# GLOSSARY

An extensive glossary of genetic terms can be found at the Web site for the National Human Genome Research Institute (*http://www.nhgri.nih.gov/DIR/VIP/Glossary/*). The entries in this glossary provide a brief written definition of the term; the user can also listen to an informative explanation of the term using RealAudio.

**algorithm**   Any sequence of actions (e.g., computational steps) that perform a particular task.

**analogous**   In phylogenetics, characters that have descended in a convergent fashion from unrelated ancestors.

**browser**   Program used to access sites on the World Wide Web. By using hypertext markup language (HTML), browsers are capable of representing a Web page the same way regardless of computer platform.

**candidate gene**   A gene that is implicated in the causation of a gene. Candidate genes lie in a region that has been identified through genetic mapping. The protein product of a candidate gene may implicate the candidate gene as being the actual disease gene being sought.

**cDNA library**   A collection of double-stranded DNA sequences that are generated by copying mRNA molecules. Because these sequences are derived from mRNAs, they contain only protein-coding DNA.

**characters and character states**   In phylogenetics, characters are homologous features in different organisms. The exact condition of that feature in a particular individual is the character state. As an example, the character ''hair color'' can have the character states ''gold,'' ''red,'' and ''yellow.'' In molecular biology, the character states can be one of the four nucleotides (A, C, T, or G) or one of the 20 amino acids. Please note that some authors define character to mean the character state as defined here.

**client** A computer, or the software running on a computer, that interacts with another computer at a remote site (server). Note the difference between client and user.

**contig** Short for ''contiguous.'' Refers to a contiguous set of overlapping DNA sequences.

**cytogenetic map** The representation of a chromosome on staining and examination by microscopy. Visually distinct light and dark bands give each chromosome a unique morphological appearance and allow for the visual tracking of cytogenetic abnormalities such as deletions or inversions.

**descriptor** Information about a sequence or set of sequences whose scope depends on its placement in a record. It is placed on a set of sequences to reduce the need to save multiple redundant copies of information.

**domain name** Refers to one of the levels of organization of the Internet and used to both classify and identify host machines. Top-level domain names usually indicate the type of site or the country in which the host is located.

**download** The act of transferring a file from a remote host to a local machine via FTP.

**E-mail** Electronic mail. Refers to messages that can be composed on the computer and transmitted via the Internet to a remote location within seconds. [*Ant*: snail mail, postal mail.]

**EST** Expressed Sequence Tags. These are usually short (300–500 bp) single reads from mRNA (cDNA) that are usually produced in large numbers. They represent a snapshot of what is expressed in a given tissue or at a given developmental stage. They represent tags (some coding, others not) of expression for a given cDNA library.

**exon** Within a gene, a region that codes for part of the gene's protein product; the ''**ex**pressed regi**on**'' of a gene.

**FAQ** Frequently asked questions. Exactly what it sounds like: a compiled list of questions and answers intended for new users of any computer-based resource, such as mailing lists or newsgroups.

**feature** Annotation on a specific location on a given sequence.

**firewall** A computer separating a company or organization's internal network from the public part, if any, of the same network. Intended to prevent unauthorized access to private computer systems.

**FTP** File transfer protocol. The method by which files are transferred between hosts.

**genetic map** Gives the relative positions of known genes and or markers. Markers must have two or more alleles that can be readily distinguished.

**genome** All of the DNA found within each of the cells of an organism. Eukaryotic genomes can be subdivided into their nuclear genome (chromosomes found within the nucleus) and their mitochondrial genome.

**Gopher** A document delivery system allowing the retrieval and display of text-based files.

**GSS**   Genome Survey Sequences. This DDBJ/EMBL/GenBank division is similar in nature to the EST division, except that its sequences are genomic in origin, rather than being cDNA (mRNA). The GSS division contains (but will not be limited to) the following types of data: random ''single pass read'' genome survey sequences, single pass reads from cosmid/BAC/YAC ends (these could be chromosome specific, but need not be), exon-trapped genomic sequences, and Alu PCR sequences.

**GUI**   Graphical user interface. Refers to software front ends that rely on pictures and icons to direct the interaction of users with the application.

**heuristic algorithm**   An economical strategy for deriving a solution to a problem for which an exact solution is computationally impractical or intractible. Consequently, a heuristic approach is not guaranteed to find the optimal or ''true'' solution.

**homologs/homologous**   In phylogenetics, particular features in different individuals that are genetically descended from the same feature in a common ancestor are termed homologous. In molecular biology, homologous is often used simply to mean similar, regardless of genetic relationship.

**homoplasy**   Similarity that has evolved independently and is not indicative of common phylogenetic origin.

**host**   Any computer on the Internet that can be addressed directly through a unique IP address.

**HTG/HTGS**   High Throughput Genome Sequences. Various genome sequencing centers worldwide are performing large-scale, systematic sequencing of human and other genomes of interest. The databases have deemed it beneficial to put the unfinished sequences resulting from such sequencing efforts in a separate division. HTG sequence entries undergo a maturation process. In Phase 0, the entry contains a single-to-few pass read of a single clone. In Phase 1, the entry contains unfinished sequence, which may be unordered, contain unoriented contigs, or a large number of gaps. In Phase 2, the entry still contains unfinished sequence but is ordered, with oriented contigs that may or may not contain gaps. In Phase 3, the entry contains finished sequence, with no gaps; at this point, the entry is moved into the appropriate primary GenBank division.

