## Data Warehouse Optimization - Report

### Piotr Sulewski, Tomasz Sankowski

### January 2025

### Contents

1	Ain	Aim of the Laboratory								
<b>2</b>	Pre	Preliminary Assumptions								
	2.1	Database Specifications	2							
	2.2	Testing Environment								
	2.3	Conditions During Testing								
3	The	eoretical Assumptions	3							
4	Tes	ting	4							
	4.1	Brief description of the queries	4							
	4.2	MOLAP Performance	6							
	4.3	HOLAP Performance	6							
	4.4	ROLAP Performance	7							
	4.5	Comparison of processing time for different cube models	8							
		4.5.1 Comparison of cube processing performance	8							
		4.5.2 Comparison of Query 2 performance	8							
		4.5.3 Comparison of Query 5 performance	9							
		4.5.4 Comparison of Query 9 performance	9							
5	Cac	the and Aggregation Settings	10							
	5.1	Aggregation Design	10							
	5.2	Cache removal	11							
6	Dis	cussion	12							
	6.1	Cube Processing	12							
	6.2									
	6.3		13							

### 1. Aim of the Laboratory

The aim of this task is to address issues related to various physical cube models and the design of aggregations. The focus is on demonstrating the impact of these models on performance and usability in a data warehouse environment.

### 2. Preliminary Assumptions

#### 2.1. Database Specifications

• Size of the Data Warehouse: 976.00 MB

• Number of rows in the monitoring fact table (main): 1,176,609

#### 2.2. Testing Environment

Measurements were taken on the following setup:

• Processor: Intel Core i7-6700K

• Memory: 16 GB RAM

• Storage: 1 TB NVMe SSD

• Operating System: Windows 10

• Tools Used:

- SQL Server Management Studio (SSMS) with SQL Server Profiler extension
- Visual Studio 2022 with Data Tools extension

#### 2.3. Conditions During Testing

To ensure reliable results, the following conditions were maintained:

- Only the necessary applications were running during measurements:
  - SQL Server Management Studio (SSMS)
  - Visual Studio 2022
  - A web browser with the instructions open

### 3. Theoretical Assumptions

Table 1: Comparison of MOLAP, HOLAP, and ROLAP

	MOLAP	HOLAP	ROLAP
Querying time	Short	Moderate <sup>1</sup> Moderate <sup>2</sup> Moderate	Long
Processing time	Long		Short
Total size	Large <sup>3</sup>		Small

<sup>&</sup>lt;sup>1</sup> Querying time can be short if well-designed aggregations are used.

Processing time is moderate, but if no aggregations are designed, it will be short.

<sup>&</sup>lt;sup>3</sup> Total size can be smaller if no aggregations are designed for measure groups.

### 4. Testing

### 4.1. Brief description of the queries

[Car Sharing];

Testing query execution times for different models, with and without defined aggregations. Testing cube processing times in the same testing settings

Figure 1: Query 2 Which car models generate the highest earnings per vehicle? WITH MEMBER [Measures]. [Vehicle Count] AS COUNT(EXISTING [Car].[Car ID].[Car ID]) MEMBER [Measures]. [Revenue Per Vehicle] AS IIF( [Measures].[Vehicle Count] > 0, [Measures].[Rental Cost] / [Measures].[Vehicle Count], NULL ) **SELECT** NON EMPTY TOPCOUNT ( [Car].[Brand].[Brand].MEMBERS \* [Car].[Model].[Model].MEMBERS, 10, [Measures].[Revenue Per Vehicle] ) ON ROWS, [Measures].[Revenue Per Vehicle] ON COLUMNS FROM

#### Figure 2: Query 5

Create a ranking of the 5 car models with the highest total mileage driven by customers in the last year.

