Concepts in Bioinformatics and Genomics

Solutions Manual

Jamil Momand / Alison McCurdy
Silvia Heubach / Nancy Warter-Perez

Chapter 1

Review of Molecular Biology

EXERCISE 1

A good resource for information on disease-related molecules in the Online Mendelian Inheritance of Man (OMIM). OMIM is linked to the National Center of Biotechnology Information (NCBI) website. To use this website optimally, it is important to become familiar with chromosomes. Chromosomes are segments of DNA bound to protein found in the nucleus of eukaryotic cells. For many organisms, their genomes are distributed as chromosomes. Humans have 23 sets of chromosomes (46 total chromosomes). Each chromosome is divided into two major areas. The short arm (p-arm) is located above the centromere, and the long arm (q-arm) is located below the centromere. Telomeres are located at the very ends of the chromosome (Figure 1-17).

To determine the location of a gene on a chromosome, chromosomes can be stained with special dyes that create a series of bands on the chromosomes. If a gene is located at 8q21.3, that would mean chromosome 8, q-arm, band 21.3. The higher the band number, the farther away from the centromere the gene is located.

a. Give the chromosome location (in terms of banding position of the gene that causes sickle cell anemia.

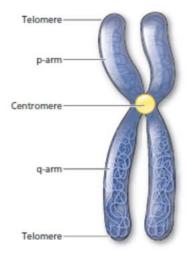


Figure 1-17. Human chromosome showing the centromere, telomeres, p-arm, and q-arm.

- b. Give the name of the gene that causes sickle cell anemia.
- c. Find the nucleotide change and amino acid change that leads to sickle cell anemia.
- d. Explain the change in distribution of spots in the peptide map (Figure 1-16) based on properties of the side chain of the mutant amino acid in the disease.

e. If sickle cell anemia is so devastating, why has it persisted in the population for such a long time? Give a molecular explanation.

Answer 1

- a. chromosomal location 11p15.4
- b. *HBB* (human beta globin)
- c. A to T transversion in the third position of the seventh codon of the hemoglobin beta gene. E6V mutation-glutamate is changed to valine at position 6 in the mature protein. The mature protein sequence is highlighted in yellow below. The glutamate at position 6 is highlighted in pink.
 - 1 MVHLTPEEKS AVTALWGKVN VDEVGGEALG RLLVVYPWTQ RFFESFGDLS TPDAVMGNPK
 - 61 VKAHGKKVLG AFSDGLAHLD NLKGTFATLS ELHCDKLHVD PENFRLLGNV LVCVLAHHFG

121 KEFTPPVQAA YQKVVAGVAN ALAHKYH

- d. The mutant peptide containing the valine is more hydrophobic and less charged. It will stay closer to the neutral point on the paper.
- e. Ongoing research suggests that sickle cells prevent the pathological effects of malaria. Red blood cells with at least one mutant copy of *HBB* are less prone to infection by the mosquito carrying the protozoan that causes malaria. Individuals heterozygous for the *HBB* sickle-cell allele do not have sickle-cell disease, and have increased malaria resistance compared to individuals homozygous for the normal *HBB* allele. The increased survival of heterozygous individuals promotes selection of the *HBB* sickle-cell allele. This is why the *HBB* sickle-cell allele persists. *Plasmodium* is the protozoan that causes malaria. There are several suggestions as to how the sickle human hemoglobin protein (Hb) operates at the molecular level to increase tolerance to *Plasmodium* infection. Here is one:

"Sickle human hemoglobin (Hb) confers a survival advantage to individuals living in endemic areas of malaria, the disease caused by *Plasmodium* infection. As demonstrated hereby, mice expressing sickle Hb do not succumb to experimental cerebral malaria (ECM). This protective effect is exerted irrespectively of parasite load, revealing that sickle Hb confers host tolerance to *Plasmodium* infection. Sickle Hb induces the expression of heme oxygenase-1 (H0-1) in hematopoietic cells, via a mechanism involving the transcription factor NF-E2-related factor 2 (Nrf2). Carbon monoxide (CO), a byproduct of heme catabolism by H0-1, prevents further accumulation of circulating free heme after *Plasmodium* infection, suppressing the pathogenesis of ECM. Moreover, sickle Hb inhibits activation and/or expansion of pathogenic CD8+T cells recognizing antigens expressed by *Plasmodium*, an immunoregulatory effect that does not involve Nrf2 and/or H0-1. Our findings provide insight into molecular mechanisms via which sickle Hb confers host tolerance to severe forms of malaria. (A. Ferreira, I. Marguti, I. Bechmann, et al. "Sickle Hemoglobin Confers Tolerance to Plasmodium Infection." *Cell* 145, no. 3 [April 2011]. 398-409. doi: 10.1016/j.cell.2011.03.049.PubMed PMID: 21529713)

EXERCISE 2

One important exercise a bioinformatician performs is to compare amino acid sequences. One reason to make comparisons is to determine the parts of the proteins that are critical for function. These regions are generally conserved within proteins that perform the same duties. Conserved regions are those that have nearly the same amino acid sequences. Proteins that perform the same duties are called homologs and can be found in different species. For example, p53 from humans and p53 from frogs perform the same functions. There are some regions within these proteins that will be similar in both humans and frogs. We call these regions conserved sequences. A multiple sequence alignment allows the bioinformatician to readily line up amino acid sequences of related proteins. The conserved regions are identified in the alignment. For this exercise, perform a multiple sequence alignment of three homologs of cytochrome C. The three homologs are from human, yeast, and dog, and the accession numbers for these sequences are AAA35732, NP 001183974, and 1YCC. The first two sequences can be found in the protein database at the NCBI. The last sequence can be found in the structure database at the Protein Data Bank. Print out your multiple sequence alignment result and attach a short paragraph explaining how the alignment gives you a clue as to which parts of the cytochrome C protein you would hypothesize are most important to its function. (The function is the same in all three organisms.)

For those of you unfamiliar with NCBI, here are specific instructions:

- Go to the NCBI website.
- Use the dropdown menu to search in the "Protein" database.
- Enter the accession number and click "Search."
- Change the format to FASTA.
- Copy the FASTA output into a sequence alignment window of a website that hosts a sequence alignment program. Make sure the header (the line with the > symbol) is placed on top of the sequence within the window.
- Repeat for each FASTA output of the remaining two proteins.
- Once all three sequences are pasted into the sequence alignment input window, run the program

Answer 2

Note that in the alignment output the asterisks (*) indicate perfect matches across the sequences, colons (:) indicate highly similar amino acids, periods (.) indicate slightly similar amino acids and no symbol indicate dissimilar amino acids.

```
CLUSTAL O(1.2.1) multiple sequence alignment
```

```
yeast NVLWDENNMSEYLTNPKKYIPGTKMAFGGLKKEKDRNDLITYLKKACE-
human GIIWGEDTLMEYLENPKKYIPGTKMIFVGIKKKEERADLIAYLKKATNE
dog GITWGEETLMEYLENPKKYIPGTKMIFAGIKKTGERADLIAYLKKATKE
: * *:.: *** ******** * *:* :* ***:**** :
```

Sequence alignment shows that these three proteins are well conserved, especially sequences $\mbox{QCHTVEKGGPHKVGPNLHGIFGRHSGQAEGYSYTDAN}$ and $\mbox{EYLTNPKKYIPGTKM}$. These may be regions that have important features for activity and structure of the protein.

Chapter 2

Information Organization and Sequence Databases

EXERCISE 1

For this exercise, you will need to access GenBank by going to the NCBI website and using the dropdown menu to search "Nucleotide." Note that the definition of the coding strand is the strand of DNA within the gene that is identical to the transcript (for genetic code, see Figure 1.11). On the other hand, the template strand is the strand that is complementary to the coding strand.

