# Sequence Alignment Tools

The development of efficient sequence alignment tools represents one of the most transformative milestones in bioinformatics. Modern biology generates vast quantities of genomic and proteomic data, and manually aligning even a handful of sequences using the original dynamic programming algorithms is computationally too expensive. To speed up this tasks, or better, to make them available for large community of molecular biologists, a series of heuristic and optimized tools were created to approximate optimal alignments with remarkable speed and accuracy. Among these, BLAST revolutionized local alignment searches by introducing a word-based heuristic that made database-scale similarity searches routine. Later, tools such as CLUSTAL, MAFFT, and MUSCLE extended these ideas to multiple sequence alignment, enabling researchers to analyze entire gene families and protein domains simultaneously. Together, these tools form the computational backbone of modern bioinformatics, supporting genome annotation, evolutionary inference, and functional analysis.

The acronym BLAST stands for Basic Local Alignment Search Tool, but its creators likely knew what they were doing. Biologists quickly began saying things like "I'll blast this sequence and see what comes up." The verb caught on instantly; few tools in science have managed to become both indispensable and grammatically productive.

Some Typical Questions We Can Answer with Alignment Tools

Sequence alignment tools are not only about matching letters in biological strings-they are instruments for asking and answering biological questions that can be expressed in computational terms. Here are some common questions that can be answered with modern sequence alignment tools:

- Does this new gene resemble any known gene? Given a newly sequenced DNA fragment, BLAST can search large genomic databases to find similar sequences. In computer science terms, this is analogous to finding approximate string matches in a massive text corpus, where the cost of substitutions and gaps is defined by a scoring function.
- Are two proteins likely to have the same function? Similar amino acid sequences often imply structural or functional similarity.

  Alignment tools quantify this similarity, much like comparing two

software functions by their abstract syntax trees or control-flow patterns.

- Which species or genes are most closely related? By aligning multiple sequences and comparing their similarities, we can infer evolutionary relationships. Conceptually, this is like clustering documents or source code repositories based on shared features or edit distances (more on this in our next chapter).
- Has a mutation occurred, and where? Comparing two versions of the same gene is similar to computing a diff between two files—alignments highlight insertions, deletions, and substitutions corresponding to biological mutations.
- Which regions of a genome are conserved across species? Multiple sequence alignment reveals conserved segments that have remained stable through evolution, similar to detecting invariant code segments across different implementations of the same algorithm.

By framing biological questions in terms of sequence similarity and alignment, these tools enable researchers to extract meaningful insights from raw genomic data, bridging the gap between biology and computation.

#### **BLAST**

BLAST (Basic Local Alignment Search Tool) identifies regions of similarity between sequences by comparing a query sequence to a database of known sequences. Unlike global alignment algorithms, which attempt to align entire sequences, BLAST uses local alignment, focusing on finding high-similarity segments (subsequences) within larger sequences. This approach is computationally efficient and wellsuited for detecting functional or evolutionary relationships across genomic and protein data.

The BLAST algorithm works in three main phases: seeding, extension, and evaluation.

1. **Seeding**: The algorithm first breaks down the query sequence into shorter words or "k-mers" of a fixed length, usually three residues for proteins, eleven for nucleotides. These k-mers are compared against all possible k-mers in the database. If the similarity between the query k-mer and a database k-mer meets a certain threshold score, it becomes a "seed" that may lead to a high-scoring alignment.

BLAST was crafted by Altschul and colleagues and first reported in the Journal of Molecular Biology (1990). The tool revolutionized bioinformatics by making large-scale sequence similarity searches both fast and practical.

Each entry in these databases is annotated with metadata—organism name, gene or protein identifier, functional description, literature references, and cross-links to resources like GenBank, UniProt, or Gene Ontology—so that when BLAST finds a similar sequence, the result can be interpreted in a biological context.

- 2. Extension: For each seed match, BLAST extends the alignment in both directions, calculating a score for each extension by adding or subtracting points for matches, mismatches, and gaps based on a scoring matrix (e.g., BLOSUM for proteins, which gives positive scores to biologically likely substitutions). When multiple highscoring segment pairs occur between the same query and database entry, BLAST combines them into a single hit under that accession, reporting the best overall score and the combined alignment coverage. This ensures the output highlights biologically meaningful relationships rather than listing every local match separately. This process continues until the alignment score drops below a threshold, meaning no further extension would improve the alignment. The resulting segments, known as high-scoring segment pairs, represent local regions of high similarity between the query and database sequence.
- 3. Evaluation and Ranking (E-value calculation): The significance of each HSP is assessed using a statistical measure called the E-value, or "expect" value. The E-value estimates the number of alignments with a score equal to or greater than the observed score that could be expected by chance when searching a database of a given size. It's calculated using the formula:

$$E = K \times m \times n \times e^{-\lambda S}$$

where:

- K and  $\lambda$  are statistical parameters, essentially a normalization constants. They dependent on the scoring system and sequence composition, Together, they describe the extreme value distribution of alignment scores expected by chance.
- *m* and *n* are the lengths of the query and database sequences, respectively.
- *S* is the raw alignment score of the high-scoring segment pairs, that is, the score of the best contiguous local alignment produced by extending a single seed in both directions.

