

# Serverless Analytics

Using Synapse Distributed Query Processor

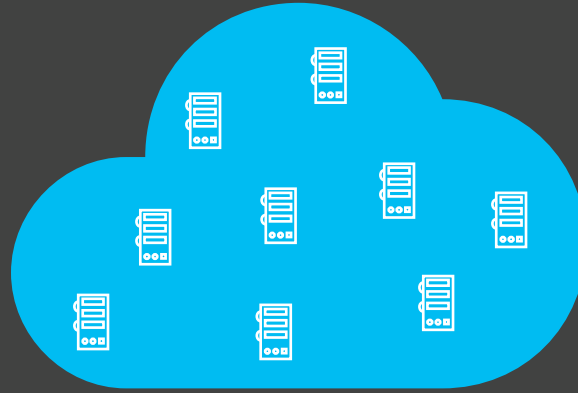
CE, CSA, MCS & Architects - Global Data & AI Team Meeting

Devin Jaiswal – Data SQL Ninja Team

Abraham Samuel – Global Black Belt Team

Today's customer, before using Azure resource, first, they need to size the Azure resources, customize it to their needs, and then use it.

Serverless is the one step closer, where they just bring the data, and start using it.



Serverless, the next gen platform



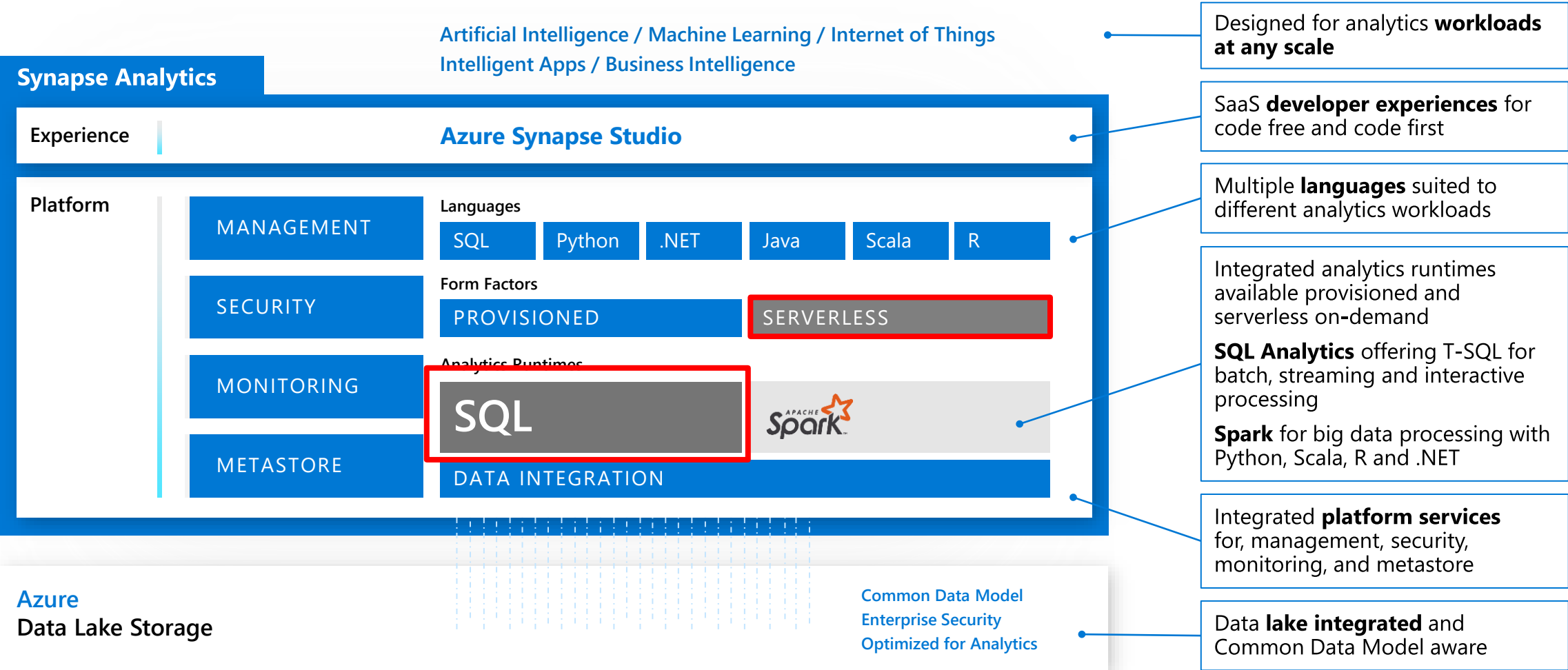
The "evolution" of database platforms

# Objectives

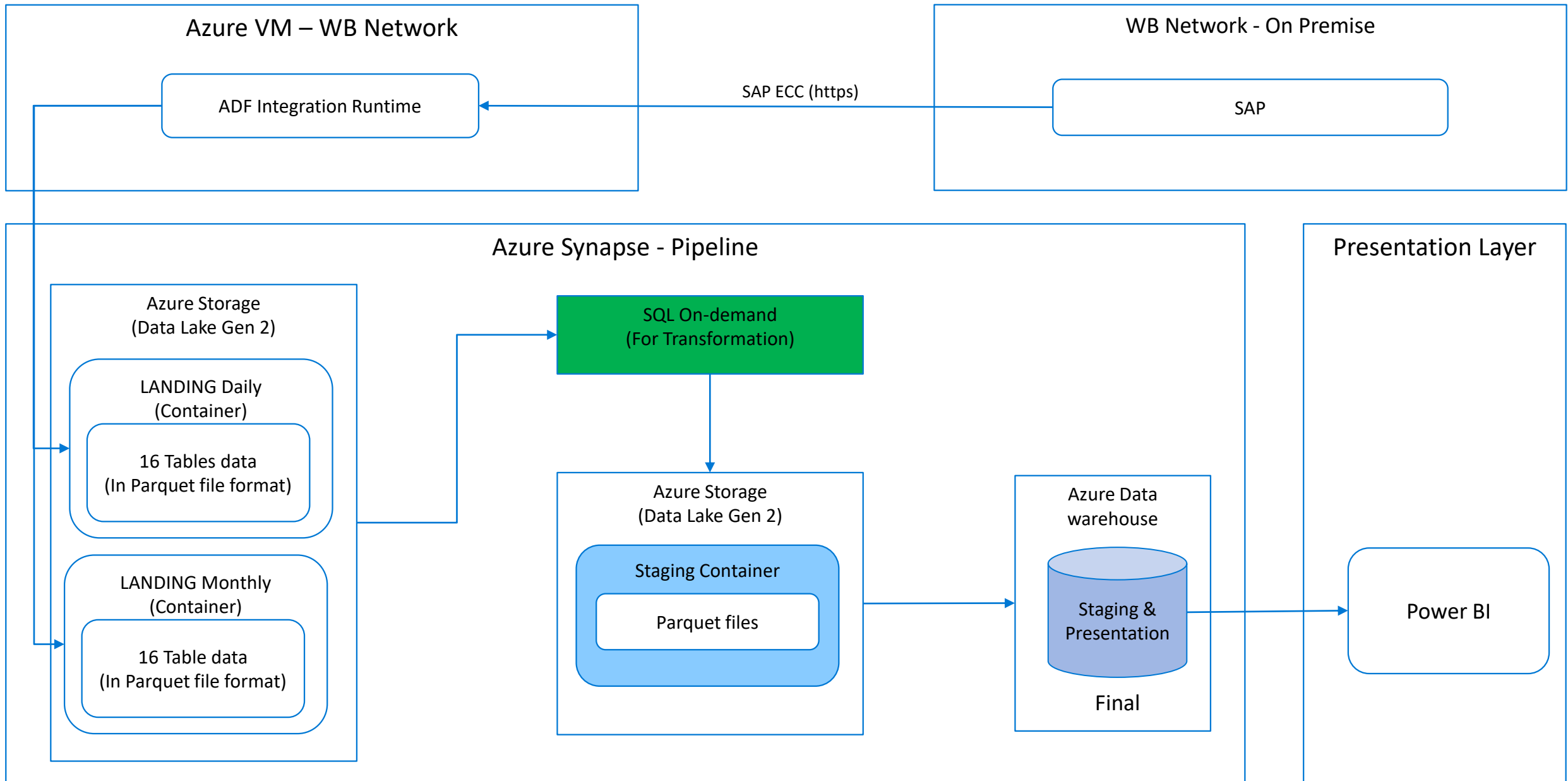
- Synapse DQP Technical deep dive
- Best practices
- Demo

