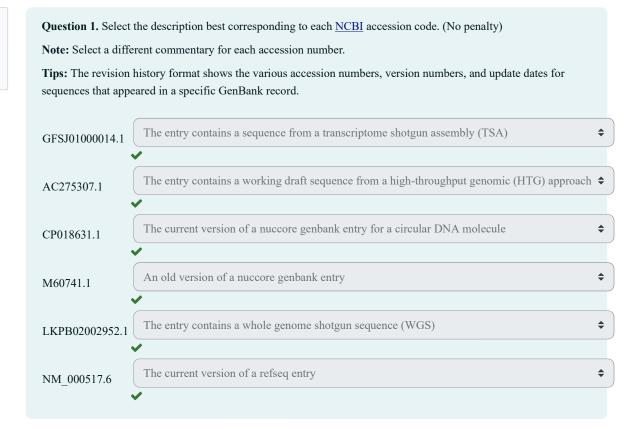
Started on	Thursday, 5 October 2023, 9:26 AM
State	Finished
Completed on	Monday, 9 October 2023, 11:08 PM
Time taken	4 days 13 hours
Marks	14.43/20.00
Grade	<b>7.21</b> out of 10.00 ( <b>72.14</b> %)

Correct

Mark 1.00 out of 1.00



The correct answer is: GFSJ01000014.1  $\rightarrow$  The entry contains a sequence from a transcriptome shotgun assembly (TSA), AC275307.1  $\rightarrow$  The entry contains a working draft sequence from a high-throughput genomic (HTG) approach, CP018631.1  $\rightarrow$  The current version of a nuccore genbank entry for a circular DNA molecule, M60741.1  $\rightarrow$  An old version of a nuccore genbank entry, LKPB02002952.1  $\rightarrow$  The entry contains a whole genome shotgun sequence (WGS), NM\_000517.6  $\rightarrow$  The current version of a refseq entry

Illumina Hi	Seq2500.
PacBio.✓	
Ion torrent.	
Oxford nand	opore.  /no answer. (without penalty)
Don't know/	no answer. (without penaity)
Mark 1.00 o	out of 1.00
The correct	answer is: PacBio.
	im of the sequencing project in which this sequence was obtained?
	e only the X and Y chromosomes in a family.
-	expression of human genes.
	the human genome for the very first time.
	rse allelic variation to the reference human genome.
Don't know/	(no answer (without penalty)
Mark 1.00 o	out of 1.00
The correct	answer is: To add diverse allelic variation to the reference human genome.
GenBank files i	usually contains this information.
ttps://www.nc	bi.nlm.nih.gov/nuccore/LKPB02002952.1
equencing Te	chnology :: PacBio
	this project will be used to improve the contiguity of the human reference add diverse allelic variation
ou can also fii	nd the second answer in the bioproject link:
ttps://www.nc	bi.nlm.nih.gov/bioproject/PRJNA288807
	'111
or the 1st one	you will have to go to ASSEMBLY from the "related information" link in the bioproject page:

And select "NA19240\_prelim\_3.0" which will send you to the entry in the Assembly database:

 $https://www.ncbi.nlm.nih.gov/assembly/GCA\_001524155.4$ 

Question 2

Correct

2.00

Mark 2.00 out of

# Question 3. Two of these search tags will help you finding sequences that are codifying at nucleotide database (CDS). Which ones? Mark 1.00 out of Tips: Learn about Search Field Descriptions for Sequence Database in this <u>link</u>. Select one or more: a. protein\_nucleotide[filter] b. All[protein] ✓ c. nucleotide\_protein[filter] ✓ ✓ d. CDS[fkey] ✓ e. protein[fkey] \_\_\_ f. protein[filter] nucleotide protein[filter] will deliver sequences that are in nucleotide (AKA: nuccore) and have a link to protein. CDS[fkey]: sequences that contain one or more CDS features annotated however you must take into account that not all annotated CDSs have a link to the protein DB The correct answers are: nucleotide\_protein[filter], CDS[fkey] Question 4. Localize all of the records in the NCBI nucleotide database containing rabbit genomic sequences that Partially correct are codifying. How many have you found? Mark 2.00 out of Optionally write here your search strategy rabbit[Organism] AND biomol\_genomic[PROP] AND CDS[fkey] (no rate) ¿How many of these sequences come from the mitochondrial compartment? Here the question asked for genomic sequences, then, biomol\_genomic[PROP]. rabbit[porgn] AND biomol\_genomic[PROP] AND CDS[fkey] This is not the right answer: rabbit[PORGN] AND biomol\_genomic[PROP] AND nucleotide\_protein[filter] the fact that their products don't have a link to the protein DB does not mean that they are not there.

