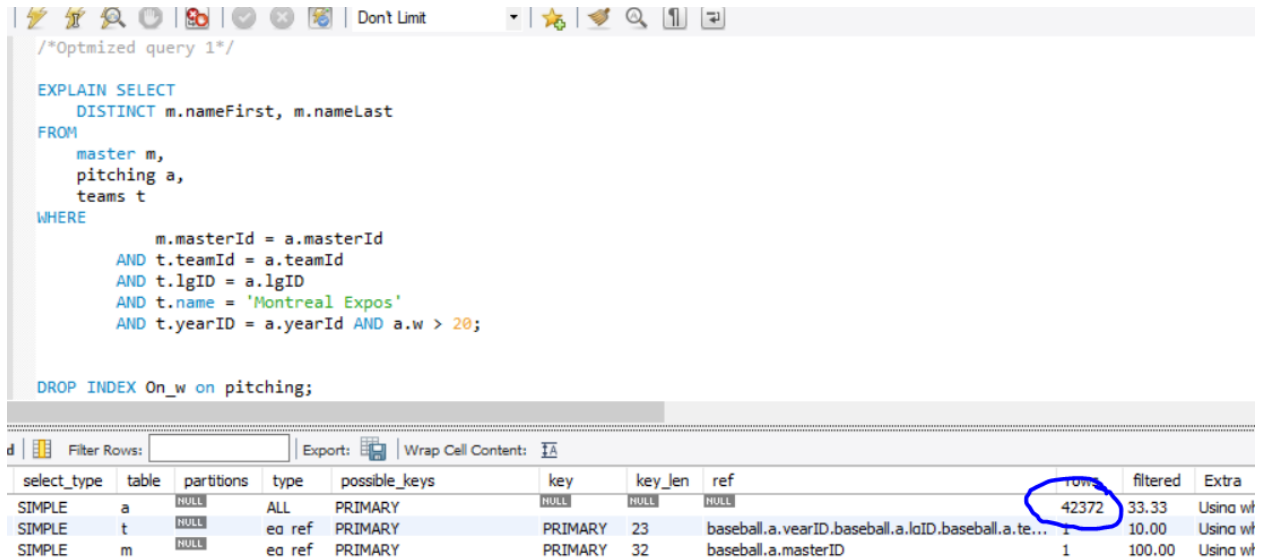


## Query 1:

I executed the given query and used the explain keyword to check the number of rows that were being read. There were 42372 rows that were being read.



```
/*Optimized query 1*/  
EXPLAIN SELECT  
  DISTINCT m.nameFirst, m.nameLast  
FROM  
  master m,  
  pitching a,  
  teams t  
WHERE  
  m.masterId = a.masterId  
  AND t.teamId = a.teamId  
  AND t.lgID = a.lgID  
  AND t.name = 'Montreal Expos'  
  AND t.yearID = a.yearID AND a.w > 20;  
  
DROP INDEX On_w on pitching;
```

select_type	table	partitions	type	possible_keys	key	key_len	ref	rows	filtered	Extra
SIMPLE	a	NONE	ALL	PRIMARY	NONE	NONE	NONE	42372	33.33	Using where
SIMPLE	t	NONE	eq ref	PRIMARY	PRIMARY	23	baseball.a.yearID,baseball.a.lgID,baseball.a.te...	1	10.00	Using where
SIMPLE	m	NONE	eq ref	PRIMARY	PRIMARY	32	baseball.a.masterID	1	100.00	Using where

Firstly, I removed the like keyword and replaced it with =, which improves the efficiency of finding matching rows by a very slight margin.

I assessed the conditions under the where clause and I see that we're trying to filter out rows using the a.w column. We can apply indexing on the "w" column of the pitching table, which can help us get lesser number of rows from the pitching table, which are then used to join with master and teams table.

I create an index using

```
CREATE INDEX Index_W on pitching(w);
```

After creating this index the number of rows we get from the pitching table decreases to 906 rows which means we have to read through almost 50 times lesser number of rows.

```

1 /*Optimized query 1*/
2
3 EXPLAIN SELECT
4 DISTINCT m.nameFirst, m.nameLast
5 FROM
6 master m,
7 pitching a,
8 teams t
9 WHERE
10 m.masterId = a.masterId
11 AND t.teamId = a.teamId
12 AND t.lgID = a.lgID
13 AND t.name = 'Montreal Expos'
14 AND t.yearID = a.yearID AND a.w > 20;
15
16
17 CREATE INDEX Index_W on pitching(w);

```

id	select_type	table	partitions	type	possible_keys	key	key_len	ref	rows	filtered	Extra
1	SIMPLE	a	INDEX	range	PRIMARY, Index W	Index W	5	INDEX	906	100.00	Using index condition: Using where
1	SIMPLE	t	INDEX	eq ref	PRIMARY	PRIMARY	23	baseball.a.yearID,baseball.a.lgID,baseball.a.te...	1	10.00	Using where
1	SIMPLE	m	INDEX	eq ref	PRIMARY	PRIMARY	32	baseball.a.masterID	1	100.00	Using where

## Query 2

I executed the below query using the explain function and it can be seen that schoolID from schoolsPlayers is what is being used to retrieve data for a school player.

