# CS411 Stage3 Improvement

## Team043-InfelPhira

Original screenshots of the explain analyze commands without creating any index.

First Analysis:



Second Analysis:

Change Second Analysis:

```
EXPLAIN ANALYZE
SELECT DISTINCT title
FROM Videos
JOIN BunchContain ON Videos.videoId = BunchContain.videoId
WHERE bunchId IN (
    SELECT bunchId
    FROM BunchContain
    GROUP BY bunchId
    HAVING COUNT(videoId) >= ALL(
        SELECT COUNT(videoId)
        FROM BunchContain
        GROUP BY bunchId
    )
)
AND categoryId >= 1
AND publishedAt >= '1970-01-01'
LIMIT 15;
CREATE INDEX idx_categoryId ON Videos (categoryId);
```

```sql
CREATE INDEX idx_publishedAt ON Videos (publishedAt);
CREATE INDEX idx_category_publishedAt ON Videos (categoryId, publishedAt);
```

Original:

CREATE INDEX idx_categoryId on Videos(categoryId);
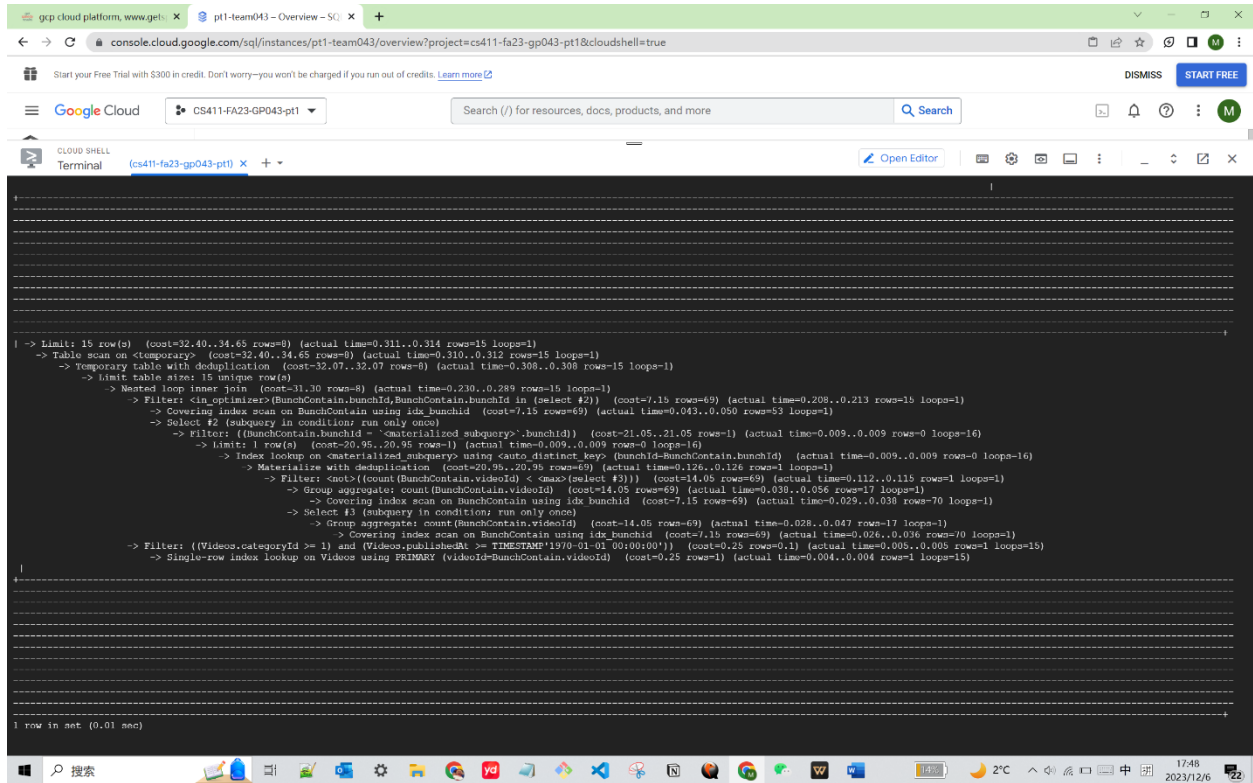


CREATE INDEX idx_publishedAt on Videos(publishedAt);

CREATE INDEX idx_category_publishedAt on Videos(categoryId, publishedAt);

**Analyses:**

**idx_categoryId on Videos(categoryId):**

This index is beneficial for the categoryId >= 1 condition. It allows the database engine to efficiently retrieve rows where the categoryId is greater than or equal to 1. There are a large number of distinct categoryId values and queries often filter based on this column, the index is appropriate

**idx_publishedAt on Videos(publishedAt):**

This index supports the publishedAt >= '1970-01-01' condition. It allows for efficient retrieval of videos published after a specific date. Since queries frequently filter videos based on the publication date, this index is beneficial.

**idx_category_publishedAt on Videos(categoryId, publishedAt):**

This composite index covers both the categoryId and publishedAt conditions used in the WHERE clause. It is helpful when filtering based on both categoryId and publishedAt, potentially improving the performance of queries with combined conditions.

**Conclusion:**

Since our biggest bunch contains data that is not enough(only hundreds), the time it executed had little differences. idx_category_publishedAt on Videos(categoryId, publishedAt) should be good when filtering based on both categoryId and publishedAt.