

Bài tập lớn BigData

Nhóm 4

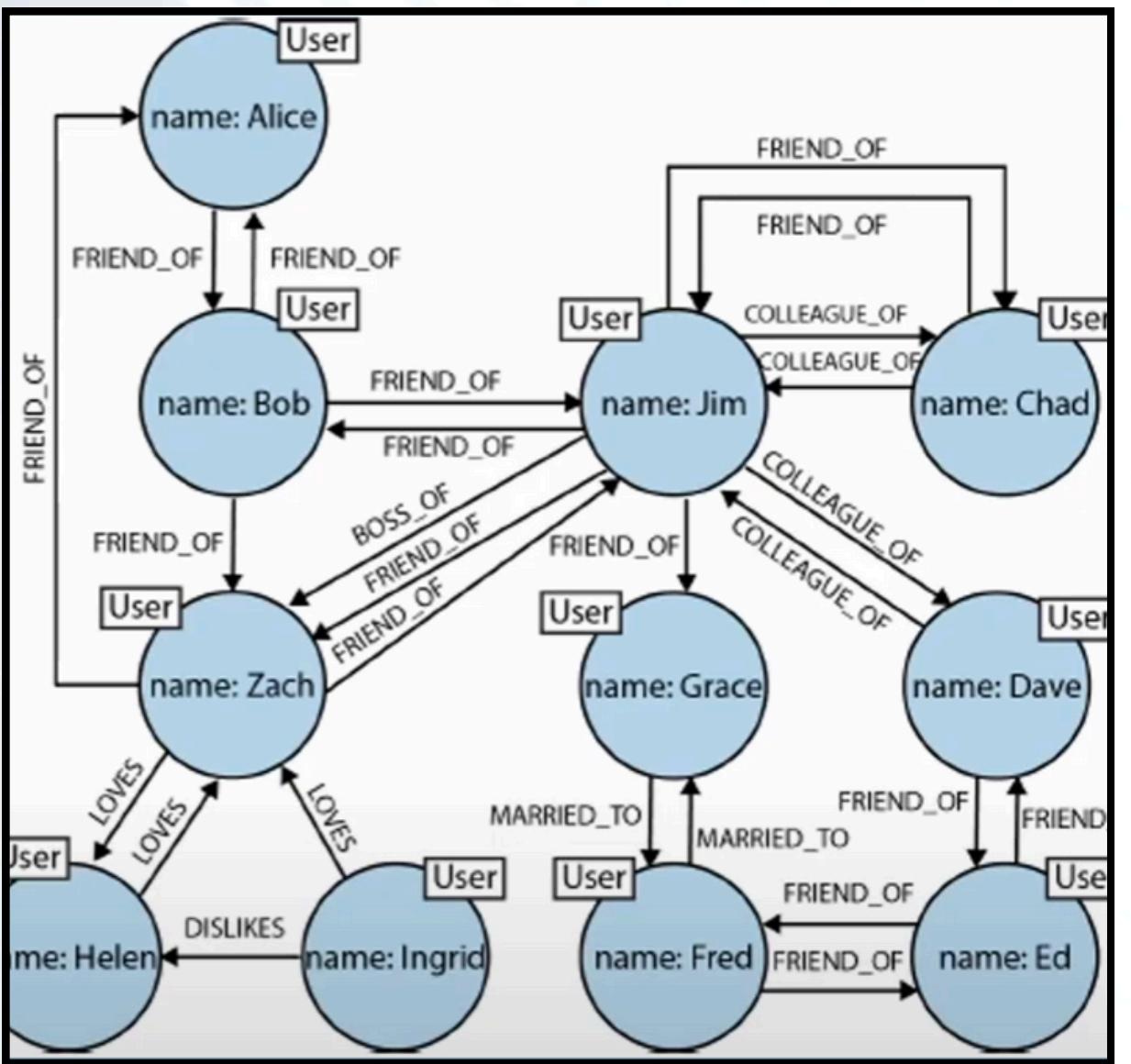


Graph Analytics With Chat Data Using Neo4j

Thành viên nhóm

- 22022522 Đàm Thái Ninh
- 22022653 Long Trí Thái Sơn
- 22022548 Hoàng Đăng Khoa

Graph DataBase



Graph Data sử dụng các nút, cạnh và thuộc tính để biểu diễn và lưu trữ dữ liệu

