



Dremio Software

Query Performance Analysis and Improvement

Introduction

This document aims to help identify the reason for poor query performance and suggest ways to improve it. The document assumes the reader is familiar with Dremio, either as an administrator or a user, can run queries and navigate around Dremio's UI (specifically the "Jobs" page).

As with any data platform, the amount of memory and the number of processors has the highest impact on the response time of a query. Those resources together make up the Dremio Cluster. Note: data storage is not part of the Dremio Cluster (the main difference between a Data Warehouse and a Data Lakehouse). Of course, IO performance is also critical to any analytics platform.

Additionally, query performance is directly related to the complexity of the query (aggregations, joins, etc), the concurrent number of users running queries on a Dremio cluster and how the data is stored versus how it's scanned.

We use Dremio's Query Profile report (referred to in the Docs as "[Raw Profile](#)") to analyze the query, identify where most of the time is spent, potential bottlenecks, and discuss ways to improve the query.

Gather Dremio Cluster Information

It will be beneficial to have the following information handy. You can refer to these details while analyzing the query in the *Profile* report.

Cluster Size	No. of Nodes	No. of Cores	Memory (GB)
Coordination nodes			
Executor nodes			

Workload Details	Details
File format(s)	
File size	
Number of concurrent queries	

Dremio Cluster Components

Here are the roles for each node type:

- Coordinator(s) - responsible for handling connections, serving the UI and query planning. There can be more than one coordinator node (scale-out coordinator).
- Executors - running the queries, joining, sorting and combining result sets, metadata and reflection refresh

Query Profile

A query profile report is generated for each executed query. It contains query metrics that can be used to analyze query performance. To view a Query Profile report, click on the Jobs icon on the left side panel in the Dremio UI, then click on the Job ID of the query you want to analyze, then click on the Raw Profile link on the top right.

If you need to provide this information to Dremio Support, you generate the zip file using the 'Download Profile' button at the bottom of the page.

Search: kamran_hussain | Start Time: All | Status: | UI, +1 | User: | Manage Columns

Job ID	User	Dataset	Query Type	Queue	Start Time	Duration	SQL
...628-cd2d-5ca80f954d00	kamran.hussain@dremio...	trips	UI (download)	High Cost User Q...	10/17/2023, 12:28:24	00:00:04	CREATE TABLE "...dataset@download"...
...a22-c86a-4f058c07cd00	kamran.hussain@dremio...	my_taxi_trips_iceberg	UI (run)	Low Cost User Q...	10/25/2023, 06:09:59	00:00:03	delete from my_taxi_trips_iceberg where passen
...0a08-4a3a-2c213f7cc300	kamran.hussain@dremio...	NYC Trips	JDBC Client	High Cost User Q...	10/17/2023, 12:24:11	00:00:03	SELECT * FROM "Business.Transportation"."NYC T
...11fe-2c50-3389dc35ad00	kamran.hussain@dremio...	NY_Trips	UI (run)	Low Cost User Q...	10/25/2023, 06:08:25	00:00:03	create table "s3"."kamran-dremio"."ny_taxi_trip
...8a09-e33e-0a1f3ed69100	kamran.hussain@dremio...	time_travel	UI (run)	Low Cost User Q...	10/25/2023, 06:15:19	00:00:01	create vds "Solution Architect".RH_gp as selec
...8db-e68e-93e4240e8800	kamran.hussain@dremio...	my_taxi_trips_iceberg	UI (run)	Low Cost User Q...	10/25/2023, 06:05:30	<1s	---SELECT CURRENT_TIMESTAMP() --2023-09-05 14:
988e-37a1-61909c251200	kamran.hussain@dremio...	trips	UI (preview)	UI Previews	10/27/2023, 09:49:24	<1s	SELECT nested_0.vendor_id AS vendor_id, nested
1ac4320b-16f5-988a-37a1-61909c251200	kamran.hussain@dremio...	my_taxi_trips_iceberg...	UI (run)	Low Cost User Q...	10/25/2023, 06:13:08	<1s	select count(*) from my_taxi_trips_iceberg at
...fec-2e3c-96e8b9ed2e00	kamran.hussain@dremio...	mytaxi	UI (preview)	UI Previews	10/17/2023, 12:21:59	<1s	SELECT * FROM mytaxi
...7b6-a341-8e2119642800	kamran.hussain@dremio...	Unavailable	UI (run)	Low Cost User Q...	10/25/2023, 06:14:37	<1s	create vds gp as select * from my_taxi_trips_1

Jobs » 1ac4320b-16f5-988a-37a1-61909c251200 | Overview | SQL | Visual Profile | **Raw Profile**

Summary

Status: COMPLETED

Total Memory: 819.96 MB

CPU Used: 0136%

Query Type: UI (preview)

Start Time: 10/27/2023 09:49:24

Duration: <1s

Wait on Client: <1s

User: kamran.hussain@dremio...