- a. Use the following accession number to access the nucleotide sequence in GenBank: CU329670
- b. Go to the FEATURES section of the record.
- c. Link to the CDS to gain access to the first 5662 nucleotides of the sequence.
- d. Name the protein product of the CDS.
- e. Write the first four amino acids (starting from the N-terminus).
- f. Write the nucleotide sequence of the coding strand that corresponds to these amino acids.
- g. Write the nucleotide sequence of the template strand that corresponds to these amino acids.
- h. Using the sequence shown in the record, give the nucleotide number range that corresponds to these amino acids.

Answer 1

- d. RecQ type DNA helicase
- e. MVVA
- f. 3'tcgctgctggta5'
- g. 5'agcgacgaccat3'
- h. 5651-5662

EXERCISE 2

Genes in eukaryotes are often organized into exons and introns, which require splicing to produce an mRNA that can be translated. The gene organization is the order of the DNA segments that comprise the gene starting with the promoter, the first exon, the first intron, the second exon, and so on. The interspersed introns can make gene identification difficult in eukaryotes—particularly in higher eukaryotes with many introns and alternative spliced mRNAs. Prediction of many genes and their organization has been based on similarity searches between genomic sequence and known protein amino acid sequences and genomic sequence and the corresponding full-length cDNAs. cDNAs are reverse-transcribed mRNAs and therefore generally do not contain intron sequences. cDNAs (i.e., copied DNA) can be considered mRNAs. A comparison of a genomic sequence (with introns) to its corresponding cDNAs will reveal where introns begin and end. GenBank will contain the genomic sequence and the cDNA sequence. To find out the structure of the gene (i.e., the arrangement of the exons and introns) we simply need to perform a sequence comparison between the genomic sequence and the cDNA sequence. Shown below is a genomic sequence from the species *C. elegans*. The Basic Local Alignment Sequence Tool (BLAST) can be used to elucidate part of the gene organization (arrangement of exons and introns) of a genomic sequence. BLAST can be used to compare genomic DNA sequence with all RNA sequences (i.e., cDNA sequences) in GenBank. The top hit of the output will be a sequence comparison between your sequence (the guery sequence) and the most similar sequence in the database (subject sequence). Subsequent hits will display sequence comparisons between the query sequence and subject sequences that are increasingly less similar. If all hits have 100% identity, use the hit with the most extensive percentage coverage to report on. Use the nucleotide BLAST tool and appropriate databases to construct a schematic diagram that shows the arrangement of introns and exons in the genomic sequence. Remember that the species source of genomic sequence is Caenorhabditis elegans.

ATCTATTTATATTTACCGAATAAATATATTCATCAATTAACCTGAAGAACAAACGAATTCGGCTAC AGGCGTCGATCAGTCTCGAATCTAGTAACAACAAGAGAGCCAATACGAAAACCGGTAAATCAATAGG GGGAAGCGAAACAGTAGGTACAAATTGGAGGGGAAGCACCAATACATTAGGTGGGGGGTACGACTTG AAAAATGAGCTGATTTTCGAATAGTTAAAGCGATGATCGTGTCCGAAAAACAGTTCATTTTTCAAG A CAACATTGAGACTGGGAGTACGGGGAAGCTCATTTACGGTGAGAGGAATTGGTGAGATCTTTAGAATATGCTTAAGGAGTTGGGGTGGCTGGAGAAGTTCCTGTAGCCTCCGTGCCGGGATTCGATGGAGA AGTCGTTGCGGCTGGTCCCTTTTCCTTCACTGGTGCTGGATCCTTGGCTGGAAGACATATGCGTGGC TTGACAGTCGATGAGGTGCGAGCCGACGAGTCCTTGTGAACTTCGTATCTGGAAATATTTTACTTAGA TAGCAAATACTAAAATTGTAAAATTACCTCAAAATCTCAGTATCCGGAATGCTCAATTTCTGCTTCA $\tt AAACCTGTCCGATGCGAGATTGACATCATCGCGAGTAGCATCACGAGTCCACAAGGAAACCTTGT$ CACCCTTTTGACGAACATTCACGACAGCTCCGCAGATGTAGTCTCCGTACTCGTCGAATTGCTCTC CAACAATAGCCATCAACAGCTCCAACCAGTAGTGATCGAGCAATTGCGTTCTTCTCTGAAGCTTCTA $\tt CCTTGAACGTTGTTGACGTCCTCCCACATTGGCTTGATTCCTTGAACAAGTAATAATCGGATCC$ CCAGTTCAATCCTCCGGCAGACTGAATGTGATTGTACAGCGACCAGAAGTCCTCGACAGTGTCGAAA ${\tt AGTGAAACCATCTGGAAAAAATCGATAAAAGACGTATTTAAAAATCTTCTACCTTCAGACAATCCTC}$ CCATTCCTTGTTACGGTCAGCTTTCAAGTACCAGAGAGCCCAGCGATTCTGGAGGGGGGTGTCTGGTGA GAAGCTCTGGAGGAACTGAAGCATCGGACGCATTCACATCGCCGGAAGCTGACAATGCTTTGTTTTCC GCTACGGATGTGCTCATTTAGCTGAAAATAGGTAATATTATATACGATTAGAGCTCGGAAAACGATA ${\tt AAATAGAGAAGTATGAATTTGGTTCAAATAACTCGGATTTTATAGGAAATTTTGTTTTACTGCAC}$ ATTTTCGGCTAGTTTCCAAGCTTTTTAGATTTTTCAAGTGTAATTGGTAACATCGGGCACAATAAAT TGATATTAAAGCTTGGAAAACAATAA

