



MLops: get your ML apps ready for production

Francesco Tamborra

Outline

- ◆ Intro: What & why

```
requirements.txt
```

- ◆ ML development VS standard development
- ◆ ML lifecycle

- ◆ **MLOps maturity levels**

- ◆ Tools overview

- ◆ `@app.get("/predictions")`

What is MLOps ?

"The ability to apply DevOps principles to Machine Learning applications"

What is MLOps ?

"The ability to apply DevOps principles to Machine Learning applications"

MLOps is an ML engineering **culture and practice** that aims at unifying ML system development (Dev) and ML system operations (Ops).

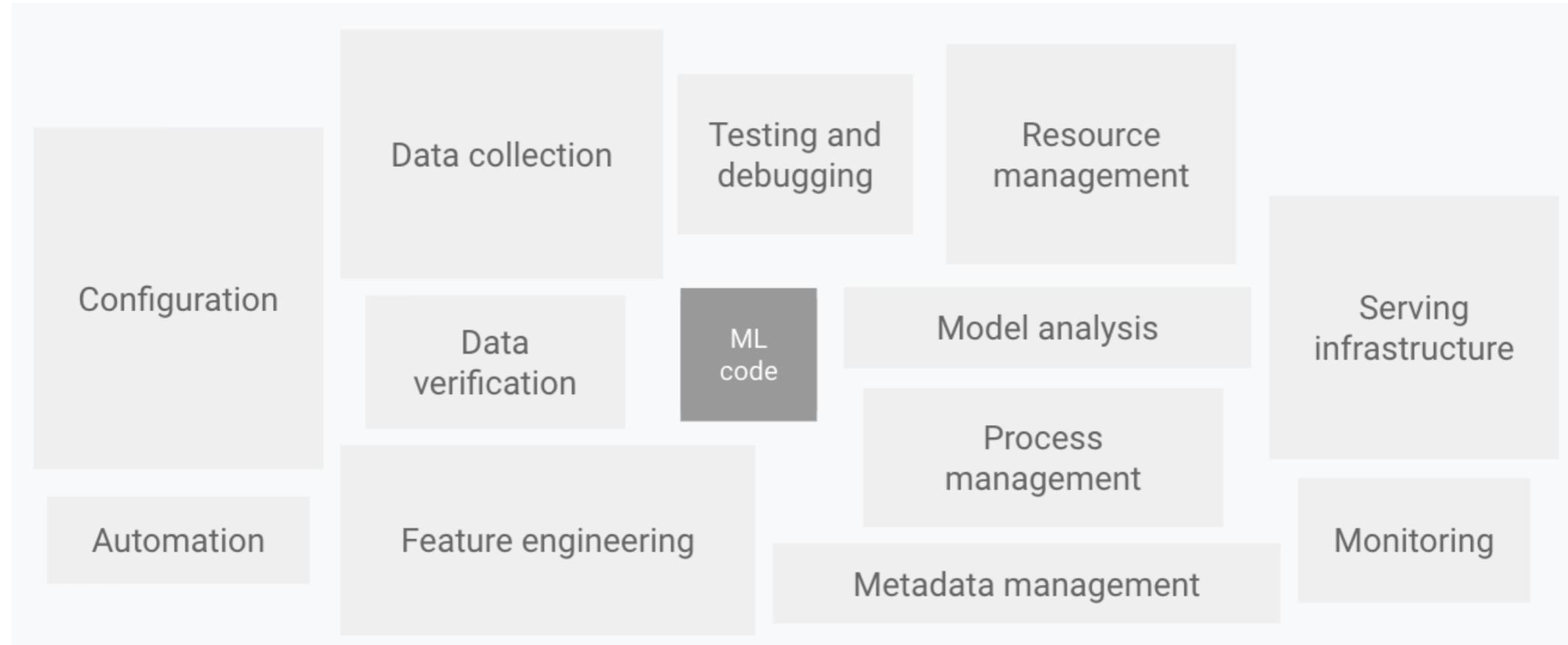
Practicing MLOps means that you advocate for automation and monitoring at all steps of ML system construction, including integration, testing, releasing, deployment and infrastructure management.

What is MLOps ?

"The ability to apply DevOps principles to Machine Learning applications"

MLOps is an ML engineering **culture and practice** that aims at unifying ML system development (Dev) and ML system operations (Ops).

Practicing MLOps means that you advocate for automation and monitoring at all steps of ML system construction, including integration, testing, releasing, deployment and infrastructure management.

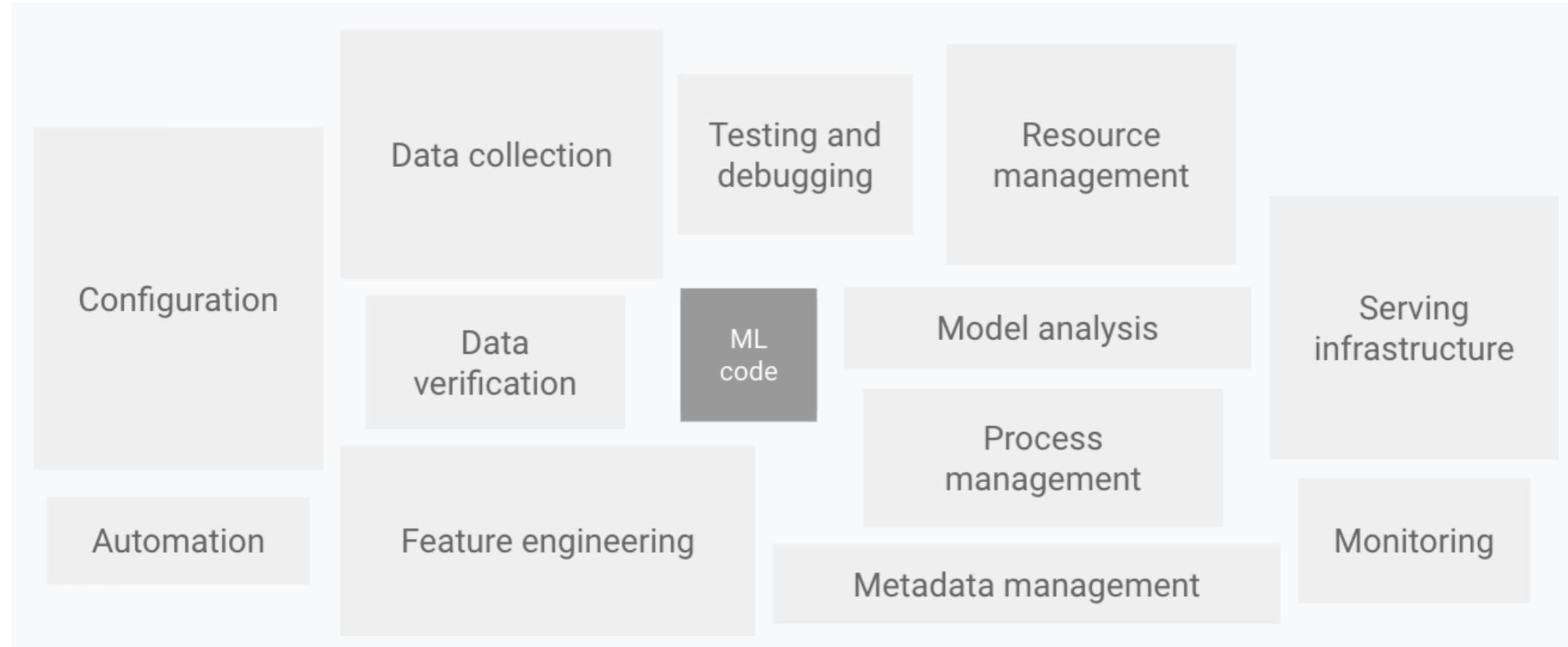


What is MLOps ?

"The ability to apply DevOps principles to Machine Learning applications"

MLOps is an ML engineering **culture and practice** that aims at unifying ML system development (Dev) and ML system operations (Ops).

Practicing MLOps means that you advocate for automation and monitoring at all steps of ML system construction, including integration, testing, releasing, deployment and infrastructure management.



Global MLOps Market Size estimated at **612 USD million** in **2021** and projected to reach **6.1 USD billion** by **2028**

Why MLOps ?

"2020 State of Enterprise ML" by Algorithmia (survey size: 750) report results:

Why MLOps ?

"2020 State of Enterprise ML" by Algorithmia (survey size: 750) report results:

55% of companies surveyed have **not** deployed
a machine learning model

Why MLOps ?

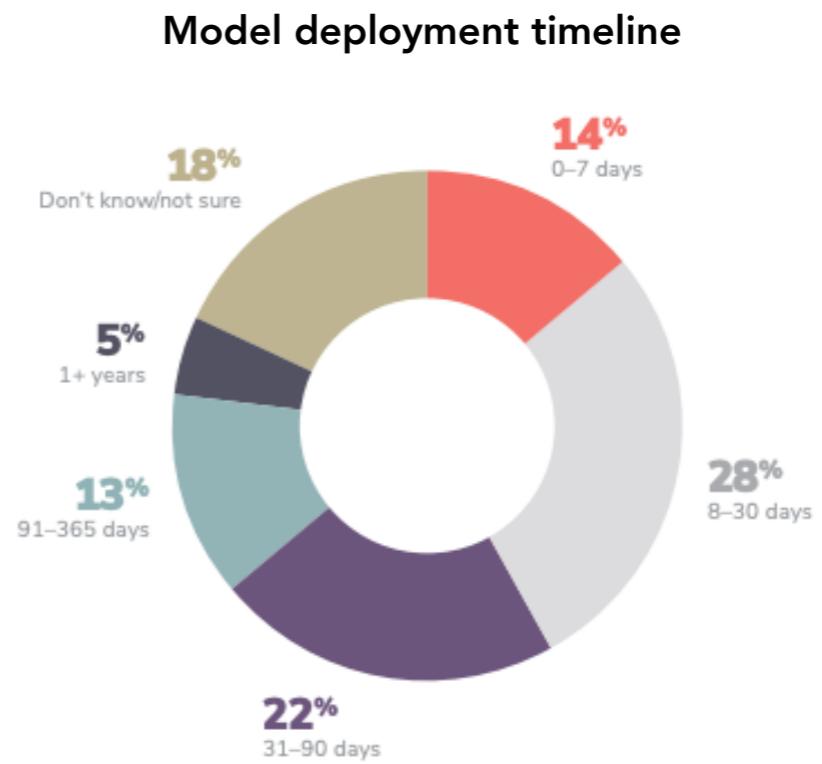
"2020 State of Enterprise ML" by Algorithmia (survey size: 750) report results:

55% of companies surveyed have **not** deployed
a machine learning model

Why MLOps ?

"2020 State of Enterprise ML" by Algorithmia (survey size: 750) report results:

55% of companies surveyed have **not** deployed a machine learning model

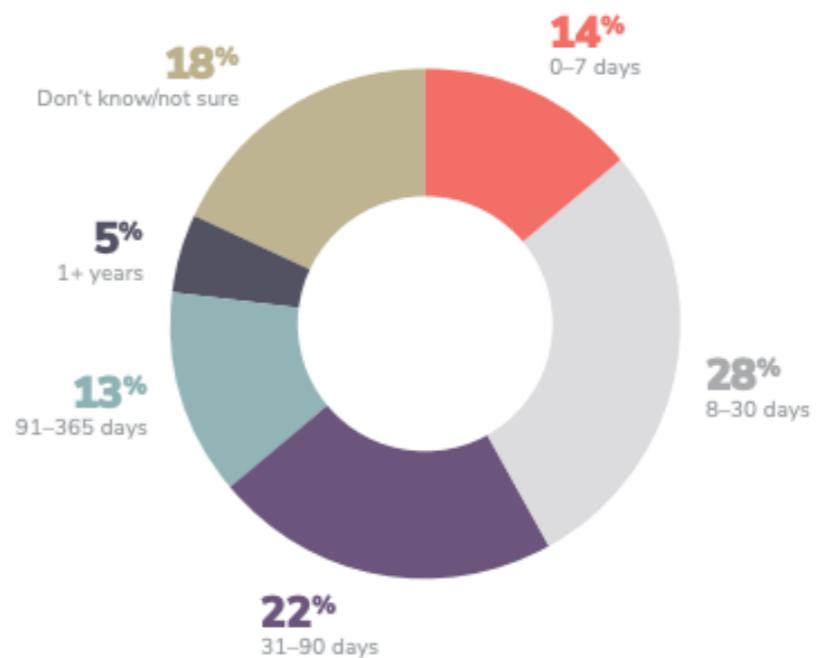


Why MLOps ?