Now, here if we assess the query, there can be multiple number of schoolIDs for a player belonging to Utah state university. In this case, once we retrieve the schoolID of Utah State University, 6147 rows are scanned to find out all the players that belong to Utah state university.

If we think about how we can optimize this, we could create an index on schoolID and schoolPlayers ID.

If we do that, then we would have much lesser rows to go through to find out schoolPlayers ID belonging to Utah State.

```

1 explain SELECT
2   h / ab as Average, h as Hits, ab as 'At Bats', nameFirst as 'First Name', nameLast as 'Last Name', batting.yearID as Year
3 FROM
4   batting,
5   master
6 WHERE
7   ab is not null
8   and batting.masterID = master.masterID
9   AND master.masterID IN (SELECT
10     masterID
11   FROM
12     schoolsplayers
13   WHERE
14     schoolID in (SELECT
15       schoolID
16     FROM
17       schools
18     WHERE
19       schoolName like '%Utah State%')
20 )) order by year;

```

	id	select_type	table	partitions	type	possible_keys	key	key_len	ref	rows	filtered	Extra
	1	SIMPLE	schoolsplayers	NULL	index	PRIMARY	PRIMARY	76	NULL	6147	94.84	Using index; Using temporary; Using where
	1	SIMPLE	schools	NULL	eq ref	PRIMARY	PRIMARY	47	baseball.schoolsplayers.schoolID	1	11.11	Using where; FirstMatch(schools)
	1	SIMPLE	master	NULL	eq ref	PRIMARY	PRIMARY	32	baseball.schoolsplayers.masterID	1	100.00	Using where
	1	SIMPLE	batting	NULL	ref	PRIMARY	PRIMARY	29	baseball.master.masterID	5	90.00	Using where

First I create an index on schoolID and test if there is an improvement in the number of rows we get.

CREATE INDEX Index\_On\_SchoolID on schoolsPlayers(schoolID);

	id	select_type	table	partitions	type	possible_keys	key	key_len	ref	rows	filtered	Extra
	1	SIMPLE	<subquery2>	NULL	ALL	NULL	NULL	NULL	NULL	1	100.00	Using temporary; Using where
	1	SIMPLE	master	NULL	eq ref	PRIMARY	PRIMARY	32	<subquery2>.masterID	1	100.00	Using where
	1	SIMPLE	batting	NULL	ref	PRIMARY	PRIMARY	29	baseball.master.masterID	5	90.00	Using where
	2	MATERIALIZED	schools	NULL	ALL	PRIMARY	PRIMARY	47	baseball.schools.schoolID	749	11.11	Using where
	2	MATERIALIZED	schoolsplayers	NULL	ref	PRIMARY, Index On Plav...	Index On SchoolID	47	baseball.schools.schoolID	8	100.00	Using index

The rows count has now been decreased to 8 because we indexed the school players by their indexes, which makes it easier to find school players now.

We could also apply another index to the schools table on schoolName column to reduce the current rows count 749.

CREATE INDEX Index\_On\_SchoolName on schools(schoolName);

```

explain SELECT
  h / ab as Average, h as Hits, ab as 'At Bats', nameFirst as 'First Name', nameLast as 'Last Name', batting.yearID as Year
FROM
  batting,
  master
WHERE
  ab is not null
  and batting.masterID = master.masterID
  AND master.masterID IN (SELECT
    masterID
  FROM
    schoolsplayers
  WHERE
    schoolID in (SELECT
      schoolID
    FROM
      schools
    WHERE
      schoolName like '%Utah State%')
  )) order by year;

```

id	select_type	table	partitions	type	possible_keys	key	key_len	ref	rows	filtered	Extra
1	SIMPLE	<subquery2>	NULL	ALL	NULL	NULL	NULL	NULL	1	100.00	Using temporary; Using filesort
2	SIMPLE	master	NULL	eq ref	PRIMARY	PRIMARY	32	<subquery2>.masterID	1	100.00	Using where
3	SIMPLE	batting	NULL	ref	PRIMARY	PRIMARY	29	baseball.master.masterID	5	90.00	Using where
4	MATERIALIZED	schools	NULL	ref	PRIMARY.Index On SchoolName	Index On SchoolName	768	const	1	100.00	Using index
5	MATERIALIZED	schoolsplayers	NULL	ref	PRIMARY.Index On PlayerID	Index On SchoolID	47	baseball.schools.schoolID	8	100.00	Using index

If you notice, it is still the same. That is because we have used schoolName like '%Utah State%'

If we know the name of the university we can improve the performance of the query by replacing it by schoolName='Utah State University' and below will be the result-

```

1 explain SELECT
2   h / ab as Average, h as Hits, ab as 'At Bats', nameFirst as 'First Name', nameLast as 'Last Name', batting.yearID as Year
3 FROM
4   batting,
5   master
6 WHERE
7   ab is not null
8   and batting.masterID = master.masterID
9   AND master.masterID IN (SELECT
10     masterID
11   FROM
12     schoolsplayers
13   WHERE
14     schoolID in (SELECT
15       schoolID
16     FROM
17       schools
18     WHERE
19       schoolName = 'Utah State University')
20   )) order by year;

