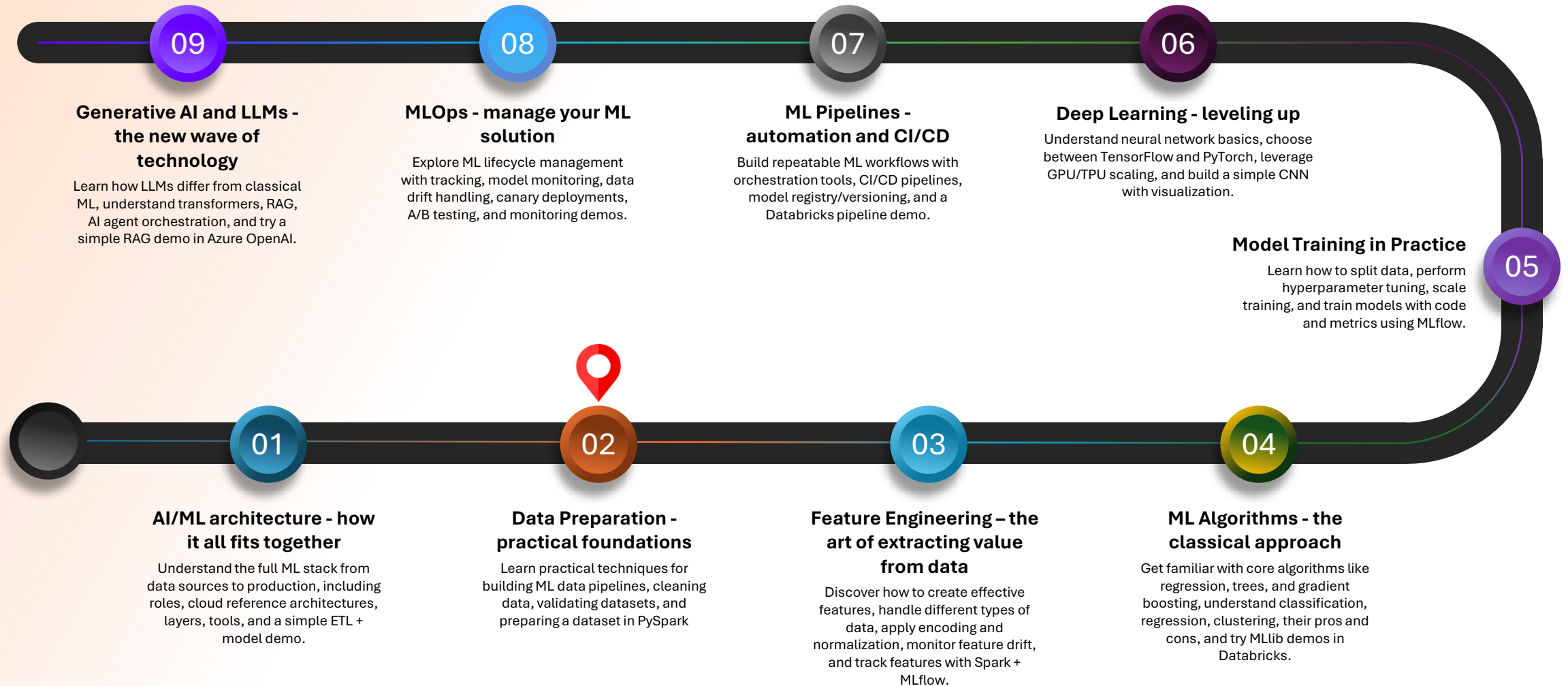


Data preparation - practical foundations

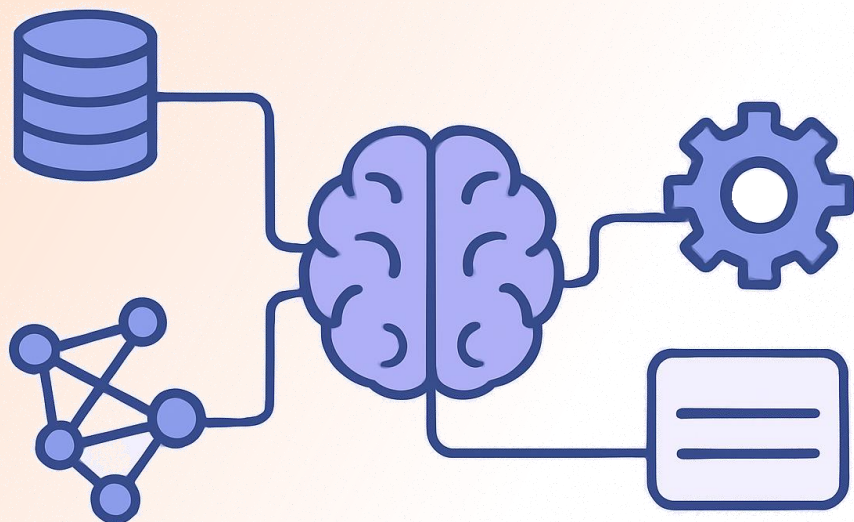
