

# TP2 : manipulation relationnelle (SQL)

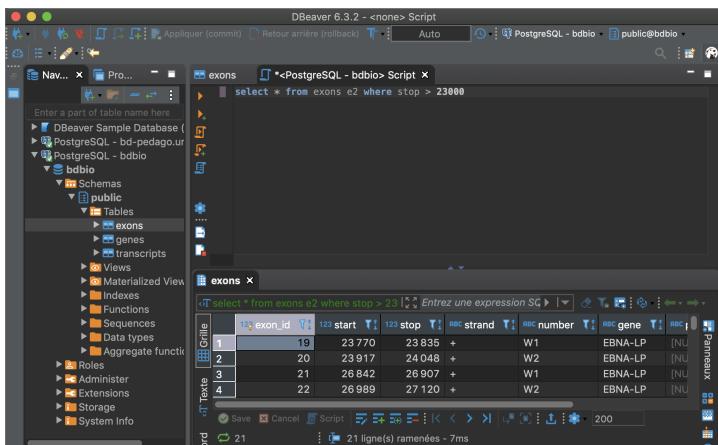
UCBL - Base de données pour la bioinformatique - 2023 / 2024

Objectif du TP : écrire des requêtes SQL sur une BD relationnelle

## 1 Préparation de l'environnement

Nous utiliserons le SGBD [PostgreSQL](#), disponible sur le serveur *bd-pedago*<sup>1</sup>. Pour interagir avec ce serveur, nous conseillons l'utilisation de l'un des outils suivants (libres et multi-plateformes) :

- [DBeaver](#) est une interface graphique (Java) pour différents SGBD (dont SQLite et PostgreSQL). Disponible sur les machines du campus via le menu du système d'exploitation ;
- [pgAdmin](#) est l'interface graphique (web) de PostgreSQL. Normalement disponible sur les machines du campus ;
- [psql](#) est l'interface en ligne de commande de PostgreSQL.



Capture d'écran de DBeaver

```
psql (12.1, server 11.5 (Ubuntu 11.5-3.pgdg18.04+1))
SSL connection (protocol: TLSv1.3, cipher: TLS_AES_256_GCM_SHA384, bits: 256)
Type "help" for help.

bdbio=> \d
          List of relations
 Schema |   Name    | Type | Owner
-----+-----+-----+
 public | exons   | table | bdbio
 public | genes   | table | bdbio
 public | transcripts | table | bdbio
(3 rows)

bdbio=> select * from exons limit 10;
      exon_id | start | stop | strand | rec_number | rec_gene | rec_product | locus_tag
-----+-----+-----+-----+-----+-----+-----+-----+
       1 | 58   | 272  | +     | 1         | W1      | EBNA-LP    | LMP2
       2 | 368  | 458  | +     | 2         | W2      | EBNA-LP    | LMP2
       3 | 540  | 788  | +     | 3         | W1      | EBNA-LP    | LMP2
       4 | 871  | 951  | +     | 4         | W2      | EBNA-LP    | LMP2
       5 | 1026 | 1196 | +     | 5         | W1      | EBNA-LP    | LMP2
       6 | 1288 | 1495 | +     | 6         | W2      | EBNA-LP    | LMP2
       7 | 1574 | 1682 | +     | 7         | W1      | EBNA-LP    | LMP2
       8 | 5498 | 5956 | +     | 8         | W2      | EBNA-LP    | LMP2
       9 | 11336| 11480| +     | 9         | C1      | EBNA       | EBNA
      10 | 11626| 11657| +     | 10        | C2      | EBNA       | LMP2
(10 rows)

bdbio=>
```

Capture d'écran de psql

Pour les **outils graphiques**, il faut créer une nouvelle connexion avec les paramètres suivants :

- Serveur : `bd-pedago.univ-lyon1.fr`
- Base de données : `p1234567` (à remplacer par votre numéro étudiant)
- Utilisateur : `p1234567` (à remplacer par votre numéro étudiant)
- Mot de passe : voir suivi Tomuss BDBIO, case `mdp_BDBIO`

Pour l'**outil psql** en ligne de commande, il faut valider la commande suivante dans un terminal, puis saisir le mot de passe demandé (voir suivi Tomuss BDBIO, case `mdp_BDBIO`).

```
psql -h bd-pedago.univ-lyon1.fr -U p1234567
```

Il est également possible de travailler localement, en installant le SGBD PostgreSQL directement sur votre machine et en y important le jeu de données (voir section suivante).

