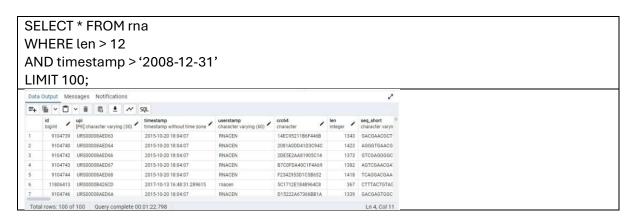# Uzair Khan

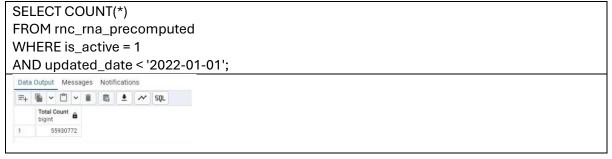# BWT - Data Engineering

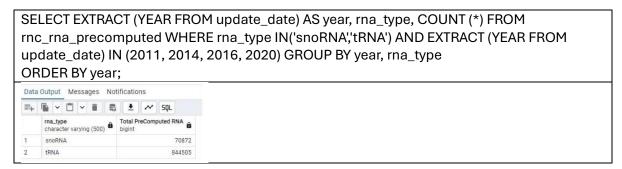# Task 02:

1. **Write a query to get data having length of Rna structures more than 12 with them being added after 2008.**

```
SELECT * FROM rna
WHERE len > 12
AND timestamp > '2008-12-31'
LIMIT 100;
```



2. **How many pre computed RNA are present that are still active and got their last release update before 2022:**

```
SELECT COUNT(*)
FROM rnc_rna_precomputed
WHERE is_active = 1
AND updated_date < '2022-01-01';
```



3. **How many total pre computed RNA records for snoRNA and tRNA were recorded in 2011, 2016, 2014, and 2020:**

```
SELECT EXTRACT (YEAR FROM update_date) AS year, rna_type, COUNT (*) FROM
rnc_rna_precomputed WHERE rna_type IN('snoRNA','tRNA') AND EXTRACT (YEAR FROM
update_date) IN (2011, 2014, 2016, 2020) GROUP BY year, rna_type
ORDER BY year;
```

4. **Can you give me the names of all databases built for RNA with minimum length other than 100, 200, 300, 400, and 15**

```
SELECT display_name
FROM rnc_database
WHERE min_length NOT IN (100, 200, 300, 400, 15);
```

Data Output  Messages  Notifications

| | display_name<br>character varying (60) |
|---|---|
| 1 | ENA |
| 2 | GENCODE |
| 3 | MGnify |
| 4 | GeneCards |
| 5 | RDP |
| 6 | snoRNA Database |
| 7 | Rfam |

5. **Can you get complete 500 records of sequences for active regions and name your column as myregions in which you are getting the region name column vakue. Then tell me what different chromosomes with exon_count we have for regions including center, east and north using the name you set for your column**

```
SELECT region_name AS myregions FROM rnc_sequence_regions WHERE region_name IS
NOT NULL LIMIT 500;

SELECT DISTINCT chromosome, exon_count FROM rnc_sequence_regions WHERE
region_name IN ('center', 'east', 'north') LIMIT 1000
```

Data Output  Messages  Notifications

| | myregions<br>text |
|---|---|
| 1 | URS00006F9F83_10020@KN672353.1/2907709-2907773:- |
| 2 | URS00006F9F83_10020@KN672353.1/2908196-2908260:+ |
| 3 | URS00006F9F83_10020@KN672353.1/3567704-3567768:- |
| 4 | URS00006F9F83_10020@KN672353.1/4727654-4727718:+ |
| 5 | URS00006F9F83_10020@KN672353.1/5119365-5119429:+ |
| 6 | URS00006F9F83_10020@KN672353.1/515849-515913:+ |