

# Récupération et intégration de lots de données depuis le GBIF

## Fonctionnement global

1. Configurer et déclencher un export ciblé via l'API du GBIF
2. Récupérer le fichier généré sur le serveur
3. Préparer les données et les stocker dans la base de données
4. Préparer les métadonnées associées
5. Intégrer ou actualiser les données dans la synthèse

## Scripts

### Configurer l'export GBIF

La récupération d'un grand volume de données pré-filtrées selon différentes variables est possible en configurant un export via l'API download du GBIF (ici en json). Le fichier de configuration est transmis à l'API du GBIF, qui met ensuite à disposition un export (fichiers) dans l'espace utilisateur du demandeur, correspondant à la requête configurée. Il est donc nécessaire d'avoir un compte GBIF actif.

Dans le fonctionnement du pôle invertébrés, les données filtrées concernent les données : \* disposant de coordonnées géographiques \* sans problème géospatial (géométries invalides etc) \* situées en Auvergne-Rhône-Alpes \* d'un ou certains jeux de données ciblés à l'avance \* des groupes (phylum) les mieux connus de la faune invertébrée

Pour connaître les "PHYLUM\_KEY" à filtrer, les correspondances avec taxref sont disponibles par défaut dans les bases de données GeoNature dans taxonomie.taxref\_liens. Par exemple pour les invertébrés :

```
SELECT DISTINCT t.cd_nom, tl.ct_sp_id
FROM taxonomie.taxref t
JOIN taxonomie.taxref_liens tl ON tl.cd_nom = t.cd_nom
WHERE tl.ct_name = 'GBIF'
AND t.id_rang = 'PH'
AND t.regne='Animalia'
AND t.phylum!='Chordata';
```

La configuration de l'export doit être stockée dans un fichier json. À titre d'exemple pour récupérer les Annélides, Mollusques et Arthropodes du jeu de données iNaturalist, le fichier query.json comporte la requête suivante :

```
{
  "creator": "USER",
  "notificationAddresses": [
```

```
"EMAIL"  
],  
"sendNotification": true,  
"format": "DWCA",  
"predicate": {  
  "type": "and",  
  "predicates": [  
    {  
      "type": "equals",  
      "key": "HAS_COORDINATE",  
      "value": "True"  
    },  
    {  
      "type": "equals",  
      "key": "HAS_GEOSPATIAL_ISSUE",  
      "value": "False"  
    },  
    {  
      "type": "equals",  
      "key": "GADM_GID",  
      "value": "FRA.1_1"  
    },  
    {  
      "type": "in",  
      "key": "DATASET_KEY",  
      "values": [ "50c9509d-22c7-4a22-a47d-8c48425ef4a7" ]  
    },  
    {  
      "type": "in",  
      "key": "PHYLUM_KEY",  
      "values": [ "42", "52", "54" ]  
    }  
  ]  
}  
}  
}
```

## Serveur : Récupérer, préparer et stocker les données en base

Une fois le fichier query.json configuré, la requête est transmise à l'API à l'aide de la commande suivante (installation du package jq pour une meilleure lisibilité des réponses) :

```
sudo apt install jq  
curl -s -silent --user USER_GBIF:PASS_GBIF --header "Content-Type: application/json" --data @query.json  
https://api.gbif.org/v1/occurrence/download/request | tail -n 1 | tr -d '\r'
```

**Cette commande retournera la clé de l'export** en question, sous une forme telle que :  
0015000-260108000000665

On peut ensuite suivre l'état de la requête (elle peut prendre quelques dizaines de minutes) avec la commande suivante :

```
curl -Ss https://api.gbif.org/v1/occurrence/download/<CLE> | jq -r '.status'
```

Une fois le statut "SUCCEEDED", se placer dans le répertoire /tmp et procéder au téléchargement des données :

```
cd /tmp
curl --location --remote-name
https://api.gbif.org/v1/occurrence/download/request/<CLE>.zip & unzip
<CLE>.zip
```

Une fois les données récupérées, les fichiers seront nettoyés des champs non utilisés pour optimiser les temps de traitement et l'espace occupé sur le serveur. Deux fichiers en particulier seront exploités : occurrence.txt et multimedia.txt

```
cut -f
1,6,8,18,25,27,30,31,32,33,34,35,40,48,63,90,98,99,100,135,149,180,188,194,2
06 /tmp/occurrence.txt > data_gbif_import.txt
cut -f 1,4 /tmp/multimedia.txt > multimedia_import.txt
```

Enfin, deux tables de destination sont créées dans la base de données (au premier usage), puis alimentées par ces données :

```
CREATE SCHEMA pinv_gbif;
```

```
CREATE TABLE pinv_gbif.tmp_gbif_data (
  gbif_id bigint primary key,
  modified timestamp,
  "references" varchar(255),
  basis_of_record varchar(255),
  recorde_by varchar(255),
  individual_count varchar(255),
  sex varchar(255),
  life_stage varchar(255),
  reproductive_condition varchar(255),
  caste varchar(255),
  behavior varchar(255),
  vitality varchar(255),
  occurrence_status varchar(255),
  occurrence_remarks text,
  event_date timestamp,
  verbatim_locality varchar(255),
  decimal_latitude numeric,
  decimal_longitude numeric,
  coordinate_uncertainty_meters varchar(255),
  identified_by varchar(255),
  scientific_name varchar(255),
  dataset_key uuid,
```

```
issue text,  
taxonKey int,  
verbatim_scientific_name varchar(255)  
);  
  
CREATE TABLE pinv_gbif.cor_multimedia (  
gbif_id int8 NULL,  
media_url text NULL  
);
```

## Postgresql : Insérer ou actualiser les données en synthèse

From:  
<http://wiki-sinp.cbn-alpin.fr/> - **CBNA SINP**

Permanent link:  
[http://wiki-sinp.cbn-alpin.fr/procedures/recuperation\\_et\\_integracion\\_de\\_donnees\\_depuis\\_le\\_gbif?rev=1769899142](http://wiki-sinp.cbn-alpin.fr/procedures/recuperation_et_integracion_de_donnees_depuis_le_gbif?rev=1769899142)

Last update: **2026/01/31 22:39**