# Azure Synapse Analytics

Integrated data platform for BI, AI and continuous intelligence



# LIS To Azure Synapse – Data flow



# Synapse SQL on-demand scenarios

## Discovery and exploration

What's in this file? How many rows are there? What's the max value?

**SQL On-demand reduces data lake exploration to the right-click!**

## Data transformation

How to transform the raw data?

**Use the full power of T-SQL to transform the data in the data lake**

# SQL Serverless – Querying on storage Demo

The image consists of two screenshots of the Microsoft Azure Synapse Analytics interface, demonstrating SQL querying on storage.

**Left Screenshot:** The 'Data' view is active. A context menu is open for a file named 'part-00133...'. The menu options include: New SQL script, New notebook, Copy ABFS path, Manage Access..., Rename..., Download, Delete, and Properties... The 'New SQL script' option is highlighted.

**Right Screenshot:** The 'SQL script 2' editor is active. The query is as follows:

```
1 SELECT
2   TOP 100 *
3 FROM
4   OPENROWSET(
5     BULK 'https://prlangaddemosa.dfs.core.windows.net/nyctlc/yellow/puYear=2015/puMonth=3/part-00133-tid-210938564719836543-aea5b543-5e83-4a7d-8d31-69f72c50b05d-15253-1.c000.snappy.parquet'
6     FORMAT='PARQUET'
7   ) AS nyc;
8
```

The 'Run' button is highlighted, and the 'SQL Analytics on-demand' option is selected in the 'Connect to' dropdown menu. The 'Results' tab is active, showing a table of taxi trip data.

VENDORID	TPEPPICKUPDATETIME	TPEPDROPOFFDATETIME	PASSENGERCOUNT	TRIPDISTANCE	PULOCATIONID	DOLOCATIONID	STARTLON	STARTLAT	ENDLON	ENDLAT
2	2015-02-28T23:5...	2015-03-01T00:0...	6	1.63	NULL	NULL	-74.000846862793	40.7306938171387	-73.977653503418	40.7631607055664
1	2015-03-28T19:2...	2015-03-28T19:2...	1	2.2	NULL	NULL	-73.977653503418	40.7631607055664	-73.96012878417...	40.7621574401855
2	2015-02-28T23:5...	2015-03-01T00:1...	5	3.23	NULL	NULL	-73.96012878417...	40.7621574401855	-73.98143005371...	40.7815055847168
1	2015-03-28T19:2...	2015-03-28T19:3...	1	2.1	NULL	NULL	-73.98143005371...	40.7815055847168	-73.98373413085...	40.7497062683105
2	2015-02-28T23:5...	2015-03-01T00:1...	1	3.52	NULL	NULL	-73.98373413085...	40.7497062683105	-73.98373413085...	40.7497062683105

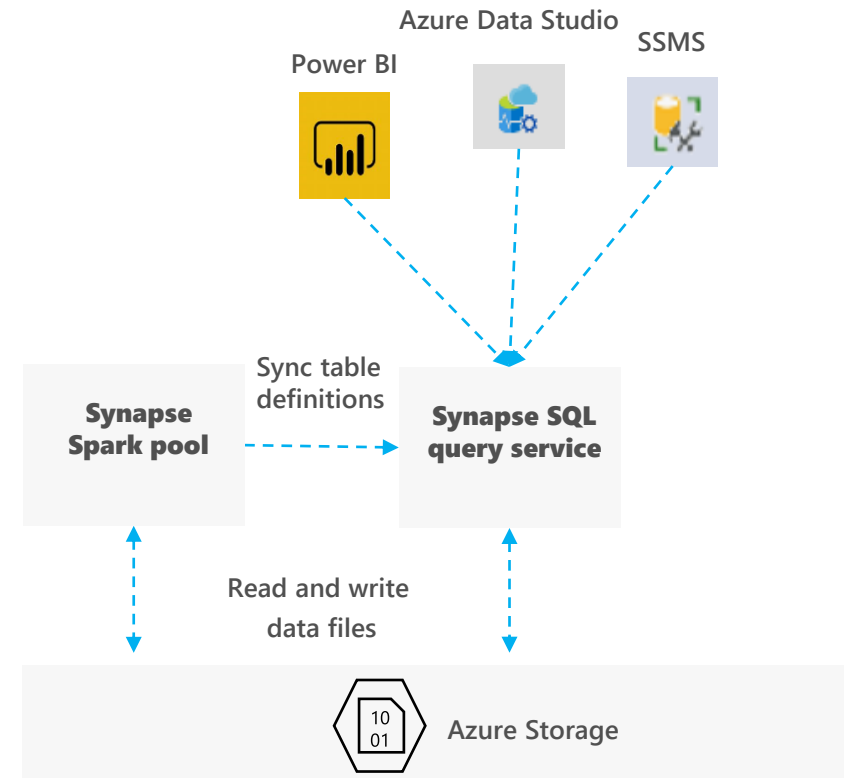
The status bar at the bottom indicates: 00:01:00 Query executed successfully.

# Synapse serverless SQL pool

An interactive query service that enables you to use standard T-SQL queries over files in Azure storage.