Note: The sequence can be found on the publisher's website.

In addition to construction of the schematic diagram, answer the following:

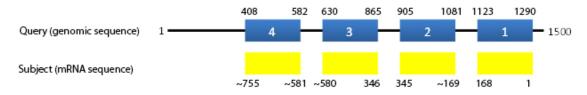
- a. Give the name and accession numbers of each distinct mRNA produced from this gene.
- b. Give the names and accession numbers of the protein product(s).
- c. Note the numbering of the sequences in the alignments. Does the database genomic sequence progress in the same direction as the database mRNA? In other words, is it the same orientation (see below):

1 114 61 98	1 0
1	1 5

- d. Consider the alignment of the query sequence and the subject sequence. What does the orientation of the sequences relative to each other tell you about the sequence that was used as the query sequence?
- e. Give the amino acid sequences separately translated from each exon sequence of the longest transcript.
- f. How many alternative splice variants are associated with this genomic sequence? List their accession numbers.
- g. Give the chromosome position numbers that denote the start and end of the *TP53* gene. The position number is the base number on chromosome 17. Calculate the length of the primary transcript. Give the lengths, in base pairs, of each exon and intron that is used for the transcription of *TP53* into mRNA isoform a. Cite the sources you used to gather your information.

Answer 2

A schematic diagram of the genomic sequence is shown below on the first line with exons shaded in blue (these exons correspond to isoform b transcript). The exon numbers at the top of diagram proceed from right to left because the given genomic sequence is the complement (also known as the negative strand or non-coding strand). These numbers mark the boundaries of the exons within the genomic sequence. A schematic diagram of the mRNA sequence is shown on the second line. The numbers at the bottom of the diagram mark the approximate nucleotide boundaries of the exons sequences (shaded yellow) where nucleotide 1 is the first nucleotide of exon sequence 1 (from BLAST analysis against RefSeq_genome sequence database—see Accession number NC_003283). Note that the numbers at the bottom may not match the output from BLAST because some nucleotides within the intron sequences at exon/intron junctions in the genomic sequence will be identical to nucleotides at the exon/exon junctions in the mRNA.



mRNA sequence is separated into exon sequences for illustration purposes. The mRNA molecule is continuous inside the cell.