Queue: UI Previews

Input: 14.19 MB / 97.2K Rows

Output: 20.15 MB / 82K Rows

Total Execution Time <1s (100%)

Pending: 1ms (0.12%)

Metadata Retrieval: 98ms (17.17%)

Planning: 81ms (19.84%)

Queued: 12ms (1.48%)

Execution Planning: 5ms (10.73%)

Starting: 4ms (10.49%)

Running: 669ms (100.19%)

Submitted SQL

```

1 SELECT nested_0.vendor_id AS vendor_id, nested_0.pickup_datetime AS pickup_datetime, nested_0.dropoff_date AS dropoff_date, nested_0.dropoff_datetime, nested_0.passenger_count AS passenger_count, nested_0.trip_distance_mi AS trip_distance_mi, nested_0.pickup_longitude AS pickup_longitude, nested_0.pickup_latitude AS pickup_latitude, nested_0.rate_code AS rate_code, nested_0.store_and_fwd_flag AS store_and_fwd_flag, nested_0.dropoff_longitude, nested_0.dropoff_latitude AS dropoff_latitude, nested_0.payment_type AS payment_type, nested_0.fare_amount AS fare, AS surcharge, nested_0.mta_tax AS mta_tax, nested_0.tip_amount AS tip_amount, nested_0.tolls_amount AS tolls_amount, nested_0.total_or AS total_or, nested_0.pickup_date AS pickup_date, join_weather_nyc."date" AS "date", join_weather_nyc.station AS station, join_weather_nyc.name AS AS aamd, join_weather_nyc.prcp AS prcp, join_weather_nyc.snow AS snow, join_weather_nyc.snowd AS snowd, join_weather_nyc.tempx AS temp join_weather_nyc.tempin AS tempin
2 FROM (

```

Queried Datasets

- trips (Preparation: trips)
- weather_nyc (pg_public:weather_nyc)

Scans

- trips_pickupDate

Acceleration

Reflections Used

Reflections Not Used

- Aggregation Reflection [AlgMatch] (Business.Transportation.NYC.Trips) - Did not cover query

Query Lifecycle - Raw Profile

From the time you hit run to the time the query returns data, it goes through several phases. Let's take a look at each phase of the query. This information is available in the Raw Profile. Right below the SQL on the Query tab of the Raw Profile is the **Job Summary** section (see screenshot below). This section of the Profile will show the times associated with each query phase, which will help you determine what to focus on to improve the performance by reducing that wait time. (start with the highest wait time).

Phase 1: **Command pool wait** - Dremio can run $n-1$ queries in parallel, where n is the number of cores on the coordinator node. If more than $n-1$ queries are executed simultaneously, then those queries go into the command pool. Command pool wait time should be at most 2-3 seconds.

Phase 2: **Metadata retrieval** - command parsing, dataset retrieval from source and KV store. Metadata is stored on the executor nodes. (starting with v18).

⚠ IMPORTANT

Metadata retrieval should be done before the query needs to be run. This step should be avoided during query execution time.

Phase 3: **Query planning time** - query planning should not take more than 500ms to 1 sec.

⚠ IMPORTANT

If you see a long planning time, it's directly related to the coordinator resources, query complexity or the number of concurrent queries.

Phase 4: **Queued** - waiting for resources.

⚠ IMPORTANT

If there is a lot of queueing, that would either mean that the concurrency is high or the queries are running for a long time. This would imply that the Dremio cluster is not correctly sized for the workload.

Phase 5: **Execution planning** - Dremio has the fully optimized plan at this point. Now, Dremio splits the plan and assigns it to the executor nodes.

Phase 6: **Starting** - propagates fragments (units of work) to the executors.

⚠ IMPORTANT

If this step takes more than a few seconds, the executors may be too busy with other work. This would imply that the Dremio cluster is not correctly sized for the workload and you need to add more executors.

Phase 7: **Running** - waits for fragment completion
Completed, Failed, or Canceled

Below is an example screenshot of the Job Summary information:

Job Summary

State:	COMPLETED
Coordinator:	dremio-master-0.dremio-cluster-pod.default.svc.cluster.local
Threads:	7
Command Pool Wait:	0ms
Total Query Time:	77,814ms
# Joins in user query:	1
# Joins in final plan:	1
Considered Reflections:	3
Matched Reflections:	0
Chosed Reflections:	0

State Durations

Pending:	0ms
Metadata Retrieval:	60ms
Planning:	127ms
Engine Start:	0ms
Queued:	12ms
Execution Planning:	5ms
Starting:	3ms
Running:	77,607ms

This query ran for over 77 seconds. The **Total Query Time of 77 seconds** consists of 60 ms of **Metadata Retrieval**, 127 ms of **Planning**, and 77,607 ms of **Running**, etc. As with any query tuning exercise, you must identify and address the state that takes the longest time.