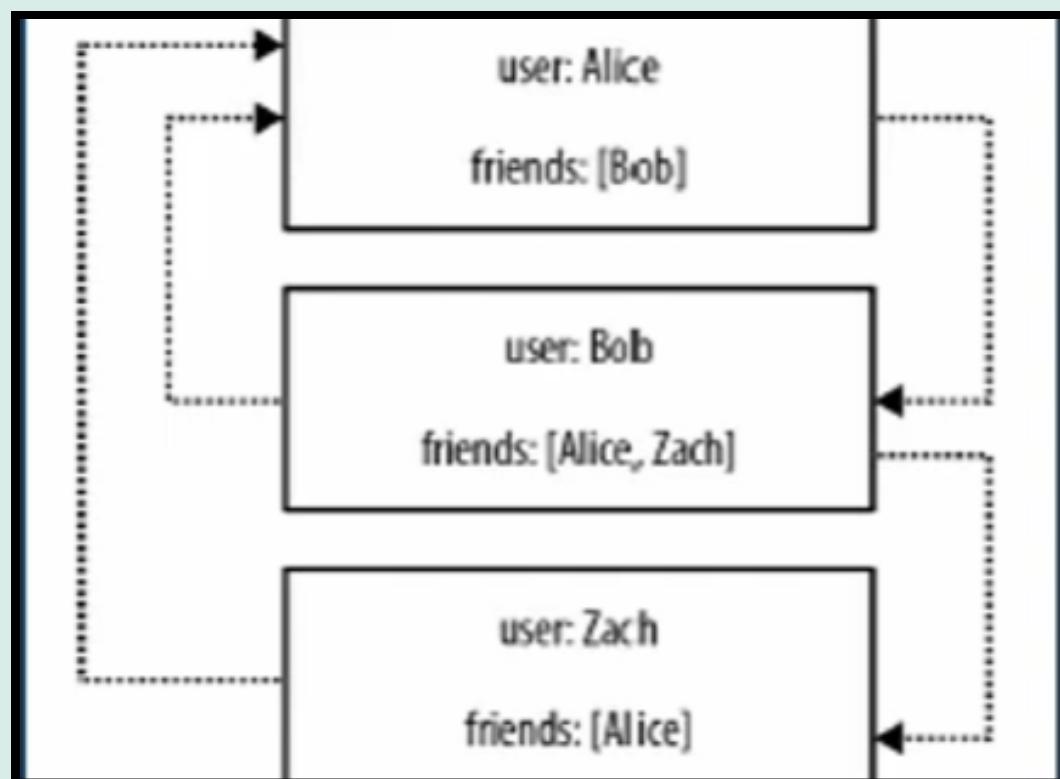
Lợi ích

- Hiệu quả
 - Linh hoạt
 - Nhanh chóng

SQL, NOSQL và Graph Database

| Person | | PersonFriend | |
|--------|--------|--------------|----------|
| ID | Person | PersonID | FriendID |
| 1 | Alice | 1 | 2 |
| 2 | Bob | 2 | 1 |
| ... | ... | 2 | 99 |
| 99 | Zach | ... | ... |
| 99 | | 99 | 1 |