Figure 3: Query 9

Considering the number of rides taken and the number of reported damages, what is the theoretical probability of causing damage within each age group of users?

```
WITH
   MEMBER [Measures].[Damaged Rides] AS
   (
        [Was Damaged].[Was Damaged].[Damaged],
        [Measures].[Rental Count]
   )
   MEMBER [Measures].[Probability of Damage] AS
   IIF(
        [Measures].[Rental Count] = 0,
        null,
        [Measures].[Damaged Rides] / [Measures].[Rental Count]
   )
   SELECT
   { [Measures].[Probability of Damage] } ON COLUMNS,
   [User].[Age Category].Children ON ROWS
   FROM
   [Car Sharing];
```

### 4.2. MOLAP Performance

Table 2: Performance Comparison of MOLAP with and without Aggregators

	MOLAP							
	Wit	hout Agg	gregators		W	ith Aggr	egators	
	Cube Processing	Query 2	Query 5	Query 9	Cube Processing	Query 2	Query 5	Query 9
	4740	45	43	43	4910	4	2	2
	4635	43	46	44	4830	4	2	2
	4567	44	44	43	4907	4	3	3
	4732	45	43	44	4879	4	2	2
	4629	43	40	43	4867	4	5	1
Average SD	$4660.6 \\ 66.05$	44 0.89	43.2 1.94	43.4 0.49	4878.6 $29.29$	4 0.0	2.8 1.17	2 0.63

### 4.3. HOLAP Performance

Table 3: Performance Comparison of HOLAP with and without Aggregators

	HOLAP							
	Wit	hout Agg	gregators	}	With Aggregators			
	Cube Processing	Query 2	Query 5	Query 9	Cube Processing	Query 2	Query 5	Query 9
	483	141	135	301	513	5	2	2
	487	135	136	295	493	4	2	2
	479	136	150	299	498	5	2	1
	481	137	132	294	513	5	2	2
	482	132	163	316	558	4	2	1
Average SD	$482.4 \\ 2.65$	136.2 2.93	143.8 11.16	301.0 7.92	515.0 22.93	4.6 0.49	2.0 0.0	1.6 0.49

### 4.4. ROLAP Performance

Table 4: Performance Comparison of ROLAP with and without Aggregators

	ROLAP							
	Witl	hout Agg	gregators	}	With Aggregators			
	Cube Processing	Query 2	Query 5	Query 9	Cube Processing	Query 2	Query 5	Query 9
	447	131	148	306	511	138	154	309
	$526 \\ 471$	$\frac{135}{134}$	$\frac{142}{140}$	$307 \\ 310$	566 $482$	$139 \\ 149$	$\frac{137}{137}$	$\frac{305}{296}$
	487	137	151	299	499	133	139	323
	502	138	132	309	486	132	145	301
Average SD	$486.6 \\ 26.84$	$135.0 \\ 2.45$	$142.6 \\ 6.62$	$306.2 \\ 3.87$	$508.8 \\ 30.37$	$138.2 \\ 6.05$	$142.4 \\ 6.50$	$306.8 \\ 9.17$

# 4.5. Comparison of processing time for different cube models

To achieve optimal results of the processing time of a cube we decided to take approximately 5 samples for each modification. The obtained results are presented in the following tables

#### 4.5.1 Comparison of cube processing performance

Table 5: Performance Comparison of Cube Processing with and without Aggregators in different cubes

	Cube Processing							
	With	out Aggreg	gators	With Aggregators				
	MOLAP	HOLAP	ROLAP	MOLAP	HOLAP	ROLAP		
Average SD	4660.6 66.05	482.4 2.65	486.6 26.64	4878.6 29.29	515.0 22.93	508.8 30.37		

### 4.5.2 Comparison of Query 2 performance

Table 6: Performance Comparison of Query 2 with and without Aggregators in different cubes

CICII CUDOS	·									
		Query 2								
	With	out Aggreg	gators	With Aggregators						
	MOLAP	HOLAP	ROLAP	MOLAP	HOLAP	ROLAP				
Average	44.0	136.2	135.0	4.0	4.6	138.2				
SD	0.89	2.93	2.45	0.0	0.49	6.05				

### 4.5.3 Comparison of Query 5 performance

Table 7: Performance Comparison of Query 5 with and without Aggregators in different cubes

	Query 5							
	With	out Aggreg	gators	With Aggregators				
	MOLAP	HOLAP	ROLAP	MOLAP	HOLAP	ROLAP		
Average SD	43.2 1.94	143.8 11.6	142.6 6.62	2.8 1.17	2.0 0.0	142.4 6.50		

