

Implementing an End-to-End Demand Forecasting Solution Through Databricks and MLflow

ORGANIZED BY  databricks



Ivana Pejeva

Cloud Solution Architect, Microsoft



Yoshi Coppens

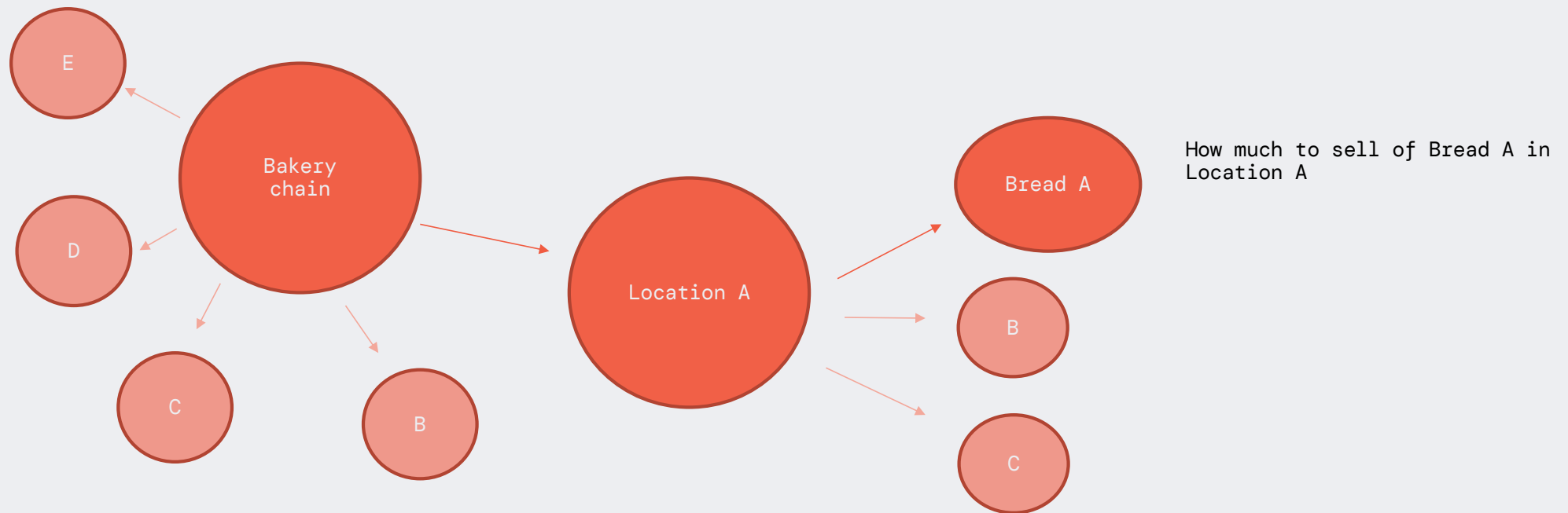
Data Engineer, element61

Outline

- Problem Setting
- Ingestion
- Cleaning, transforming, enriching
- Model selection and training
- Scoring
- Feeding it back to the client's system
- Running the model in production

