National Institute of Allergy and Infectious Diseases

**Phylogenetics and Sequence Analysis**

# Lecture 2
# BLAST and Sequence Alignment

**Kurt Wollenberg, PhD**
Phylogenetics Specialist
Bioinformatics and Computational Biosciences Branch
Office of Cyber Infrastructure and Computational Biology

Fall 2015

NIH National Institute of Allergy and Infectious Diseases

NIAID

---

# We Are BCBB!

NIH National Institute of Allergy and Infectious Diseases

- Group of 37
  - Bioinformatics Software Developers
  - Computational Biologists
  - Project Management & Analysis Professionals

NIH National Institute of Allergy and Infectious Diseases

NIAID

2

---

# Course Organization

- Building a clean sequence
- **Collecting homologs**
- **Aligning your sequences**
- Building trees
- Further analyses

NIH National Institute of Allergy and Infectious Diseases

NIAID

4

## Previously

- Hierarchical and genealogical data
- Comparative sequence analysis
- Generating clean sequence
  - Trim vector contamination
  - Trim low-quality ends
  - Align fragment overlap to build contig
  - Export contig (consensus)

NIAID

5

## Today…

**Pairwise sequence alignment**
- How does it work?

**BLAST**
- How does it work?
- The many flavors of BLAST
- Demo

**Multiple Sequence Alignment**
- How does it work?
- Demo
- Inspect and correct your MSA

NIAID

6

## PAIRWISE ALIGNMENT

and **BLAST**: **B**asic **L**ocal **A**lignment **S**earch **T**ool

- Sequence Alignment: Assigning homology to sites among a group of known sequences
- BLAST: Alignment of one sequence with many unknown sequences

NIAID

7

**HOMOLOGY vs. ANALOGY**

common ancestry          convergence



**PAIRWISE ALIGNMENT**

Pairing of <u>sites</u> based on an assessment of homology

<u>Homology</u> assessed using Substitution Matrices



**PAIRWISE ALIGNMENT**

```
HBA_HUMAN  GSAQVKGHGKKVADALTNAVAHVDDMPNALSALSDLHAHKL
           G+ +VK+HGKKV  A+++++AH+D++ +++++LS+LH  KL
HBB_HUMAN  GNPKVKAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKL


HBA_HUMAN  GSAQVKGHGKKVADALTNAVAHV---D--DMPNALSALSDLHAHKL
           ++ ++++H+ KV   + +A  ++           +L+ L+++H+ K
LGB2_LUPLU NNPELQAHAGKVFKLVYEAAIQLQVTGVVVTDATLKNLGSVHVSKG


HBA_HUMAN  GSAQVKGHGKKVADALTNAVAHVDDMPNALSALSD----LHAHKL
           GS+ + G +   +D L ++ H+ D+  A +AL D    ++AH+
F11G11.2   GSGYLVGDSLTFVDLL--VAQHTADLLAANAALLDEFPQFKAHQE
```
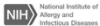
## PAIRWISE ALIGNMENT

Substitution Matrices
➡ Derived mathematically
➡ Derived from data

"A substitution matrix (even one derived by arbitrarily assigning probabilities to pairs) is a statement of the probability of observing these pairs in real alignment."

NIAID

National Institute of Allergy and Infectious Diseases

11

## PAIRWISE ALIGNMENT

DNA Substitution Matrices

• Single parameter - Jukes-Cantor
  - Equal base frequencies
  - Uniform rates of change

• Two parameter - Kimura
  - Equal base probabilities
  - Two rates of change

NIAID

National Institute of Allergy and Infectious Diseases

12

## PAIRWISE ALIGNMENT

DNA Substitution Matrices

• More parameters - HKY
  - Unequal base frequencies
  - Two rates of change

• Fully parameterized - GTR
  - Unequal base probabilities
  - Six rates of change

NIAID

National Institute of Allergy and Infectious Diseases

13

## PAIRWISE ALIGNMENT

Jukes-Cantor Substitution Probabilities

$$P_{ij}(t) = \begin{cases} \dfrac{1}{4} + \dfrac{3}{4} e^{-4\mu t} & i = j \\ \dfrac{1}{4} - \dfrac{1}{4} e^{-4\mu t} & i \neq j \end{cases}$$