### 4.5.4 Comparison of Query 9 performance

Table 8: Performance Comparison of Query 9 with and without Aggregators in different cubes

		Query 9							
	With	out Aggreg	gators	With Aggregators					
	MOLAP	HOLAP	ROLAP	MOLAP	HOLAP	ROLAP			
Average SD	43.4 0.49	301.0 7.92	306.2 3.87	2.0 0.63	1.6 0.49	306.8 9.17			

### 5. Cache and Aggregation Settings

### 5.1. Aggregation Design

The aggregations were designed with the "Perforfamous gain reaches" option set to 100%, which gave us 7 aggregations with an optimization level of 40% (4.4MB).

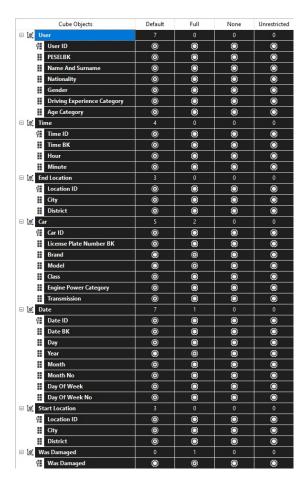


Figure 4: Aggregation Usage

### 5.2. Cache removal

Figure 5: Cache

```
Removing cache
```

### 6. Discussion

#### 6.1. Cube Processing

As observed, MOLAP required nearly ten times more time to process the cube compared to the ROLAP and HOLAP approaches, both with and without aggregations. This is because the MOLAP server needed to copy the entire warehouse storage into the cube, significantly slowing down the processing time. Additionally, adding aggregations only slightly affected the cube processing time, increasing it by approximately 10%. The higher processing time is due to the need to group measures by dimension members during cube processing. Our results confirm the assumption that MOLAP's cube processing time is considerably slower than ROLAP's and HOLAP's.

#### 6.2. Query Execution

MOLAP query execution was, on average and depending on the specific query, 3 to 8 times faster than executing queries on the ROLAP server. The query execution time on the HOLAP server was similar to that of ROLAP. This supports our assumption that, despite requiring more time for cube processing, the MOLAP server significantly outperforms the other types of data warehouses in query execution because it already contains the necessary data within the cube.

Moreover, adding aggregations had no noticeable effect on the ROLAP server's query execution time. However, it made MOLAP and HOLAP almost equally efficient, with aggregations speeding up MOLAP query execution by about 20 times and HOLAP by approximately 70 to 150 times. This highlights the effectiveness of using aggregations. Aggregations accelerate query execution because the cube doesn't need to spend time grouping measures by dimension attributes; this grouping is already done during cube processing. HOLAP's query execution time was significantly better because, during the cube processing, aggregations are transferred to the cube, eliminating the need to retrieve data from the relational warehouse during aggregation.

#### 6.3. General Conclusions

The tests on MOLAP, ROLAP, and HOLAP reveal key performance differences:

- 1. Cube Processing Time: MOLAP is slower to process the cube, as it requires transferring data into the multidimensional structure. ROLAP, however, directly accesses the relational database, making it faster. HOLAP offers a hybrid approach, with performance similar to ROLAP.
- 2. Query Execution Time: MOLAP accelerates query execution by storing data in the cube, significantly outperforming ROLAP. HOLAP, with aggregations, also sees major query performance improvements, as aggregated data is stored within the cube, bypassing the need for relational data access.
- 3. Impact of Aggregations: Adding aggregations enhances query performance in MOLAP and HOLAP, reducing query time by eliminating the need for grouping measures during execution. ROLAP shows no improvement as it still relies on the relational database for query execution.
- 4. Overall Performance: MOLAP excels in query speed but has slower processing times. ROLAP is faster for cube processing but slower in query execution. HOLAP offers a balanced solution, with faster query performance when aggregations are used, making it ideal for systems requiring both efficient data storage and quick queries.

In conclusion, the choice between OLAP models depends on specific needs: MOLAP is best for fast queries, ROLAP for quicker processing, and HOLAP offers a middle ground.