```

id	select_type	table	partitions	type	possible_keys	key	key_len	ref	rows	filtered	Extra
1	SIMPLE	<subquery2>	NULL	ALL	NULL	NULL	NULL	NULL	1	100.00	Using temporary; Using filesort
1	SIMPLE	master	NULL	eq ref	PRIMARY	PRIMARY	32	<subquery2>.masterID	1	100.00	Using where
2	SIMPLE	batting	NULL	ref	PRIMARY	PRIMARY	29	baseball.master.masterID	5	90.00	Using where
1	MATERIALIZED	schools	NULL	ref	PRIMARY.Index On SchoolName	Index On SchoolName	768	const	1	100.00	Using index
2	MATERIALIZED	schoolsplayers	NULL	ref	PRIMARY.Index On PlayerID	Index On SchoolID	47	baseball.schools.schoolID	8	100.00	Using index

### Query 3

The below query is quite slow uses block nested loop to run. The appearances table is being read in a nested way for the tables jeter and jeterTY and both tables are being read for around 97k rows,

97000\*97000 amounts to a huge amount of rows and that is why this query ends up taking about 6 seconds to execute.

```
Explain SELECT distinct jeter.masterID, jeterT.masterID, jeterTY.masterID, jeterTT.masterID
FROM
  master m,
  appearances jeter,
  appearances jeterT,
  appearances jeterTY,
  appearances jeterTT
WHERE
  m.masterID = jeter.masterID
  AND m.nameLast = 'Jeter'
  AND m.nameFirst = 'Derek'
  AND jeter.teamID = jeterT.teamID AND jeter.yearID = jeterT.yearID AND jeter.lgID = jeterT.lgID AND jeter.masterID <> jeterT.masterID
  AND jeterT.masterID = jeterTY.masterID
  AND jeterTY.teamID = jeterTT.teamID and jeterTY.yearID = jeterTT.yearID AND jeterTY.lgID = jeterTT.lgID AND jeterTY.masterID <> jeterTT.masterID
  AND jeterTT.masterID <> jeter.masterID
  AND jeter.teamID <> jeterTY.teamID
```

select_type	table	partitions	type	possible_keys	key	key_len	ref	rows	filtered	Extra
SIMPLE	jeter	HULL	ALL	PRIMARY	HULL	HULL	HULL	97811	100.00	Using temporary
SIMPLE	m	HULL	eq ref	PRIMARY	PRIMARY	32	baseball.jeter.mast...	1	5.00	Using where
SIMPLE	jeterTY	HULL	ALL	PRIMARY	HULL	HULL	HULL	97811	90.00	Using where: Using join buffer (Block Nested Lo...
SIMPLE	jeterT	HULL	eq ref	PRIMARY	PRIMARY	44	baseball.jeter.vearI...	1	10.00	Using where
SIMPLE	jeterTT	HULL	ref	PRIMARY	PRIMARY	15	baseball.jeterTY.ve...	36	8.10	Using where

We can create an index and the masterID column in the appearances table as it is being heavily used to join the jeter and jeterTY table and improve the cost.

```
CREATE INDEX masterID_Index ON appearances(masterID);
```

The cost improvement after doing this is pretty good-

id	select_type	table	partitions	type	possible_keys	key	key_len	ref	rows	filtered	Extra
1	SIMPLE	m	HULL	ALL	PRIMARY	HULL	HULL	HULL	18224	1.00	Using where: Using temporary
1	SIMPLE	jeter	HULL	ref	PRIMARY.masterID Index	masterID Index	29	baseball.m.masterID	5	100.00	Using index condition
1	SIMPLE	jeterT	HULL	ref	PRIMARY.masterID Index	PRIMARY	15	baseball.jeter.vearID.baseball...	36	9.00	Using where
1	SIMPLE	jeterTY	HULL	ref	PRIMARY.masterID Index	masterID Index	29	baseball.jeterT.masterID	5	90.00	Using index condition
1	SIMPLE	jeterTT	HULL	ref	PRIMARY	PRIMARY	15	baseball.jeterTY.vearID.baseb...	36	8.10	Using where

Now that we have improved the cost for reading rows from the appearances table. We can notice that the master table(m) has quite a lot of cost we can improve on.

We can apply another index on the nameFirst and nameLast column as it is being used in the where clause to find the player named Jason Derek in query.

```
CREATE INDEX name_Index ON master(nameFirst,nameLast);
```

Result Grid											
Filter Rows: [ ] Export: [ ] Wrap Cell Contents: [ ]											
id	select_type	table	partitions	type	possible_keys	key	key_len	ref	rows	filtered	Extra
1	SIMPLE	m	NONE	ref	PRIMARY, name Index	name Index	306	const,const	1	100.00	Using index; Using temporary
1	SIMPLE	teier	NONE	ref	PRIMARY, masterID Index	masterID Index	29	baseball.m.masterID	5	100.00	Using index condition
1	SIMPLE	teierT	NONE	ref	PRIMARY, masterID Index	PRIMARY	15	baseball.teierT.yearID,baseball...	36	9.00	Using where
1	SIMPLE	teierTY	NONE	ref	PRIMARY, masterID Index	masterID Index	29	baseball.teierTY.masterID	5	90.00	Using index condition
1	SIMPLE	teierTT	NONE	ref	PRIMARY	PRIMARY	15	baseball.teierTY.yearID,baseb...	36	8.10	Using where

This results in significant improvement. As we can see, the number of rows it has to read now is way lesser than we had earlier.

