

# Basic Local Alignment Search Tool

## Alignments

- used to uncover homologies between sequences
- combined with phylogenetic studies
  - can determine *orthologous* and *paralogous* relationships

## Local Alignment

- uses a subset of a sequence and attempts to align it to subset of other sequences
- computationally less expensive than other methods

## Local Alignment Example

- a small seed is uncovered

*The initial seed for the alignment:*

```
      TAT
      |||
AAGCGAATAATATATTTATACTCAGATTATTGCGCG
```

- And now the extended alignment:

```
      TATATATTAGTA
      ||||| ||
AAGCGAATAATATATTTATACTCAGATTATTGCGCG
```

# BLAST

## Basic Local Alignment Search Tool

Altschul et al. (1990. J Mol Biol 215:403)

- Set of alignment algorithms
- Use the same search protocol to:
  - Find a short fragment of a query sequence
  - That aligns with a fragment of a subject sequence found in a database

### **General Concept for Original BLAST Program**

- Sequence (query) is broken into words of *length W*
- Align all words with sequences in the database
- Calculate *score T* for each word that aligns with a sequence in the database using a substitution matrix
- Discard words whose T value is below a *neighborhood score threshold*
- Extend words in both directions until score falls by *dropoff value X* when compared to previous best score



## BLAST Words

- Three characters in length for proteins
- Compiled by using a sliding window

|-----|

|-----|

|-----|

|-----|

M E N G G P A P E S



## Build Alignment

### 1. Original alignment: T Score = 19

			G	G	P						
1	-1	-2	6	6	7	4	7	5	4	<b>BLOSUM 62 Score</b>	
I	P	A	G	G	P	A	P	E	S		

### 2. Extend one amino acid in each direction: T Score = 21

			N	G	G	P	A			
1	-1	-2	6	6	7	4	7	5	4	<b>BLOSUM 62 Score</b>
I	P	A	G	G	P	A	P	E	S	

### 3. Stop when next extension drops off below value X compared to previous score

## Points to Remember

- The T score is converted into a bits score by a complicated formula
- The X value is based on the bit score

## **BLAST Statistics**

### **Score (bits)**

- A statistical conversion of the score derived by summing using the substitution matrix

### **E value of $-10$ ( $=1 \times 10^{-10}$ )**

- Unlikely that random chance lead to this current alignment compared to an alignment with an e value of 1
- Often considered to be a probability

### **Rules of thumb:**

- **E value of  $-30$  or less**
  - Sequences are homologous
- **E values of  $-5$** 
  - Often considered significant enough when annotating a genome

## BLAST2

(1997. Nucleic Acids Research 25:3389-3402)

Takes a different (and three-times faster) approach than the original BLAST algorithm

- Same word search
- Lower T value
- **Neighboring words discovered**
  - Must be at a distance less than A (default 40)
- Alignment extended from the neighboring words

Gap penalties

- New in BLAST2
- **Allow for better alignments**
- Default for amino acid search
- Introducing a gap
  - -11
- Extending that gap
  - -1

## BLAST Algorithms

Search	Query	Database
blastn	nucleotide	nucleotide
blastx	translated nucleotide in all six frames	protein
tblastx	translated nucleotide in all six frames	translated nucleotide in all six frames
blastp	protein	protein



# Homology, Orthology, Paralogy, Identity, Similarity

How do we define the relationship between two nucleotide or protein sequences?

## Homology

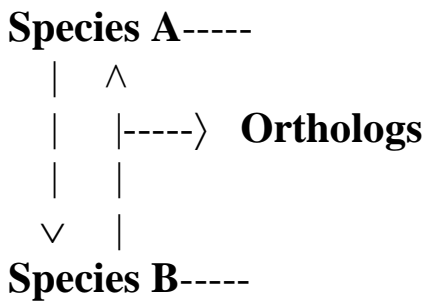
Two sequences are said to show homology if they are identical by descent in a taxonomic lineage or the result of within species duplication

## Homologs

Two sequences that exhibit homology

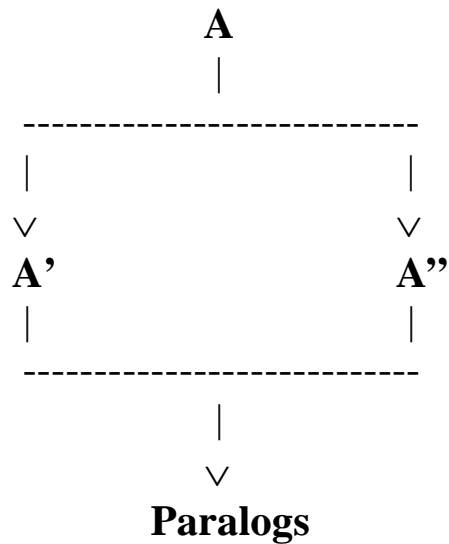
## Orthologs

Two sequences related by descent in a species lineage



## Paralogs

Two sequences related by a duplication event within a species



Describing the relationship between orthologs and paralogs

## Important note

Two sequences are homologous or not homologous; **there is no percentage of homology**

## Identity

The % of exact nucleotide or amino acid matches between two sequences

## Similarity

The % of identical or similar amino acids between two protein sequences

## Examples of similar amino acids

Valine/Leucine

Valine/Isoleucine

Threonine/Serine