

Optimization

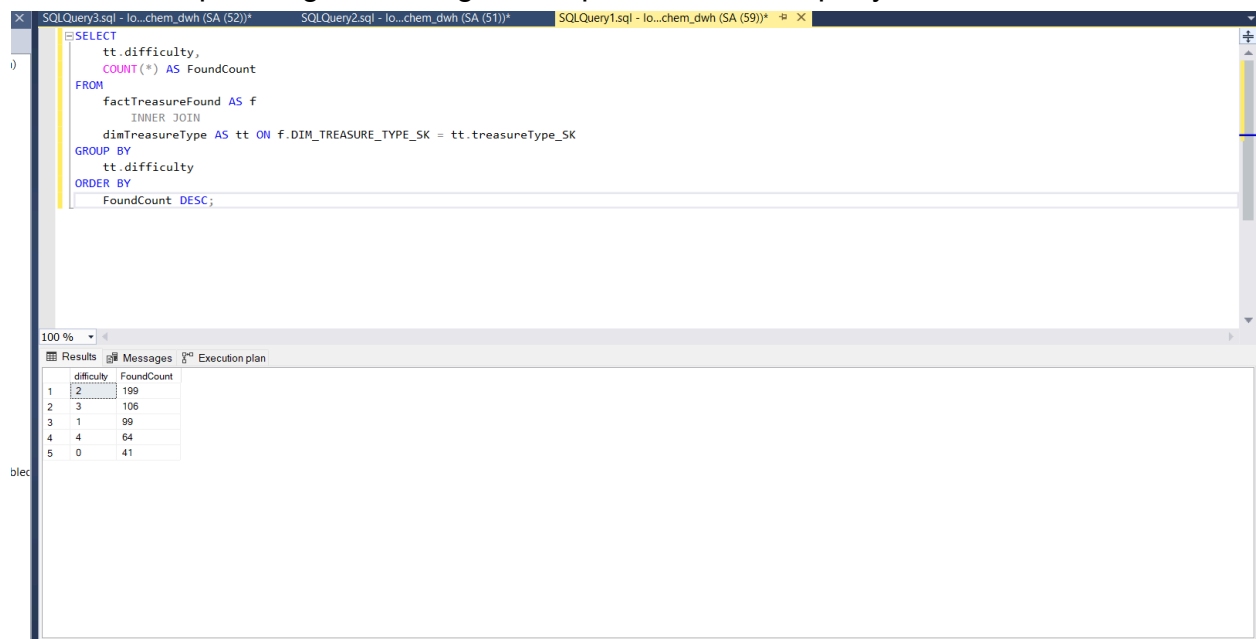
[S1]: Create a logical index based on a research question

This is the research question:

```
-- Are there caches which are more popular at the moment and  
sought after by various people?
```

Below is the output of the said research question

Photo from sql manager showing the output of the basic query:



The screenshot shows a SQL Server Enterprise Manager window with three tabs. The active tab is 'SQLQuery1.sql - lo...chem_dwh (SA (59))'. The query editor displays the following SQL query:

```
SELECT  
    tt.difficulty,  
    COUNT(*) AS FoundCount  
FROM  
    factTreasureFound AS f  
    INNER JOIN  
    dimTreasureType AS tt ON f.DIM_TREASURE_TYPE_SK = tt.treasureType_SK  
GROUP BY  
    tt.difficulty  
ORDER BY  
    FoundCount DESC;
```