Problem setting

Scope of Project



Ingestion

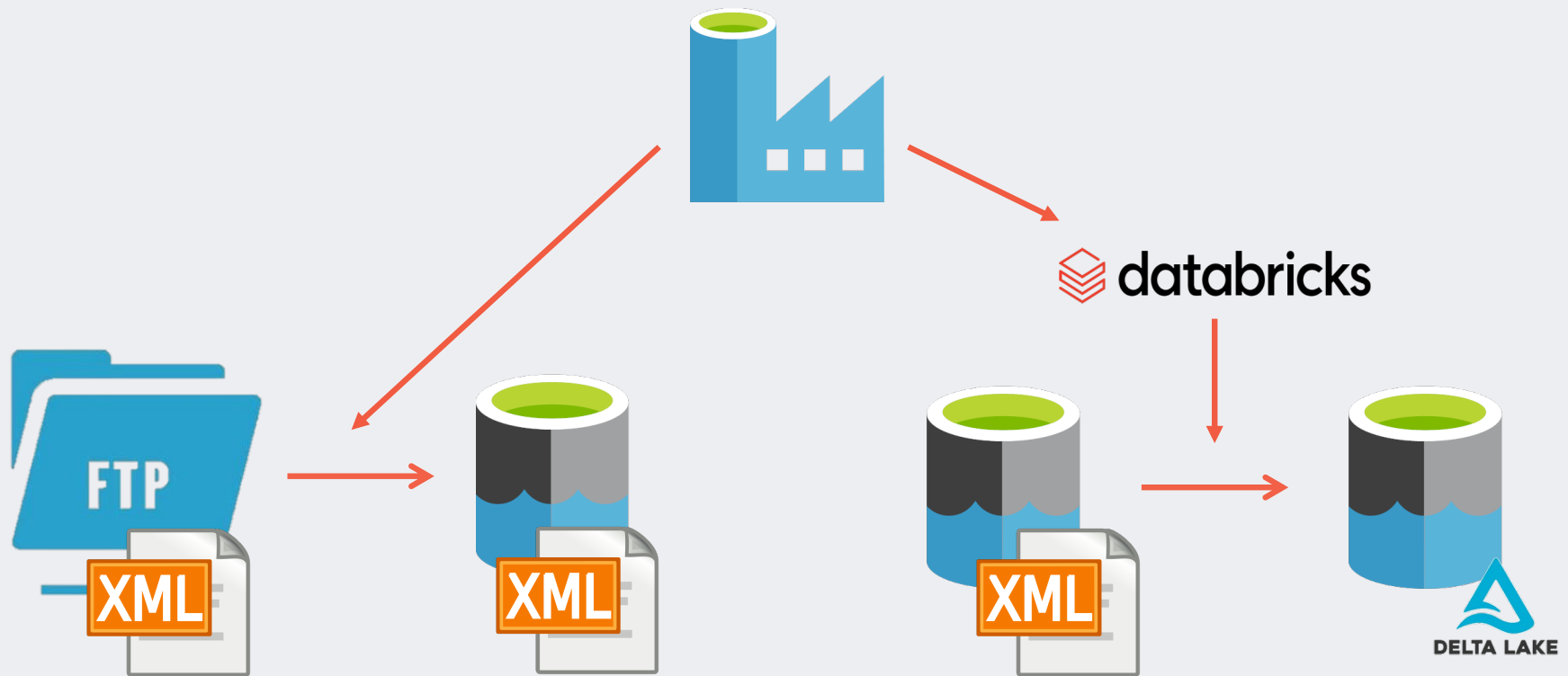
Ingesting the files

XML files on an FTP server??

- The provided salesdata is made available every day on an FTP-server
- File format is XML
- Azure Data Factory used to copy the binary files from FTP server
- Use spark-xml package in Databricks to translate from XML
https://mvnrepository.com/artifact/com.databricks/spark-xml_2.12

Ingesting the files

Structure of Data Lake and ingestion



Ingesting the files

Learnings

Ideally the data gets fed into your systems in an optimal way but when you cannot choose the setup, you have to make it work, ideally not having to incorporate more and more different tools

Combining Azure Data Factory with Azure Databricks
gives you quite a lot of power

Cleaning,
transforming,
enriching

Using extra data sources

Client based



Store Information



Product Information



Store Holiday
Information



Competitor Holiday
Information

Using extra data sources

External sources



Feature Engineering

- Feedback from client is important
 - 'felt that model was too slow with incorporating a big shift in demand trend'
 - incorporate **Sales of last week**
- Analyzing performance is important
 - 'we saw that there seemed to be a bias to underestimate usually'
 - when they sell 10 but only had an inventory of 10 vs selling 10 when there is 20 delivered, huge difference => incorporating a **SoldOut** factor

Model Selection and Training

Picking a ML model

Learnings

ARIMA vs ML

- Very dependent on data and problem when one is performing better than the other
- Big advantage in ML => make a joint model from multiple time series

Picking LightGBM

- Focus was not on complicated ML due to time constraints
- We wanted some level of **interpretability** for the client

Model Granularity

1 Model for all stores

- Features such as Adoption Level, Location, Population, Events etc.
- We did not get the required information

1 Model for all products

- Difficult, pain au chocolat vs croissant vs éclair?