SQL

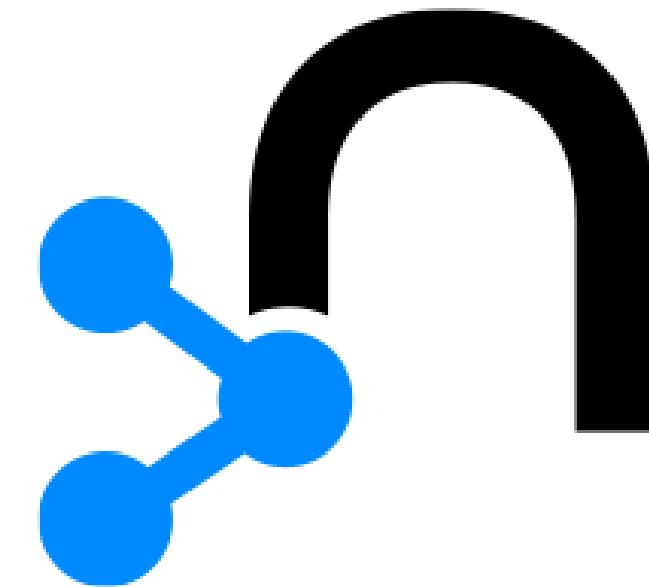


NOSQL

| Depth | RDBMS execution time(s) | Neo4j execution time(s) | Records returned |
|-------|-------------------------|-------------------------|------------------|
| 2 | 0.016 | 0.01 | ~2500 |
| 3 | 30.267 | 0.168 | ~110,000 |
| 4 | 1543.505 | 1.359 | ~600,000 |
| 5 | Unfinished | 2.132 | ~800,000 |

Source: Neo4j in Action

Neo4j



Là gì?

Neo4j là một cơ sở dữ liệu đồ thị, cho phép lưu trữ và truy vấn mối quan hệ giữa các dữ liệu một cách hiệu quả

Neo4j Browser

Thực hiện truy vấn

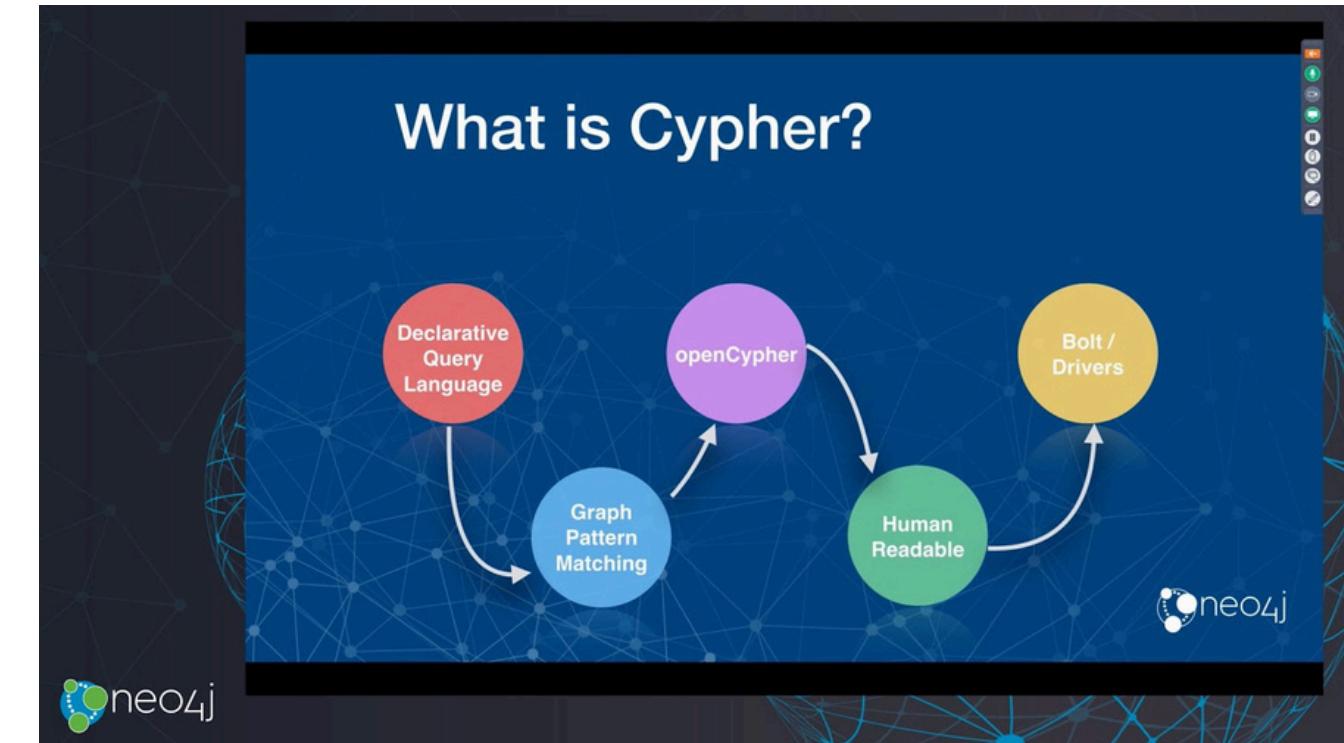
Cách dùng

Có thể tạo ràng buộc duy nhất cho các nút và cạnh, sau đó tải dữ liệu từ các tệp CSV và xây dựng đồ thị

