

Database Design

Database Implementation

Connection:

```
Welcome to Cloud Shell! Type "help" to get started.
Your Cloud Platform project in this session is set to watchdojo.
tiancheng0115@cloudshell:~ (watchdojo)$ gcloud sql connect watchdojo --user=root --quiet
Allowlisting your IP for incoming connection for 5 minutes...done.
Connecting to database with SQL user [root].Enter password:
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 41
Server version: 8.0.26-google (Google)

Copyright (c) 2000, 2022, Oracle and/or its affiliates.

Oracle is a registered trademark of Oracle Corporation and/or its
affiliates. Other names may be trademarks of their respective
owners.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

mysql> use watchdojo
Reading table information for completion of table and column names
You can turn off this feature to get a quicker startup with -A

Database changed
```

DDL Commands:

```
CREATE TABLE Login (
    username VARCHAR(255) PRIMARY KEY,
    password VARCHAR(255)
);
```

```
CREATE TABLE Review (
    username VARCHAR(255) PRIMARY KEY,
    show_id INT PRIMARY KEY,
    reviews VARCHAR(4095),
    FOREIGN KEY (username) REFERENCES Login (username),
    FOREIGN KEY (show_id) REFERENCES Movie_TV (show_id)
);
```

```
CREATE TABLE Watch_list (  
    list_id INT PRIMARY KEY,  
    username VARCHAR(255),  
    FOREIGN KEY (username) REFERENCES Login (username)  
);
```

```
CREATE TABLE Existed_in (  
    list_id INT,  
    show_id INT,  
    PRIMARY KEY (list_id, show_id),  
    FOREIGN KEY (list_id) REFERENCES Watch_list(list_id),  
    FOREIGN KEY (show_id) REFERENCES Movie_TV(show_id)  
);
```

```
CREATE TABLE Movie_TV (  
    show_id INT PRIMARY KEY,  
    title VARCHAR(255),  
    year_released INT,  
    category VARCHAR(255),  
    country VARCHAR(255),  
    rating VARCHAR(20),  
    duration VARCHAR(255)  
);
```

```
CREATE TABLE Cast (  
    cast_id VARCHAR(255) PRIMARY KEY,  
    first_name VARCHAR(255),  
    last_name VARCHAR(255)  
);
```

```
CREATE TABLE Director (  
    director_id VARCHAR(255) PRIMARY KEY,  
    first_name VARCHAR(255),  
    last_name VARCHAR(255)  
);
```

```
CREATE TABLE Directed_by (  
    director_id VARCHAR(255),  
    show_id INT,  
    FOREIGN KEY (director_id) REFERENCES Director (director_id),  
    FOREIGN KEY (show_id) REFERENCES Movie_TV (show_id),  
    PRIMARY KEY (director_id, show_id)  
);
```

```
CREATE TABLE Casted_by (  
    cast_id VARCHAR(255),  
    show_id INT,  
    FOREIGN KEY (cast_id) REFERENCES Cast (cast_id),
```

FOREIGN KEY (show_id) REFERENCES Movie_TV (show_id),
PRIMARY KEY (cast_id, show_id)
);

Inserting at least 1000 rows:

```
mysql> SELECT COUNT(*) FROM Movie_TV;
+-----+
| COUNT(*) |
+-----+
|      8807 |
+-----+
1 row in set (0.14 sec)
```

```
mysql> SELECT COUNT(*) FROM Cast;
+-----+
| COUNT(*) |
+-----+
|    29443 |
+-----+
1 row in set (0.23 sec)
```

```
mysql> SELECT COUNT(*) FROM Director;
+-----+
| COUNT(*) |
+-----+
|     4985 |
+-----+
1 row in set (0.00 sec)
```

```
mysql> SELECT COUNT(*) FROM Casted_by;
+-----+
| COUNT(*) |
+-----+
|    53053 |
+-----+
1 row in set (0.72 sec)
```

```
mysql> SELECT COUNT(*) FROM Directed_by;
+-----+
| COUNT(*) |
+-----+
|     6976 |
+-----+
1 row in set (0.09 sec)
```

