

Sequence alignment Sekwence alignment Sequence alinement

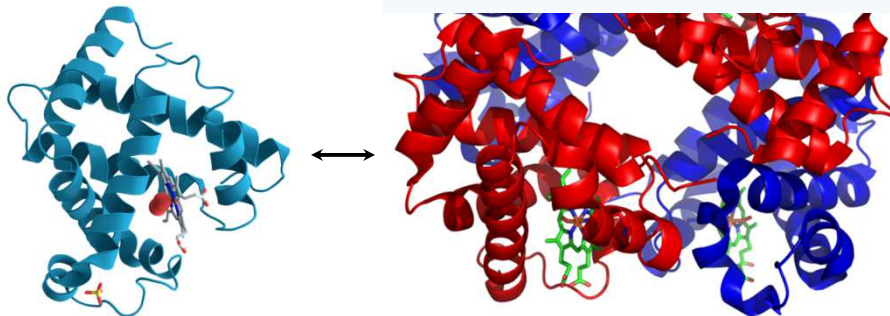
(Drawing heavily from Durbin *et al.*, *Biological Sequence Analysis*)

Systems Biology / Bioinformatics
Edward Marcotte, Univ of Texas at Austin

Typically, to be “biologically related” means to share a common ancestor. In biology, we call this *homologous*.

Two proteins sharing a common ancestor are said to be *homologs*. Homology often implies structural similarity & sometimes (not always) sequence similarity. **A statistically significant sequence or structural similarity can be used to infer homology (common ancestry).**

e.g., Myoglobin & Hemoglobin

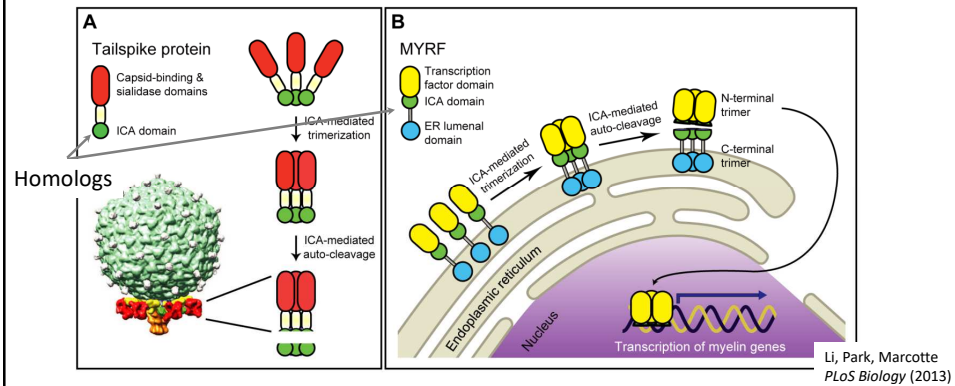


“X-Ray data suggest that the globin chain has the same configuration in the myoglobins and haemoglobins of all vertebrates.”
Kendrew, Perutz, HC Watson. JMB (1965) 13, 669-678

http://en.wikipedia.org/wiki/File:Myoglobin.png&File:1GZX_Haemoglobin.png

In practice, searching for sequence or structural similarity is one of the most powerful computational approaches to discover a gene's function. We can often gain insight about a protein from its homologs.

For example, my lab discovered that myelinating the neurons in your brain reuses the same biochemical mechanism that phage use to make capsids. The key breakthrough was recognizing that the human and phage proteins contained homologous domains.

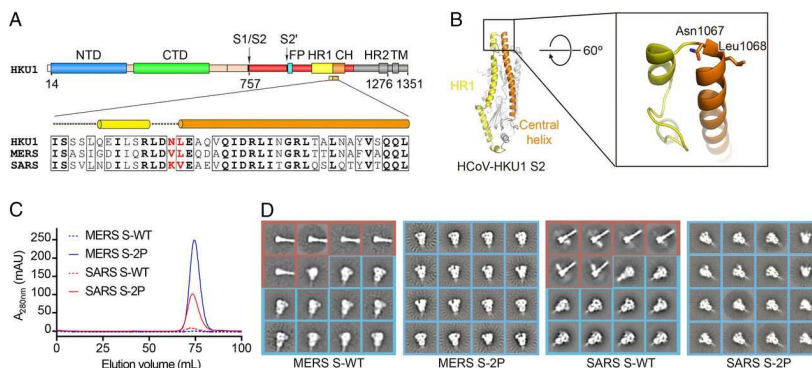


& here's the "trillion dollar" paper from the McLellan lab that the SARS-CoV-2 vaccines are designed from based on homology to MERS and SARS spike antigens

Immunogenicity and structures of a rationally designed prefusion MERS-CoV spike antigen

Jesper Pallelesen, Nanshuang Wang, Kizmekia S. Corbett, Daniel Wrapp, Robert N. Kirchdoerfer, Hannah L. Turner, Christopher A. Cottrell, Michelle M. Becker, Lingshu Wang, Wei Shi, Wing-Pui Kong, Erica L. Andres, Arminia N. Vothschach, Mark R. Denison, James D. Chappell, Barney S. Graham, Andrew B. Ward, and Jason S. McLellan

PNAS August 29, 2017 114 (35): E7348-E7357; first published August 14, 2017; <https://doi.org/10.1073/pnas.1707304114>



Sequence alignment algorithms such as BLAST, PSI-BLAST, FASTA, MMSeqs2, the Needleman–Wunsch & Smith-Waterman algorithms are arguably some of the most important driver technologies of modern biology and underlie the sequencing revolution.

So, let's start learning bioinformatics algorithms by learning how to align two protein sequences.

Live demo:

http://blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastp&BLAST_PROGRAMS=blastp&PAGE_TYPE=BlastSearch&SHOW_DEFAULTS=on&LINK_LOC=blasthome

MVLSPADKTNVKAAWGKVGGAHAGEYGAEALERMFLSFPTTKTYFPHFDLSHGSAQVKGHG
KKVADALTNAAVHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLAAHLPAEFTP
AVHASLDKFLASVSTVLTSKYR

The next few slides show the data from searching this dbase
(#'s may be a bit different from the live version):

Title: clustered nr

Description: **ClusteredNR** is derived from the protein nr database by clustering sequences at 90% identity and 90% length. more...

