

Assignment Number 2  
-Deepika Kini

**Q.1 please check comments below**

- **Title (id, type, title, originalTitle, startYear, endYear, runtime, avgRating, numVotes)**

Preprocessing: no changes to primary key(unlike previous assignment)

Joined with ratings table to get extra columns

File: preprocess.py

Time for preprocessing and loading: 5 mins

- **Genre (id, genre)**

Preprocessing: Take unique genres from title and break it down, add a unique identifier

Time for preprocessing and loading: 1 mins

- **Title\_Genre (genre, title)**

- genre FK Genre(id)

- title FK Title(id)

Preprocessing: map title id with many genre ids(break into many rows using explode)

Time for preprocessing and loading: 2 mins

- **Member (id, name, birthYear, deathYear)**

Preprocessing: Used name.basics file

Time for preprocessing and loading: 3 mins

- **Title\_Actor (actor, title)**

- actor FK Member(id)

- title FK Title(id)

- **Title\_Producer (producer, title)**

- producer FK Member(id)

- title FK Title(id)

Preprocessing: Used principals file to filter on actor/actress/producer category

Time for preprocessing and loading: 55 mins

- **Title\_Director (director, title)**

- director FK Member(id)

- title FK Title(id)

- **Title\_Writer (writer, title)**

- writer FK Member(id)

- title FK Title(id)

Preprocessing: Used crew file to get columns of writer and director and then spilt them into rows(contained arrays of ids)

Time for preprocessing and loading: 45 mins

- **Character (id, character)**

**Note: this table should contain individual characters**

*Tables are created in python(unlike my previous implementation of hw1) . here i've use psycopg2 library*

- **Actor\_Title\_Character (actor, title, character)**
  - **(actor,title) FK Title\_Actor (actor, title)**
  - **character FK Character(id)**

Preprocessing: Used explode to create these tables

Time for preprocessing and loading: 35 mins

**Note:** For rows in title\_actor, ... they should be present in both title and principal files. Hence inner joins are run

*Tables are created in python(unlike my previous implementation of hw1) . here i've use psycopg2 library*

*Below queries in psql used to load data quicker:(due to slowness later on)*

```
\copy producer(title, producer) FROM
'/Users/deepika/Desktop/CSCI620/Assignment_1/IMDB/ModifiedData/producer_mapping_project2.csv'(format csv, delimiter ',', header true);
```

```
\copy ACTOR(title, actor) FROM
'/Users/deepika/Desktop/CSCI620/Assignment_1/IMDB/ModifiedData/actor_mapping_project2.csv'(format csv, delimiter ',', header true);
Time: 10mins
```

Wasn't able to load director and writer table(with primary keys applied) since they had duplicates in primary keys which I wasn't able to handle in python.

**RUNTIME: 4.25 hours approx**

**Loading using copy command took time hence pushed data first into sql and added constraints:**

SQL:

```
create table atc
(
    title varchar(20),
```

```
        actor varchar(20),
character integer
    )
```

```
insert into atc(
select actor,title,character from actor_title_character)
```

```
create table actor_title_character1(title, actor, character) as
select * from actor_t_c group by title, actor, id
```

```
Alter table a_t_c ADD PRIMARY KEY(id, ACTOR, TITLE) ;
```

```
Alter table a_t_c ADD FOREIGN KEY(character) REFERENCES Character(id);
```

Not run:

```
Alter table a_t_c ADD FOREIGN KEY(TITLE, ACTOR) REFERENCES title_ACTOR(TITLE,
ACTOR); since extra rows present in this table
(wasn't able to disable constraints using the code:
```

```
sql = "set constraints all deferred;"
cursor.execute(sql)
conn.commit() )
```

**Please refer to main\_assignment2\_preprocessing.py and hw2.sql for preprocessing and create table codes. Other python files (other than queries\_script.py) are used for preprocessing**

-----

## **Q.2.**

2.1. Number of invalid Title\_Actor relationships with respect to characters. (That is, entries in Title\_Actor which do not appear in Actor\_Title\_Character.)

```
select title, actor from title_actor where concat(title, actor) not in (select concat(title, actor) from
atc )
```