NIAID

14

## PAIRWISE ALIGNMENT

Jukes-Cantor Substitution Probabilities

μt = 0.25

|   | A | C | G | T |
|---|---|---|---|---|
| A | 0.5259 | 0.1580 | 0.1580 | 0.1580 |
| C | 0.1580 | 0.5259 | 0.1580 | 0.1580 |
| G | 0.1580 | 0.1580 | 0.5259 | 0.1580 |
| T | 0.1580 | 0.1580 | 0.1580 | 0.5259 |

NIAID

15

## PAIRWISE ALIGNMENT

Kimura Two-Parameter Substitution Model

If the probability of transitions (A ⇔ G, C ⇔ T) is different from the probability of transversions (A ⇔T, G ⇔T, A ⇔C, G ⇔C), then there are two relative rate parameters expressed as the transition/transversion rate ratio κ

NIAID

16

## PAIRWISE ALIGNMENT

Kimura Two-Parameter Substitution Probabilities

$$P_{ij}(t) = \begin{cases} \dfrac{1}{4} - \dfrac{1}{4}e^{-4\mu t} & i \neq j, transversion \\ \dfrac{1}{4} + \dfrac{1}{4}e^{-4\mu t} - \dfrac{1}{2}e^{-2(\kappa+1)\mu t} & i \neq j, transition \\ \dfrac{1}{4} + \dfrac{1}{4}e^{-4\mu t} + \dfrac{1}{2}e^{-2(\kappa+1)\mu t} & i = j \end{cases}$$

NIAID

National Institute of Allergy and Infectious Diseases

17

## PAIRWISE ALIGNMENT

Kimura Two-Parameter Substitution Probabilities

$\mu t = 0.25 \quad \kappa = 2.0$

|   | A | C | G | T |
|---|---|---|---|---|
| A | 0.4535 | 0.1580 | 0.2304 | 0.1580 |
| C | 0.1580 | 0.4535 | 0.1580 | 0.2304 |
| G | 0.2304 | 0.1580 | 0.4535 | 0.1580 |
| T | 0.1580 | 0.2304 | 0.1580 | 0.4535 |

NIAID

National Institute of Allergy and Infectious Diseases

18

## PAIRWISE ALIGNMENT

HKY Substitution Probabilities

$$P_{ij}(t) = \begin{cases} \pi_j + \pi_j\left(\dfrac{1}{\Pi_j} - 1\right)e^{-\mu t} + \left(\dfrac{\Pi_j - \pi_j}{\Pi_j}\right)e^{-\mu t A} & (i = j) \\ \pi_j + \pi_j\left(\dfrac{1}{\Pi_j} - 1\right)e^{-\mu t} + \left(\dfrac{\pi_j}{\Pi_j}\right)e^{-\mu t A} & (i \neq j, transition) \\ \pi_j\left(1 - e^{-\mu t}\right) & (i \neq j, transversion) \end{cases}$$

NIAID

National Institute of Allergy and Infectious Diseases

19

## PAIRWISE ALIGNMENT

HKY Substitution Probabilities

$$\Pi_j = \pi_A + \pi_G \quad \text{if } j \text{ is a purine}$$

$$\Pi_j = \pi_C + \pi_T \quad \text{if } j \text{ is a pyrimidine}$$

$$A = 1 + \Pi_j(\kappa - 1)$$

NIAID

20

## Substitution Models

**Jukes-Cantor**
(One substitution type, equal nucleotide frequencies)

Independent nucl. freq.          Two substitution types  .

**F81/TN82**                    **Kimura 2-Parameter (K2P)**

Two substitution types      Indep. nucl. freq.      Three substitution types   .

**HKY85/F84**                   **Kimura 3 subst. type (K3ST)**

Three substitution types          Six substitution types  .

**Tamura-Nei (TrN)**            **Symmetric (SYM)**

Six substitution types      Independent nucl. freq.   .