Molecule Type: Protein

Update date: 2025/12/27

Number of sequences: 470,748,714

BLAST searches a pre-clustered dataset to give more taxonomic diversity to your results

Clusters

Graphic Summary

Alignments

Taxonomy

Clusters producing significant alignments

Download

Select columns

Show 5000

select all 4502 clusters selected

GenPept

Graphics

Distance tree of results

Multiple alignment

MSA Viewer

Cluster Composition	Cluster Ancestor	Cluster Representative Sequence	Max Score	Total Score	Query Cover	E value	Per. Ident	Acc. Len	Accession																																		
<div> <div>1 member(s), 1 organism(s)</div> <div>human</div> <div>Chain A, HEMOGLOBIN-BASED BLOOD SUBSTITUTE (Ho...</div> <div>284 567</div> <div>99%</div> <div>2e-94</div> <div>100.00%</div> <div>283</div> <div>1ABW_A</div> </div>																																											
<div> <div>1ABW_A (human) - 6 member(s)</div> <div> <div>Reports</div> <div>Cluster Members</div> <div>Cluster Taxonomy</div> <div>Download</div> </div> <div> <div>BLAST Alignment</div> <div>Tree View</div> <div>MSA Viewer</div> <div>Multiple alignment</div> </div> <div> <div>1 - 6 cluster member(s)</div> <table> <thead> <tr> <th>Accession</th> <th>Scientific Name</th> <th>Common Name</th> <th>Taxid</th> <th>Additional Proteins</th> </tr> </thead> <tbody> <tr> <td>1ABW_A</td> <td>Homo sapiens</td> <td>human</td> <td>9606</td> <td>4 or more identical proteins</td> </tr> <tr> <td>1C7D_A</td> <td>Homo sapiens</td> <td>human</td> <td>9606</td> <td>N/A</td> </tr> <tr> <td>1O1J_A</td> <td>Homo sapiens</td> <td>human</td> <td>9606</td> <td>N/A</td> </tr> <tr> <td>1O1L_A</td> <td>Homo sapiens</td> <td>human</td> <td>9606</td> <td>N/A</td> </tr> <tr> <td>1O1M_A</td> <td>Homo sapiens</td> <td>human</td> <td>9606</td> <td>N/A</td> </tr> <tr> <td>1O1N_A</td> <td>Homo sapiens</td> <td>human</td> <td>9606</td> <td>N/A</td> </tr> </tbody> </table> </div> </div>	Accession	Scientific Name	Common Name	Taxid	Additional Proteins	1ABW_A	Homo sapiens	human	9606	4 or more identical proteins	1C7D_A	Homo sapiens	human	9606	N/A	1O1J_A	Homo sapiens	human	9606	N/A	1O1L_A	Homo sapiens	human	9606	N/A	1O1M_A	Homo sapiens	human	9606	N/A	1O1N_A	Homo sapiens	human	9606	N/A								
Accession	Scientific Name	Common Name	Taxid	Additional Proteins																																							
1ABW_A	Homo sapiens	human	9606	4 or more identical proteins																																							
1C7D_A	Homo sapiens	human	9606	N/A																																							
1O1J_A	Homo sapiens	human	9606	N/A																																							
1O1L_A	Homo sapiens	human	9606	N/A																																							
1O1M_A	Homo sapiens	human	9606	N/A																																							
1O1N_A	Homo sapiens	human	9606	N/A																																							
<div> <div>10 member(s), 6 organism(s)</div> <div>bat</div> <div>PREDICTED_hemoglobin subunit alpha [Miniopterus natalen...</div> <div>254 254</div> <div>100%</div> <div>1e-84</div> <div>86.62%</div> <div>142</div> <div>XP_016071731.1</div> </div>																																											
<div> <div>56 member(s), 32 organism(s)</div> <div>even-toed ungulates & w...</div> <div>hemoglobin subunit alpha [Bos taurus]</div> <div>252 252</div> <div>100%</div> <div>7e-84</div> <div>88.03%</div> <div>142</div> <div>NP_001070890.2</div> </div>																																											
<div> <div>4 member(s), 4 organism(s)</div> <div>odd-toed ungulates</div> <div>RecName: Full=Hemoglobin subunit alpha; AltName: Full=Al...</div> <div>251 251</div> <div>99%</div> <div>1e-83</div> <div>87.23%</div> <div>141</div> <div>P09906.1</div> </div>																																											
<div> <div>1 member(s), 1 organism(s)</div> <div>eastern chipmunk</div> <div>RecName: Full=Hemoglobin subunit alpha [Tamias striatus]</div> <div>251 251</div> <div>99%</div> <div>2e-83</div> <div>87.23%</div> <div>141</div> <div>B3EWD9.1</div> </div>																																											

Clusters

Graphic Summary

Alignments

Taxonomy

hover to see the title

click to show alignments

Show Conserved Domains

Alignment Scores

< 40

40 - 50

50 - 80

80 - 200

>= 200

4502 clusters selected

Putative conserved domains have been detected, click on the image below for detailed results.

Query

heme binding site

tetramer interface

Hb-alpha-like

Globin-like

Distribution of the top 4673 Blast Hits on 4502 subject clusters

Query

1 20 40 60 80 100 120 140

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

4

Clusters
Graphic Summary
Alignments
Taxonomy

Alignment view
Pairwise
Restore defaults

4502 clusters selected

Download
GenPept
Graphics
sort by: E value

Chain A, HEMOGLOBIN-BASED BLOOD SUBSTITUTE [Homo s
Sequence ID: 1ABW_A Length: 283 Number of Matches: 2

Range 1: 143 to 283
GenPept
Graphics

Score	Expect	Method	Identities	Positiv
284 bits(726)	2e-94	Compositional matrix adjust.	141/141(100%)	141/
Query 2	VLSPADKTNVKAAGKVGAGHAGEYGAEALERMFLSFPTTKTYFPHFDLSH			
Sbjct 143	VLSPADKTNVKAAGKVGAGHAGEYGAEALERMFLSFPTTKTYFPHFDLSH			
Query 62	KVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLA			
Sbjct 203	KVADALTNAVAHVDDMPNALSALSDLHAHKLRVDPVNFKLLSHCLLVTLA			
Query 122	VHASLDKFLASVSTVLTSKYR 142			
Sbjct 263	VHASLDKFLASVSTVLTSKYR 283			
Sbjct 121	AVHASLDKFLASVSTVLTSKYR 142			

If you're curious why the top hit for hemoglobin is a "blood substitute"...

LETTERS TO NATURE

A human recombinant haemoglobin designed for use as a blood substitute

Douglas Looker, Debbie Abbott-Brown, Paul Cozart, Steven Durfee, Stephen Hoffman, Antony J. Mathews, Jeanne Miller-Roechrich, Steven Shoemaker, Stephen Trimble, Giulio Fermi*, Noboru H. Komiyama*, Kiyoshi Nagai* & Gary L. Stetler†

Somatogen Inc., 5797 Central Avenue, Boulder, Colorado 80301, USA
MRC Laboratory of Molecular Biology, Hills Road, Cambridge CB2 2QH, UK

THE need to develop a blood substitute is now urgent because of the increasing concern over blood-transmitted viral and bacterial pathogens¹. Cell-free haemoglobin solutions^{2,3} and human haemoglobin synthesized in *Escherichia coli*⁴ and *Saccharomyces cerevisiae*⁵ have been investigated as potential oxygen-carrying substitutes for red blood cells. But these haemoglobins cannot be used as a blood substitute because (1) the oxygen affinity in the absence of 2,3-bisphosphoglycerate is too high to allow unloading of enough oxygen in the tissues⁶, and (2) they dissociate into $\alpha\beta$ dimers⁷ that are cleared rapidly by renal filtration⁸⁻¹⁰, which can result in long-term kidney damage¹¹. We have produced a human haemoglobin using an expression vector containing one gene encoding a mutant β -globin with decreased oxygen affinity and one duplicated, tandemly fused α -globin gene. Fusion of the two