Neo4j Bloom

Visualize dữ liệu

Cypher – Ngôn ngữ truy vấn của Neo4j



Tạo Node

```
CREATE (n:Person {name: 'Alice', age: 30})
```

Tạo Relationship

```
MATCH (a:Person {name: 'Alice'}), (b:Person {name: 'Bob'})  
CREATE (a)-[:FRIEND_OF]->(b)
```

Truy vấn Node

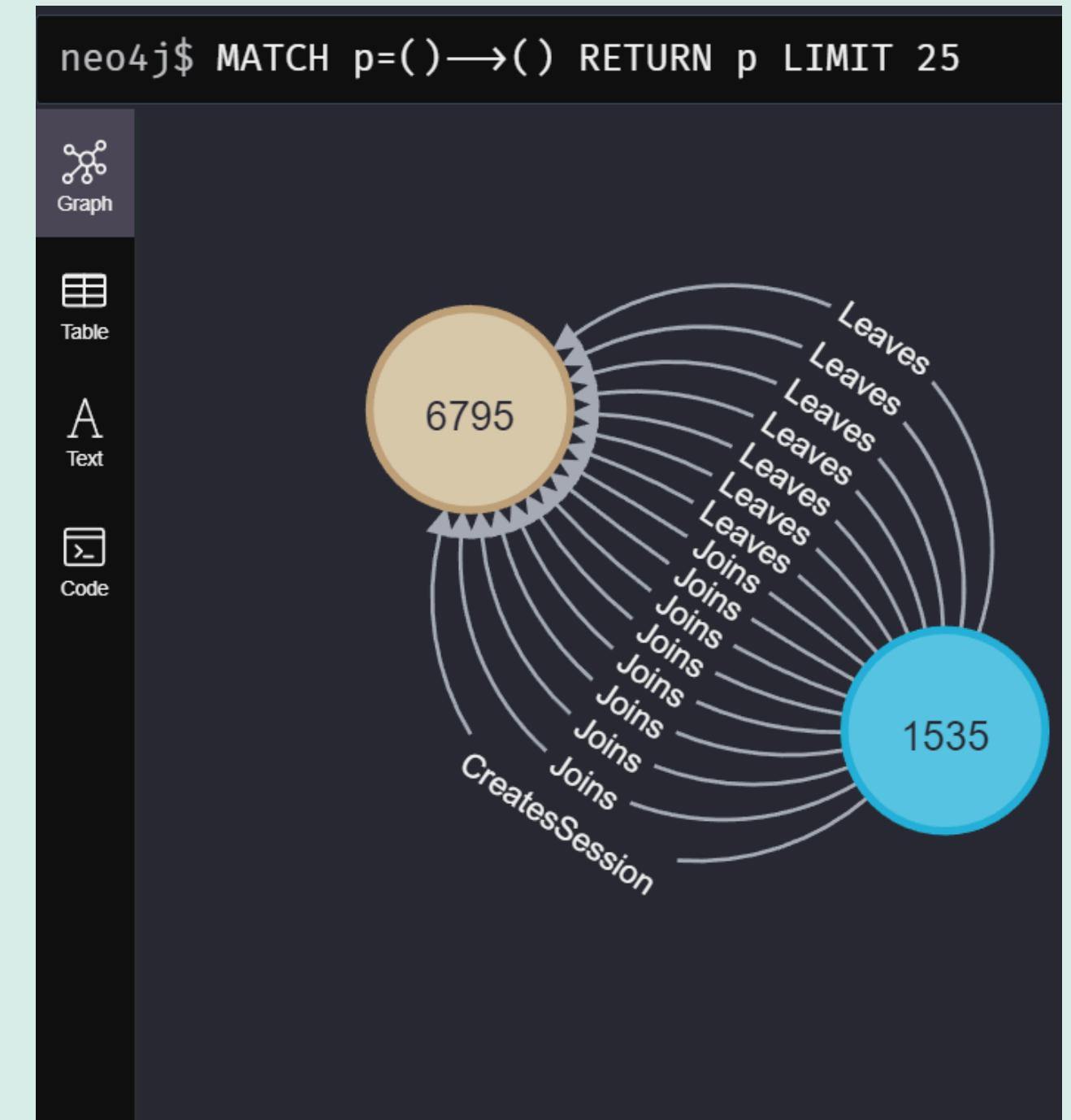
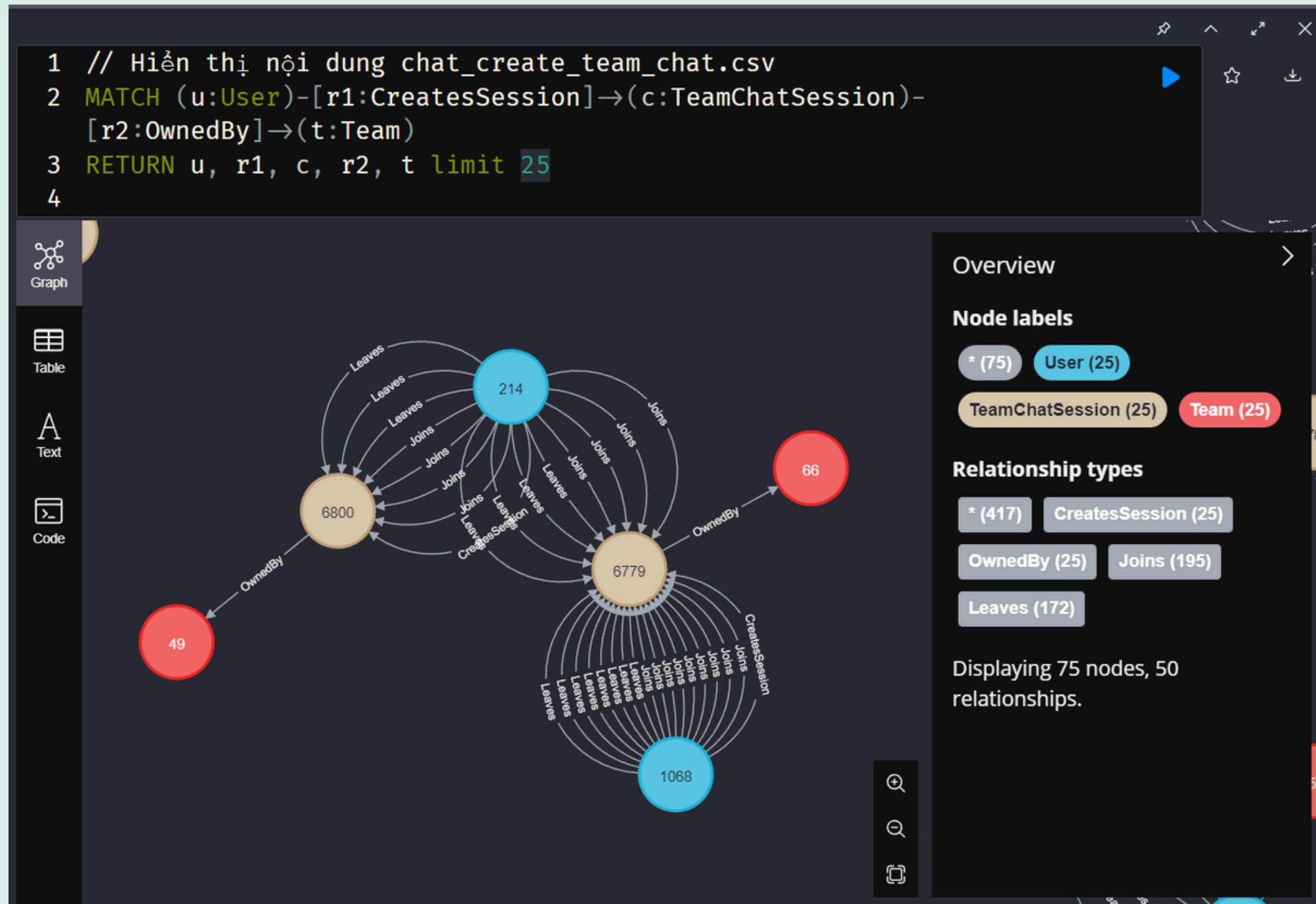
```
MATCH (n:Person)  
RETURN n
```

