# Building scalable and automated Databricks-oriented Data Platform from scratch

Wojciech Pratkowiecki

**DATUMO** **SQL Day**

**17th edition**
**SQLDay Conference**
12-14 May 2025, WROCŁAW + ONLINE

─────── Platinum sponsors ───────

**DATUMO**     **Woodler**

─────── Gold sponsors ───────

**lingaro**     **VOLVO**     TECHNOLOGY INNOVATION DATA KNOWLEDGE **tidk**

**software one**     **Dataedo**

─────── Silver sponsors ───────

Capgemini     DATA RISE LAB     OM     elitmind     C&F     kursy FABRIC

# Agenda

- The story

- The challenge

- The solution

- Lessons learned

# The story

**Leader of CEE-pharmaceutical market with departments and subsidiaries across Europe**

Production department

Regulatory department

Strategy department

R&D department

**1** Unify and centralize data pipelines and analysis

**2** Standardize data management and processing

**3** Provide tailored Data Platform to each tenant
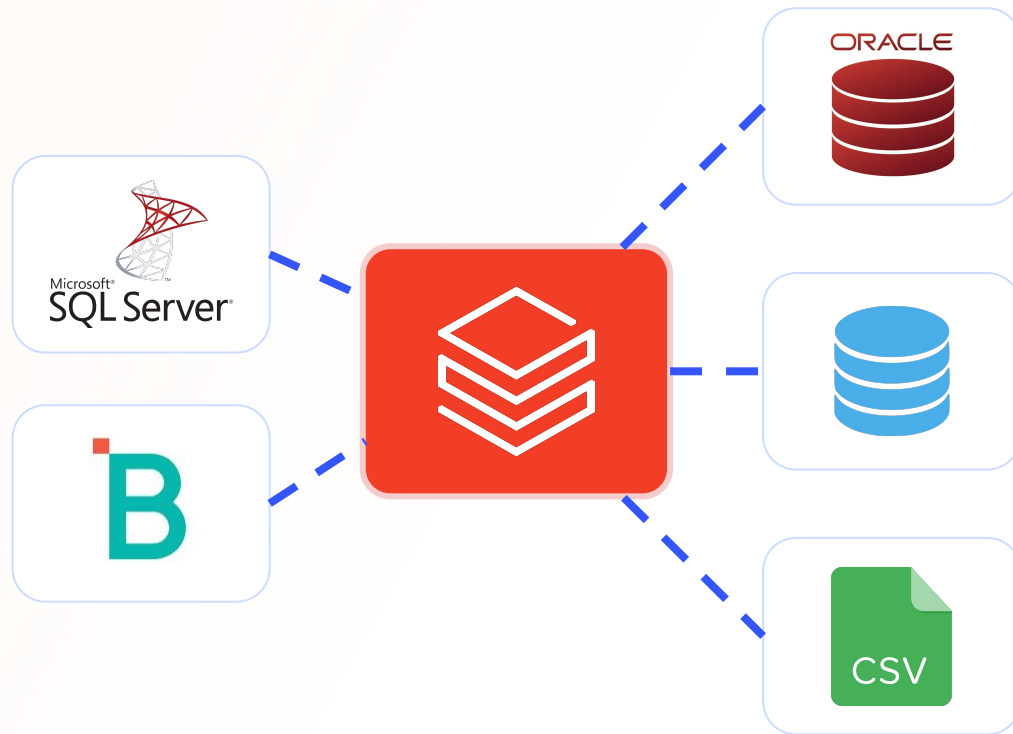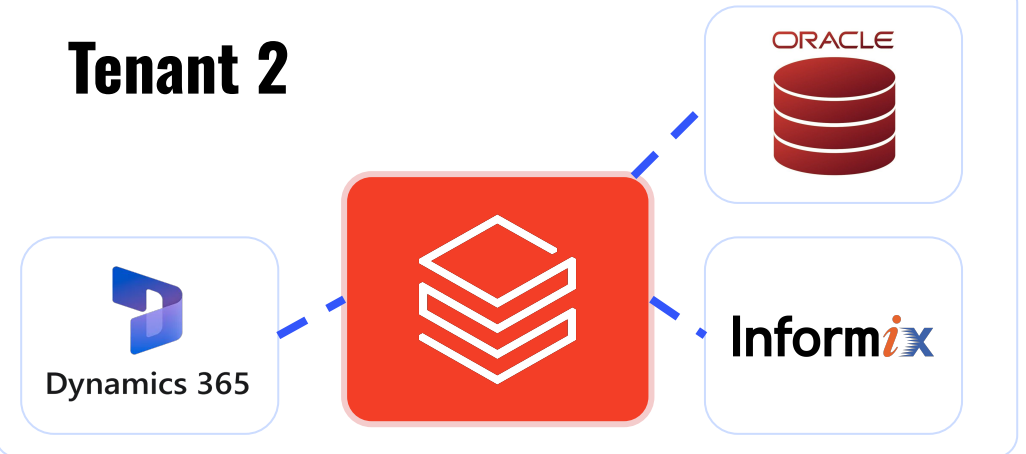
**4** Data mesh approach for datasets exchange