*No outcome, just runs*

2.2. Alive actors whose name starts with “Phi” and did not participate in any title in 2014.

```
select m.name, t.title, t.startYear from title_actor a
inner join member m on a.actor = m.id
inner join title t on a.title = t.id
where m.deathYear = 0 and t.startYear != 2014 and endYear != 2014
and lower(m.name) like 'phi%';
Time: 19.2sec
```

2.3. Producers who have produced the most talk shows in 2017 and whose name contains “Gill”. (Hint: talk show is a genre)

```
select m.name, count(t.id) from member m
join title_producer p on m.id = p.producer
  join title t on t.id = p.title
  join title_genre tg on tg.title = t.id
join genre g on tg.genre = g.id
where m.name like '%gill%' and t.startYear = 2017 and lower(g.genre) like '%talk-show%'
group by m.name
order by count(t.id) desc
Time:21 secs
```

2.4. Alive producers ordered by the greatest number of long-run titles produced (runtime greater than 120 minutes).

```
select m.name, count(t.title)
from title_producer p
join member m on p.producer = m.id
join title t on t.id = p.title
where t.runtime > 120 and m.deathYear != 0
group by m.name
order by count(t.title) desc
Time :23 secs
```

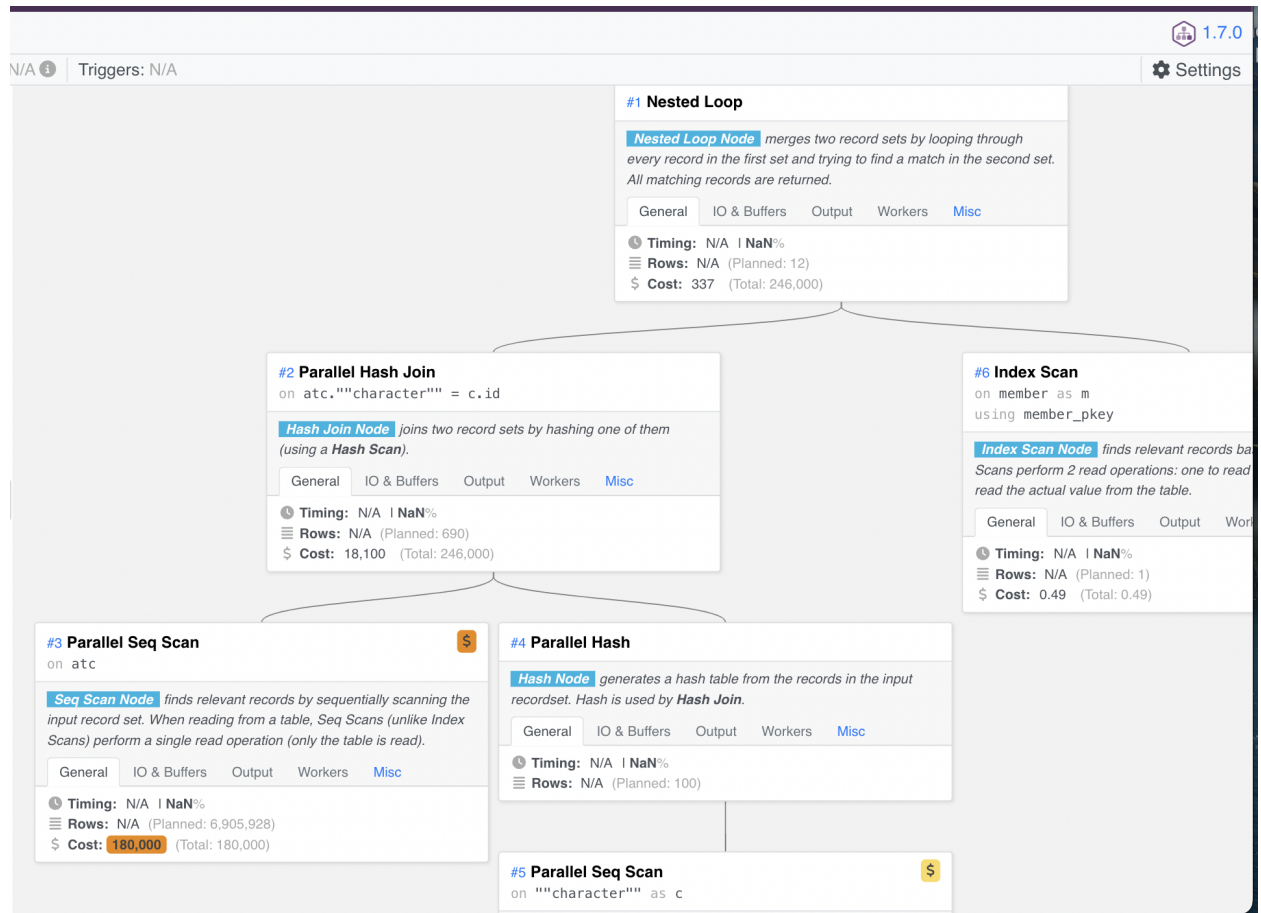
2.5. Alive actors who have portrayed Jesus Christ (simply look for a character with this specific name).

```
select m.name, atc.character, c.character from member m  
join atc atc on atc.actor = m.id  
join character c on c.id = atc.character  
where lower(c.character) like '%jesus christ%' and m.deathYear != 0  
Time: 1min 20 sec
```