Advanced Queries

Query 1:

Find the actors and actresses cast the top 15 largest number of movies that are released after 2008.

```
SELECT first_name, last_name, COUNT(show_id)
FROM Movie_TV NATURAL JOIN Casted_by NATURAL JOIN Cast
WHERE year_released >= 2008 AND category = "Movie"
GROUP BY cast_id
ORDER BY COUNT(show_id) DESC
LIMIT 15
```

result:

first_name	last_name	COUNT(show_id)
Julie	Tejwani	28
Rupa	Bhimani	27
Anupam	Kher	24
Rajesh	Kava	21
Nawazuddin	Siddiqui	19
Jigna	Bhardwaj	19
Naseeruddin	Shah	18
Akshay	Kumar	18
Boman	Irani	18
James	Franco	17
Andrea	Libman	17
Fred	Tatasciore	16
Om	Puri	16
Rajesh	Sharma	16
David	Spade	16

Query 2:

Find how many different ratings of American movies and shows each director have

```
SELECT first_name, last_name, rating, COUNT(show_id)
FROM Movie_TV NATURAL JOIN Directed_by db NATURAL JOIN Director d
WHERE country = "United States"
GROUP BY rating, d.director_id
ORDER BY COUNT(show_id) DESC
LIMIT 15
```

result:

first_name	last_name	rating	COUNT(show_id)
Marcus	Raboy	TV-MA	12
Jay	Karas	TV-MA	11
Jay	Chapman	TV-MA	10
Shannon	Hartman	TV-MA	8
Lance	Bangs	TV-MA	7
Ryan	Polito	TV-MA	7
William	Lau	TV-Y7	6
Martin	Scorsese	R	6
Robert	Rodriguez	PG	6
Quentin	Tarantino	R	5
John G.	Avildsen	PG	4
Mike	Gunther	R	4
Kevin	Smith	R	4
Brian	Levant	PG	4
Paul Thomas	Anderson	R	4

INDEX DESIGN & PERFORMANCE

Query 1:

Default Index: cast_id, show_id:

```
mysql> EXPLAIN ANALYZE (SELECT first_name last_name COUNT(show_id)
-> FROM MOVIE_TV NATURAL JOIN Casted by NATURAL JOIN Cast
-> WHERE year released >= 2008 AND category = "Movie"
-> GROUP BY cast_id
-> ORDER BY COUNT(show_id) DESC);
```

```
|
```

```
+-----+
| |
+-----+
```

```
|
```

```
Sort: `COUNT(show_id)` DESC (actual time=217.654..220.198 rows=18426 loops=1)
```

```
-> Stream results (cost=26475.08 rows=2794) (actual time=0.126..208.683 rows=18426 loops=1)
```

```
-> Group aggregate: count(Movie_TV.show_id) (cost=26475.08 rows=2794) (actual time=0.123..197.357 rows=18426 loops=1)
```

```
-> Nested loop inner join (cost=26195.68 rows=2794) (actual time=0.078..181.047 rows=30139 loops=1)
```

```
-> Nested loop inner join (cost=23217.80 rows=2794) (actual time=0.069..130.449 rows=30139 loops=1)
```

```
-> Index scan on Casted by using PRIMARY (cost=5660.15 rows=55879) (actual time=0.046..17.967 rows=53053 loops=1)
```

```
-> Filter: ((Movie_TV.category = 'Movie') and (Movie_TV.year released >= 2008)) (cost=0.25 rows=0) (actual time=0.002..0.002 rows=1 loops=53053)
```

```
-> Single-row index lookup on Movie_TV using PRIMARY (show_id=Casted by.show_id) (cost=0.25 rows=1) (actual time=0.001..0.001 rows=1 loops=53053)
```

```
-> Single-row index lookup on Cast using PRIMARY (cast_id=Casted by.cast_id) (cost=0.25 rows=1) (actual time=0.001..0.001 rows=1 loops=30139)
```

```
|
```

```
+-----+
```

```
1 row in set (0.23 sec)
```