## The Serverless Pool uses separation of Compute and State:

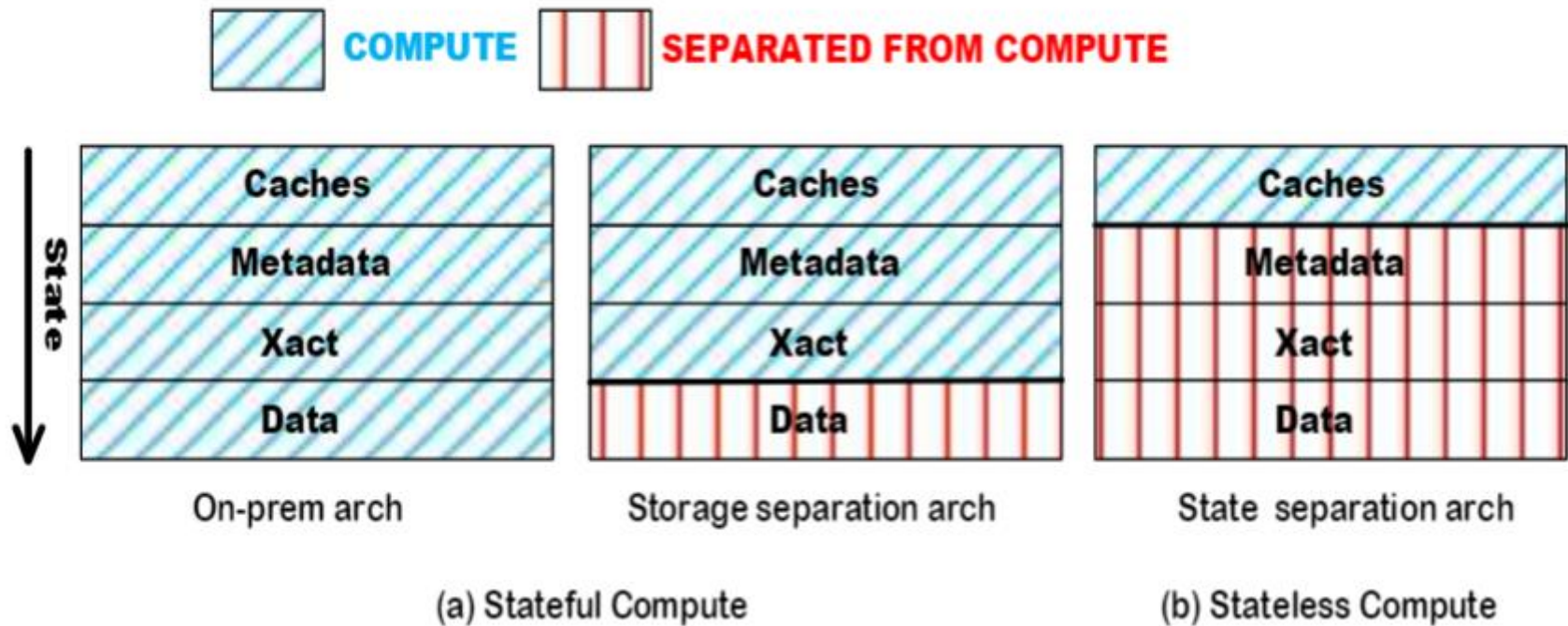
- **Stateless, resilient compute**
- **A logical database storage model, rather than physical**
- **Seamless scalability**



