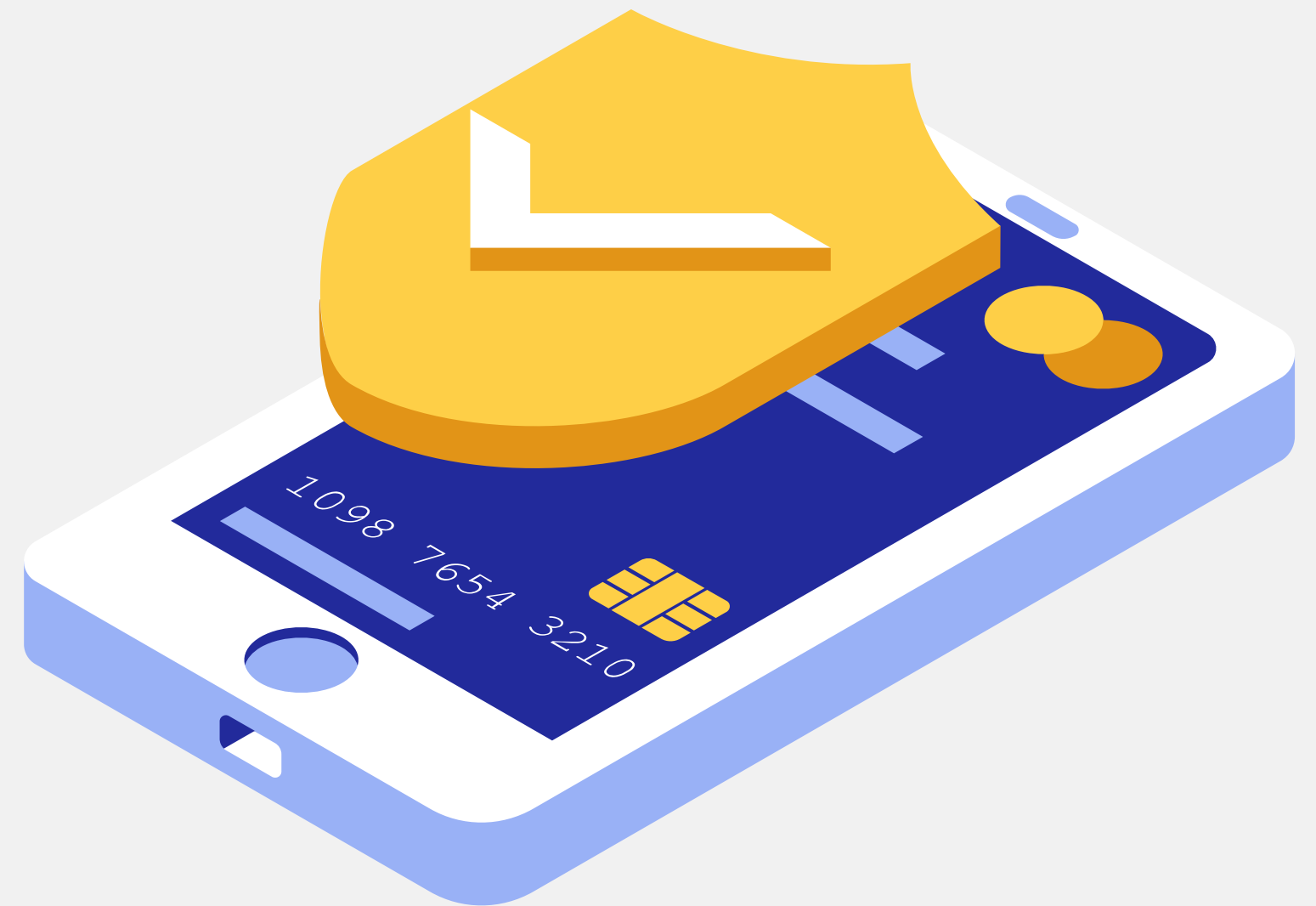


# Data Processing Pipelines

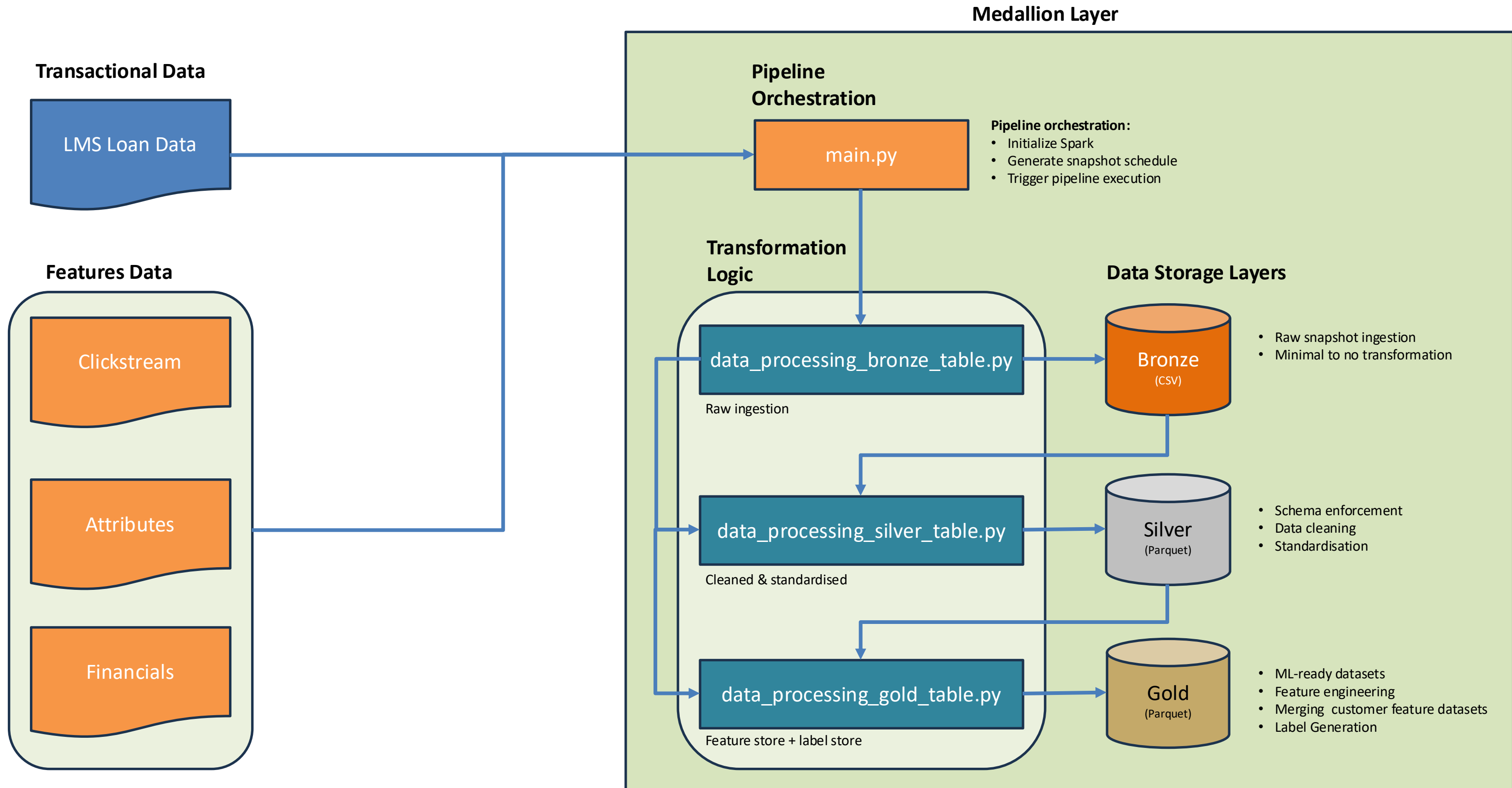
# Assignment Objective

To build an end to end data processing pipeline that transforms raw loan and customer data into ML-ready feature and label stores using the Medallion Architecture.

The pipeline applies data ingestion, cleaning, schema enforcement, feature engineering and modular PySpark processing across **bronze, silver and gold** layers. It also uses Docker and GitHub to support a more reproducible machine learning engineering workflow.



# Medallion Architecture Design



# Bronze Layer Implementation: Raw Data Ingestion

Raw LMS and customer feature datasets are ingested using PySpark filtered by snapshot\_date and stored as monthly CSV partitions in the bronze datamart. Minimal transformation is applied to preserve source data integrity and maintain traceability for downstream silver and gold processing.



