

BUILD PIPELINE ON CLOUD

MENTOR: NGUYỄN HOÀNG VŨ THÂN
TẶNG THỊ THÚY VÂN





BUILD PIPELINE ON CLOUD

CONTEXT

01.

BASIC KNOWLEDGE

02.

CLOUD COMPUTING

03.

MICROSOFT AZURE

04.

DEMO



01. BASIC KNOWLEDGE

PYTHON

- Concept
- Compiler & Interpreter

SQL

- Concept
- OLTP & OLAP

PYTHON

DEFINITION

- High-level
- Open source
- Cross-platform

LIBRARIES

- Faker
- Pandas
- Datetime

```
card.py > ...
8 # Khởi tạo đối tượng Faker
9 fake = Faker()
10
11 # Kết nối đến cơ sở dữ liệu SQL Server
12 sql_server_host = r'INTERN-NTNGHI-L\SQLSEVER'
13 sql_server_database = 'projectsale'
14 sql_server_user = 'sa'
15 sql_server_password = '1234567'
16
17 # Tạo SQLAlchemy engine
18 engine = create_engine(f'mssql+pyodbc://{sql_server_user}:{sql_server_password}@{sql_server_host}/{sql_server_database}?driver=ODBC')
19
20 # Tạo DataFrame chứa dữ liệu
21 data = {'card_number': [], 'cc_name': [], 'pin': [], 'b_street': [], 'b_city': [], 'b_state': [], 'b_zip': []}
22
23 # Tạo một tập hợp để lưu trữ các card_number đã sử dụng
24 existing_card_numbers = set()
25
26 # Tạo 300 dòng dữ liệu
27 for _ in range(300):
28     card_number = ''.join(random.choices('0123456789', k=16))
29
30     # Kiểm tra xem card_number đã tồn tại chưa
31     while card_number in existing_card_numbers:
32         card_number = ''.join(random.choices('0123456789', k=16))
33
34     # Thêm card_number vào tập hợp đã sử dụng
35     existing_card_numbers.add(card_number)
36
37     cc_name = fake.name()
38     pin = ''.join(random.choices('0123456789', k=4))
39     b_street = fake.street_address()
40     b_city = fake.city()[:20]
41     b_state = fake.state_abbr()
42     b_zip = fake.zipcode()
```