### **Q. 3.**

2.1.query  
No outcome

2.2



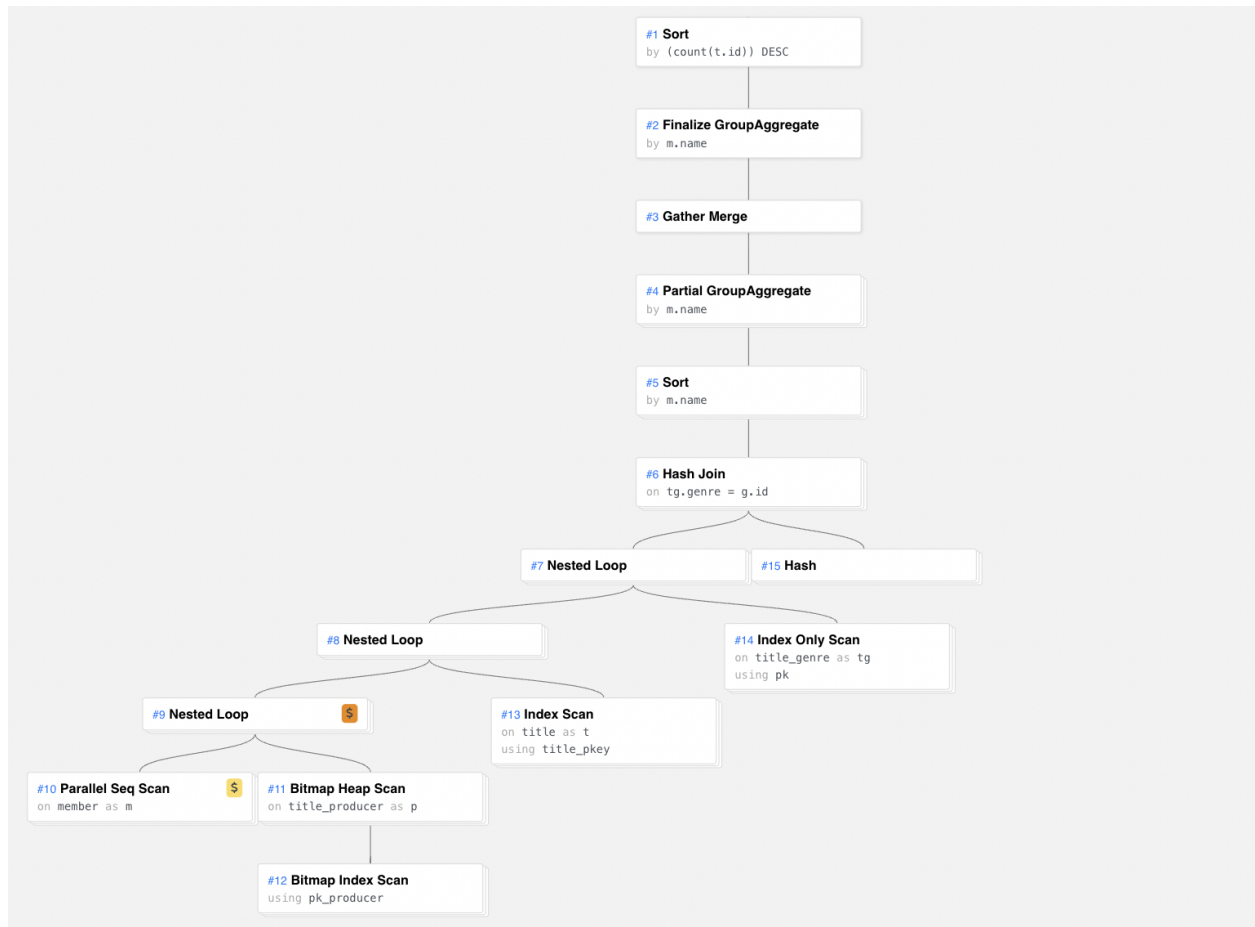
The nested loop is for joins. Member table uses table to filter out year information  
Index of title\_actor is used to filter out on title\_actor for join.