**General time-reversible (GTR)**

NIAID

## PAIRWISE ALIGNMENT

Protein Score Matrices
Similarity of Amino Acids



**Amino Acids**
A alanine (ala)
R arginine (arg)
N asparagine (asn)
D aspartic acid (asp)
C cysteine (cys)
Q glutamine (gln)
E glutamic acid (glu)
G glycine (gly)
H histidine (his)
I isoleucine (ile)
L leucine (leu)
K lysine (lys)
M metioneine (met)
F phenyalanine (phe)
P proline (pro)
S serine (ser)
T threonine (thr)
W trytophan (trp)
Y tyrosine (tyr)

From Esquivel RO, et al.. 2013. Advances in Quantum Mechanics, Chapter 27 InTech.

NIAID

22

## PAIRWISE ALIGNMENT

Protein Score Matrices

- Derived from empirical data
- Account for depth of relationship among the data
- Expressed as log-odds ratio:
  - Logarithm of the ratio of the probabilities of two residues being aligned due to homology versus random chance

NIAID

23

## PAIRWISE ALIGNMENT

Protein Score (Substitution) Matrices

The log-odds ratio:
$$s(a,b) = \log(p_{ab}/q_a q_b)$$

$q_a$ = frequency of residue a in the data

$p_{ab}$ = probability that residues a and b have been derived from a common ancestor

NIAID

24

## PAIRWISE ALIGNMENT

Protein Substitution Matrices

- PAM250: Based on phylogenies where all sequences differ by no more than 15%.

- BLOSUM62: Based on clusters of sequences with greater than 62% identical residues.

NIAID

25

## Protein Substitution Matrices

PAM250

```
C  12
S   0   2
T  -2   1   3
P  -3   1   0   6
A  -2   1   1   1   2
G  -3   1   0  -1   1   5
N  -4   1   0  -1   0   0   2
D  -5   0   0  -1   0   1   2   4
E  -5   0   0  -1   0   0   1   3   4
Q  -5  -1  -1   0   0  -1   1   2   2   4
H  -3  -1  -1   0  -1  -2   2   1   1   3   6
R  -4   0  -1   0  -2  -3   0  -1  -1   1   2   6
K  -5   0   0  -1  -1  -2   1   0   0   1   0   3   5
M  -5  -2  -1  -2  -1  -3  -2  -3  -2  -1  -2   0   0   6
I  -2  -1   0  -2  -1  -3  -2  -2  -2  -2  -2  -2  -2   2   5
L  -6  -3  -2  -3  -2  -4  -3  -4  -3  -2  -2  -3  -3   4   2   6
V  -2  -1   0  -1   0  -1  -2  -2  -2  -2  -2  -2  -2   2   4   2   4
F  -4  -3  -3  -5  -4  -5  -4  -6  -5  -5  -2  -4  -5   0   1   2  -1   9
Y   0  -3  -3  -5  -3  -5  -2  -4  -4  -4   0  -4  -4  -2  -1  -1  -2   7  10
W  -8  -2  -5  -6  -6  -7  -4  -7  -7  -6  -3   2  -3  -4  -5  -2  -6   0   0  17
    C   S   T   P   A   G   N   D   E   Q   H   R   K   M   I   L   V   F   Y   W
```
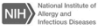
26

## Protein Substitution Matrices

BLOSUM62

```
C   9
S  -1   4
T  -1   1   5
P  -3  -1  -1   7
A   0   1   0  -1   4
G  -3   0  -2  -2   0   6
N  -3   1   0  -2  -2   0   6
D  -3   0  -1  -1  -2  -1   1   6
E  -4   0  -1  -1  -1  -2   0   2   5
Q  -3   0  -1  -1  -1  -2   0   0   2   5
H  -3  -1  -2  -2  -2  -2   1  -1   0   0   8
R  -3  -1  -1  -2  -1  -2   0  -2   0   1   0   5
K  -3   0  -1  -1  -1  -2   0  -1   1   1  -1   2   5
M  -1  -2  -1  -2  -1  -3  -2  -3  -2   0  -2  -1  -1   5
I  -1  -2  -1  -3  -1  -4  -3  -3  -3  -3  -3  -3  -3   1   4
L  -1  -2  -1  -3  -1  -4  -3  -4  -3  -2  -3  -2  -2   2   2   4
V  -1  -2   0  -2   0  -3  -3  -3  -2  -2  -3  -3  -2   1   3   1   4
F  -2  -2  -2  -4  -2  -3  -3  -3  -3  -3  -1  -3  -3   0   0   0  -1   6
Y  -2  -2  -2  -3  -2  -3  -2  -3  -2  -1   2  -2  -2  -1  -1  -1  -1   3   7
W  -2  -3  -2  -4  -3  -2  -4  -4  -3  -2  -2  -3  -3  -1  -3  -2  -3   1   2  11
    C   S   T   P   A   G   N   D   E   Q   H   R   K   M   I   L   V   F   Y   W
```