- a. 1. Caenorhabditis elegans Eukaryotic translation initiation factor 4E-3 (ife-3), partial mRNA (longest mRNA) accession number: NM_070722; 2. Caenorhabditis elegans Eukaryotic translation initiation factor 4E-3 (ife-3), partial mRNA mRNA accession number: NM_171920; 3. Caenorhabditis elegans Eukaryotic translation initiation factor 4E-3 (ife-3), partial mRNA mRNA accession number: NM_070723 (shortest mRNA)
- b. 1. Eukaryotic translation initiation factor 4E-3 [Caenorhabditis elegans] accession number: NP_503123 (251 amino acids); 2. Eukaryotic translation initiation factor 4E-3 [Caenorhabditis elegans] accession number: NP_741502 (250 amino acids); 3. Eukaryotic translation initiation factor 4E-3 [Caenorhabditis elegans] accession number: NP_503124 (248 amino acids)
- c. Opposite orientation
- d. The query sequence matches the template strand of the gene.
- e. 1-168 Exon 1: Met S T S V A E N K A L S A S G D V N A S D A S V P P E L L T R H P L Q N R W A L W Y L K A D R N K E W E D C L K

 167-345 Exon 2: Met V S L F D T V E D F W S L Y N H I Q S A G G L N W G S D Y Y L F K E G I K P Met W E D V N N V Q G G R W L V V V D K Q

 346-581 Exon 3: K L Q R R T Q L L D H Y W L E L L Met A I V G E Q F D E Y G D Y I C G A V V N V R Q K G D K V S L W T R D A T R D D V N L R I G Q V L K Q K L S I P D T E I L

 408-580 Exon 4: R Y E V H K D S S A R T S S T V K P R I C L P A K D P A P V K
- f. 3 splice variants: NM_070722, NM_171920, NM_070723.
- g. This answer may change slightly depending on the version of the human genome used.

EKGPAATTSPSNPGTEATGTSPATPTP**Stop**

Chromosome position number at end of TP53 gene: 7590863 Chromosome position number at beginning of TP53 gene 7571720 Primary strand length = 7590863-7571720 + 1 = 19,144 bases

The coding strand of TP53 is on the complementary (negative) strand.

```
Gene segment Ex.1 In.1 Ex.2 In.2 Ex.3 In.3 Ex.4 In.4 Ex.5 In.5 Ex.6 (# of base pairs) (169) (10,754) (102) (117) (22) (109) (279) (757) (184) (81) (113)

Gene segment In.6 Ex.7 In.7 Ex.8 In.8 Ex.9 In.9 Ex.10 In.10 Ex.11 (# of base pairs) (567) (110) (343) (137) (92) (74) (2819) (107) (919) (1289)
```

Primary strand length = 19,144

How to answer this question: Obtained reference sequence information for p53 isoform a mRNA from Entrez gene with keywords TP53 and *Homo sapiens* (accession number: NP_000537.3). Used annotations in Refseq record of isoform a to obtain exon positions and lengths in the mRNA. Used BLAST against RefGenome database *Homo sapiens* with each exon sequence to obtain chromosome positions that correspond to beginning and end of each exon. The absolute value of the difference of the chromosome position corresponding to end of one exon and beginning of next exon minus 1 gives the number of bases in the intron. One final note: There are likely to be several approaches to arriving at the correct answer to this question.

Chapter 3

Molecular Evolution

EXERCISE 1

Obtain the human p63 isoform 1 protein sequence and the human p53 isoform a protein sequence. Use the Needleman-Wunsch program (http://www.ebi.ac.uk/Tools/psa/) to align the two proteins. Use the following parameters when you run the Needleman-Wunsch program:

Matrix: BLOSUM62; Gap Open Penalty: 10; Gap Extend Penalty: 0.5; End Gap Penalty: False; End Gap Open Penalty: 10; End Gap Extend Penalty: 0.5

- a. Which amino acid sequences within p53 are conserved in p63? (Hint: Find a region longer than 50 fifty amino acids with no gaps spanning longer than 5 five amino acids.)
- b. A second paralog of p53 is p73. Which amino acid sequences within p53 isoform a are conserved in p73 isoform a?
- c. Is the particular domain common within p53, p63, and p73 associated with a specific function?

Answer 1

- a. p53 and p63 are conserved from aa 97 through aa 312 within human p53.
- b. p53 and p73 are conserved from aa 97 through aa 291 within human p53.
- c. This domain binds DNA in a sequence-specific fashion.

