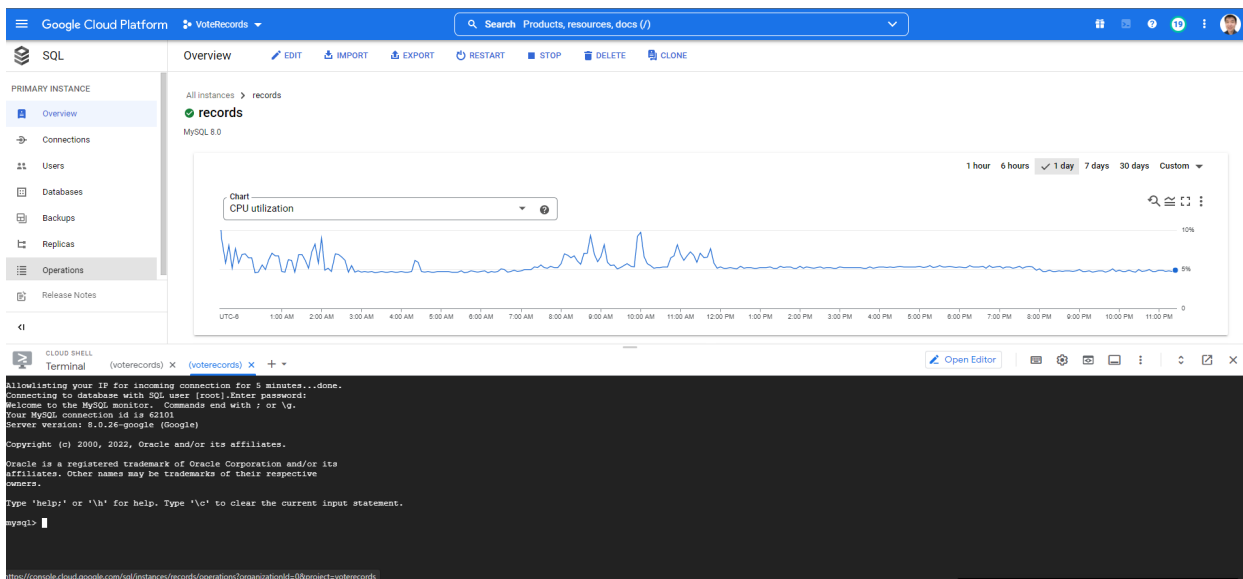


2.1 Screenshot of connection to GCP



2.2 DDL Commands

```
CREATE TABLE Senators(  
    SenatorID INTEGER NOT NULL,  
    Name VARCHAR(255) NOT NULL,  
    BirthYear INTEGER,  
    PRIMARY KEY (SenatorID)  
);  
CREATE TABLE Wikipedia(  
    PageTitle VARCHAR(255) NOT NULL,  
    PageURL VARCHAR(1024),  
    PRIMARY KEY (PageTitle)  
);  
CREATE TABLE Bills(  
    BillID INTEGER NOT NULL,  
    Date VARCHAR(255),  
    Results VARCHAR(255),  
    Description VARCHAR(1024),  
    PRIMARY KEY (BillID)  
);  
CREATE TABLE Parties(  
    PartyName VARCHAR(255) NOT NULL,  
    YearFounded INTEGER,  
    PRIMARY KEY (PartyName)  
);  
CREATE TABLE States(  
    StateID CHAR(2) NOT NULL,  
    StateName VARCHAR(32),  
    DominantParty VARCHAR(255),  
    PRIMARY KEY (StateID)  
);
```

```

CREATE TABLE Vote(
  SenatorID INTEGER NOT NULL,
  BillID INTEGER NOT NULL,
  VoteType VARCHAR(32),
  PRIMARY KEY (SenatorID, BillID),
  FOREIGN KEY (SenatorID) REFERENCES Senators(SenatorID),
  FOREIGN KEY (BillID) REFERENCES Bills(BillID)
);
CREATE TABLE AffiliatedTo(
  SenatorID INTEGER NOT NULL,
  PartyName VARCHAR(255) NOT NULL,
  PRIMARY KEY (SenatorID, PartyName),
  FOREIGN KEY (SenatorID) REFERENCES Senators(SenatorID),
  FOREIGN KEY (PartyName) REFERENCES Parties(PartyName)
);
CREATE TABLE FromState(
  SenatorID INTEGER NOT NULL,
  StateID CHAR(2) NOT NULL,
  PRIMARY KEY (SenatorID, StateID),
  FOREIGN KEY (SenatorID) REFERENCES Senators(SenatorID),
  FOREIGN KEY (StateID) REFERENCES States(StateID)
);
CREATE TABLE LooksLike(
  SenatorID INTEGER,
  PageTitle VARCHAR(255),
  PRIMARY KEY (SenatorID, PageTitle),
  FOREIGN KEY (SenatorID) REFERENCES Senators(SenatorID),
  FOREIGN KEY (PageTitle) REFERENCES Wikipedia(PageTitle)
);