The Dedicated Pool uses separation of Compute and Storage



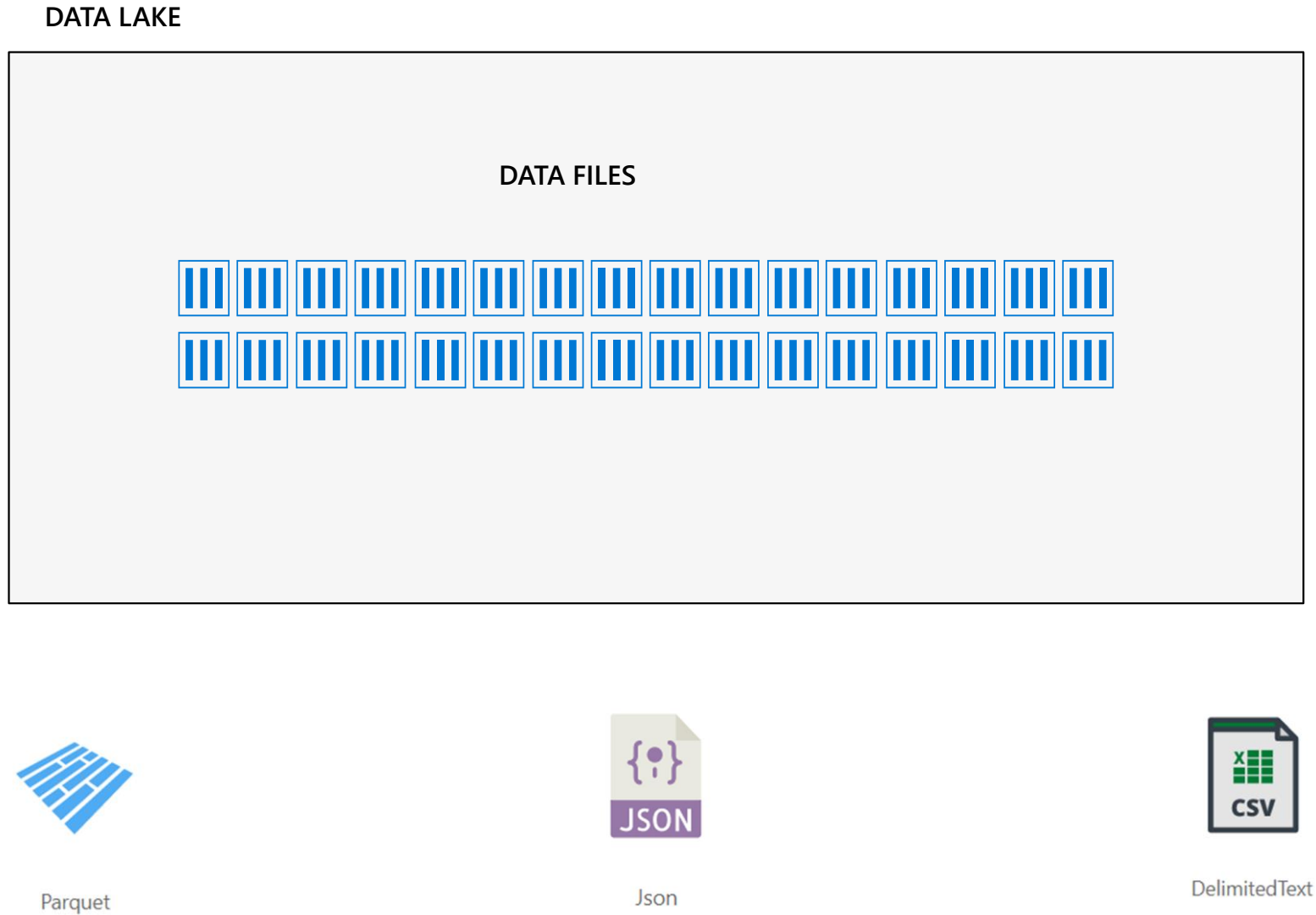
# Separating Compute and State



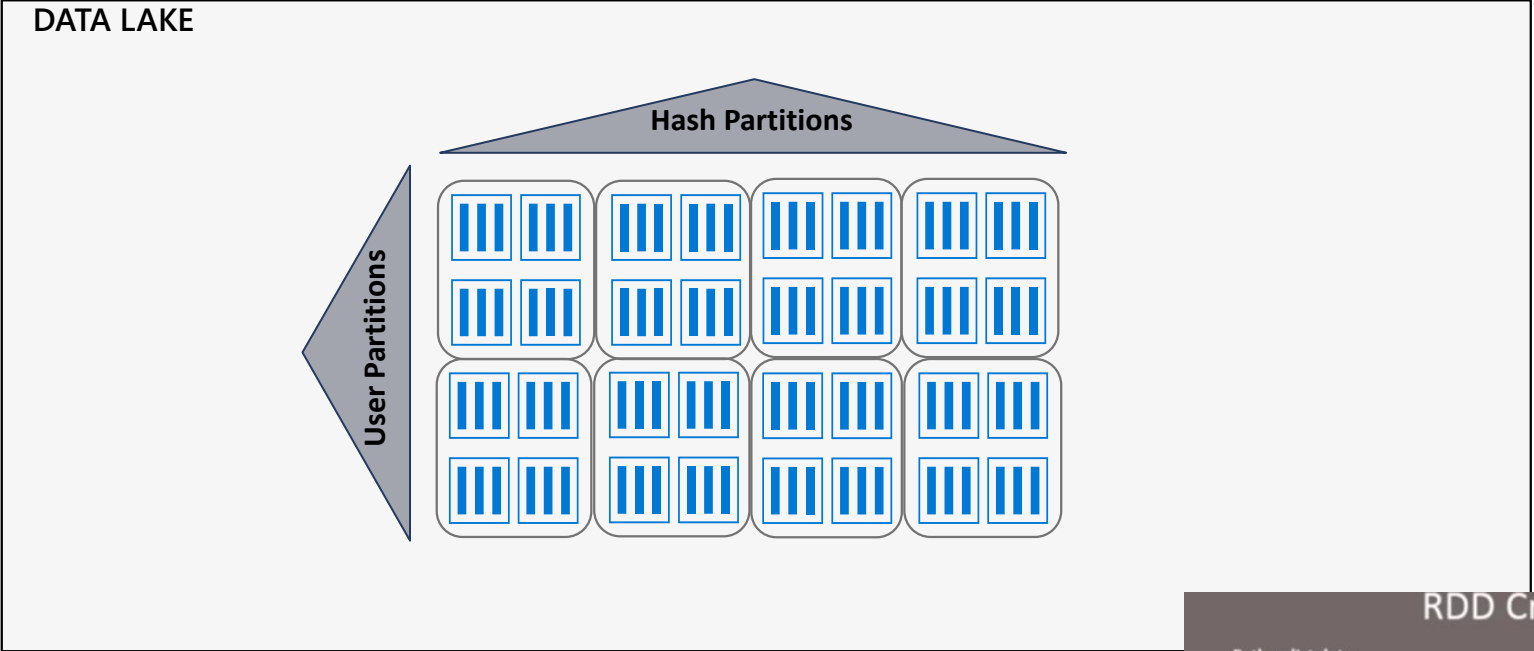
Evolution of data warehouse architectures over the years, illustrating how state has been coupled with compute.

# Data to Data Cells

A collection (e.g., table) of data objects (e.g., rows) in serverless pool can be logically abstracted as a collection of cells



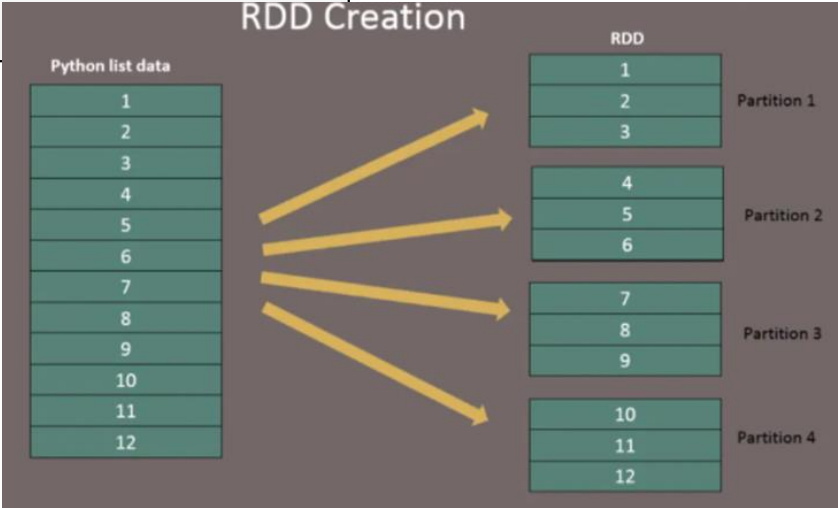
# Distribute the data



Dataset must be uniformly distributed across a large number of cells

We achieve that using Hash Partition, The hash-distribution  $h$  is used to map cells to compute nodes.

Similar concept as Spark RDD



The metadata and transactional log state is off-loaded to centralized services.