The metadata retrieval, planning and queued times are within the expected range; however, the Running time is unacceptable. Further investigation needs to be done to determine why the query took 77 seconds in the **Running** state. This paper's [Optimization Concepts](#) section will cover some possible reasons for extended run times.

Please refer to the [Reading Dremio Job Profiles](#) paper to learn more about the information in Dremio profiles and where to find it.

Common reasons for slow performance

Given the various states of query execution mentioned in the previous section, here we will briefly examine some common reasons for slow performance for a few of those states.

If the **Command Pool Wait** is more than 1-2 seconds, then either there were too many queries initiated at the same time (for the cores available on the coordinator), or the query planning is taking too long, which is consuming the cores, due to the complexity of the query. Please review the resources available on the Coordinator node(s) and take corrective actions. Dremio supports scale-out coordinator nodes, which may be needed to reduce this wait.

Metadata Retrieval should not happen at query runtime (aka inline refresh). If the duration of this phase is high, the metadata is either unavailable or possibly expired. Please review the configuration for the Source and adjust it accordingly:

Click on the Datasets icon > right-click on the Source in question > Settings > Metadata tab on the left. Optionally, review the Doc “Administration > [Refreshing Metadata](#)” to avoid inline metadata refresh.

Details of the **Planning** phase are in the Raw Profile > Planning tab:



One possible reason for high **Planning** time is too many reflections on the datasets in the query. The planner has to analyze all the reflections, calculate the cost, compare the cost and choose the best reflection. This can take significant time, sometimes more than running the query and fetching the data. Scroll down in the Planning tab and look at the values of ‘Normalization’ and ‘Logical Planning’ (see screenshot below). If these are high, more than 10-20 seconds, then it points to having too many reflections. Please review how many reflections you have and if they are needed.

Starting with version 24.2, [reflection hints](#) are now available to control which reflections are considered or excluded. You can consider hints to accelerate reflection selection.

Let's briefly review a few of the steps of the **Planning** state:

1. Validate and Convert To Rel - Dremio looks at the physical and virtual datasets by pulling this information into the planning phase. These should typically be in milliseconds and not the cause for any slowness.
2. Logical Planning - reflection substitution, partition pruning, filter pushdown
3. Final Physical Plan (Optimized Plan) - Dremio determines how to divide the work into Phases (or units of work).

Scroll further in the report to get to the **Threads** section.

Threads

Overview																				
Phase	Weight	Thread Reporting	First Start	Last Start	First End	Last End	Min First-run	Avg First-run	Max First-run	Min Wall-clock	Avg Wall-clock	Max Wall-clock	Min Sleep	Avg Sleep	Max Sleep	Min Blocked	Avg Blocked	Max Blocked	Last Update	Progress
00-xx-xx	2	1 / 1	0.151s	0.151s	0.813s	0.813s	0.004s	0.004s	0.004s	0.662s	0.662s	0.662s	0.000s	0.000s	0.000s	0.653s	0.653s	0.653s	14:49:25	14:49:25
01-xx-xx	2	8 / 8	0.151s	0.163s	0.810s	0.822s	0.002s	0.003s	0.004s	0.657s	0.657s	0.659s	0.001s	0.008s	0.024s	0.607s	0.626s	0.638s	14:49:25	14:49:25
02-xx-xx	1	8 / 8	0.151s	0.163s	0.715s	0.822s	0.002s	0.006s	0.010s	0.560s	0.584s	0.659s	0.003s	0.005s	0.011s	0.002s	0.376s	0.510s	14:49:25	14:49:25
03-xx-xx	2	1 / 1	0.163s	0.163s	0.725s	0.725s	0.003s	0.003s	0.003s	0.562s	0.562s	0.562s	0.000s	0.000s	0.000s	0.494s	0.494s	0.494s	14:49:25	14:49:25
04-xx-xx	1	8 / 8	0.152s	0.164s	0.708s	0.744s	0.005s	0.009s	0.019s	0.555s	0.562s	0.580s	0.003s	0.009s	0.021s	0.009s	0.462s	0.549s	14:49:25	14:49:25
05-xx-xx	1	1 / 1	0.153s	0.153s	0.708s	0.708s	0.003s	0.003s	0.003s	0.555s	0.555s	0.555s	0.001s	0.001s	0.001s	0.497s	0.497s	0.497s	14:49:25	14:49:25

Each Phase can have different degrees of parallelism - how this is derived is discussed in the next section. The example above has 6 phases: 00, 01, ..., 05. Phases 00, 03 and 05 are single threaded. Look at the 'Thread Reporting' column, which shows 1/1. Phases 01, 02, and 04 have eight threads each.