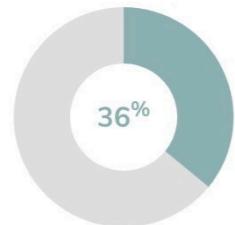
"2020 State of Enterprise ML" by Algorithmia (survey size: 750) report results:

55% of companies surveyed have **not** deployed a machine learning model

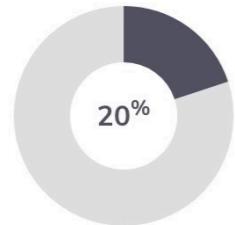
Model deployment timeline



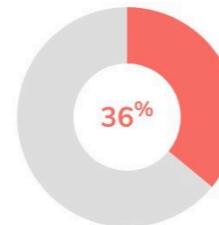
What percentage of your data scientists' time is spent deploying ML models?



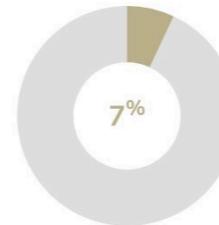
36% of survey participants said their data scientists spend **a quarter** of their time deploying ML models



20% of survey participants said their data scientists spend **half to three-quarters** of their time deploying ML models



36% of survey participants said their data scientists spend **a quarter to half** of their time deploying ML models



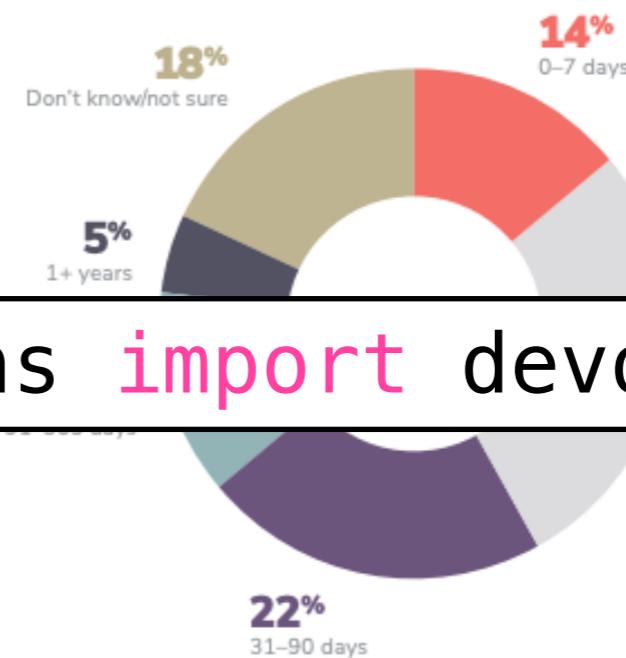
7% of survey participants said their data scientists spend **more than three-quarters** of their time deploying ML models

Why MLOps ?

"2020 State of Enterprise ML" by Algorithmia (survey size: 750) report results:

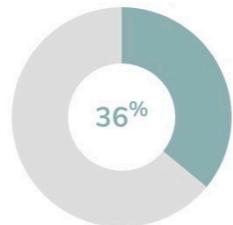
55% of companies surveyed have **not** deployed a machine learning model

Model deployment timeline

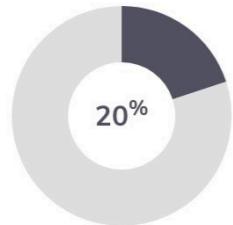


from Development Operations import devops

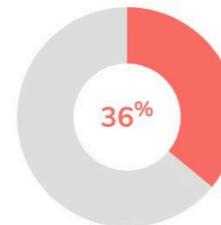
What percentage of your data scientists' time is spent deploying ML models?



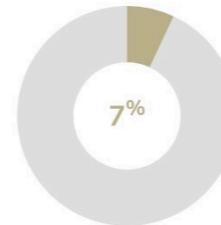
36% of survey participants said their data scientists spend **a quarter** of their time deploying ML models



20% of survey participants said their data scientists spend **half to three-quarters** of their time deploying ML models



36% of survey participants said their data scientists spend **a quarter to half** of their time deploying ML models

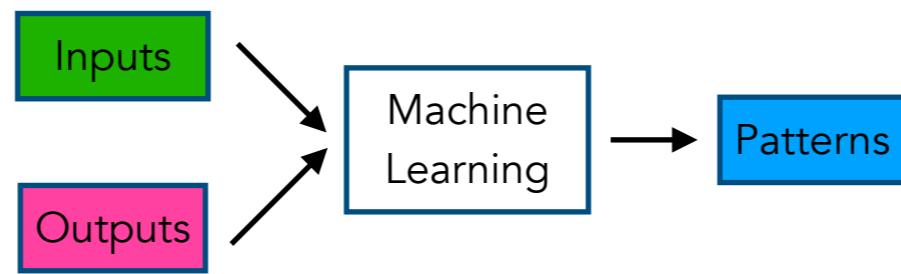


7% of survey participants said their data scientists spend **more than three-quarters** of their time deploying ML models

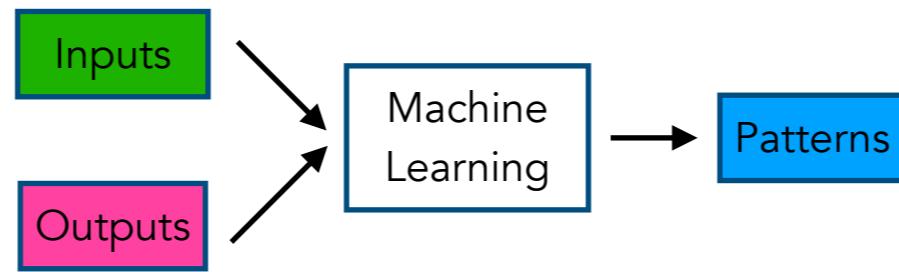
ML Dev VS standard SW Dev



ML Dev VS standard SW Dev



ML Dev VS standard SW Dev



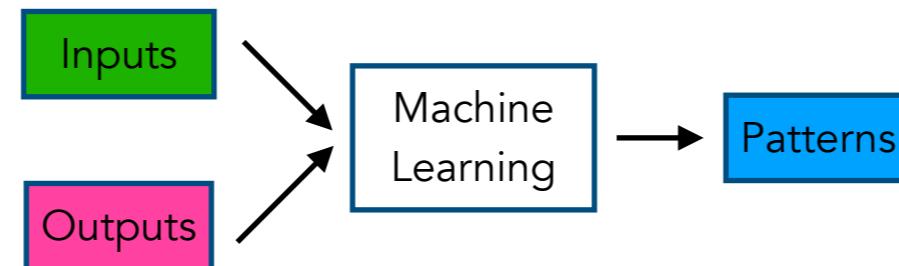
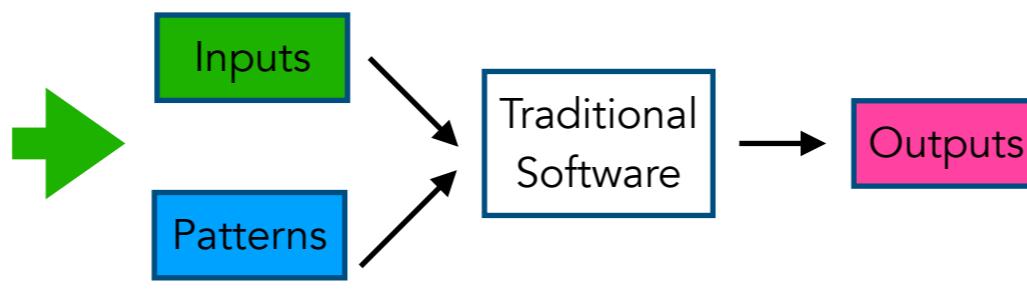
Andrej Karpathy

Nov 11, 2017 · 9 min read · Listen

Software 2.0

ML Dev VS standard SW Dev

- (Static data)
- Code (*Testing & Versioning*)



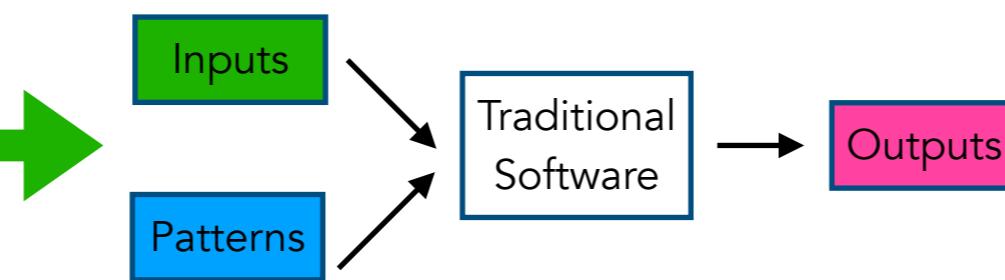
Andrej Karpathy

Nov 11, 2017 · 9 min read · Listen

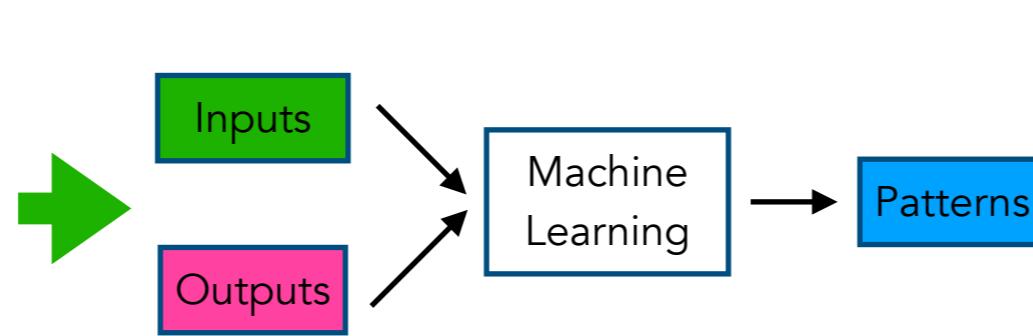
Software 2.0

ML Dev VS standard SW Dev

- (Static data)
- Code (*Testing & Versioning*)



- **Big** dynamical **data**
- Code
(*Testing & Versioning both*)