1 Model per product group

- Yes, can distinguish some fungible groups, such as Bread, Confiserie, Small Cakes

Making use of parallelism

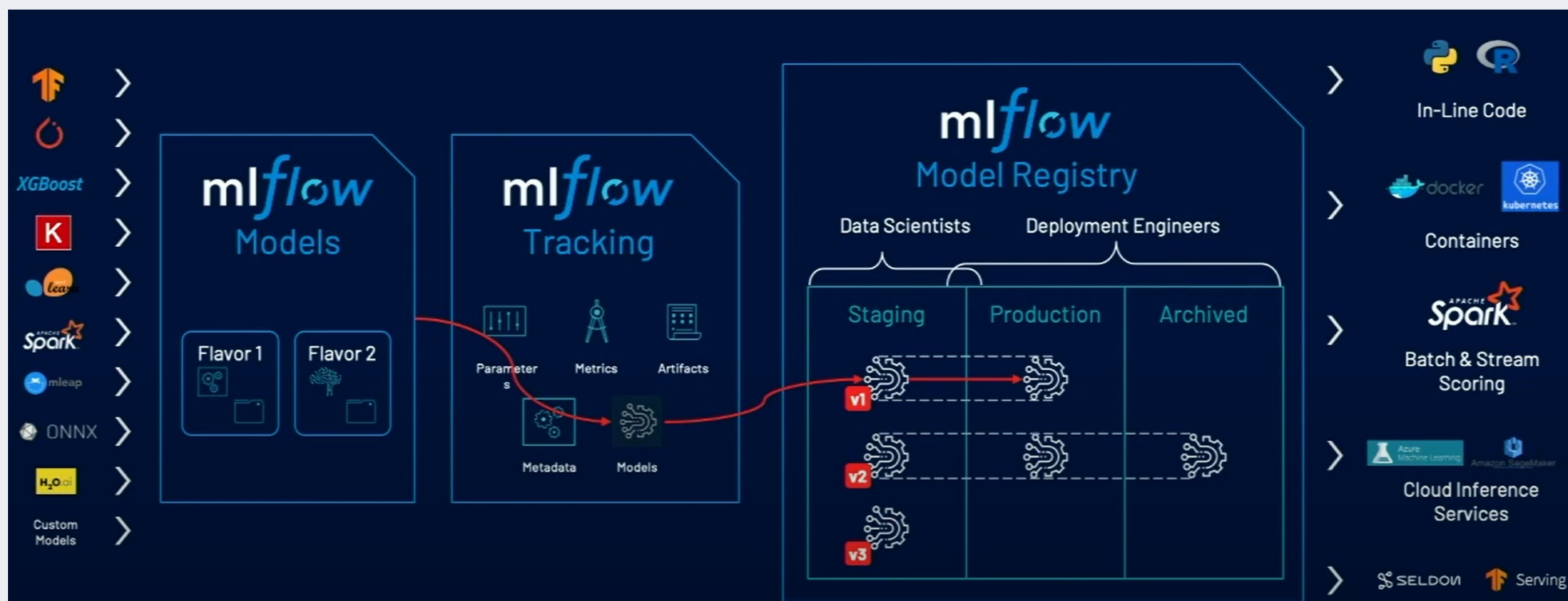
Training multiple models at once using Pandas UDFs

- We have multiple productgroups and multiple stores => **248 models**
- Spark MLLib is not made for training multiple models at once
- A for loop of course is also very unperformant
- **Grouped Map Pandas UDFs! For both Training and Scoring!**

```
df  
  .groupBy(['StoreKey', 'ProductGroupKey'])  
  .applyInPandas(trainingUDF, schema = schema_training)
```


Tracking Performance and experiments

Introducing MLFlow



Source: Spark AI Summit - Richard Zang & Denny Lee

Tracking Performance and experiments

Continued MLFlow - code for training

```
def trainUdf (pandasDf:object):  
    ...  
    with mlflow.start_run(experiment_id=experiment_id):  
        gbm = lgb.train(...)   
        X_test['pred'] = gbm.predict(...)   
        ...  
        mlflow.log_metric("MAPE", MAPE)   
        ...  
        mlflow.lightgbm.log_model(lgb_model = gbm,   
        artifact_path=f"{storeKey}_{productGroupKey}_model",   
        registered_model_name=f"{storeKey}_{productGroupKey}_model")   
    ...
```