Table scan is basically seq scan and index scan uses primary key as index

Explaining tree structure:

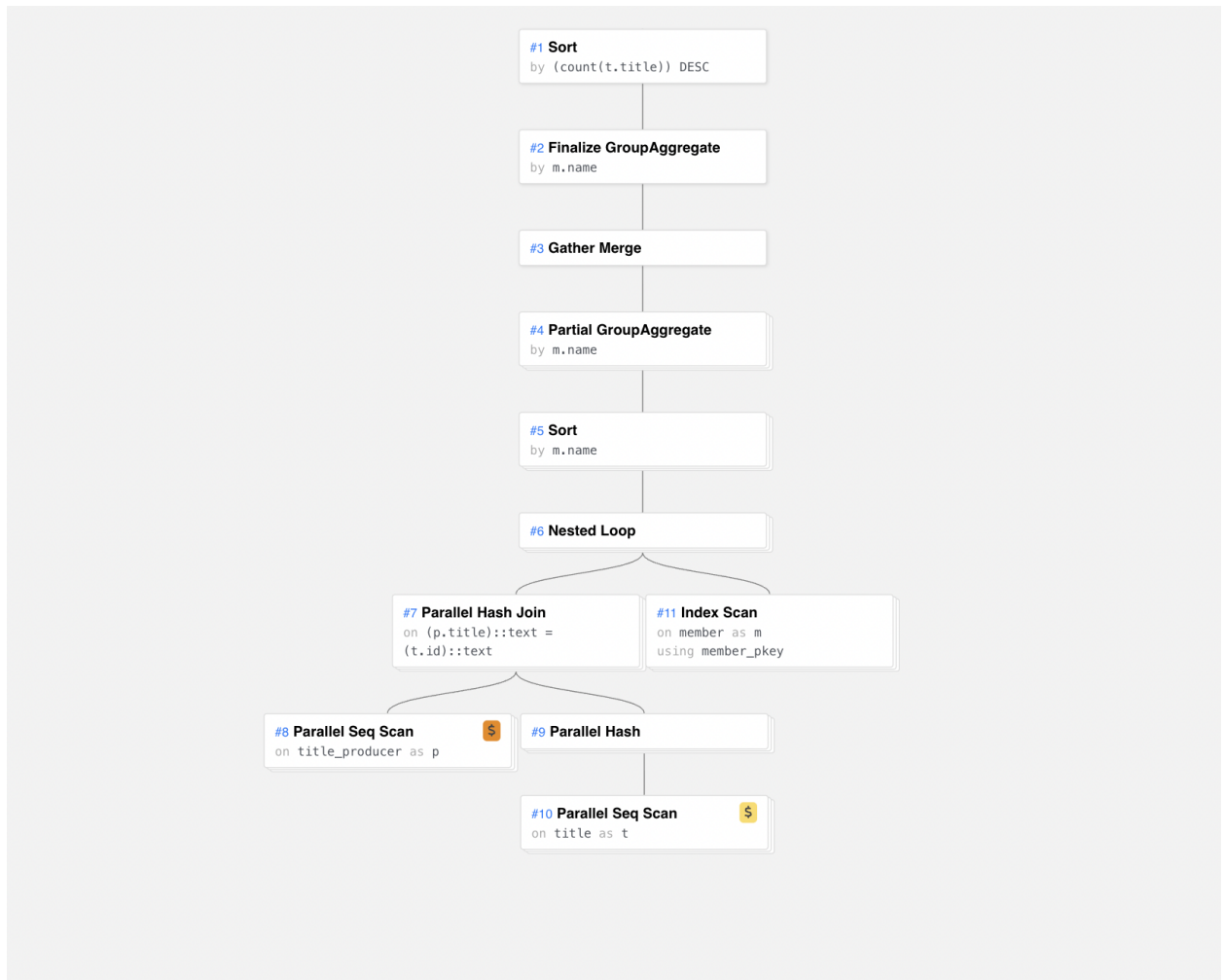
member(pulled by data scan) and title\_Actor(scanned on index) join is carried out in the leaf node of tree and nested loop is applied to check on equal values.(line 1 of sql)