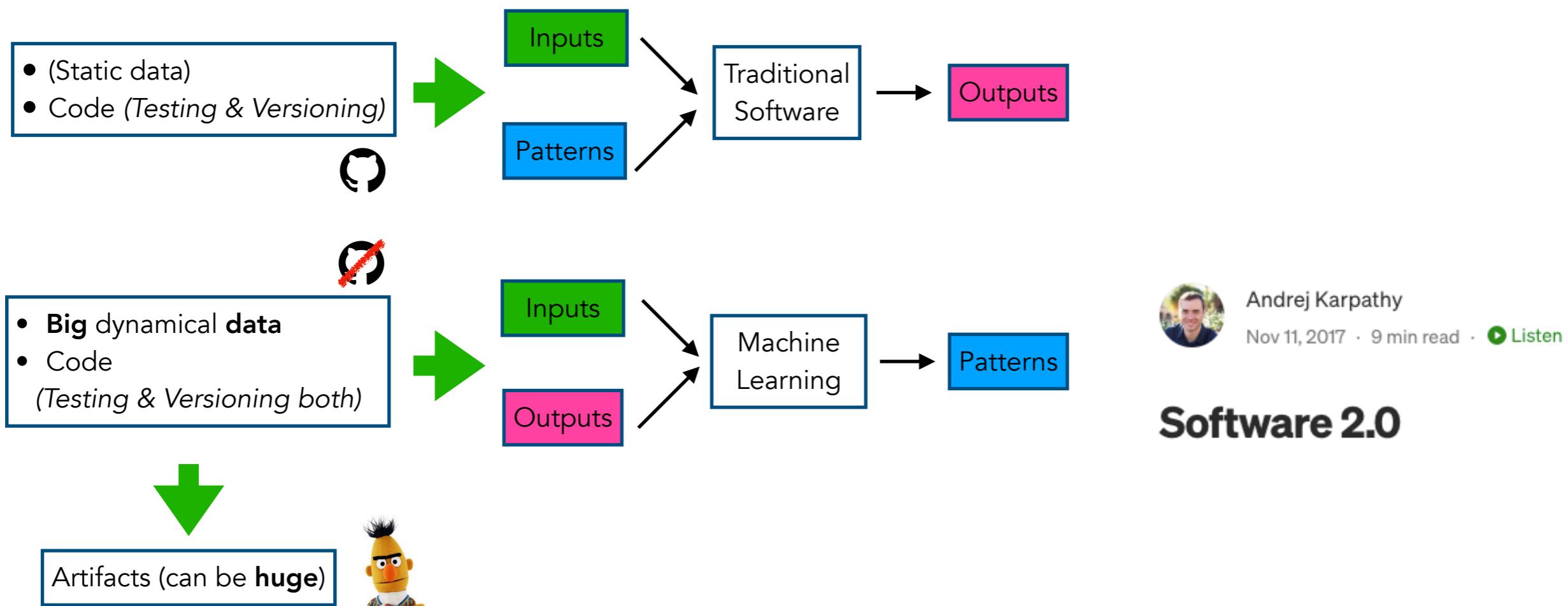


Andrej Karpathy

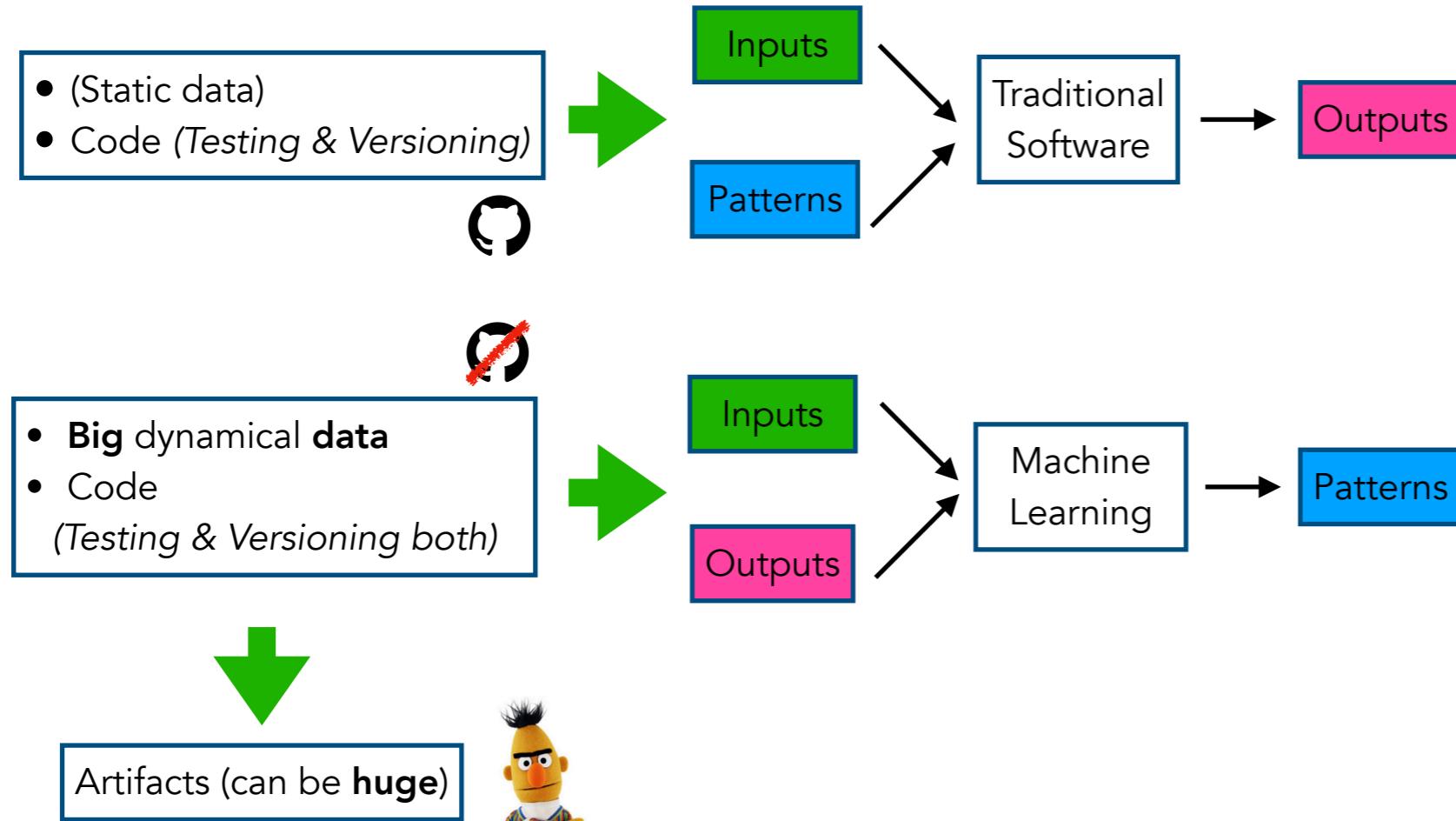
Nov 11, 2017 · 9 min read · Listen

Software 2.0

ML Dev VS standard SW Dev



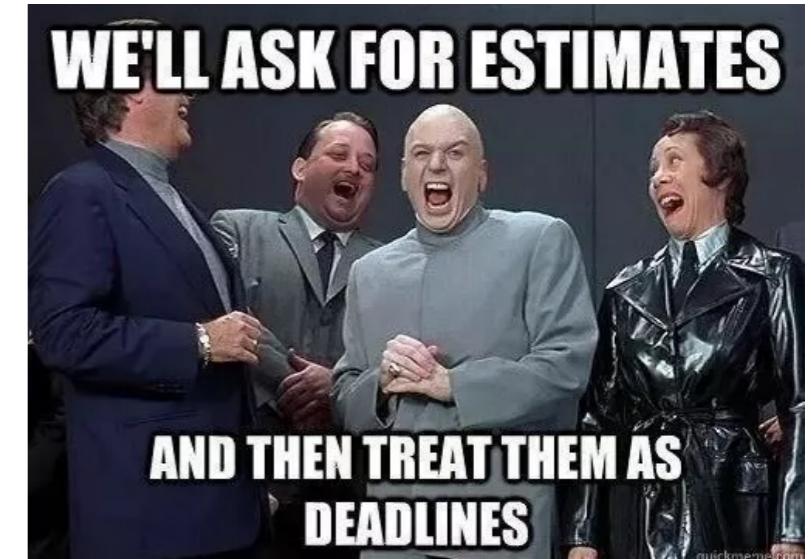
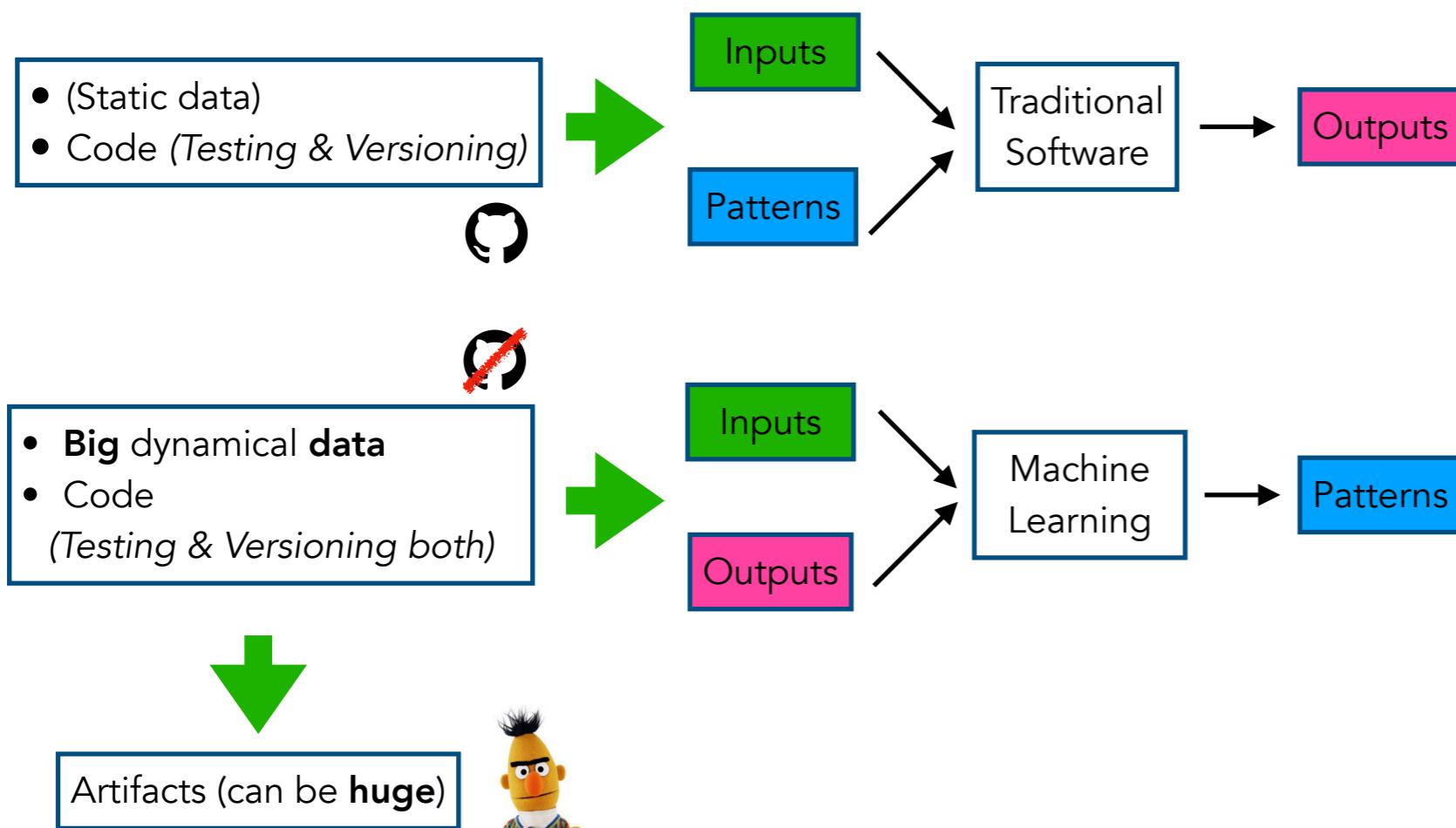
ML Dev VS standard SW Dev



Other differences:

Development Experimental in nature: **tracking** what worked and what didn't is essential to maintain **reproducibility**

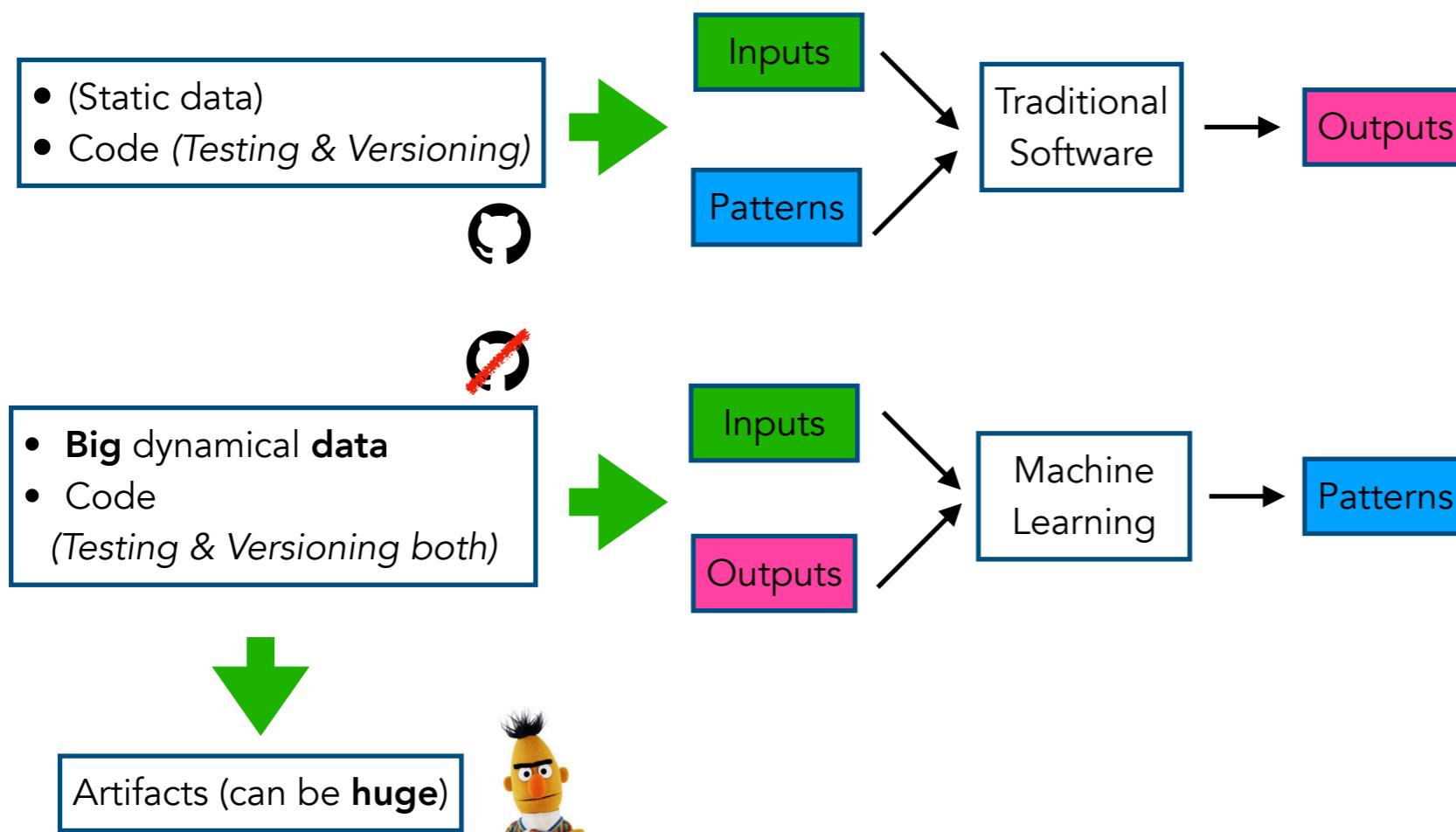
ML Dev VS standard SW Dev



Other differences:

Development Experimental in nature: **tracking** what worked and what didn't is essential to maintain **reproducibility**

ML Dev VS standard SW Dev

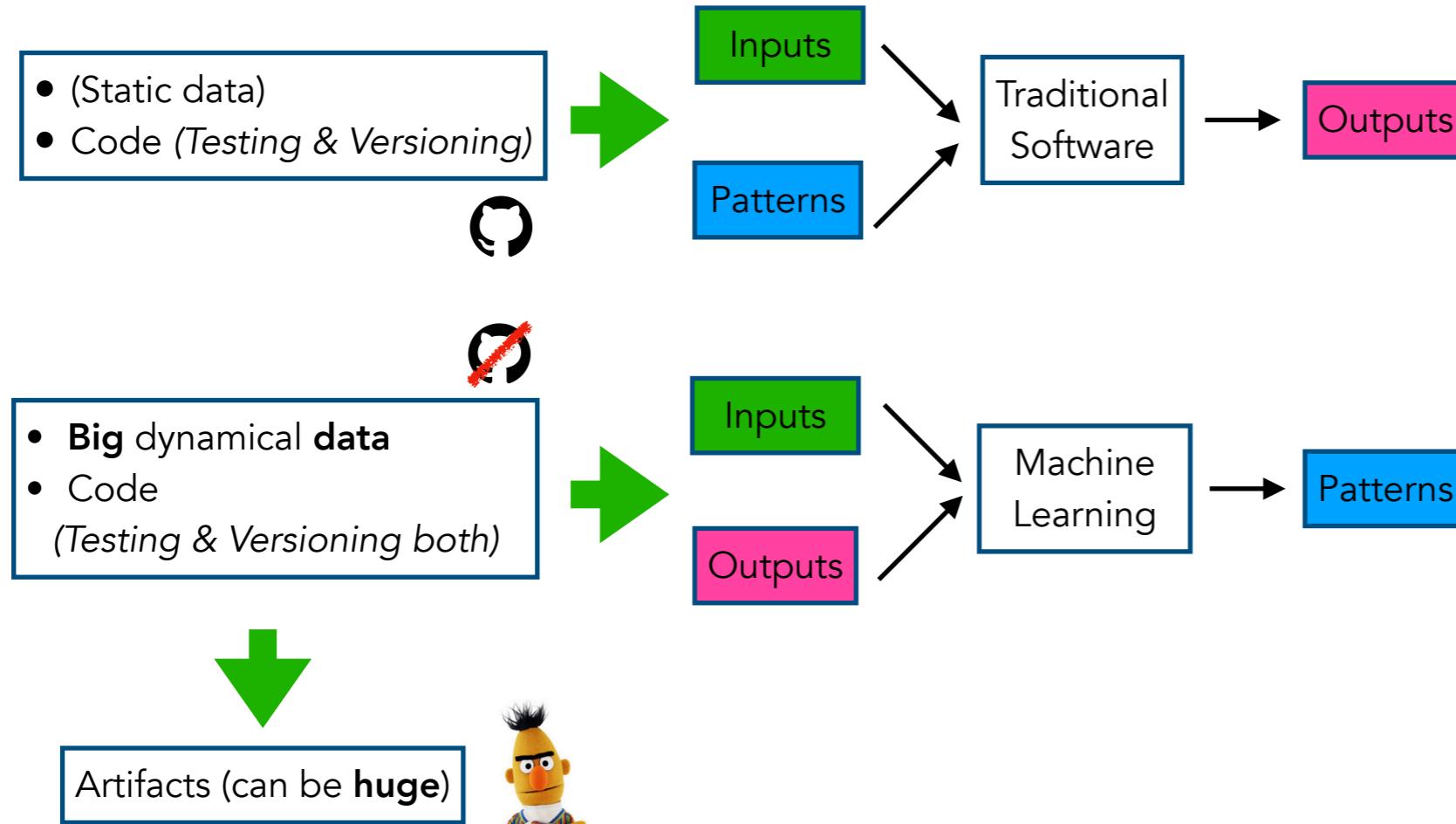


Other differences:

Development Experimental in nature: **tracking** what worked and what didn't is essential to maintain **reproducibility**

Testing In addition to typical unit and integration tests, you need **data validation**, **trained model quality evaluation**, and **model validation**

ML Dev VS standard SW Dev



Other differences:

Development	Experimental in nature: tracking what worked and what didn't is essential to maintain reproducibility
Testing	In addition to typical unit and integration tests, you need data validation , trained model quality evaluation , and model validation
Production	Models can decay in more ways than conventional software systems, and you need to consider this degradation by tracking summary statistics of your data and monitor the online performance

DevOps VS MLOps in a nutshell

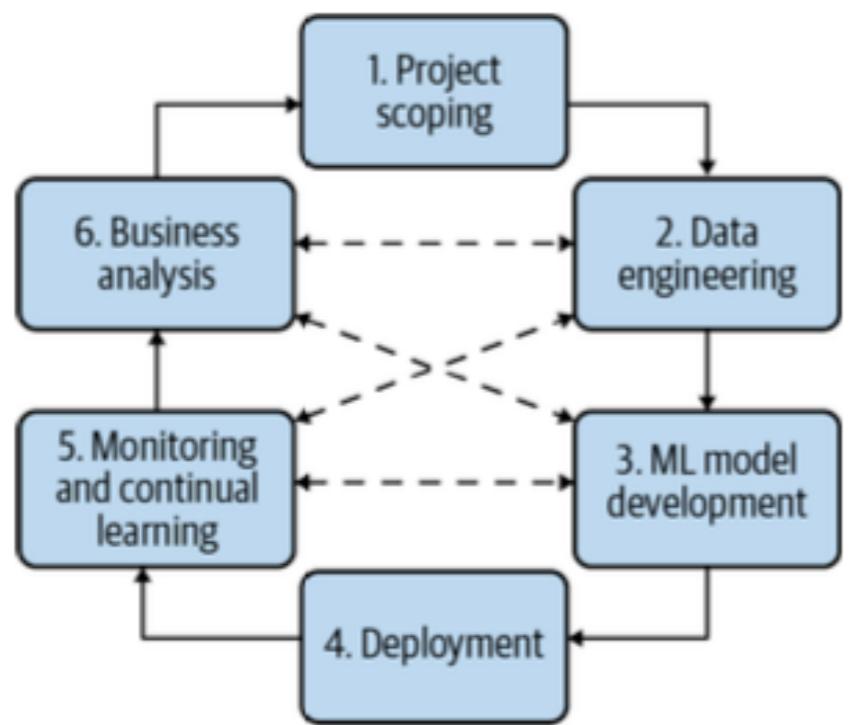
DevOps VS MLOps in a nutshell

CI is no longer only about testing and validating code and components, but also testing and validating data, data schemas, and models.

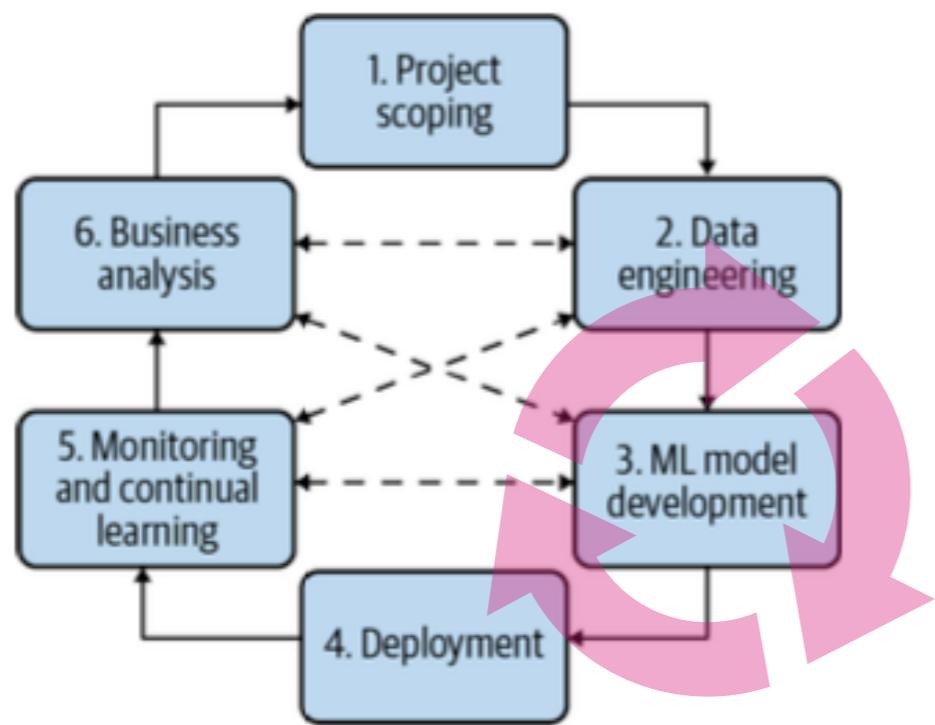
CD is no longer about a single software package or a service, but a system (an ML training pipeline) that should automatically deploy another service (model prediction service).

CT is a new property, unique to ML systems, that's concerned with automatically retraining and serving the models.

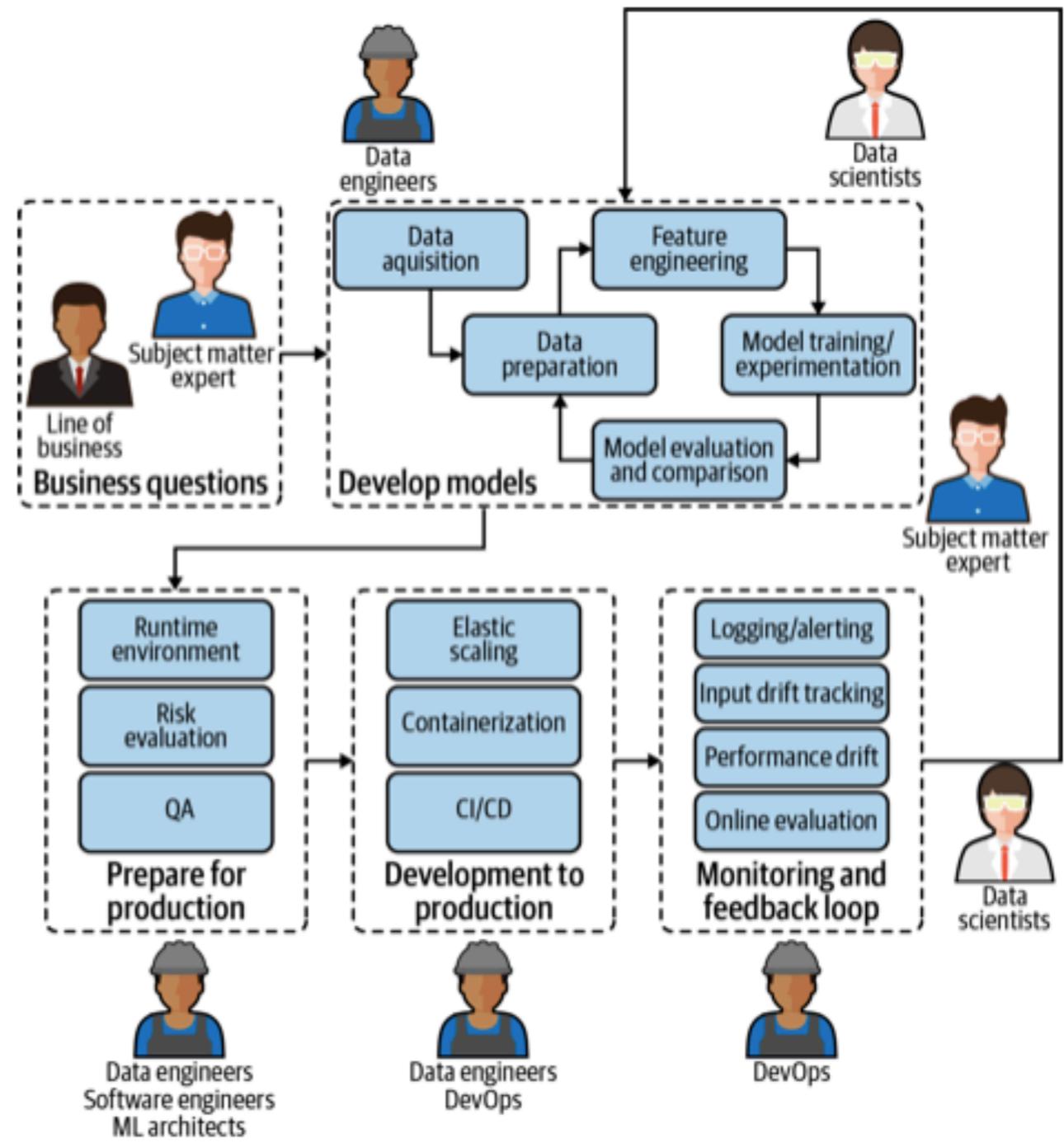
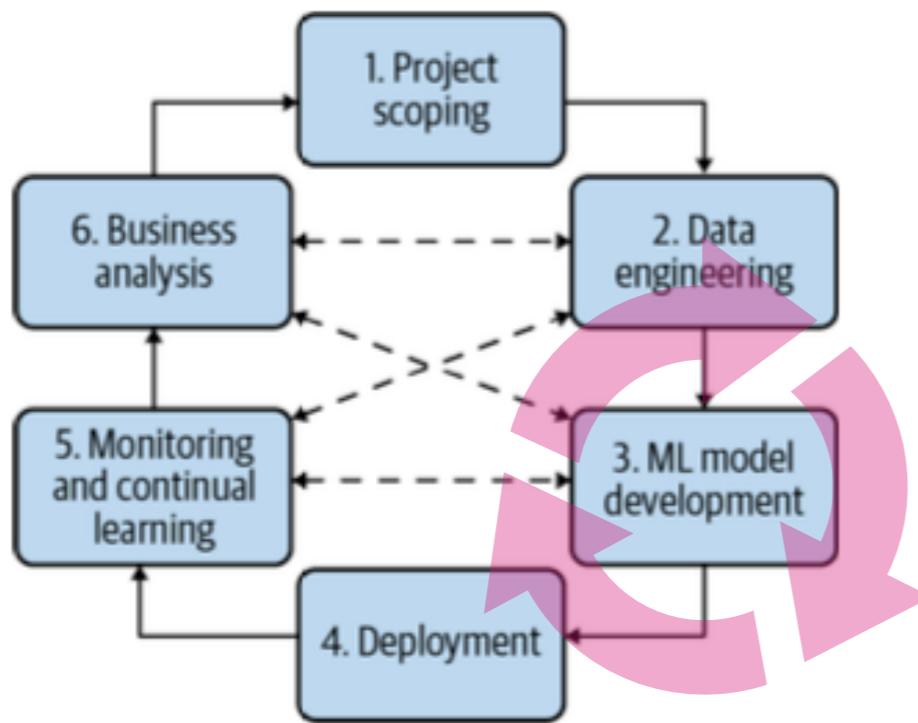
ML Lifecycle