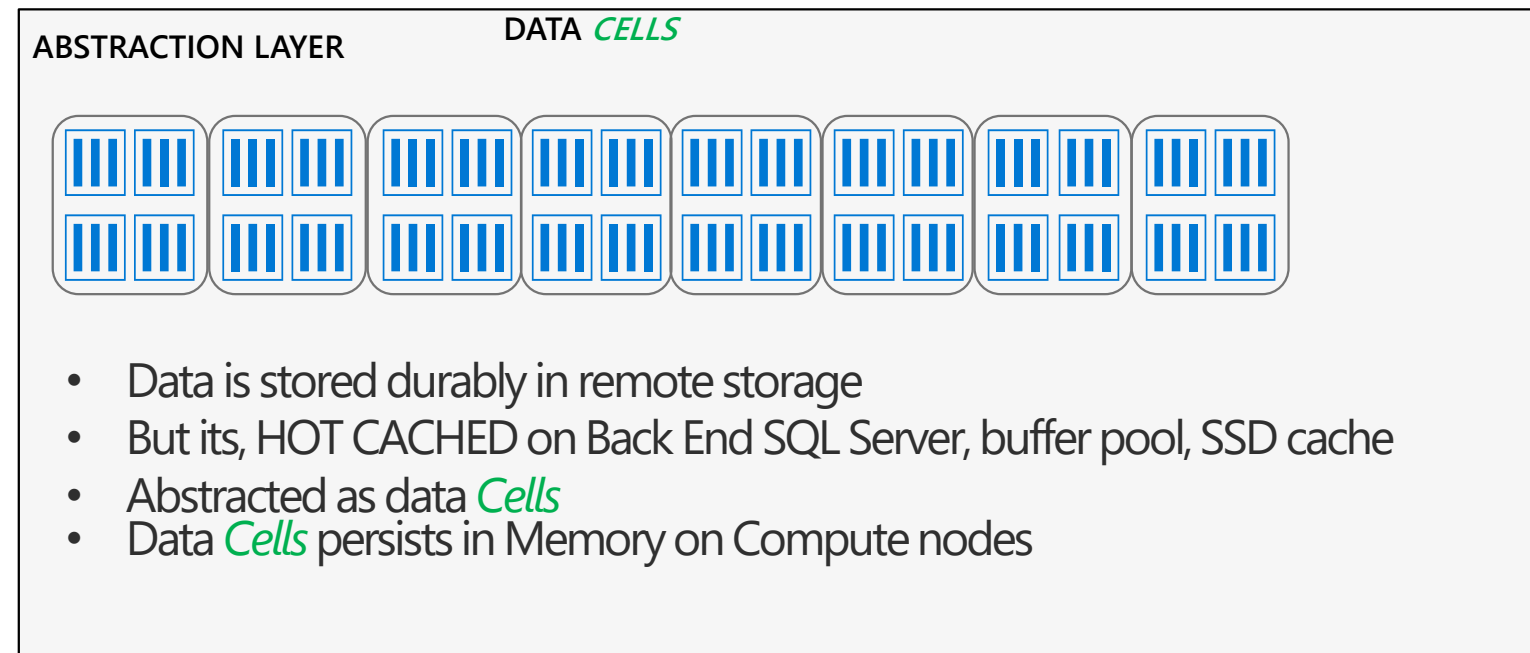
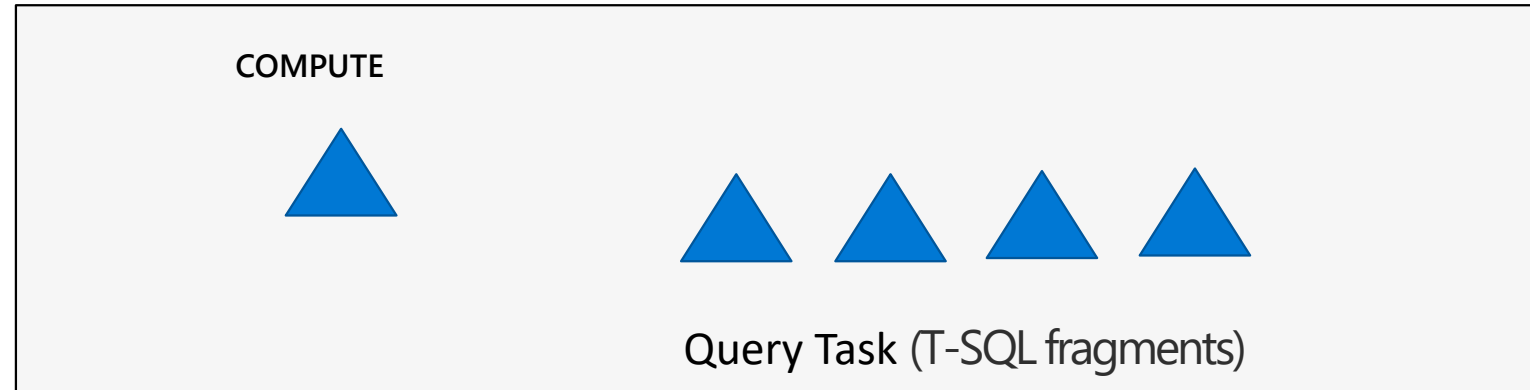
The SQL Server Front End is the service responsible for authentication, and metadata (bind metadata to data cells)



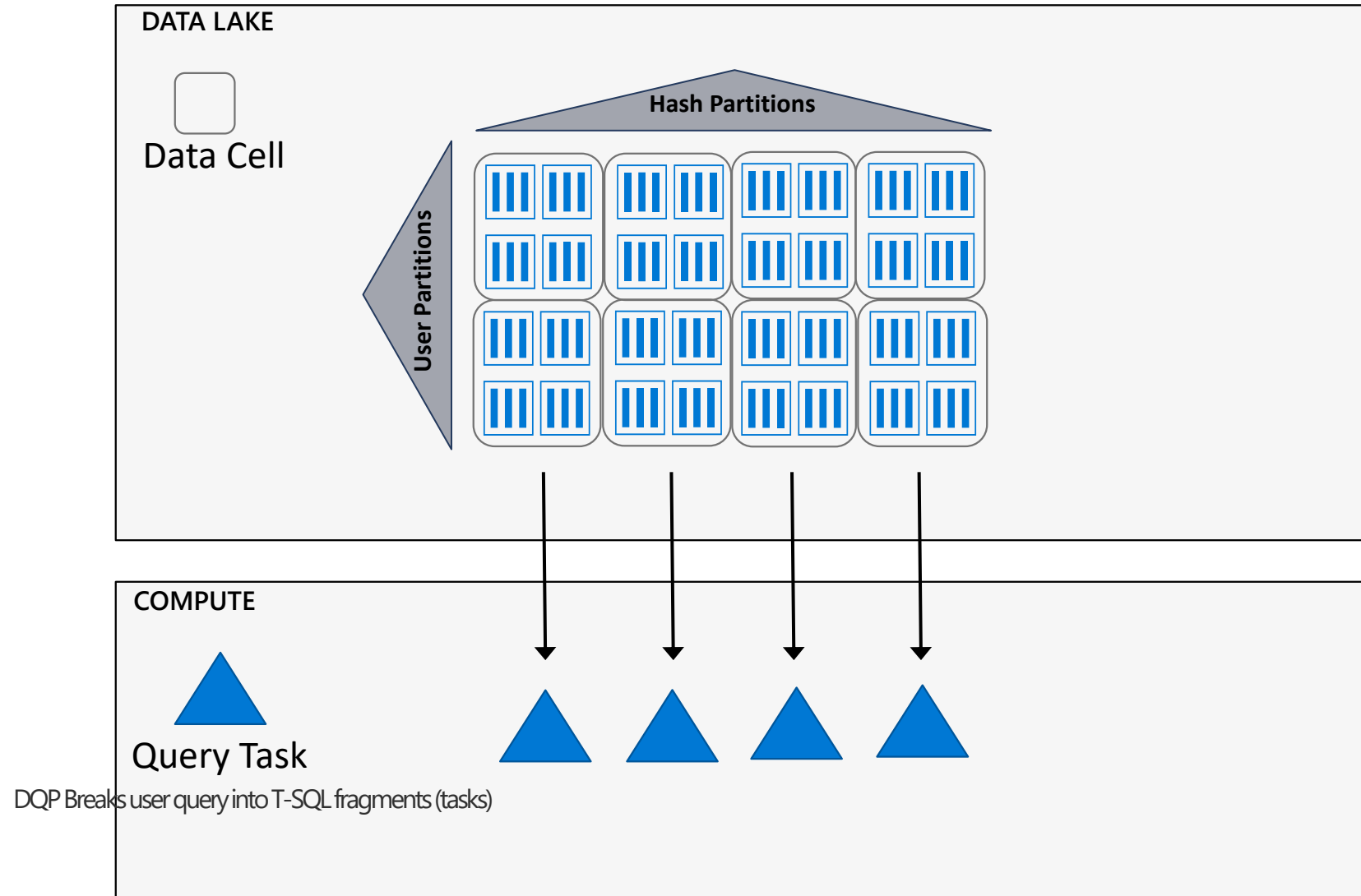
Data is stored durably in remote storage, and HOT CACHED on Compute Server. Abstracted as data cells.