BUILD PIPELINE ON CLOUD



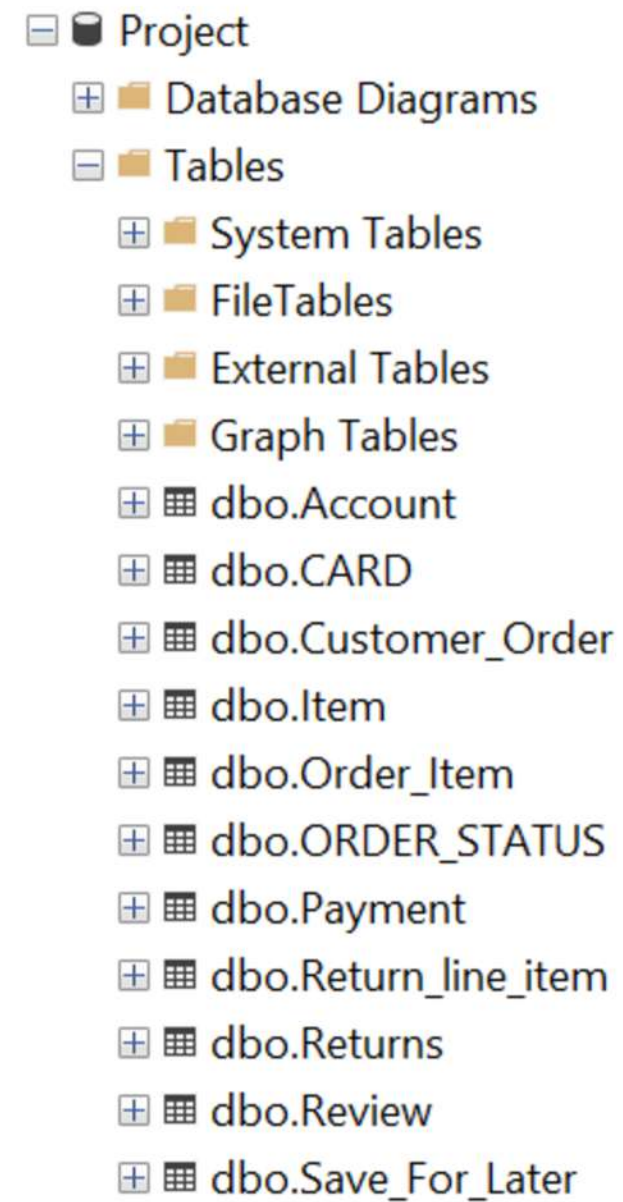
COMPILER & INTERPRETER

Criteria	Compiler	Interpreter
Input	Entire program	Single line of code
Output	Intermediate object code	No intermediate object code is generated
Mechanism	Compilation happens before execution	Compilation and execution happen simultaneously
Speed	Faster	Slower
Memory	Requires more memory due to intermediate object code generation	Requires less memory since it does not generate intermediate object code
Errors	Displays all errors after compilation, all at once	Displays errors one line at a time
Error Detection	Very difficult	Relatively easy
Programming Languages	C, C++, C#, Scala, typescript	PHP, Perl, Python, Ruby