27

## Protein Substitution Matrices

```
W  -8  -2  -5  -6  -6  -7  -4  -7  -7  -6  -3   2  -3  -4  -5  -2  -6   0   0  17   P250
W  -2  -3   2  -4  -3  -2  -4  -4  -3  -2  -2  -3  -3  -1  -3  -2  -3   1   2  11   B62
    C   S   T   P   A   G   N   D   E   Q   H   R   K   M   I   L   V   F   Y   W
```
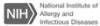
28

## BLAST and Sequence Alignment

How do two sequences get "aligned"?

- Global alignment (Needleman-Wunsch)
    - Assign homology across the entire sequence
    - Clustal

- Local alignment (Smith-Waterman)
    - Assign homology for subsequences
    - MUSCLE and BLAST
    - Good for aligning very divergent sequences

NIAID

NIH National Institute of Allergy and Infectious Diseases

29

## SEQUENCE ALIGNMENT

**HEAGAWGHEE ⇔ PAWHEAE**

Build a matrix of score values for all site pairs

PAM250                          BLOSUM62

```
    H  E  A  G  A  W  G  H  E  E          H  E  A  G  A  W  G  H  E  E
P   0 -1  1  0  1 -5  0  0 -1 -1     P   -2 -1 -2 -1 -4 -2 -2 -1 -1
A  -1  0  2  1  2 -6  1 -1  0  0     A   -2 -1  4  0  4 -3  0 -2 -1 -1
W  -3 -7 -6 -7 -6 17 -7 -3 -7 -7     W   -2 -3 -3 -2 -3 11 -2 -2 -3 -3
H   6  1 -1 -2 -1 -3 -2  6  1  1     H    8  0 -2 -2 -2 -2 -2  8  0  0
E   1  4  0  0  0 -7  0  1  4  4     E    0  5 -1 -2 -1 -3 -2  0  5  5
A  -1  0  2  1  2 -6  1 -1  0  0     A   -2 -1  4  0  4 -3  0 -2 -1 -1
E   1  4  0  0  0 -7  0  1  4  4     E    0  5 -1 -2 -1 -3 -2  0  5  5
```

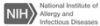NIH National Institute of Allergy and Infectious Diseases

NIAID

30

## SEQUENCE ALIGNMENT

What about gaps?

- Score penalty for opening
- Score penalty for extending

Penalties are log probabilities of a gap of a specific length

NIH National Institute of Allergy and Infectious Diseases

NIAID

31

## SEQUENCE ALIGNMENT

Standard gap costs

| Substitution Matrix | Gap Costs (Open, Extend) |
|---|---|
| PAM30 | (9,1) |
| PAM70 | (10,1) |
| BLOSUM80 | (10,1) |
| BLOSUM62 | (10,1) |
| BLOSUM45 | (15,2) |

NIH National Institute of Allergy and Infectious Diseases

NIAID

32

---

## SEQUENCE ALIGNMENT

Dynamic Programming:
Calculate a matrix of alignment scores

BLOSUM62

```
     H   E   A
P  -2  -1  -1
A  -2  -1   4
W  -2  -3  -3
```

```
            H     E     A
     0    -8   -16   -24
P   -8    -2    -9   -17
A  -16   -10    -3    -5
W  -24   -18   -11    -6
```