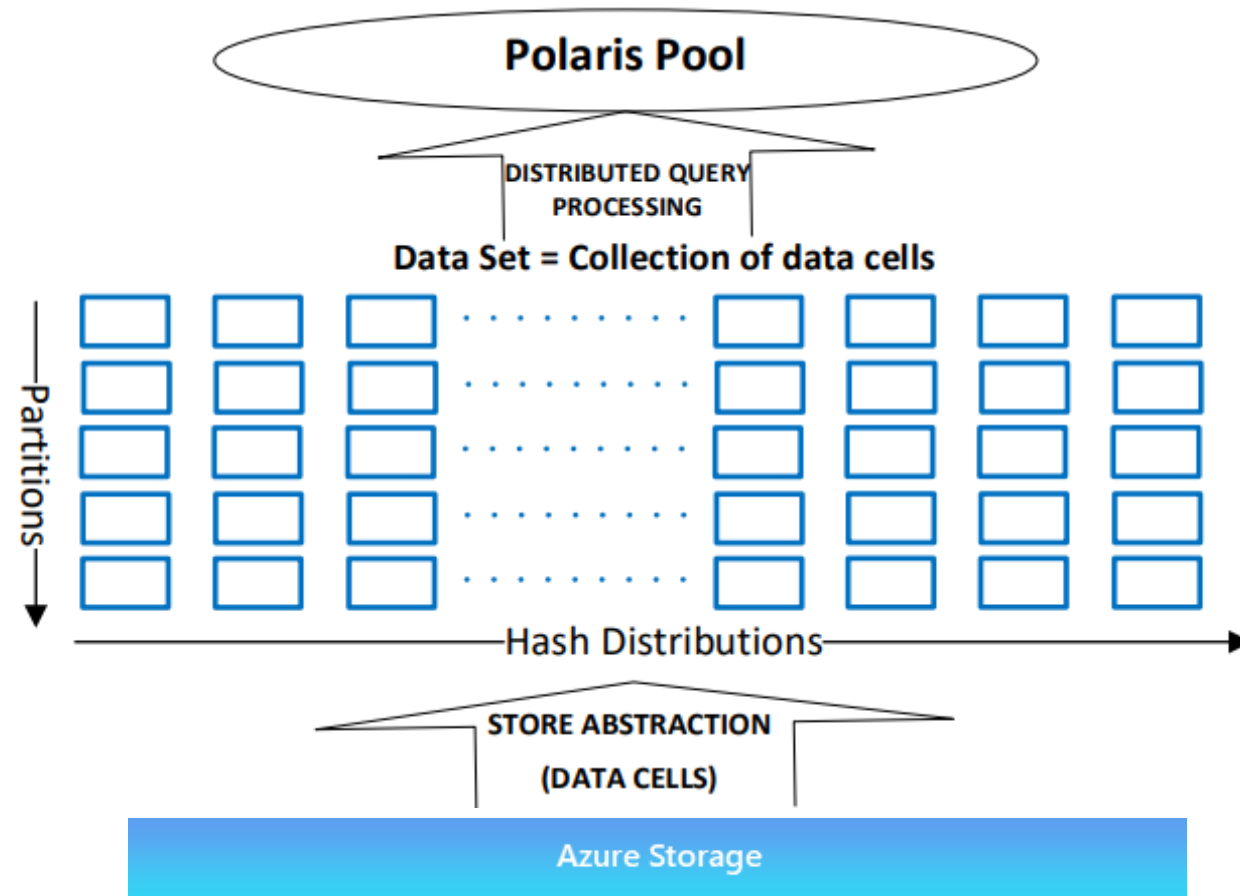
# Data Cells and Tasks



# Task are mapped to cell

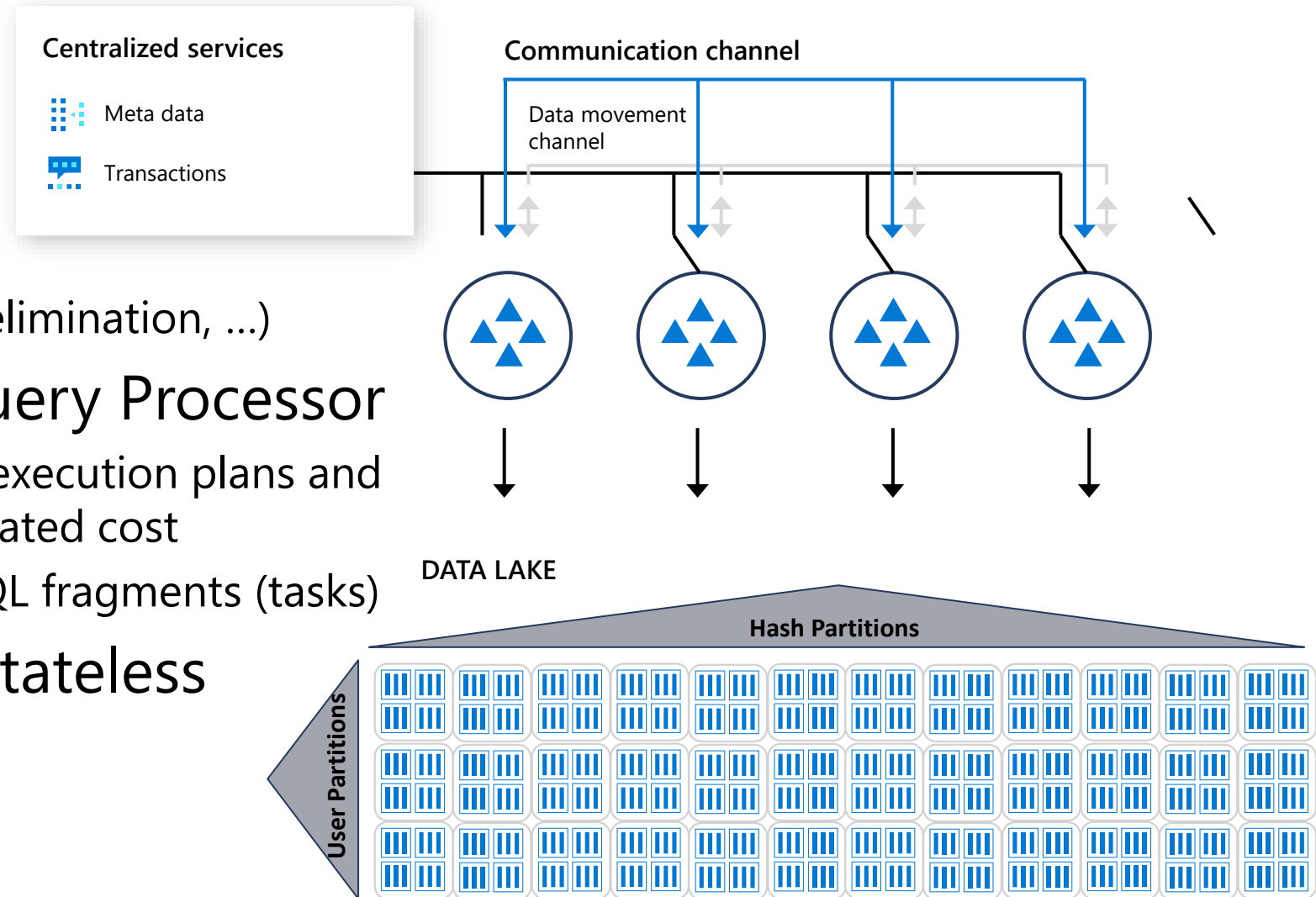


# Abstraction



# Distributed query execution flow

- SQL Frontend
  - Metadata
  - Security
  - Query simplification (filter pushdown, partition elimination, ...)
- DQP – Distributed Query Processor
  - Explores viable distributed execution plans and picks one with lowest estimated cost
  - Breaks user query into T-SQL fragments (tasks)
- SQL Backend – fully stateless
  - Executes tasks
  - Propagate results to parent





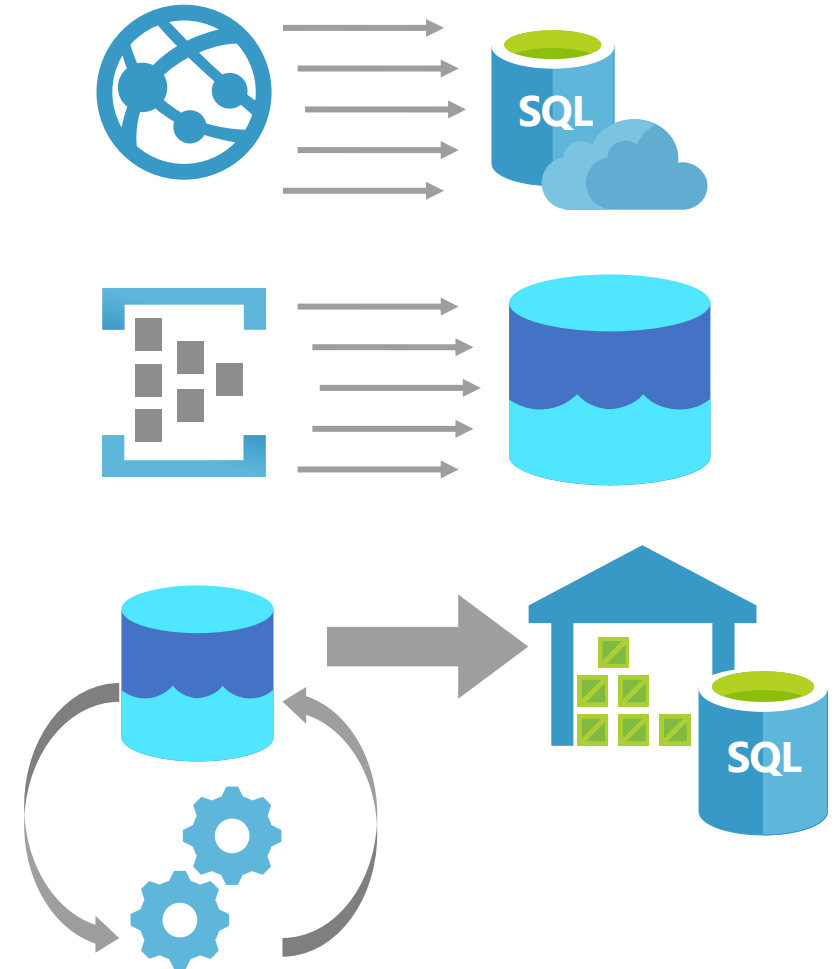
# What workloads are NOT suitable?

## Operational workloads (OLTP)

- High frequency reads and writes.
- Large numbers of singleton selects.
- High volumes of single row inserts.
- Critical SLA Reporting and Queries.
- Not recommended for sub-second performance.

## Data Preparations

- Row by row processing needs.
- Incompatible formats (XML).



# What Workloads are Suitable?

Store large volumes of data.

Consolidate disparate data into a single location.

Shape, model, transform and aggregate data.

Batch/Micro-batch loads.

Perform query analysis across large datasets.

Ad-hoc reporting across large data volumes (relaxed SLA's).

All using simple SQL constructs.

# Observations

- External tables (CETAS)
  - Statistics Matter!
  - Serverless SQL pool relies on statistics to generate optimal query execution plans.
  - Parquet files - automatically created