SQL

DEFINITION

- Programming language
- Storing and processing
- Relational databases

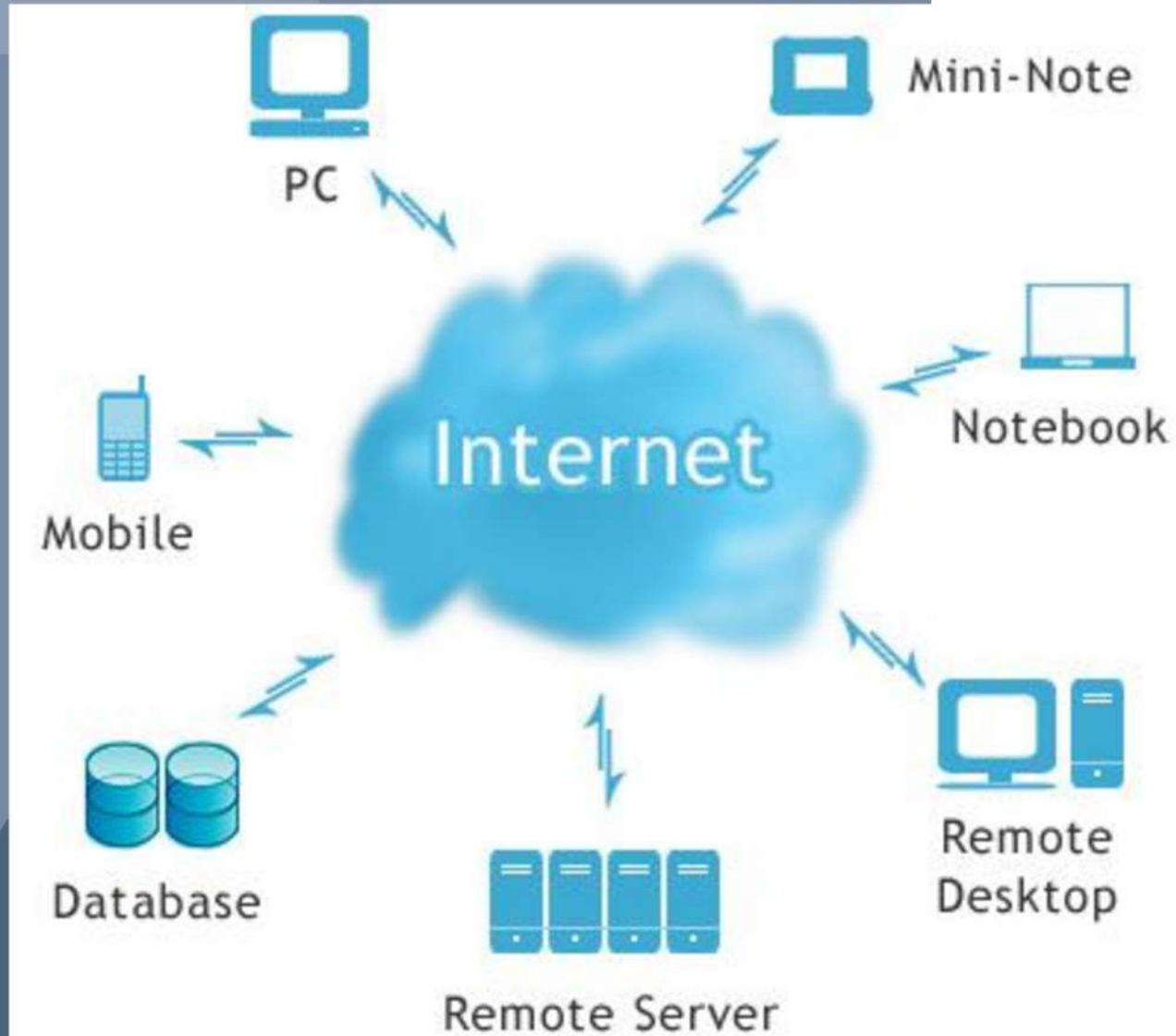


OLTP & OLAP

Feature	OLAP	OLTP
Purpose	OLAP helps you analyze large volumes of data to support decision-making	OLTP helps you manage and process real-time transactions
Data source	OLAP uses historical and aggregated data from multiple sources	Uses real-time and transactional data from a single source
Data structure	OLAP uses multidimensional (cube) or relational databases	Uses relational databases
Data model	Uses star schema, snowflake schema, or other analytical models	Uses normalized or denormalized models
Data volume	Requires large storage capacity. Think terabytes (TB) and petabytes (PB).	Has relatively small storage requirements. Think gigabytes (GB).
Response time	Has longer response times, typically in seconds or minutes	Has shorter response times, typically in milliseconds
Example applications	Analyze trends, predict customer behavior, and identify profitability	Process payments, manage customer data, and fulfill orders

02. CLOUD COMPUTING

- Providing services
- Through the Internet



• CLOUD SERVICE TYPES




- IaaS
- PaaS
- SaaS

• CLOUD MODELS

- Private Cloud
- Public Cloud
- Hybrid Cloud

RESPONSIBILITIES OF EACH CLOUD SERVICE TYPE

	Responsibility	SaaS	PaaS	IaaS	On-prem
Responsibility always retained by the customer	Information and data	Customer	Customer	Customer	Customer
	Devices (Mobile and PCs)	Customer	Customer	Customer	Customer
	Accounts and identities	Customer	Customer	Customer	Customer
Responsibility varies by type	Identity and directory infrastructure	Shared	Shared	Customer	Customer
	Applications	Microsoft	Shared	Customer	Customer
	Network controls	Microsoft	Shared	Customer	Customer
	Operating system	Microsoft	Microsoft	Customer	Customer
Responsibility transfers to cloud provider	Physical hosts	Microsoft	Microsoft	Microsoft	Customer
	Physical network	Microsoft	Microsoft	Microsoft	Customer
	Physical datacenter	Microsoft	Microsoft	Microsoft	Customer

 Microsoft  Customer  Shared

CLOUD MODELS

Public Cloud	Private Cloud	Hybrid Cloud
No capital expenditures to scale up	Complete control	Provides the most flexibility
Quickly provisioned and deprovisioned	Data is not collocated with other organizations' data	Determine where to run their applications
Pay only for what they use	Hardware must be purchased for startup and maintenance	Control security, compliance, or legal requirements
Insufficient control over resources and security	Responsible for hardware maintenance and updates	No mention