Note: Phase 00 is always single-threaded; this is the phase that sends the results back to the client.

Degree of Parallelism

Based on the format of the data and the data source, the degree of parallelism in Dremio is determined in slightly different ways.

Parquet, ORC (read-only) and AVRO file formats fall into Dremio's unlimited splits flow, where the metadata is stored in the distributed store. During the Planning phase, the coordinator cannot access the number of splits, and partition pruning happens during the first part of query execution. In this case, the degree of parallelism is determined during the Execution Planning phase. It is calculated as the minimum of these two values:

- Number of cores on the executors (70% of executor cores for each query)
- Estimated number of rows. You can see each dataset's estimated number of rows in the Final Planning section. (search for 'rowcount')

For all other file and table formats and external sources, the metadata still resides in the Master Coordinator's KVstore and pruning happens during the Planning phase. In this case, the degree of parallelism is calculated as the minimum of these three values:

- Number of splits to be scanned (in Final Physical plan)
- Number of cores on the executors (70% of executor cores for each query)
- Estimated number of rows. You can see each dataset's estimated number of rows in the Final Planning section. (search for 'rowcount')

Final Physical Transformation (8 ms)

```

+): rowcount = 784414.6833656224, cumulative cost = {1.1867924750484336E7 rows, 8.277115391290028E7 cpu, 1.5733495667312449E7 i
+): rowcount = 784414.6833656224, cumulative cost = {1.1083510067118714E7 rows, 8.277107547143194E7 cpu, 1.5733495667312449E7 i
10612
. VARCHAR(65536) uploadcommenttext, DOUBLE CAST, INTEGER +, ANY E_X_P_R_H_A_S_H_F_I_E_L_D): rowcount = 784414.6833656224, cumula
INTEGER +): rowcount = 784414.6833656224, cumulative cost = {8730266.017021846 rows, 7.022028365464531E7 cpu, 1.5733495667312449
00000002 memory}, id = -1678540615

IT runid, DOUBLE CAST): rowcount = 211388.17778794913, cumulative cost = {1691105.422303593 rows, 1.3518314133293128E7 cpu, 4227

commenttext, BIGINT runid, DOUBLE CAST): rowcount = 211388.17778794913, cumulative cost = {1268329.0667276948 rows, 1.1827187572

3712'), =($2, '0452'), =($2, '2715'), =($2, '3798'), =($2, '2716'), =($2, '3815'), =($2, '1696'), =($2, '1897'), =($2, '3807'),
rtid, VARCHAR(65536) originalsourcesystemidentifier, VARCHAR(65536) sourcesystemidentifier, DECIMAL(38, 9) spotbalanceusdeamount
IT runid, DOUBLE CAST): rowcount = 573026.5055776733, cumulative cost = {4584212.044621387 rows, 3.664515390672827E7 cpu, 1.1460

```

Optimization Concepts

Optimization concept #1: Runtime Filtering

The remainder of this paper will focus on the Physical Plan. We're looking for optimization patterns that will ensure efficient execution.

```
Final Physical Transformation (1 ms)
00-00 Screen : rowType = RecordType(VARCHAR(65536) Fragment, BIGINT Records, VARCHAR(65536) Path, VARBINARY(65536) Metadata, INTEGER Partition, I
00-01 Project(Fragment=[0], Records=[1], Path=[2], Metadata=[3], Partition=[4], FileSize=[5], IcebergMetadata=[6], fileschema=[7], Par
00-02 WriterCommitter(final=[/opt/dremio/data/results/1ac4320b-18f5-988a-37a1-6f9092251200]) : rowType = RecordType(VARCHAR(65536) Fragment,
00-03 UnionExchange : rowType = RecordType(VARCHAR(65536) Fragment, BIGINT Records, VARCHAR(65536) Path, VARBINARY(65536) Metadata, INTEGEI
01-01 Writer : rowType = RecordType(VARCHAR(65536) Fragment, BIGINT Records, VARCHAR(65536) Path, VARBINARY(65536) Metadata, INTEGER Par
01-02 Project(vendor_id=[0], pickup_datetime=[1], dropoff_date=[2], dropoff_datetime=[3], passenger_count=[4], trip_distance_mi=[
01-03 Project(vendor_id=[0], pickup_datetime=[2], dropoff_date=[3], dropoff_datetime=[4], passenger_count=[5], trip_distance_mi
01-04 HashJoin(condition=[(=$1, $22)], joinType=[full]) : rowType = RecordType(VARCHAR(65536) vendor_id, DATE pickup_date, TIMESTAM
01-06 Project(vendor_id=[0], pickup_date=[1], pickup_datetime=[2], dropoff_date=[3], dropoff_datetime=[4], passenger_count=
01-08 HashToRandomExchange(dist0=[[$1]]) : rowType = RecordType(VARCHAR(65536) vendor_id, DATE pickup_date, TIMESTAMP(3) picku
02-01 Project(vendor_id=[0], pickup_date=[1], pickup_datetime=[2], dropoff_date=[3], dropoff_datetime=[4], passenger_co
02-02 Project(vendor_id=[0], pickup_date=[1], pickup_datetime=[2], dropoff_date=[3], dropoff_datetime=[4], passenger_
02-03 SelectionVectorRemover : rowType = RecordType(VARCHAR(65536) vendor_id, DATE pickup_date, TIMESTAMP(3) pickup_date
02-04 Filter(condition=[AND(>=($6, 0), <=($6, 360))] ) : rowType = RecordType(VARCHAR(65536) vendor_id, DATE pickup_da
02-05 Limit(offset=[0:BIGINT], fetch=[10000:BIGINT]) : rowType = RecordType(VARCHAR(65536) vendor_id, DATE pickup_da
02-06 TableFunction(columns=[vendor_id` `pickup_date` `pickup_datetime` `dropoff_date` `dropoff_datetime` `n
```

