

# Azure Support for 4 Projects.

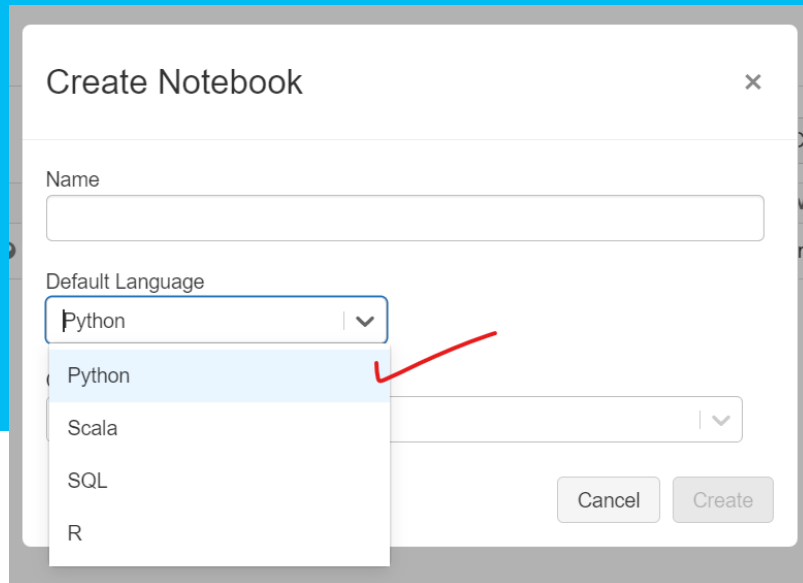
- Python for Apache Spark
- Apache Kafka
- Apache Beam
- Apache Cassandra



# PySpark in Azure

1

1. Use Azure Databricks.
2. Import Python Notebook into Databricks Workspace or code from scratch.
3. PySpark (Native Python API for Apache Spark Programming)
4. Supports:



Create Notebook

Name

Default Language

Python

Python

Scala

SQL

R

Cancel Create

2

1. Use Spark Pools in Synapse.
2. Synapse Notebook web interface.
3. Exports in .ipynb
4. Supports:

%%pyspark	Python ✓	Execute a <b>Python</b> query against Spark Context.
%%spark	Scala	Execute a <b>Scala</b> query against Spark Context.
%%sql	SparkSQL	Execute a <b>SparkSQL</b> query against Spark Context.
%%csharp	.NET for Spark C#	Execute a <b>.NET for Spark C#</b> query against Spark Context.

R currently NA (2021)

Python API for  
Apache Spark

# Apache Kafka in Azure

1

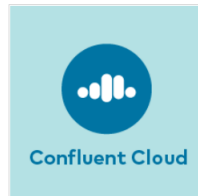
## Use Confluent on Azure

- Original creators of Apache Kafka.
- Fully Managed Apache Kafka on Azure.
- Created by Confluent; as a Solution on Azure. Comes with managed connectors with Blob, ADLS Gen2, SQL Server etc.

Try on Azure Marketplace > Click [here](#).  
Explore Pricing [here](#): Basic, Standard & Dedicated.

### **Better to use if you're using OR wanting to use features like:**

- Kafka [client-side end-to-end block compression](#).
- [Log compaction](#) (evicting all except last record of each key from a partition)
- Kafka [Streams client library](#) for Java/Scala apps integration
- Confluent's "[ksqlDB](#)" product in your solution architecture



2

## Use [Azure Event Hub for Kafka](#) (using Kafka protocol)

- Event Hubs end-point for Apache Kafka compatible producer & consumer API.
- Supports API clients at v1.0+.

Key difference between: Apache Kafka & Event Hubs:

1. Kafka you either manage on your OWN, or via Confluent. Event Hubs is a managed service from Microsoft.
2. Scale in Event Hub is managed by TU.

### **Better to use if you're using OR wanting to use features like:**

- [Competing-consumer](#) queue pattern (including Idempotency Patterns)
- Pub/Sub at a level that allows subs access to incoming messages based on server-evaluated rules other than plain offsets.
- Tracking of lifecycle of a job initiated by a message
- Sidelining faulty messages into a "dead-letter" queue

Excellent document is: [Asynchronous messaging options](#) in Azure.



