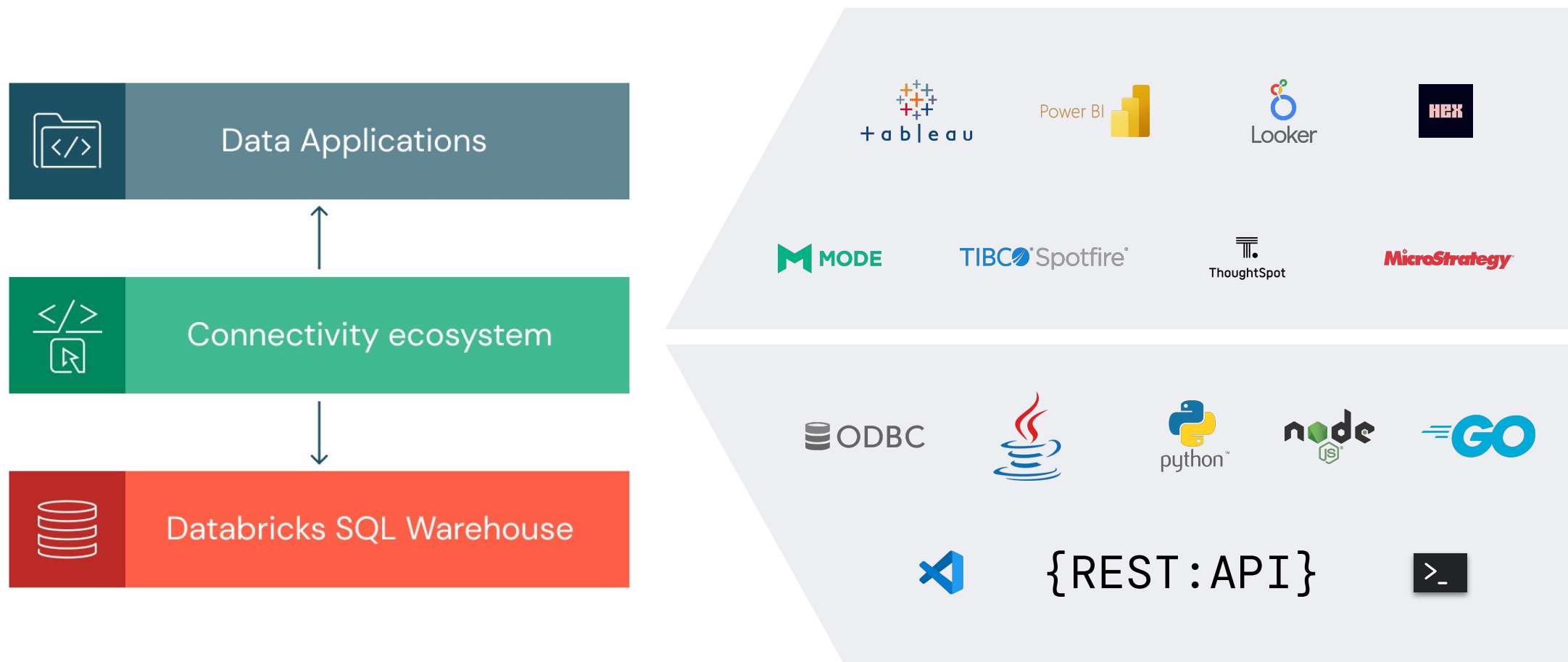# Databricks SQL

## Run SQL on the Lakehouse with your tools of choice



- – Connect to the lakehouse using well established BI tools, e.g., Power BI, Tableau, or Looker.

- – Easily ingest and transform data in-place using your favorite tools like Fivetran or dbt.