Protein sequence alignment

Two biologically related proteins with similar sequences:

```

FlgA1 EAGNVKLKRGRLDTLPPTVLDINQLVDAISLRDLSPDQPIQLTQFRQAWRVKAGQRVNVIASGD
      ++K+K+GRLDTLPP +L+ N A+SLR ++ QP+ R+ W +KAGQ V +A G+
FlgA2 TLQDIKMKQGRDLTLPPGALLEPNFAQGAVSLRQINAGQPLTRNMLRRLWI IKAGQDVQVLALGE

```

Also biologically related (& fold up into the same 3D protein structure):

```

FlgA1 EAGNVKLKRGRLDTLPPTVLDINQLVDAISLRDLSPDQPIQLTQFRQAWRVKAGQRVNVIASGD
      A + P +L I+ R L P + I R+AW V+ G V V
FlgA3 LAALKQVTTLIAGKHKPDAMATHAEELQGKIAKRTLLPGRYIPTAAIREAWLVEQGAQVFFIAG

```

But these are biologically unrelated (& fold up into unrelated structures):

```

FlgA1 AGNVKLKRGRLDTLPPTVLDINQLVDAISLRDLSPDQPIQLTQFRQA-WRVKAGQRVNVIASGD
      AG+V K G + + PRT ++ I+ P PI +++A WRV A + V V+ GD
HvcPP AGHV--KNGTMRIVGPRTCSNVWNGTFPINATTGPSIPI PAPNYKKALWRVSATEYVEVVRVGD

```

(FYI, we'll draw examples from Durbin *et al.*, *Biological Sequence Analysis*, Ch. 1 & 2).

To align two sequences, we need to perform 3 steps:

- 1. We need some way to decide which alignments are better than others.
For this, we'll invent a way to give the alignments a "score" indicating their quality.**
- 2. Align the two proteins so that they get the best possible score.**
- 3. Decide if the score is "good enough" for us to believe the alignment is biologically significant.**

To align two sequences, we need to perform 3 steps:

- 1. We need some way to decide which alignments are better than others.
For this, we'll invent a way to give the alignments a "score" indicating their quality.**
- 2. Align the two proteins so that they get the best possible score.**
- 3. Decide if the score is "good enough" for us to believe the alignment is biologically significant.**

We'll treat mutations as independent events.

This allows us to create an ***additive scoring scheme***.

The score for a sequence alignment will be the sum of the scores for aligning each of the individual positions in two sequences.

What kind of mutations should we expect?

Substitutions, insertions and deletions.

Insertions and deletions can be treated as equivalent events by considering one or the other sequence as the reference, and are usually called ***gaps***.

Diagram illustrating sequence alignment and mutation types:

```
AGNVKLKRG
AG+V    K G
AGHV- -KNG
```

Arrows point from the labels ***substitution*** and ***gap*** to the corresponding positions in the alignment:

- substitution*** points to the 'V' in the second sequence (AG+V) and the 'V' in the third sequence (AGHV).
- gap*** points to the '-' in the third sequence (AGHV- -KNG).

Let's consider two models:

First, a **random** model, where amino acids in the sequences occur independently at some given frequencies.

The probability of observing an alignment between x and y is just the product of the frequencies (q) with which we find each amino acid.

We can write this as:

$$P(x, y | R) = \prod_i q_{x_i} \prod_j q_{y_j}$$

What does the capital pi mean?

What's this mean?

What's this mean?

i is just a counter indicating the sequence position

Here's our pair of proteins from before:

FlgA1 EAGNVKLKRGRDLTPPRTVLVDINQLVDAISLRDLSPDQPIQLTQFRQAWRVKAGQRVNVIASGD
 ++K+K+GRDLTLPP +L+ N A+SLR ++ QP+ R+ W +KAGQ V V+A G+
FlgA2 TLQDIKMKQGRDLTPPGALLEPNFAQGAVSLRQINAGQPLTRNMLRRLWI IKAGQDVQVLALGE

So, our random model is:

$$P(x, y | R) = \prod_i q_{x_i} \prod_j q_{y_j} = \underbrace{f(E)*f(A)*f(G)*\dots*f(G)*f(D)*f(T)*f(L)*f(Q)*\dots*f(G)*f(E)}_{\text{frequencies of each amino acid in protein 1 \& 2}}$$

Second, a **match** model, where amino acids at a given position in the alignment arise from some common ancestor with a probability given by the joint probability p_{ab} .

So, under this model, the probability of the alignment is the product of the probabilities of seeing the individual amino acids aligned.

We can write that as:

What does the capital pi mean again?

$$P(x, y | M) = \prod_i p_{x_i, y_i}$$

What's this mean?

What's this mean?

Here's our pair of proteins from before:

FlgA1 EAGNVKLKRGRDLTPPRTVLVDINQLVDAISLRDLSPDQPIQLTQFRQAWRVKAGQRVNVIASGD
 ++K+K+GRDLTLPP +L+ N A+SLR ++ QP+ R+ W +KAGQ V V+A G+
FlgA2 TLQDIKMKQGRDLTPPGALLEPNFAQGAVSLRQINAGQPLTRNMLRRLWIIKAGQDVQVLALGE

So, our match model is:

$$P(x, y | M) = \prod_i p_{x_i, y_i} = \underbrace{f(\text{E aligned with T}) * f(\text{A aligned with L}) * \dots f(\text{D aligned with E})}_{\text{frequencies of the aligned residue pairs}}$$

To decide which model better describes an alignment, we'll take the ratio:

$$\frac{P(x, y | M)}{P(x, y | R)} = \frac{\prod_i P_{x_i y_i}}{\prod_i q_{x_i} \prod_j q_{y_j}} = \prod_i \frac{P_{x_i y_i}}{q_{x_i} q_{y_i}}$$

What did these mean again?

Such a ratio of probabilities under 2 different models is called an **odds ratio**.

Where else have you heard odds ratios used?

Basically: if the ratio > 1, model *M* is more probable
if < 1, model *R* is more probable.

Now, to convert this to an additive score *S*, we can simply take the logarithm of the odds ratio (called the **log odds ratio**):

$$S = \sum_i s(x_i, y_i)$$

This is just the score for aligning one amino acid with another amino acid:

$$s(a, b) = \log \left(\frac{P_{ab}}{P_a P_b} \right)$$

Here written *a* and *b* rather than *x_i* and *y_i* to emphasize that this score reflects the inherent preference of the two amino acids (*a* and *b*) to be aligned.

Almost done with step 1...