Similarly the outcome of this is combined with a join on title\_key (line 2 of sql) which gives the 2nd level data. At the end the data is gathered (the columns are chosen)



*Bitmap index filters data based on index of the table. Similar to above query outcome*

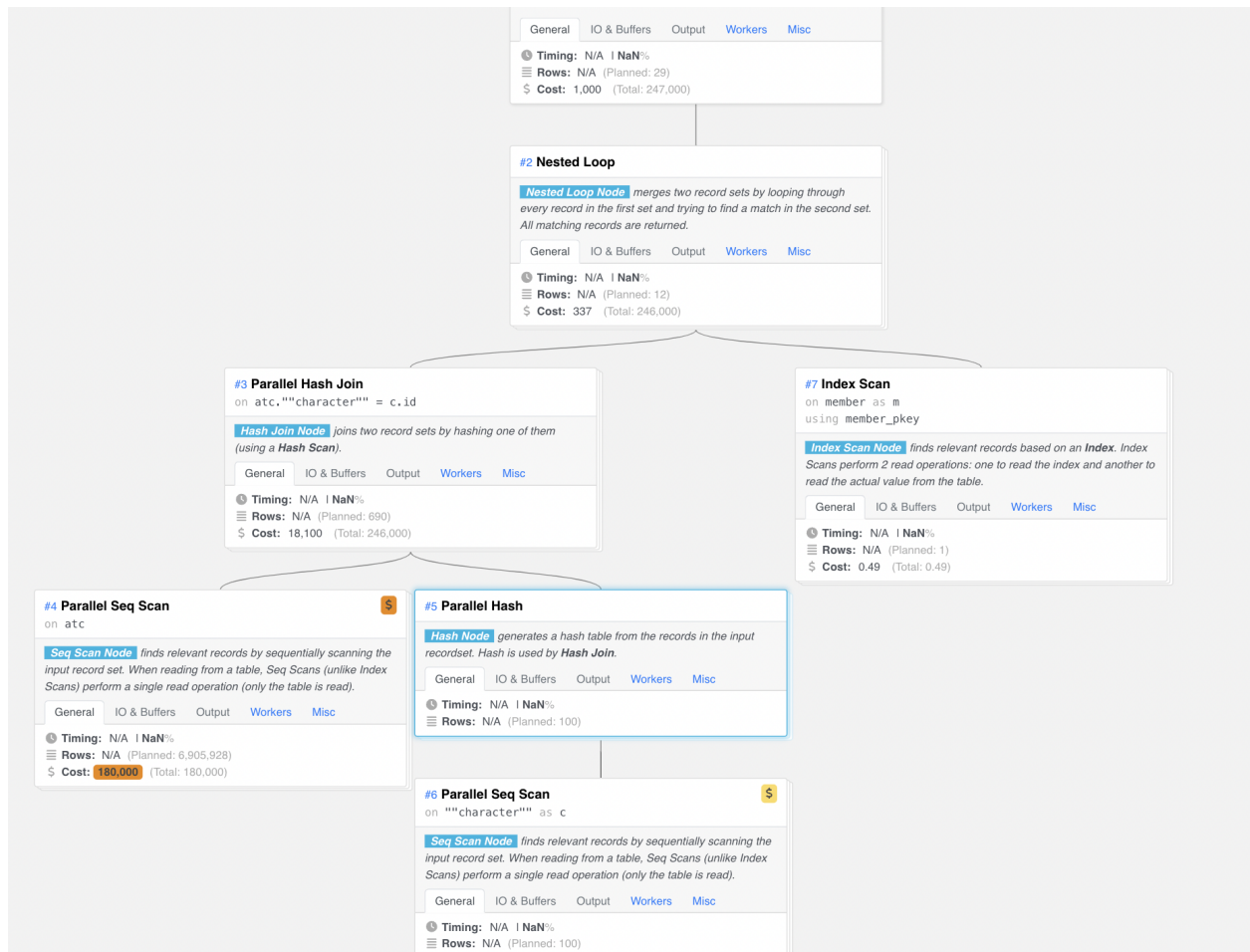
2.4



Where and join takes place first, then aggregate works, and then sort (based on count). For **text columns**, **hash** is used as seen above. M.name is sorted first and then applied the filter to



# GATHER



Q.4.

Below are the relational algebra queries corresponding to sql queries:

(2.1)  $\pi_{\text{count}}(\pi_{\text{title, actor}}(\text{title\_actor})) - \pi_{\text{title, actor}}(\text{title\_actor\_character})$

(2.2)

$\pi_{\text{m.name}}(\sigma_{\text{death\_year} \neq 0 \wedge \text{start\_year} \neq 2014 \wedge \text{end\_year} \neq 2014 \wedge \text{lower}(\text{name}) \text{ like } 'Phi-.'}(\text{title\_actor} \times \text{title} \times \text{member}))$