- – **Leverage existing applications to find insights or build data apps with tools and languages you already know.**
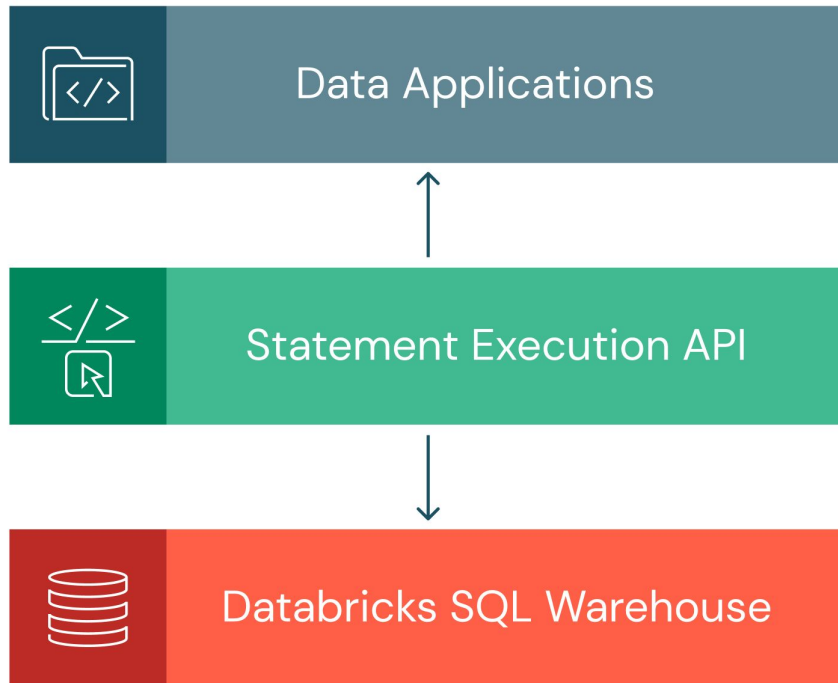
# Build apps powered by the Lakehouse

## Benefit from a rich ecosystem

# SQL Statement Execution API

## Access and manage data by executing SQL statements over HTTP

| | |
|---|---|
| Data Applications | |
| Statement Execution API | |
| Databricks SQL Warehouse | |

- Build custom data applications

- Integrate with a wide range of applications and computing devices

- Create a generic integration layer for enterprise services

- Create client libraries for your programming language of choice

# SQL Statement Execution API

## Access and manage data by executing SQL statements over REST



Business user

Data Application
</>

SQL Warehouse

Submit SQL statement for execution
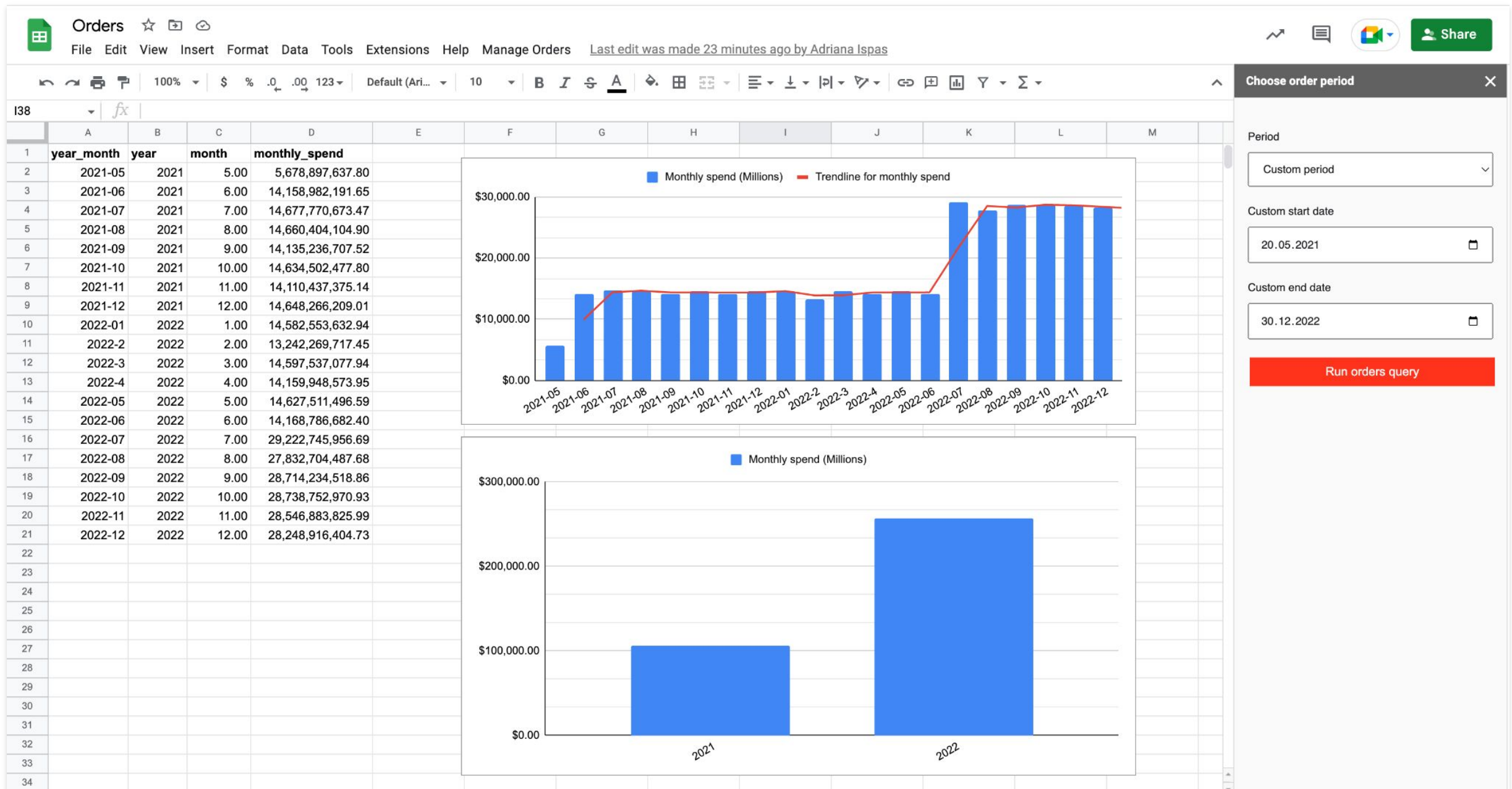`POST /sql/statements`