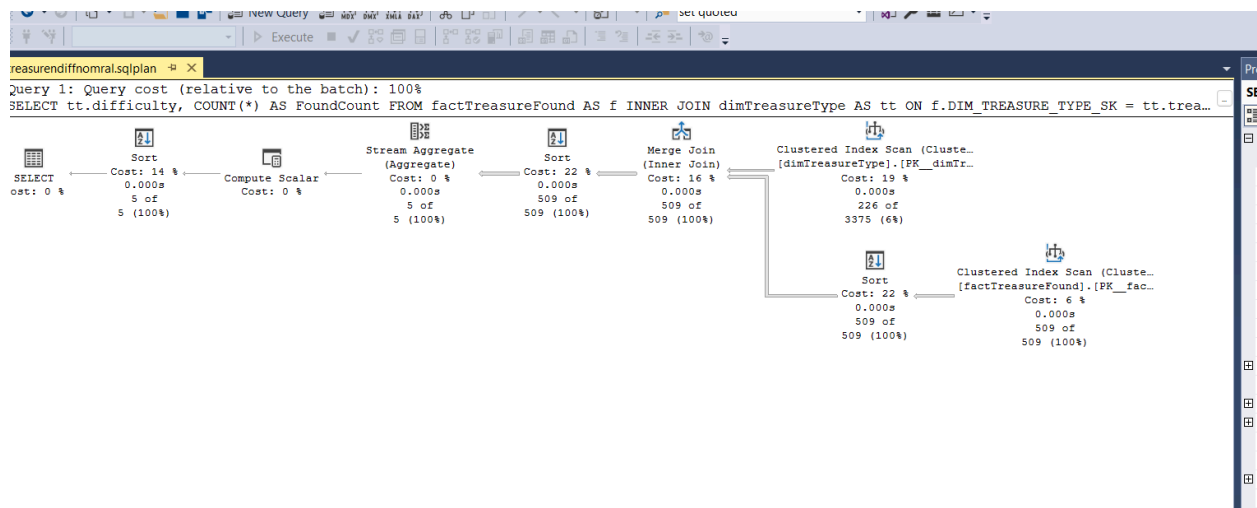
Below the query editor, the 'Results' tab is selected, showing the output of the query. The results are as follows:

	difficulty	FoundCount
1	2	199
2	3	106
3	1	99
4	4	64
5	0	41

Based on the provided data, it appears that caches with a difficulty level of 2 are the most popular at the moment, with 199 found counts.

-- This suggests that caches with a moderate difficulty level are sought after by various people. The popularity decreases as the difficulty level increases,
-- with fewer caches found for difficulty levels 3, 1, 4, and 0, respectively.

Here is the execution plan for the research question query:



Before the index optimization, the query had many joins and sorting operations, leading to a complex execution plan with numerous links. These joins and sorting operations increase computational overhead and make the query execution less efficient.

Then we will be going to tuning advisor and login with same credentials and will first run the query if you have not done so yet, so that the tuning advisor knows from which database and table to run the recommendations on, for us we will be choosing the index options in the tuning options then we start the analysis.

Recommendation	Details	Partition Scheme	Size (KB)	Definition
stat_370100359_1_2				(treasureType_SK] [difficulty]
[_cta_mv_0]				SELECT [dbo].[dimTreasureType].[difficulty] as _col_1, count(*) as _col_2 FROM [cta_mv_0]
index_cta_mv_0_c_7_1390627997_K1	clustered, unique		8	([_col_1] asc)

We will be choosing both the create view and the clustered index on said view: which are shown below

```
SQLQuery3.sql - lo...chem_dwh (SA (52))*  SQLQuery2.sql - lo...chem_dwh (SA (51))*  SQLQuery1.sql - lo...chem_dwh (SA (59))*
CREATE VIEW [dbo].[_dta_mv_0] WITH SCHEMABINDING
AS
SELECT [dbo].[dimTreasureType].[difficulty] as _col_1, count_big(*) as _col_2 FROM [dbo].[dimTreasureType], [dbo].[factTreasureFound] WHERE [dbo].[dimTreasureType].[treasure_id] = [dbo].[factTreasureFound].[treasure_id]
```

Above we create the view on the original query parameters

Then create a unique clustered index on that view:

```
SQLQuery3.sql - lo...chem_dwh (SA (52))*  SQLQuery2.sql - lo...chem_dwh (SA (51))*  SQLQuery1.sql - lo...chem_dwh (SA (59))*
SET ARITHABORT ON
SET CONCAT_NULL_YIELDS_NULL ON
SET QUOTED_IDENTIFIER ON
SET ANSI_NULLS ON
SET ANSI_PADDING ON
SET ANSI_WARNINGS ON
SET NUMERIC_ROUNDABORT OFF

CREATE UNIQUE CLUSTERED INDEX [_dta_index__dta_mv_0_c_7_1390627997__K1] ON [dbo].[_dta_mv_0]
(
    [_col_1] ASC
) WITH (SORT_IN_TEMPDB = OFF, DROP_EXISTING = OFF, ONLINE = OFF) ON [PRIMARY]
```