Question 3

Question 4

2.00

Correct

1.00

In the left margin you can see that 159 are in the Mithocondrion OR use this search: rabbit[porgn] AND biomol genomic[PROP] AND CDS[fkey] AND mitochondrion[filter]

Partially correct

Mark 1.43 out of 5.00

Question 5. Do you think you have found in Question 4 the NCBI nucleotide reference entry for the complete rabbit mitochondrial chromosome? Yes ❖ ✓

If yes, what is its refseq access number? AJ001588.1

How many proteins does it code? 13 ✓

Between what coordinates do we find the tRNA for tryptophan? From 4971 ✓ to 5037 ✓

Click on the filter "Mithochondiron" in the question 4 results and select the only refseq result refering to the Mithocondrion complete genome "Oryctolagus cuniculus mitochondrion, complete genome". NC\_001913.1 (https://www.ncbi.nlm.nih.gov/nuccore/NC\_001913.1). You can recognize it is a code from RefSeq database due to its starting characters "NC\_".

Or perform a search like this:

rabbit[porgn] AND biomol genomic[prop] AND CDS[fkey] AND refseq[filter] AND mitochondrion[filter]

Use "Related information" box and click in protein: (https://www.ncbi.nlm.nih.gov/protein?LinkName=nuccore\_protein&from\_uid=5835526) to obtain results for the proteins encoded in this entry.

In the gbff file (https://www.ncbi.nlm.nih.gov/nuccore/NC\_001913.1) search for "tRNA-Trp"

tRNA 4971..5037
/product="tRNA-Trp"

### Question 6

Partially correct
Mark 1.00 out of
1.00

Question 6. ¿How many complete mammalian mitochondrial chromosomes can be found at the NCBI nucleotide database? 77718 ✓

Optionally write here your search strategy

mammals[filter] AND "complete genome"[title] also select mitochondrion in (no rate)

**Tips**: Mitochondria has its own genome. Researchers and/or curators use to add to the title of entries containing complete genomic sequences the phrase "complete genome" for fully sequenced chromosomes or genomes of eukaryotic organelles.

mammalia[porgn] AND mitochondrion[filter] AND "complete genome"[TI]

mitochondrion[filter] and gene in mitochondrion[PROP] both give the same result.

txid40674 and mammalia are synonyms.

"complete genome" should be added using [title] or [word] as search tag or simply nothing. This is an approximation to the best way to get the most results. There is not a search tag to filter for the status of a genome sequencing project at NCBI nucleotide.

Correct

Mark 1.00 out of 1.00

**Question 7.** One of the proteins encoded in the mitochondrial chromosomes is cytochrome b. The best strategy to search for sequences encoding for cytochrome b proteins in NCBI nucleotide database is to include in a query the expression:

#### Select one:

- a. AND (cytochrome b[protein name]) ✓
- O b. AND "cytochrome b"[title]
- O c. NOT cytochrome\*
- od. AND cytochrome\*[title]
- oe. Don't know/no answer (without penalty)

in this case you should evaluate what would happen if you add each of the phrases in a search strategy:

AND "cytochrome b"[title]; with this phrase you are including other proteins like cytochrome b reductase (e.g. NM 001011954.1).

AND (cytochrome b[protein name]); this strategy will retrieve all nucleotide sequences encoding for a protein named exactly as "cytochrome b", including complete bacterial genomes and eukaryotic mitochondrial genomes.

AND cytochrome\*[title]; this strategy includes other variants like cytochrome c.

NOT cytochrome\*; this strategy will exclude cytochrome b from your search

The correct answer is: AND (cytochrome b[protein name])

## Ouestion 8

Correct

Mark 2.00 out of 2.00

Question 8. Find out the NCBI genomic reference sequence encoding for human (us) cytochrome b protein.

Indicate the NCBI accession numbers for both the genomic sequence

NC 012920.1

and

the encoded protein YP\_003024038

#### Strategy

"human" [Primary Organism] AND "cytochrome b" [protein name] AND refseq[filter]

Sequence with accession number <u>NC\_012920.1</u> is the reference sequence in refseq for the Homo sapiens mitochondrion complete genome. This is the reference DNA sequence encoding, among other proteins, for the human cytochrome b protein (YP 003024038.1).

AC\_000021 is an obsolete version of NC\_012920, however, if you have answred "NC\_012920, AC\_000021" this is counted as correct. Please notify me if you have introduced a different variation

Question 9 Correct	<b>Question 9.</b> By querying the NCBI nucleotide database investigate if bacteria genomes encode for cytochrome b protein.					
Mark 1.00 out of 1.00	Tips: <u>Enterobacteria</u> and <u>Cyanobacteria</u> are two different groups of <u>bacteria</u> .					
	Select one:					
	a. Yes, cytochrome b is exclusive of enterobacteria.					
	b. No, bacteria do not have cytochrome b protein.					
	c. Don't know/no answer (without penalty).					