<sup>1</sup>Plus de détails sur la [documentation du serveur bd-pedago](#).

## 2 Creation du jeu de donnees

Commencez par telecharger le [script PostgreSQL](#).

Le jeu de donnees represente les concepts de genes, exons et transcripts utilises en bio-informatique. Il provient a l'origine d'un fichier au format GenBank<sup>2</sup>, qui a ete converti en SQL.

Copiez-coller le script SQL et executez-le dans l'diteur SQL de votre outil, ou directement au niveau du *prompt* de psql<sup>3</sup>. Vous devez obtenir 152 genes, 39 exons et 95 transcripts.

Schema relationnel de la base de donnees :

```
GÈNES  (gene_id, start, stop, strand, gene, GeneID, locus_tag, gene_synonym)
EXONS  (exon_id, start, stop, strand, number, #gene, product, locus_tag)
TRANSCRIPTS (CDS_id, start, stop, strand, codon_start, protein_id, product, #gene, translation, UniProtKBSwissProt, InterPro, GOA, GI, PDB, GeneID, locus_tag, note, UniProtKBTrEMBL, function)
```

Questions sur la modelisation :

1. Quel probleme y a t-il dans cette base de donnees, en particulier au niveau des cles etrangeres ?
2. Cette base de donnees est-elle en troisime forme normale ?

## 3 Execution de requetes SQL

Pour ecrire des requetes SQL, vous devez lancer l'diteur SQL de l'outil, ou directement dans la console sur psql. Il est conseille d'ecrire les requetes dans un fichier texte (sauvegarde) et de les copier-coller dans l'outil.

Documentations : [SQL sur PostgreSQL](#), [DBeaver](#), [pgAdmin](#), [SQL.sh](#).

Traduire les requetes suivantes en SQL :

1. Nom des proteines de la table transcripts, que l'on obtient avec la requete `SELECT product FROM Transcripts`. Pourquoi obtient-on des doublons ? Modifiez la requete pour obtenir 84 tuples.
2. Informations sur les genes. Les attributs *start*, *stop* et *strand* seront concatenes et spares par des virgules au sein d'un seul attribut nomme *coords* (152 tuples resultat).
3. Informations sur les genes qui possedent un *locus\_tag*. Le resultat sera ordonne par nom de gene de croissant (45 tuples resultat)
4. Identifiant des transcripts qui codent pour une proteine dont l'existence est supposee, i.e. contenant *hypothetical*. L'identifiant (le *gene\_id* de GENES) et le nom du gene associe a ce transcript seront galement donnes (65 tuples resultat)
5. Identifiant des transcripts dont le *locus\_tag* est celui d'un gene (45 tuples resultat)
6. Nom des genes dont le *locus\_tag* contient un *P* ou qui sont associes a un exon codant pour une proteine (5 tuples resultat)
7. Nombre de transcripts appartenant au meme gene. Les resultats seront ordonnes par nombre de croissant (85 tuples resultat)
8. Noms des genes qui ont au moins deux *geneid* differents. On affichera aussi ce nombre de *geneid* (34 tuples resultat)
9. Longueurs minimale et maximale des translations, mais uniquement pour les transcripts dont le gene est galement associe a un exon (1 tuple resultat : 378 et 944)
10. Identifiant et *locus\_tag* des paires de transcripts qui possedent le meme complement / *strand* et dont les *locus\_tags* sont celui d'un gene commenant par *BALF* (1 tuple resultat)

<sup>2</sup>Genome EBV de GenBank, [http://www.ncbi.nlm.nih.gov/nuccore/NC\\_007605.1](http://www.ncbi.nlm.nih.gov/nuccore/NC_007605.1)

<sup>3</sup>Sur psql, ne pas oublier le point-virgule en fin de requete. Tapez \? pour la liste des commandes.