100 %

Messages Execution plan

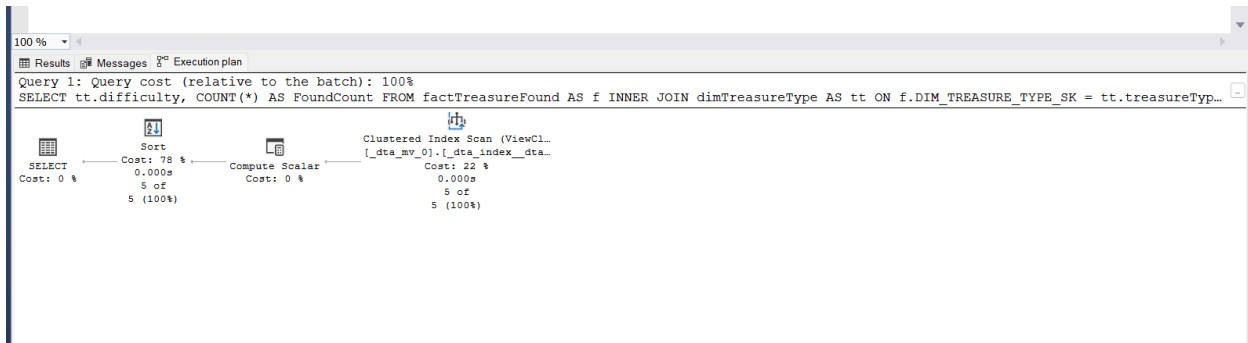
SET ANSI_WARNINGS ON

T-SQL

We create it because:

Creating a unique clustered index on the view derived from the original query parameters boosts query performance by organizing data efficiently, speeding up data retrieval, enhancing join operations, and reducing I/O operations. This index assists the query optimizer in generating efficient execution plans, resulting in overall better database performance.

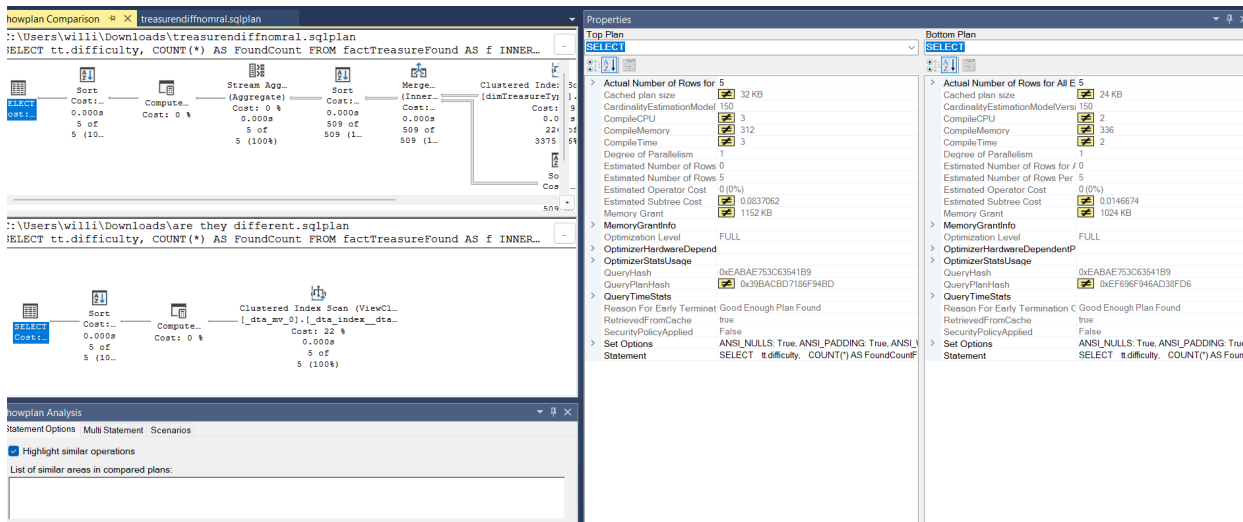
And run the query again:



Incorporating a clustered index materialized view into the original query simplifies execution by reducing the need for complex joins, sorting, and merging operations. With pre-computed data stored in the materialized view, the query optimizer can streamline execution, resulting in a more efficient query plan with fewer computational steps

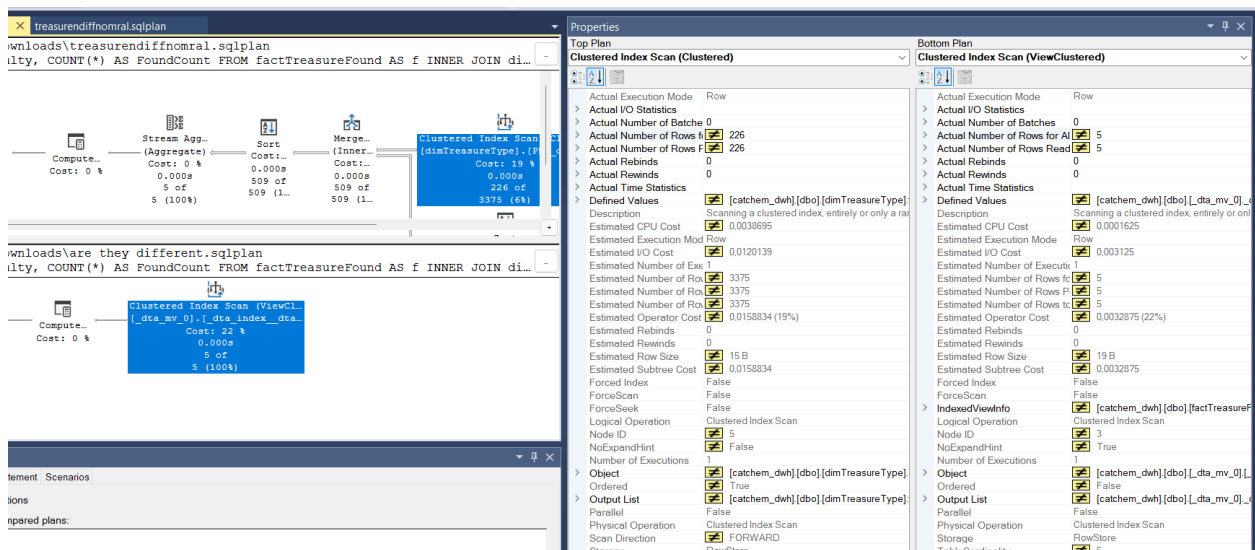
Also, the logical index on the original query helps streamline data access by providing a structured organization of the queried columns. This aids in faster data retrieval and improved query performance without physically altering the database schema.

Execution Plan comparison:



Based on this analysis, it appears that the query has been optimized effectively. By employing joins, indexes, and clusters, the query execution has notably improved in terms of speed and efficiency, resulting in a lower subtree cost.

Comparing the performance metrics before and after optimization, there's a significant reduction in the subtree cost with the optimized query. For instance, the subtree cost has decreased from 0.08 to 0.01 with the utilization of an index. This indicates a substantial enhancement in query execution efficiency.



The decrease in I/O from 0.01 to 0.003 signifies improved efficiency in data input/output operations. This optimization reduces the amount of data read from or written to storage, leading to faster query execution and lower resource usage.

Similarly, the decrease in CPU cost from 0.003 to 0.0001 indicates enhanced efficiency in CPU utilization during query execution. This optimization results in reduced processing time and resource consumption, contributing to overall faster query performance and improved system responsiveness.

In summary, the integration of joins, indexes, and clusters has significantly optimized the query, resulting in notable improvements in performance metrics such as subtree cost, CPU utilization, and I/O operations. The reduction in subtree cost highlights the effectiveness of these optimization techniques in enhancing overall query efficiency, while the decrease in CPU cost and I/O reflects improved resource utilization and faster data access.

Decision:

Considering the significant decreases observed in subtree cost, CPU cost, and I/O operations after implementing the logical index, it's evident that the index has effectively optimized the query's performance. The improvements in these key metrics indicate enhanced query efficiency and resource utilization.

Therefore, based on the observed benefits of the logical index, it is recommended to use the indexed version of the query for improved performance and resource efficiency.

[S2] Column Storage Optimization

I chose the following research question and query for column storage optimization.