Index on (first_name, last_name) of Cast:

```
mysql> CREATE INDEX name_idx ON Cast (first_name, last_name);  
Query OK, 0 rows affected (0.35 sec)  
Records: 0 Duplicates: 0 Warnings: 0  
  
mysql> EXPLAIN ANALYZE (SELECT first_name, last_name, COUNT(show_id) FROM Movie_TV NATURAL JOIN Casted_by NATURAL JOIN Cast WHERE year_released >= 2008 AND category = "Movie" GROUP BY cast_id ORDER BY COUNT(show_id) DESC);  
+-----+  
| EXPLAIN |  
+-----+  
|  
+-----+  
| Sort: 'COUNT(show_id)' DESC (actual time=218.826..221.948 rows=18426 loops=1)|  
-> Stream results (cost=26475.08 rows=2794) (actual time=0.076..209.609 rows=18426 loops=1)|  
-> Group aggregate: count(Movie_TV.show_id) (cost=26475.08 rows=2794) (actual time=0.074..197.828 rows=18426 loops=1)|  
-> Nested loop inner join (cost=26135.68 rows=2794) (actual time=0.060..181.357 rows=30139 loops=1)|  
-> Nested loop inner join (cost=25217.80 rows=2794) (actual time=0.063..130.577 rows=30139 loops=1)|  
-> Index scan on Casted by using PRIMARY (cost=5660.15 rows=55979) (actual time=0.032..17.725 rows=53053 loops=1)|  
-> Filter: ((Movie_TV.category = 'Movie') and (Movie_TV.year_released >= 2008)) (cost=0.25 rows=0) (actual time=0.002..0.002 rows=1 loops=53053)|  
-> Single-row index lookup on Movie_TV using PRIMARY (show_id=Casted_by.show_id) (cost=0.25 rows=1) (actual time=0.001..0.001 rows=1 loops=53053)|  
-> Single-row index lookup on Cast using PRIMARY (cast_id=Casted_by.cast_id) (cost=0.25 rows=1) (actual time=0.001..0.001 rows=1 loops=30139)|  
+-----+
```

The default index of entity set Cast is cast_id, which is the combination of first name and last name without white space. In this case, we are actually indexing on a single word, which is a precise while time-consuming and space-consuming way of storage. We choose to use both

Index on year_released of Movie_TV:

We try `year_released` as an additional index here because it is used in the `WHERE` statement in the query, and thus we are expecting an increase of the search speed compared with using other attributes as index. However, we find the result is almost the same as our original expectation, and it's probably because there are not enough number of `year_released`, resulting in the inefficiency of the searching.

```
mysql> CREATE INDEX category_id ON Movie_TV (category);
Query OK, 0 rows affected (0.13 sec)
Records: 0 Duplicates: 0 Warnings: 0

mysql> EXPLAIN ANALYZE (SELECT first_name, last_name, COUNT(show_id) FROM Movie_TV NATURAL JOIN Casted_by NATURAL JOIN Cast WHERE year_released >= 2008 AND category = "Movie" GROUP BY cast_id ORDER BY COUNT(show_id) DESC);
+-----+
| EXPLAIN |
+-----+
|
|
+-----+
|> Sort: 'COUNT(show_id)' DESC (actual time=227.528..230.063 rows=18426 loops=1)
|>   -> Stream results (cost=91081.54 rows=13031) (actual time=0.080..217.598 rows=18426 loops=1)
|>     -> Group aggregate: count(Movie_TV.show_id) (cost=91081.54 rows=13031) (actual time=0.078..206.464 rows=18426 loops=1)
|>       -> Nested loop inner join (cost=29778.49 rows=13031) (actual time=0.063..189.331 rows=30139 loops=1)
|>         -> Nested loop inner join (cost=63217.80 rows=13031) (actual time=0.056..135.932 rows=30139 loops=1)
|>           -> Index scan on Casted_by using PRIMARY (cost=5660.15 rows=58873) (actual time=0.040..18.306 rows=53053 loops=1)
|>             -> Filter: ((Movie_TV.category = 'Movie') AND (Movie_TV.year_released >= 2008)) (cost=0.25 rows=0) (actual time=0.002..0.002 rows=1 loops=53053)
|>               -> Single-row index lookup on Movie_TV using PRIMARY (show_id=Casted_by.show_id) (cost=0.25 rows=1) (actual time=0.001..0.001 rows=1 loops=53053)
|>                 -> Single-row index lookup on Cast using PRIMARY (cast_id=Casted_by.cast_id) (cost=0.25 rows=1) (actual time=0.001..0.001 rows=1 loops=30139)
|>
|
+-----+
1 row in set (0.24 sec)
```