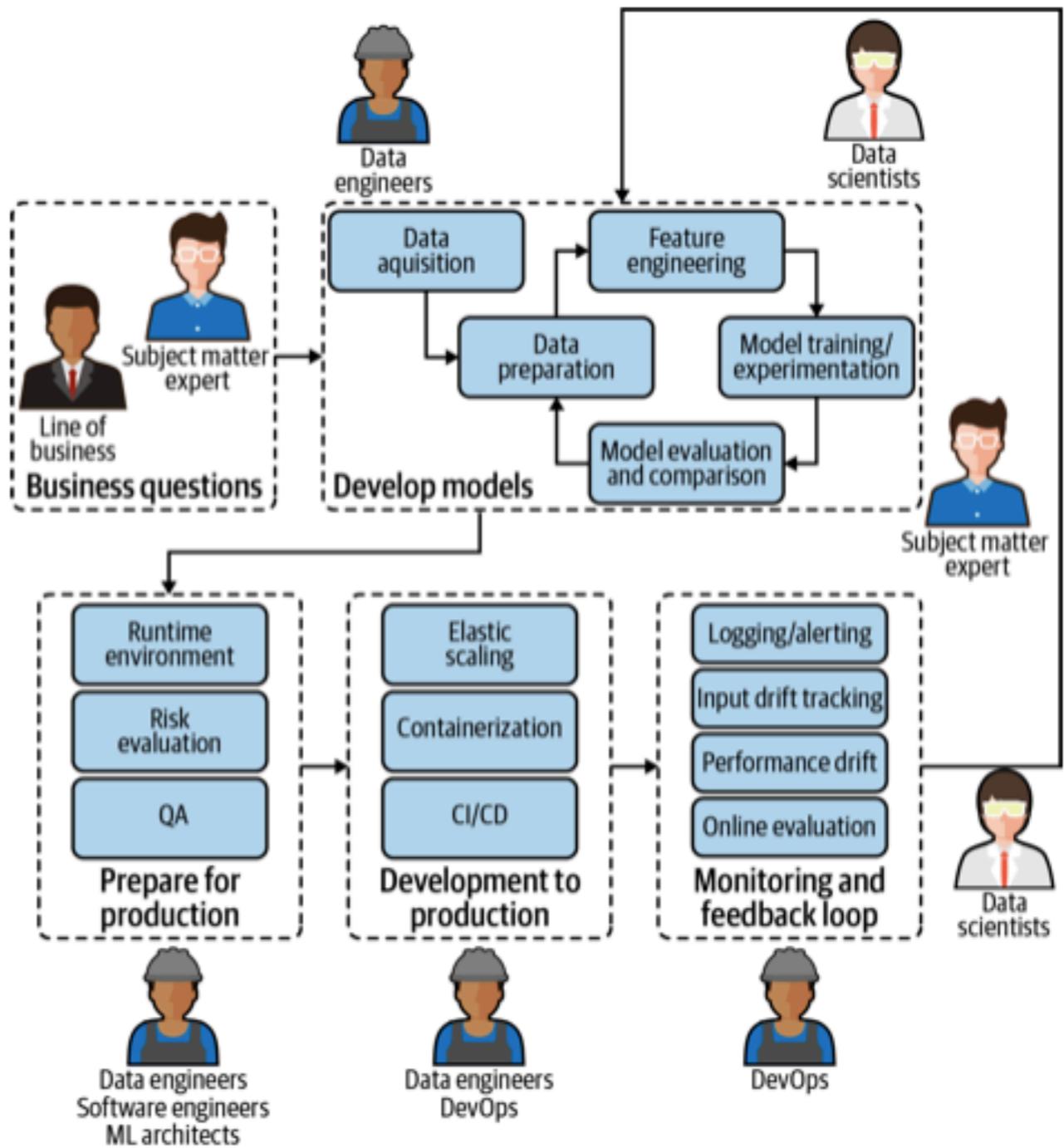
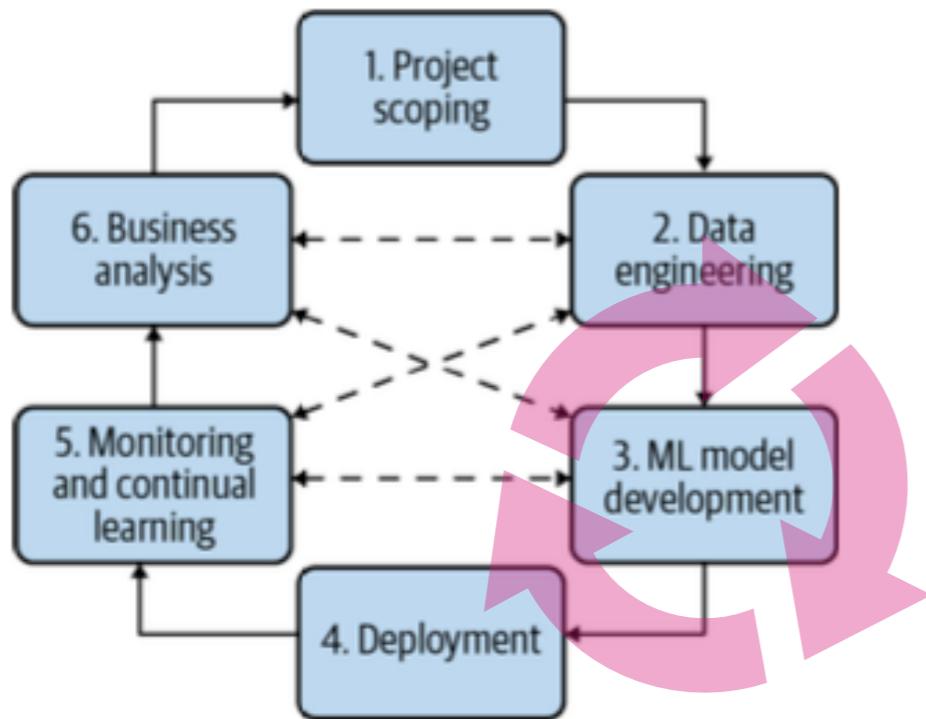
ML Lifecycle



ML Lifecycle

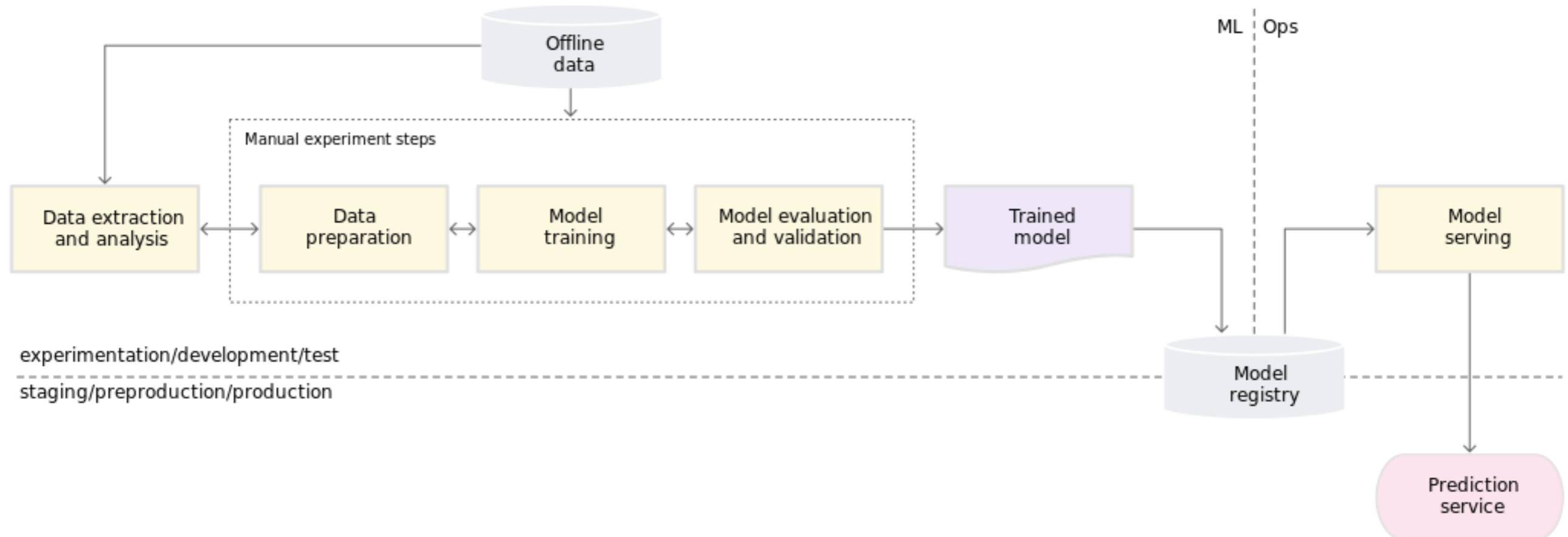


ML Lifecycle

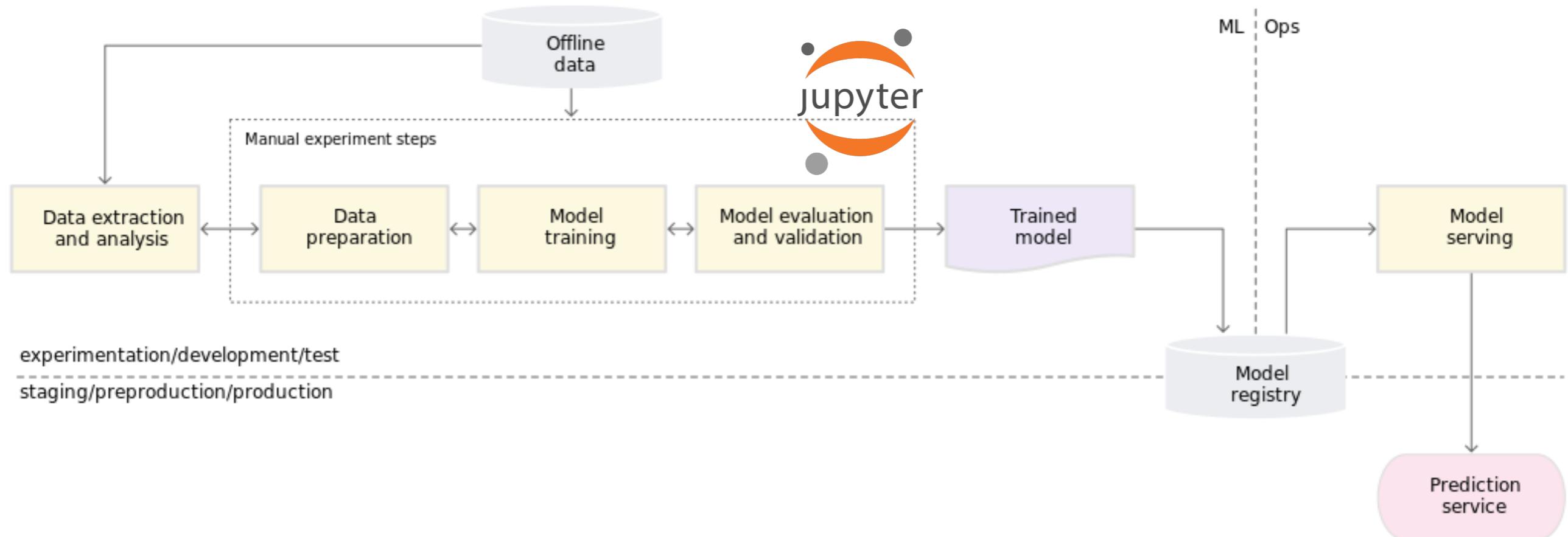


The **level of automation** of these steps defines the **maturity** of the ML process, which reflects the **velocity** of training new models given new data or new implementations

