

Indraprastha Institute of Information Technology Delhi (IIITD)
Department of Computational Biotechnology

BIO213 – Introduction to Quantitative Biology

Rubrics_Quiz-3 (October 21, 2025)

Total time: 50 mins, Total marks: 40

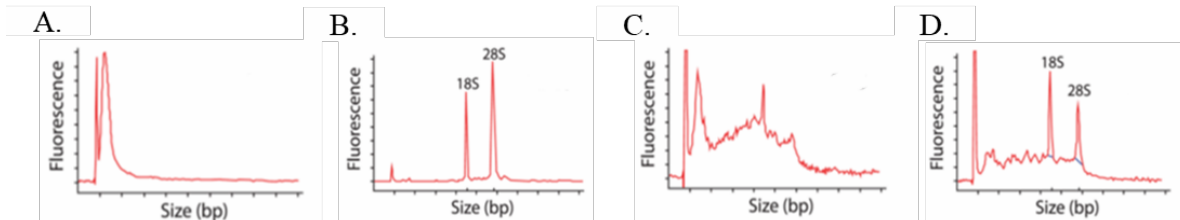
Name: _____

Roll No. _____

Choose the correct answer for the following questions (1-12):

1. In RNA-seq, biological replicates are important because **(1 mark)**
 - a. Technical variability is higher than biological variability
 - b. They increase statistical power by capturing biological variability
 - c. They increase sequencing read length
 - d. They remove GC bias
2. Loop modelling is difficult because loops **(1 mark)**
 - a. Are always missing from templates
 - b. Have high conformational variability and depend on long-range packing
 - c. Contain no hydrogen bonds
 - d. Break templates
3. A statistically significant p-value (<0.05) in Differential Expression analysis means **(1 mark)**
 - a. The null hypothesis is definitely false
 - b. There is less than 5% probability that the observed difference occurred by chance
 - c. The alternative hypothesis is false
 - d. The experiment must always be repeated
4. Adapter contamination appears in the quality check as: **(1 mark)**
 - a. Uniform peaks across cycles
 - b. Sudden drop in sequence complexity
 - c. Overrepresented short sequences
 - d. Increased GC-bias
5. Why do *ab initio* models frequently become trapped in incorrect conformations? **(1 mark)**
 - a. No phi (ϕ) and psi (ψ) angles are allowed
 - b. Energy landscapes have many local minima
 - c. Templates are missing
 - d. Hydrogen atoms are ignored
6. During DE analysis, log fold change is preferred over raw fold change because **(1 mark)**
 - a. RNA-seq cannot c
 - b. compute raw fold change
 - c. Log scaling stabilises variance and is symmetric for up/down-regulated genes
 - d. It removes GC-bias
 - e. It increases counts

7. Which of the following statements about the Ramachandran plot is INCORRECT? (1 mark)
- It plots phi (ϕ) and psi (ψ) dihedral angles of amino acid residues
 - Glycine has the most restricted conformational space due to its side chain
 - Sterically allowed regions indicate no atomic clashes
 - Proline has limited phi angle rotation due to its cyclic structure
8. In RNA-seq analysis, what is the PRIMARY reason for performing multiple testing corrections? (1 mark)
- To increase statistical power
 - To reduce false positives when testing thousands of genes simultaneously
 - To normalise for sequencing depth
 - To account for biological variability
9. A Phred quality score of 30 indicates a base call accuracy of? (1 mark)
- 90%
 - 99%
 - 99.9%
 - 99.99%
10. Which statement about microarrays vs RNA-seq is CORRECT? (1 mark)
- Microarrays have better dynamic range than RNA-seq
 - RNA-seq requires prior knowledge of the transcriptome
 - RNA-seq can detect novel transcripts and gene fusions
 - Microarrays are better for detecting single-nucleotide polymorphisms
11. Which of the following is true with respect to RNA-seq differential gene expression analysis? (2 marks)
- The chances of finding false positives increase with the increasing number of genes.
 - For 30000 genes and a p-value of 0.05, one would expect 1500 genes to have a p-value < 0.05 by chance.
 - p-value is adjusted using methods like false discovery rate.
- 1 and 2
 - 1 and 3
 - 2 and 3
 - All of these
 - None of these
12. Which of the following indicates the best and the worst quality RNA? (2 marks)



- Best- A, Worst- D
- Best- B, Worst- A
- Best- B, Worst- C
- Best- A, Worst- D

Answer the following question:

13. A researcher sequences two RNA samples with the following results: Sample A has 10 million reads, Sample B has 50 million reads. Both have 5000 reads mapping to Gene X (length 2 kb). Calculate the RPKM for Gene X in both samples and explain if there's differential expression.