  - CSV files - you should create statistics manually for columns
  - Where? Columns - particularly used in DISTINCT, JOIN, WHERE, ORDER BY and GROUP BY
  - `sys.sp_create_openrowset_statistics [ @stmt = ] N'statement_text'`
- [Create and update statistics using Azure Synapse SQL resources - Azure Synapse Analytics | Microsoft Docs](#)

# Observations

- External tables (CETAS)

- Data Types Matter!
- Recommend using the right data types for the CREATE EXTERNAL TABLE at all times
- Don't use varchar(8000) for all columns 😊. If the max length for a column is 25 use varchar(25).
- If the data is numeric/decimal use int, bigint, decimal, float, etc. Don't use varchar type.
  - [Synapse SQL Data Types - CREATE TABLE \(Azure Synapse Analytics\) - SQL Server | Microsoft Docs](#)
- Correct collations also matter in terms of performance. Avoid unexpected conversions.
  - Details: [Always use UTF-8 collations to read UTF-8 text in serverless SQL pool - Microsoft Tech Community](#)

# Observations

- [Guidance : Best practices for serverless SQL pool - Azure Synapse Analytics | Microsoft Docs](#)
- External tables (CETAS)
  - CETAS is a parallel operation that creates external table metadata and exports the SELECT query results to a set of files in your storage account.
  - As CETAS generates Parquet files, statistics will be automatically created when the first query targets this external table, resulting in improved performance for subsequent queries targeting table generated with CETAS.