## Input Sources

- lms\_loan\_daily.csv
- feature\_clickstream.csv
- features\_attributes.csv
- features\_financials.csv



## Implementation Steps

1. Read raw CSV files using PySpark.
2. Filter records by monthly snapshot\_date.
3. Save each monthly snapshot as a separate bronze file.
4. Store outputs under the bronze datamart folders.



## Output Tables

- bronze\_loan\_daily\_YYYY\_MM\_DD.csv
- bronze\_clickstream\_YYYY\_MM\_DD.csv
- bronze\_attributes\_YYYY\_MM\_DD.csv
- bronze\_financials\_YYYY\_MM\_DD.csv



## Design Rationale

- Preserves raw source data for traceability.
- Keeps bronze layer close to the original source.
- Enables repeatable monthly backfill processing.
- Foundation for downstream silver processing.

# Silver Layer Implementation: Cleaning and Standardisation

The silver layer reads bronze monthly snapshots, enforces schema, removes duplicates, cleans invalid numeric values, converts date fields and derives repayment indicators such as mob and dpd.

The cleaned outputs are stored as Parquet files and used as the foundation for gold-level feature and label generation.



## Input Sources

- bronze\_loan\_daily\_YYYY\_MM\_DD.csv
- bronze\_clickstream\_YYYY\_MM\_DD.csv
- bronze\_attributes\_YYYY\_MM\_DD.csv
- bronze\_financials\_YYYY\_MM\_DD.csv



## Key Processing Steps

1. Load monthly bronze loan and feature files.
2. Enforce column data types.
3. Remove duplicate feature records.
4. Clean invalid characters such as \_ in numeric fields.
5. Extreme or unrealistic values were replaced with nulls.
6. Convert snapshot\_date into date format.
7. Create loan repayment **features\*** such as: dpd, mob  
installments\_missed, first\_missed\_date



## Output Tables

- silver\_loan\_daily\_YYYY\_MM\_DD.parquet
- silver\_clickstream\_YYYY\_MM\_DD.parquet
- silver\_attributes\_YYYY\_MM\_DD.parquet
- silver\_financials\_YYYY\_MM\_DD.parquet



## Design Rationale

- Improves data quality before feature engineering.
- Standardises schema for downstream joins.
- Converts raw CSV data into Parquet for scalable storage and analytics.
- Creates reliable intermediate tables for gold feature and label stores.

*features\* are further elaborated in slide 7*

# Gold Layer Implementation

The gold layer converts cleaned silver tables into final ML-ready outputs. It creates a default label based on 30dpd\_6mob, joins customer feature datasets, engineers financial risk indicators, and stores the final feature and label stores as Parquet files for future model training.

Label store and feature store are stored separately to support reusable training pipelines.



## Input Sources

- silver\_loan\_daily\_YYYY\_MM\_DD.parquet
- silver\_clickstream\_YYYY\_MM\_DD.parquet
- silver\_attributes\_YYYY\_MM\_DD.parquet
- silver\_financials\_YYYY\_MM\_DD.parquet



## Output Tables

- gold\_label\_store\_YYYY\_MM\_DD.parquet
- gold\_feature\_store\_YYYY\_MM\_DD.parquet

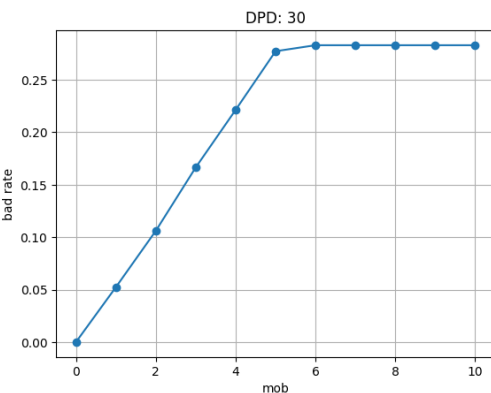


## Design Rationale

- The gold layer separates labels and features into reusable ML-ready stores.
- This supports future model training while keeping feature engineering and label generation traceable within the Medallion Architecture.

