

## Glossary of Terms

Note: Words in **bold** in a definition indicate terms also defined in this Glossary

Term	Definition
<b>3'</b>	"3 prime"; Refers to carbon 3 of the nucleic acid sugar component (either ribose in RNA or deoxyribose in DNA) to which additional <b>nucleotides</b> may be added by polymerase, often used to refer to that end of a single-stranded DNA or RNA molecule where the 3' carbon retains its hydroxyl group (-OH) and no further nucleotides are bonded.
<b>5'</b>	"5 prime"; Refers to carbon 5 of the nucleic acid sugar component (either ribose in RNA or deoxyribose in DNA), to which the triphosphate is attached in a <b>nucleotide</b> triphosphate, often used to refer to that end of a single-stranded DNA or RNA molecule where the 5' carbon's phosphate group(s) is/are unattached to a preceding nucleotide.
<b>alternative splicing</b>	The inclusion or exclusion of certain <b>exons</b> in the <b>splicing</b> reactions that determine the sequences included in the final <b>mRNA</b> product. This mechanism is utilized to generate a series of closely related protein <b>isoforms</b> , which differ by the inclusion or exclusion of the particular protein regions encoded by those exons. Alternative splicing is directed by RNA-binding proteins that may block, or stimulate, utilization of a particular splice site.
<b>amino acid</b>	The basic building block of proteins, a small molecule with a -C-C- core, an amine group (-NH <sub>2</sub> ) at one end and a carboxylic acid group (-COOH) at the other end. The general structure can be represented as NH <sub>2</sub> -CHR-COOH, where R can be any of 20 different functional groups of acidic, basic, or nonpolar character.
<b>annotation</b>	Gene annotation is the process of notating the location, structure, and identity of genes in a genome. As initial attempts may be based on incomplete information, gene annotations are constantly changing as further data becomes available. Gene annotation databases are updated regularly, and different databases may refer to the same gene/protein by different names, reflecting new knowledge and improved understanding of protein function.
<b>base</b>	Although formally incorrect (the nitrogenous base that defines A, C, G, T and U is only part of the whole <b>nucleotide</b> ), this is often used as a synonym for "nucleotide" in referring to the A, C, G, T, and U components of DNA and RNA.
<b>base pair/base pairing</b>	The hydrogen bonding of one of the <b>bases</b> (A, C, G, T, U) with another, as dictated by optimal hydrogen bond formation in DNA (A-T and C-G) or in RNA (A-U and C-G). Two polynucleotide strands, or regions thereof, in which all the <b>nucleotides</b> form such base pairs are said to be complementary. In achieving complementarity, each strand of DNA can serve as a template for synthesis of its partner strand - the secret of DNA replication's extremely high accuracy and thereby of inheritance.

<b>canonical</b>	In agreement with existing principles and standards generated from data and evidence. For example, the canonical <b>splice donor sequence</b> is GT. In rare instances, however, the sequence GC is used instead; since GC is not the "standard", it would be referred to as <b>non-canonical</b> .
<b>cDNA</b>	"complementary DNA"; a double-stranded DNA molecule prepared <i>in vitro</i> ("outside of the body"; i.e., in a test tube) by employing an RNA molecule as a template to synthesize DNA using reverse transcriptase. The RNA component of the resulting RNA-DNA hybrid is enzymatically degraded, and the complementary strand then synthesized by DNA polymerase. The resulting double-stranded DNA can be used for cloning and analysis.
<b>CDS</b>	"coding sequence"; that part of the DNA sequence of a gene that is translated into protein.
<b>coding exon</b>	In a gene, any <b>exon</b> that contains some part of the <b>CDS</b> ; in contrast, an exon that has no part translated into protein is called a " <b>non-coding exon</b> ."
<b>coding strand/ positive strand</b>	In a gene, the DNA strand that has the sequence found in the RNA molecule. Also called the sense, positive, or non-template strand.
<b>codon</b>	The sequence of three nucleotides in DNA or RNA that specifies a particular <b>amino acid</b> .
<b>coordinate</b>	Numerical position within a biological sequence; for example, the first base in a DNA sequence would have the coordinate "1".
<b>downstream</b>	Refers to the genomic region that comes after the <b>feature</b> being examined.
<b>exon</b>	In eukaryotes, a contiguous segment of DNA that corresponds to a portion of the mature (processed) RNA product of that gene. Exons in eukaryotic genomes are often, but not always, separated by <b>introns</b> . Although exons are transcribed with the introns, the latter are spliced out during RNA processing and degraded.
<b>feature</b>	Any region of defined structure/sequence in a genomic fragment of DNA. Inherent features would include genes, pseudogenes, and repetitive elements. A feature may also be predicted by computational algorithms, such as those aimed at identifying protein-coding genes.
<b>intron</b>	Non-coding section of a eukaryotic nucleic acid sequence found between exons. Introns are removed ("spliced out") from the primary transcript/pre-mRNA after transcription and before the molecule is exported to the cytoplasm for translation.
<b>isoforms</b>	Potentially different versions of a protein encoded by a single gene. Isoforms result from alternative splicing of a particular <b>pre-mRNA</b> , and/or the use of a different transcription start site.
<b>mRNA</b>	Mature messenger RNA that has been completely processed and is ready for translation; it has a 7-methylguanosine cap at its <b>5'</b> end, a <b>poly(A) tail</b> at its <b>3'</b> end, and has all its <b>introns</b> spliced out.