```

2.3 1000 rows on four tables

```

Database changed
mysql> SELECT COUNT(SenatorID)
-> FROM Senators;
+-----+
| COUNT(SenatorID) |
+-----+
|          1016 |
+-----+
1 row in set (0.02 sec)

mysql> SELECT COUNT(PageTitle)
-> FROM Wikipedia;
+-----+
| COUNT(PageTitle) |
+-----+
|          1016 |
+-----+
1 row in set (0.02 sec)

```

```

mysql> SELECT COUNT(BillID)
-> FROM Bills;
+-----+
| COUNT(BillID) |
+-----+
|          1313 |
+-----+
1 row in set (0.01 sec)

mysql> SELECT COUNT(*)
-> FROM Vote;
+-----+
| COUNT(*) |
+-----+
|       58142 |
+-----+
1 row in set (0.02 sec)

mysql> 

```

3.1 SQL Query #1

```
SELECT BillID, COUNT(SenatorID) as YesCount, Results, Date
FROM Vote NATURAL JOIN Bills
WHERE VoteType = 1
GROUP BY BillID
HAVING YesCount > 50;
-- This query returns the BillID, number of "yea" votes, vote results, and date,
-- for all bills that has amajority yea votes and their information.
-- This helps us see the bills that were passed.
```

Top 15 rows

```
mysql> SELECT BillID, COUNT(SenatorID) as YesCount, Results, Date
-> FROM Vote NATURAL JOIN Bills
-> WHERE VoteType = 1
-> GROUP BY BillID
-> HAVING YesCount > 50 LIMIT 15;
```

BillID	YesCount	Results	Date
1173	85	Nomination Confirmed	1/20/2021
1174	70	Bill Passed	1/21/2021
1175	94	Nomination Confirmed	1/22/2021
1176	85	Nomination Confirmed	1/25/2021
1177	79	Nomination Confirmed	1/26/2021
1179	83	Resolution Agreed to	1/26/2021
11711	87	Nomination Confirmed	2/2/2021
11714	90	Amendment Agreed to	2/4/2021
11716	100	Amendment Agreed to	2/4/2021
11718	58	Amendment Agreed to	2/4/2021
11719	99	Amendment Agreed to	2/4/2021
11721	100	Amendment Agreed to	2/4/2021
11722	98	Amendment Agreed to	2/4/2021
11723	52	Motion Rejected	2/4/2021
11725	52	Motion Rejected	2/4/2021

15 rows in set (0.07 sec)

3.2 SQL Query #2

```
Select BillID, tmp.NoCount, Results, Date
FROM Bills NATURAL JOIN
  (Select BillID, COUNT(SenatorID) as NoCount
   FROM Vote NATURAL JOIN Bills
   WHERE SenatorId IN
     (SELECT SenatorID
      FROM AffiliatedTo NATURAL JOIN Parties
      WHERE PartyName = 'democrat')
   AND VoteType BETWEEN 4 AND 6
   GROUP BY BillID) as tmp
WHERE Results LIKE "%Agree%" OR Results LIKE "%Confirm%" OR Results LIKE "%Pass%"
ORDER BY tmp.NoCount DESC;
-- This query returns the BillID, number of "nay" votes, vote results, and date,
-- for all bills that were passed. Sorted in descending order by number of "nay" votes from democrats.
-- This helps us see what bills were really disliked by the democrats but were still passed.
```

Note: An explanation of VoteType is included in the proposal. (4: Announced Nay, 5: Paired Nay, 6: Nay)

Top 15 rows

```
--> GROUP BY BillID) as tmp
--> WHERE Results LIKE "%Agree%" OR Results LIKE "%Confirm%" OR Results LIKE "%Pass%"
--> ORDER BY tmp.NoCount DESC
--> LIMIT 15;
```

BillID	NoCount	Results	Date
11743	47	Amendment Agreed to	2/5/2021
11744	47	Amendment Agreed to	2/5/2021
11776	47	Amendment Agreed to	3/5/2021
117336	47	Amendment Agreed to	8/10/2021
117342	47	Amendment Agreed to	8/11/2021
11732	46	Amendment Agreed to	2/4/2021
117349	46	Amendment Agreed to	8/11/2021
117489	46	Joint Resolution Passed	12/8/2021
11748	45	Amendment Agreed to	2/5/2021
117332	45	Amendment Agreed to	8/10/2021
117346	45	Amendment Agreed to	8/11/2021
117330	44	Amendment Agreed to	8/10/2021
117351	44	Amendment Agreed to	8/11/2021
11728	41	Amendment Agreed to	2/4/2021
117323	41	Amendment Agreed to	8/10/2021

