

CENTRALESUPÉLEC

INFRASTRUCTURES MODERNES

TP1 - Bases de Données Relationnelles et Graphes

Élève:

Daniel STULBERG HUF

Encadrants:

Francesca BUGIOTTI

Thierry RAPATOUT

Luc VO VAN

Table des matières

1	Base de données relationnelle - Azure SQL Database	1
2	Base de données graphe - Neo4j	6
3	Export des données vers un modèle graphe	6
4	Requêtes graphe (Cypher)	8
5	Requêtes graphe (Gremlin)	18
A	bashrc.sh	29
B	export-neo4j.py	29

1 Base de données relationnelle - Azure SQL Database

Exercice 0 : Décrivez les tables et les attributs.

Table	Column	Type	Null	Description
tArtist	idArtist	nvarchar	False	Primary key
tArtist	primaryName	nvarchar	False	Artist name
tArtist	birthYear	smallint	True	Artist birth year
tFilm	idFilm	nvarchar	False	Primary key
tFilm	primaryTitle	nvarchar	False	Movie name
tFilm	startYear	smallint	True	Release year
tFilm	runtimeMinutes	smallint	True	Movie duration in minutes
tFilmGenre	idFilm	nvarchar	False	Foreign key to tFilm
tFilmGenre	idGenre	nvarchar	True	Foreign key to tGenre
tGenre	idGenre	nvarchar	False	Primary key
tGenre	genre	nvarchar	True	Genre
tJob	idArtist	nvarchar	False	Foreign key to tArtist
tJob	category	nvarchar	True	Role in the movie
tJob	idFilm	nvarchar	True	Foreign key to tFilm

Exercice 1 ($\frac{1}{4}$ pt) : Visualisez l'année de naissance de l'artiste Nicolas Cage.

```

1 SELECT birthYear
2 FROM [dbo].[tNames]
3 WHERE primaryName = 'Nicolas Cage'

```

Query

birthYear
1964

Result

This SQL query retrieves the birth year of the artist named Nicolas Cage from the table '[dbo].[tNames]' with a simple select operation and a filter.

Exercice 2 ($\frac{1}{4}$ pt) : Comptez le nombre d'artistes présents dans la base de donnée.

```
1 SELECT COUNT(*) AS NumberOfArtists
2 FROM [dbo].[tArtist];
```

Query

NumberOfArtists
5227

Result

This SQL query counts the total number of artists in the table '[dbo].[tArtist]' by simply counting the total number of rows of such table.

Exercice 3 ($\frac{1}{4}$ pt) : Trouvez les noms des artistes nés en 1960, affichez ensuite leur nombre.

```
1 SELECT primaryName
2 FROM [dbo].[tArtist]
3 WHERE birthYear = 1960;
```

Query 1

primaryName
Antonio Banderas
John Leguizamo
Oliver Platt
...

Result 1

This SQL query retrieves the names of artists who were born in the year 1960 from the table '[dbo].[tArtist]' with a simple select operation and a filter.

```
1 SELECT COUNT(*) AS NumberOfArtistsBornIn1960
2 FROM [dbo].[tArtist]
3 WHERE birthYear = 1960;
```

Query 2

NumberOfArtistsBornIn1960
20

Result 2

This SQL query counts the number of artists born in the year 1960 in the table '[dbo].[tArtist]' by simply counting the total number of rows of such table.

Exercice 4 (1 pt) : Trouvez l'année de naissance la plus représentée parmi les acteurs (sauf 0 !), et combien d'acteurs sont nés cette année là.

```
1 SELECT TOP 1 birthYear, COUNT(*) AS NumberOfActors
2 FROM [dbo].[tArtist]
3 WHERE birthYear <> 0
4 GROUP BY birthYear
5 ORDER BY COUNT(*) DESC;
```

Query

birthYear	NumberOfActors
1979	40

Result

This SQL query finds the most common birth year (excluding the year 0) among artists in the table '[dbo].[tArtist]' and shows how many artists were born in that year. The query groups the artists from the table by birth year, counts how many artists are in each group, then sorts these groups so the one with the most artists is first, and finally, it selects the top group.

Exercice 5 ($\frac{1}{2}$ pt) : Trouvez les artistes ayant joué dans plus d'un film

```
1 SELECT a.idArtist, a.primaryName, COUNT(j.idFilm) AS NumberOfFilms
2 FROM [dbo].[tArtist] AS a
3 JOIN [dbo].[tJob] AS j ON a.idArtist = j.idArtist
4 WHERE j.category = 'acted in'
5 GROUP BY a.idArtist, a.primaryName
6 HAVING COUNT(j.idFilm) > 1
7 ORDER BY NumberOfFilms DESC;
```

Query

idArtist	primaryName	NumberOfFilms
nm5401689	Miles Jonn-Dalton	7
nm0562210	Antonio Mayans	6
nm12302814	Ellen Wing	6
...

Result

This SQL query lists artists who have acted in more than one film, along with the number of films they acted in. The query joins the tables '[dbo].[tArtist]' and '[dbo].[tJob]' based on the artist's ID, filters to include only those jobs where the artist's role was 'acted in', groups the result by artist ID and name, counts the number of films for each artist, includes only those artists who have acted in more than one film, and finally, orders the results by the number of films in descending order.

Exercice 6 ($\frac{1}{2}$ pt) : Trouvez les artistes ayant eu plusieurs responsabilités au cours de leur carrière (acteur, directeur, producteur...).

```
1 SELECT a.idArtist, a.primaryName, COUNT(DISTINCT j.category) AS
   NumberOfDistinctCategories
2 FROM [dbo].[tArtist] AS a
3 JOIN [dbo].[tJob] AS j ON a.idArtist = j.idArtist
4 GROUP BY a.idArtist, a.primaryName
5 HAVING COUNT(DISTINCT j.category) > 1;
```

Query

idArtist	primaryName	NumberOfDistinctRoles
nm0083851	Sam Bisbee	2
nm0164703	Christopher Clarke	2
nm0168002	Glen Coburn	2
...

Result

This SQL query lists artists who have had more than one distinct type of job in their career, showing each artist's ID, name, and the number of different job types they've had. The query joins the tables '[dbo].[tArtist]' and '[dbo].[tJob]' on the artist's ID, groups the result by artist ID and name, counts the distinct types of job categories for each artist, and finally, includes only those artists who have more than one distinct type of job.