The last trick:

Take a big set of pre-aligned protein sequence alignments (that are correct!) and measure all of the pairwise amino acid substitution scores (the $s(a,b)$'s). Put them in a 20x20 **amino acid substitution matrix** :

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	5	-2	-1	-2	-1	-1	-1	0	-2	-1	-2	-1	-1	-3	-1	1	0	-3	-2	0
R	-2	7	-1	-2	-4	1	0	-3	0	-4	-3	3	-2	-3	-3	-1	-1	-3	-1	-3
N	-1	-1	7	2	-2	0	0	0	1	-3	-4	0	-2	-4	-2	1	0	-4	-2	-3
D	-2	-2	2	8	-4	0	2	-1	-1	-4	-4	-1	-4	-5	-1	0	-1	-5	-3	-4
C	-1	-4	-2	-4	13	-3	-3	-3	-3	-2	-2	-3	-2	-2	-4	-1	-1	-5	-3	-1
Q	-1	1	0	0	-3	7	2	-2	1	-3	-2	2	0	-4	-1	0	-1	-1	-1	-3
E	-1	0	0	2	-3	2	6	-3	0	-4	-3	1	-2	-3	-1	-1	-1	-3	-2	-3
G	0	-3	0	-1	-3	-2	-3	8	-2	-4	-4	-2	-3	-4	-2	0	-2	-3	-3	-4
H	-2	0	1	-1	-3	1	0	-2	10	-4	-3	0	-1	-1	-2	-1	-2	-3	2	-4
I	-1	-4	-3	-4	-2	-3	-4	-4	-4	5	2	-3	2	0	-3	-3	-1	-3	-1	4
L	-2	-3	-4	-4	-2	-2	-3	-4	-3	2	5	-3	3	1	-4	-3	-1	-2	-1	1
K	-1	3	0	-1	-3	2	1	-2	0	-3	-3	6	-2	-4	-1	0	-1	-3	-2	-3
M	-1	-2	-2	-4	-2	0	-2	-3	-1	2	3	-2	7	0	-3	-2	-1	-1	0	1
F	-3	-3	-4	-5	-2	-4	-3	-4	-1	0	1	-4	0	8	-4	-3	-2	1	4	-1
P	-1	-3	-2	-1	-4	-1	-1	-2	-2	-3	-4	-1	-3	-4	10	-1	-1	-4	-3	-3
S	1	-1	1	0	-1	0	-1	0	-1	-3	-3	0	-2	-3	-1	5	2	-4	-2	-2
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	2	5	-3	-2	0
W	-3	-3	-4	-5	-5	-1	-3	-3	-3	-3	-2	-3	-1	1	-4	-4	-3	15	2	-3
Y	-2	-1	-2	-3	-3	-1	-2	-3	2	-1	-1	-2	0	4	-3	-2	-2	2	8	-1
V	0	-3	-3	-4	-1	-3	-3	-4	-4	4	1	-3	1	-1	-3	-2	0	-3	-1	5

This is the **BLOSUM50** matrix.

(The numbers are scaled & rounded off to the nearest integer):

What's the score for aspartate (D) aligning with itself?

How about aspartate with phenylalanine (F)? Why?

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	5	-2	-1	-2	-1	-1	-1	0	-2	-1	-2	-1	-1	-3	-1	1	0	-3	-2	0
R	-2	7	-1	-2	-4	1	0	-3	0	-4	-3	3	-2	-3	-3	-1	-1	-3	-1	-3
N	-1	-1	7	2	-2	0	0	0	1	-3	-4	0	-2	-4	-2	1	0	-4	-2	-3
D	-2	-2	2	8	-4	0	2	-1	-1	-4	-4	-1	-4	-5	-1	0	-1	-5	-3	-4
C	-1	-4	-2	-4	13	-3	-3	-3	-3	-2	-2	-3	-2	-2	-4	-1	-1	-5	-3	-1
Q	-1	1	0	0	-3	7	2	-2	1	-3	-2	2	0	-4	-1	0	-1	-1	-1	-3
E	-1	0	0	2	-3	2	6	-3	0	-4	-3	1	-2	-3	-1	-1	-1	-3	-2	-3
G	0	-3	0	-1	-3	-2	-3	8	-2	-4	-4	-2	-3	-4	-2	0	-2	-3	-3	-4
H	-2	0	1	-1	-3	1	0	-2	10	-4	-3	0	-1	-1	-2	-1	-2	-3	2	-4
I	-1	-4	-3	-4	-2	-3	-4	-4	-4	5	2	-3	2	0	-3	-3	-1	-3	-1	4
L	-2	-3	-4	-4	-2	-2	-3	-4	-3	2	5	-3	3	1	-4	-3	-1	-2	-1	1
K	-1	3	0	-1	-3	2	1	-2	0	-3	-3	6	-2	-4	-1	0	-1	-3	-2	-3
M	-1	-2	-2	-4	-2	0	-2	-3	-1	2	3	-2	7	0	-3	-2	-1	-1	0	1
F	-3	-3	-4	-5	-2	-4	-3	-4	-1	0	1	-4	0	8	-4	-3	-2	1	4	-1
P	-1	-3	-2	-1	-4	-1	-1	-2	-2	-3	-4	-1	-3	-4	10	-1	-1	-4	-3	-3
S	1	-1	1	0	-1	0	-1	0	-1	-3	-3	0	-2	-3	-1	5	2	-4	-2	-2
T	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	2	5	-3	-2	0
W	-3	-3	-4	-5	-5	-1	-3	-3	-3	-3	-2	-3	-1	1	-4	-4	-3	15	2	-3
Y	-2	-1	-2	-3	-3	-1	-2	-3	2	-1	-1	-2	0	4	-3	-2	-2	2	8	-1
V	0	-3	-3	-4	-1	-3	-3	-4	-4	4	1	-3	1	-1	-3	-2	0	-3	-1	5

Using this matrix, we can score any alignment as the sum of scores of individual pairs of amino acids.

For example, the top alignment in our earlier example:

```
FlgA1 EAGNVKLKRGRDTPPRTVLDINQLVDAISLRDLSPDQPIQLTQFRQAWRVKAGQRVNVIASGD
      ++K+K+GRDTPPP +L+ N   A+SLR ++  QP+      R+ W +KAGQ V V+A G+
FlgA2 TLQDIKMKQGRDTPPGALLEPNFAQGAVSLRQINAGQPLTRNMLRRLWIIKAGQDVQVLALGE
```

gets the score:

$$S(\text{FlgA1}, \text{FlgA2}) = -1 - 2 - 2 + 2 + 4 + 6 + \dots = 186$$

We also need to penalize gaps. For now, let's just use a constant penalty d for each amino acid gap in an alignment, *i. e.*:

the penalty for a gap of length $g = -g*d$

PAM



Margaret Dayhoff (1925-1983)
Developed point accepted mutation
matrices = PAM matrices
(& also made the 1 letter aa codes!)

Calibrated for different evolutionary times
PAM- n = n substitutions per 100 residues
e.g. matrices from PAM1 to PAM250
measure PAM1,
calculate higher PAMs from that

Explicit model of evolution
(calculated using a phylogenetic tree)

vs.

