (different from profiles)
- New "hits" are used to recalculate PSSM and redo the search.
- Searches continue to be done iteratively until no new hits are found.

Problems with iterative searching with Psi-Blast

- You do not make the multiple alignment. There could be errors in the alignment.
- Can pick up significant hits due to shared domains.



Summary

- Blastp may miss some significant matches.
- Using evolutionary information in the form of a multiple alignment increases sensitivity of database searching.

1. What kind of search does BlastX perform? Why would you use BlastX?
1. Describe a scoring matrix (i.e. PAM250).
2. A blast score has an E-value of e^{-3} , is this statistically significant? Is it biologically significant?
3. I want to search for EST sequences, should I search the DNA nr database?
4. BlastP replaced part of my sequence with X's, why did it do that?
5. Name of few databases you would search using BlastN.
6. Describe the type of search that is performed by BlastN, BlastP, and BlastX?
7. What does an E-value from blast mean?
8. Define homologous, orthologous, and paralogous.
9. You have a protein coding sequence and want to compare it to the database. Should you search the protein or DNA sequence? Why?
10. You are told a sequence alignment has a raw score of 1000. Is this a significant match?
11. What does the filter do in a Blast search?
12. The bit score includes the parameters k and lamda (λ). Why are these parameters in the equation to calculate bit scores?
13. The E-value score includes parameters m and n. Why are these parameters in the equation to calculate E-values?
14. What is a position specific scoring matrix?
15. How does Psi-Blast differ from BLASTP?