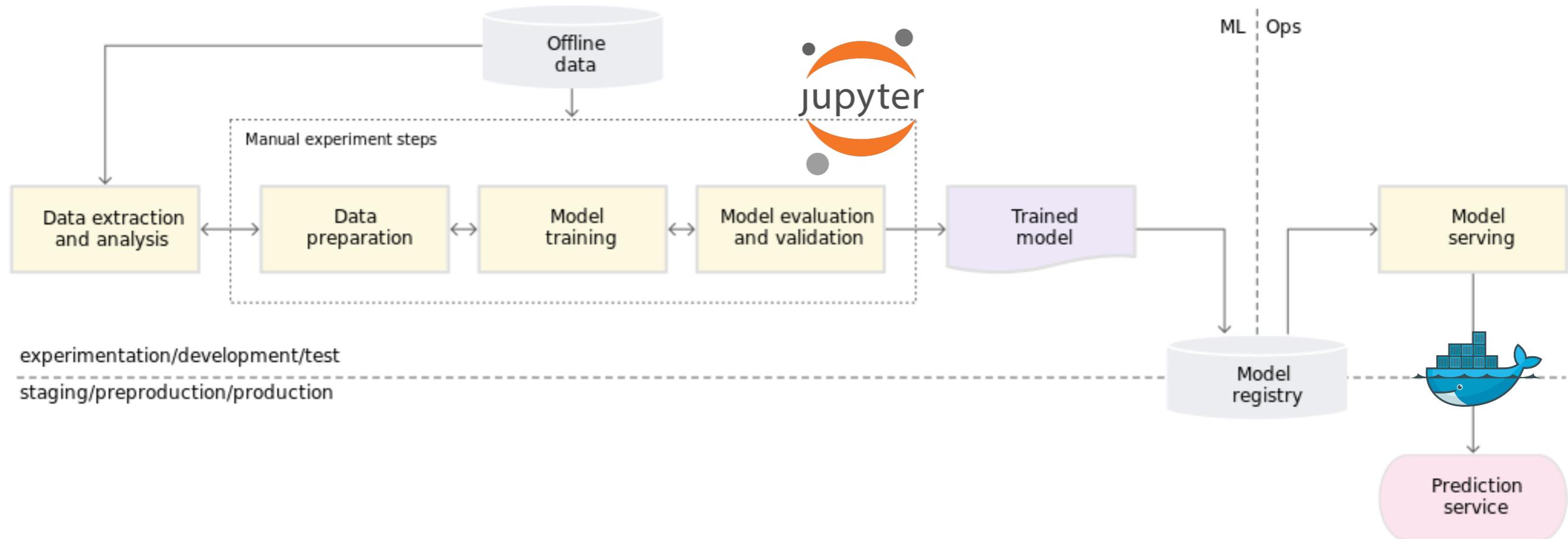
MLOps level 0



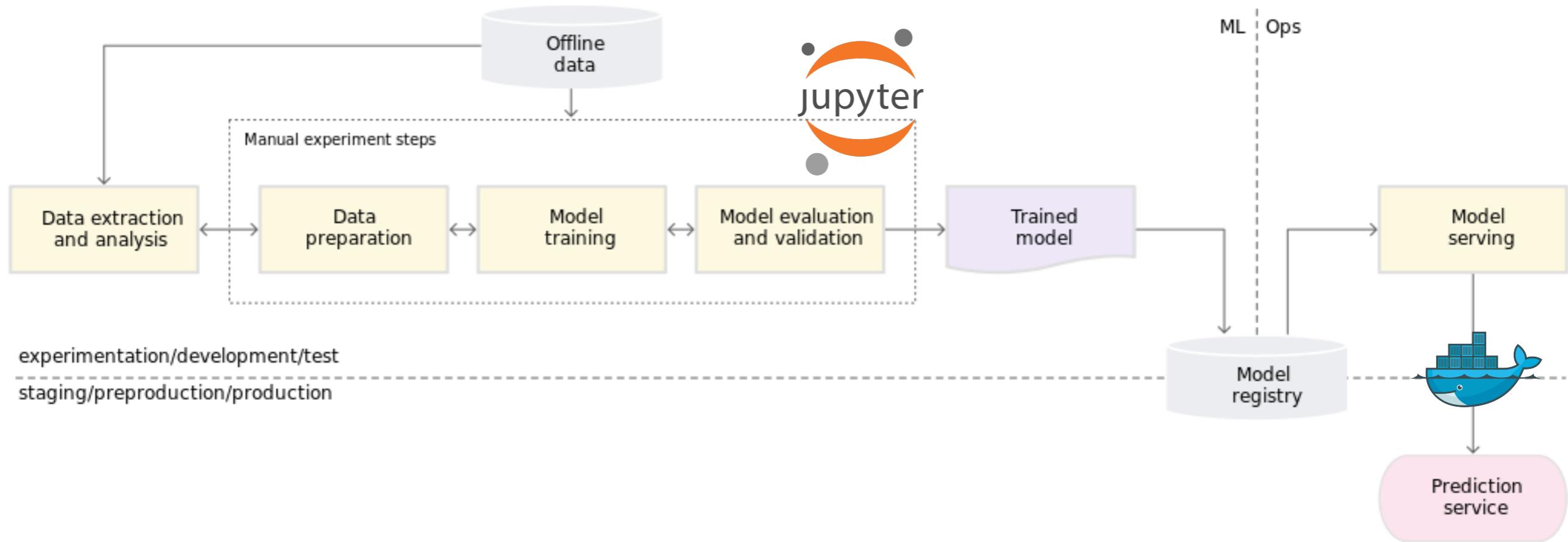
MLOps level 0



MLOps level 0

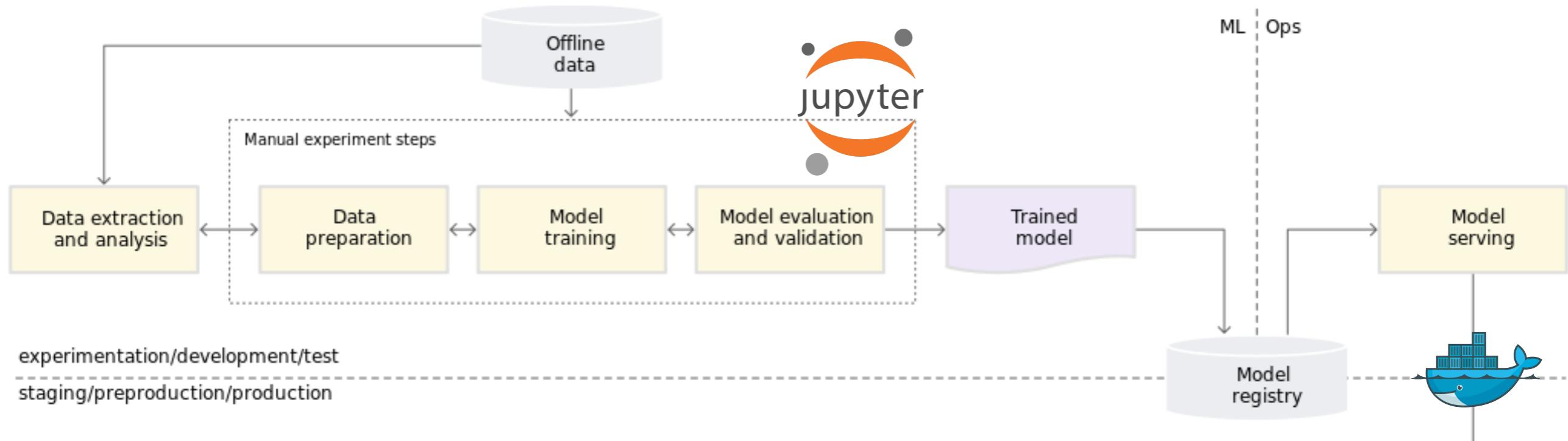


MLOps level 0



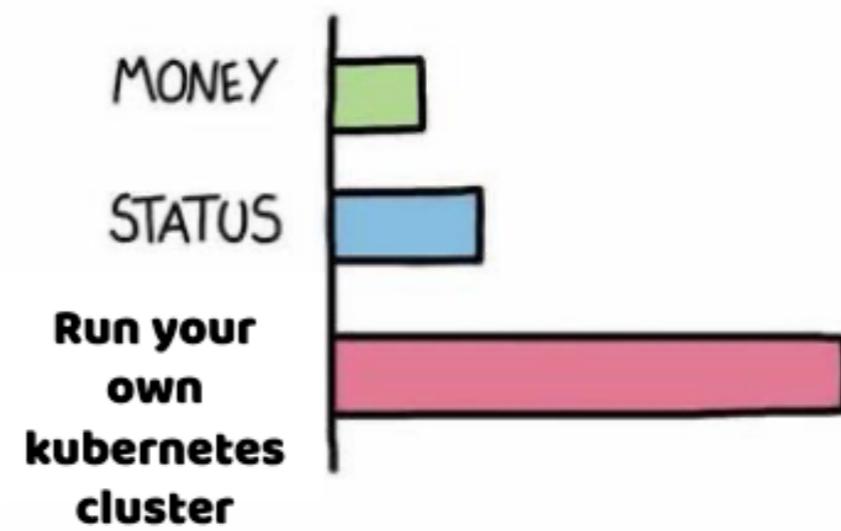
	Dev	Prod
Scale	<ul style="list-style-type: none"> • usually one instance • no need to worry about autoscaling 	<ul style="list-style-type: none"> • multiple instances/nodes • you most certainly want autoscaling
State	stateful: reproducible but inflexible <ul style="list-style-type: none"> • install dependencies once and forget • can persist data in dedicated storage 	inherently stateless: flexible but hard to reproduce <ul style="list-style-type: none"> • need to install dependencies on any new instance • need to figure out how to persist data/state across instances

