

## ASSESSMENT 9 DATABASES

### Scenario 2

RNA-seq is a powerful means to assess the quantity of RNA in a cell. By determining the “transcriptome” in cells at a given time or under particular conditions, we can uncover a great amount of information about the cells’ activities. Furthermore, by comparing quantifications for different conditions, we can potentially identify the genes that are involved in mechanisms that are important in the response to these changing conditions. Your group is trying to understand the development of pancreatic cancer cells in order to develop targeted treatments. As a first step, you will find data on how gene expression is affected in these cells compared to healthy cells. You would also like to know whether any mutations have occurred in transcription factors or in their binding sites in pancreatic cancer cells. Finally, in order to assess the possible functional implications of this dysregulation of expression, you will determine whether the affected transcripts encode proteins that interact with each other. You must now determine what data you would need and how you would integrate that data from disparate sources in order to complete your analysis. Assume that all the required data is publicly available.

**As a first step, you will find data on how gene expression is affected in these cells(*pancreatic cancer cells*) compared to healthy cells.**

**You would also like to know whether any mutations have occurred in transcription factors or in their binding sites in pancreatic cancer cells.**

**Finally, in order to assess the possible functional implications of this dysregulation of expression, you will determine whether the affected transcripts encode proteins that interact with each other.**

**You must now determine what data you would need and how you would integrate that data from disparate sources in order to complete your analysis.**

**Assume that all the required data is publicly available.**

- 1. Read the scenario**
- 2. Identify all types of data that are required in order to complete the analysis.**
- 3. Assign a data type to each member of the group and keep the implementation details of each data source independent and secret until all of them have been completed**
- 4. For each data type, describe an imaginary data source. This can be a published data set or online database/service. If you are already familiar**

**with an appropriate, real data source, feel free to describe it instead of an imaginary one. Do not worry about specifying technical details. Avoid the temptation to produce a data source that is exactly tailored to your needs!**

**Please consider the following questions about the data source.**

- **Is it from a single publication or is it a repository of multiple data sets?**
  - **If it is from a single publication, in what year was it published?**
  - **If it is a repository, how frequently is it updated? When was the last update?**
- **Is the data produced experimentally, computationally or both? If multiple data types are provided, specify which are experimental and which are computationally predicted.**
- **Which species is or are represented in the data?**
- **What is the main entity or entities of the data source? That is, what would the user primarily search for in order to find the data that they need (e.g., genes, proteins, domains, etc.)?**
- **How are these entities identified? For example, if it were a genome database, what would the individual gene identifiers (unique names) look like? Is this identifier format unique to this data source? Does the data source provide mappings to other common identifier formats?**
- **Transcriptomic Data:**
  - **Type:** RNA-seq data for pancreatic cancer cells and healthy cells.
  - **Imaginary Data Source:** This data source for this would come from a single publication, published in 2020. This data was produced both experimentally and computationally. For this publication, the only species represented were humans, mice and horses. The main entity of the data source were genes, whose unique identifiers are the official gene symbol provided by HGNC (for example BRCA1, TP53 or APOE) as can be found in other databases like NCBI. The output is a RNA nucleotide sequence with uracil instead of thymine.
- **Mutation Data:**
  - **Type:** Information about mutations in transcription factors or their binding sites in pancreatic cancer cells.
  - **Imaginary Data Source:** This data source would come from a repository of multiple data sets, updated every six months, whose last update was June 2023. As this is a repository with many data sources,