Check the execution status and retrieve results
`GET   /sql/statements/{statement_id}`

Cancel the execution of a SQL statement
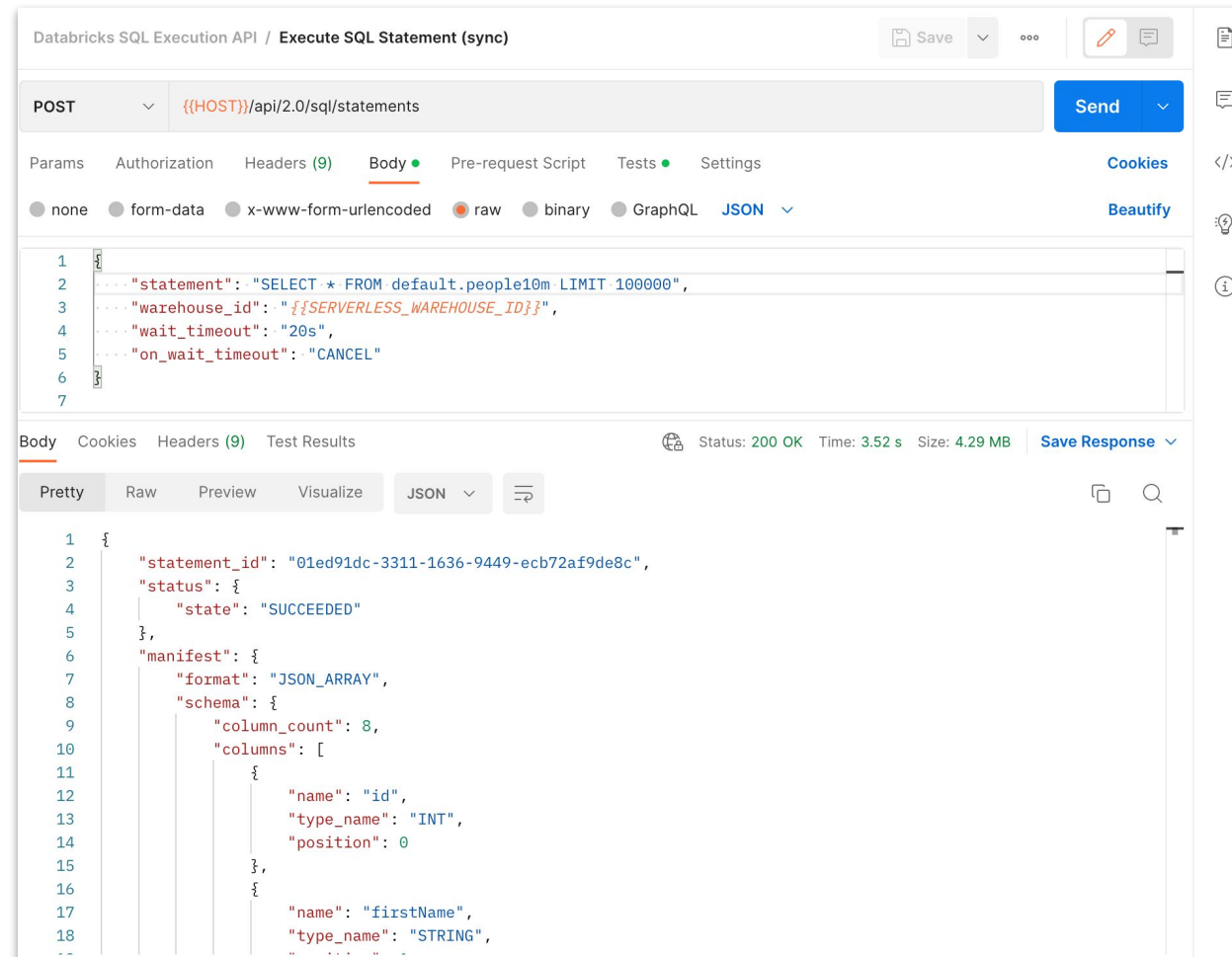`POST /sql/statements/{statement_id}/cancel`