Also, the query after this optimization runs in 1.1 seconds as compared to earlier which was ~6.5 seconds.

## Query 4

The below query cost comes out to be 2745 rows for the team table. Here I can try the following approaches-

1. Apply an index on W since it is being used in the where clause
2. Apply a composite index on yearID, W since we first use yearID to find a max W then we use W in the query.
3. Try out a join instead of a subquery and then apply a composite index on W, yearid.

Result Grid											
Filter Rows: [ ] Export: [ ] Wrap Cell Contents: [ ]											
id	select_type	table	partitions	type	possible_keys	key	key_len	ref	rows	filtered	Extra
1	PRIMARY	T	NONE	ALL	NONE	NONE	NONE	NONE	2745	100.00	Using where
2	DEPENDENT	v	NONE	ref	PRIMARY	PRIMARY	4	baseba...	19	100.00	Using index

Approach 1 – This didn't work out because it still has to go through all the rows in the subquery to find out a valid yearID, and then it finds the max. The index did not improve any performance here.

Approach 2 – This slightly improved the performance. The performance improved by 40% on average.

Approach 3 – This significantly improved the performance. The performance differences are as below

Original Query-

Automatic context help is disabled. Use the toolbar to manually get help for the current caret position or to toggle automatic help.

```

1  EXPLAIN SELECT
2  name, yearId, W
3  FROM
4  teams T
5  WHERE
6  W = (SELECT
7        MAX(W)
8      FROM
9        teams y
10     WHERE
11       t.yearID = y.yearID
12     );

```

Result Grid

name	yearId	W
Philadelphia Athletics	1871	21
Boston Red Stockings	1872	39
Boston Red Stockings	1873	43
Boston Red Stockings	1874	52
Boston Red Stockings	1875	71
Chicago White Stockings	1876	52
Boston Red Caps	1877	42
Boston Red Caps	1878	41
Providence Grays	1879	59

Output

#	Time	Action	Message	Duration / Fetch
903	18:07:33	SELECT T.name,T.yearId,T.W FROM teams T JOIN...	154 row(s) returned	0.016 sec / 0.000 sec
904	18:07:34	SELECT T.name,T.yearId,T.W FROM teams T JOIN...	154 row(s) returned	0.000 sec / 0.000 sec
905	18:10:17	DROP INDEX IndexOn_W_yearID on teams	Error Code: 1091. Can't DROP 'IndexOn_W_yearID'; check that column/key exists	0.000 sec
906	18:10:45	SELECT name, yearId, W FROM teams T WH...	154 row(s) returned	0.032 sec / 0.000 sec
907	18:10:49	SELECT name, yearId, W FROM teams T WH...	154 row(s) returned	0.031 sec / 0.000 sec
908	18:10:50	SELECT name, yearId, W FROM teams T WH...	154 row(s) returned	0.031 sec / 0.000 sec
909	18:10:50	SELECT name, yearId, W FROM teams T WH...	154 row(s) returned	0.047 sec / 0.000 sec

Result Grid

id	select_type	table	partitio	type	possible_keys	key	key_le	ref	rows	filtered	Extra
1	PRIMARY	T	NULL	ALL	NULL	NULL	NULL	NULL	2745	100.00	Using where
2	DEPENDEN...	v	NULL	ref	PRIMARY.IndexOn W yearID	IndexOn W yearID	4	baseba...	19	100.00	Using index

## Modified Query

```

13
14  EXPLAIN SELECT T.name,T.yearId,T.W FROM teams T JOIN
15  (SELECT
16    MAX(W) w ,yearId
17  FROM
18    teams
19  where w>=0
20  group by yearID
21  ) MW on MW.w=T.W AND MW.yearId=T.yearId;
22
23  SELECT max(W) FROM teams;
24
25  CREATE INDEX IndexOn_W_yearID on teams(yearID,W);
26  DROP INDEX IndexOn_W_yearID on teams;
27
28  CREATE INDEX IndexOn_YearID on teams(yearID) using btree;

```

Result Grid

id	select_type	table	partitio	type	possible_keys	key	key_le	ref	rows	filtered	Extra
1	PRIMARY	<derived2>	NULL	ALL	NULL	NULL	NULL	NULL	47	100.00	Using where
1	PRIMARY	T	NULL	ref	PRIMARY.IndexOn W yearID	IndexOn W yearID	9	MW,ve...	1	100.00	Using index
2	DERIVED	teams	NULL	range	PRIMARY.IndexOn W yearID	IndexOn W yearID	4	NULL	144	33.33	Using where: Using index for aro

Output

#	Time	Action	Message	Duration / Fetch
915	18:14:10	SELECT T.name,T.yearId,T.W FROM teams T JOIN...	154 row(s) returned	0.000 sec / 0.000 sec
916	18:14:14	SELECT T.name,T.yearId,T.W FROM teams T JOIN...	154 row(s) returned	0.000 sec / 0.000 sec
917	18:14:14	SELECT T.name,T.yearId,T.W FROM teams T JOIN...	154 row(s) returned	0.000 sec / 0.000 sec
918	18:14:15	SELECT T.name,T.yearId,T.W FROM teams T JOIN...	154 row(s) returned	0.000 sec / 0.000 sec
919	18:14:16	SELECT T.name,T.yearId,T.W FROM teams T JOIN...	154 row(s) returned	0.016 sec / 0.000 sec
920	18:14:16	SELECT T.name,T.yearId,T.W FROM teams T JOIN...	154 row(s) returned	0.000 sec / 0.000 sec
921	18:14:22	SELECT T.name,T.yearId,T.W FROM teams T JOIN...	154 row(s) returned	0.000 sec / 0.000 sec

So, using a join instead of a subquery and then creating a composite index on yearId,W helped improve the query performance overall.