**HTML**   Hypertext markup language. The standard, text-based language used to specify the format of World Wide Web documents. HTML files are translated and rendered through the use of Web browsers.

**haplotype**   Sets of alleles that are usually inherited together.

**hyperlink**   A graphic or text within a World Wide Web document that can be selected using a mouse. Clicking on a hyperlink transports the user to another part of the same Web page or to another Web page, regardless of location.

**hypertext**   Within a Web page, text that is differentiated either by color or by underlining, which functions as a hyperlink.

**indel**   Acronym for ''**in**sertion or **del**etion.'' Applied to length-variable regions of a multiple alignment when it is not specified whether sequence length differences have been created by insertions or deletions.

**Internet**   A system of linked computer networks used for the transmission of files and messages between hosts.

**IP address**   The unique, numeric address of a computer host on the Internet.

**Intranet**   A computer network internal to a company or organization. Intranets are often not connected to the Internet or are protected by a firewall.

**Java**   A programming language developed by Sun Microsystems that allows small programs (applets) to be run on any computer. Java applets are typically invoked when a user clicks on a hyperlink.

**LAN**   Local Area Network. A network that connects computers in a small, defined area, such as the offices in a single wing or a group of buildings.

**LOD score**   For "**log od**ds," a statistical estimate of the linkage between two loci on the same chromosome.

**molecular clock**   The hypothesis that nucleotide or amino acid substitutions occur at more or less a fixed rate over evolutionary time, like the slow ticking of a clock. It has been proposed that given a calibration date and a constant molecular clock, the amount of sequence divergence can be used to calculate the time that has elapsed since two molecules diverged.

**mutation studies**   In Sequin, a set of sequences for the same gene in the same species, perhaps the same individual, in which several different induced mutations are isolated and sequenced.

**oligo**   For oligonucleotide, a short, single-stranded DNA or RNA. Most often used as probes for the detection of complementary DNA or RNA.

**orthologs/orthologous**   Homologous sequences are said to be orthologous when they are direct descendants of a sequence in the common ancestor, i.e., without having undergone a gene duplication event. See also homologs and paralogs.

**PAM matrix**   PAM (percent accepted mutation) and BLOSUM (blocks substitution matrix) are matrices that define scores for each of the 210 possible amino acid substitutions. The scores are based on empirical substitution frequencies observed in alignments of database sequences and in general reflect similar physicochemical properties (e.g., a substitiution of leucine for isoleucine, two amino acids of similar hydrophobicity and size, will score higher than a substitution of leucine for glutamate.)

**paralogs/paralogous**   Homologous sequences in two organisms A and B that are descendants of two different copies of a sequence created by a duplication event in the genome of the common ancestor. See also homologous and orthologs.

**pedigree**   A tree representation of a family (cohort) showing the relationships between members and the pattern of inheritance of a given trait.

**phylogenetic studies**   In Sequin, a set of sequences for the same gene in individuals of different species. The presumption is that the individuals cannot interbreed. Sequin does not allow a single organism name but expects the organism to be encoded in the definition line. It does, however, present a control for setting the proper genetic code.

**physical map**   A genome map showing the exact location of genes and markers. The highest-resolution physical map is the DNA sequence itself.

**platform**   Properly, the operating system running software on a computer, e.g., UNIX or Windows 95. More often used to refer to the type of computer, such as a Macintosh or PC-compatible.

**polymorphism**   A gene or locus that can have one or more alleles or haplotypes.

**population studies**   In Sequin, a set of sequences for the same gene in individuals of a single species. The presumption is that the individuals can interbreed. Sequin allows entry of a single organism name, although it would expect that some distinguishing source information, such as strain, clone, or isolate, be entered for each sequence.

**positional cloning**   Relies on the identification of a gene through pedigree analysis, genetic and physical mapping, and mutation analysis. Does not require extensive knowledge of the biochemistry of the disease to determine the gene responsible for the disease. [*Ant*: functional cloning.]

**protein name**   In a sequence record, the preferred field for a protein feature.

**protein description**   In a sequence record, used if the protein name is not known.

**server**   A computer that processes requests issued from remote locations by client machines.

**site**   An individual column of residues in an amino acid or nucleotide alignment. The residues at a site are presumed to be homologous.

**spam**   Postings to newsgroups or mail broadcast to a large number of E-mail accounts that usually are wholly irrelevant or not of interest to the recipients. Analogous to postal junk mail.

**STS**   Sequenced Tagged Sites. An operationally unique sequence that identifies the combination of primer pairs used in a PCR assay, generating a reagent that maps to a single position within the genome. STS sequences are usually on the order of 200–500 bases in length. dbSTS is a division of GenBank devoted to STS sequences; it is intended to facilitate cross-comparison of STSs with sequences in other divisions for the purpose of correlating map positions of anonymous sequences with known genes.

**Telnet**   An Internet protocol or application that allows users to connect to computers at remote locations and use these computers as if they were physically sitting at that computer.

**URL**   Uniform resource locator. Used within Web browsers, URLs specify both the type of site being accessed (FTP, Gopher, or Web) and the address of the Web site.

**user**   The person using client-server or other types of software.

**World Wide Web**   A document delivery system capable of handling various types of non-text-based media.

# INDEX

**457**