03. MICROSOFT AZURE



- Azure Blob Storage
- Azure Data Factory
- Azure Databricks

03. MICROSOFT AZURE

AZURE BLOB STORAGE

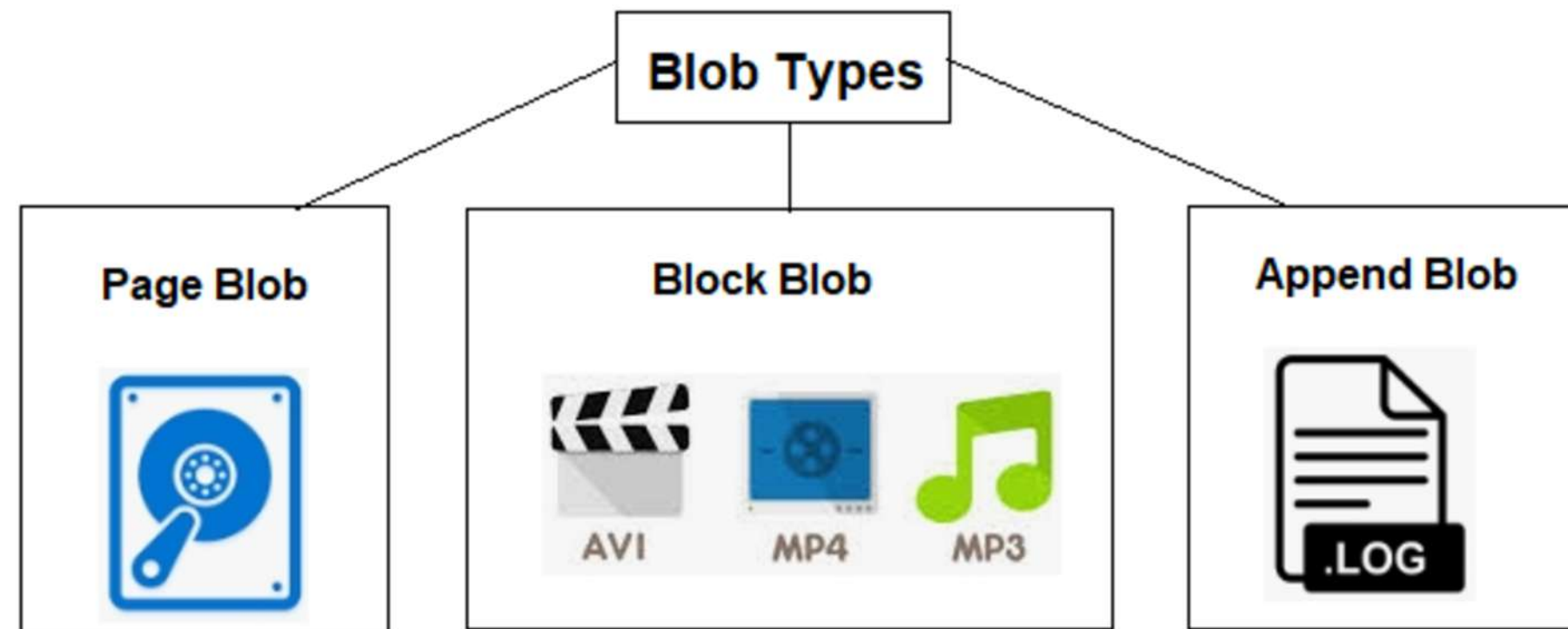
- Microsoft's object storage solution
- Store large amounts of unstructured data

