

Data Contracts in action

powered by Python open source ecosystem

Agenda

- About me
- My challenges before data contracts
- Data contracts primer
- Streaming solution with a data contract
- Batch solution with a data contract
- Talk to me

About me

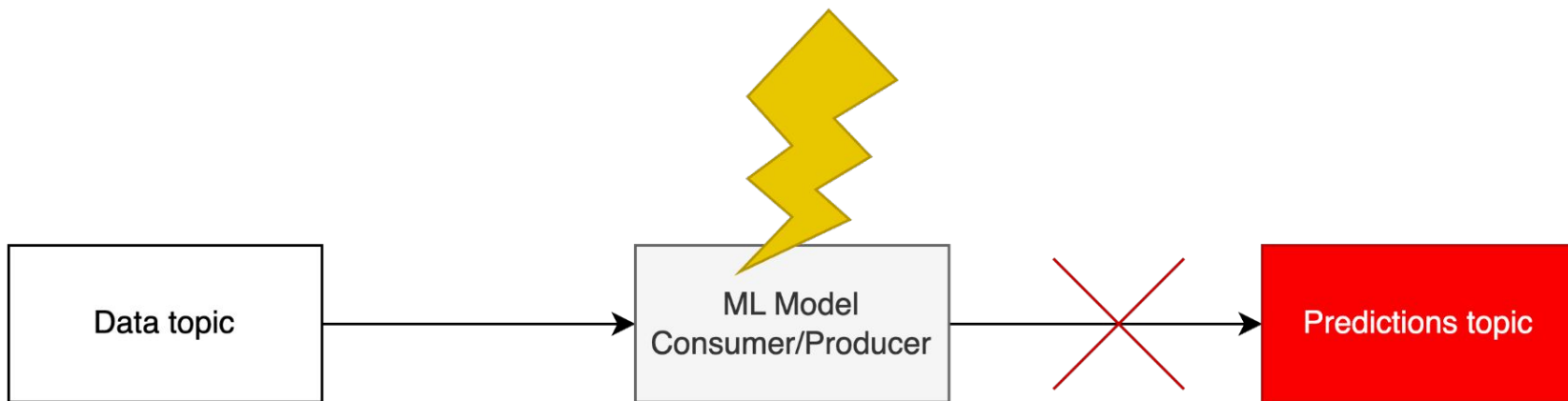
Principal MLOps and Data Engineer at Thoughtworks

Organiser of PyLadies Amsterdam

Microsoft AI MVP since 2021

My challenges before data contracts

Streaming challenge



My challenges before data contracts

Batch challenge



Value 1	Value 2	Value 3
Value 4	Value 5	Value 6
Value 7	Value 8	Value 9



Prediction 1	Prediction 2	Prediction 3
Prediction 4	Prediction 5	Prediction 6
Prediction 7	Prediction 8	Prediction 9

Data contracts primer

"A data contract is an agreed **interface** between the generators of data and its consumers. It sets the **expectations** around that data, defines how it should be **governed**, and facilitates the **explicit** generation of quality data that meets the business requirements."

Andrew Jones, Driving Data Quality with Data Contracts

Data contracts primer

interface

abstraction

agreement

confidence

Andrew Jones, Driving Data Quality with Data Contracts

Data contracts primer

expectations

the structure/schema

data quality checks

SLOs (completeness, timeliness, availability)

ownership and responsibilities

Andrew Jones, Driving Data Quality with Data Contracts

Data contracts primer

governed

PII
sensitivity

access

retention

anonymization

lineage

Andrew Jones, Driving Data Quality with Data Contracts

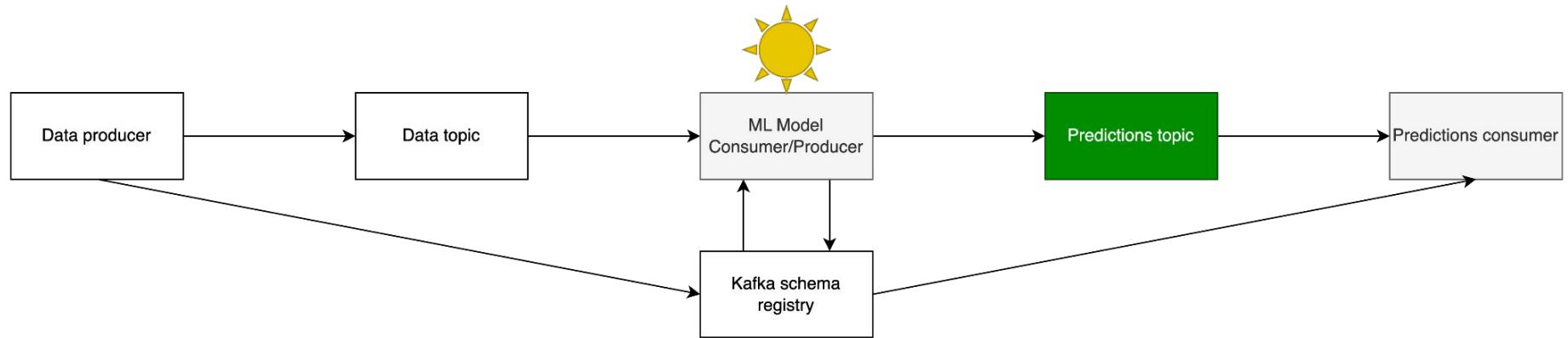
Data contracts primer

explicit

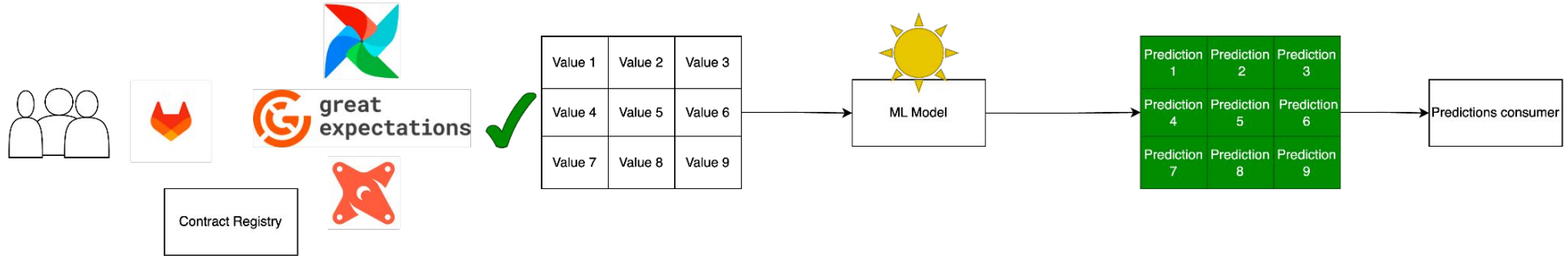
built for consumption

Andrew Jones, Driving Data Quality with Data Contracts

Streaming solution with a data contract



Batch solution with a data contract



Talk to me



About me



About PyLadies Amsterdam