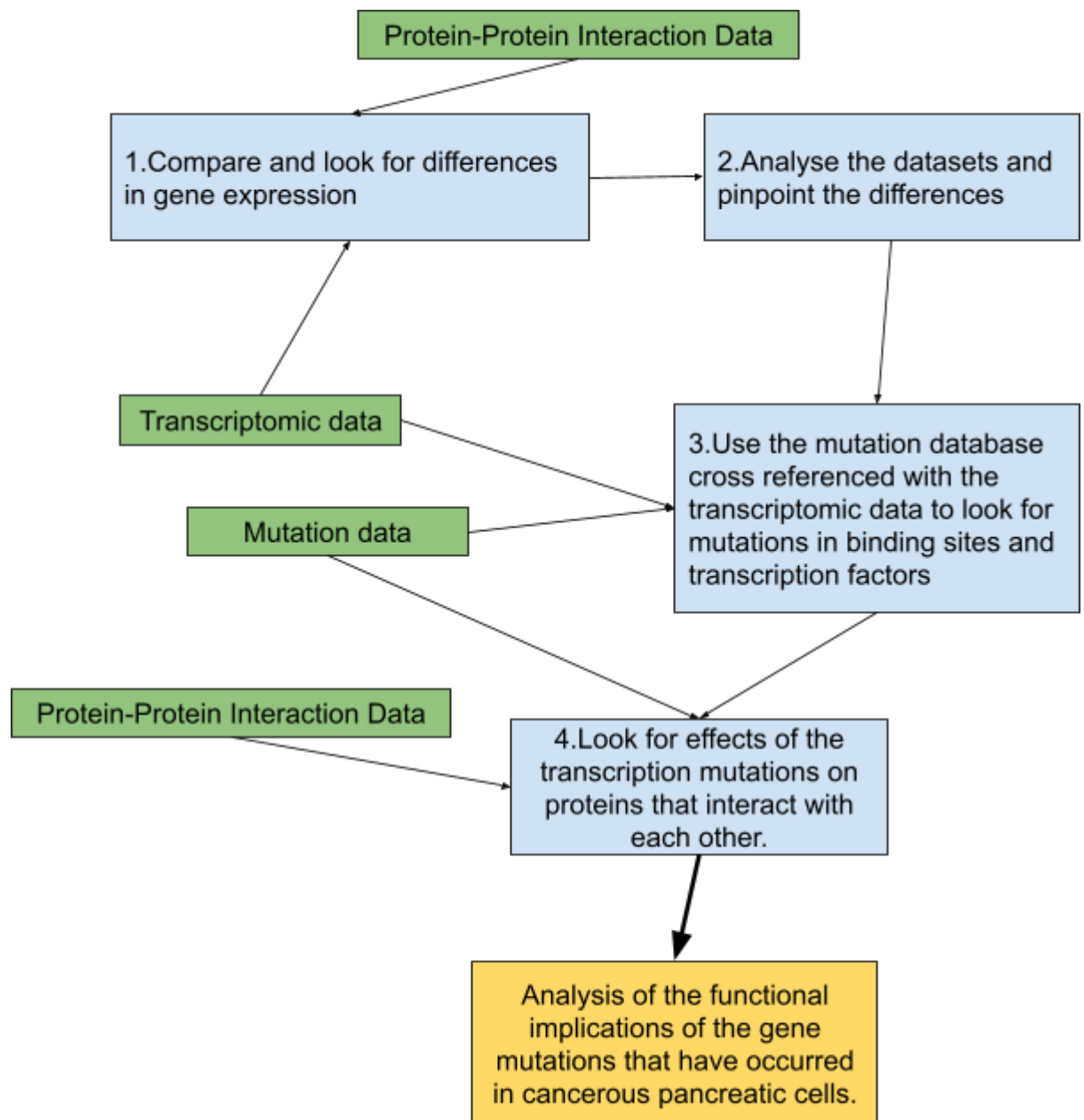
there are data sources that are obtained experimentally, computationally or both. There are many species represented in this imaginary repository. The identifiers used are gene symbols in a style of their own creation. provided the entry of every gene symbol is the gene suffering the mutation explained in that entry. The imaginary database offers mapping to other databases that describe the gene and the protein encoded in both the wild type and the mutated kind.

- **Protein-Protein Interaction Data:**

- **Type:** Data on protein-protein interactions.
- **Imaginary Data Source:** This data source would be a public repository containing multiple data sets. This repository is updated monthly as many data is introduced daily, the last update then was last month. In this repository the only data accepted is the data that is provided both experimentally and computationally predicted, as this ensures the most confident data. There's many species represented. The main entity of the data source are proteins, and its identifiers are UniProt ID of every protein, thus the database is mapped to UniProt directly. The output is an amino acid sequence

**5. With the data sources now completed, work together to draw a schematic of how your analysis “pipeline” would look.**

The schematic would look like:



The databases(green) used in each step(blue) are indicated by the arrows pointing from those databases to that step. The last step is in the yellow box

**6. Are there any incompatibilities between the data sources? How would you resolve them?**

Indeed, there are incompatibilities between the data sources, primarily arising from Mutated data, which possesses its own set of IDs and gene symbols. However, if we incorporate a mechanism within the databases that allows us to translate the IDs or gene symbols to the specific ones used in the Mutated data, the issue can be effectively resolved. This translation feature would facilitate seamless integration and correlation between the various datasets.

**7. Do the data from any of the data sources need to be transformed in order to be useful for your analysis?**

Yes, the data from Transcriptomic data, for instance, is an RNA sequence, so we would need to transform it into a DNA sequence in order to have thymine instead of uracil, so it can be analysed as a DNA sequence. This in order to analyse it with other DNA sequences

Also, in the Protein-Protein interaction data, the data is an amino acid sequence, so we would need to transform it into nucleotides in order to analyse the sequence as DNA. This in order to analyse it with other DNA sequences