EXERCISE 2

Retinoblastoma is a rare childhood cancer of the retina. The retina is located in the back of the inside of the eye and is responsible for detecting light. Alfred Knudson, at the University of Texas at Houston, observed that retinoblastomas fall into two classes. In Class I, the retinoblastoma is first diagnosed in children at the mean age of three, and the tumors are often found in both eyes. In Class II, the retinoblastoma is diagnosed at the mean age of five, and a single tumor is detected in one eye. Given what you know about germ cell mutations and somatic cell mutations, give a plausible explanation for the observation of two classes of retinoblastoma.

Answer 2

Class I cases are characterized by a germ cell mutation in one allele of the retinoblastoma gene (Rb). As the retinal cells of the infant proliferate during the first few years of life, mutations in the second Rb allele occur in more than one cell. This is a somatic cell mutation. The independent somatic cell mutations in separate cells cause several retinal cells to form tumors. In other words, there are multiple tumors in both eyes. This accounts for the observation of tumors in both eyes of the patients

and for their relative early onset. Class II cases are characterized by two somatic cell mutations in one retinal cell. Because it takes longer for two independent mutations to occur (rather than one mutation) in the retinal cell, only one tumor is created per patient and thus only one eye shows the cancer. This requirement for two independent mutations also accounts for the relatively late onset of the disease. These observations caused Knudson to formulate the "two-hit hypothesis" and led to the prediction for the existence of tumor suppressor genes.

EXERCISE 3

Peyton Rous (1879–1970) was a relatively young man when, in 1911, he discovered a virus that causes sarcomas in chickens. The virus was named Rous sarcoma virus (RSV). Later, it was found that RSV is a retrovirus that contains an oncogene, *v-src*, in its RNA genome. Perform pairwise global alignments with the Needleman-Wunsch algorithm between v-src protein and the following proteins: chicken c-src, human c-src, and mouse c-src proteins. The "c-" prefix is short for cellular. Sometimes proto-oncogenes are distinguished from viral oncogenes with the prefix "c-" and "v-," respectively. Report the identities for each pairwise alignment. Given what you know about the origin of v-src, does the result match your expectations? Use the Schmidt-Ruppin A strain of RSV as your source for the v-src sequence.

Answer 3

Chicken: 94% identity Human: 88% identity Mouse: 87% identity

The high identity shared between v-src and chicken c-src (94%) is consistent with the idea that RSV was first generated in a chicken and not a human or a mouse. Here is a possible scenario: a retrovirus captured a chicken c-src proto-oncogene in a mechanism similar to what is described in Figure 3-9. The resulting virus, called RSV, was able to out-compete the parental retrovirus for survival. RSV confers a growth advantage onto cells it infects because the c-src promotes cell growth. As the virus continues to regenerate itself, it creates a few mutations that make v-src a more powerful cell growth inducer than its c-src counterpart.

EXERCISE 4

In the Li-Fraumeni syndrome *TP53* is often mutated. Compare the mRNAs from a Li-Fraumeni patient to transcript variant 1 of wild-type p53 with the Needleman-Wunsch global alignment software program.

Matrix: DNAfull; Gap Open Penalty: 10; Gap Extend Penalty: 0.5; End Gap Penalty: False; End Gap Open Penalty: 10; End Gap Extend Penalty: 0.5

Classify the mutations in the areas that overlap. List the nucleotide location (using the wild-type p53 nucleotide numbering system as the reference) of the mutation. Indicate whether it is a point mutation or indel, transversion or transition, missense, or other type of DNA mutation. Give the amino acid substitution(s), if there is any, that occurs in the p53 protein. The p53 transcript sequence from a Li-Fraumeni patient can be obtained from GenBank with accession number BT019622.1. For a codon table, please see Chapter 1.

Answer 4Accession number of *Homo sapiens* tumor protein p53 (TP53), transcript variant 1, mRNA: NM_000546.5

	Nucleotide	Point or	Transition/	Missense, silent	-
Location(nt)	change	indel	transversion	or other	Amino acid change
417	C→G	Point	Transversion	Missensea	P72→R72
764	C→A	Point	Transversion	Missenseb	L187→M187
1384	- →A	Indel		Silent	None

^aPolymorphism

Below are segments from sequence alignment with nucleotide changes in boldface. "Li" is mRNA from Li-Fraumeni individual and "wt" is wild-type tumor protein p53 transcript variant 1.