Maciej Kępa

Roadmap

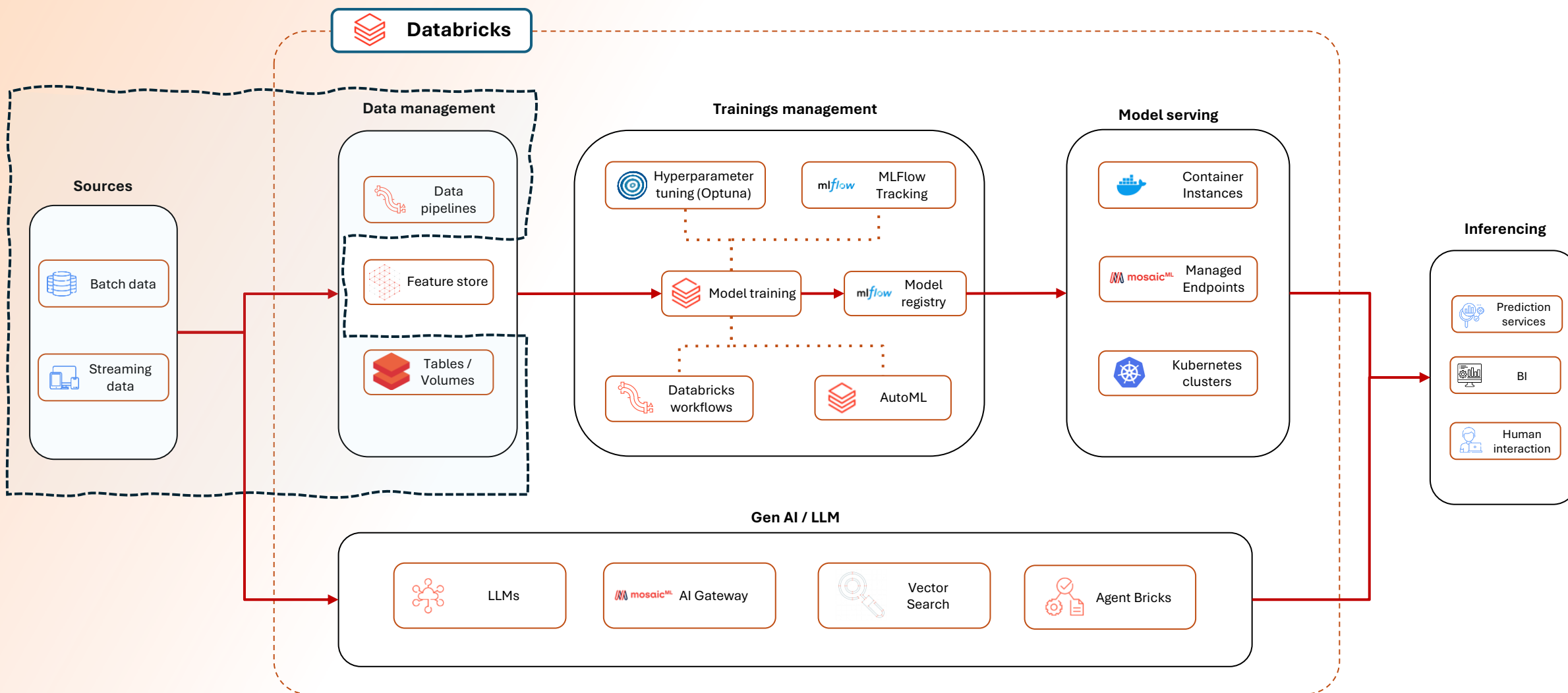


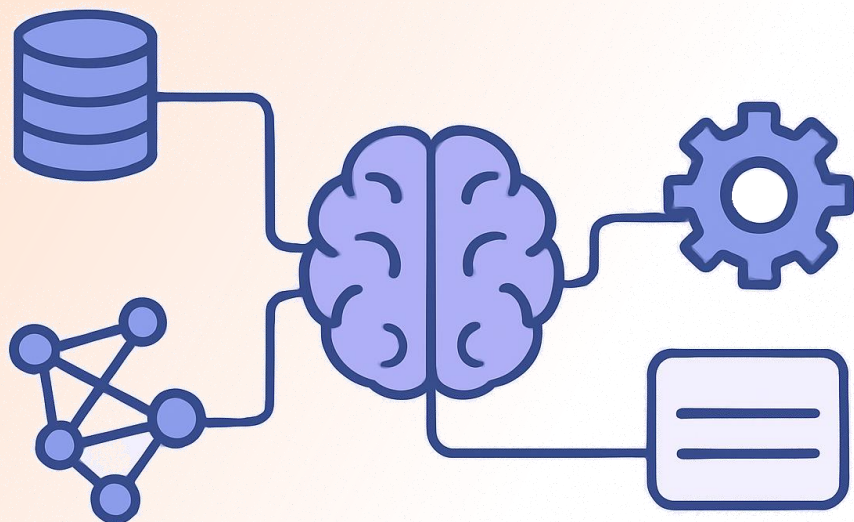
Agenda

1. Introduction
2. From raw to model-ready data
3. Workshop: data cleaning
4. Workshop: data validation
5. Workshop: data pipeline
6. Best practices