NOTE
If your query does an inner or right outer join - your query should use runtime filtering.

Runtime filtering improves query performance by reducing the amount of data that will be projected when joining fact and dimension tables. Dremio will quickly determine the join keys from the dimension table and project only the matching data in the fact table, making this step much faster. This is only applicable for inner joins or right outer joins.

Let's now find out if the runtime filter was used (or not used). Click on the Profile tab > go to the Planning tab > and search for 'runtime' (see highlighted row in the screenshot below).

```
Final Physical Transformation (108 ms)
01-05 HashAgg(group=[{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, :
01-06 Project(TIME_INV_DESC=[0], FACILITY_ID=[1], CLNT_UCN=[2], CLNT_NM=[3], RCV_PAYMENT_CCY=[4], PAY
01-07 HashToRandomExchange(dist0=[[$0]], dist1=[[$1]], dist2=[[$2]], dist3=[[$3]], dist4=[[$4]], dist5=[
02-01 UnorderedMuxExchange : rowType = RecordType(VARCHAR(65536) TIME_INV_DESC, DECIMAL(28, 0) FACILIT
03-01 Project(TIME_INV_DESC=[0], FACILITY_ID=[1], CLNT_UCN=[2], CLNT_NM=[3], RCV_PAYMENT_CCY=[4]
03-02 Project(TIME_INV_DESC=[60], FACILITY_ID=[0], CLNT_UCN=[62], CLNT_NM=[61], RCV_PAYMENT_CC
03-03 Project(FACILITY_ID=[0], CALCULATION_EXCEPTION=[2], SOURCE_SYSTEM=[3], RECALC_MTM_AMT=[
03-04 HashJoin(condition=[(=$1, $66)], joinType=[inner], runtimeFilter=[314-03 dir0]) : rowT
03-06 Project(FACILITY_ID=[0], TIME_INV_KEY0=[1], CALCULATION_EXCEPTION=[2], SOURCE_SYSTE
03-07 SelectionVectorRemover : rowType = RecordType(DECIMAL(28, 0) FACILITY_ID, VARCHAR(65
03-08 Filter(condition=[OR(IS NULL($58), CAST(AND(IS NOT NULL($68), <>($67, 0))):BOOLEAN
03-09 Project(FACILITY_ID=[0], TIME_INV_KEY0=[1], CALCULATION_EXCEPTION=[2], SOURCE
03-10 Project(FACILITY_ID=[0], TIME_INV_KEY0=[1], CALCULATION_EXCEPTION=[3], SOUR
03-11 HashJoin(condition=[AND(=( $4, $71), =( $2, $70))] ), joinType=[inner], runtimeF
03-13 Project(FACILITY_ID=[0], TIME_INV_KEY0=[1], TIME_INV_KEY1=[2], CALCULAT
03-15 HashToRandomExchange(dist0=[[$4]], dist1=[[$2]]) : rowType = RecordType(I
```

This example shows that the runtime filter is being used; it is a HashJoin, join type=[inner] join, and the column name used to filter is dir0. The scan that benefited from this filter was in Phase

314-03. (we'll use this phase to drill deeper into the operation to see the benefits of runtime filtering).

Referring to the screenshots below, scroll down to the **Operators** section > Find and expand the phase from the 'runtimeFilter' > then Expand **Operator Metrics** (in this example, it's Phase 314-03). Scroll to the right until you see NUM_PARTITION_PRUNED. Because of runtime filtering, the query is optimized. Dremio can skip this many row groups (approx. 30 per phase x number of phases = 80 in this example).