# Integrate with Google Sheets

Blog post: https://www.databricks.com/blog/2023/03/07/databricks-sql-statement-execution-api-announcing-public-preview.html

# Demo

## Basics of the API using Postman

# Submit SQL statement for execution

## Three modes: synchronous, asynchronous and hybrid

**Default**: `wait_timeout` ∈ `[5, 50]s`. Then, continue async & fetch results in subsequent calls via ID

```
POST /sql/statements
    statement : "SELECT * FROM my_table"
    wait_timeout : "15s"
    on_wait_timeout: "CONTINUE"
```
→
```
# wait up to 15s ...
statement_id: "ID123"
status: { state: "RUNNING" }
```

**Asynchronous**: `wait_timeout = 0` → Execute async & fetch results in subsequent calls via ID

```
POST /sql/statements
    statement : "SELECT * FROM my_table"
    wait_timeout : "0s"
```
→
```
# no wait
statement_id: "ID123"
status: { state: "PENDING" }
```

**Synchronous**: `wait_timeout` ∈ `[5, 50]s` and return results in the same call. Otherwise, cancel.

```
POST /sql/statements
    statement : "SELECT * FROM my_table"
    wait_timeout : "15s"
    on_wait_timeout: "CANCEL"
```
→
```
# wait up to 15s, then cancel
statement_id: "ID123"
status: { state: "SUCCEEDED" }
manifest: { ... }
result: { ... }
```

# Fetching results

## Two modes: INLINE or EXTERNAL_LINKS

**Inline**: Results returned as payload, limited to 16 MiB, usually chunked, formats: JSON/CSV

```
POST /sql/statements
    ...
    disposition: "INLINE"
```

→

```
...
result:
  chunk_index: 0,
  row_offset: 0,
  row_count: 1000,
  data_array: [["1234","3.14159"],...]
```

**External links**: Results returned via resolved pre-signed URLs; 100 GB, formats: JSON/CSV/Arrow

```
POST /sql/statements
    ...
    disposition :
    "EXTERNAL_LINKS"
```

→

```
...
result:
  external_links:
  - chunk_index: 0,
    row_offset: 0,
    row_count: 257500,
    next_chunk_index": 1,
    next_chunk_internal_link: "/api/2.0/sql/statements/.../result/chunks/1?row_offset=..."
    external_link: "https://cloud.store/path/chunk00_abc?token=YYZ"
    expiration: "2022-09-22T19:21:03Z"
```

# Retrieve results w/ EXTERNAL_LINKS (1)

Check execution status using the handle & retrieve the 1st result chunk if ready

| `GET /sql/statements/ID123` | → | total_chunk_count: 3<br>chunks:<br>  -   chunk_index: 0<br>      row_offset: 0<br>      row_count: 257500<br>      ...<br>result:<br>  external_links:<br>  - chunk_index: 0,<br>    row_offset: 0,<br>    row_count: 257500,<br>    next_chunk_index": 1,<br>    next_chunk_internal_link:<br>"/api/2.0/sql/statements/.../result/chunks/1?row_offset=..."<br>    **external_link: "https://cloud.store/path/chunk00_abc?token=YYZ"**<br>    expiration: "2022-09-22T19:21:03Z" |
| `GET https://cloud.store/path/chunk00_abc?token=YYZ` | → | `[["4444","2.0"],...]` |