MLOps level 0

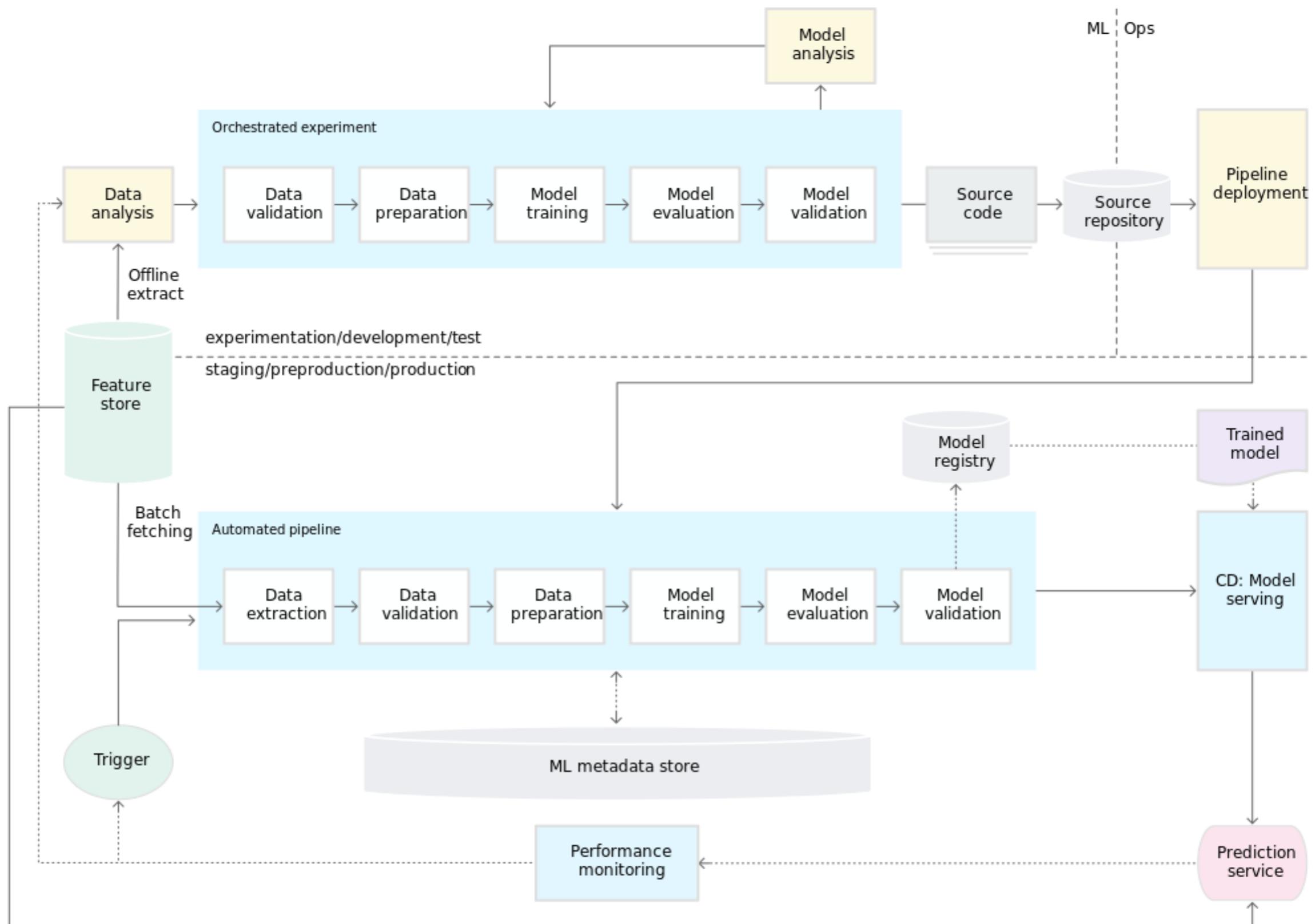


	Dev	
Scale	<ul style="list-style-type: none"> • usually one instance • no need to worry about autoscaling 	<ul style="list-style-type: none"> • multiple instances • you most likely have to pay for them
State	<p>stateful: reproducible but inflexible</p> <ul style="list-style-type: none"> • install dependencies once and forget • can persist data in dedicated storage 	<p>inherently stateless</p> <ul style="list-style-type: none"> • need to install dependencies • need to find instances