Exercice 7 ($\frac{3}{4}$ pt) : Trouver le nom du ou des film(s) ayant le plus d'acteurs (i.e. uniquement acted in).

```
1 SELECT f.primaryTitle, COUNT(*) AS NumberOfActors
2 FROM tFilm AS f
3 JOIN tJob AS j ON f.idFilm = j.idFilm
4 WHERE j.category = 'acted in'
5 GROUP BY f.primaryTitle
6 ORDER BY COUNT(*) DESC;
```

Query

primaryTitle	NumberOfActors
Carole & Grey	10
Full Clip for a Snitch	10
Noches de Sol	10
...	...

Result

This SQL query lists films along with the number of actors in each, ordering first by the greatest number of actors who participated in it. The query joins the tables '[dbo].[tFilm]' and '[dbo].[tJob]' based on the film ID, filters to include only those jobs where the category is 'acted in', groups the result by the film title, counts the number of actors for each film, and finally, orders the films by the number of actors in descending order.

Exercice 8 (1 pt) : Montrez les artistes ayant eu plusieurs responsabilités dans un même film (ex : à la fois acteur et directeur, ou toute autre combinaison) et les titres de ces films.

```
1 SELECT a.primaryName, f.primaryTitle
2 FROM tArtist AS a
3 JOIN tJob AS j ON a.idArtist = j.idArtist
4 JOIN tFilm AS f ON j.idFilm = f.idFilm
5 GROUP BY a.idArtist, a.primaryName, f.idFilm, f.primaryTitle
6 HAVING COUNT(DISTINCT j.category) > 1;
```

Query

primaryName	primaryTitle
Jeffrey Schneider	The Accident
Lenni Uitto	Permafrost
Jeff Proffitt	Street Crimes Unit
...	...

Result

This SQL query lists the names of artists who had more than one type of job in the same film, along with the titles of those films. The query joins the tables '[dbo].[tArtist]', '[dbo].[tJob]', and '[dbo].[tFilm]' on artist ID and film ID, groups the results by both artist and film, counts the distinct job categories for each artist in each film, and finally, includes only the pairs of artist-film where the artist had more than one type of job in the same film.

2 Base de données graphe - Neo4j

Modifiez la fin du fichier `/.bashrc` les variables d'environnement nécessaires à la connexion à votre base Neo4.

The fragment of the `bashrc` file containing the necessary environment variables is attached to this delivery.

3 Export des données vers un modèle graphe

Dans le dossier `tp`, complétez le programme `export-neo4j.py` aux endroits notés **A COMPLETER. N'hésitez pas à déboguer en ajoutant des `print`, créer des programmes de test etc. Utilisez les fonctions `create_nodes` et `create_relationships` de `py2neo`.**

The completed `export-neo4j.py` script is attached to this delivery.

(optionnel, 2 points bonus) Décrivez de quelle manière l'environnement de développement a été préconfiguré pour vous :

- **Environnement Python**
- **Informations de connexion aux bases de données**

The provided script is designed to migrate data from a SQL database to a Neo4j graph database. It demonstrates how to preconfigure a development environment for such a task, which includes setting up the Python environment and connecting to both SQL and Neo4j databases.

Python Environment

This script uses `py2neo`, which is a client toolkit for working with Neo4j from Python applications. It provides functionalities to interact with the graph database, such as creating nodes and relationships, executing Cypher queries, and handling transactions. For connecting to the SQL database, the script relies on `pyodbc`, a Python open database connectivity interface that allows connecting to any database using the ODBC drivers installed on the system. Finally, the script relies on environment variables to securely manage the database connection details, which avoids hardcoding sensitive data (like usernames and passwords) directly in the script.

Database Connection

A connection string is constructed with the details provided by the environment variables to connect to the SQL database and perform queries to fetch data that will be migrated to the Neo4j graph database. The `Graph` object from `py2neo` is instantiated with the environment variables, enabling the script to interact with the Neo4j database, in order to create nodes, relationships, and run Cypher commands for data manipulation and schema updates.

Workflow

The script starts by deleting existing nodes and relationships in the Neo4j database to prepare for fresh data import. The script then fetches data from the tables `tFilm`, `tArtist`, and `tJob` in the SQL database, and creates the corresponding `Film` and `Artist`

nodes in Neo4j, along with their properties. The relationships (ACTED_IN, DIRECTED, etc) between Artist and Film nodes are created based on the roles defined in the tJob table. In addition, the indexes on idFilm for Film nodes and idArtist for Artist nodes are created to optimize query performance in Neo4j.

Data Migration

The script processes data in batches (that are defined by BATCH_SIZE) to efficiently handle large datasets. For each entity type (Film, Artist), the script fetches a batch of records from the SQL database, converts each record into a Node object with appropriate labels and properties, and uses 'create_nodes' to bulk-import these nodes into Neo4j. For relationships, the script constructs tuples representing the relationships from the tJob table and uses 'create_relationships' to bulk-import these relationships into Neo4j, therefore linking artists and films.

4 Requetes graphe (Cypher)

Exercice 1 ($\frac{1}{4}$ pt) : Ajoutez une personne ayant votre prénom et votre nom dans le graphe. Verifiez que le noeud a bien été créé.

```
1 CREATE (n: Artist { primaryName: "Daniel Stulberg Huf", birthYear:
   ↳ 2002 })
2 RETURN n
```

Query

```
1 {
2   "identity": 6473,
3   "labels": [
4     "Artist"
5   ],
6   "properties": {
7     "birthYear": 2002,
8     "primaryName": "Daniel Stulberg Huf"
9   },
10  "elementId": "6473"
11 }
```

Result

This query creates a new node in the graph database with the label Artist and sets the properties primaryName with the value "Daniel Stulberg Huf" and birthYear with the value 2002. Then, it returns the newly created node.

Exercice 2 ($\frac{1}{4}$ pt) : Ajoutez un film nommé L’histoire de mon 20 au cours Infrastructure de donnees

```
1 CREATE (n: Film { primaryTitle: "L’histoire de mon 20 au cours
   ↳ Infrastructure de donnees", startYear: 2024 })
2 RETURN n
```

Query

```
1 {
2   "identity": 6474,
3   "labels": [
4     "Film"
5   ],
6   "properties": {
7     "primaryTitle": "L’histoire de mon 20 au cours Infrastructure de
   ↳ donnees",
```

```
8     "startYear": 2024
9   },
10   "elementId": "6474"
11 }
```

Result

This query creates a new node in the graph database with the label Film and sets the properties primaryTitle with the value "L'histoire de mon 20 au cours Infrastructure de donnees" and startYear with the value 2024. Then, it returns the newly created node.