NIH National Institute of Allergy and Infectious Diseases

NIAID

33

---

## SEQUENCE ALIGNMENT

Dynamic Programming

1) Calculate a full matrix
2) Traceback to get the Global Alignment

```
         H    E    A    G    A    W    G    H    E    E
              -24  -32  -40  -48  -56  -64  -72  -80
P   -8   -2   -9        -33  -41  -49  -57  -65  -73
A  -16  -10   -3   -5  -13       -29  -37  -45  -53  -61
W  -24  -18  -11   -6   -7  -15       -18  -26  -34  -41
H  -32  -16  -18  -13   -8   -9  -17       -10  -18  -26
E  -40  -24  -11  -19  -15   -9  -12  -19        -5  -13
A  -48  -32  -19   -7  -15  -11  -12  -12  -20        -6
E  -58  -40  -27  -15   -9  -16  -14  -14  -12  -15
```

```
H E A G A W G H E E
- - P - A W H E A E
```

NIH National Institute of Allergy and Infectious Diseases

NIAID

34

# SEQUENCE ALIGNMENT

## Local Alignment

- Alignment of subsequences
- Good for aligning very divergent sequences

## Score Calculation

- Minimum score is zero
- Traceback begins at the highest score
- Score = 0 ➔ End of subsequence

NIH National Institute of Allergy and Infectious Diseases

35

---

# SEQUENCE ALIGNMENT

## Local Alignment

|   | H | E | A | G | A | W | G | H | E | E |
|---|---|---|---|---|---|---|---|---|---|---|
|   | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| P | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| A | 0 | 0 | 0 | 4 | 0 |   | 0 | 0 | 0 | 0 |
| W | 0 | 0 | 0 | 0 | 0 |   |   | 0 | 0 | 0 |
| H | 0 | 8 | 0 | 0 | 0 | 0 | 7 | 13 |   | 7 | 0 |
| E | 0 | 0 | 13 | 5 | 0 | 0 | 0 | 5 | 13 |   | 12 |
| A | 0 | 0 | 5 | 17 | 9 | 4 | 0 | 0 | 5 | 12 | 17 |
| E | 0 | 0 | 5 | 9 | 15 | 8 | 0 | 0 | 0 | 10 | 17 |

```
        A W G H E
Repeat Match    A W - H E        Overlap Match
H E A G A W G H E e              H E A G A W G H E e
p a w H E A e                            p A W - H E a e
        p A W - H E a e
```

NIH National Institute of Allergy and Infectious Diseases

36

---

# SEQUENCE ALIGNMENT

### Scoring alignments and expect values

**Score** := Value in the dynamic programming matrix where the traceback began.

Expect (**E**) value := Number of matches expected due to chance, with a score greater than **S**, based on a stochastic sequence model.

**P** value := Probability of finding at least one match with score ≥ **S**

$$P = 1-e^{-E(S)}$$

NIH National Institute of Allergy and Infectious Diseases

37

## BLAST
**(Basic Local Alignment Search Tool)**

### How does BLAST work?

- Create a list of query sequence "words"
  - Word lengths: 11 nucleotides, 3 amino acids
- Create a list of neighborhood words
  - Similar to query words and above a score threshold
- Search for matches in the database
- Extend matches
  - Below threshold? Discard!
  - Above threshold? Keep it!
- Format and output maximally extended matches

NIH National Institute of Allergy and Infectious Diseases

NIAID

38

---

## BLAST
**(Basic Local Alignment Search Tool)**

How does BLAST work?

How does BLAST evaluate matches?

It uses (local) alignment scores

NIH National Institute of Allergy and Infectious Diseases

NIAID

39

---

## BLAST

The Many Flavors of BLAST

- BLASTn and BLASTp
- short, nearly-exact match BLAST
- Translated BLAST
  - BLASTx    nt → aa ⇨ protein db
  - tBLASTn   aa ⇨ protein db ← DNA db
  - tBLASTx   nt → aa ⇨ protein db ← DNA db
- PSI-BLAST (Position-Specific Iterated BLAST)
- PHI-BLAST (Pattern Hit Initiated BLAST)
- bl2seq

NIH National Institute of Allergy and Infectious Diseases

NIAID

40

## BLAST

### short, nearly-exact match BLAST

- Increase Expect threshold
- Reduce word size (7 for nt, 2 for aa)
- Turn off low complexity filter
- Protein: Use a more stringent substitution matrix