15 rows in set (0.01 sec)

4.1.a Indexing Analysis for Query #1

Performance Before Indexing

```
mysql> EXPLAIN ANALYZE
-> SELECT BillID, COUNT(SenatorID) as YesCount, Results, Date
-> FROM Vote NATURAL JOIN Bills
-> WHERE VoteType = 1
-> GROUP BY BillID
-> HAVING YesCount > 50;

+-----+
| EXPLAIN |
+-----+

+-----+
| -> Filter: (YesCount > 50) (actual time=81.506..81.624 rows=384 loops=1)
  -> Table scan on <temporary> (actual time=0.001..0.041 rows=581 loops=1)
    -> Aggregate using temporary table (actual time=81.503..81.577 rows=581 loops=1)
      -> Nested loop inner join (cost=10196.38 rows=7488) (actual time=0.058..60.316 rows=34696 loops=1)
        -> Filter: (Vote.VoteType = 1) (cost=7575.75 rows=7488) (actual time=0.045..21.876 rows=34696 loops=1)
          -> Table scan on Vote (cost=7575.75 rows=74875) (actual time=0.040..14.261 rows=58142 loops=1)
            -> Single-row index lookup on Bills using PRIMARY (BillID=Vote.BillID) (cost=0.25 rows=1) (actual time=0.001..0.001 rows=1 loops=34696)
        |
      +-----+
    +-----+
  +-----+
1 row in set (0.09 sec)
```

This is our Query #1; without indexing we get a baseline performance of 0.09 seconds.

1. CREATE INDEX idx1 ON Bills (BillID);

```
mysql> EXPLAIN ANALYZE
-> SELECT BillID, COUNT(SenatorID) as YesCount, Results, Date
-> FROM Vote NATURAL JOIN Bills
-> WHERE VoteType = 1
-> GROUP BY BillID
-> HAVING YesCount > 50;

+-----+
| EXPLAIN |
+-----+

+-----+
| -> Filter: (YesCount > 50) (actual time=82.596..82.718 rows=384 loops=1)
  -> Table scan on <temporary> (actual time=0.001..0.041 rows=581 loops=1)
    -> Aggregate using temporary table (actual time=82.593..82.671 rows=581 loops=1)
      -> Nested loop inner join (cost=10196.38 rows=7488) (actual time=0.045..60.788 rows=34696 loops=1)
        -> Filter: (Vote.VoteType = 1) (cost=7575.75 rows=7488) (actual time=0.035..21.859 rows=34696 loops=1)
          -> Table scan on Vote (cost=7575.75 rows=74875) (actual time=0.029..14.217 rows=58142 loops=1)
            -> Single-row index lookup on Bills using PRIMARY (BillID=Vote.BillID) (cost=0.25 rows=1) (actual time=0.001..0.001 rows=1 loops=34696)
        |
      +-----+
    +-----+
  +-----+
1 row in set (0.08 sec)
```

This is our Query #1, with indexing based on the id of a bill (BillID). We use this for our index because this is one of the attributes returned by our query, and there are hundreds of bills in our dataset. We believe using this index should lead to a faster performance, as there are a hundred senators voting for a given bill; we should therefore easily retrieve the votes for a given bill. However, using this index only slightly increased our performance down to 0.08 seconds. The reason for this could be because our default query is already fast with the data set that we are using.

2. CREATE INDEX idx2 ON Vote (SenatorID);

```
mysql> EXPLAIN ANALYZE
-> SELECT BillID, COUNT(SenatorID) as YesCount, Results, Date
-> FROM Vote NATURAL JOIN Bills
-> WHERE VoteType = 1
-> GROUP BY BillID
-> HAVING YesCount > 50;
+-----+
| EXPLAIN |
+-----+
-> Filter: (YesCount > 50) (actual time=82.415..82.531 rows=384 loops=1)
-> Table scan on <temporary> (actual time=0.001..0.041 rows=581 loops=1)
-> Aggregate using temporary table (actual time=82.411..82.484 rows=581 loops=1)
-> Nested loop inner join (cost=10196.38 rows=7488) (actual time=0.070..60.858 rows=34696 loops=1)
-> Filter: (Vote.VoteType = 1) (cost=7575.75 rows=7488) (actual time=0.056..22.194 rows=34696 loops=1)
-> Table scan on Vote (cost=7575.75 rows=74875) (actual time=0.049..14.334 rows=58142 loops=1)
-> Single-row index lookup on Bills using PRIMARY (BillID=Vote.BillID) (cost=0.25 rows=1) (actual time=0.001..0.001 rows=1 loops=34696)
+-----+
1 row in set (0.08 sec)
```

This is our Query #1, with indexing based on the id of a senator (SenatorID). We use this for our index because this is another one of the attributes returned by our query, and there are a hundred senators. We believe using this index should lead to a faster performance, as there are a hundred senators voting for a given bill; we should therefore easily retrieve the votes for a given bill. However, using this index only slightly increased our performance down to 0.08 seconds again. The reason for this could be because our default query is already fast with the data set that we are using.