## Key Processing Steps

1. Load cleaned silver loan and feature tables.
2. Generate default labels from loan repayment behaviour.
3. Filter loans at mob = 6.
4. Assign default label using dpd  $\geq$  30.
5. Join attributes, financials, and clickstream features by Customer\_ID and snapshot\_date.
6. Create engineered **features\***: debt\_to\_income\_ratio, emi\_to\_salary\_ratio, high\_credit\_utilization\_flag



*features\* are further elaborated in slide 7*

# Feature Engineering & Label Design

Layer	Feature/ Label	Formula/ Logic	Business Meaning
Silver	dpd	Days between snapshot date and first missed payment	Measures delinquency severity
	mob	Instalment number	Loan repayment progression
	installments_missed	Ceiling of Overdue_amt/ due_amnt	Rounded up number of missed installment
	first_missed_date	snapshot_date - installments_missed months	Capture the 1 <sup>st</sup> time period where the payment was due
Gold	label	$dpd \geq 30$ at $mob = 6$	Binary loan default prediction target
	debt_to_income_ratio	Outstanding Debt / Annual Income	Debt burden risk
	emi_to_salary_ratio	EMI / Monthly Salary	Repayment affordability
	high_credit_utilization_flag	Utilization > 50%	Indicator of high credit usage risk

# Data Quality & Validation Checks ✓

<p><b>Duplicate Checks</b></p>	<p>Duplicate records were removed from feature datasets using dropDuplicates() before downstream joins and feature engineering.</p> <p>Validation Logic</p> <ul style="list-style-type: none"> <li>• Checked duplicate Customer_ID and snapshot_date combinations.</li> <li>• Prevented duplicate customer feature records in the gold feature store.</li> </ul>
<p><b>Schema enforcement</b></p>	<p>Column data types were explicitly enforced using PySpark casting to standardise loan and feature datasets.</p>
<p><b>Invalid Value Cleaning (Silver Layer)</b></p>	<ul style="list-style-type: none"> <li>• Removed _ characters from numeric financial columns.</li> <li>• Converted cleaned values into numeric datatypes.</li> </ul> <p>Extreme or unrealistic values were replaced with nulls based on cutoff thresholds derived from binning the latest cumulative dataset.</p> <ul style="list-style-type: none"> <li>• Age values &lt; 0 or &gt; 100 → NULL.</li> <li>• Num_Bank_Accounts, Num_Credit_Card values &lt;0 or &gt; 15 → NULL.</li> <li>• Interest_Rate values &lt; 0 or &gt; 100 → NULL.</li> <li>• Delay_from_due_date, Num_Credit_Inquiries values &lt;0 or &gt; 15 → NULL.</li> </ul>
<p><b>Derived Repayment Validation</b></p>	<p>Repayment behaviour features were engineered and validated using loan overdue information.</p> <ul style="list-style-type: none"> <li>• mob</li> <li>• installments_missed</li> <li>• first_missed_date</li> <li>• dpd</li> </ul>



# ML Compataility and Leakage Control

## ML Compatibility



- Gold feature and label stores are structured for downstream model training.
- Feature and label stores are kept separate to prevent leakage, improve reusability and support different model training use cases (can be joined later during model training using Customer\_ID & snapshot\_date).
- Parquet format improves scalability and analytical efficiency.

## Leakage Control



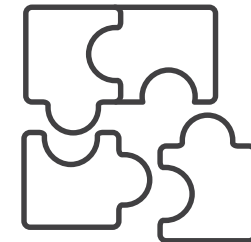
- Labels are generated using repayment behaviour at mob = 6.
- Feature joins use historical customer information available at snapshot time.
- Future repayment outcomes are used only for label generation, not feature creation.

# Limitations and Improvements



## Current Limitations

- Pipeline currently runs locally using local[\*] Spark mode.
- Batch processing only; no real-time streaming ingestion. Train, validation, test, and OOT data split are not yet implemented.
- Feature store metadata cataloging is not implemented.
- No workflow orchestration tool such as Airflow is used yet.
- No automated model monitoring or retraining pipeline.



## Future Improvements

- Deploy pipeline on distributed cloud infrastructure.
- Integrate Airflow for automated scheduling and orchestration.
- Add train, validation, test, and OOT dataset generation.
- Introduce online feature store for real-time inference.
- Add model monitoring and data drift detection.
- Implement CI/CD pipeline and container registry deployment.