Note: the disposition = EXTERNAL_LINKS is specified when submitting the execution request

# Retrieve results w/ EXTERNAL_LINKS (2)

Retrieve a specific chunk in the result set

| | | |
|---|---|---|
| `GET /sql/statements/ID123/result/chunks/1/?row_offset=...` | → | `external_links:`<br>`- chunk_index: 1,`<br>`  row_offset: 257500,`<br>`  row_count: 257500,`<br>`  next_chunk_index": 2,`<br>`  next_chunk_internal_link:`<br>`"/api/2.0/sql/statements/.../result/chunks/2?row_offset=..."`<br>`  external_link: "https://cloud.store/path/chunk00_abc?token=YYZ"`<br>`  expiration: "2022-09-22T19:21:03Z"` |
| `GET https://cloud.store/path/chunk01_abc?token=YYZ` | → | `[["4444","2.0"],...]` |

Note: the disposition = EXTERNAL_LINKS is specified when submitting the execution request

# Parameterized SQL statements

## Improved security and reusability

```
POST /api/2.0/sql/statements HTTP/1.1
Host: <base_HOST>
Authorization: Bearer <personal_access_token>
Content-Type: application/json
{
    "statement": "SELECT * FROM stores WHERE store_id = :store_id",
    "warehouse_id": "<warehouse_ID>",
    "parameters": [
        {
         "name":"store_id",
         "type":"INT",
         "Value": 1234
        }
    ]
}
```

# Build a data app

# Acme, Inc's Data API

## Manage stores and their sales

### Get all stores

```
GET /stores
```

→

```
state: "SUCCEEDED"
stores: [["123", "Acme, Inc", …], ["456", "Databricks", …], …]
```

### Get sales for a store

```
GET /stores/<store_id>/sales
    request_id: Optional[token]
    limit: Optional[int]
    format: "CSV"
```

→

```
request_id: "ID456"
state: "PENDING"
links: Optional[Array]
```

### Create new sale

```
POST /stores/<store_id>/sales
    date: "2023-06-29"
    quantity: 10
    price: 2.50
    item_id: 1234
```

→

```
sale_id: "ID789"
```

15

# List Stores

## Synchronous mode, inline small data

**Acme Inc's API Request**

```
GET /stores
```

→

**SQL Statement Execution API Request**

```
POST /sql/statements
    statement: "SELECT * FROM stores"
    wait_timeout: "50s"
    on_wait_timeout: "CANCEL"
```

**Acme Inc's API Response**

```
state: "SUCCEEDED"
stores: [
  ["123", "Acme, Inc", …],
  ["456", "Databricks", …]
]
```

←

**SQL Statement Execution API Response**

```
statement_id: "ID123"
status: { state: "SUCCEEDED" }
manifest: { ... }
result: {
  data_array: [
    ["123", "Acme, Inc", …],
    ["456", "Databricks", …]
  ]
}
```

# Download Sales for a Store

## Asynchronous mode, large data with external links

**Acme Inc's API Request**

```
GET /stores/123/sales
    format: "CSV"
```

→

**SQL Statement Execution API Request**

```
POST /sql/statements
    statement: "SELECT * FROM stores where store_id = :store_id"
    parameters: [
        { name: "store_id", value: "123", type: "INT" }
    ]
    disposition: "EXTERNAL_LINKS"
    wait_timeout: "0s"
    on_wait_timeout: "CONTINUE"
```

**Acme Inc's API Response**

```
request_id: "ID123"
state: "RUNNING"
```

←

**SQL Statement Execution API Response**

```
statement_id: "ID123"
status: { state: "RUNNING" }
```

17

# Download Sales for a Store

## Asynchronous mode, large data with external links

**Acme Inc's API Request**

```
GET /stores/123/sales
     request_id: "ID123"
```

→

**SQL Statement Execution API Request**

```
GET /sql/statements/ID123
```

**Acme Inc's API Response**

```
request_id: "ID123"
state: "SUCCEEDED"
links: [
  "https://cloud.store/path/chunk00_abc?token=YYX",
  "https://cloud.store/path/chunk01_def?token=YYW",
  "https://cloud.store/path/chunk02_ghi?token=YYZ"
]
```