Tracking Performance and experiments

Continued MLFlow - code for scoring

```
def predictUdf (pandasDf:object):  
    ...  
    model_name = f"{storeKey}_{productGroupKey}_model"  
    model_uri = f"models:{model_name}/Production"  
    model = mlflow.lightgbm.load_model(model_uri)  
    predictions = model.predict(...)  
    ...
```

Tracking Performance and experiments

Databricks - Experiments

The screenshot displays the Databricks Experiments page for Experiment ID: 3558968592689005. The interface includes a sidebar with navigation icons, a top bar with the experiment ID, and a main content area. The main content area has a 'Description' section with an 'Edit' link, a row of action buttons (Refresh, Compare, Delete, Download CSV), and a search bar with a filter dropdown set to 'All time'. Below these are view toggles (list/table), a 'Columns' button, a toggle for 'Only show differences', and a search query: 'metrics.rmse < 1 and params.model = "tree"'. The table shows 5 matching runs. The 'Metrics' section of the table is circled in red.

Experiment ID: 3558968592689005

► Description [Edit](#)

[Refresh](#) [Compare](#) [Delete](#) [Download CSV](#) [Start Time](#) [All time](#)

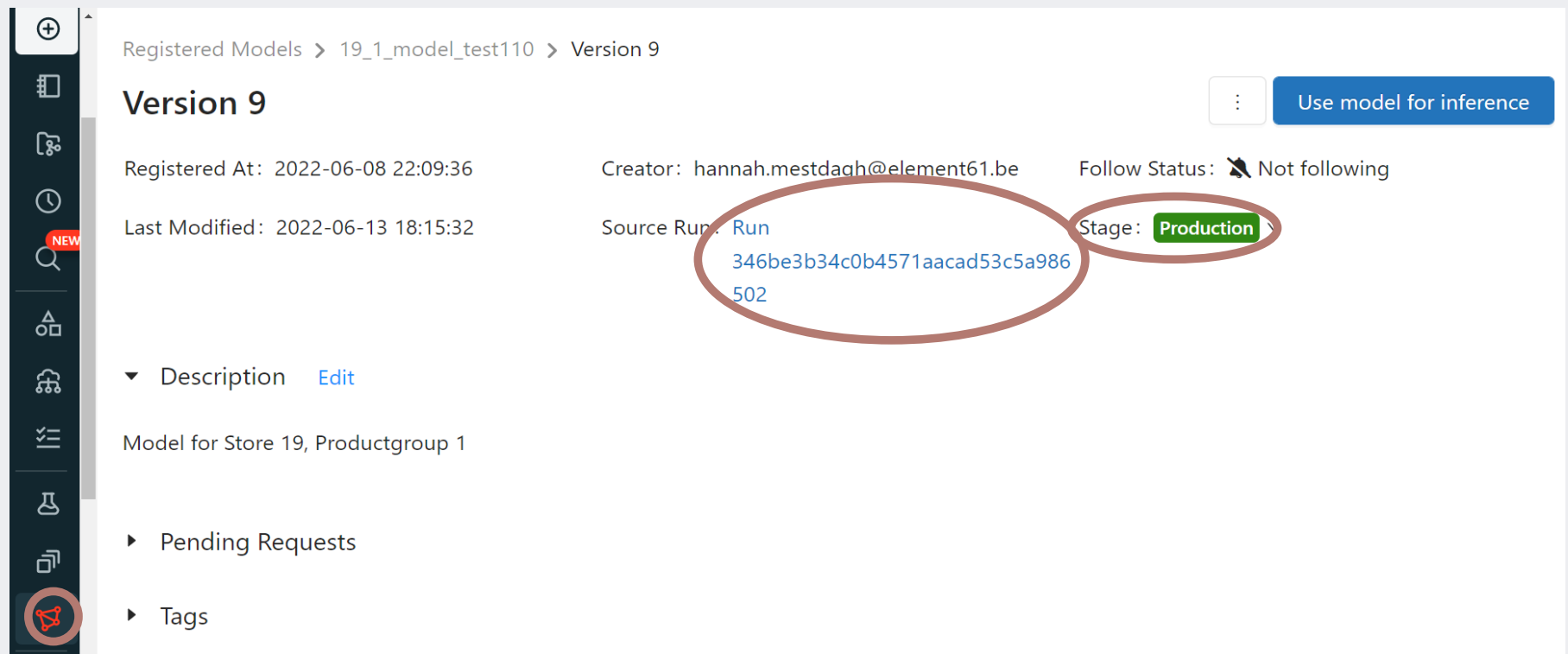
[Columns](#) Only show differences ☒ [Search](#) [Filter](#) [Clear](#)

