# Spike: Evaluate Spark Engine in Azure Fabric

## Description 🔗

As a data Engineer, I want to run transformation jobs in Spark Engine in Fabric to evaluate the following:

- Integrate with MVE Data Generator
- Flexibility of Job Scheduling
- Integration with ADLS/OneLake
- Performance of Spark Engine
- CI/CD of Spark Notebooks
- Data Observability
- Integration of Great Expectations for Data Validation Checks
- KeyVault Integration in Spark Engine
- Cost Estimation for Workloads

## Summary 🔗

| Feature | Support Feature |
|---|---|
| Flexibility of Job Scheduling | Partial |
| Integration with ADLS/OneLake | Yes |
| Performance of Spark Engine | Yes |
| CI/CD of Spark Notebooks | Limited |
| Data Observability | Limited |
| Integration of Great Expectations for Data Validation Checks | Yes |
| KeyVault Integration in Spark Engine | Yes |
| Cost Estimation for Workloads | Needs Review |

## Findings 🔗

### Flexibility of Job Scheduling 🔗

In short Apache Spark job scheduling can be done at the Cluster level or at Spark Application level.

Cluster Level: Here a job refers to an multiple Spark Application and this refers to scheduling jobs on the same cluster based on the resources available in the cluster. Here each job requests for resources to the cluster manager for processing the job. Once the request is approved, these resources are locked and won't be available for subsequent job until they are released.

Application Level: In Spark Action triggers a Job, In this case this is what we are talking about. There might be case where there are multiple jobs request that comes in, job scheduling refers to scheduling these jobs accordingly.

Spark job definition comparison

Concurrency Throttling and Queueing:

- In Microsoft Fabric Spark, job submission is based on the purchased Fabric capacity SKUs.
- The queueing mechanism operates as a simple FIFO-based queue, allowing jobs to be submitted once capacity becomes available.
- When the capacity is fully utilized due to concurrent running jobs, new submissions are throttled with the message: "Unable to submit this request because all the available capacity is currently being used. Cancel a currently running job, increase your available capacity, or try again later."
- Notebook jobs and Spark job definitions are automatically retried when capacity becomes available. The queue expiration is set to 24 hours from the job submission time1.

Bursting:

- Fabric capacities support bursting, which allows you to consume extra compute cores beyond what you've purchased.
- For Spark workloads, bursting enables users to submit jobs with a total of 3X the Spark VCores purchased.
- Note that bursting increases the total number of Spark VCores for concurrency but doesn't increase the maximum cores per job.
- Users cannot submit a job that requires more cores than what their Fabric capacity offers.

Optimistic Job Admission:

- Recently introduced, Optimistic Job Admission for Fabric Spark provides more flexibility for concurrency usage.
- With this feature, you can run up to ~24 jobs concurrently with the same configuration when submitted concurrently.
- It prevents job starvation and significantly improves concurrency (in some cases, up to ~12X increase)2.

Custom Apache Spark Pools:

- You can create custom Spark pools in Fabric.
- Enable or disable autoscaling for your custom Spark pools.
- When autoscaling is enabled, the pool dynamically acquires new nodes up to the maximum node limit specified by the user and retires them after job execution.
- This dynamic resource adjustment ensures better performance based on job requirements3.

Considerations and limitations:

- Livy API and how to submit and manage Spark jobs: Livy API is in the roadmap but not exposed yet in Fabric. You must create notebooks and Spark job definitions with the Fabric UI. (9)
- Managed identity: Currently, Fabric doesn't support running notebooks and Spark job definitions using the workspace identity or managed identity for Azure KeyVault in notebooks. (9)

## Performance of Spark Engine 🔗

Spark engine in Fabric spins up significantly faster than the Synapse engine in 2 to 4 seconds against 2 minutes. (1)

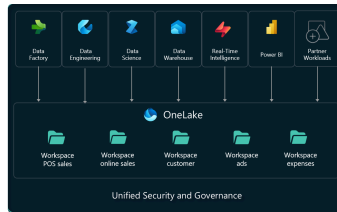Compare Fabric Data Engineering and Azure Synapse Spark

Node size family: Fabric Spark pools only support Memory Optimized node size family for now. If you're using a GPU-accelerated SKU Spark pool in Azure Synapse, they aren't available in Fabric.