Exercice 3 ($\frac{1}{2}$ pt) : Ajoutez la relation ACTED_IN qui modélise votre participation à ce film en tant qu'acteur/actrice

```
1 MATCH (a:Artist), (b:Film)
2 WHERE a.primaryName = "Daniel Stulberg Huf" AND b.primaryTitle =
   ↳ "L'histoire de mon 20 au cours Infrastructure de donnees"
3 CREATE (a)-[r:ACTED_IN]->(b)
4 RETURN type(r)
```

Query

```
1 type(r)
2 "ACTED_IN"
```

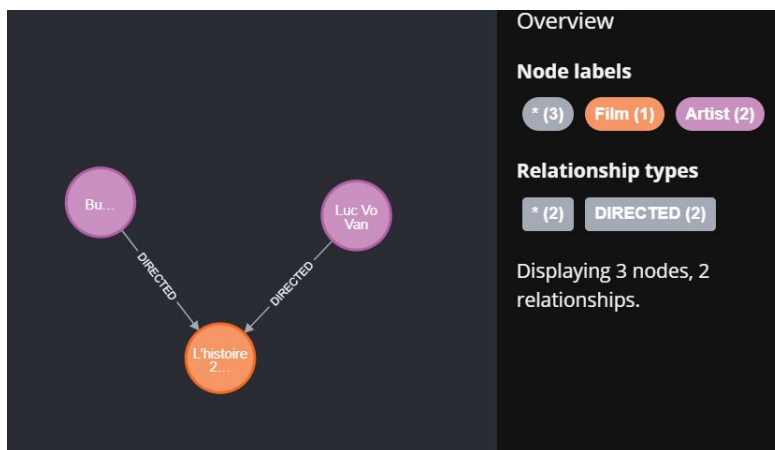
Result

This query first finds two existing nodes in the graph, one with the label Artist and a primaryName property of "Daniel Stulberg Huf", and another with the label Film and a primaryTitle property of "L'histoire de mon 20 au cours Infrastructure de donnees". It then creates a relationship of type ACTED_IN from the Artist node to the Film node. Finally, it returns the type of the created relationship, which in this case will be "ACTED_IN".

Exercice 4 ($\frac{1}{2}$ pt) : Ajoutez deux de vos professeurs/enseignants comme réalisateurs/réalisatrices de ce film.

```
1 MATCH (a: Film)
2 WHERE a.primaryTitle = "L'histoire de mon 20 au cours Infrastructure
   ↳ de donnees"
3 CREATE (b: Artist { primaryName: "Luc Vo Van" })-[:DIRECTED]->(a)
4 CREATE (c: Artist { primaryName : "Francesca Bugiotti"
   ↳ })-[:DIRECTED]->(a)
5 RETURN a, b, c
```

Query



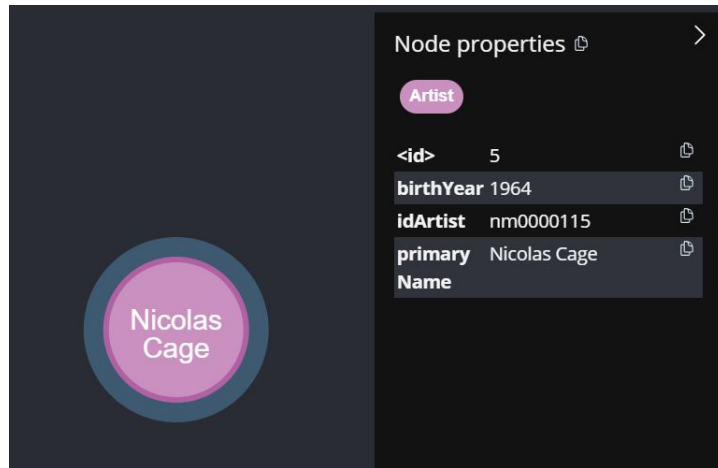
Result

This query adds two new artists as directors to an existing film, creating both the artist nodes and their directing relationships to the film in a single operation. It first searches for a node labeled as Film with the primaryTitle property equal to "L'histoire de mon 20 au cours Infrastructure de donnees", then it creates a new Artist node with the primaryName property set to "Luc Vo Van" and establishes a DIRECTED relationship from this new artist node to the previously matched film node. The same relation is established for the new Artist node with the primaryName property set to "Francesca Bugiotti". Finally, the query returns the film node 'a' and the two newly created artist nodes 'b' and 'c'.

Exercice 5 ($\frac{1}{2}$ pt) : Affichez le noeud représentant l'acteur nommé Nicolas Cage, et visualisez son année de naissance.

```
1 MATCH (a: Artist)
2 WHERE a.primaryName = "Nicolas Cage"
3 RETURN a
```

Query



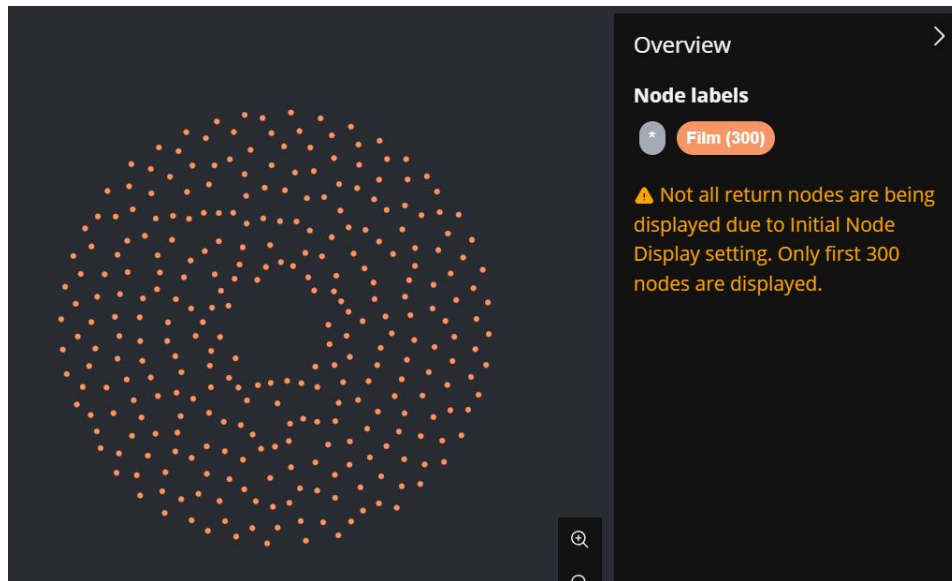
Result

This query searches for a node with the label Artist that has a primaryName property equal to "Nicolas Cage" and returns that node. Essentially, it finds and returns the node representing the artist named Nicolas Cage from the graph database.