Showing 5 matching runs

						Metrics >		
	↓ Start Time	Duration	User	Source	Models	MAPE	MSE	RMSE
<input type="checkbox"/>	✓ 4 days ago	4.6min	hannah.mes...	2. Model Ti	31_5_model.../5	0.642	5.729	2.393
<input type="checkbox"/>	✓ 11 days ago	11.6s	hannah.mes...	2. Model Ti	31_5_model.../4	0.652	5.769	2.402
<input type="checkbox"/>	✓ 1 month ago	11.5s	hannah.mes...	2. Model Ti	31_5_model.../3	0.662	5.319	2.306
<input type="checkbox"/>	✓ 1 month ago	10.9s	hannah.mes...	2. Model Ti	31_5_model.../2	0.822	4.158	2.039

Tracking Performance and experiments

Databricks - Models



The screenshot shows the Databricks Models interface. On the left is a sidebar with navigation icons, including a red circle icon at the bottom. The main content area displays details for 'Version 9' of a model. The breadcrumb path is 'Registered Models > 19_1_model_test110 > Version 9'. The title 'Version 9' is prominently displayed. To the right of the title is a menu icon and a blue button labeled 'Use model for inference'. Below the title, the 'Registered At' timestamp is '2022-06-08 22:09:36' and the 'Last Modified' timestamp is '2022-06-13 18:15:32'. The 'Creator' is listed as 'hannah.mestdagh@element61.be' and the 'Follow Status' is 'Not following'. The 'Source Run' is a blue link that has been circled in red, showing the ID '346be3b34c0b4571aacad53c5a986502'. The 'Stage' is a green button labeled 'Production', also circled in red. Below this, there is a 'Description' section with an 'Edit' link, showing the text 'Model for Store 19, Productgroup 1'. At the bottom, there are expandable sections for 'Pending Requests' and 'Tags'.

Registered Models > 19_1_model_test110 > Version 9

Version 9

Registered At: 2022-06-08 22:09:36 Creator: hannah.mestdagh@element61.be Follow Status: Not following

Last Modified: 2022-06-13 18:15:32 Source Run: [Run](#) Stage: **Production**

[346be3b34c0b4571aacad53c5a986502](#)

▼ Description [Edit](#)

Model for Store 19, Productgroup 1

► Pending Requests

► Tags

Scoring

Which metrics to use?

Can't go wrong with MAPE

- MAPE = Mean Absolute Percentage Error
- Discussion with client: Overestimating is better than underestimating!
 - Option 1: Change our metric to punish missing sales harder than having leftovers (using production cost and opportunity loss)
 - **Option 2: Instead of using the quantity sold, use quantity + wanted margin**