Main components of Azure Blob Storage

Storage Account

Container

Blob



03. MICROSOFT AZURE

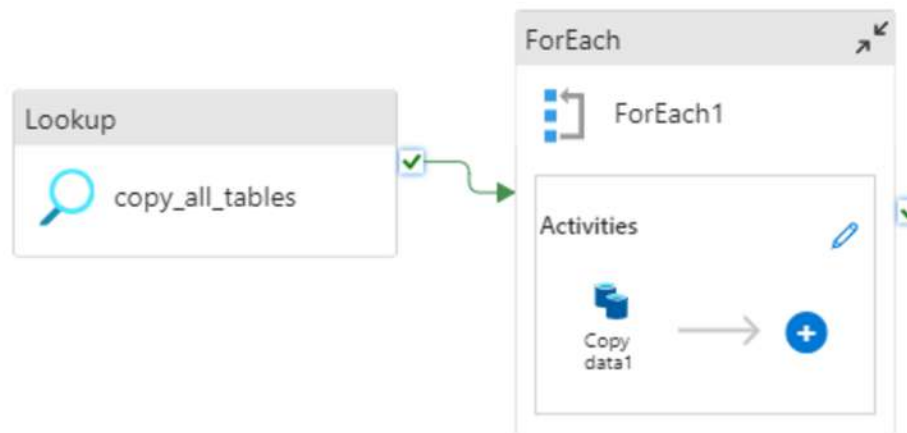
AZURE DATA FACTORY

Key features

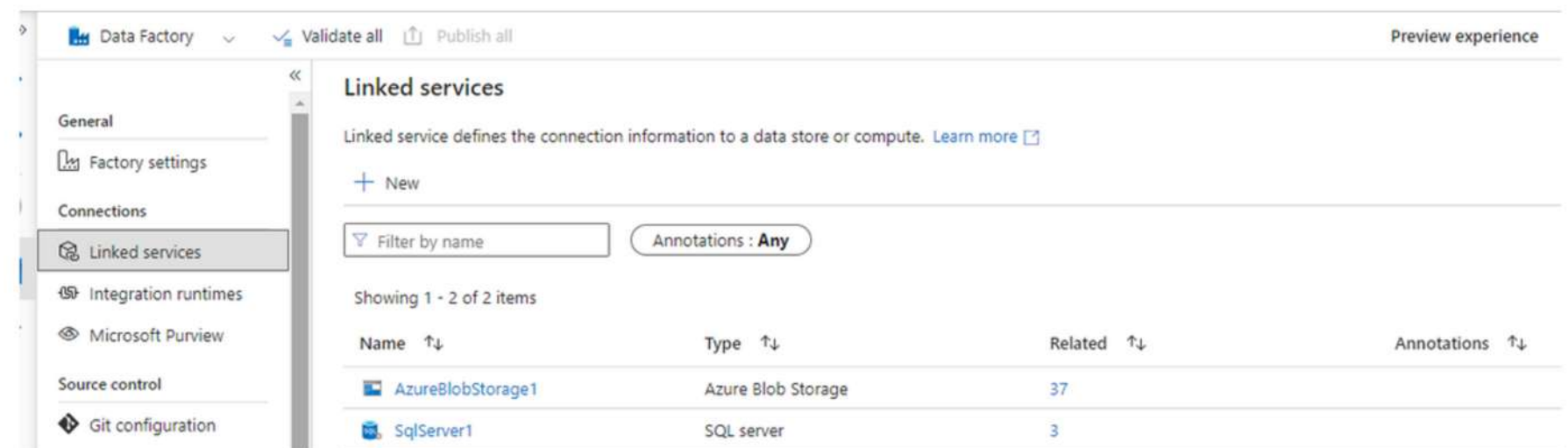
- Data Integration Service
- Data Transformation
- Data Movement
- Automation and Scheduling