## Query 5

In this query, I observe the following things-

Result Grid												
Filter Rows: <input type="text"/>												
Export: <input type="button" value=""/> Wrap Cell Content: <input checked="" type="checkbox"/>												
	id	select_type	table	partitio	type	possible_keys	key	key_le	ref	rows	filtered	Extra
	1	PRIMARY	teams	<a href="#">NULL</a>	ALL	PRIMARY	<a href="#">NULL</a>	<a href="#">NULL</a>	<a href="#">NULL</a>	2745	100.00	<a href="#">NULL</a>
	1	PRIMARY	<derived5>	<a href="#">NULL</a>	ref	<auto kev0>	<auto kev0>	25	baseba...	32	100.00	<a href="#">NULL</a>
	1	PRIMARY	<derived2>	<a href="#">NULL</a>	ref	<auto kev0>	<auto kev0>	25	baseba...	3625	100.00	Usina where
	5	DERIVED	battino	<a href="#">NULL</a>	ALL	<a href="#">NULL</a>	<a href="#">NULL</a>	<a href="#">NULL</a>	<a href="#">NULL</a>	98745	90.00	Usina where: Usina temporarv: Usina f
	2	DERIVED	<derived4>	<a href="#">NULL</a>	ALL	<a href="#">NULL</a>	<a href="#">NULL</a>	<a href="#">NULL</a>	<a href="#">NULL</a>	88870	100.00	Usina where: Usina temporarv: Usina f
	2	DERIVED	<derived3>	<a href="#">NULL</a>	ref	<auto kev0>	<auto kev0>	25	A.team...	336	33.33	Usina where
	4	DERIVED	battino	<a href="#">NULL</a>	ALL	<a href="#">NULL</a>	<a href="#">NULL</a>	<a href="#">NULL</a>	<a href="#">NULL</a>	98745	90.00	Usina where: Usina temporarv: Usina f
	3	DERIVED	battino	<a href="#">NULL</a>	ALL	<a href="#">NULL</a>	<a href="#">NULL</a>	<a href="#">NULL</a>	<a href="#">NULL</a>	98745	100.00	Usina temporarv: Usina filesort

The use of AB and H is excessive from the batting table. Although we're grouping by teamID, yearID and lgID, even if we apply a composite index on these three, the performance degrades, because it has to look through the batting table anyway to group the tuples.

What we can do here is reduce the number of tuples before we group them. We are summing up AB and H and also using 'AB is not null' in the where clause. Now, while grouping them, we can ignore all the tuples that do not contribute to this aggregation by applying a where clause in such a way that we don't have to go through all these non-contributing tuples.

Here is the re-written query-



```

7 FROM
8 (SELECT
9     count(masterID) as cnt, A.yearID, A.teamID, A.lgID
10    FROM
11    (select
12        masterID,
13        teamID,
14        yearID,
15        lgID,
16        sum(AB),
17        sum(H),
18        sum(H) / sum(AB) as avg
19    FROM
20        batting
21    WHERE AB>=0 and H>=0
22    GROUP BY teamID , yearID , lgID , masterID) B
23    INNER JOIN
24    (select
25        teamID,
26        yearID,
27        lgID,
28        sum(AB),
29        sum(H),
30        sum(H) / sum(AB) as avg
31    FROM
32        batting
33    WHERE AB>=0 and H>=0
34    GROUP BY teamID , yearID , lgID) A
35    ON
36        A.avg >= B.avg AND A.teamID = B.teamID
37        AND A.yearID = B.yearID
38        AND A.lgID = B.lgID
39    GROUP BY teamID , yearID , lgID) C
40    INNER JOIN
41    (SELECT
42        count(masterID) as cnt, teamID,yearID, lgID
43    FROM
44        batting
45    WHERE ab>=0
46    GROUP BY teamID,yearID,lgID) D
47    INNER JOIN
48    teams
49    ON
50        C.cnt / D.cnt >= 0.75
51        AND C.yearID = D.yearID
52        AND C.teamID = D.teamID
53        AND C.lgID = D.lgID
54        AND teams.yearID = C.yearID
55        AND teams.lgID = C.lgID
56        AND teams.teamID = C.teamID;

```

I've highlighted the changes I've made in yellow. These changes help in reducing the number of tuples we need to group.