NIAID

NIH National Institute of Allergy and Infectious Diseases

41

## BLAST

### PSI-BLAST
(Position-Specific Iterated BLAST)

- Perform initial BLASTp search
- Generate a sequence profile from results
- BLASTp using the profile
- Iterate until no new sequences are found
- Convergence

NIAID

NIH National Institute of Allergy and Infectious Diseases

42

## BLAST

### PHI-BLAST
(Pattern Hit Initiated BLAST)

Sequence Profile

[LIVMF]-G-E-x-[GAS]-[LIVM]-x(5,11)-R-[STAQ]-A-x-[LIVMA]-x-[STACV]

[ ] = Any of the residues within the brackets
- = spacer separating sites in the profile
x = Any residue
x(a,b) = Any residues a to b in length

    VGERGLEEDKRKRSAWMQC

    MGETALRRRKKEDEERTANVYT

    FGEAAMPGGPHQSRSAFAWV

NIAID

NIH National Institute of Allergy and Infectious Diseases

43

## BLAST

### Access to BLAST

- NCBI
- Your own computer
- NIAID HPC cluster

NIAID

National Institute of
Allergy and
Infectious Diseases

44

## Multiple Sequence Alignment

### Multiple Sequence Alignment
### The Progressive Alignment Algorithm



From RC Edgar 2004, NAR 32: 1792-1797

NIAID

National Institute of
Allergy and
Infectious Diseases

45

## Multiple Sequence Alignment

### Programs

- Clustal
  - Your own computer
  - Web Server
  - NIAID HPC cluster
- MUSCLE
  - Your own computer
  - Web Server
  - NIAID HPC cluster
- MAFFT
  - Web Server

NIAID

National Institute of
Allergy and
Infectious Diseases

46

## Multiple Sequence Alignment

**NEVER**
directly input the output of a MSA program into an analysis program!

**ALWAYS**
inspect the alignment to improve it.

NIH National Institute of Allergy and Infectious Diseases

NIAID

47

## Multiple Sequence Alignment

Multiple Sequence Alignment Editors

- MacVector
  - Commercial software
- MegAlign (Lasergene)
  - Commercial software
- AliView
  - Public domain
- GeneDoc
  - Public domain
- BioEdit
  - Public domain

NIH National Institute of Allergy and Infectious Diseases

NIAID

48

## Web Resources

**ClustalW**
http://www.clustal.org/

**Muscle**
http://www.drive5.com/muscle/download3.6.html

**MAFFT**
http://mafft.cbrc.jp/alignment/server/

**AliView**
http://www.ormbunkar.se/aliview/

**GeneDoc**
http://www.nrbsc.org/downloads/

**BioEdit**
http://www.mbio.ncsu.edu/BioEdit/bioedit.html

NIH National Institute of Allergy and Infectious Diseases

NIAID

49

## Recapitulation

- BLAST search for contig0001 homologs

- Download selected sequence records

- Align sequence records with Clustal2

NIAID

NIH National Institute of Allergy and Infectious Diseases

50

## Seminar Follow-Up Site

- For access to past recordings, handouts, slides visit this site from the NIH network: http://collab.niaid.nih.gov/sites/research/SIG/Bioinformatics/



Recommended Browsers:
- IE for Windows,
- Safari for Mac (Firefox on a Mac is incompatible with NIH Authentication technology)

Login
- If prompted to log in use "NIH\" in front of your username

NIAID

NIH National Institute of Allergy and Infectious Diseases

51

## Retrieving Slides/Handouts



NIAID

NIH National Institute of Allergy and Infectious Diseases

52

## Retrieving Slides/Handouts

**BCBB Seminar Follow-Up Site**

This lecture

These slides

53

## Questions?

Consultation & Advice | Software Development | Biocomputing Resources

ScienceApps@niaid.nih.gov

Bioinformatics and Computational Biosciences Branch
NIAID Office of Cyber Infrastructure and Computational Biology

54

## Next Lecture

The Count

Grover

Oscar    Elmo    Ernie

Bert

Cookie

Herry

55