

## **GLOSSARY OF TERMS AND ABBREVIATIONS**

---

**BLAST** Basic Local Alignment Search Tool—a tool for identifying sequences in a database that match a given query sequence. Statistical analysis is applied to judge the significance of each match. Matching sequences may be homologous to, or related to, the query sequence. There are several versions of BLAST:

- BLASTP** compares an amino acid query sequence against a protein sequence database
- BLASTN** compares a nucleotide query sequence against a nucleotide sequence database
- BLASTX** compares a nucleotide query sequence translated in all reading frames against a protein sequence database
- TBLASTN** compares a protein query sequence against a nucleotide sequence database dynamically translated in all reading frames
- TBLASTX** compares the six-frame translations of a nucleotide query sequence against the six-frame translations of a nucleotide sequence database.

**BLAT** BLAST-Like Alignment Tool. BLAT might superficially appear to be like BLAST, also being a tool for detecting subsequences that match a given query sequence, however BLAT and BLAST have a number of differences. BLAT was developed at the UCSC; it searches the human genome by keeping an index of the entire genome in memory. The index consists of all non-overlapping 11-mers except for repeat sequences. A BLAT search of the human genome will quickly find sequences of 95% and greater similarity of length 40 bases or more. It may miss more divergent or shorter sequence alignments (see the UCSC FAQ for more details on this tool — <http://genome.ucsc.edu/FAQ.html>).

**CDS** Coding sequence.

**Contig Map** A map depicting the relative order of overlapping (contiguous) clones representing a complete genomic or chromosomal segment.

**DAS (Distributed Annotation System)** DAS is a protocol for browsing and sharing genome sequence annotations across the Internet, allowing users to search and compare annotations from several sources. Ensembl provides a DAS reference server giving access to a wide range of specialist annotations of the human genome (see <http://www.ensembl.org/das/> for more detail).

**Data Mining** The ability to query very large databases in order to satisfy a hypothesis (“top-down” data mining); or to interrogate a database in order to generate new hypotheses based on rigorous statistical correlations (“bottom-up” data mining).

**Domain (protein)** A region of special biological interest within a single protein sequence. However, a domain may also be defined as a region within the three-dimensional structure of a protein that may encompass regions of several distinct protein sequences that accomplishes a specific function. A domain class is a group of domains that share a common set of well-defined properties or characteristics.

**Electronic PCR (ePCR)** An electronic process analogous to lab based PCR. Two primers are used to map a sequence feature (e.g. a SNP). To validate the position both primers must map in the same vicinity spanning a defined distance, effectively producing an electronic PCR product.

**Expressed Sequence Tag (EST)** A short sequence read from an expressed gene derived from a cDNA library. Databases storing large numbers of ESTs can be used to gauge the relative abundance of different transcripts in cDNA libraries and the tissues from which they are derived. An EST can also act as a physical tag for the identification, cloning and full length sequencing of the corresponding cDNA or gene.

**FASTA format** FASTA format, originally devised for Lipman & Pearson's FASTA (Fast-All) sequence alignment algorithm, is one of the simplest and most widely accepted formats for sequences, taking the form of a simple header preceded by a ">" sign and sequence on the following line, e.g.

```
>sequence_id
gataggctgagcgcgatgctagctagctagc
```

**Golden Path** The golden path is a term applied to the first and subsequent assemblies of the human genome.

**Hidden Markov model (HMM)** A joint statistical model for an ordered sequence of variables. The result of stochastically perturbing the variables in a Markov chain (the original variables are thus "hidden"), where the Markov chain has discrete variables which select the "state" of the HMM at each step. The perturbed values can be continuous and are the "outputs" of the HMM. A Hidden Markov Model is equivalently a coupled mixture model where the joint distribution over states is a Markov chain. Hidden Markov models are valuable in bioinformatics because they allow a search or alignment algorithm to be trained using unaligned or unweighted input sequences; and because they allow position-dependent scoring parameters such as gap penalties, thus more accurately modelling the consequences of evolutionary events on sequence families.

**Homology** (strict) Two or more biological species, systems or molecules that share a common evolutionary ancestor. (general) Two or more gene or protein sequences that share a significant degree of similarity, typically measured by the amount of identity (in the case of DNA), or conservative replacements (in the case of protein), that they register along their lengths. Sequence "homology" searches are typically performed with a query DNA or protein sequence to identify known genes or gene products that share significant similarity and hence might inform on the ancestry, heritage and possible function of the query gene.

**in silico (biology)** (Lit. computer mediated). The use of computers to simulate, process, or analyse a biological experiment.

**NCBI** National Center for Biotechnology Information, Washington, D.C., USA.

**Open reading frame (ORF)** Any stretch of DNA that potentially encodes a protein. Open reading frames start with a start codon, and end with a termination codon. No termination codons may be present internally. The identification of an ORF is the first indication that a segment of DNA may be part of a functional gene.

- Ortholog/Paralog** Paralogs are genes related by duplication within a genome. Orthologs retain the same function in the course of evolution, whereas paralogs evolve new functions, even if these are related to the original one.
- PERL** PERL is the short form acronym for Practical Extraction and Report Language. Perl is relatively straightforward up to a certain level—this has encouraged its development as the primary language of biological computing.
- Relational Database** A database that follows E. F. Codd's 11 rules, a series of mathematical and logical steps for the organization and systemization of data into a software system that allows easy retrieval, updating, and expansion. A relational database management system (RDBMS) stores data in a database consisting of one or more tables of rows and columns. The rows correspond to a record (tuple); the columns correspond to attributes (fields) in the record. RDBMSs use Structured Query Language (SQL) for data definition, data management, and data access and retrieval. Relational and object-relational databases are used extensively in bioinformatics to store sequence and other biological data.
- Secondary structure (protein)** The organization of the peptide backbone of a protein that occurs as a result of hydrogen bonds e.g. alpha helix, Beta pleated sheet.
- Sequence Tagged Site (STS)** A unique sequence from a known chromosomal location that can be amplified by PCR. STSs act as physical markers for genomic mapping and cloning.
- Single Nucleotide polymorphism (SNP)** A DNA sequence variation resulting from substitution of one nucleotide for another.
- SQL** Structured Query Language. A type of programming language used to construct database queries and perform updates and other maintenance of relational databases, SQL is not a fully-fledged language that can create standalone applications, but it is powerful enough to create interactive routines in other database programs.
- Substitution matrix** A model of protein evolution at the sequence level resulting in the development of a set of widely used substitution matrices. These are frequently called Dayhoff, MDM (Mutation Data Matrix), BLOSUM or PAM (Percent Accepted Mutation) matrices. They are derived from global alignments of closely related sequences. Matrices for greater evolutionary distances are extrapolated from those for lesser ones.
- Tertiary structure (protein)** Folding of a protein chain via interactions of its side-chain molecules including formation of disulphide bonds between cysteine residues.
- UCSC** University of California Santa Cruz
- UTR** Untranslated region. The non coding region of an mRNA transcript flanking either side of the open reading frame.