Below are the performance statistics-

Original Query

Result Grid

Filter Rows:

Export:

Wrap Cell Content:

	id	select_type	table	partitio	type	possible_keys	key	key_le	ref	rows	filtered	Extra
	1	PRIMARY	teams	NULL	ALL	PRIMARY	NULL	NULL	NULL	2745	100.00	NULL
	1	PRIMARY	<derived5>	NULL	ref	<auto kev0>	<auto kev0>	25	baseba...	32	100.00	NULL
	1	PRIMARY	<derived2>	NULL	ref	<auto kev0>	<auto kev0>	25	baseba...	3625	100.00	Using where
	5	DERIVED	battino	NULL	ALL	NULL	NULL	NULL	NULL	98745	90.00	Using where: Using temporary: Using filesort
	2	DERIVED	<derived4>	NULL	ALL	NULL	NULL	NULL	NULL	88870	100.00	Using where: Using temporary: Using filesort
	2	DERIVED	<derived3>	NULL	ref	<auto kev0>	<auto kev0>	25	A.team...	336	33.33	Using where
	4	DERIVED	battino	NULL	ALL	NULL	NULL	NULL	NULL	98745	90.00	Using where: Using temporary: Using filesort
	3	DERIVED	battino	NULL	ALL	NULL	NULL	NULL	NULL	98745	100.00	Using temporary: Using filesort

810

17:33:10

SELECT

C.yearID as year, name as teamName,...

93 row(s) returned

14.969 sec / 0.000 sec

## Modified Query

<

Result Grid

Filter Rows:

Exports:

Wrap Cell Content:

	id	select_type	table	partitio	type	possible_keys	key	key_le	ref	rows	filtered	Extra
	1	PRIMARY	teams	NULL	ALL	PRIMARY	NULL	NULL	NULL	2745	100.00	Using where
	1	PRIMARY	<derived5>	NULL	ref	<auto kev0>	<auto kev0>	25	baseba...	11	100.00	NULL
	1	PRIMARY	<derived2>	NULL	ref	<auto kev0>	<auto kev0>	25	baseba...	57	100.00	Using where
	5	DERIVED	battino	NULL	ALL	NULL	NULL	NULL	NULL	98745	33.33	Using where: Using temporary: Using filesort
	2	DERIVED	<derived3>	NULL	ALL	NULL	NULL	NULL	NULL	10969	100.00	Using where: Using temporary: Using filesort
	2	DERIVED	<derived4>	NULL	ref	<auto kev0>	<auto kev0>	25	B.team...	43	33.33	Using where
	4	DERIVED	battino	NULL	ALL	NULL	NULL	NULL	NULL	98745	11.11	Using where: Using temporary: Using filesort
	3	DERIVED	battino	NULL	ALL	NULL	NULL	NULL	NULL	98745	11.11	Using where: Using temporary: Using filesort

812 17:35:02 SELECT C.yearID as year, name as teamName,... 93 row(s) returned

1.781 sec / 0.000 sec

I also tried applying index on AB and H so that it is easier to search the when we use the where clause but the performance didn't speed up. This is one of the cases where using an index does not increase performance.

As it can be seen from the screenshots, the query performance is 14 times better now.

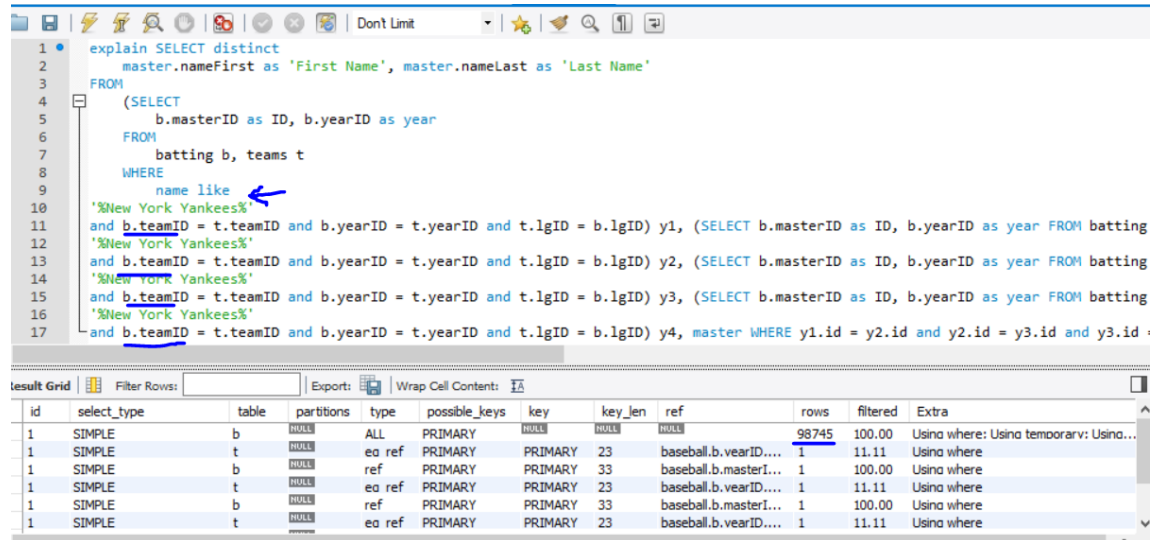
## Query 6

In this query, I observe three things-

One, that the cost of scanning the batting table b is quite a lot(98745 rows). We can try adding an index on the teamID, yearID and leagueID columns, to improve cost.