Lower E-values indicate more significant alignments. For example, an E-value close to zero suggests the alignment is unlikely to be random, hinting at a meaningful biological relationship.

Intuitively, the E-value indicates how many matches of this quality would be expected to occur by chance in a database of the same size. For example, an E-value of  $10^{-5}$  means that about one in 100,000

random searches would produce an alignment this good. Small Evalues therefore suggest that the match is unlikely to be random and is likely to be biologically meaningful.

By efficiently identifying and ranking local similarities, BLAST allows researchers to filter out biologically relevant hits from random alignments. Its computational efficiency and effectiveness have made it a cornerstone in bioinformatics, used for everything from annotating gene sequences to exploring evolutionary relationships.

## Accessing BLAST and BLAST Variants

BLAST is accessible through various bioinformatics websites, like that of the National Center for Biotechnology Information (NCBI)<sup>4</sup>, and accessed on site, through application program interface, or run on a local computer. There are several BLAST variants tailored for different types of sequence comparisons:

- blastn: Compares a nucleotide query sequence against a nucleotide database. It is commonly used for identifying similar DNA or RNA sequences, locating genes within genomes, or verifying cloned sequences.
- blastp: Compares an amino acid query sequence against a protein database. This is typically used to infer possible protein function or homology when the corresponding gene sequence is known.
- blastx: Translates a nucleotide query sequence in all six reading frames and compares the resulting protein sequences against a protein database. This is particularly useful when analyzing DNA sequences that may encode proteins but where the correct reading frame is unknown.
- tblastn: Uses a protein query to search a nucleotide database translated in all six reading frames. This allows detection of potential protein-coding regions in genomic DNA.
- tblastx: Translates both the nucleotide query and the nucleotide database sequences in all six reading frames and compares the resulting proteins. It is the most computationally intensive variant but can detect distant relationships that may be missed at the nucleotide level.

Each variant uses the same core BLAST algorithm but is optimized for different biological contexts, depending on whether the input and database sequences are nucleotides or proteins.

4 https://blast.ncbi.nlm.nih.gov/ Blast.cgi

#### BLAST Example: Identifying an Unknown Gene Fragment

Suppose a molecular biologist sequences a short fragment of DNA from the soil amoeba Dictyostelium discoideum, but the gene's identity is unknown. The fragment is:

> 1 atggaaacaa ttcaatcagt tattacagaa tggagtgatt caaaaagttg ggatcattta 61 tttcaacata atttcaaaga ttcaaactgg tcagaattat ttgacccagt aaatttcaaa 121 ttcaaatttg gtacaacacc attttctcaa ttccaaattc ttccatcagt tatttcctta

This sequence looks suspiciously regular—perhaps it encodes a conserved protein motif. We can use blastx to translate the nucleotide sequence in all six reading frames and search against the NCBI protein database.

The BLAST output shows a top hit with, perhaps (try it on your own to see what actually happens), the following summary:

Top hit: cAMP-dependent protein kinase catalytic subunit (PKA-C) **Organism:** Dictyostelium discoideum **E-value:**  $5 \times 10^{-87}$  **Percent identity:** 98% Accession: XM\_640288.1

The alignment reveals that the query fragment encodes a portion of the catalytic domain of the pkac gene, a well-studied protein kinase involved in cell signaling and developmental regulation.

#### BLAST Example: Finding Protein Homologs Across Mammals

Suppose we have obtained a short amino acid sequence fragment from a human protein sample and wish to identify it and find homologous proteins in other mammals. The fragment is:

#### **MVHLTPEEKSAVTALWGKVNV**

Running this fragment through blastp (protein-protein search) at the NCBI BLAST server may yield the following top hit:

Top hit: Hemoglobin subunit beta (Homo sapiens) Accession: NP\_-000509.1 E-value:  $2 \times 10^{-29}$  Percent identity: 100% over 21 amino acids

The BLAST report also lists numerous homologs with similarly high scores:

Species	Protein	Identity
Pan troglodytes (chimpanzee)	Hemoglobin subunit beta	100%
Mus musculus (mouse)	Hemoglobin beta chain	90%
Canis lupus familiaris (dog)	Hemoglobin beta chain	86%
Bos taurus (cow)	Hemoglobin beta chain	85%

These results show that the query sequence corresponds to the N-terminal region of human hemoglobin's beta subunit and that highly similar homologs are present in all tested mammals. The high conservation reflects the protein's essential role in oxygen transport.