Key components

- Pipeline
- Activity
- Dataset
- Trigger



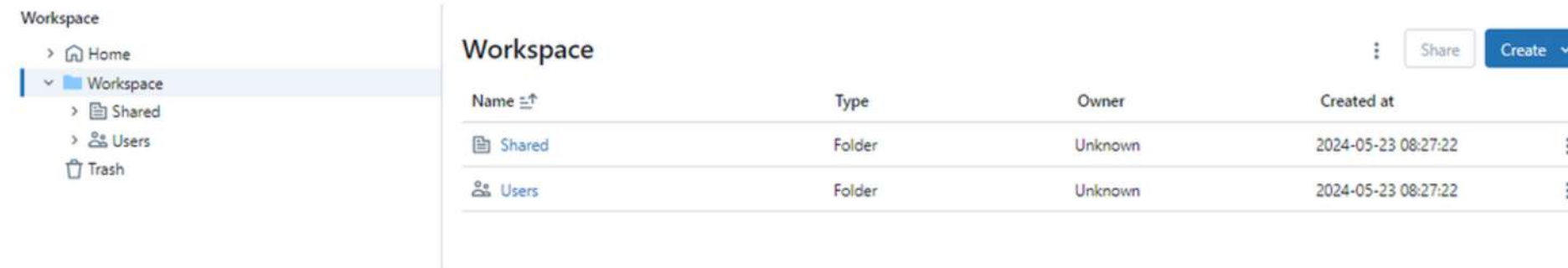
- Linked Service