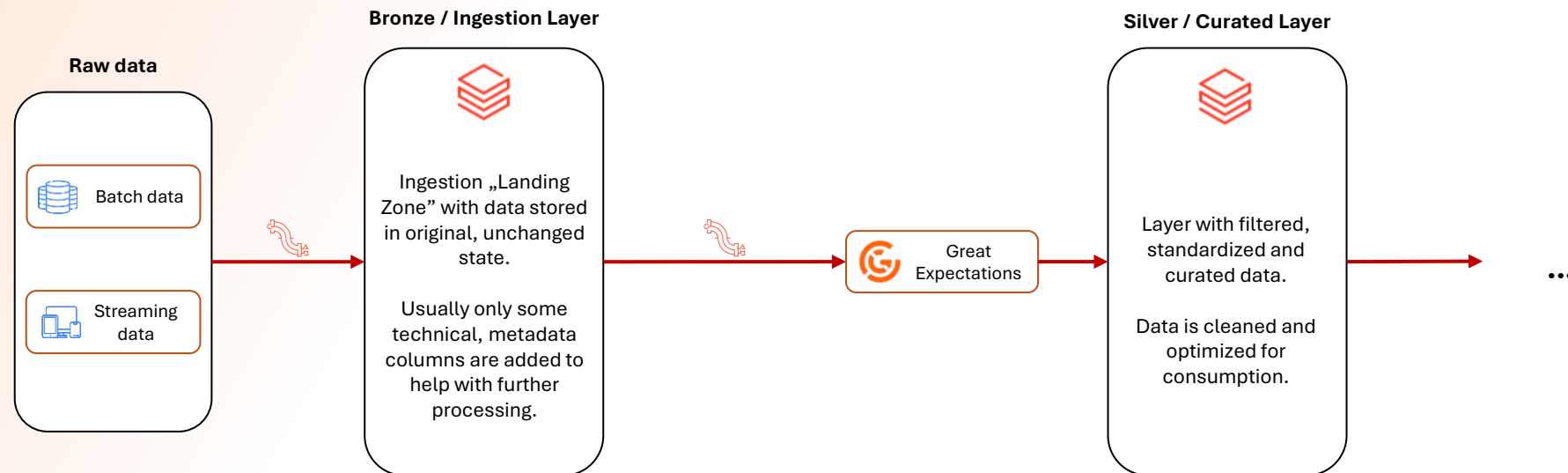
Introduction





**From raw to
model-ready data**

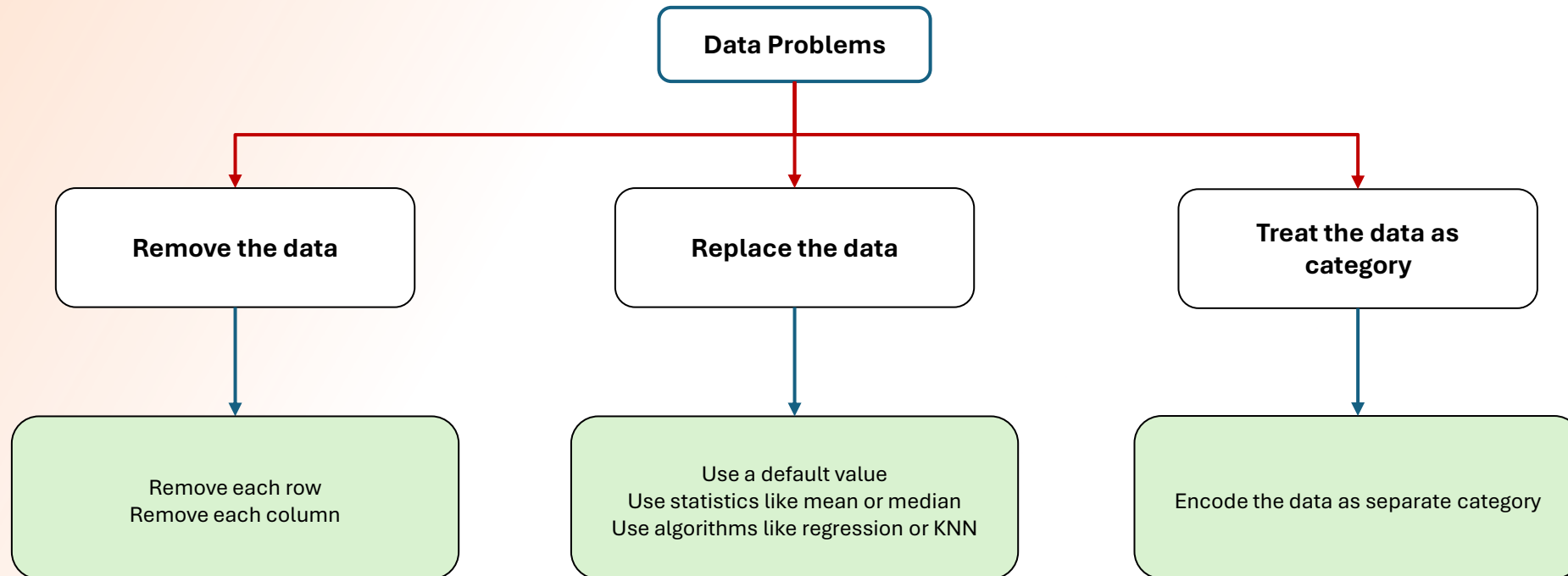
Data pipeline



Most common data problems

1	Missing values	Data points are absent due to collection issues or system gaps
2	Inconsistent data types	Numeric fields stored as strings or mixed formats breaking transformations
3	Duplicates	Repeated records inflate counts and distort model training
4	Outliers or impossible values	Extreme values that skew distributions and can mislead model behavior
5	Inconsistent categorical labels	Variations like "USA", "Usa", "United States" cause incorrect grouping