Open-source Distributed  
Event-Streaming

# Apache Beam in Azure

A pipeline reads from an external source (file/db), writes output to a sink.

- Built-in I/O Transforms supported by Beam is > [Built-in I/O Transforms \(apache.org\)](#)

**Databases:**

## Java SDK

**Database**

These I/O connectors are used to connect to database systems.

Name	Description	Javadoc
CassandraIO	An IO to read from <a href="#">Apache Cassandra</a> .	<a href="#">org.apache.beam.sdk.io.cassandra.CassandraIO</a>
HadoopFormatIO (guide)	Allows for reading data from any source or writing data to any sink which implements <a href="#">Hadoop</a> InputFormat or OutputFormat.	<a href="#">org.apache.beam.sdk.io.hadoop.format.HadoopFormatIO</a>
HBaseIO	A bounded source and sink for <a href="#">HBase</a> .	<a href="#">org.apache.beam.sdk.io.hbase.HBaseIO</a>
HCatalogIO (guide)	HCatalog source supports reading of HCatRecord from a <a href="#">HCatalog</a> -managed source, for example <a href="#">Hive</a> .	<a href="#">org.apache.beam.sdk.io.hcatalog.HCatalogIO</a>
KuduIO	A bounded source and sink for <a href="#">Kudu</a> .	<a href="#">org.apache.beam.sdk.io.kudu</a>
SolrIO	Transforms for reading and writing data from/to <a href="#">Solr</a> .	<a href="#">org.apache.beam.sdk.io.solr.SolrIO</a>
ElasticsearchIO	Transforms for reading and writing data from/to <a href="#">Elasticsearch</a> .	<a href="#">org.apache.beam.sdk.io.elasticsearch.ElasticsearchIO</a>
BigQueryIO (guide)	Read from and write to <a href="#">Google Cloud BigQuery</a> .	<a href="#">org.apache.beam.sdk.io.gcp.bigquery.BigQueryIO</a>
BigTableIO	Read from and write to <a href="#">Google Cloud Bigtable</a> .	<a href="#">org.apache.beam.sdk.io.gcp.bigtable.BigtableIO</a>
DatastoreIO	Read from and write to <a href="#">Google Cloud Datastore</a> .	<a href="#">org.apache.beam.sdk.io.gcp.datastore.DatastoreIO</a>
SnowflakeIO (guide)	Experimental Transforms for reading from and writing to <a href="#">Snowflake</a> .	<a href="#">org.apache.beam.sdk.io.snowflake.SnowflakeIO</a>
SpannerIO	Experimental Transforms for reading from and writing to <a href="#">Google Cloud Spanner</a> .	<a href="#">org.apache.beam.sdk.io.gcp.spanner.SpannerIO</a>
JdbcIO	IO to read and write data on <a href="#">JDBC</a> .	<a href="#">org.apache.beam.sdk.io.jdbc.JdbcIO</a>
MongoDbIO	IO to read and write data on <a href="#">MongoDB</a> .	<a href="#">org.apache.beam.sdk.io.mongodb.MongoDbIO</a>
MongoDbGridFSIO	IO to read and write data on <a href="#">MongoDB GridFS</a> .	<a href="#">org.apache.beam.sdk.io.mongodb.MongoDbGridFSIO</a>
RedisIO	An IO to manipulate a <a href="#">Redis</a> key/value database.	<a href="#">org.apache.beam.sdk.io.redis.RedisIO</a>
DynamoDBIO	Read from and write to <a href="#">Amazon DynamoDB</a> .	<a href="#">org.apache.beam.sdk.io.aws.dynamodb.DynamoDBIO</a> <a href="#">org.apache.beam.sdk.io.aws2.dynamodb.DynamoDBIO</a>
ClickHouseIO	Transform for writing to <a href="#">ClickHouse</a> .	<a href="#">org.apache.beam.sdk.io.clickhouse.ClickHouseIO</a>

Open-source unified model  
for Batch + Streaming

## Python SDK

**Database**

These I/O connectors are used to connect to database systems.