Profile: 1 ▾ Download

Overview Profile Visual Profile *

Thread	SETUP_NS	NUM_READERS	NUM_REMOTE_READERS	NUM_ROW_GROUPS	NUM_VECTORIZED_COLUMNS	NUM_NON_VECTORIZED_COLUMNS	COPY_NS	FILTER_NS	PARQUET_EXEC_PATH	FILTER_EXIST
314-85-03	0.008s	0.007s	0.105s	7	23,808	3MB	iaasn00297746.svr.us.jpurchase.net			3,290,722
314-86-03	0.005s	0.010s	0.209s	11	36,412	736KB	iaasn00297790.svr.us.jpurchase.net			3,705,965
314-87-03	0.074s	0.072s	0.062s	36	130,263	3MB	iaasn00297754.svr.us.jpurchase.net			1,811,787

Operator Metrics

Thread	SETUP_NS	NUM_READERS	NUM_REMOTE_READERS	NUM_ROW_GROUPS	NUM_VECTORIZED_COLUMNS	NUM_NON_VECTORIZED_COLUMNS	COPY_NS	FILTER_NS	PARQUET_EXEC_PATH	FILTER_EXIST
314-00-03		2		2	2	0				2
314-01-03		8		8	2	0				2
314-02-03		9		9	2	0				2
314-03-03		2		2	2	0				2

Profile: 1 ▾ Download

Overview Profile Visual Profile *

IS	NUM_IO_READ	NUM_HIVE_PARQUET_DECIMAL_COERCIONS	NUM_ROW_GROUPS_TRIMMED	NUM_COLUMNS_TRIMMED	NUM_PARTITIONS_PRUNED	NUM_BOOSTED_FILE_READS	MAX_BOOSTED_FILE_READ_TIME_NS
			1	111	34	0	0
			7	777	28	0	0
			8	888	26	0	0
			1	111	34	0	0
			1	111	35	0	0
			1	111	33	0	0
			24	2,664	7	0	0
			4	444	29	0	0

Optimization concept #2: Partition Pruning:

Every query should do partition pruning. If it is not, then it needs to be investigated.

To see if the query is doing partition pruning, let's look at the physical splits for the dataset and then compare that to the number of splits scanned to retrieve the data for this query.

To get the physical splits for the dataset, look at the raw metadata in the 'Convert To Rel' section in the **Planning** tab. Search for 'splits' - in the example below, the dataset has 214,109 physical splits.

Query and Planning

Query Visualized Plan **Planning** Acceleration

Validation (6 ms)

Convert To Rel (43 ms)

```

=[$30])

counterpartylegalentity', 'currency', 'bookinglegalentity', 'ledgeraccount', '_record_key', 'runid'], splits=[214109])

id'], splits=[149704])

```

Scroll left and capture the table for the dataset.

Query and Planning

Query Visualized Plan **Planning** Acceleration

Validation (6 ms)

Convert To Rel (43 ms)

```

LogicalFilter(condition=[($32, 10237, -($32, 100977))
ExpansionNode(path=[EPS_FRI.com_jpmc_ct_mis_balancesheetaverages_base_v])
LogicalProject(as_of_date=[$34], affected_feed=[$38], _abexid=[$0], _runtimeid=[$1], accrualtypecode=[$2], adjustmentidentifier=[$3],
LogicalJoin(condition=[($30, $31)], joinType=[inner])
LogicalUnion(all=[true])
LogicalProject(_abexid=[$0], _runtimeid=[$1], accrualtypecode=[$2], adjustmentidentifier=[$3], adjustmentreasonext=[$4],
ExpansionNode(path=[EPS_FRI.source.com_jpmc_ct_mis_balancesheetaverages_cold])
LogicalProject(_abexid=[$0], _runtimeid=[$1], accrualtypecode=[$2], adjustmentidentifier=[$3], adjustmentreasonext=[$4],
ScanCrel(table=["SFP-HIVE".db_103118_eps_cons_priv.com_jpmc_ct_mis_balancesheetaverages_cold], columns=['_abexid',
LogicalProject(_abexid=[$0], _runtimeid=[$1], accrualtypecode=[$2], adjustmentidentifier=[$3], adjustmentreasonext=[$4],
ExpansionNode(path=[EPS_FRI.source.com_jpmc_ct_mis_balancesheetaverages_hot])
LogicalProject(_abexid=[$0], _runtimeid=[$1], accrualtypecode=[$2], adjustmentidentifier=[$3], adjustmentreasonext=[$4],
ScanCrel(table=["SFP-HIVE".db_103118_eps_cons_priv.com_jpmc_ct_mis_balancesheetaverages_hot], columns=['_abexid', '
ExpansionNode(path=[EPS_FRI.source.com_jpmc_ct_mis_balancesheetaverages_metadata])
LogicalProject(runid=[$0], ingest_mode=[$1], dlp_event_time=[$2], as_of_date=[CAST($3):DATE], validations=[$4], meta_info=[$5],
ScanCrel(table=["SFP-HIVE".db_103118_eps_cons_priv.com_jpmc_ct_mis_balancesheetaverages_metadata], columns=['runid', 'i
ExpansionNode(path=[EPS_FRI.com_innc_ct_refdata.businesswithiereserve.ctb]ok ul)