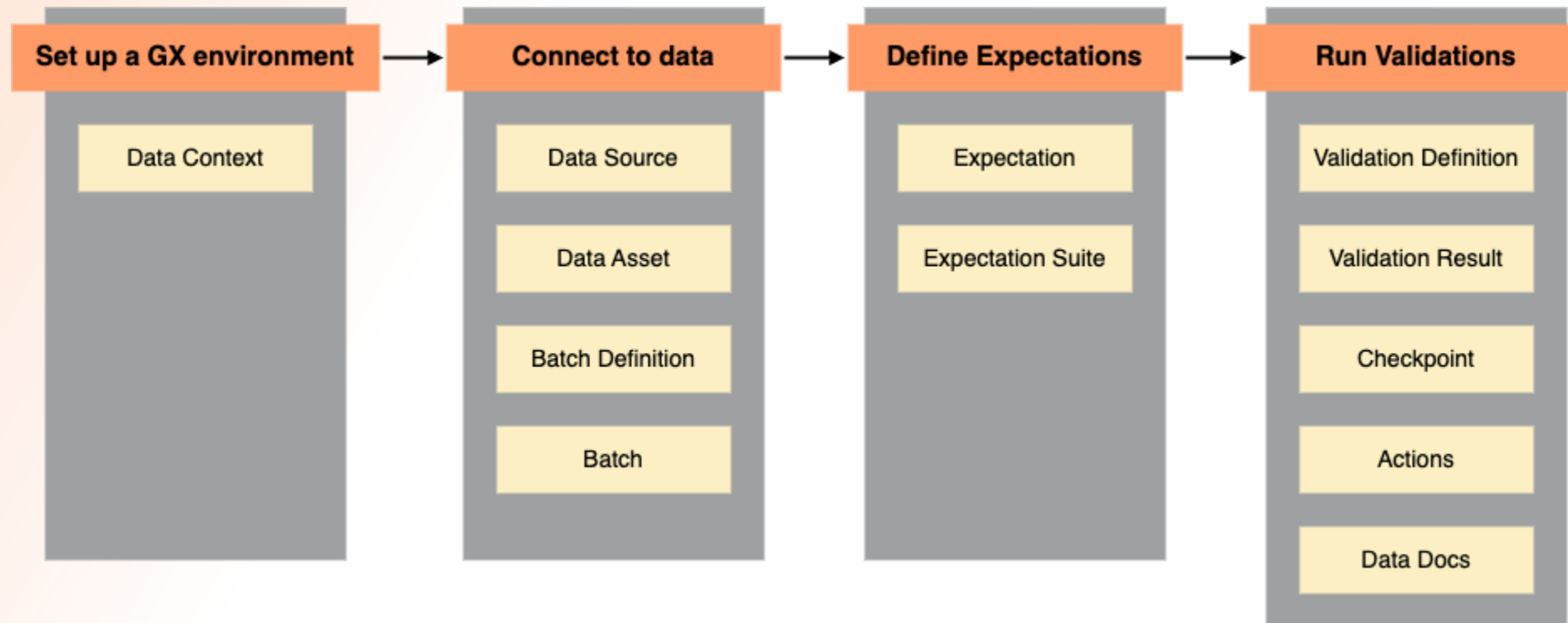
Handling data quality



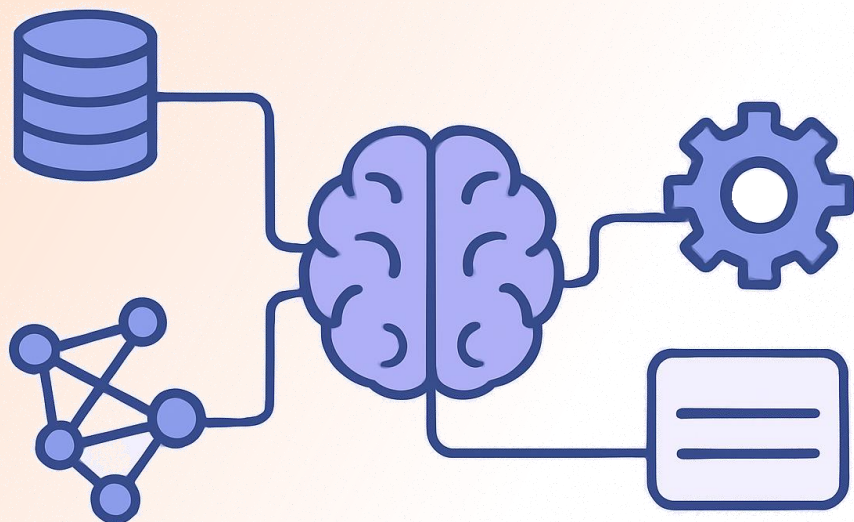
Rules for handling data quality:

- There is no „best” way to do it
- Understand the reasons
- Consider and validate the impact and potential bias
- Assess the amount and pattern of missing data

Great Expectations



DEMO



Summary

Takeaways

- **There is no universal cleaning recipe:** every dataset has its own problems, so methods must be chosen based on context, not habit.
- **You can't clean what you don't understand:** meaningful preprocessing starts with understanding semantics, business logic, and how data is generated.
- **Validate early, validate always:** automated checks catch data problems before they hit training or production.
- **Consistency beats cleverness:** stable schemas, clear types, and predictable formats are more beneficial than the most complicated feature engineering.

Thank you!

Contact:



<https://www.linkedin.com/in/maciej-kepa>



<https://github.com/maciejkepa/ai-ml-in-practice>