# Observations

- Execution Plans – DAGS
- DMV's - Coming soon

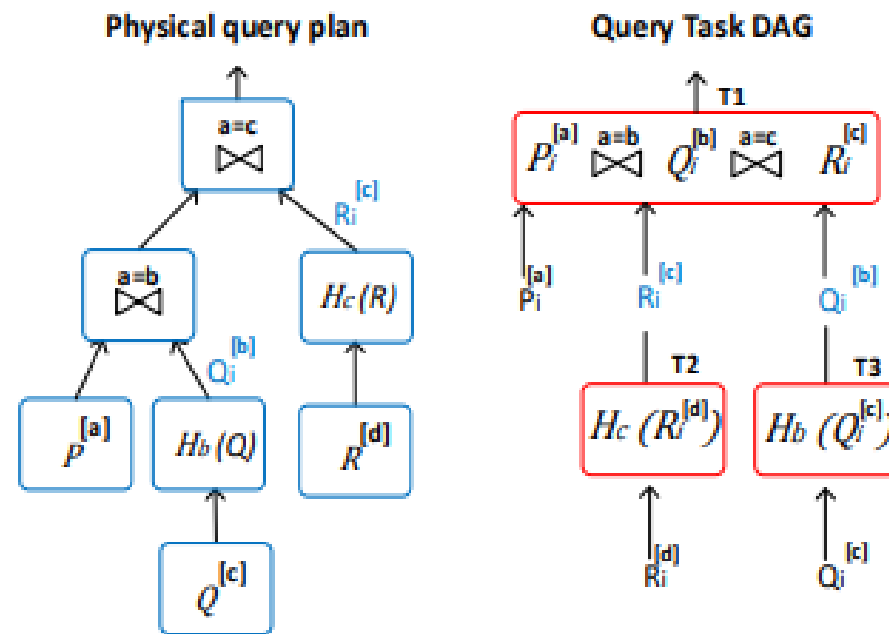


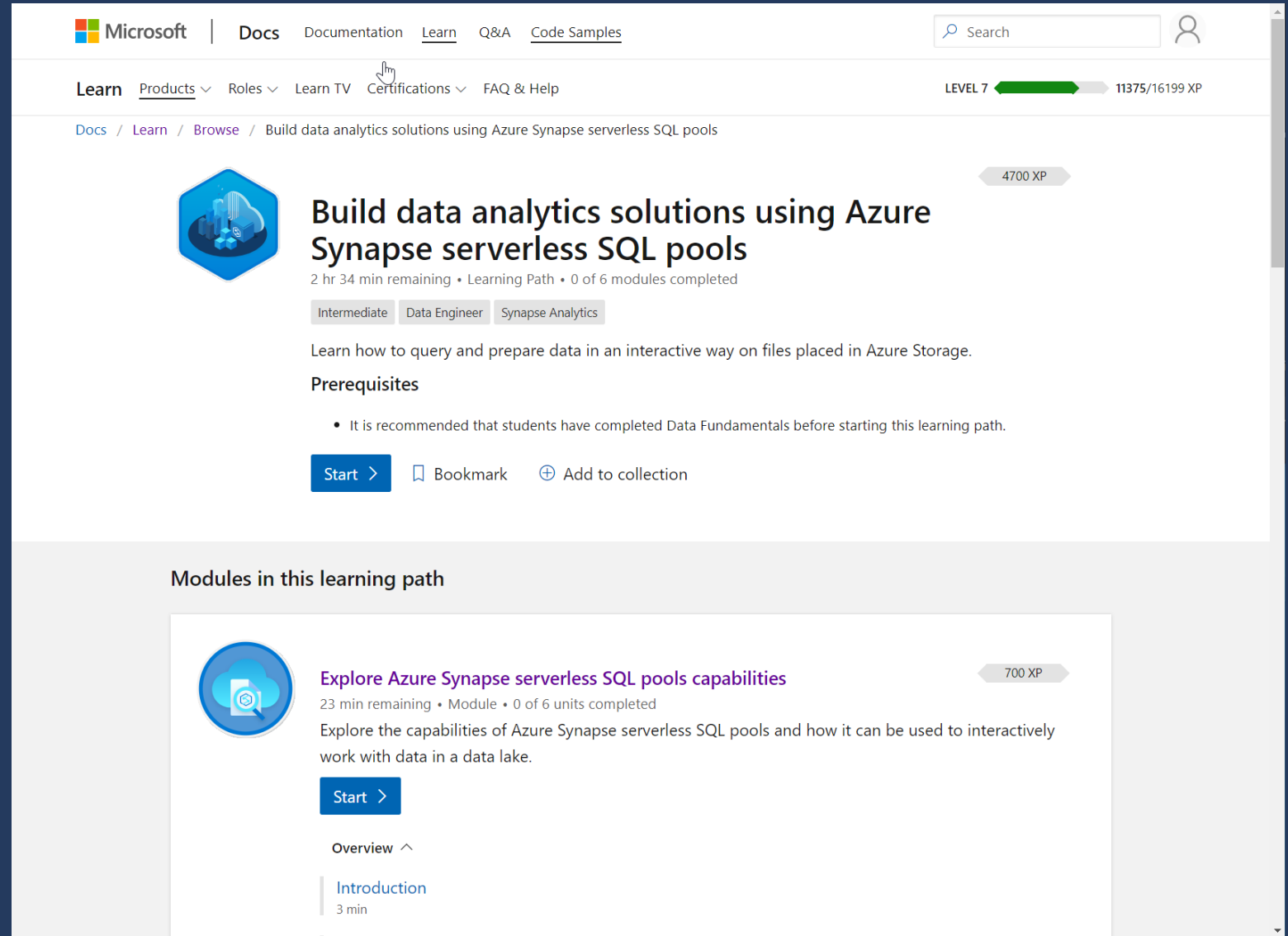
Figure 6. The Query Task DAG

# Demo: Run interactive queries using Azure Synapse serverless SQL pools

# /learn alert

Complete interactive learning exercises, watch videos, and practice and apply your new skills.

[aka.ms/mslearnasaserverless](https://aka.ms/mslearnasaserverless)



The screenshot shows the Microsoft Learn interface. At the top, there's a navigation bar with 'Microsoft', 'Docs', 'Documentation', 'Learn', 'Q&A', and 'Code Samples'. A search bar and a user profile icon are on the right. Below this, a secondary navigation bar includes 'Learn', 'Products', 'Roles', 'Learn TV', 'Certifications', and 'FAQ & Help'. A progress indicator shows 'LEVEL 7' with a green bar and '11375/16199 XP'.

The main content area displays a learning path titled 'Build data analytics solutions using Azure Synapse serverless SQL pools'. It features a blue hexagonal icon with a database and cloud. The title is in large, bold text. Below the title, it says '2 hr 34 min remaining • Learning Path • 0 of 6 modules completed'. There are three tags: 'Intermediate', 'Data Engineer', and 'Synapse Analytics'. A description states: 'Learn how to query and prepare data in an interactive way on files placed in Azure Storage.' Under 'Prerequisites', it lists: 'It is recommended that students have completed Data Fundamentals before starting this learning path.' At the bottom of this section are three buttons: 'Start >', 'Bookmark', and 'Add to collection'.

Below this, a section titled 'Modules in this learning path' shows a list of modules. The first module is 'Explore Azure Synapse serverless SQL pools capabilities', with a blue circular icon. It indicates '23 min remaining • Module • 0 of 6 units completed'. The description is: 'Explore the capabilities of Azure Synapse serverless SQL pools and how it can be used to interactively work with data in a data lake.' It has a 'Start >' button. Below the module title, there's an 'Overview' section with a dropdown arrow, and an 'Introduction' section with a '3 min' duration.





# Azure Synapse serverless SQL Pools

Synapse Serverless has two Analytics Runtimes

- SQL engine – Per TB
- Spark engine – Per Hour