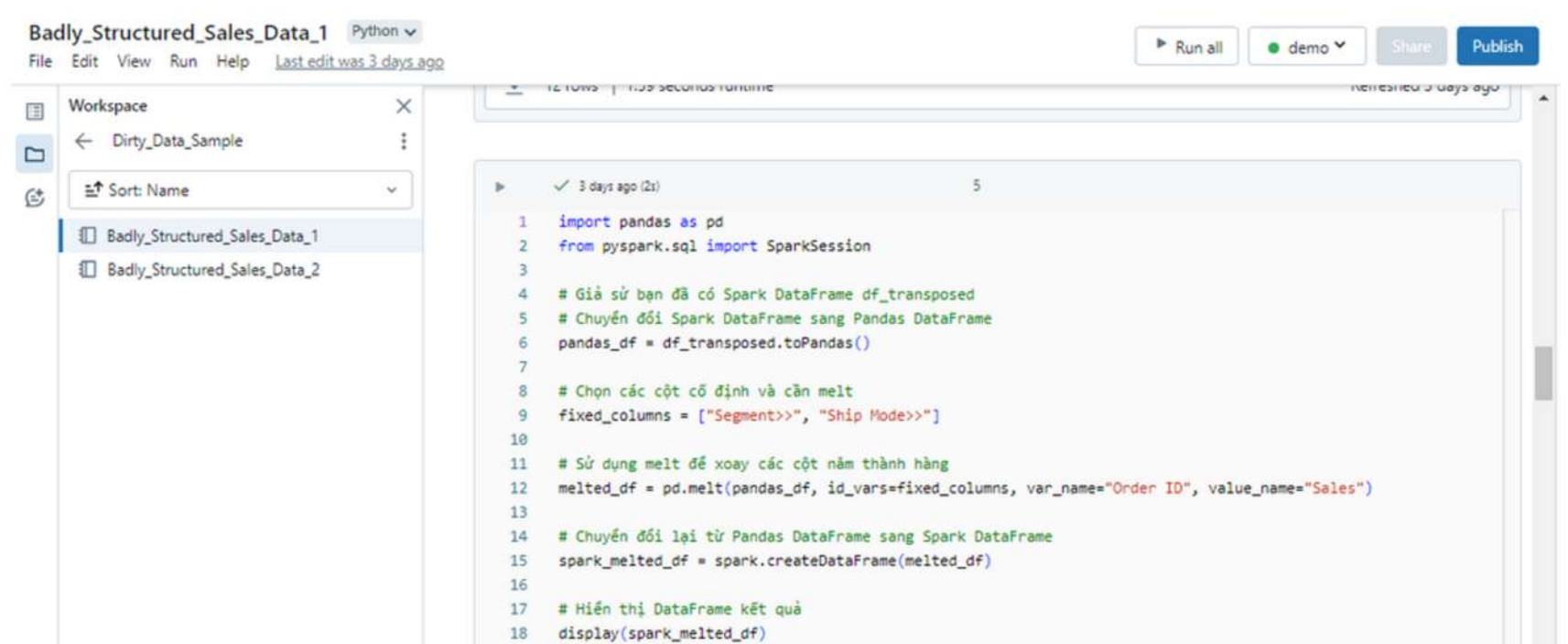
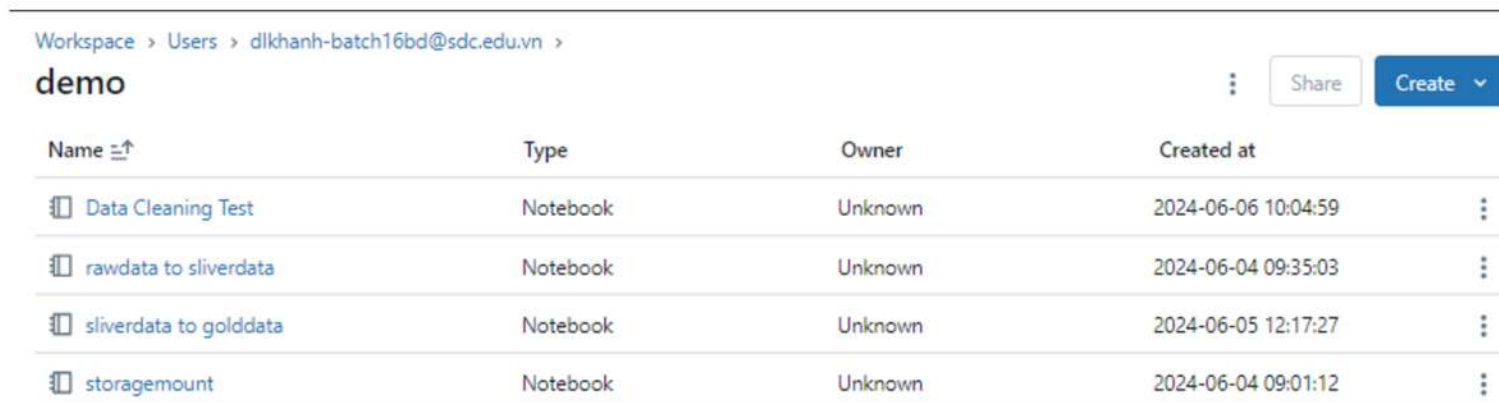
03. MICROSOFT AZURE