Graph Analytics of Catch the Pink Flamingo Chat Data Using Neo4j

| File Name | Description | Fields |
|---------------------------|-------------------|---|
| chat_create_team_chat.csv | userid | the user id assigned to the user |
| | teamid | the id of the team |
| | teamChatSessionID | a unique id for the chat session |
| | timestamp | a timestamp denoting when the chat session created |
| chat_item_team_chat.csv | userid | the user id assigned to the user |
| | teamchatsessionid | a unique id for the chat session |
| | chatitemid | a unique id for the chat item |
| | timestamp | a timestamp denoting when the chat item created |
| chat_join_team_chat.csv | userid | the user id assigned to the user |
| | teamChatSessionID | a unique id for the chat session |
| | timestamp | a timestamp denoting when the user join in a chat session |

| | | |
|----------------------------|-------------------|---|
| chat_leave_team_chat.csv | userid | the user id assigned to the user |
| | teamchatsessionid | a unique id for the chat session |
| | timestamp | a timestamp denoting when the user leave a chat session |
| chat_mention_team_chat.csv | ChatItemId | the id of the ChatItem |
| | userid | the user id assigned to the user |
| | timeStamp | a timestamp denoting when the user mentioned by a chat item |
| chat_respond_team_chat.csv | chatid1 | the id of the chat post 1 |
| | chatid2 | the id of the chat post 2 |
| | timestamp | a timestamp denoting when the chat post 1 responds to the chat post 2 |

Graph Analytics



Graph Analytics

```
1 // Question 1-a
2 // Find the longest conversation chain in the chat data using the
3 // "ResponseTo" edge label. This question has two parts
4 MATCH p=(start:ChatItem)-[:ResponseTo*]→(end:ChatItem)
5 RETURN p, length(p) AS longestChainLength
6 ORDER BY longestChainLength DESC
7 LIMIT 1;
8
```

p

| longestChainLength |
|--------------------|
| 9 |

```
1 // Question 1-a
2 // Find the longest conversation chain in the chat data using the
3 // "ResponseTo" edge label. This question has two parts
4 MATCH p=(start:ChatItem)-[:ResponseTo*]→(end:ChatItem)
5 RETURN p, length(p) AS longestChainLength
6 ORDER BY longestChainLength DESC
7 LIMIT 1;
8
```

Graph

Overview

Node labels

- * (10)
- ChatItem (10)