Exercice 6 ($\frac{1}{2}$ pt) : Visualisez l'ensemble des films.

```
1 MATCH (a: Film)
2 RETURN a
```

Query



Result

This query searches for all nodes in the graph database that have the label Film and returns those nodes, therefore retrieving every node representing a film.

Exercice 7 ($\frac{1}{2}$ pt) : Trouvez les noms des artistes nés en 1960, affichez ensuite leur nombre.

Query 1 :

```
1 MATCH (a: Artist)
2 WHERE a.birthYear = 1960
3 RETURN a
```

Query 1

```
1 {
2   "identity": 0,
3   "labels": [
4     "Artist"
5   ],
6   "properties": {
7     "idArtist": "nm0000104",
8     "birthYear": 1960,
9     "primaryName": "Antonio Banderas"
10  },
11  "elementId": "0"
12 }
13 {
14   "identity": 23,
15   "labels": [
16     "Artist"
```

```

17 ],
18 "properties": {
19   "idArtist": "nm0000491",
20   "birthYear": 1960,
21   "primaryName": "John Leguizamo"
22 },
23 "elementId": "23"
24 }
25 {
26   "identity": 55,
27   "labels": [
28     "Artist"
29 ],
30   "properties": {
31     "idArtist": "nm0001624",
32     "birthYear": 1960,
33     "primaryName": "Oliver Platt"
34   },
35   "elementId": "55"
36 }
37 ...

```

Result 1

This query searches for all nodes in the graph database that have the label Artist and a birthYear property equal to 1960, and then returns those nodes, therefore returning all artist nodes representing artists born in the year 1960.

Query 2 :

```

1 MATCH (a: Artist)
2 WHERE a.birthYear = 1960
3 RETURN COUNT(a)

```

Query 2

```

1 COUNT(a)
2 20

```

Result 2

This query searches for all nodes in the graph database that have the label Artist with a birthYear property equal to 1960. Instead of returning the nodes themselves, it returns the count of these nodes, which represents the number of artists born in the year 1960.

Exercice 8 (1 pt) : Trouver l'ensemble des acteurs (sans entrées doublons) qui ont joué dans plus d'un film.

```
1 MATCH (n: Artist)
2 WHERE SIZE([(n)-[:ACTED_IN]->(f) | f]) > 1
3 RETURN n
```

Query

```
1 {
2   "identity": 6,
3   "labels": [
4     "Artist"
5   ],
6   "properties": {
7     "idArtist": "nm0000128",
8     "birthYear": 1964,
9     "primaryName": "Russell Crowe"
10  },
11  "elementId": "6"
12 }
13 {
14   "identity": 12,
15   "labels": [
16     "Artist"
17   ],
18   "properties": {
19     "idArtist": "nm0000199",
20     "birthYear": 1940,
21     "primaryName": "Al Pacino"
22  },
23  "elementId": "12"
24 }
25 {
26   "identity": 27,
27   "labels": [
28     "Artist"
29   ],
30   "properties": {
31     "idArtist": "nm0000616",
32     "birthYear": 1956,
33     "primaryName": "Eric Roberts"
34  },
35   "elementId": "27"
36 }
37 ...
```


Result

This query searches for artist nodes (:Artist) who have acted in more than one film. It uses a pattern comprehension inside the SIZE function to count the distinct films (f) each artist (n) has acted in. If an artist has acted in more than one film, the artist node is returned.

Exercice 9 (1 pt) : Trouvez les artistes ayant eu plusieurs responsabilités au cours de leur carrière (acteur, directeur, producteur...).

```
1 MATCH (n: Artist)-[r]->(:Film)
2 WITH COUNT(DISTINCT type(r)) as nb_roles, n
3 WHERE nb_roles > 1
4 RETURN n
```

Query

```
1 {
2   "identity": 1960,
3   "labels": [
4     "Artist"
5   ],
6   "properties": {
7     "idArtist": "nm7700848",
8     "birthYear": 0,
9     "primaryName": "Maja Prelog"
10  },
11  "elementId": "1960"
12 }
13 {
14   "identity": 717,
15   "labels": [
16     "Artist"
17   ],
18   "properties": {
19     "idArtist": "nm6366616",
20     "birthYear": 0,
21     "primaryName": "Jeffrey Tipton"
22  },
23  "elementId": "717"
24 }
25 {
26   "identity": 1116,
27   "labels": [
28     "Artist"
29   ],
30   "properties": {
```

```

31     "idArtist": "nm6815189",
32     "birthYear": 0,
33     "primaryName": "Martin W. Payne"
34 },
35     "elementId": "1116"
36 }
37 ...

```

Result

This query finds artists (:Artist nodes) who have more than one type of relationship (r) with any film (:Film nodes), indicating they have had multiple roles in relation to films. The query first matches artists that have any kind of relationship to films, then by using the WITH clause, it counts the distinct types of relationships (type(r)) each artist has with films, and aliases this count as nb_roles. After that, the query filters these results to include only those artists where nb_roles is greater than 1. Finally, it returns the artist nodes that meet this criterion.

Exercice 10 (1 pt) : Montrez les artistes ayant eu plusieurs responsabilités dans un même film (ex : à la fois acteur et directeur, ou toute autre combinaison) et les titres de ces films.

```

1  MATCH (a:Artist)-[r]->(f:Film)
2  WITH a, f, COLLECT(DISTINCT type(r)) AS roles
3  WHERE SIZE(roles) > 1
4  RETURN a.primaryName AS Artist, f.primaryTitle AS Film, roles

```

Query

Artist	Film	roles
Maja Prelog	Cent'anni	["DIRECTED", "ACTED_IN"]
Jeffrey Tipton	Father Fauci : The Covid Story	["DIRECTED", "ACTED_IN"]
Ben Petrie	The Heirloom	["DIRECTED", "ACTED_IN"]
...