BLOSUM



Steve and Jorja Henikoff
Developed BLOSUM matrices

Calibrated for different % identity sequences
BLOSUM- n = for sequences of about n % identity
averages substitution probabilities over
sequence clusters, gives better estimates
for highly divergent cases

Implicit model of evolution
(calculated from blocks of aligned sequences)

To align two sequences, we need to perform 3 steps:

1. We need some way to decide which alignments are better than others.
For this, we'll invent a way to give the alignments a "score" indicating their quality.
2. **Align the two proteins so that they get the best possible score.**
3. Decide if the score is "good enough" for us to believe the alignment is biologically significant.

A sense of scale:

There are $\binom{2n}{n} \approx \frac{2^{2n}}{\sqrt{\pi n}}$ possible global alignments between two sequences of length n if we use gaps

So, with 2 sequences of length 100, that's $> 10^{60}$ possible alignments

We'll use something called **dynamic programming**.

This is **mathematically guaranteed** to find the best scoring alignment, and uses **recursion**. This means problems are broken into sub-problems, which are in turn broken into sub-problems, etc, until the simplest sub-problems can be solved.

We're going to find the best **local** alignment—the best matching internal alignment—without forcing all of the amino acids to align (i.e. to match **globally**).

i.e., this \longrightarrow $\begin{matrix} \text{ATGCAT} \\ \text{ATGCAT} \end{matrix}$

Not this \longrightarrow $\begin{matrix} \text{ACGTTATGCATGACGTA} \\ \text{-C---ATGCAT-----T-} \end{matrix}$

Here's the main idea:

We'll make a **path matrix**, showing the possible alignments and their scores. There are simple rules for how to fill in the matrix.

This will test all possible alignments & give us the top-scoring alignment between the two sequences.

		$i=0$					x					$i=n$
		H	E	A	G	A	W	G	H	E	E	
	0											
P	$\leftarrow j=0$											
A												
W												
y H												
E												
A												
E	$\leftarrow j=m$											

The path matrix will be filled from the top left to the bottom right

Here are the rules:

For a given square in the matrix $F(i,j)$, we look at the squares to its left $F(i-1,j)$, top $F(i,j-1)$, and top-left $F(i-1,j-1)$. Each should have a score.

We consider **3 possible events** & **choose the one scoring the highest**:

(1) x_i is aligned to y_j

$$F(i-1,j-1) + s(x_i, y_j)$$

(2) x_i is aligned to a gap

$$F(i-1,j) - d$$

(3) y_j is aligned to a gap

$$F(i,j-1) - d$$

For this example, we'll use $d = 8$. We also set the left-most & top-most entries to zero.

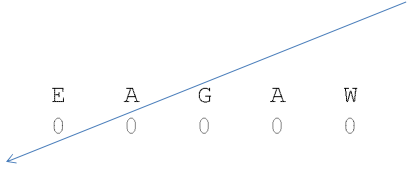
Just two more rules:

If the score is negative, set it equal to zero.

At each step, we also keep track of which event was chosen by **drawing an arrow from the cell we just filled back to the cell which contributed its score to this one.**

That's it! Just repeat this to fill the entire matrix.

Here we go! Start with the borders & the first entry.



		H	E	A	G	A	W	G	H	E	E
	0	0	0	0	0	0	0	0	0	0	0
P	0	0									
A	0										
W	0										
H	0										
E	0										
A	0										
E	0										

Why is this zero?

What's the score from our BLOSSUM matrix for substituting H for P?

Next round!

		H	E	A	G	A	W	G	H	E	E
	0	0	0	0	0	0	0	0	0	0	0
P	0	0	0								
A	0	0	0								
W	0										
H	0										
E	0										
A	0										
E	0										

Terrible! Again, none of the possible give positive scores.

We have to go a bit further in before we find a positive score...

A few more rounds, and a positive score at last!

		H	E	A	G	A	W	G	H	E	E
	0	0	0	0	0	0	0	0	0	0	0
P	0	0	0	0							
A	0	0	0	5							
W	0	0	0								
H	0										
E	0										
A	0										
E	0										

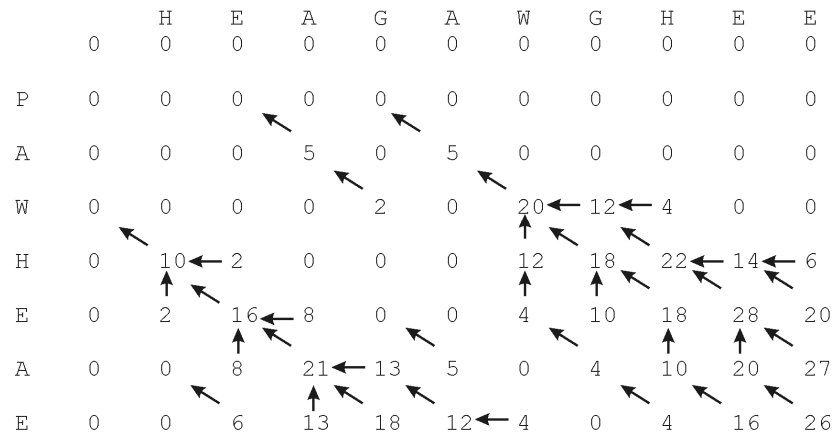
How did we get this one?

& a few more rounds...

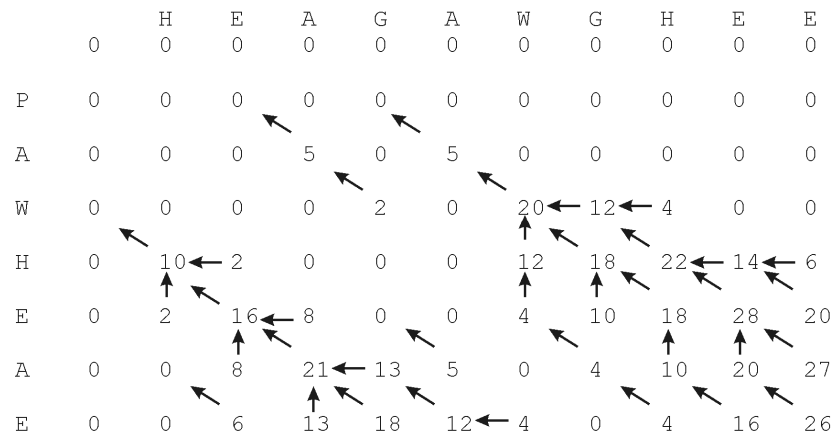
		H	E	A	G	A	W	G	H	E	E
	0	0	0	0	0	0	0	0	0	0	0
P	0	0	0	0	0						
A	0	0	0	5	0						
W	0	0	0	0	2						
H	0	10	2	0	0						
E	0										
A	0										
E	0										

What does this mean?

The whole thing filled in!