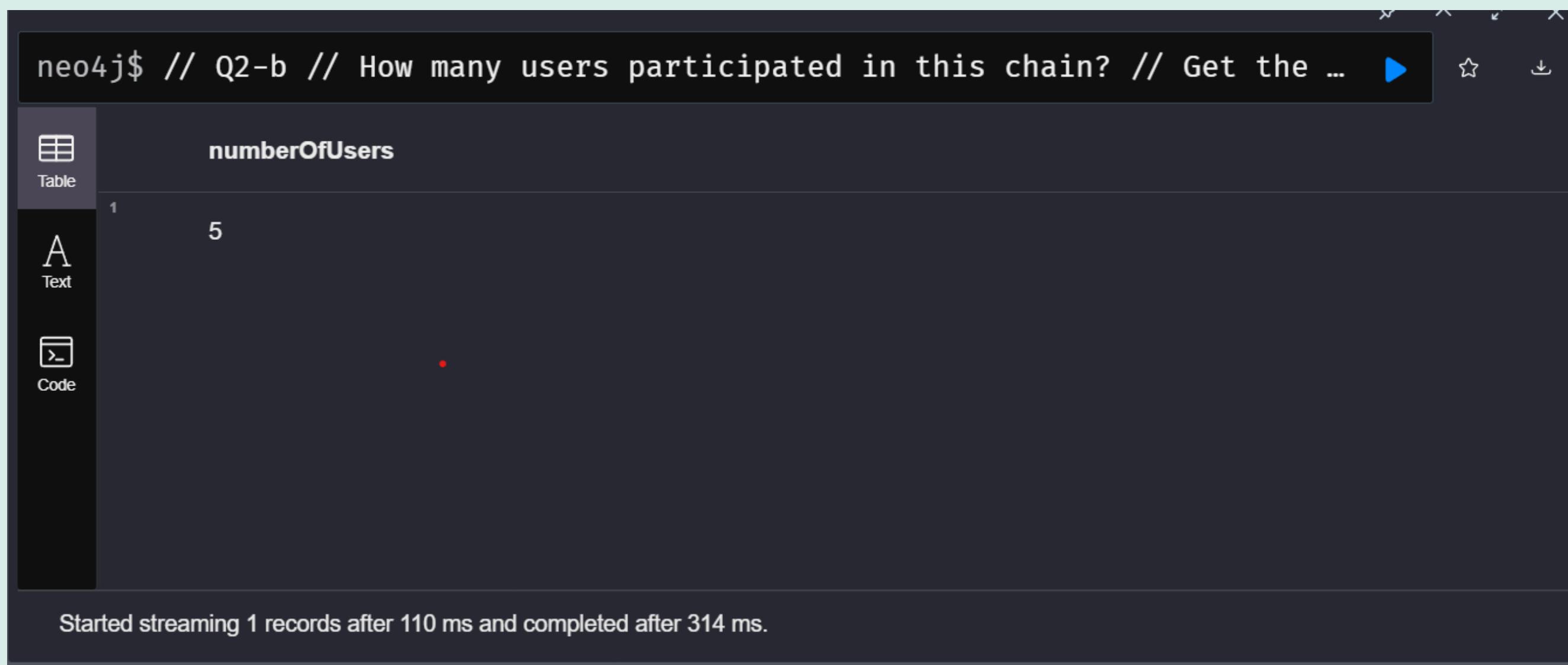
Relationship types

- * (9)
- ResponseTo (9)

Displaying 10 nodes, 9 relationships.

Started streaming 1 records after 8 ms and completed after 96 ms.

Graph Analytics



A screenshot of a terminal window with a dark background. The window title bar says "neo4j\$ // Q2-b // How many users participated in this chain? // Get the ...". Below the title bar is a toolbar with icons for copy, star, and download. On the left is a sidebar with three items: "Table" (selected), "Text" (with icon A), and "Code" (with icon Σ). The main area shows a table with one row. The table has a single column labeled "numberOfUsers" with the value "5". At the bottom of the window, a status message reads "Started streaming 1 records after 110 ms and completed after 314 ms."

| numberOfUsers |
|---------------|
| 5 |

Started streaming 1 records after 110 ms and completed after 314 ms.

Graph Analytics

Xác định xem những người chat nhiều nhất có tham gia team tạo nhiều chat nhất không?

```
1 // Q2-b
2
3 // Chattiest Teams
4
5 // Query to find the top 10 chattiest teams
6 MATCH (ci:ChatItem)-[:PartOf]→(tcs:TeamChatSession)-[:OwnedBy]→
  (t:Team)
7 RETURN t.id AS teamId, count(ci) AS chatCount
8 ORDER BY chatCount DESC
9 LIMIT 10;
10
11
```

| teamId | chatCount |
|--------|-----------|
| 82 | 1324 |
| 185 | 1036 |
| 112 | 957 |
| 18 | 844 |
| 194 | 836 |
| 129 | 814 |

▶ ☆ 4

```
5 // Chattiest Users
6
7 // Query to find the top 10 chattiest users
8 MATCH (u:User)-[:CreatesChat]→(c:ChatItem)
9 RETURN u.id AS userId, count(c) AS chatCount
10 ORDER BY chatCount DESC
11 LIMIT 10;
12
```

Table

| userId | chatCount |
|--------|-----------|
| 394 | 115 |
| 2067 | 111 |
| 1087 | 109 |
| 209 | 109 |
| 554 | 107 |
| 1627 | 105 |

Text

Code

Started streaming 10 records after 7 ms and completed after 40 ms.