Result

This query identifies artists who have had multiple distinct roles in the same film, then returns the names of these artists, the titles of the films, and the distinct roles they held. The query first finds patterns where an artist (a) has any type of relationship (r) with a film (f), then for each artist-film pair, it collects the distinct types of relationships into a list named roles. The query filters these pairs to only include those where the artist has more than one distinct role in the same film (SIZE(roles) > 1), and finally, returns the primary name of the artist, the primary title of the film, and the list of distinct roles that the artist had in that film.

Exercice 11 (2 pt) : Trouver le nom du ou des film(s) ayant le plus d'acteurs.

```
1 MATCH (n: Artist)-[:ACTED_IN]->(t: Film)
2 WITH COUNT(DISTINCT n) as NumberOfActors, t
3 ORDER BY NumberOfActors DESC
4 LIMIT 3
5 RETURN t.primaryTitle, NumberOfActors
```

Query

primaryTitle	NumberOfActors
Full Clip for a Snitch	10
Noches de Sol	10
Carole & Grey	10
...	...

Result

This query finds films with the most distinct actors, limiting the results to the top 3 films with the highest number of actors for simplicity reasons, and returns the titles of these films along with the count of distinct actors in each. The query first matches patterns where an artist (n with the label Artist) acted in (:ACTED_IN relationship) a film (t with the label Film), then for each film t, it counts the number of distinct actors n that have acted in it, grouping by the film. It orders the films by the count of distinct actors in descending order (DESC) and limits the results to the top 3 films. Finally, the query returns the primary title of each of these films and the count of distinct actors that acted in them as NumberOfActors.

5 Requêtes graphe (Gremlin)

Exercice 1 ($\frac{1}{4}$ pt) : Ajoutez une personne ayant votre prénom et votre nom dans le graphe. Vérifiez que le noeud a bien été créé.

```
1 g.addV("Artist")
2   .property("primaryName", "Daniel Stulberg Huf")
3   .property("birthYear", 2002)
4   .property("pk", "pk_dh")
```

Query

```
1 [
2   {
3     "id": "cbb9c809-10f8-4718-a250-013826b84150",
4     "label": "Artist",
5     "type": "vertex",
6     "properties": {
7       "primaryName": [
8         {
9           "id": "79e2e5ce-6a16-44fa-9d55-dc4f05f4eb0d",
10          "value": "Daniel Stulberg Huf"
11        }
12      ],
13      "birthYear": [
14        {
15          "id": "404b2e2b-b8a4-4dfe-966a-d4203684652d",
16          "value": 2002
17        }
18      ],
19      "pk": [
20        {
21          "id": "cbb9c809-10f8-4718-a250-013826b84150|pk",
22          "value": "pk_dh"
23        }
24      ]
25    }
26  }
27 ]
```

Result

This query creates a new vertex in the graph database labeled "Artist" and assigns it three properties : 'primaryName' set to "Daniel Stulberg Huf"; 'birthYear' set to 2002; and 'pk' set to "pk_dh," which serves as a primary key for this vertex. The query initiates the addition of a vertex with `g.addV("Artist")`, followed by chaining `.property()` steps to assign each of the properties.

Exercice 2 ($\frac{1}{4}$ pt) : Ajoutez un film nommé L’histoire de mon 20 au cours Infrastructure de donnees

```
1 g.addV("Film")
2   .property("primaryTitle", " L histoire  de mon 20 au cours
3     Infrastructure de donnees")
4   .property("startYear", 2024)
5   .property("pk", "pk_film")
```

Query

```
1 [
2   {
3     "id": "69e4c697-3191-4957-962e-62abdf9661ec",
4     "label": "Film",
5     "type": "vertex",
6     "properties": {
7       "primaryTitle": [
8         {
9           "id": "67594b6f-a3cb-4029-9df6-b0f75f61efe8",
10          "value": " L histoire  de mon 20 au cours Infrastructure
11          de donnees"
12        }
13      ],
14      "startYear": [
15        {
16          "id": "74496039-3ca2-41fc-b045-a520d92bcd8a",
17          "value": 2024
18        }
19      ],
20      "pk": [
21        {
22          "id": "69e4c697-3191-4957-962e-62abdf9661ec|pk",
23          "value": "pk_film"
24        }
25      ]
26    }
27 ]
```

Result

This query creates a new vertex in the graph database labeled "Film" and assigns it three properties : ‘primaryTitle’ set to "L’histoire de mon 20 au cours Infrastructure de donnees"; ‘startYear’ set to 2024; and ‘pk’ set to "pk_film," which serves as a primary key for this vertex. The query initiates the addition of a vertex with `g.addV("Film")`, followed by chaining `.property()` steps to assign each of the properties.

Exercice 3 ($\frac{1}{4}$ pt) : Ajoutez la relation **ACTED_IN** qui modélise votre participation à ce film en tant qu'acteur/actrice

```

1 g.V().has("Artist", "primaryName", "Daniel Stulberg Huf").as("artist")
2   .V().has("Film", "primaryTitle", " L histoire de mon 20 au cours
   Infrastructure de donnees").as("film")
3   .addE("ACTED_IN")
4   .from("artist")
5   .to("film")

```

Query

```

1 [
2   {
3     "id": "7eada684-32c0-4dfd-b428-a5815389ef1d",
4     "label": "ACTED_IN",
5     "type": "edge",
6     "inVLabel": "Film",
7     "outVLabel": "Artist",
8     "inV": "69e4c697-3191-4957-962e-62abdf9661ec",
9     "outV": "cbb9c809-10f8-4718-a250-013826b84150"
10  }
11 ]

```

Result

This query creates a relationship by linking the artist named "Daniel Stulberg Huf" to the film titled "L'histoire de mon 20 au cours Infrastructure de donnees" with an edge labeled "ACTED_IN". It begins by searching for a vertex labeled "Artist" with the primaryName property matching "Daniel Stulberg Huf" and aliases this vertex as "artist." It simultaneously searches for another vertex labeled "Film" with the primaryTitle property of "L'histoire de mon 20 au cours Infrastructure de donnees" and aliases it as "film." Then, the query creates an "ACTED_IN" edge from the artist vertex to the film vertex.