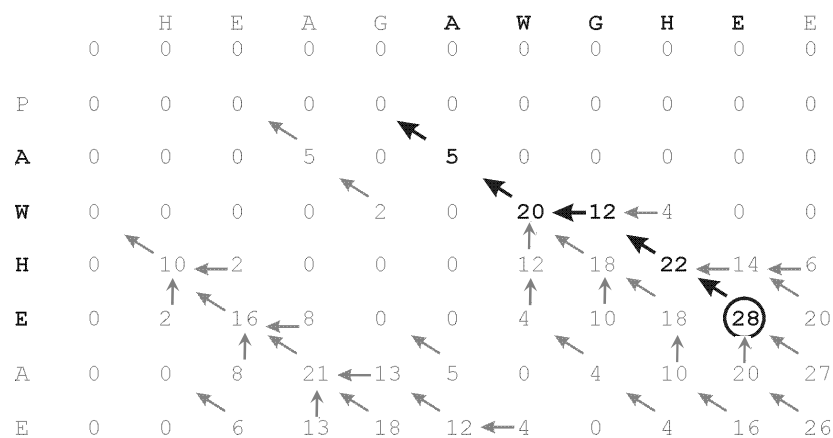


Now, find the optimal alignment using a **traceback** process:
Look for the highest score, then follow the arrows back.
The alignment “grows” from right to left



This gives the following alignment: **AWGHE**
 AW-HE

(Note: for gaps, the arrow points to the sequence that gets the gap)



To align two sequences, we need to perform 3 steps:

1. We need some way to decide which alignments are better than others.
 For this, we'll invent a way to give the alignments a "score" indicating their quality.
2. Align the two proteins so that they get the best possible score.
3. Decide if the score is "good enough" for us to believe the alignment is biologically significant.

This algorithm always gives the best alignment.

Every pair of sequences can be aligned in some fashion.

So, when is a score “good enough”?

How can we figure this out?

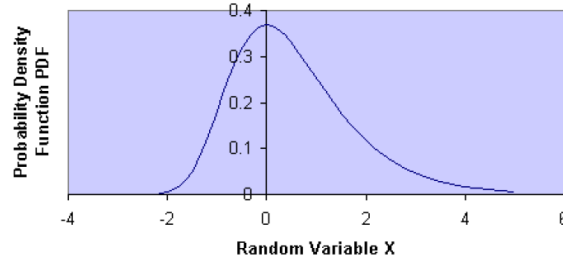
Here's one approach:

**Shuffle one sequence. Calculate the best alignment & its score.
Repeat 1000 times.**

If we never see a score as high as the real one, we say the real score has < 1 in a 1000 chance of happening just by luck.

But if we want something that only occurs < 1 in a million, we'd have to shuffle 1,000,000 times...

Luckily, alignment scores follow a well-behaved distribution, the **extreme value distribution**, so we can do a few trials & fit to this.



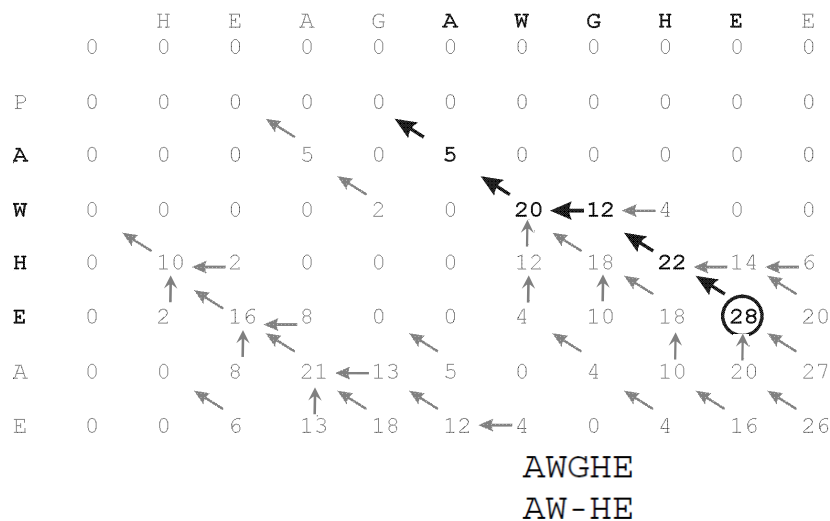
random trials & their average score

$$p(\text{max score} \leq X) \approx e^{-kNe^{\lambda(X-\mu)}}$$

This p-value gives the significance of your alignment.
But, if we search a database and perform many alignments, we still need something more (next time).

Describe the shape & can be fit from a few trials

Some extensions: Local vs. global alignments
How might you force the full sequences to align?



Some extensions: Local vs. global alignments

How might you force the full sequences to align?

A few tiny changes:

Initialize only the top left cell of the path matrix to zero
(not all top and left cells).

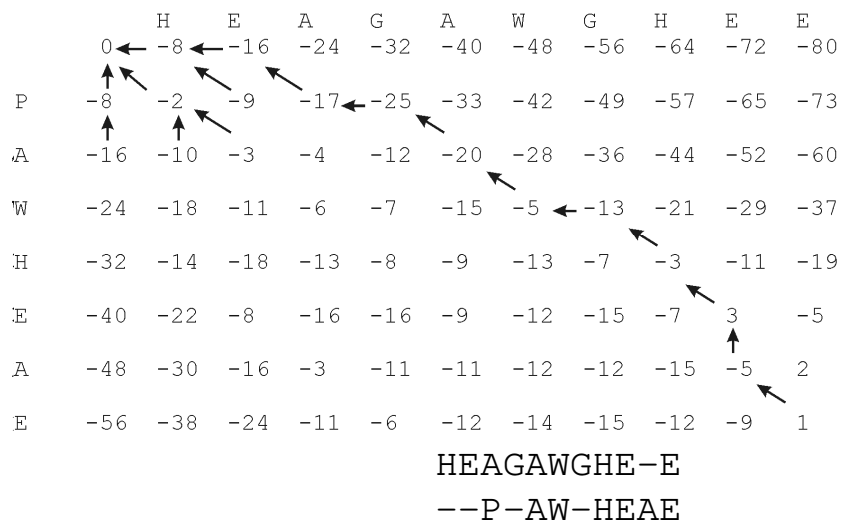
Leave the negative values (don't set them to zero).

The optimal alignment should start at the top left cell and finish at the bottom right cell of the path matrix.

Start the trace-back at the bottom right cell

Some extensions: Local vs. global alignments

How might you force the full sequences to align?



How can you try this yourself using BioPython?