DATABRICKS

- Workspace



- Notebook



03. MICROSOFT AZURE

DATABRICKS

- Clusters

State	Name	Runtime	Active memory	Active cores	Active DBU / h	Source	Creator	Notebooks	
	demo	12.2	15 GB	2 cores	1	UI	dlkhanh-batch16bd@sd...	-	

- Delta Lake

Database TablesDBFS

Upload

/delta/heart_rate_data

Prefix search

heart_rate_data

Prefix search

_delta_log

part-00000-5b24c8bc-4757-41d4-...

part-00000-6eff063f-084c-4368-97...

part-00000-72244e75-1895-4583-...

part-00000-9ea588f8-55a0-42d6-b...

04. DEMO

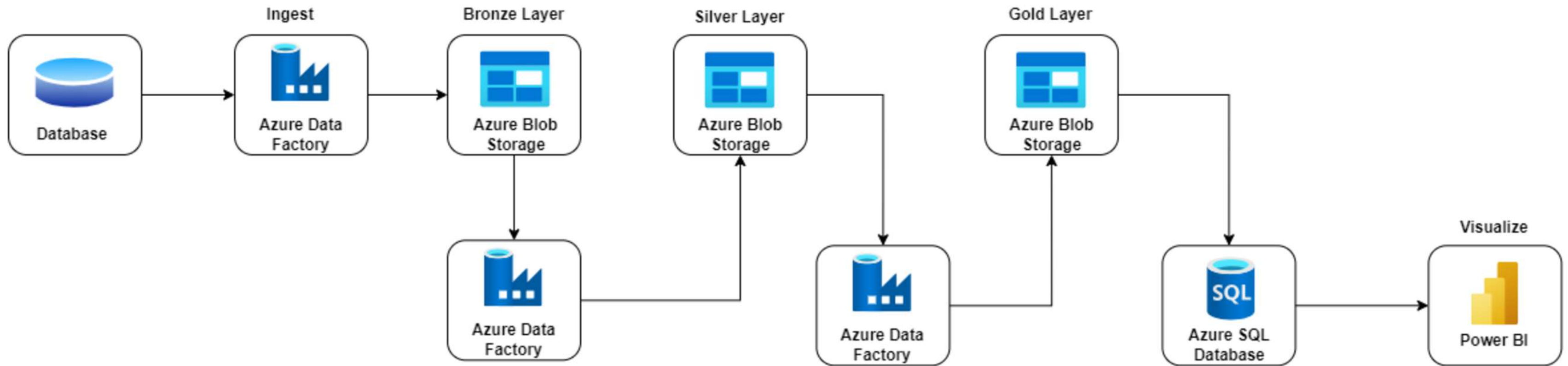


Azure
Data Factory

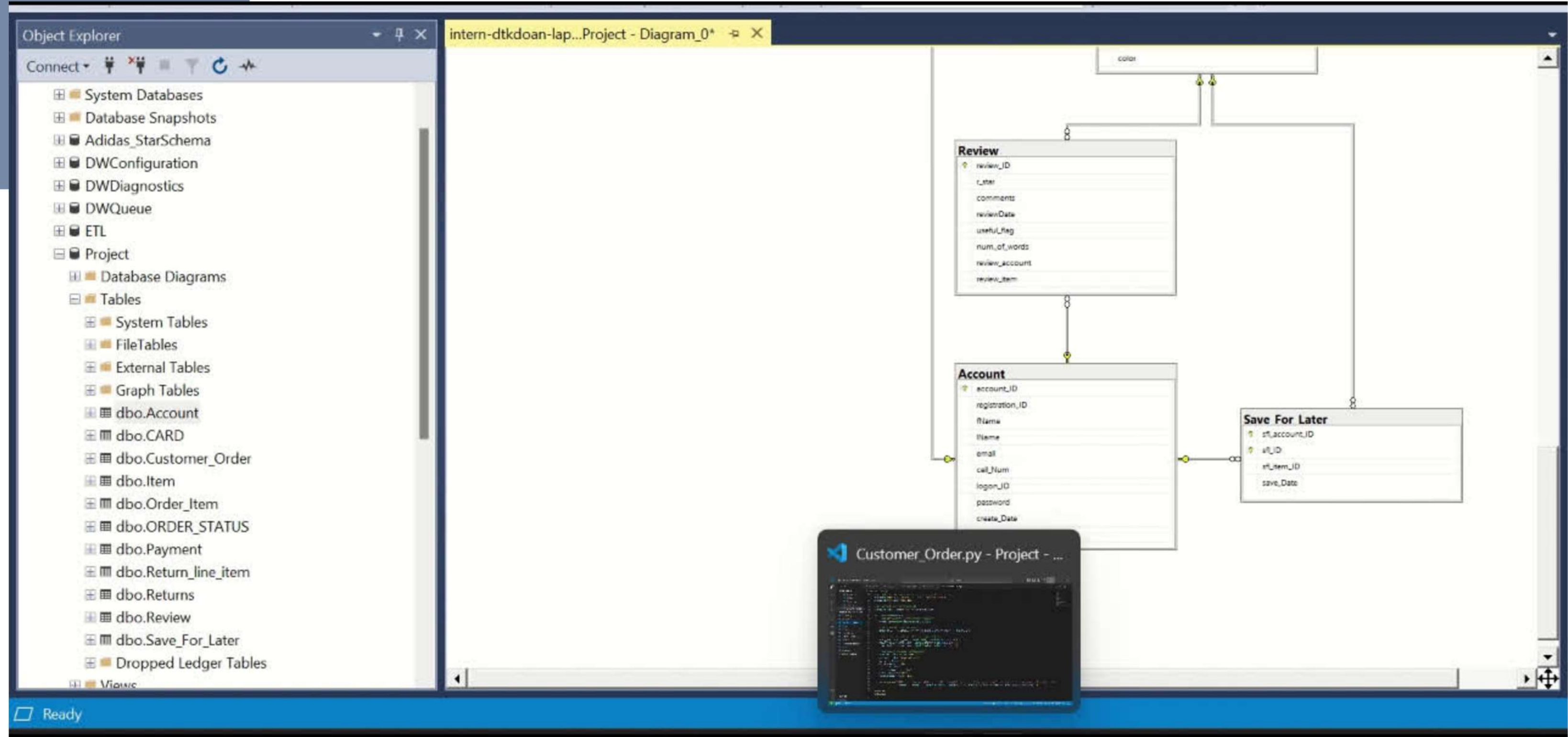


BUILD PIPELINE ON CLOUD

DATA FLOW



04. DEMO



04. DEMO



databricks



BUILD PIPELINE ON CLOUD

*Thank
You*