## Integration with ADLS/OneLake 🔗

### OneLake - To rule them all 🔗

OneLake Overview OneLake is a single, unified, logical data lake for your whole organization. Like OneDrive, OneLake comes automatically with every Microsoft Fabric tenant and is designed to be the single place for all your analytics data. OneLake brings customers:

- One data lake for the entire organization
- One copy of data for use with multiple analytical engines

OneLake is built on top of Azure Data Lake Storage (ADLS) Gen2 and can support any type of file, structured or unstructured. All Fabric data items like data warehouses and lakehouses store their data automatically in OneLake in Delta Parquet format. If a data engineer loads data into a lakehouse using Spark, and then a SQL developer uses T-SQL to load data in a fully transactional data warehouse, both are contributing to the same data lake. OneLake stores all tabular data in Delta Parquet format.

OneLake supports the same ADLS Gen2 APIs and SDKs to be compatible with existing ADLS Gen2 applications, including Azure Databricks. You can address data in OneLake as if it's one big ADLS storage account for the entire organization. Every workspace appears as a container within that storage account, and different data items appear as folders within those containers.

## OneLake file explorer for Windows 🔗

OneLake is the OneDrive for data. Just like OneDrive, you can easily explore OneLake data from Windows using the OneLake file explorer for Windows. You can navigate all your workspaces and data items, easily uploading, downloading, or modifying files just like you do in Office. The OneLake file explorer simplifies working with data lakes, allowing even nontechnical business users to use them.

For more information, see OneLake file explorer.

## Shortcuts connect data across domains without data movement 🔗

Shortcuts allow your organization to easily share data between users and applications without having to move and duplicate information unnecessarily. When teams work independently in separate workspaces, shortcuts enable you to combine data across different business groups and domains into a virtual data product to fit a user's specific needs.

A shortcut is a reference to data stored in other file locations. These file locations can be within the same workspace or across different workspaces, within OneLake or external to OneLake in ADLS, S3, or Dataverse — with more target locations coming soon. No matter the location, shortcuts make files and folders look like you have them stored locally.

## CI/CD of Spark Notebooks 🔗

Git Integration is in Preview (3) 🟦 Lifecycle management tutorial - Microsoft Fabric

Only works with Azure DevOps and Azure Repos (3) (4)

In order to edit the workspace without interfering with other team members' changes, each team member creates their own isolated workspace to work in until they're ready to share their changes with the team.

Azure DevOps integration (notebook) (3) Deployment pipelines (notebook) (3)

## Supported items 🔗

The following items are currently supported: (7)

Data pipelines Dataflows Gen1 Datamarts Lakehouse Notebooks Paginated reports Reports (except reports connected to semantic models hosted in Azure Analysis Services, SQL Server Analysis Services or reports exported by Power BI Desktop that depend on semantic models hosted in MyWorkspace) Semantic models (except push datasets, live connections, model v1, and semantic models created from the Data warehouse/lakehouse.) If the workspace or Git directory has unsupported items, it can still be connected, but the unsupported items are ignored. They aren't saved or synced, but they're not deleted either. They appear in the source control pane but you can't commit or update them. (3)

## Considerations and limitations 🔗

- Currently, only Git in Azure Repos with the same tenant as the Fabric tenant is supported.
- If the workspace and Git repo are in two different geographical regions, the tenant admin must enable cross-geo exports.