Similar to the last try, we want to use category as the additional index here because it is used in the WHERE statement in the query. Unfortunately, we still find the result is almost the same as our original expectation, and what's worse is there are only two elements in our category column, so the search speed is still a loss.

Query 2:

Default Index: director_id, show_id:

[illegible]

Index on (first_name, last_name) of Director:

[illegible]

The default index of entity set Director is director_id, which is the combination of first name and last name without white space. In this case, we are actually indexing on a single word, which is a precise while time-consuming and space-consuming way of storage. We choose to use both first and last name as indexes so that the search can utilize search in both groups. For example, if first name starts with “a” and last name starts with “b”, we can

Index on rating of Movie_TV:

We try to use rating in Movie_TV table as an index because we group the data by rating in our query. By using rating as an index, we think that the speed of searching different ratings can be improved. The result shows that the time increased a little bit. However, we are not sure if this small increase is caused by the new index. The increase is really small, maybe because there are not enough number of different ratings, resulting in the inefficiency of the searching.

```
mysql> CREATE INDEX country_idx ON Movie_TV (country);
Query OK, 0 rows affected (0.12 sec)
Records: 0 Duplicates: 0 Warnings: 0

mysql> EXPLAIN ANALYZE (SELECT first_name, last_name, rating, COUNT(show_id) FROM Movie_TV NATURAL JOIN Directed by db NATURAL JOIN Director d WHERE country = "United States" GROUP BY rating, d.director_id ORDER BY COUNT(show_id) DESC);
```

```
+-----+
|EXPLAIN|
+-----+
|       |
+-----+
|<br>-> Sort: `COUNT(show_id)` DESC (actual time=23.811..24.087 rows=2000 loops=1)|
|<br>-> Table scan on <Temporary> (actual time=0.003..0.329 rows=2000 loops=1)|
|<br>-> Aggregate using temporary table (actual time=22.847..23.337 rows=2000 loops=1)|
|<br><br>-> Nested loop inner join (cost=2503.56 rows=3185) (actual time=0.104..18.006 rows=2394 loops=1)|
|<br><br><br>-> Nested loop inner join (cost=1388.96 rows=3185) (actual time=0.097..12.109 rows=2394 loops=1)|
|<br><br><br><br>-> Index lookup on Movie TV using country_idx (country='United States') (cost=354.55 rows=2818) (actual time=0.088..5.231 rows=2818 loops=1)|
|<br><br><br><br>-> Index lookup on db using show_id (show_id=Movie_TV.show_id) (cost=0.25 rows=1) (actual time=0.002..0.002 rows=1 loops=2818)|
|<br><br><br><br>-> Single-row index lookup on d using PRIMARY (director_id=db.director_id) (cost=0.25 rows=1) (actual time=0.002..0.002 rows=1 loops=2394)|
|<br><br><br><br>|
|<br><br><br><br>+-----+
```

```
| row in set (0.03 sec)
```

We try to use country in Movie_TV table as an index because we try to filter the data by "United States" in our query. By using country as an index, we think that the speed of searching different entries can be improved. The results fluctuate around by only an extremely small amount and have no significant difference comparing to using default index.

Reason of inefficient index:

We suspect there are couple factors that are affecting how the indexes are influencing the query performance. First, we don't have a lot of unique values in columns such as category. When we search by one category, almost half our dataset isn't being selected and it's not big enough to make a big impact on the query performance. That's why we see the extremely small amount of change on the time. We also suspect that most of the search in our query uses primary keys which are set to default index for search. Even if we add more index on top of the default indexes, there won't be a big improvement.