<b>non-coding strand/ negative strand</b>	Also called the anti-sense, template, or non-coding strand. This strand of the DNA sequence of a single gene is the complement of the 5' to 3' DNA strand known as the positive, sense, non-template, or <b>coding strand</b> . The term loses meaning for longer DNA sequences with genes on both strands.
<b>nucleotide</b>	The basic building block of DNA (A, C, G, T) and RNA (A, C, G, U). Nucleotides consist of a nitrogenous base, a 5-carbon sugar (either ribose in RNA or deoxyribose in DNA), and phosphate group(s).
<b>ORF</b>	"Open reading frame"; a long stretch of <b>codons</b> in the same reading <b>frame</b> uninterrupted by <b>termination codons</b> ; an ORF may reflect the presence of a gene.
<b>phase</b>	The phase describes the number of <b>bases</b> between the end of the <b>exon</b> (defined by the splice site) and the full <b>codon</b> nearest that splice site. The number of bases between the adjacent full codon and an exon/splice site can be 0, 1 or 2. The phase of an <b>upstream</b> exon will determine which <b>frame</b> is translated in the <b>downstream</b> exon by indicating how many bases after the splice acceptor site are needed to create a full codon of 3 bases.
<b>poly(A) tail</b>	About 250 adenine <b>nucleotides</b> that are post-transcriptionally added by poly(A) polymerase to the <b>3'</b> end of eukaryotic transcripts, following cleavage of the newly synthesized RNA ~20 nucleotides <b>downstream</b> of an AAUAAA polyadenylation signal sequence.
<b>pre-mRNA (primary transcript)</b>	The initial transcript from a protein-coding gene that contains both <b>introns</b> and <b>exons</b> . Pre-mRNA requires the addition of a <b>5'</b> cap and <b>3' poly (A) tail</b> and the removal of introns to produce the final <b>mRNA</b> molecule containing joined exons.
<b>promoter</b>	A segment of DNA to which RNA polymerase binds to initiate <b>transcription</b> of the <b>downstream</b> gene(s).
<b>putative</b>	Something that may be predicted or inferred but that requires more evidence to confirm or refute.
<b>read</b>	A raw DNA sequence.
<b>reading frame/frame</b>	A frame is a single series of adjacent nucleotide triplets in DNA or RNA: one frame would have bases at positions 1, 4, 7, etc. as the first base of sequential codons. There are three possible reading frames in an mRNA strand and six in a double stranded DNA molecule due to the two strands from which transcription is possible. Different computer programs number these frames differently, so care should be taken when comparing designated frames from different programs. One common way is to refer to the three possible left-to-right reading frames as +1, +2, and +3 and the three possible right-to-left reading frames as -1, -2, and -3.
<b>splicing</b>	The process by which <b>introns</b> are removed and <b>exons</b> are joined to produce a mature, functional RNA ( <b>mRNA</b> ) from a <b>primary transcript</b> . Some RNAs are self-splicing, but most require a specific ribonucleoprotein complex to catalyze the reaction.

<b>splice acceptor site</b>	The <b>splicing</b> site at the <b>3'</b> end of an <b>intron</b> , at the boundary between an intron and the <b>exon</b> immediately <b>downstream</b> . The <b>canonical</b> splice acceptor site dinucleotide sequence is AG.
<b>splice donor site</b>	The splicing site at the 5' end of an <b>intron</b> , at the boundary between an intron and the <b>exon</b> immediately <b>upstream</b> . The <b>canonical</b> splice donor site dinucleotide sequence is GT; in rare cases, the <b>non-canonical</b> sequence GC is used instead.
<b>splice junction</b>	Either a <b>splice acceptor site</b> or a <b>splice donor site</b> .
<b>start codon (initiation codon)</b>	The first codon of a <b>CDS</b> . In eukaryotes this is almost always ATG, which codes for methionine (one of the 20 <b>amino acids</b> ).
<b>stop codon (termination codon)</b>	A codon that specifies the termination of protein synthesis; sometimes called a "nonsense codon" since it does not specify an <b>amino acid</b> .
<b>transcription</b>	The process of copying one strand of a DNA double helix by RNA polymerase, creating a complementary strand of RNA called the transcript.
<b>translation</b>	The process by which codons in an mRNA are "read" by the ribosome and tRNAs to direct protein synthesis.
<b>TSS (transcription start site)</b>	The location in DNA, generally <b>upstream</b> of a gene's coding sequence, where RNA polymerase begins <b>transcription</b> .
<b>UTR</b>	"Untranslated region"; a segment of DNA (or RNA) that is transcribed and present in the mature <b>mRNA</b> but is not translated into protein. UTRs may be found at either or both of the <b>5'</b> and <b>3'</b> ends of a gene or transcript.
<b>upstream</b>	Refers to the genomic region prior to the <b>feature</b> being examined.