```

Next, scroll down to the Physical Plan (Final Physical Transformation) and search for that table name to get the number of splits scanned for that table. Scroll to the right and see the value of 'splits':

Final Physical Transformation (8 ms)

```
adjustmentreasontext, VARCHAR(65536) bookingbusinessunitidentifier, VARCHAR(65536) bookinglegalentityidentifier,
VARCHAR(65536) originalsourcesystemidentifier, VARCHAR(65536) sourcesystemidentifier, VARCHAR(65536) uploadcommenttext,
VARCHAR(65536) sourcesystemidentifier, VARCHAR(65536) uploadcommenttext, BIGINT runid, DOUBLE CAST): rowcount = 784414.68:
VARCHAR(65536) adjustmentreasontext, VARCHAR(65536) bookingbusinessunitidentifier, VARCHAR(65536) bookinglegalentity:
VARCHAR(65536) sourcesystemidentifier, VARCHAR(65536) uploadcommenttext, BIGINT runid, DOUBLE CAST): rowcount = 211388.:
RecordType(VARCHAR(65536) adjustmentreasontext, VARCHAR(65536) bookingbusinessunitidentifier, VARCHAR(65536) bookinglegalentity:
VARCHAR(65536) sourcesystemidentifier, DECIMAL(38, 9) spotbalanceusdeamount, VARCHAR(65536) uploadcommenttext, VARCHAR(65536)
'3208'), =($2, '3982'), =($2, '0847'), =($2, '1697'), =($2, '3767'), =($2, '3725'), =($2, '3760'), =($2, '0449'
ntifier', `spotbalanceusdeamount`, `uploadcommenttext`, `runid`) splits=[36], mode=[NATIVE_PARQUET] : rowType
VARCHAR(65536) adjustmentreasontext, VARCHAR(65536) bookingbusinessunitidentifier, VARCHAR(65536) bookinglegalentity:
VARCHAR(65536) sourcesystemidentifier, VARCHAR(65536) uploadcommenttext, BIGINT runid, DOUBLE CAST): rowcount = 573026.:
RecordType(VARCHAR(65536) adjustmentreasontext, VARCHAR(65536) bookingbusinessunitidentifier, VARCHAR(65536) bookinglegalentity:
VARCHAR(65536) sourcesystemidentifier, DECIMAL(38, 9) spotbalanceusdeamount, VARCHAR(65536) uploadcommenttext, VARCHAR(65536)
'3208'), =($2, '3982'), =($2, '0847'), =($2, '1697'), =($2, '3767'), =($2, '3725'), =($2, '3760'), =($2, '0449'
ntifier', `spotbalanceusdeamount`, `uploadcommenttext`, `runid`], splits=[1], mode=[NATIVE_PARQUET]) : rowType =
```

In this example, Dremio scanned only 36 splits out of 214,109 physical splits, resulting in a much faster query response time.

Optimization concept #3: Filter Push Down

Filter Push Down is an optimization technique that minimizes the data that must be scanned, improving query performance. It's a feature of Dremio and nothing explicit has to be done by the data engineer to use this feature. Filter push-down is not as optimal as the other optimization methods discussed earlier - because Dremio will still need to open the splits and read the footer.

Let's look at a scenario when you are querying a dataset using a filter column, but the data is not partitioned by that column. In this case, the number of splits scanned will not be reduced, but the number of parquet bytes read should be less.

To identify a 'filter push down' in a query Profile, search for 'filters=' in the Planning tab > Final Physical Plan section. Look at the number of splits, which are the splits to be scanned, and compare it to the number of splits from the 'Convert to Rel' section for the same table, which shows the actual number of splits for that table. In the example below, you can see it's 920 in both cases. This means that there was no reduction in splits.

This section shows the actual number of splits for that table.

Convert To Rel (3,855 ms)

```

DERIVS_DIM_TIME_LOV))
DESC=[{2}, TIME_INV_KEY=[{3}], ACTUAL_TIME_INV_KEY=[{4}], ACTUAL_TIME_INV_DESC=[{5}], SNAPSHOT_TYPE=[{6}], DIR0=[{7}]

columns=[`SLN0`, `DATAMART`, `TIME_INV_DESC`, `TIME_INV_KEY`, `ACTUAL_TIME_INV_KEY`, `ACTUAL_TIME_INV_DESC`, `SNAPSHOT_TYPE`, `dir0`], splits=[920])
))

%'}})
LOV], columns=[`SLN0`, `DATAMART`, `TIME_INV_DESC`, `TIME_INV_KEY`, `ACTUAL_TIME_INV_KEY`, `ACTUAL_TIME_INV_DESC`, `SNAPSHOT_TYPE`, `dir0`], splits=!
```

```

DERIVS_DIM_TIME_LOV))
DESC=[{2}, TIME_INV_KEY=[{3}], ACTUAL_TIME_INV_KEY=[{4}], ACTUAL_TIME_INV_DESC=[{5}], SNAPSHOT_TYPE=[{6}], DIR0=[{7}]
```

This section shows the number of scanned splits for that table.