Second, instead of using the like keyword on the name column in teams table, we can use =

Third, Since we are trying to find a specific team name, we could use an index on the name column of the teams table.



```
1 explain SELECT distinct
2   master.nameFirst as 'First Name', master.nameLast as 'Last Name'
3 FROM
4   (SELECT
5     b.masterID as ID, b.yearID as year
6   FROM
7     batting b, teams t
8   WHERE
9     name like
10    '%New York Yankees%'
11    and b.teamID = t.teamID and b.yearID = t.yearID and t.lgID = b.lgID) y1, (SELECT b.masterID as ID, b.yearID as year FROM batting
12    '%New York Yankees%'
13    and b.teamID = t.teamID and b.yearID = t.yearID and t.lgID = b.lgID) y2, (SELECT b.masterID as ID, b.yearID as year FROM batting
14    '%New York Yankees%'
15    and b.teamID = t.teamID and b.yearID = t.yearID and t.lgID = b.lgID) y3, (SELECT b.masterID as ID, b.yearID as year FROM batting
16    '%New York Yankees%'
17    and b.teamID = t.teamID and b.yearID = t.yearID and t.lgID = b.lgID) y4, master WHERE y1.id = y2.id and y2.id = y3.id and y3.id =
```

id	select_type	table	partitions	type	possible_keys	key	key_len	ref	rows	filtered	Extra
1	SIMPLE	b	NONE	ALL	PRIMARY	NONE	NONE	NONE	98745	100.00	Using where: Using temporary: Using...
1	SIMPLE	t	NONE	eq ref	PRIMARY	PRIMARY	23	baseball.b.yearID...	1	11.11	Using where
1	SIMPLE	b	NONE	ref	PRIMARY	PRIMARY	33	baseball.b.masterI...	1	100.00	Using where
1	SIMPLE	t	NONE	eq ref	PRIMARY	PRIMARY	23	baseball.b.yearID...	1	11.11	Using where
1	SIMPLE	b	NONE	ref	PRIMARY	PRIMARY	33	baseball.b.masterI...	1	100.00	Using where
1	SIMPLE	t	NONE	eq ref	PRIMARY	PRIMARY	23	baseball.b.yearID...	1	11.11	Using where

## What I did

Created an index on teamID, yearID and lgID columns in the batting table.

Created an index on teamID, yearID and lgID columns in the teams table.

Removed the like keyword and replaced it with =

Created a hash index on teamName column in the teams table.

## Improvement

The cost came down from 98745 rows to 101 rows.

```

1 explain SELECT distinct
2   master.nameFirst as 'First Name', master.nameLast as 'Last Name'
3 FROM
4   (SELECT
5     b.masterID as ID, b.yearID as year
6   FROM
7     batting b, teams t
8   WHERE
9     name=
10    'New York Yankees'
11 and b.teamID = t.teamID and b.yearID = t.yearID and t.lgID = b.lgID) y1, (SELECT b.masterID as ID, b.yearID as year FROM batting b, tea
12 'New York Yankees'
13 and b.teamID = t.teamID and b.yearID = t.yearID and t.lgID = b.lgID) y2, (SELECT b.masterID as ID, b.yearID as year FROM batting b, tea
14 'New York Yankees'
15 and b.teamID = t.teamID and b.yearID = t.yearID and t.lgID = b.lgID) y3, (SELECT b.masterID as ID, b.yearID as year FROM batting b, tea
16 'New York Yankees'
17 and b.teamID = t.teamID and b.yearID = t.yearID and t.lgID = b.lgID) y4, master WHERE y1.id = y2.id and y2.id = y3.id and y3.id = y4.id
18 CREATE INDEX teamName_Index ON teams(name) using hash;

```

select_type	table	partition	type	possible_keys	key	key_len	ref	rows	filtered	Extra
SIMPLE	t	NULL	ref	PRIMARY, newIndexBYL, teamName_Index	teamName_Index	153	const	101	100.00	Using index; Using temporary; Using filesort
SIMPLE	b	NULL	ref	PRIMARY, newIndexBYL	newIndexBYL	25	baseba...	34	100.00	Using index
SIMPLE	b	NULL	ref	PRIMARY, newIndexBYL	PRIMARY	33	baseba...	1	100.00	Using where
SIMPLE	t	NULL	eq ref	PRIMARY, newIndexBYL, teamName_Index	PRIMARY	23	baseba...	1	5.00	Using where
SIMPLE	b	NULL	ref	PRIMARY, newIndexBYL	PRIMARY	33	baseba...	1	100.00	Using where
SIMPLE	t	NULL	eq ref	PRIMARY, newIndexBYL, teamName_Index	PRIMARY	23	baseba...	1	5.00	Using where
SIMPLE	b	NULL	ref	PRIMARY, newIndexBYL	PRIMARY	33	baseba...	1	100.00	Using where
SIMPLE	t	NULL	eq ref	PRIMARY, newIndexBYL, teamName_Index	PRIMARY	23	baseba...	1	5.00	Using where
SIMPLE	m...	NULL	eq ref	PRIMARY	PRIMARY	32	func	1	100.00	Using where

104 x Read Only

## Query 7

In the given query, the current cost is 23956 rows. This cost stems from the grouping that is being applied on yearID, teamID and lgID 2 times to find out the sum of the salaries.