(5 marks)

Ans: Sample A total reads $N_A = 10,000,000$.
Sample B total reads $N_B = 50,000,000$.
Reads mapping to Gene X $C = 5,000$.
Gene length $L = 2 \text{ kb} = 2,000 \text{ bp}$.

RPKM formula: $RPKM = (10^9 \times C) / (N \times L)$.

Calculate Sample A step-by-step:

Numerator: $10^9 \times 5,000 = 5 \times 10^{12}$.

Denominator: $10,000,000 \times 2000 = 10^7 \times 2 \times 10^3 = 2 \times 10^{10}$.

RPKM (A) = $5 \times 10^{12} \div 2 \times 10^{10} = (5/2) \times 10^2 = 2.5 \times 10^2 = 250$.

Calculate Sample B:

Numerator: 5×10^{12} (same).

Denominator: $50,000,000 \times 2000 = 5 \times 10^7 \times 2 \times 10^3 = 1 \times 10^{11}$.

RPKM (B) = $5 \times 10^{12} \div 1 \times 10^{11} = 50$.

RPKM for Gene X — Sample A = 250, Sample B = 50.

Interpretation: Gene X appears 5-fold higher in Sample A than in Sample B ($250/50 = 5$; $\log \text{fold-change} \approx 2.32$).

14. Which of the following is not true about the template Selection Step? Justify your answer.
- The first step in protein structural modelling is to select appropriate structural templates
 - This forms the foundation for the rest of the modelling process.
 - There is no use of heuristic alignment search programs.
 - The template selection involves searching the non-redundant protein sequence database for homologous proteins to be used as templates.

(5 marks)

Ans: **C. NOT TRUE (Incorrect Statement)**

Template selection **heavily depends** on heuristic alignment search programs such as:

- BLAST
- PSI-BLAST
- HMMER / HHsearch / HHblits

These programs are essential for:

- Rapid similarity searches
- Identifying homologous proteins
- Detecting distant evolutionary relationships

Therefore, the statement "*There is no use of heuristic alignment search programs*" is **false**.

D. NOT TRUE (Incorrect Statement)

Template selection **does not typically use the non-redundant (NR) protein sequence database**.

Instead, it uses **structural databases** such as:

- **PDB (Protein Data Bank)**
- **SCOP / CATH**

- **PDB70 / PDB100 chains** (non-redundant structural sets)
NR is a *sequence* database, not a *structural* one.
Template selection requires **structurally solved proteins**, not just sequences.

15. Use the information given in the table to find out the following value: **(8 marks)**

	β-Sheet	Helix	Others	Total
G	200	150	30	380
V	100	400	180	680
All residues	600	2100	500	3200

a. $P(SS = \beta\text{-Sheet} \mid aa = G)$

Ans: $P(SS = \beta\text{-Sheet} \mid aa = G) = 200 / 380$
Simplify: divide numerator & denominator by 20 $\rightarrow 10 / 19$
Decimal: $10 \div 19 = 0.526315$ (6 d.p.)
Per cent: 52.6316% (4 d.p.)

b. $P(SS = \text{Helix} \mid aa = V)$

Ans: $P(SS = \text{Helix} \mid aa = V) = 400 / 680$
Simplify: divide by 20 $\rightarrow 20 / 34 = \text{divide by } 2 \rightarrow 10 / 17$
Decimal: $10 \div 17 = 0.588235$ (6 d.p.)
Per cent: 58.8235% (4 d.p.)

c. $P(SS = \sim \beta\text{-Sheet} \mid aa = G)$

Ans: $P(SS \neq \beta\text{-Sheet} \mid aa = G) = 1 - P(\beta \mid G) = (380 - 200) / 380 = 180 / 380$
Simplify: divide by 20 $\rightarrow 9 / 19$
Decimal: $9 \div 19 = 0.473684$ (6 d.p.)
Per cent: 47.3684% (4 d.p.)