(2.3)  $\pi_{\text{m.name}}(\sigma_{\text{count}(\text{t.title})}(\sigma_{\text{t.start\_year} = 2014 \wedge \text{lower}(\text{g.genre}) \text{ like } 'talk-show.'}(\text{member} \times \text{title} \times \text{producer} \times \text{title\_genre} \times \text{genre})))$

(2.4)  $\pi_{\text{m.name}}(\sigma_{\text{count}(\text{t.title})}(\sigma_{\text{t.start\_year} > 1920 \wedge \text{death\_year} \neq 0}(\text{title} \times \text{producer} \times \text{member} \times \text{title})))$

(2.5)  $\pi_{\text{m.name}}(\sigma_{\text{death\_year} \neq 0 \wedge \text{lower}(\text{c.character}) \text{ like } 'J.'}(\text{title\_actor\_character}))$   
 $(\text{member} \times \text{atc} \times \text{character})$

atc  $\rightarrow$  actor - title - character

2.5 with different jesus and christ terms:

$\pi_{\text{distinct}(\text{m.name})}(\sigma_{\text{death\_year} \neq 0 \wedge \text{lower}(\text{character}) \text{ like } 'J.'}(\text{title\_actor\_character}))$   
 $\wedge (\sigma_{\text{death\_year} \neq 0 \wedge \text{lower}(\text{character}) \text{ like } 'Christ.'}(\text{title\_actor\_character}))$

$(A \cup B) - (A \cap B)$

unique

Diagram showing two overlapping sets J and C, with an arrow pointing from J to C.