- Azure DevOps on-prem isn't supported.
- Sovereign clouds aren't supported. (3)
- Direct Query and composite models on Power BI Datasets and Analysis Services aren't supported at this time. (5)
- DirectLake semantic models aren't supported at this time. (5)

## Data Observability (6) 🔗

Spark Advisor Built-in monitoring pools and jobs (through Monitoring hub) Spark history server No support for:

- Prometheus/Grafana
- Log Analytics
- Storage Account
- Event Hubs

## Integration of Great Expectations for Data Validation Checks 🔗

Tutorial: Validate data using SemPy and Great Expectations (GX) Data Validation with Great Expectations on Microsoft Fabric

## KeyVault Integration in Spark Engine 🔗

Azure Key Vault integration is support via Microsoft Spark Utilities (MSSparkUtils)

The Credential Utilities specifically reference getting access tokens and managing secrets in an Azure Key Vault.

## Cost Estimation for Workloads 🔗

Fabric SKU Types (8)

- Azure - Billed per second with no commitment.
  - Pay as you go with no time commitment.
  - You can scale your capacity up or down using the Azure portal.
  - You can pause and resume your capacity as needed. This feature is designed to save money when the capacity isn't in use.
  - Microsoft Cost Management.
  - Azure Monitor Metrics.
- Microsoft 365 - Billed monthly or yearly, with a monthly commitment (8)
  - Microsoft 365 SKUs, also known as P SKUs, are Power BI SKUs that also support Fabric when it's enabled on top of your Power BI subscription. Power BI EM SKUs don't support Microsoft Fabric.

Understand your Azure bill on a Fabric capacity

Microsoft Fabric pricing

In terms of Power BI licenses, if you go for a F64 or a higher capacity, read users will not need Power BI Pro license only users who create content. Users will be able to create data warehouses, use notebooks, and manage their capacity without a Power BI Pro license just as it was possible to do with a Power BI Premium per Capacity license. For those who opt for a lower capacity version than F64, one Power BI Pro license per user will be required to consume the content.

Price for 1 CU pay-as-you-go is $0.18/hour (USD) for EAST US 2. F2 offers 2 CUs, F4 offers 4 CUs, F8 offers 8 CUs, F16. .. and F2048 offers 2048 CUs. So, for a scenario with an F64 SKU that is open for 730 hours (average number of hours per month), will cost $8,409.60/month ($0.18 per hour * 730h * 64 CUs). It is important to note, as mentioned previously, that cost optimization strategies can be put in place to optimize and reduce performance according to busy periods.

In addition to Compute's costs, there are OneLake storage costs which will be detailed around 0.023$/GB/month (again depending on the EAST US 2 region). To these charges, inter-region data transfer network charges may apply depending on the source/destination of each storage access.

A concrete example of cost calculation A company X wants to implement Microsoft Fabric, it has 500 active consumers and 100 users authoring in Power BI and wants to use all the functionalities of Fabric. She estimates that she will have around 1000GB of data in her OneLake and estimates that half the time of the day (outside working hours), half the capacity will be needed to support her activities.

Compute Power (CUs) • F64 for 365h per month @$0.18/CU/hour: $4,204.80 • F32 for 365h per month @$0.18/CU/hour: $2,102.40

Data Storage • 1000GB @ $0.023/GB/month: $23

Power BI Pro Licenses • 100 users @$10/user/month: $4000

For a monthly total of $11,330.20/month.

Note: Prices are displayed in USD for the EAST US 2 region.

[Microsoft Fabric pricing](#)

## Resources: 🔗

(1) 🔴 [Microsoft Fabric vs Synapse: A Comparative Study](#) .

(2) ⊞ [Comparison between Fabric and Azure Synapse Spark. - Microsoft Fabric](#)

(3) ⊞ [Overview of Fabric Git integration - Microsoft Fabric](#)

(4) ⊞ [Get started with Git integration - Microsoft Fabric](#)

(5) ⊞ [Git integration process - Microsoft Fabric](#)

(6) ⊞ [Apache Spark monitoring overview - Microsoft Fabric](#)

(7) ⊞ [The Microsoft Fabric deployment pipelines process - Microsoft Fabric](#)

(8) ⊞ [Buy a Microsoft Fabric subscription - Microsoft Fabric](#)

(9) ⊞ [Comparison between Fabric and Azure Synapse Spark. - Microsoft Fabric](#)