d. $P(SS = \sim \text{Helix} \mid aa = V)$

Ans: $P(SS \neq \text{Helix} \mid aa = V) = 1 - P(\text{Helix} \mid V) = (680 - 400) / 680 = 280 / 680$
Simplify: divide by 20 $\rightarrow 14 / 34 = \text{divide by } 2 \rightarrow 7 / 17$
Decimal: $7 \div 17 = 0.411765$ (6 d.p.)
Per cent: 41.1765% (4 d.p.)

e. $P(SS = \beta\text{-Sheet})$

Ans: $P(SS = \beta\text{-Sheet}) = \beta_{\text{total}} / \text{grand total} = 600 / 3200$
Simplify: divide by 200 $\rightarrow 3 / 16$
Decimal: $3 \div 16 = 0.187500$ (6 d.p.)
Per cent: 18.7500% (4 d.p.)

f. $P(SS = \sim \beta\text{-Sheet})$

Ans: $P(SS \neq \beta\text{-Sheet}) = 1 - P(\beta) = (3200 - 600) / 3200 = 2600 / 3200$
Simplify: divide by 200 $\rightarrow 13 / 16$
Decimal: $13 \div 16 = 0.812500$ (6 d.p.)
Per cent: 81.2500% (4 d.p.)

g. $P(SS = \text{Helix})$

Ans: $P(SS = \text{Helix}) = \text{Helix_total} / \text{grand total} = 2100 / 3200$
Simplify: divide by 100 $\rightarrow 21 / 32$
Decimal: $21 \div 32 = 0.656250$ (6 d.p.)
Per cent: 65.6250% (4 d.p.)

h. $P(SS = \sim\text{Helix})$

Ans: $P(SS \neq \text{Helix}) = 1 - P(\text{Helix}) = (3200 - 2100) / 3200 = 1100 / 3200$
Simplify: divide by 100 $\rightarrow 11 / 32$
Decimal: $11 \div 32 = 0.343750$ (6 d.p.)
Per cent: 34.3750% (4 d.p.)

16. State whether the following statements are correct or incorrect? In case of an incorrect statement, justify your answer. **(5 marks)**

- a. Conformational search algorithm in *ab initio* protein structure modelling explores the potential energy surface and locates the local minimum.

Ans: **Incorrect.**

Justification: Ab-initio methods (e.g., Monte Carlo, simulated annealing, fragment assembly) explore the protein's potential energy landscape to find low-energy conformations; these searches typically locate local minima (ideally the global minimum, but often local minima).

- b. Logs put positive and negative values of fold changes on a symmetric scale.

Ans: **Correct.**

Justification: Log transformation (commonly log₂) converts fold changes so up- and down-regulation are symmetric (e.g., FC=4 \rightarrow +2; FC=0.25 \rightarrow -2), which stabilises variance and simplifies interpretation.

- c. The correct sequence of steps involved in RNA-seq analysis includes library preparation, sequencing, quantification, read mapping, and differential expression analysis.

Ans: **Incorrect.**

Justification: The correct order is **library preparation \rightarrow sequencing \rightarrow read mapping \rightarrow quantification \rightarrow differential expression analysis**. Quantification (counting reads per feature) requires mapped reads, so it comes after read mapping.

- d. Technical replicates generally increase statistical power more than biological replicates.

Ans: **Incorrect.**

Justification: Biological replicates capture true biological variability and therefore increase statistical power for detecting real differences. Technical replicates only reduce measurement noise and have much less impact on power for biological hypotheses.

- e. Total RNA extracted from the cells can be directly used for sequencing.

Ans: **Incorrect.**

Justification: Total RNA must be processed into a sequencing library (e.g., rRNA depletion or

poly(A) selection, fragmentation, cDNA synthesis, adapter ligation). Raw total RNA cannot be directly loaded onto sequencers.

17. Differentiate between diagnostic and prognostic biomarkers with the help of an appropriate example. **(3 marks)**

Ans: **Diagnostic Biomarkers**

- These biomarkers are used to detect or confirm the presence of a disease.
- They help in early identification or accurate diagnosis.
Example: Elevated troponin levels are a diagnostic biomarker for myocardial infarction (heart attack).

Prognostic Biomarkers

- These biomarkers provide information about the likely course, outcome, or progression of a disease.
- They predict disease severity, recurrence, or survival.
Example: HER2 overexpression in breast cancer is a prognostic biomarker indicating aggressive tumour behaviour and poorer prognosis.