So, the two derived tables can be optimized by the following approach-

1. Use where before the group by clause. Apply where clause on salary. This way we will have lesser number of rows to group.
2. Try creating a composite index using yearID, teamID, lgID and salary.
3. Apply indexing on the salary column of the salaries table, so that when we use the where clause, we can reduce the cost by reducing the number of rows we need to traverse.

```

32
33 • explain SELECT
34     name,
35     A.lgID,
36     A.S as TotalSalary,
37     A.yearID as Year,
38     B.S as PreviousYearSalary,
39     B.yearID as PreviousYear
40 FROM
41     (SELECT
42         sum(salary) as S, yearID, teamID, lgID
43     FROM
44         salaries
45     group by yearID , teamID , lgID) A,
46     (SELECT
47         sum(salary) as S, yearID, teamID, lgID
48     FROM
49         salaries
50     group by yearID , teamID , lgID) B,
51     teams
52 WHERE
53     A.yearID = B.yearID + 1

```

id	select_type	table	partitio	type	possible_keys	key	key_le	ref	rows	filtered	Extra
1	PRIMARY	<derived3>	HULL	ALL	HULL	HULL	HULL	HULL	23956	100.00	HULL
2	PRIMARY	teams	HULL	ea ref	PRIMARY	PRIMARY	23	func.B....	1	100.00	Using where
3	PRIMARY	<derived2>	HULL	ref	<auto kev0>	<auto kev0>	23	func.B....	10	100.00	Using where
4	DERIVED	salaries	HULL	index	PRIMARY	PRIMARY	52	HULL	23956	100.00	HULL
5	DERIVED	salaries	HULL	index	PRIMARY	PRIMARY	52	HULL	23956	100.00	HULL

I modified the query and this is how it looks like now-

```

• Explain SELECT
    name,
    A.lgID,
    A.S as TotalSalary,
    A.yearID as Year,
    B.S as PreviousYearSalary,
    B.yearID as PreviousYear
FROM
    (SELECT
        sum(salary) as S, yearID, teamID, lgID
    FROM
        salaries
    where salary>=0
    group by yearID , teamID , lgID) A,
    (SELECT
        sum(salary) as S, yearID, teamID, lgID
    FROM
        salaries
    where salary>=0
    group by yearID , teamID , lgID) B,
    teams
WHERE
    A.yearID = B.yearID + 1
    AND (A.S * 2) <= (B.S)
    AND A.teamID = B.teamID
    AND A.lgID = B.lgID
    AND teams.yearID = A.yearID
    AND teams.lgID = A.lgID
    AND teams.teamID = A.teamID;

```

Now this is the optimization result after I amend the query with where clause mentioned in the 1<sup>st</sup> approach-

id	select_type	table	partitio	type	possible_keys	key	key_le	ref	rows	filtered	Extra
1	PRIMARY	<derived3>	HULL	ALL	HULL	HULL	HULL	HULL	7984	100.00	HULL
2	PRIMARY	teams	HULL	ea ref	PRIMARY	PRIMARY	23	func.B....	1	100.00	Using where
3	PRIMARY	<derived2>	HULL	ref	<auto kev0>	<auto kev0>	23	func.B....	10	100.00	Using where
4	DERIVED	salaries	HULL	index	PRIMARY	PRIMARY	52	HULL	23956	33.33	Using where
5	DERIVED	salaries	HULL	index	PRIMARY	PRIMARY	52	HULL	23956	33.33	Using where

This performs better but now we add the 2nd approach of creating a composite index. These are the results when we do that-

Result Grid

Filter Rows:

Export:

Wrap Cell Contents:

	id	select_type	table	partitio	type	possible_keys	key	key_le	ref	rows	filtered	Extra
	1	PRIMARY	<derived3>	NULL	ALL	NULL	NULL	NULL	NULL	7984	100.00	NULL
	1	PRIMARY	teams	NULL	ea ref	PRIMARY	PRIMARY	23	func.B....	1	100.00	Usina where
	1	PRIMARY	<derived2>	NULL	ref	<auto kev0>	<auto kev0>	23	func.B....	10	100.00	Usina where
	3	DERIVED	salaries	NULL	index	PRIMARY	PRIMARY	52	NULL	23956	33.33	Usina where
	2	DERIVED	salaries	NULL	index	PRIMARY	PRIMARY	52	NULL	23956	33.33	Usina where

The cost remains the same. Perhaps, because we already have an index created on the primary keys of the salaries table. So, its already using an index.

Now we try the 3rd approach of creating an index on the salary column. These are the results when we do that-

Result Grid											
Filter Rows:		Export:		Wrap Cell Content:							
id	select_type	table	partitio	type	possible_keys	key	key_le	ref	rows	filtered	Extra
1	PRIMARY	<derived3>	NULL	ALL	NULL	NULL	NULL	NULL	11978	100.00	NULL
1	PRIMARY	teams	NULL	eq ref	PRIMARY	PRIMARY	23	func.B....	1	100.00	Using where
1	PRIMARY	<derived2>	NULL	ref	<auto key0>	<auto key0>	23	func.B....	10	100.00	Using where
3	DERIVED	salaries	NULL	range	PRIMARY,newIndexTeams	newIndexTeams	9	NULL	11978	100.00	Using where: Using index: Using ter
2	DERIVED	salaries	NULL	range	PRIMARY,newIndexTeams	newIndexTeams	9	NULL	11978	100.00	Using where: Using index: Using ter

We have halved the cost, which means 50% improvement.

So, we optimized this query by two ways- Restricting the number of tuples that need to be grouped by using a where clause before group by and applying an index on the column that is being used in the where clause before the group by clause.