Q. Are more difficult caches done on weekends?

```
SELECT
    dt.difficulty,
    dd.Weekday,
    COUNT(ftf.TreasureFoundID) AS total_caches_searched
FROM
    catchem_dwh.dbo.factTreasureFound ftf
    JOIN
    catchem_dwh.dbo.dimTreasureType dt ON ftf.DIM_TREASURE_TYPE_SK = dt.treasureType_SK
    JOIN
    catchem_dwh.dbo.dimDay dd ON ftf.DIM_DAY_SK = dd.day_SK
GROUP BY
    dt.difficulty, dd.Weekday
ORDER BY
    dt.difficulty, dd.Weekday;
```

because it has aggregate expression in a select clause and joins with two other tables that will deal with potentially large amounts of datasets.

After I ran the query for the first time, the Turning Advisor gave the recommendation plans.

The screenshot shows the Database Engine Tuning Advisor (DETA) window. The 'Recommendations' tab is active, displaying an 'Estimated improvement: 74%'. Below this, the 'Index Recommendations' section lists two recommendations for the 'catchem_dwh' database:

Database Name	Object Name	Recommendation	Target of Recommendation	Details	Partition Scheme	Size (KB)	Definition
catchem_dwh	[dbo].[dimDay]	create	_dta_index_dimDay_7_738101670_col_				
catchem_dwh	[dbo].[dimDay]	create	_dta_stat_738101670_8_1				[(Weekday], [day_SK])

I implemented the following query to create column storage and create statistics for dimDay.

```
CREATE NONCLUSTERED COLUMNSTORE INDEX [_dta_index_dimDay_7_738101670__col__] ON [dbo].[dimDay]
(
    [day_SK],
    [Date],
    [DayOfMonth],
    [Month],
    [Year],
    [DayOfWeek],
    [DayOfYear],
    [Weekday],
    [MonthName],
    [Season]
)WITH (DROP_EXISTING = OFF, COMPRESSION_DELAY = 0) ON [PRIMARY]

CREATE STATISTICS [_dta_stat_738101670_8_1] ON [dbo].[dimDay]([Weekday], [day_SK])
```

Top plan refers to the query **after optimization** and **Bottom Plan** is the **original** query. Estimated Subtree Cost has been reduced from **0.143501** to **0.0265457**.

The top plan used 3 seconds of CPU time and 593KB of memory during compilation, whereas the bottom plan used 41 seconds of CPU time and 416 KB of memory, which indicates the top plan compiled faster and more efficiently.

Overall, the top plan seems to be more efficient than the bottom plan due to lower estimated costs, and faster compile time.

Top Plan	Bottom Plan
SELECT	SELECT
Actual Number of Rows for 0	Actual Number of Rows for 0
BatchModeOnRowStoreUs True	Cached plan size 64 KB
Cached plan size 136 KB	CardinalityEstimationMode 150
CardinalityEstimationMode 150	CompileCPU 41
CompileCPU 3	CompileMemory 416
CompileMemory 592	CompileTime 43
CompileTime 3	Estimated Number of Rows 0
Estimated Number of Rows 0	Estimated Number of Rows 29.4978
Estimated Number of Rows 29.4978	Estimated Operator Cost 0 (0%)
Estimated Operator Cost 0 (0%)	Estimated Subtree Cost 0.143501
Estimated Subtree Cost 0.0265457	MemoryGrantInfo
MemoryGrantInfo	Optimization Level FULL
Optimization Level FULL	OptimizerHardwareDependent
OptimizerHardwareDependent	OptimizerStatsUsage
OptimizerStatsUsage	QueryHash 0x421D5D60DC2B8DCE
QueryHash 0x421D5D60DC2B8DCE	QueryPlanHash 0xFBC9107803DADB13
QueryPlanHash 0xE55A109415536D99	Reason For Early Termination Good Enough Plan Found
Reason For Early Termination Time Out	RetrievedFromCache true
RetrievedFromCache true	SecurityPolicyApplied False
SecurityPolicyApplied False	Set Options
Set Options	ANSI_NULLS: True, ANSI_PADDING: True
Statement	Statement
SELECT dt.difficulty, dd.Weekday, CO	SELECT dt.difficulty, dd.Weekday, CO