BioPython can perform a wide variety of sequence alignments, DNA/protein, local/global, dynamic programming, BLAST, different scoring schemes, etc, & is a great environment to learn and play with these approaches. Here's a minimal use case to start you off:

```
1 # Here's how to perform pairwise alignments using BioPython,
2 # excerpted from https://biopython.org/DIST/docs/tutorial/Tutorial.html
3
4 # To generate pairwise alignments, first create a PairwiseAligner object:
5 from Bio.Align import PairwiseAligner
6 aligner = PairwiseAligner() # this will use a very minimal default scoring method
7 # However, BioPython knows about more sophisticated schemes
8 # e.g. uncomment the next line to use the BLASTN substitution matrix & gap penalties, which is good for nucleotides:
9 # aligner = PairwiseAligner(scoring="blastn")
10 # other options include megablast (for nucs) and blastp (for proteins)
11
12 aligner.mode = "local" # alternatively, use "global" for a global alignment
13 target = "AGAACTC"
14 query = "GAAGT"
15 score = aligner.score(target, query) # Use aligner.score to calculate the alignment score between 2 sequences:
16 print(score)
17
18 alignments = aligner.align(target, query)
19 for alignment in alignments:
20     print(alignment)
21
22 # BioPython will perform Smith-Waterman for local alignments, Needleman-Wunsch for global
23 # you can confirm which algorithm you used by typing:
24 aligner.algorithm
25
5.0
target      1 GAAGT 6
            0 |||| 5
query       0 GAAGT 5

'Smith-Waterman'
```

Using BioPython, you can change every aspect of the scoring & substitution matrices, as well as run BLAST locally or in the cloud.

e.g. here's the BLOSUM62 matrix, along w/ many others that BioPython knows about:

```
1 from Bio.Align import substitution_matrices
2 substitution_matrices.load()
3 [ 'BENNER22', 'BENNER6', 'BENNER74', 'BLASTN', 'BLASTP', 'BLOSUM45', 'BLOSUM50', 'BLOSUM62', ..., 'TRANS' ]
4 matrix = substitution_matrices.load("BLOSUM62")
5 print(matrix)

# Matrix made by matblas from blosum62.iij
# * column uses minimum score
# BLOSUM Clustered Scoring Matrix in 1/2 Bit Units
# Blocks Database = /data/blocks_5.0/blocks.dat
# Cluster Percentage: >= 62
# Entropy = 0.6979, Expected = -0.5209
A   R   N   D   C   Q   E   G   H   I   L   K   M   F   P   S   T   W   Y   V   B   Z   X   *
A 4.0 -1.0 -2.0 -2.0 0.0 -1.0 -1.0 0.0 -2.0 -1.0 -1.0 -1.0 -1.0 -2.0 -1.0 1.0 0.0 -3.0 -2.0 0.0 -2.0 -1.0 0.0 -4.0
R -1.0 5.0 0.0 -2.0 -3.0 1.0 0.0 -2.0 0.0 -3.0 -2.0 2.0 -1.0 -3.0 -2.0 -1.0 -1.0 -3.0 -2.0 -3.0 -1.0 0.0 -1.0 -4.0
N -2.0 0.0 6.0 1.0 -3.0 0.0 0.0 0.0 1.0 -3.0 -3.0 0.0 -2.0 -3.0 -2.0 1.0 0.0 -4.0 -2.0 -3.0 3.0 0.0 -1.0 -4.0
D -2.0 -2.0 1.0 6.0 -3.0 0.0 2.0 -1.0 -1.0 -3.0 -4.0 -1.0 -3.0 -3.0 -1.0 0.0 -1.0 -4.0 -3.0 -3.0 4.0 1.0 -1.0 -4.0
C 0.0 -3.0 -3.0 -3.0 9.0 -3.0 -4.0 -3.0 -3.0 -1.0 -1.0 -3.0 -1.0 -2.0 -3.0 -1.0 -1.0 -2.0 -2.0 -1.0 -3.0 -3.0 -2.0 -4.0
Q -1.0 1.0 0.0 0.0 -3.0 5.0 2.0 -2.0 0.0 -3.0 -2.0 1.0 0.0 -3.0 -1.0 0.0 -1.0 -2.0 -1.0 -2.0 0.0 3.0 -1.0 -4.0
E -1.0 0.0 0.0 2.0 -4.0 2.0 5.0 -2.0 0.0 -3.0 -3.0 1.0 -2.0 -3.0 -1.0 0.0 -1.0 -3.0 -2.0 -2.0 1.0 4.0 -1.0 -4.0
G 0.0 -2.0 0.0 -1.0 -3.0 -2.0 -2.0 6.0 -2.0 -4.0 -4.0 -2.0 -3.0 -3.0 -2.0 0.0 -2.0 -2.0 -3.0 -3.0 -1.0 -2.0 -1.0 -4.0
H -2.0 0.0 1.0 -1.0 -3.0 0.0 0.0 -2.0 8.0 -3.0 -3.0 -1.0 -2.0 -1.0 -2.0 -1.0 -2.0 -2.0 2.0 -3.0 0.0 0.0 -1.0 -4.0
I -1.0 -3.0 -3.0 -3.0 -1.0 -3.0 -3.0 -4.0 -3.0 4.0 2.0 -3.0 1.0 0.0 -3.0 -2.0 -1.0 -3.0 -1.0 3.0 -3.0 -3.0 -1.0 -4.0
L -1.0 -2.0 -3.0 -4.0 -1.0 -2.0 -3.0 -4.0 -3.0 2.0 4.0 -2.0 2.0 0.0 -3.0 -2.0 -1.0 -2.0 -1.0 1.0 -4.0 -3.0 -1.0 -4.0
K -1.0 2.0 0.0 -1.0 -3.0 1.0 1.0 -2.0 -1.0 -3.0 -2.0 5.0 -1.0 -3.0 -1.0 0.0 -1.0 -3.0 -2.0 -2.0 0.0 1.0 -1.0 -4.0
M -1.0 -1.0 -2.0 -3.0 -1.0 0.0 -2.0 -3.0 -2.0 1.0 2.0 -1.0 5.0 0.0 -2.0 -1.0 -1.0 -1.0 -1.0 1.0 -3.0 -1.0 -1.0 -4.0
F -2.0 -3.0 -3.0 -3.0 -2.0 -3.0 -3.0 -3.0 -1.0 0.0 0.0 -3.0 0.0 6.0 -4.0 -2.0 -2.0 1.0 3.0 -1.0 -3.0 -3.0 -1.0 -4.0
P -1.0 -2.0 -2.0 -1.0 -3.0 -1.0 -1.0 -2.0 -2.0 -3.0 -3.0 -1.0 -2.0 -4.0 7.0 -1.0 -1.0 -4.0 -3.0 -2.0 -2.0 -1.0 -2.0 -4.0
S 1.0 -1.0 1.0 0.0 -1.0 0.0 0.0 0.0 -1.0 -2.0 -2.0 0.0 -1.0 -2.0 -1.0 4.0 1.0 -3.0 -2.0 -2.0 0.0 0.0 0.0 -4.0
T 0.0 -1.0 0.0 -1.0 -1.0 -1.0 -1.0 -2.0 -2.0 -1.0 -1.0 -1.0 -1.0 -2.0 -1.0 1.0 5.0 -2.0 -2.0 0.0 -1.0 -1.0 0.0 -4.0
W -3.0 -3.0 -4.0 -4.0 -2.0 -2.0 -3.0 -2.0 -2.0 -3.0 -2.0 -3.0 -1.0 1.0 -4.0 -3.0 -2.0 11.0 2.0 -3.0 -4.0 -3.0 -2.0 -4.0
Y -2.0 -2.0 -3.0 -2.0 -1.0 -2.0 -3.0 2.0 -1.0 -1.0 -2.0 -1.0 3.0 -3.0 -2.0 -2.0 2.0 7.0 -1.0 -3.0 -2.0 -1.0 -4.0
V 0.0 -3.0 -3.0 -3.0 -1.0 -2.0 -2.0 -3.0 -3.0 3.0 1.0 -2.0 1.0 -1.0 -2.0 -2.0 0.0 -3.0 -1.0 4.0 -3.0 -2.0 -1.0 -4.0
B -2.0 -1.0 3.0 4.0 -3.0 0.0 1.0 -1.0 0.0 -3.0 -4.0 0.0 -3.0 -3.0 -2.0 0.0 -1.0 -4.0 -3.0 -3.0 4.0 1.0 -1.0 -4.0
Z -1.0 0.0 0.0 1.0 -3.0 3.0 4.0 -2.0 0.0 -3.0 -3.0 1.0 -1.0 -3.0 -1.0 0.0 -1.0 -3.0 -2.0 -2.0 1.0 4.0 -1.0 -4.0
X 0.0 -1.0 -1.0 -1.0 -2.0 -1.0 -1.0 -1.0 -1.0 -1.0 -1.0 -1.0 -1.0 -2.0 0.0 0.0 -2.0 -1.0 -1.0 -1.0 -1.0 -1.0 -4.0
* -4.0 -4.0 -4.0 -4.0 -4.0 -4.0 -4.0 -4.0 -4.0 -4.0 -4.0 -4.0 -4.0 -4.0 -4.0 -4.0 -4.0 -4.0 -4.0 -4.0 -4.0 1.0
```