←

**SQL Statement Execution API Response**

```
statement_id: "ID123"
status: { state: "SUCCEEDED" }
manifest: { total_chunk_count: 3 }
result: {
  external_links: [
    {
      external_link:
        "https://cloud.store/path/chunk00_abc?token=YYZ"
    }
  ]
}


GET /sql/statements/ID123/result/chunks/1
GET /sql/statements/ID123/result/chunks/2
```

# Create a new sale
## Hybrid, DML

**Acme Inc's API Request**

```
POST /stores/123/sales
    date: "2023-06-29"
    quantity: 10
    price: 2.50
    item_id: 1234
```

→

**SQL Statement Execution API Request**

```
POST /sql/statements
    statement: "INSERT INTO …"
    Parameters: [ … ]
    wait_timeout: "50s"
    on_wait_timeout: "CONTINUE"
```

**Acme Inc's API Response**

```
sale_id: "ID456"
state: "SUCCEEDED"
```

←

**SQL Statement Execution API Response**

```
statement_id: "ID123"
status: { state: "SUCCEEDED" }
```

# Create a new sales order
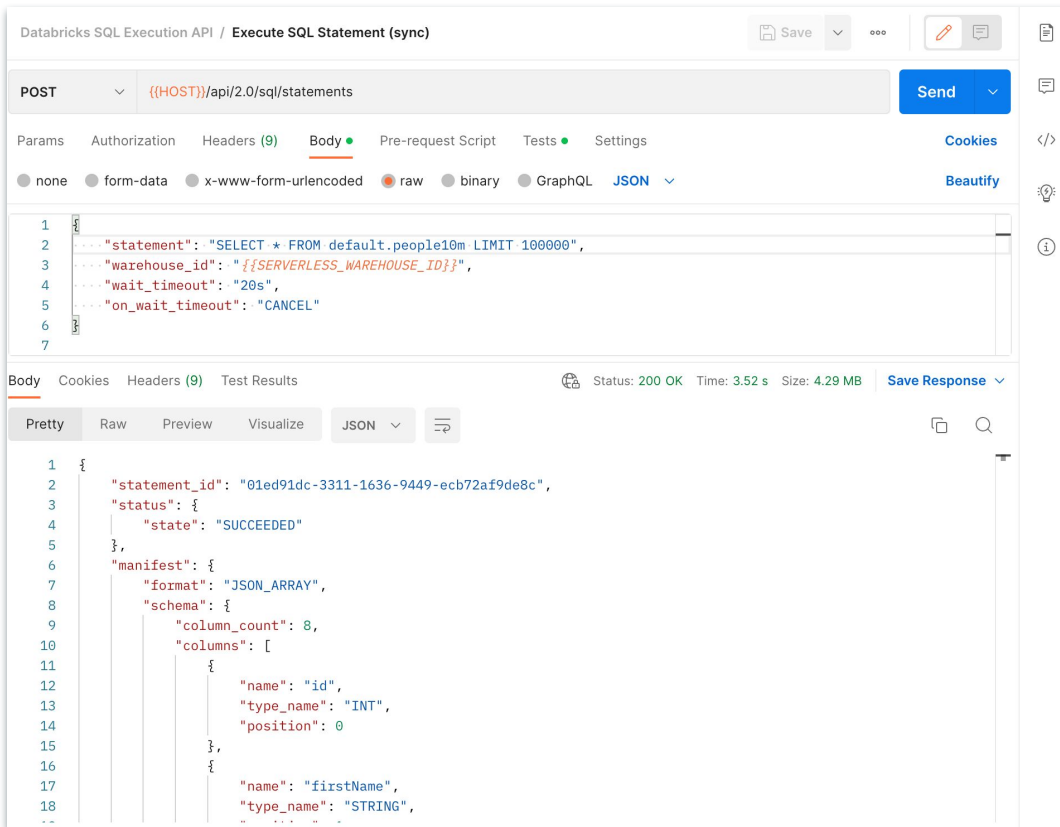
## Use **parameters** for improved security

```
POST /api/2.0/sql/statements HTTP/1.1
Host: <base_HOST>
Authorization: Bearer <personal_access_token>
Content-Type: application/json
{
    "statement":
      "INSERT INTO store_sales (ss_sold_date_sk, ss_ticket_number, ss_store_sk, ss_item_sk, ss_quantity, ss_sales_price)
       VALUES (:sold_date, :sale_id, :store_id, :item_id, :quantity, :sales_price)",
    "parameters": [
        { "name": "sold_date", "type": "DATE", "value": "2023-06-29" },
        { "name": "sale_id", "type": "BIGINT", "value": "1234" },
        { "name": "store_id", "type": "INT", "value": "567" },
        { "name": "item_id", "type": "INT", "value": "890" },
        { "name": "quantity", "type": "INT", "value": "10" },
        { "name": "sales_price", "type": "DECIMAL(7,2)", "value": "2.50" }
    ]
}
```