Exercice 4 ($\frac{1}{4}$ pt) : Ajoutez deux de vos professeurs/enseignants comme réalisateurs/réalisatrices de ce film.

```

1 g.V().as("film").has("Film", "primaryTitle", " L histoire de mon 20
   au cours Infrastructure de donnees")
2   .addV("Artist")
3   .as("director1")
4   .property("primaryName", "Luc Vo Van")
5   .property("pk", "pk_lv")
6   .addE("DIRECTED").from("director1").to("film")
7   .addV("Artist")
8   .as("director2")
9   .property("primaryName", "Francesca Bugiotti")
10  .property("pk", "pk_fb")
11  .addE("DIRECTED").from("director2").to("film")

```

```
12 .select("film", "director1", "director2");
```

Query

```
1 [
2   {
3     "film": {
4       "id": "69e4c697-3191-4957-962e-62abdf9661ec",
5       "label": "Film",
6       "type": "vertex",
7       "properties": {
8         "primaryTitle": [
9           {
10            "id": "67594b6f-a3cb-4029-9df6-b0f75f61efe8",
11            "value": " L histoire de mon 20 au cours Infrastructure
12            de donnees"
13          }
14        ],
15        "startYear": [
16          {
17            "id": "74496039-3ca2-41fc-b045-a520d92bcd8a",
18            "value": 2024
19          }
20        ],
21        "pk": [
22          {
23            "id": "69e4c697-3191-4957-962e-62abdf9661ec|pk",
24            "value": "pk_film"
25          }
26        ]
27      },
28      "director1": {
29        "id": "3c2fe96a-e056-4307-b428-b683c84af627",
30        "label": "Artist",
31        "type": "vertex",
32        "properties": {
33          "primaryName": [
34            {
35              "id": "bc0d089c-5ccb-4ff0-8752-207fa91a1f89",
36              "value": "Luc Vo Van"
37            }
38          ],
39          "pk": [
40            {
41              "id": "3c2fe96a-e056-4307-b428-b683c84af627|pk",
42              "value": "pk_lv"
43            }
44          ]
45        }
46      },
47      "director2": {
```

```

48     "id": "be1087bd-f45b-401a-a6cb-0d5410b74098",
49     "label": "Artist",
50     "type": "vertex",
51     "properties": {
52         "primaryName": [
53             {
54                 "id": "51b23b9d-c4e8-4814-81a2-2221fbce69e4",
55                 "value": "Francesca Bugiotti"
56             }
57         ],
58         "pk": [
59             {
60                 "id": "be1087bd-f45b-401a-a6cb-0d5410b74098|pk",
61                 "value": "pk_fb"
62             }
63         ]
64     }
65 }
66 }
67 ]

```

Result

This query begins by identifying a film vertex based on its "primaryTitle" property being "L'histoire de mon 20 au cours Infrastructure de donnees" and aliases it as "film". It then adds two new vertices labeled "Artist", representing directors Luc Vo Van and Francesca Bugiotti, assigning each properties for their names and unique identifiers, and aliases them as "director1" and "director2" respectively. For each director, the query creates a "DIRECTED" edge from the director to the film, signifying their roles as directors in that film. Finally, the query selects and returns the film and the two director vertices.

Exercice 5 ($\frac{1}{4}$ pt) : Affichez le noeud représentant l'acteur nommé Nicolas Cage, et visualisez son année de naissance.

```

1 g.V().has("Artist", "primaryName", "Nicolas Cage")
2   .project("name", "year")
3     .by("primaryName")
4     .by("birthYear");

```

Query

```

1 [
2   {
3     "name": "Nicolas Cage",
4     "year": 1960
5   }
6 ]

```

Result

This query searches the graph for vertices labeled "Artist" that have a primaryName property equal to "Nicolas Cage." For each vertex that matches this criterion, it constructs a map projection with two keys : "name" and "year." The .by("primaryName") specifies that the value associated with the "name" key in the output map should be the value of the primaryName property of the vertex. Similarly, the .by("birthYear") specifies that the value associated with the "year" key should be the value of the birthYear property of the vertex. This results in the matched vertex detailing Nicolas Cage's name and birth year.

Exercice 6 ($\frac{1}{4}$ pt) : Visualisez l'ensemble des films.

```
1 g.V().hasLabel("Film")
```

Query

```
1  [
2  {
3    "id": "tt15125820",
4    "label": "Film",
5    "type": "vertex",
6    "properties": {
7      "primaryTitle": [
8        {
9          "id": "af2afe0a-bdb9-40fe-9c80-98084edaa392",
10         "value": "Bionica"
11       }
12     ],
13     "startYear": [
14       {
15         "id": "faf1de08-91b0-4b6c-8b0c-34695cfab504",
16         "value": 2024
17       }
18     ],
19     "averageRating": [
20       {
21         "id": "59b18de1-2ad2-484b-8031-4c194579d91c",
22         "value": 0
23       }
24     ],
25     "pk": [
26       {
27         "id": "tt15125820|pk",
28         "value": "0"
29       }
30     ]
31   }
32 },
33 {
34   "id": "tt30874320",
35   "label": "Film",
36   "type": "vertex",
```

```

37  "properties": {
38    "primaryTitle": [
39      {
40        "id": "e257af60-15f0-4785-a0a4-1176fe13bbe0",
41        "value": "Cent'anni"
42      }
43    ],
44    "startYear": [
45      {
46        "id": "baef7fec-6b6c-4c20-8e09-ff03700abc63",
47        "value": 2024
48      }
49    ],
50    "averageRating": [
51      {
52        "id": "988a5ea5-11d9-457a-a9b2-90e7649e0b33",
53        "value": 0
54      }
55    ],
56    "pk": [
57      {
58        "id": "tt30874320|pk",
59        "value": "0"
60      }
61    ]
62  }
63 },
64 {
65   "id": "tt30877290",
66   "label": "Film",
67   "type": "vertex",
68   "properties": {
69     "primaryTitle": [
70       {
71         "id": "cdb6d0e6-b66f-40c2-bde7-d5e5d41e8f17",
72         "value": "Velky Pan"
73       }
74     ],
75     "startYear": [
76       {
77         "id": "670cd0aa-ddd2-4d83-8909-9f2e7d573b98",
78         "value": 2024
79       }
80     ],
81     "averageRating": [
82       {
83         "id": "076b3338-fa8c-4578-9a45-f62c5887596a",
84         "value": 0
85       }
86     ],
87     "pk": [
88       {

```

```
89         "id": "tt30877290|pk",
90         "value": "0"
91     }
92 ]
93 }
94 },
95 ...
96 ]
```

Result

This query searches through the graph for all vertices that have the label "Film". It starts by looking at all vertices in the graph and then filters these vertices to only include those that are labeled with "Film" using the `hasLabel("Film")` step.