⊚ e. Yes, several bacterial groups encode for a protein named cytochrome b.

Od. Yes, but only in the group of Cyanobacteria.

bacteria[porgn] AND "cytochrome b"[protein] You will get more than 128000 entries in almost all bacterial groups.

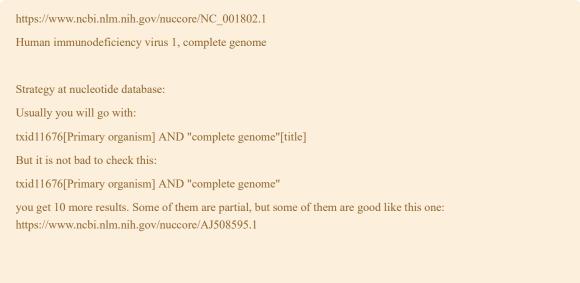
Enterobacteria and Cyanobacteria are two groups of bacteria. If you indicate the following search strategies you will find fewer number of entries.

cyanobacteria[porgn] AND "cytochrome b"[protein]

enterobacteria[porgn] AND "cytochrome b"[protein]

The correct answer is: Yes, several bacterial groups encode for a protein named cytochrome b.

## Question 10 Question 10. Complete microbial genomes at NCBI. Partially correct 10.1. Open entry with accession number NC 001802.1. What does this sequence represent? Mark 2.00 out of 2.00 A partial sequence for a viral genome. A complete genomic sequence for a human gene. A reference sequence for a complete viral genome. ✓ The reference sequence for the human genome. O Don't know/no answer (without penalty) Mark 1.00 out of 1.00 The correct answer is: A reference sequence for a complete viral genome. 10.2. It is possible to find complete genome sequences in the NCBI nucleotide database for any organism. Create a search strategy trying to find at the nucleotide database all the complete or nearly complete genomic sequences for Human immunodeficiency virus type 1. How many entries have you found? Optionally write here your search strategy txid11676[Organism] AND "complete genome"[Title] (no rate) Tips: • Use the taxid instead of the full name for the Human immunodeficiency virus type 1 • Researchers and/or curators use to add to the title of entries containing complete or nearly complete genome sequences the phrase "complete genome"



Incorrect

Mark 0.00 out of 2.00

Question 11. There are several genomes at the NCBI without annotations. That means none of the genetic elements like CDSs, genes or any other features annotated as "misc\_feature", have been indicated in the sequence. Search the records in the NCBI nucleotide database for HIV-1 (Human immunodeficiency virus type 1) complete genomes without any annotation.

How many entries have you found? 4491

Indicate here the accesion.version number of only one of them.

AB253704.1

Optionally indicate here your search strategy

((txid11676[Organism]) AND "complete genome"[Title]) NOT "misc feature (no rate)

#### Option 1.

txid11676[porgn] AND "complete genome"[title] NOT CDS[fkey] NOT Gene[fkey] NOT misc feature[fkey]

Only 5 entries fulfill these criteria. JN571034.1, AY781127.1, AY781125.1, AY781128.1 and AY781126.1.

UPDATE: now it is 10!

JN571034.1, AY781127.1, AY781125.1, AY781128.1, AY781126.1, ON245430.1, ON245428.1, ON245429.1, ON245427.1, ON245431.1

so both answers are considered right.

However, there is a second option:

txid11676[Primary organism] AND complete genome NOT CDS[fkey] NOT Gene[fkey] NOT misc feature[fkey]

Here we are not forcing the words "complete" and "genome" neither to the title nor to be together so we can view more results

We use to get 21 good entries (2LDL A is not a complete genome) with accession numbers:

JX503079.1 M93259.1 M93258.1 JX503077.1 JX503078.1

JX503080.1 JX503073.1 JX503074.1

JX503082.1

JX503071.1 JN571034.1

JX503072.1

JX503075.1

JX503081.1

JX503076.1

JX503083.1

AY781127.1

AY781125.1

AY781128.1

AY781126.1 MK457954.1

UPDATE: now is 26 (again 2LDL\_A is not a complete genome)

JX503079.1			
M93259.1			
M93258.1			
JX503077.1			
JX503078.1			
JX503080.1			
JX503073.1			
JX503074.1			
JX503082.1			
JX503071.1			
JN571034.1			
JX503072.1			
JX503075.1			
JX503081.1			
JX503076.1			
JX503083.1			
AY781127.1			
AY781125.1			
AY781128.1			
AY781126.1			
ON245430.1			
ON245428.1			
MK457954.1			
ON245429.1			
ON245427.1			
ON245431.1			