Working with reliability buckets

- We want to know for which products we can have a reliable predictions

What can we predict
well

What is okay

What is not
necessarily great

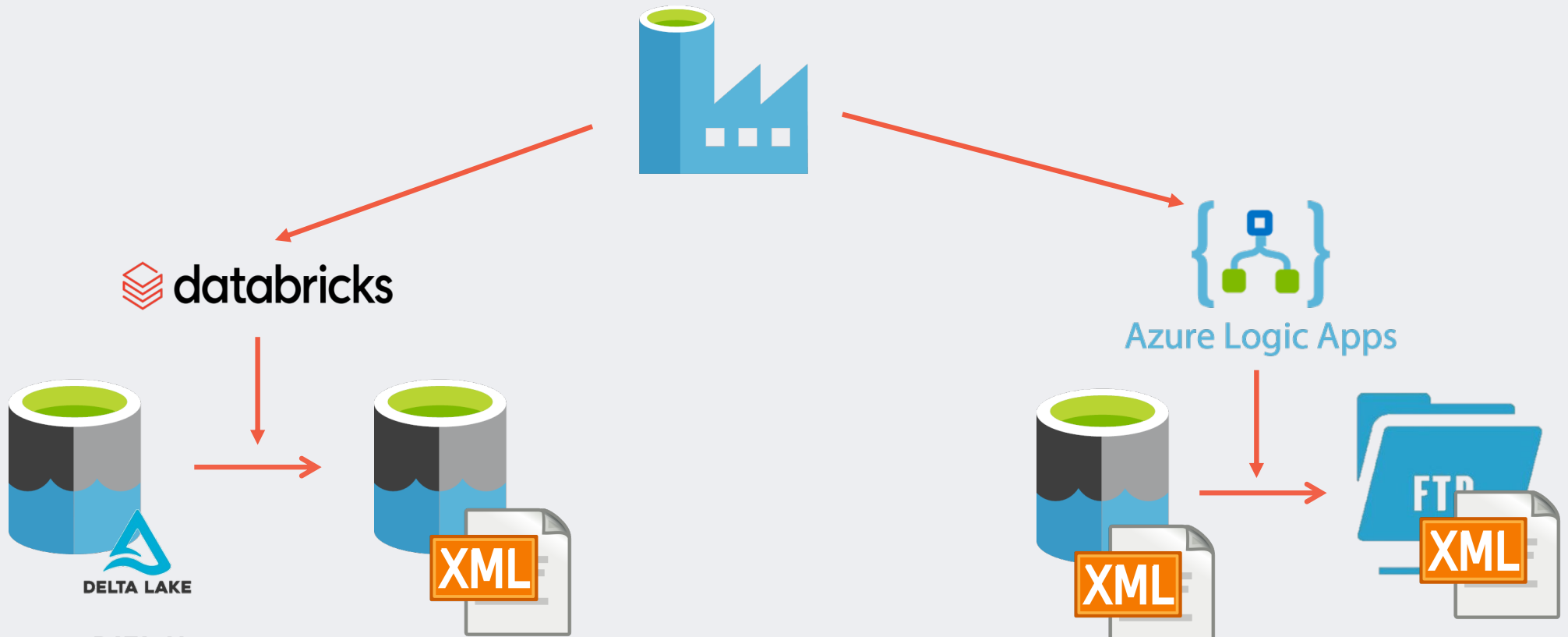
Feeding it back
into the
client's systems

Feeding it back into the client's systems

XML and FTP again!

- Their ERP system needs XML files again!
 - Translate Delta Lake tables back to XML using spark-xml package in Databricks
- XML files need to be transferred back to the FTP server
 - Azure Data Factory actually cannot do this “out-of-the-box”, can only copy to SFTP
 - Introduce **Azure Logic App** which can use FTP as a sync and has a lot more data connectors

Feeding it back into the client's systems



Running the model in production

Frequency of training and scoring

Multiple factors matter

Data can be weeks late

- Internal system can have big delays
- Forecast is needed every day
- Weather has a big impact so **scoring every day!**

Keep control of cost

- The cost of training is not negligible
- Rather more frequent scoring than more frequent training

Importance of new data

- For old stores, new data gets less and less important
- For new stores, it is the opposite
- They are expanding and for consistency's sake **training every week!**

Monitoring and Alerting

Mailing through Logic App, orchestrated in Data Factory

Data missing

- Due to operational difficulties, extensive checks needed => Mail
- Client input data not getting updated => Mail
- Follow up for new stores being added

Monitoring models

- Big improvement in MAPE => Mail
- Deteriorating MAPE vs last => Mail
- Future: Model Drift Detection

Conclusion

Conclusion



- In a couple of days we can get a Demand Forecasting model into production
- **Pandas UDFs can help us scale** hugely in the model training
- **MLFlow is a great MLOps tool** for tracking, registering Machine Learning models

DATA+AI
SUMMIT 2022

Thank you

Ivana Pejeva & Yoshi Coppens