# The story

# The challenge

- Datasets are spread across various systems

- Introduced platform must be standardized and scalable

- All the resources should be implemented in modular approach

- Tenant-specific adjustments to each of platform's elements are required

- All of the Data Platform resources and processes need to be automatized

- Significant value for technical and non-technical users is expected

# The solution



**Functional Framework**

**Infrastructure Framework**

**Datastructure Framework**

**Pipelines Framework**

**CI/CD**

# The solution

DATUMO    SQL Day

**Functional Framework**
(Python / Apache Spark)

- ADLS integration with Autoloader and native Spark

- JDBC sources integration in incremental/full mode

- Tenant-specific implementation for integrated sources

- Common functions used for Silver/Gold layer creation

- Utilities for analysis and notebooks development

- Data Quality facilities

# The solution

**DATUMO** | **SQL Day**

**Infrastructure Framework**

- Terraform IaC implementation
- Blueprints of multi-service platform components
- Modules defining adjusted templates for services
- Simple reproduction across environments
- Access management
- Automated management of all platform resources

# The solution



**Datastructure Framework**

- Terraform modules for Unity Catalog resources management
- Modules defining templates for catalogs, schemas, tables, external locations etc.
- UC resources access management
- Definitions of particular catalogs/schemas
- Table's schema evolution management
- Automated management of all UC resources

# The solution

DATUMO  SQL Day

**Pipelines Framework**

- Data pipelines as Databricks Workflows
- Implemented and managed with Databricks Asset Bundles
- CLI and templates simplifying pipeline creation - single YAML file to define a pipeline
- Management of cross-environment deployment
- Custom configuration of jobs/clusters across environments
- Automated deployment of Workflows to Workspaces

# The solution



Developer implementing a new feature → **Push** → Technical repository → **Trigger** → CI/CD with tests and deployment → **Publish** → Packages artifactory ← **Use** ← Databricks job

Data analyst defining SQLs for an ETL pipeline → **Push** → Pipelines repository → **Trigger** → CI/CD deploying pipeline → **Deployment** → Databricks Workspace

# The solution

DATUMO    SQL Day



```python
class BaseDataIngestor(ABC):
    """
    Abstract base class for data ingestion.
    """

    def __init__(self, config: BaseIngestionConfig, spark: SparkSession):
        self.config = config
        self.spark = spark

    @cached_property
    def _target_table(self) -> Table:
        return Table(
            catalog_name=self.config.target_catalog,
            schema_name=self.config.target_schema,
            table_name=self.config.target_table,
        )

    def _verify_table_existence(self) -> None:
        if not self.spark.catalog.tableExists(self._target_table.full_name):
            raise ConfigurationValidationError(
                f"Table {self._target_table.full_name} does not exist."
            )

    def _add_technical_columns(self, df: DataFrame) -> DataFrame:
        current_timestamp = F.current_timestamp()
        modified_by = self.config.modified_by or self.spark.conf.get("spark.app.name")

        return (
            df.withColumn(TechnicalColumns.LOADED_AT, current_timestamp)
            .withColumn(TechnicalColumns.MODIFIED_AT, current_timestamp)
            .withColumn(TechnicalColumns.MODIFIED_BY, F.lit(modified_by))
            .withColumn(TechnicalColumns.IS_PROCESSED, F.lit(False))
            .withColumn(TechnicalColumns.ROW_ID, F.expr("uuid()"))
        )

    @abstractmethod
    def _ingest_from_source_to_bronze(self) -> None:
        raise NotImplementedError

    def ingest(self) -> None:
        self._verify_table_existence()
        self._ingest_from_source_to_bronze()
```

# The solution

# The solution

# The solution

# The solution

# The solution

# The solution

## Short Time To Market for each tenant

Deployment of core cloud resources, setup of Unity Catalog → Provisioning of Azure & Databricks infrastructure → Adjustments to tenant-specific sources → Definition of Unity Catalog resources →

Quick integration to Bronze layer with metadata-based approach → Implementation of jobs creating Silver/Gold layer → Automation of transformations with Databricks Workflow → Deployment to higher environments
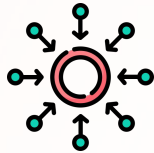
# The solution


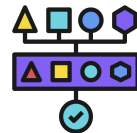
## Provided value

Centralized data platform enabling cross-source system analysis

Standardized data processing and management across numerous tenants

Automated and extendable solution allowing quick integration of new sources

Defining advanced data pipelines utilizing powerful Databricks features made easy

Modular and scalable solution enabling quick provisioning of adjusted Data Platform

# Lessons learned

- Define the data usage patterns with target users at the very beginning

- Short TTM introducing numerous novel services can be overwhelming

- Infrastructure and data structure automated management pays off

- Lack of date-based scheduling in Databricks hurts, but can be overcomed

- Strict schema management can be cumbersome, however smart scripts/notebooks can be an answer

Q&A

# Thank you for your attention