3. CREATE INDEX idx3 ON Vote (SenatorID, BillID);

```
mysql> CREATE INDEX idx3 ON Vote (SenatorID, BillID);
Query OK, 0 rows affected (0.21 sec)
Records: 0 Duplicates: 0 Warnings: 0

mysql> EXPLAIN ANALYZE
-> SELECT BillID, COUNT(SenatorID) as YesCount, Results, Date
-> FROM Vote NATURAL JOIN Bills
-> WHERE VoteType = 1
-> GROUP BY BillID
-> HAVING YesCount > 50;
+-----+
| EXPLAIN |
+-----+
-> Filter: (YesCount > 50) (actual time=81.242..81.360 rows=384 loops=1)
-> Table scan on <temporary> (actual time=0.001..0.043 rows=581 loops=1)
-> Aggregate using temporary table (actual time=81.238..81.314 rows=581 loops=1)
-> Nested loop inner join (cost=10196.38 rows=7488) (actual time=0.060..60.528 rows=34696 loops=1)
-> Filter: (Vote.VoteType = 1) (cost=7575.75 rows=7488) (actual time=0.047..21.867 rows=34696 loops=1)
-> Table scan on Vote (cost=7575.75 rows=74875) (actual time=0.041..14.334 rows=58142 loops=1)
-> Single-row index lookup on Bills using PRIMARY (BillID=Vote.BillID) (cost=0.25 rows=1) (actual time=0.001..0.001 rows=1 loops=34696)
+-----+
1 row in set (0.08 sec)
```

This is our Query #1, with indexing based on the id of a senator (SenatorID) and the id of the bill (BillID). We use this for our index because both of these attributes are returned by our query, there are a hundred senators, and hundreds of bills. We believe using this index should lead to a faster performance, as there are a hundred senators voting for a given bill; we should therefore easily retrieve the votes for a given bill. However, using this index only slightly increased our performance down to 0.08 seconds again. The reason for this could be because our default query is already fast with the data set that we are using.

Our final recommendation based on the above information is using an index based on both SenatorId and BillId. The reason for this is because it leads to a (albeit slight) performance increase versus without indexing. Our query tries to retrieve bills based on the bill id and the count of the senator ids'. Therefore, using both of these for our indexing makes the most sense.

Performance Before Indexing

```
1. CREATE INDEX idx1 ON Bills (BillID);
```

```
2. CREATE INDEX idx2 ON Bills (BillID, Results);
```

This is our Query #2, with indexing based on the id of the Bill (Bill ID). We use Bill Id for our index because in the code, we group by Bill ID. So we believe this should sort the Bill Id for us and reduce the amount of time to search. However, the runtime does not change at all. The reason for this might be because even though we sort the Bill Id because we use group by, what we really need is the SenatorID in that table created by groupby. So the runtime stay the same. Another possibility is that the performance was actually improved by a small amount, but the new runtime still got rounded up to 0.02 seconds, therefore it may seem like the performance did not change. This is possible since our query is already quite fast, and a 25% decrease in runtime may still not be visible due to rounding.

This is our Query #2, with indexing based on the id of the Bill (Bill ID) and the id of the senator from Vote. The reason why we choose BillID and SenatorID is because we believe that it would run faster if we sort the BillID and SenatorID for the subquery using the Vote table because we want to choose those from thousands of Bill ID and SenatorID. However, the runtime does not change. Since we have to search thousands of results and there is only one output, it would explain that indexing might not help decrease the runtime in this case. Another possibility is that the performance was actually improved by a small amount, but the new runtime still got rounded up to 0.02 seconds, therefore it may seem like the performance did not change. This is possible since our query is already quite fast, and a 25% decrease in runtime may still not be visible due to rounding.

We preferred the second indexing design, since theoretically it would improve the performance of the subquery. However, due to the fact that there are only several different results possible, in case of a large database, the performance increase is not significant. Therefore, we decided to use the default indexing because the runtime does not change by a noticeable amount, and this minimal performance improvement does not justify the resources it takes to implement this indexing design.