# Using BLAST for Genomic Sequence Annotation

**Jeremy Buhler**

**Adapted by Wilson Leung**
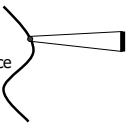
Last Update: 08/16/2023

1

## Overview

- What is comparative annotation?

- How to measure similarity between biosequences

- How to decide whether two sequences are "similar enough"

2

## Comparative Annotation

- Identify functional elements in DNA sequence

- Uses comparison to databases of sequences with known function

New DNA sequence

Probably *CFTR*

3

## Why Does It Work?

- Functional sequences are under negative selection → fewer mutations

- More conservation → greater similarity

- **BLAST software recognizes similarity.**

4

## Caveats w/Similarity Evidence

- Similarity without conservation
  - random chance

- Conservation without selective pressure
  - slow mutation
  - recent divergence

- Similar selective pressures, but seqs have two distinct functions

5

## Overview

- What is comparative annotation?

- How to measure similarity between biosequences

- How to decide whether two sequences are "similar enough"

6

## What is **Similarity**?

- How to measure similarity of two DNA seqs?

- Mutations happen…
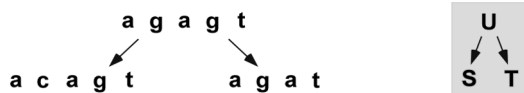
- Measure should reflect desired evolutionary inference

7

## Mutational Model

- Sequences change by series of events of (only) three types:

  - **Substitution** of one base ACG ➡ ATG

  - **Insertion** of one base ACG ➡ ACAG

  - **Deletion** of one base ACG ➡ AG

8

## Sequence History (1/2)

- Suppose seqs S, T diverged from a common ancestral sequence U…

  a g a g t

  a c a g t        a g a t

  U
  S  T

9

## Sequence History (2/2)

- Draw lines between bases of S and T that come from same base of U.

  a g a g t

  a c a g t        a g a ⌴ t

- This is a "tree alignment" of S,T,U.

10

## Sequence Alignment

- Now elide the ancestor…
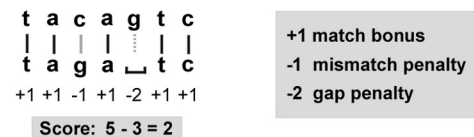
  S     a  c  a  g  t
        |  |  |     |
  T     a  g  a  ⌴  t

- Result is correspondence between bases of S, T – a sequence alignment

11

## Similarity Score of Alignment

- Fewer mutations → more conservation

  t  a  c  a  g  t  c
  |  |  |     |  |  |
  t  a  g  a  ⌴  t  c
  +1 +1 -1 +1 -2 +1 +1

  +1 match bonus
  -1 mismatch penalty
  -2 gap penalty

  Score: 5 - 3 = 2

- Give bonus for matches, penalties for substitutions and gaps

12

## One Small Problem…

- Do you own a time machine?

- If not, how do you know
  - ancestral sequence U?
  - history of mutation?

- Hence, how to get correct alignment?

13

## What We Do In Practice

- Guess an alignment that minimizes # of hypothesized mutations

```
a a g c c – – – – –    S
– – – – – a a t c c    T
```

```
a a g c c    S
| |   | |
a a t c c    T
```

- (more precisely, maximizes score)

14

## Overview

- What is comparative annotation?

- How to measure similarity between biosequences

- How to decide whether two sequences are "similar enough"

15

## Deciding What to Report

- Any two sequences can be aligned with **some** score.

- Higher scores are better…

- When is score high enough to be evidence of conservation?

16

## Idea: Test a Null Hypothesis

- Suppose two DNA seqs S, T are **completely unrelated**.

- What is probability that best alignment between S, T has score at least Θ?

- If score(S,T) is unlikely to occur by chance, then report (S,T)

17

## Null Model Assumptions

- Bases of seqs S, T generated independently at random

random process **S**

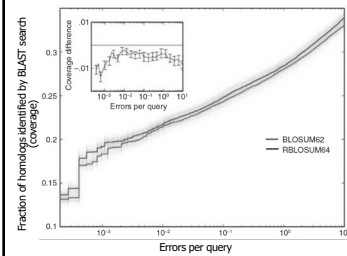random process **T**

18

## Scoring Systems



**BLOSUM62 Matrix**

Amino acids grouped by the Dayhoff classification scheme

"BLOSUM". Wikipedia. Ppgardne.

- Eddy SR. Where did the BLOSUM62 alignment score matrix come from? Nat Biotechnol. 2004 Aug;22(8):1035-6.

19

## Errors in the BLOSUM Matrices



- Benchmark: detect distant homologs in the ASTRAL database

- The "correct" matrix (**RBLOSUM64**) performs worse than the **BLOSUM62** matrix

- Subsequent "fix" of the BLOSUM matrices:
  - Hess M, et al. Addressing inaccuracies in BLOSUM computation improves homology search performance. BMC Bioinformatics. 2016 Apr 27;17:189.

Styczynski MP, et al. BLOSUM62 miscalculations improve search performance. Nat Biotechnol. 2008 Mar;26(3):274-5.

20

## P-values

- Given random seqs S', T' with same base distributions as S, T

- Karlin-Altschul theory tells us probability that S', T' align with score at least $\Theta$

- If $p(\Theta)$ is small, report alignment of S,T

21

## E-values

- For computational reasons, BLAST reports not $p(\Theta)$ but rather $E(\Theta)$

- $E(\Theta)$ = expected # times alignment with score at least $\Theta$ happens by chance in current search

- If $E(\Theta) < 1$, then score $\Theta$ is interesting

22

## Caveats about E-values

- Model from which E-values are computed is too simple for real bioseqs

- Large margin of safety is wise

- Be very skeptical of "matches" with $E > 10^{-5}$

23

## Explanation for E-value = 0.0

- E-value is less than **1.0e-180**



NCBI Home
IEB Home
C++ Toolkit docs
C Toolkit source browser
C Toolkit source browser (2)

**NCBI C++ Toolkit Cross Reference**

**c++/src/algo/blast/api/blast_seqalign.cpp**

```
0056  BEGIN_NCBI_SCOPE
0057  USING_SCOPE(objects);
0058  BEGIN_SCOPE(blast)
0059
0060  #ifndef SMALLEST_EVALUE
0061  /// Threshold below which e-values are saved as 0
0062  #define SMALLEST_EVALUE 1.0e-180
```

24

## Summary

- Comparative annotation with BLAST uses similarity as evidence for conserved function.

- Similarity score based on hypothesized evolutionary relations among sequences.

- E-values indicate whether scores are high enough to be real biological conservation.

25

## Additional Resources

- Introduction to Dynamic Programming
  - Overview of the algorithms for calculating global, semiglobal, and local alignments

- From Smith-Waterman to BLAST
  - Discuss the heuristics used by BLAST to reduce the search space and quickly report high-scoring local alignments

26