# Multiple Sequence Alignment

A popular technique and a software tool that implements a heuristic approach to multiple sequence alignment is CLUSTAL. The term CLUSTAL derived from "cluster alignment," reflecting the software's purpose of clustering and aligning multiple sequences. The name emphasizes its function in performing multiple sequence alignment by progressively clustering sequences based on their similarity. CLUSTAL works through a series of steps, typically following a progressive alignment approach:

- Pairwise Alignment: First, CLUSTAL calculates pairwise alignment scores for all pairs of sequences using methods such as Needleman-Wunsch (for global alignment) or Smith-Waterman (for local alignment). These scores quantify how similar each sequence pair is.
- 2. Guide Tree Construction: CLUSTAL then uses these pairwise alignment scores to build a guide tree, often by a method such as the neighbor-joining algorithm (more on this in the next chapter). This tree reflects the evolutionary relationships between sequences based on their similarity scores, determining the order in which sequences will be aligned.
- 3. **Progressive Alignment**: Using the guide tree, CLUSTAL progressively aligns sequences from most to least similar, starting with closely related pairs and gradually aligning less similar ones. At each step, previously aligned groups are treated as single units or "profiles" to ensure consistency across the alignment.

An important advancement came with **CLUSTAL W**. This version significantly improved alignment sensitivity by introducing sequence weighting, position-specific gap penalties, and flexible scoring matrices. These innovations made CLUSTAL W one of the most accurate and widely adopted tools for multiple sequence alignment and established the foundation for many modern alignment programs.

#### MAFFT and MUSCLE

While CLUSTAL established the foundation for multiple sequence alignment, more recent tools such as MAFFT and MUSCLE have

CLUSTAL was introduced by Higgins and Sharp in *Computer Applications in the Biosciences* (1988), the journal that is now know as *Bioinformatics* and became one of the most widely used tools for multiple sequence alignment, setting the standard for comparative sequence analysis.

The CLUSTAL W algorithm (*Nucleic Acids Research*, 1994) greatly improved multiple sequence alignment through sequence weighting, position-specific gap penalties, and refined substitution matrices, setting a new benchmark for accuracy and robustness.

greatly improved speed, scalability, and alignment quality.

MAFFT (Multiple Alignment using Fast Fourier Transform) accelerates the computation of sequence similarity by representing sequences as numerical profiles rather than strings of symbols. Each residue (nucleotide or amino acid) is replaced by a vector encoding its biochemical properties—for example, hydrophobicity or charge. These numerical profiles are then transformed into frequency space using the Fast Fourier Transform (FFT), which efficiently computes correlations between sequences. This allows MAFFT to identify regions of local similarity by comparing frequency patterns of encoded properties instead of performing direct symbol-by-symbol comparisons, greatly reducing computation time for large datasets.

**MUSCLE** (Multiple Sequence Comparison by Log-Expectation) focuses on progressively refining the alignment through three main stages:

- 1. Draft alignment: MUSCLE quickly estimates pairwise distances using short word matches (k-mers) instead of full alignments and constructs an initial guide tree that determines the order in which sequences will be aligned.
- 2. Improved guide tree: the algorithm recomputes pairwise distances using the more accurate alignment-based similarities, reconstructs a new guide tree, and realigns the sequences accordingly.
- 3. Iterative refinement: MUSCLE divides the alignment into subgroups, realigns them, and keeps changes only if they increase the overall alignment score, measured by a log-expectation function that estimates how likely aligned residues are to be related.

Through this cycle of alignment, reevaluation, and refinement, MUS-CLE efficiently converges toward a high-quality multiple sequence alignment without requiring exhaustive computation.

### What's Next?

From alignments, we move toward the broader picture of evolution. Once we can align sequences and quantify their similarity, we can begin to ask deeper questions: Which genes share common ancestry? How closely related are different species? And how can we represent these relationships visually? In the next chapter, we will build on the foundations of sequence alignment to group similar sequences through clustering, reconstruct evolutionary trees using algorithms such as neighbor joining, and visualize these results to reveal patterns of divergence and common origin. In doing so, we will transMAFFT was introduced by Katoh, Misawa, Kuma, and Miyata in Nucleic Acids Research (2002), and MUSCLE by Robert C. Edgar in Nucleic Acids Research (2004). Both became standard tools for large-scale multiple sequence alignment.

MAFFT introduced the use of FFT for biological sequence alignment, transforming residue properties into frequency space to accelerate similarity detection.

form raw sequence data into a map of life's history—an evolutionary network connecting genes, proteins, and species.