Exercice 7 ($\frac{1}{2}$ pt) : Trouvez les noms des artistes nés en 1960, affichez ensuite leur nombre.

```
1 g.V().hasLabel("Artist")
2 .has("birthYear", 1960)
3 .values("primaryName");
```

Query 1

```
1 [
2   "Joan Chen",
3   "Desmond Eastwood",
4   "Alicia Witt",
5   ...
6 ]
```

Result 1

This query searches the graph for vertices labeled as "Artist" who were born in the year 1960, and then extracts the `primaryName` property from each of these vertices. It begins by selecting all vertices and narrows down this selection to those with the label "Artist" using `.hasLabel("Artist")`. It further filters these artist vertices to only those with a `birthYear` property equal to 1960 through `.has("birthYear", 1960)`. Finally, for each vertex that matches these criteria, it retrieves the value of the `primaryName` property using `.values("primaryName")`, listing the names of all artists born in 1960.

```
1 g.V().hasLabel("Artist").has("birthYear", 1960).count()
```

Query 2

```
1 [
2   5289
3 ]
```

Result 2

This query counts the number of vertices in the graph that are labeled as "Artist" and have a birthYear property equal to 1960. It starts by selecting all vertices in the graph, then filters these to only include those with the label "Artist" (`hasLabel("Artist")`). It then narrows down the selection to those artists who were born in 1960 (`has("birthYear", 1960)`). Finally, it applies the `count()` method, which counts the number of vertices that meet these criteria, resulting in the total number of artists born in that year.

Exercice 8 (1,5 pt) : Trouver l'ensemble des acteurs (sans entrées doublons) qui ont joué dans plus d'un film.

```
1 g.V().hasLabel("Artist")
2   .where(__.out("acted in").count().is(gt(1)))
3   .dedup()
4   .values("primaryName")
```

Query

```
1 [
2   "Tim Neff",
3   "Lilly Singh",
4   "Olwen Fouere",
5   ...
6 ]
```

Result

This query finds and lists the unique names of artists who have acted in more than one film. It starts by selecting all vertices labeled "Artist" from the graph. Then, for each artist, it checks if there are more than one outgoing edges labeled "acted in," indicating the artist has acted in multiple films. This is done using a where step that contains a sub-traversal : it looks at outgoing "acted in" edges, counts them, and filters for artists with a count greater than one using `.count().is(gt(1))`. To ensure each artist's name is listed only once, even if they've acted in many films, the `dedup()` step removes any duplicates. Finally, it extracts and returns the `primaryName` property of these filtered artist vertices, giving a list of artists who have acted in more than one film.

Exercice 9 (2 pt) : Trouvez les artistes ayant eu plusieurs responsabilités au cours de leur carrière (acteur, directeur, producteur...).

```
1 g.V().hasLabel('Artist')
2   .where(outE().groupCount().by(label).unfold().
3   select(keys).dedup().count().is(gt(1)))
4   .values('primaryName')
```

Query

```
1 [
2   "Carley Harrison",
3   "Josef Hader",
4   "Pavel Nikolajev",
```

```
5   ...
6 ]
```

Result

This query identifies and retrieves the names of artists who have more than one type of relationship to other entities in the graph, indicating they've held multiple roles. It starts by selecting vertices labeled "Artist". For each artist, it examines their outgoing relationships, counting the occurrences of each relationship type using `group-Count().by(label)`. This count is then unfolded to work with individual entries, from which it extracts the relationship types (`select(keys)`), removes duplicates with `dedup()`, and counts the distinct types of relationships. If an artist has more than one distinct type of outgoing relationship (`count().is(gt(1))`), it implies they've had multiple roles. Finally, it extracts and returns the `primaryName` property of these artists, providing a list of artists with multiple roles in their careers.

Exercice 10 (2 pt) : Montrez les artistes ayant eu plusieurs responsabilités dans un même film (ex : à la fois acteur et directeur, ou toute autre combinaison) et les titres de ces films.

```
1 g.V().hasLabel('Film').has("primaryTitle")
2   .as('film')
3   .inE().as('roles')
4   .outV().hasLabel('Artist').has("primaryName")
5   .as('artist')
6   .select('artist', 'film')
7     .by('primaryName')
8     .by('primaryTitle')
9   .group().by(select('artist', 'film'))
10    .by(__.select('roles').label().dedup().count())
11  .unfold()
12  .where(select(values).is(gt(1)))
13  .select(keys)
```

Query

This query finds artists who have had more than one type of role in the same film, then groups and lists these artists along with the films. It begins by selecting all vertices labeled "Film" that have a "primaryTitle", marking these vertices as "film". Then, it looks for incoming edges to these film vertices, marking them as "roles", which represent the different roles artists have had in these films. Next, it follows those edges to vertices labeled "Artist" that have a "primaryName", marking these vertices as "artist". The query constructs a pair of "artist" and "film" and uses this pair as a key to group the data. For each artist-film pair, it counts the distinct types of roles (using deduped edge labels to ensure uniqueness) associated with the pair. After grouping, it unfolds the results to filter out artist-film pairs for only those with more than one distinct role. Finally, it selects the keys from this filtered set, which include the artist's name and the film title.

However, this query did not work properly as it lasted forever to be completed. That could have happened due to the complexity of several traversal steps, or because of

the volume of data. One possible solution for the problem would be simplifying the query by avoiding the use of traversal-based calculations within the `group()` step, by performing some of the data processing application-side after fetching more basic data from the database. Other approach would be breaking down the query into multiple simpler queries by firstly retrieving artist-film pairs, and then, in subsequent queries, process these pairs to count distinct roles.