Li	199 CCAGAGGCTGCTCCCCGCGCGCCCTGCACCAGCAGCTCCTACACCGGC	248
wt	401 CCAGAGGCTGCTCCCCCGTGGCCCCTGCACCAGCAGCTCCTACACCGGC	450
Li	549 AGATAGCGATGGT A TGGCCCCTCCTCAGCATCTTATCCGAGTGGAAGGAA	598
wt	751 AGATAGCGATGGTCTGGCCCCTCCTCAGCATCTTATCCGAGTGGAAGGAA	800
Li	1149 CATGTTCAAGACAGAAGGGCCTGACTCAGACT A G	1182
wt	1351 CATGTTCAAGACAGAAGGGCCTGACTCAGACT-GACATTCTCCACTTCTT	1399

Below are segments from sequence alignment with amino acid changes in boldface. Query is protein sequence is from Li-Fraumeni individual and subject is wild-type tumor protein p53 isoform a.

Query 61	DEAPRMPEAAP R VAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAK	120
	DEAPRMPEAAP VAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAK	
Sbjct 61	DEAPRMPEAAPPVAPAPAAPTPAAPAPAPSWPLSSSVPSQKTYQGSYGFRLGFLHSGTAK	120
Query 181	RCSDSDG M APPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGSDCTTIHYNYMCNS	240
	RCSDSDG+APPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGSDCTTIHYNYMCNS	
Sbjct 181	RCSDSDGLAPPQHLIRVEGNLRVEYLDDRNTFRHSVVVPYEPPEVGSDCTTIHYNYMCNS	240

bMissense mutation in the DNA binding domain leads to loss of tumor suppressor function.

EXERCISE 5

The sequence of two different forms of a gene starting with ATG codon is shown below. Which of the base differences in the second sequence are synonymous changes, and which are non-synonymous changes?

```
Form 1: ATGTCTCATGGACCCCTTCGTTTG
Form 1: ATGTCTCAAAGACCACATCGTCTG
```

This exercise is adapted from Hartwell et al. (2008).

Answer 5

Recall that synonymous changes are those nucleotide changes that do not change the amino acids. Nonsynonymous changes are nucleotide changes that do change the amino acids. The sequences are rewritten with corresponding amino acids.

```
Form 1:
        ATG TCT
                CAT
                     GGA CCC
                               CTT
                                   CGT
                                        TTG
        Met Ser
                 His
                     Gly
                          Pro
                               Leu
                                   Arq
                                        Leu
Form 2: ATG TCT CAA
                     AGA
                          CCA
                               CAT
                                   CGT
                                        CTG
        Met Ser Gln Arg Pro
                               His
                                   Arg Leu
```

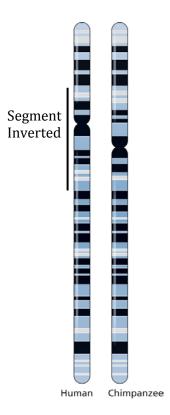
Underlined nucleotides are changed. Those in blue font are synonymous changes (no change in amino acid) and those in red font are non-synonymous changes (change in amino acid).

EXERCISE 6

Synteny is the conservation of the physical order of genetic loci in the genomes of two species. For example, a segment of mouse chromosome 2 is similar to the entirety of human chromosome 20 (chromosome 20 is one of the smaller chromosomes in humans). Cytogeneticists use chemicals to stain chromosomes to produce banding patterns. The banding patterns mark where sections of genes are located. Similar banding patterns denote regions of synteny. The figure shows the banding patterns on human chromosome 5 and the chimpanzee chromosome 5. Describe the region(s) where an inversion event occurred. Exercise adapted from Lesk (2008).



Answer 6



The region surrounding the centromere is inverted. Interestingly, many inversion mutations can account for the differences between humans and other primates.