Name	Description	pydoc
BigQueryIO (guide)	Read from and write to <a href="#">Google Cloud BigQuery</a> .	<a href="#">apache_beam.io.gcp.bigquery</a>
BigTableIO	Read from and write to <a href="#">Google Cloud Bigtable</a> .	<a href="#">apache_beam.io.gcp.bigtableio module</a>
DatastoreIO	Read from and write to <a href="#">Google Cloud Datastore</a> .	<a href="#">apache_beam.io.gcp.datastore.v1new.datastoreio</a>
MongoDbIO	IO to read and write data on <a href="#">MongoDB</a> .	<a href="#">apache_beam.io.mongodbio</a>

Use Beam to connect with Azure Cosmos DB's API for MongoDB via MongoDbIO connector.

Example Project > <https://bit.ly/3zl9Jqe>

Connectivity is one thing; Performance is another. Need to test Performance in PoC.

# Apache Cassandra in Azure

## Open-source NoSQL Distributed Database

1

### Inside Azure VM.

Benefits from more mem. Recommendations come from Performance Experiments, access in [GitHub](#).

- [Standard\\_DS14\\_v2](#) OR, Standard\_DS13\_v2 VMs, OR [Standard\\_L16s\\_v2](#).
- Data & Commit Logs stored on a striped set of 2/4 1-TB [P30](#).
- Use 1-2 TB data/VM with enough free space for compaction.
- For highest IOPS using premium managed disks, create a stripe-set instead of larger single disks.
- [Azure Ultra Disk](#) can also be evaluated for Cassandra workloads that require smaller disk capacity.
- [Accelerated Networking](#) on NIC of Cassandra node and VMs running client apps.
- For achieving "low" random-access disk latency for Cassandra read workloads, use Azure managed disks with [ReadOnly](#) caching enabled.
- Set Linux block dev read-ahead setting = 8 KB.
- Set strip set (e.g., /dev/md0) to 8 KB read-ahead.
- Commit logs should be on premium managed disks.
- You maintain IaaS and DB.

Additional [Best Practices](#) for Cassandra in VM.

2

### Managed DBaaS.

Azure [Managed Instance](#) for Cassandra.

- Currently in Public Preview (Dt. 26-08-2021).
- Automated deployment, management (patching & node health).
- Support for Hybrid scenarios.
- Support creation of a Multi-Region Cluster.
- Automated scaling of Nodes in Cluster.
- Pricing is flexible, on-demand, with no licensing fees.
- Cassandra repairs are done automatically for you via [Cassandra-repair.io](#).
- SLA: Currently in Public Preview, hence, does not come with SLA. Once comes out in "GA", SLA will be available.

3

### DataStax Astra on Azure.

DataStax Astra built on Apache Cassandra.

- OSS, serverless, DBaaS.
- [Astra](#), now in the Microsoft Azure Marketplace!
- Service plan is Pay-Per-Use.
- DataStax [Enterprise](#) also available via Marketplace. Accelerates cloud-native and bare-metal performance with specialist workloads including graph, search, analytics etc. Starter bundle (3-node) starts from \$5K/mth.

4

### Cosmos DB's API for Cassandra DB

Use apps written for Apache Cassandra.

- Use existing Apache Drivers compliant with CQLv4.
- The Cassandra API enables you to interact with data stored in Azure Cosmos DB using the Cassandra Query Language (CQL), Cassandra-based tools (like cqlsh) and Cassandra client drivers that you're already familiar with.
- No Ops.
- Cosmos DB guarantees low latency reads and writes at 99<sup>th</sup> percentile.
- Backed by SLAs.
- Use existing Code & tools.
- Throughput across all regions.
- Ability to globally distribute data across all Azure regions. SLA of 99.99% HA within a single region, 99.999% read and write HA across multiple-regions.
- Event Sourcing: Cassandra API provides access to a persistent change log, the Change Feed, which can facilitate event sourcing directly from the database.



DataStax