Putting it all together, here's the example alignment we did manually

```

1 from Bio import Align
2 from Bio.Align import substitution_matrices
3
4 aligner = Align.PairwiseAligner()
5 aligner.mode = "local"
6 matrix = substitution_matrices.load("BLOSUM50")
7
8 aligner.substitution_matrix = matrix
9 aligner.target_internal_open_gap_score = -8.000000
10 aligner.target_internal_extend_gap_score = -8.000000
11 aligner.target_left_open_gap_score = -8.000000
12 aligner.target_left_extend_gap_score = -8.000000
13 aligner.target_right_open_gap_score = -8.000000
14 aligner.target_right_extend_gap_score = -8.000000
15 aligner.query_internal_open_gap_score = -8.000000
16 aligner.query_internal_extend_gap_score = -8.000000
17 aligner.query_left_open_gap_score = -8.000000
18 aligner.query_left_extend_gap_score = -8.000000
19 aligner.query_right_open_gap_score = -8.000000
20 aligner.query_right_extend_gap_score = -8.000000
21
22 target = "HEAGAWGHEE"
23 query = "PAWHEAE"
24 score = aligner.score(target, query)
25 print(score)
26
27 alignments = aligner.align(target, query)
28 for alignment in alignments:
29     print(alignment)

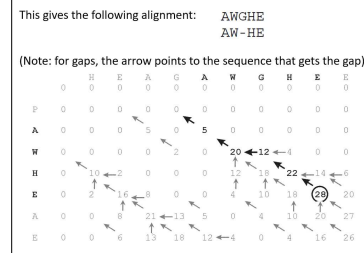
```

```

28.0
target      4 AWGHE 9
            0 ||-|| 5
query      1 AW-HE 5

```

Here was our earlier version:



You can read more about using BioPython for sequence analyses & get example code at:

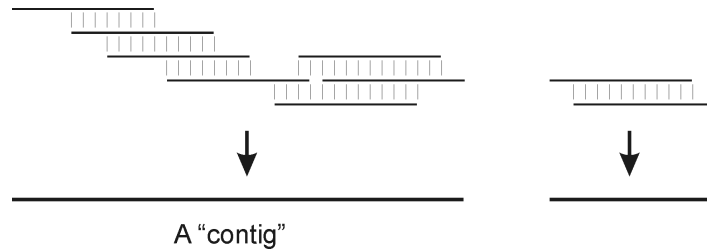
<https://biopython.org/DIST/docs/tutorial/Tutorial.html>

Chapter 7 is all about how to perform pairwise sequence alignments

Some extensions:

What about overlapping sequences?

e.g. as in 'shotgun sequencing' genomes where
'contigs' are built up from overlapping sequences



Some extensions:

What about overlapping sequences?

Modify global alignment to not penalize overhangs:

The optimal alignment should start at the top or left edge
and finish at the bottom or right edge of the path matrix.

Set these boundary conditions :

$$F(i,0) = 0 \text{ for } i=1 \text{ to } n$$

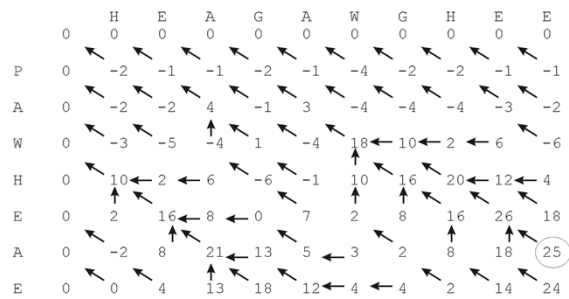
$$F(0,j) = 0 \text{ for } j=1 \text{ to } m$$

Start the traceback at the cell with the highest score on the
right or bottom border

Some extensions:

What about overlapping sequences?

e.g. as in 'shotgun sequencing' genomes where
'contigs' are built up from overlapping sequences



(overhang = HEA)

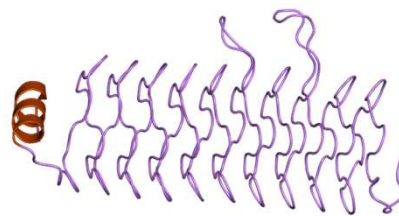
GAWGHEE

PAW-HEA

(overhang = E)

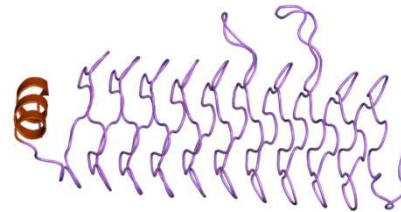
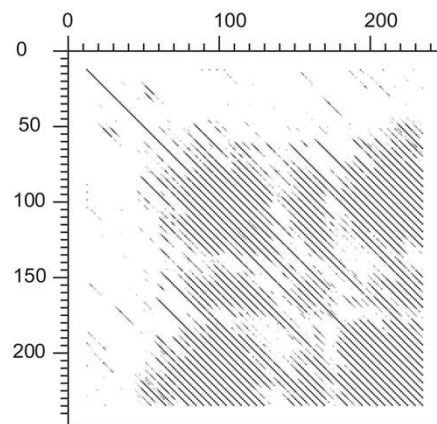
Some extensions:

How might you find repetitive sequences?



Structure of the pentapeptide
repeat protein HetL
(from wiki, PMID18952182)

Align the sequence to itself and ignore the diagonal (optimal) alignment
→ High-scoring off-diagonal alignments will be repeats



Structure of the pentapeptide
repeat protein HgIK
(from wiki, PMID18952182)

Dot plot (quick visualization of
sequence similarity)
of the pentapeptide repeat
protein HgIK protein vs. itself
(http://en.wikipedia.org/wiki/Pentapeptide_repeat)