Exercice 11 (5 pt BONUS) : Trouver le nom du ou des film(s) ayant le plus d'acteurs.

```
1 g.V().hasLabel('Film').has("primaryTitle")
2   .group()
3     .by('primaryTitle')
4     .by(__.in("acted in").count())
5   .unfold()
6   .order().by(select(values), decr)
7   .limit(3)
```

Query

```
1 [
2   {
3     "L'histoire de mon 20 au cours Infrastructure de donnees": 19
4   },
5   {
6     "Noches de Sol": 10
7   },
8   {
9     "Carole & Grey": 10
10  }
11 ]
```

Result

This query identifies the top three films with the highest number of actors and retrieves their titles along with the count of actors for each. It starts by selecting vertices labeled "Film" that also have a "primaryTitle" property to ensure consistency (some films did not have this property correctly written and that was raising an error). The query then groups these films by their titles and counts the number of incoming "acted in" edges for each film, which represents the number of actors associated with that film. After grouping, the query unfolds the results to break down the grouped structure into individual items that can be sorted. It then orders these items in descending order based on the count of actors (using `order().by(select(values), decr)`), ensuring the films with the most actors are listed first. Finally, the query limits the output to the top three films.

A bashrc.sh

```
1 export TPBDD_NEO4J_SERVER="bolt://44.214.180.114:7687"
2 export TPBDD_NEO4J_USER="neo4j"
3 export TPBDD_NEO4J_PASSWORD="commands-atom-sheets"
```

B export-neo4j.py

```
1 from py2neo import Graph
2 from py2neo.bulk import create_nodes, create_relationships
3 from py2neo.data import Node
4 import os
5 import pyodbc
6
7 server = os.environ["TPBDD_SERVER"]
8 database = os.environ["TPBDD_DB"]
9 username = os.environ["TPBDD_USERNAME"]
10 password = os.environ["TPBDD_PASSWORD"]
11 driver= '{ODBC Driver 18 for SQL Server}'
12
13 neo4j_server = os.environ["TPBDD_NEO4J_SERVER"]
14 neo4j_user = os.environ["TPBDD_NEO4J_USER"]
15 neo4j_password = os.environ["TPBDD_NEO4J_PASSWORD"]
16
17 graph = Graph(neo4j_server, auth=(neo4j_user, neo4j_password))
18
19 BATCH_SIZE = 10000
20
21 print("Deleting existing nodes and relationships...")
22 graph.run("MATCH ()-[r]->() DELETE r")
23 graph.run("MATCH (n:Artist) DETACH DELETE n")
24 graph.run("MATCH (n:Film) DETACH DELETE n")
25
26 with pyodbc.connect('DRIVER='+driver+';SERVER=tcp:'+server+';PORT
    =1433;DATABASE='+database+';UID='+username+';PWD='+ password) as
    conn:
27     cursor = conn.cursor()
28
29     # Films
30     exportedCount = 0
31     cursor.execute("SELECT COUNT(1) FROM tFilm")
32     totalCount = cursor.fetchval()
33     cursor.execute("SELECT idFilm, primaryTitle, startYear FROM tFilm
    ")
34     while True:
35         importData = []
36         rows = cursor.fetchmany(BATCH_SIZE)
37         if not rows:
38             break
39
```

```
40     i = 0
41     for row in rows:
42         # Créer un objet Node avec comme label Film et les
proprietes adequates
43         n = Node("Film", idFilm=row[0], primaryTitle=row[1],
startYear=row[2])
44         importData.append(n)
45         i += 1
46
47     try:
48         create_nodes(graph.auto(), importData, labels={"Film"})
49         exportedCount += len(rows)
50         print(f"{exportedCount}/{totalCount} film records
exported to Neo4j")
51     except Exception as error:
52         print(error)
53
54     # Artists
55     # En vous basant sur ce qui a été fait dans la section précédente
, exportez les données de la table tArtist
56     # A COMPLETER
57     exportedCount = 0
58     cursor.execute("SELECT COUNT(1) FROM tArtist")
59     totalCount = cursor.fetchall()
60     cursor.execute("SELECT idArtist, primaryName, birthYear FROM
tArtist")
61     while True:
62         importData = []
63         rows = cursor.fetchmany(BATCH_SIZE)
64         if not rows:
65             break
66
67     i = 0
68     for row in rows:
69         # Créer un objet Node avec comme label Film et les
proprietes adequates
70         # A COMPLETER
71         n = Node("Artist", idArtist=row[0], primaryName=row[1],
birthYear=row[2])
72         importData.append(n)
73         i += 1
74
75     try:
76         create_nodes(graph.auto(), importData, labels={"Artist"})
77         exportedCount += len(rows)
78         print(f"{exportedCount}/{totalCount} artist records
exported to Neo4j")
79     except Exception as error:
80         print(error)
81
82     try:
83         print("Indexing Film nodes...")
```



```
84     graph.run("CREATE INDEX FOR (n:Film) ON n.idFilm")
85     print("Indexing Artist nodes...")
86     graph.run("CREATE INDEX FOR (n:Artist) ON n.idArtist")
87 except Exception as error:
88     print(error)
89
90
91 # Relationships
92 exportedCount = 0
93 cursor.execute("SELECT COUNT(1) FROM tJob")
94 totalCount = cursor.fetchval()
95 cursor.execute(f"SELECT idArtist, category, idFilm FROM tJob")
96 while True:
97     importData = { "acted in": [], "directed": [], "produced":
98 [], "composed": [] }
99     rows = cursor.fetchmany(BATCH_SIZE)
100     if not rows:
101         break
102
103     for row in rows:
104         relTuple=(row[0], {}, row[2])
105         importData[row[1]].append(relTuple)
106
107     try:
108         for cat in importData:
109             # Utilisez la fonction create_relationships de py2neo
110             # pour créer les relations entre les noeuds Film et Artist
111             # (les tuples nécessaires ont déjà été créés ci-
112             dessus dans la boucle for précédente)
113             # https://py2neo.org/2021.1/bulk/index.html
114             # ATTENTION: remplacez les espaces par des _ pour
115             nommer les types de relation
116             create_relationships(graph.auto(), importData[cat],
117 cat.upper().replace(" ", "_"), start_node_key=("Artist", "idArtist
118 "), end_node_key=("Film", "idFilm"))
119             exportedCount += len(rows)
120             print(f"{exportedCount}/{totalCount} relationships
121 exported to Neo4j")
122         except Exception as error:
123             print(error)
```