Final Physical Transformation (108 ms)

```

unionexchange : rowType = RecordType(VARCHAR(65536) dir0): rowcount = 92.0, cumulative cost = 1104.0
  ParquetScan(table=[CRRTDW.PUBLIC_USERS_FLAT_VERSION], columns=[`dir0`], splits=[91]) : rowType = Re
  ve = RecordType(VARCHAR(65536) ACTUAL_TIME_INV_KEY): rowcount = 1159091.9999999998, cumulative cost = {2457521
  rowType = RecordType(VARCHAR(65536) ACTUAL_TIME_INV_KEY): rowcount = 13171.499999999996, cumulative cost = {12
  ve = RecordType(VARCHAR(65536) ACTUAL_TIME_INV_KEY): rowcount = 131715.0, cumulative cost = {1166714.760869565
  INV_KEY=[{0}]} : rowType = RecordType(VARCHAR(65536) ACTUAL_TIME_INV_KEY): rowcount = 131715.0, cumulative cos
  =[{1}, {2}], joinType=[inner], runtimeFilter=[{247-05 dir0}]) : rowType = RecordType(VARCHAR(65536) ACTUAL_T
  TIME_INV_KEY=[{1}], dir0=[{2}]) : rowType = RecordType(VARCHAR(65536) ACTUAL_TIME_INV_KEY, VARCHAR(65536) dir0):
  ble=[CRRTDW.DERIVS_DIM_TIME_LOV], columns=[`TIME_INV_KEY`, `ACTUAL_TIME_INV_KEY`, `dir0`], splits=[920], fi
  le : rowType = RecordType(VARCHAR(65536) EXPR$0): rowcount = 88.0, cumulative cost = {508051.76086956525 rows,
  ip=[{}], EXPR$0=[MAX({0})]} : rowType = RecordType(VARCHAR(65536) EXPR$0): rowcount = 1.0, cumulative cost = {5
  ie : rowType = RecordType(VARCHAR(65536) EXPR$0): rowcount = 1.0, cumulative cost = {507962.76086956525 rows,
  group=[{}], EXPR$0=[MAX({0})]} : rowType = RecordType(VARCHAR(65536) EXPR$0): rowcount = 1.0, cumulative cost
  can(table=[CRRTDW.DERIVS_DIM_TIME_LOV], columns=[`dir0`], splits=[887]) : rowType = RecordType(VARCHAR(65536)
  _EXCEPTION=[{1}], CLIENT_ID=[{2}], LE_CLIENT_ID=[{3}], SOURCE_SYSTEM=[{4}], COB_DATE=[{5}], RECALC_MTM_AMT=[{6}], S
  N_EXCEPTION=[{1}], CLIENT_ID=[{2}], LE_CLIENT_ID=[{3}], SOURCE_SYSTEM=[{4}], COB_DATE=[{5}], RECALC_MTM_AMT=[{6}],
```

However the query is still optimized because Dremio looks at the parquet file footer to see if that column exists in that file and only scans that file. Dremio will also know which row groups

have that column and only scan those row groups. You can see this in the Profile by looking at the number of parquet bytes read. Determine the phase where the filter happens > go to Operators > for that scan, expand the Operator Metrics and see the value of `parquet_bytes_read`.

Operator Metrics						
<code>_NON_VECTORIZED_COLUMNS</code>	<code>COPY_NS</code>	<code>FILTER_NS</code>	<code>PARQUET_EXEC_PATH</code>	<code>FILTER_EXISTS</code>	<code>PARQUET_BYTES_READ</code>	<code>NUM_ASYNC_STREAMS</code>
0			2	1	19,720,055	8

It's essential to understand the impact of a filter on a non-partition column if there are **too many small files**. In this case, Dremio has to look at the footer of all the small files, which will adversely impact performance.

Conclusion

In this document, we introduced some key concepts to analyze and tune queries in Dremio. We started by looking at the metrics in the query profile report and took action accordingly.

Remember, looking at the whole cluster is essential, as other queries may impact the query. To minimize the impact of other workloads, you can set up engines to isolate workloads. Please review the [Workload Management](#) section in the Docs.

We recommend having a Dremio Professional Services resource do a thorough health check of your cluster at least once yearly. Please reach out to your account executive to request a health check.