Graph Analytics

Xác định xem những người chat nhiều nhất có tham gia team tạo nhiều chat nhất không?



The screenshot shows the Neo4j browser interface with a query results table. The table has two columns: 'User' and 'Team'. There is one row of data with values 999 and 52 respectively. The table is displayed in a dark-themed environment.

| User | Team |
|------|------|
| 999 | 52 |

Started streaming 1 records after 25 ms and completed after 48 ms.

```
neo4j$ // Q2-c // Step 1: Identify top 10 chattiest users WITH [394, 206...]
```

Graph Analytics

```
1
2 MATCH
3 (u:User)-[c:CreatesChat]→()
4 WITH u, COUNT(c) as Chats
5 ORDER BY Chats DESC LIMIT 10 WITH [u] as ChattiestUsers
6 //Getting the neighbours of all Users and the count
7 MATCH (u1:User)-[:InteractsWith]→(u2:User)
8 WHERE u1 in ChattiestUsers
9 WITH u1.id AS UserID, COLLECT(DISTINCT u2.id) AS Neighbours RETURN
UserID, Neighbours, SIZE(Neighbours) AS k
```

| | UserID | Neighbours | k |
|---|--------|-----------------------------------|---|
| 1 | 394 | [1012, 2011, 1997, 1782] | 4 |
| 2 | 2067 | [63, 516, 1265, 1672, 209, 1627] | 6 |
| 3 | 1087 | [1311, 426, 929, 772, 1879, 1098] | 6 |
| 4 | 209 | [63, 516, 1627, 2067, 1672] | 5 |
| 5 | 554 | [2018, 1959, 1687, 1096, 1010] | 5 |
| 6 | 1627 | [516, 2067, 209, 63, 1672, 1265] | 6 |
| 7 | | | |

Started streaming 10 records after 11 ms and completed after 74 ms.

Xác định những người hoạt động
tích cực nhất

Những khó khăn gặp phải

- Vấn đề khi chạy trên WSL
- Lạ lẫm với Cypher
- Gặp nhiều vấn đề khi tích hợp với Spark

Dự định

- Tiếp tục tìm cách xử lý Neo4j với Spark
- Có thể sẽ sử dụng Docker

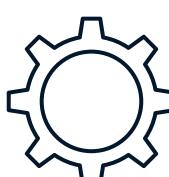
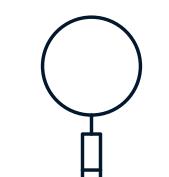
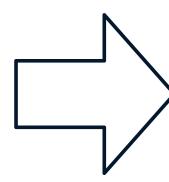
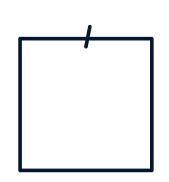
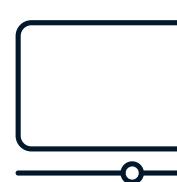
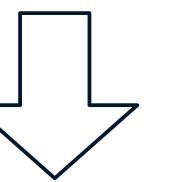
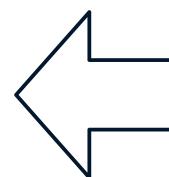
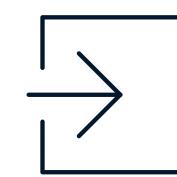
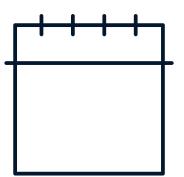
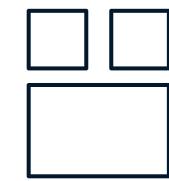
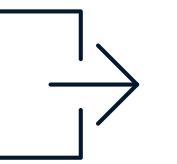




Demo

- Các câu lệnh đơn giản
- Graph Analytics of Catch the Pink Flamingo Chat Data Using Neo4j

QA



Xin chân thành cảm ơn