# Further code samples

## Check out our [Git repo](): Postman, cURL, notebooks, etc.

[github.com/databricks-demos/dbsql-rest-api]()

# SQL Statement Execution API
## Simplified data access using a programming language of your choice

- Removes the need to install drivers and manage Cloud infrastructure, or manage connections

- Allows querying and manipulating data, or defining data objects (DDL, DML, DQL, DCL)

- Allows different execution modes: synchronous, asynchronous, or hybrid

- Allows efficient access to large data sets with EXTERNAL_LINKS

- Leverages authentication options supported by Databricks REST APIs.

# Build apps powered by the Lakehouse

## Learn more and get started

→ Learn more about our connectors and tools

   https://docs.databricks.com/dev-tools/index-driver.html

→ Learn about apps in the marketplace

   https://www.databricks.com/blog/introducing-lakehouse-apps

23

# Related talks at Data+AI Summit

## Customer Talks

- **Akamai |** [Internet-Scale Analytics: Migrating a Mission Critical Product to the Cloud](#)

- **AT&T |** [Building and Managing Data Platform for 13+ PB Delta Lake and 1000s of Users: AT&T's Story](#)

- **S&P GLOBAL |** [Using Databricks to Power Insights and Visualizations on the S&P Global Marketplace](#)

- **Land O'Lakes |** [Self-Service Geospatial Analysis Leveraging Databricks, Apache Sedona, And R](#)

- **American Airlines |** [Making Travel More Accessible For Customers Bringing Mobility Devices](#)

- **Collins Aerospace |** [Jet Streaming Data and Predictive Analytics: How the Lakehouse and Apache Spark™ Enable Collins Aerospace to Keep Aircraft Flying](#)

- **Banco Bradesco |** [Data Democratization with Lakehouse: An Open Banking Application Case](#)

- **Michelin |** [Data Democratization at Michelin](#)

- **Zurich Insurance |** [dbt Labs | Modernizing the Data Stack: Lessons Learned From the Evolution at Zurich Insurance](#)

- **Rec Room |** [How Rec Room Processes Billions of Events Per Day with Databricks and RudderStack](#)

- **RaceTrac Inc. |** [Unlocking the Power of Real-Time Data to Maximize Data Insights](#)

- **dbt Labs |** [Modernizing the Data Stack: Lessons Learned From the Evolution at Zurich Insurance](#)

# Related talks at Data+AI Summit

## Product Deep Dive

- **Databricks |** [Databricks SQL: Why The Best Serverless Data Warehouse Is A Lakehouse](#)
- **Databricks |** [What's New In Databricks SQL -- With Live Demos](#)
- **Databricks |** [Databricks SQL Serverless Under the Hood: How We Use ML to Get the Best Price/Performance](#)
- **Databricks |** [Best Practices For Setting Up Databricks SQL At Enterprise Scale](#)
- **Databricks |** [Building Apps on the Lakehouse with Databricks SQL](#)
- **Databricks |** [Unlock The Next Evolution Of The Modern Data Stack With The Lakehouse Revolution -- With Live Demos](#)
- **Databricks |** [Unleashing Large Language Models with Databricks SQL's AI Functions](#)
- **Databricks |** [Under the Hood: Intelligent Workload Management](#)
- **Databricks |** [Going Beyond SQL: Python UDFs in Unity Catalog for all your Lakehouse](#)