Q.5

Applied index to columns which are used frequently in either join, where or select clauses. Q2, 3 ,4 and 5 run faster when indexed manually

Query 2.1:

```
select title, actor from title_actor where concat(title, actor) not in (select
concat(title, actor) from atc )
```

*Taking time to run*

Query 2.2:

<i>name</i>	<i>title</i>	<i>startyear</i>
<i>Philip Martin Brown</i>	<i>The Birdman</i>	<i>1994</i>
<i>Philip Bretherton</i>	<i>The Car Washer</i>	<i>2015</i>
<i>Philippa Baker</i>	<i>Episode #1.97</i>	<i>1972</i>
<i>Philipp Seiser</i>	<i>Die Beförderung</i>	<i>1983</i>
<i>Philipp Stadler</i>	<i>Mirakel</i>	<i>1990</i>
<i>Phill Jupitus</i>	<i>The Great Water</i>	<i>2000</i>
<i>Phil Brooks</i>	<i>Episode #10.53</i>	<i>2009</i>
<i>Philip Hawthorn</i>	<i>Episode #5.28</i>	<i>2000</i>
<i>Phillip J. Roth</i>	<i>Chimes of Freedom</i>	<i>2019</i>

***Time: 12.2 sec***

Query 2.3:

***output:***

<i>"Name"</i>	<i>"count"</i>
<i>"Robert Shergill"</i>	<i>1</i>

***Time: 3 seconds***

Query 2.4 :

***Sample rows:***

<i>"name"</i>	<i>"count"</i>
---------------	----------------

"Aachi Manorama"	172
"N. T. Rama Rao"	164
"Nagesh"	157
"Shivaji Ganesan"	141

**Time: 17.9 sec**

**Query 2.5:**

**Time: 7.3 seconds**

**Sample rows:**

name	character	character-2
Claudio Brook	1598433	["Jesucristo (Jesus Christ)"]
Jon Shepodd	228330	["Jesus Christ"]
Nelson Leigh	228330	["Jesus Christ"]
Peter Chown	228330	["Jesus Christ"]
Arsenio Corsellas	228330	["Jesus Christ"]
Hans-Reinhard Müller	988461	["Jesus Christus"]
Michael Gwynn	228330	["Jesus Christ"]
Sydney Ayres	228330	["Jesus Christ"]
Nelson Leigh	1248861	["Jesus", "Jesus Christ"]
Nelson Leigh	228330	["Jesus Christ"]

**When run in python code: queries\_script.py, it takes 54 secs for indexed data and 72 sec for non indexed (not manually indexed by me) data**



```
CSCI-620 ~/PycharmProjects/assignment2
venv
├── bin
├── include
├── lib
├── assignment2
│   ├── connection.py
│   ├── duplicates.py
│   ├── main_assignment2_preprocessing.py
│   ├── member.py
│   └── producer.py
└── python3.10
    ├── .gitignore
    ├── basics.py
    ├── imdb.py
    ├── pyvenv.cfg
    ├── test.py
    └── test1.py
crew.py
main.py
main.py
duplicates.py
main_assignment2_preprocessing.py
queries()

186 --q3
187 select m.name, count(t.id) from member m
188 join title_producer p on m.id = p.producer
189 join title t on t.id = p.title
190 join title_genre tg on tg.title = t.id
191 join genre g on tg.genre = g.id
192 where m.name like '%gill%' and t.startYear = 2017 and lower(g.genre) like '%talk-show%'
193 group by m.name
194 order by count(t.id) desc;
195 --q4
196 select m.name, count(t.title)
197 from title_producer p
198 join member m on p.producer = m.id
199 join title t on t.id = p.title
200 where t.runtime > 120 and m.deathYear != 0
201 group by m.name

/Users/deepika/PycharmProjects/CSCI-620/venv/bin/python /Users/deepika/PycharmProjects/CSCI-620/venv/lib/assignment2
54.89199429098517
```

*Index sql statements:(in hw2.sql file)*

*create index q2 on member(name)*

*create index q3 on title(id, endYear, startYear)*

*create index q4 on genre(genre)*

*create index q5 on title\_actor(title, actor)*

*create index q6 on character(id, character)*

*create index q7 on actor\_title\_character(actor, character)*

-----END-----