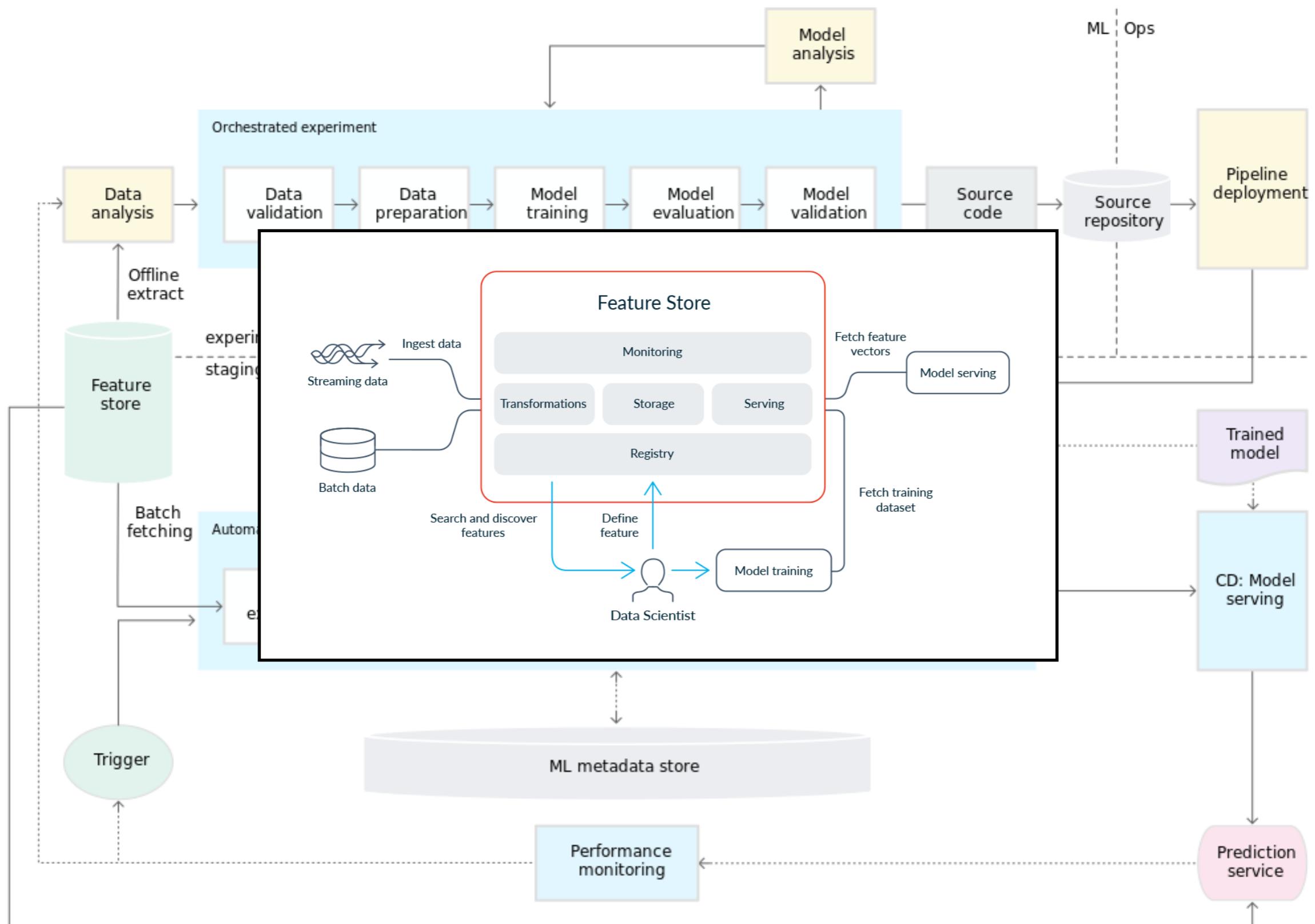
WHAT GIVES PEOPLE FEELINGS OF POWER



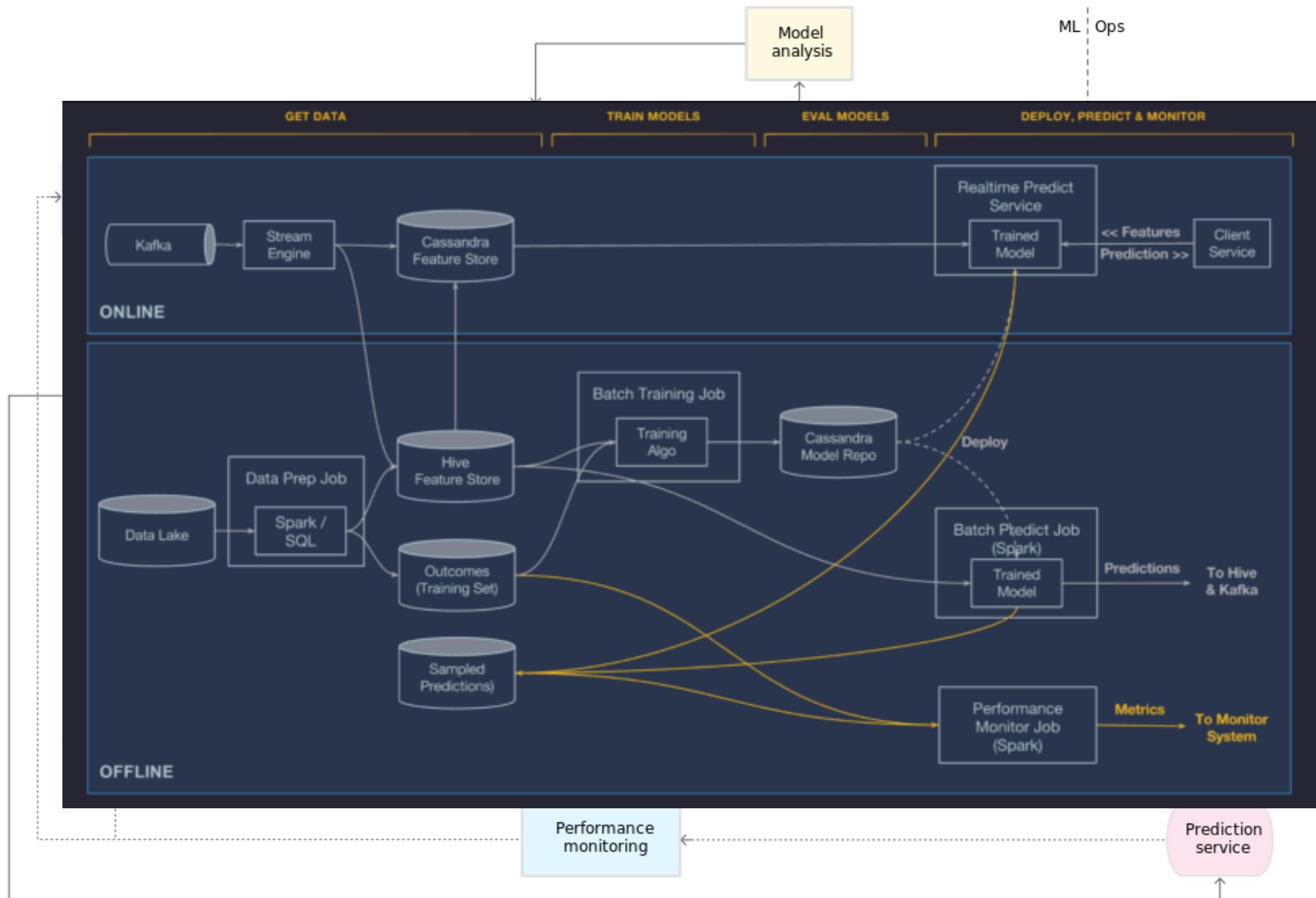
MLOps level 1



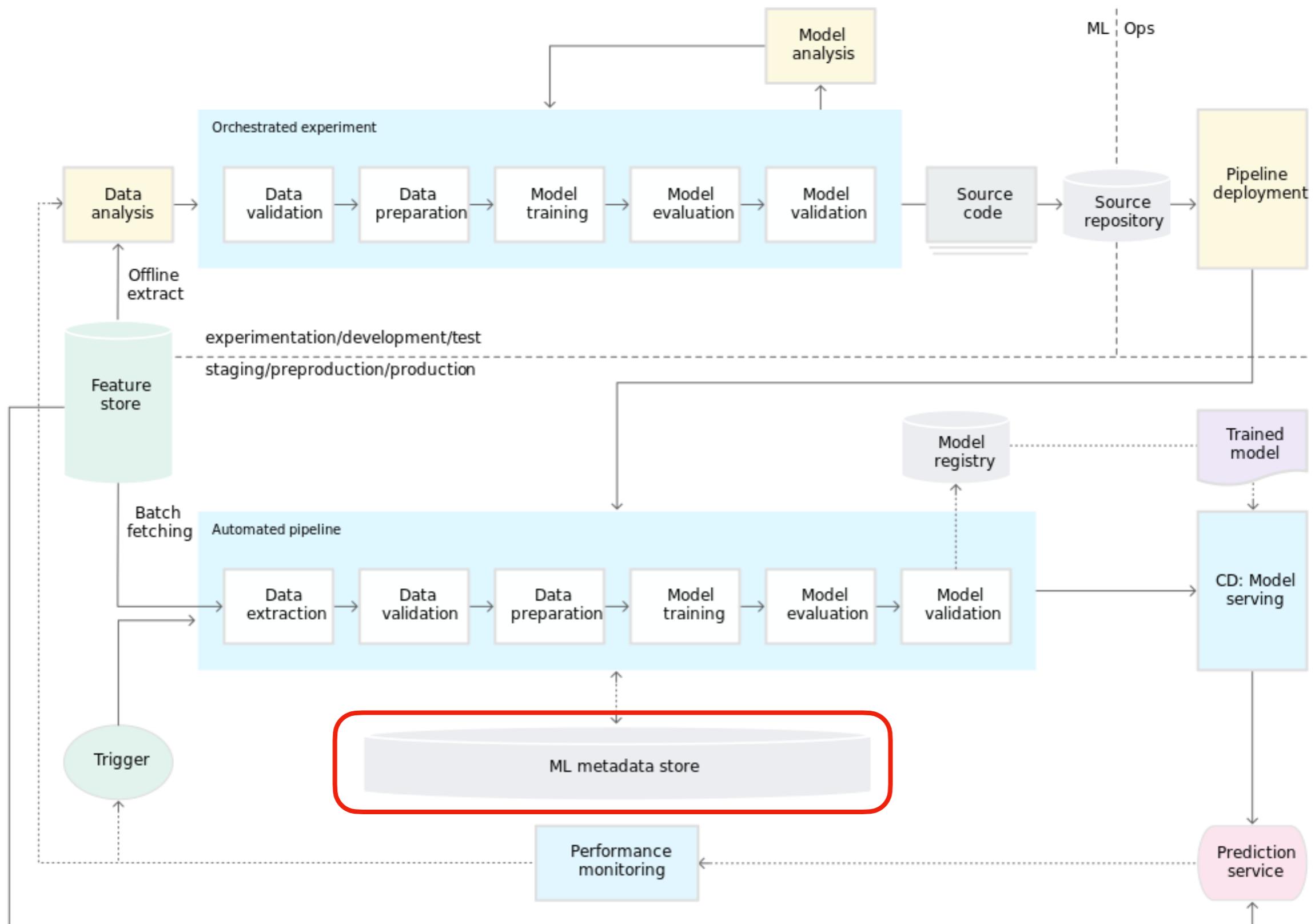
MLOps level 1



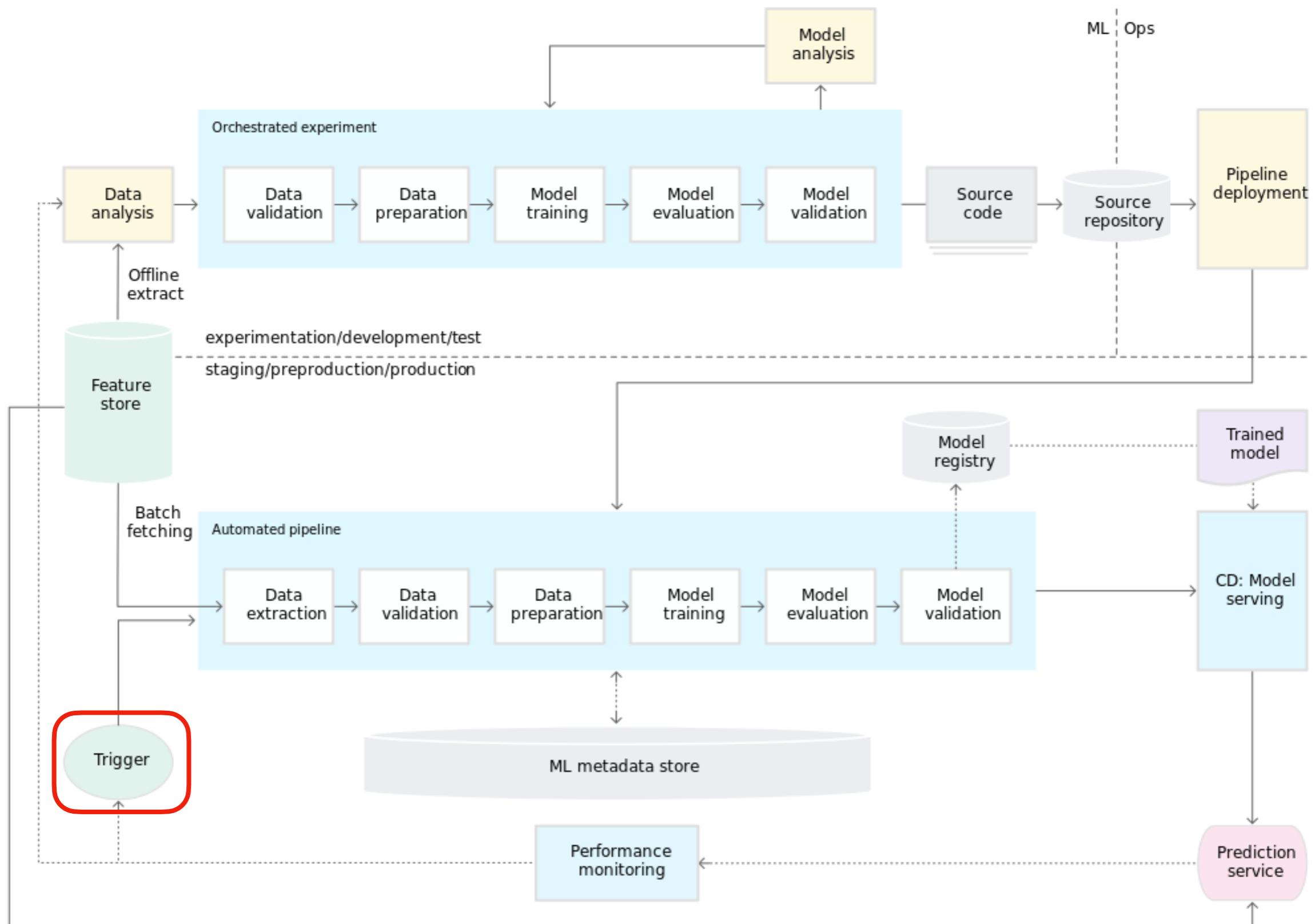
MLOps level 1



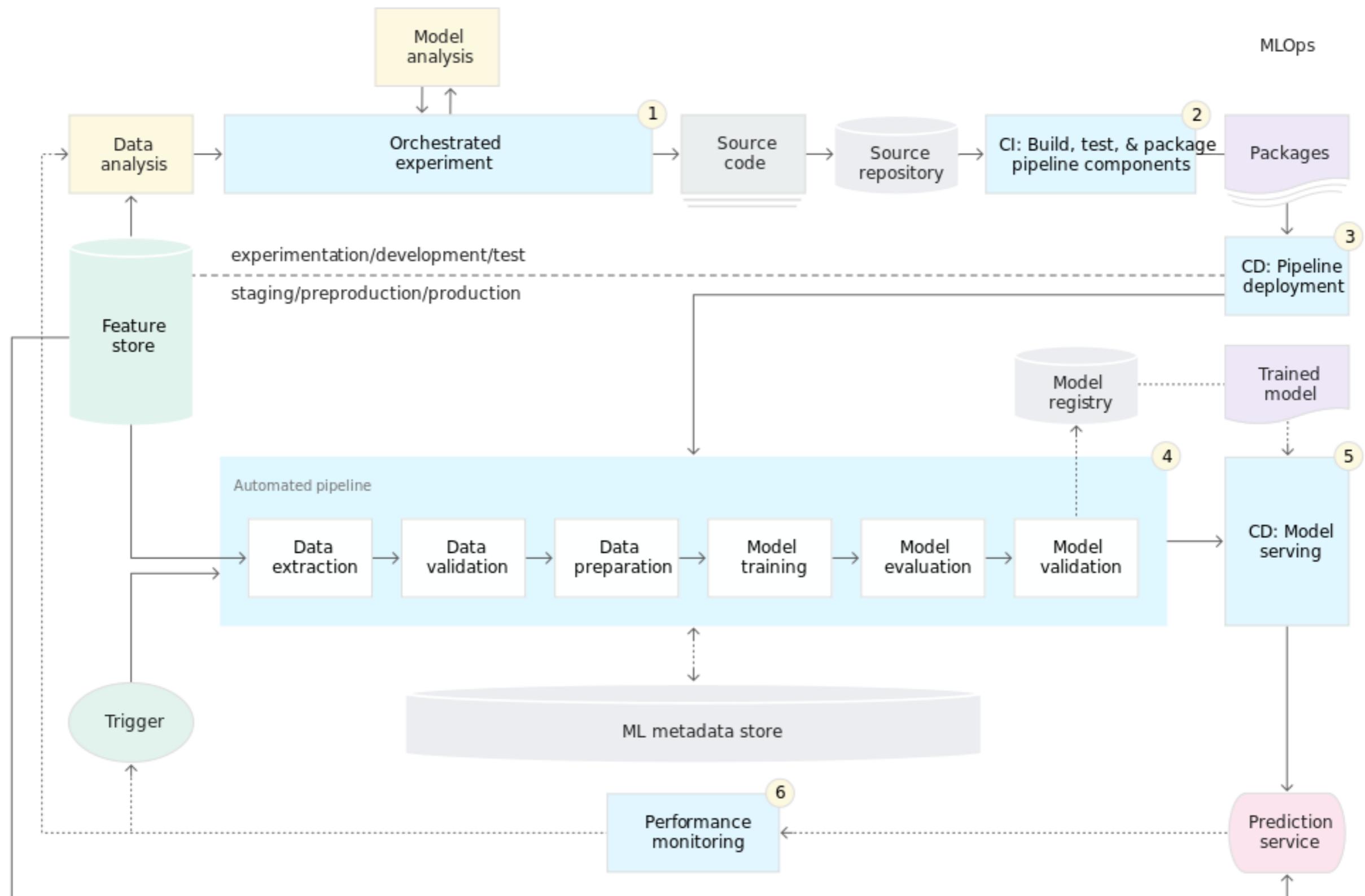
MLOps level 1



MLOps level 1



MLOps level 2



Summary

Summary

Implementing ML in a production environment doesn't **only** mean deploying your model as an API **for prediction** but rather deploying an ML pipeline that can **automate the retraining and deployment** of new models.

Setting up a **CI/CD** system enables you to **automatically test and deploy** new pipeline implementations.

This system lets you **cope with rapid changes** in your data and business environment.

You don't have to immediately move all of your processes from **one level to another**.

MLOps tools

MLOps tools

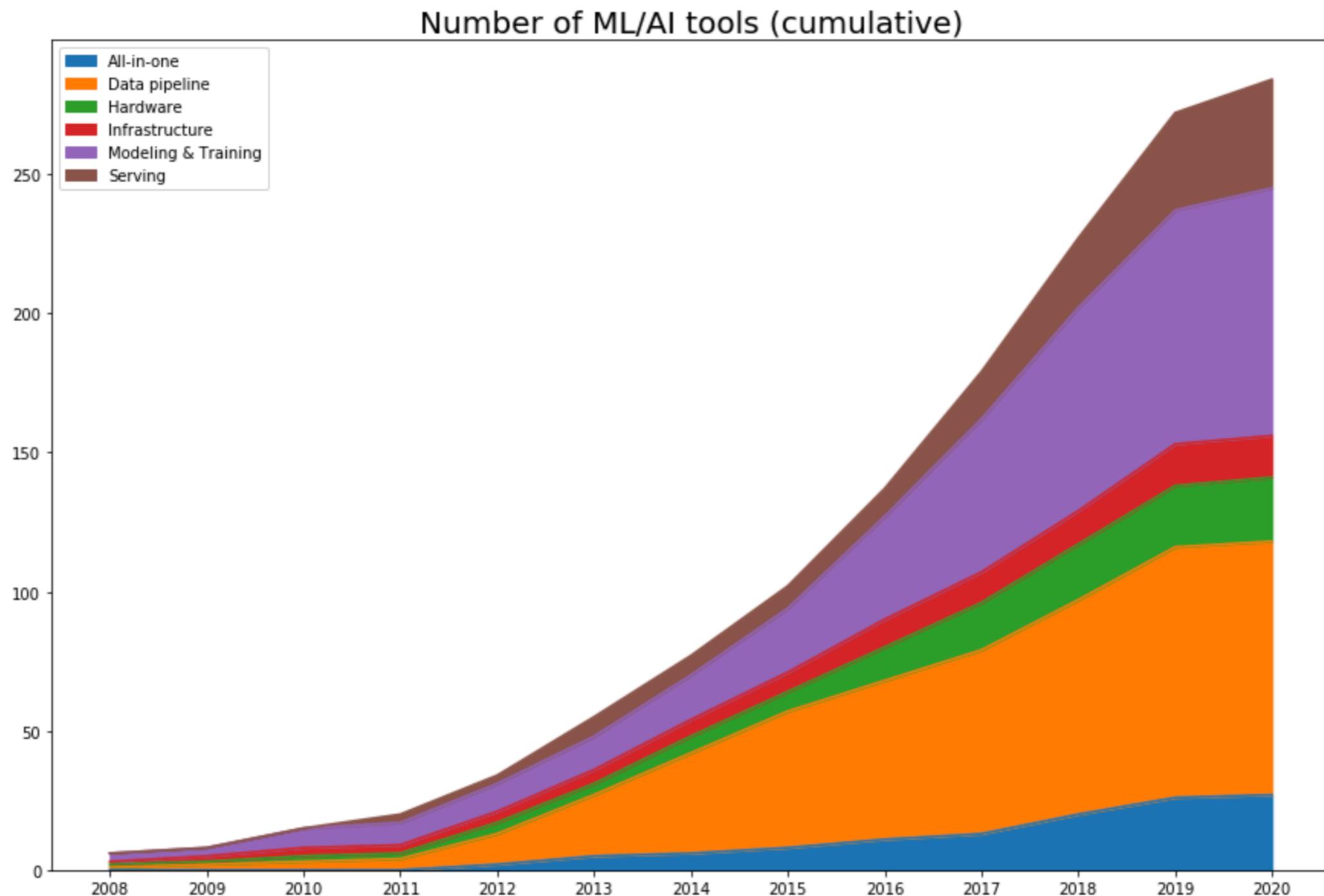
BIG DATA & AI LANDSCAPE 2018

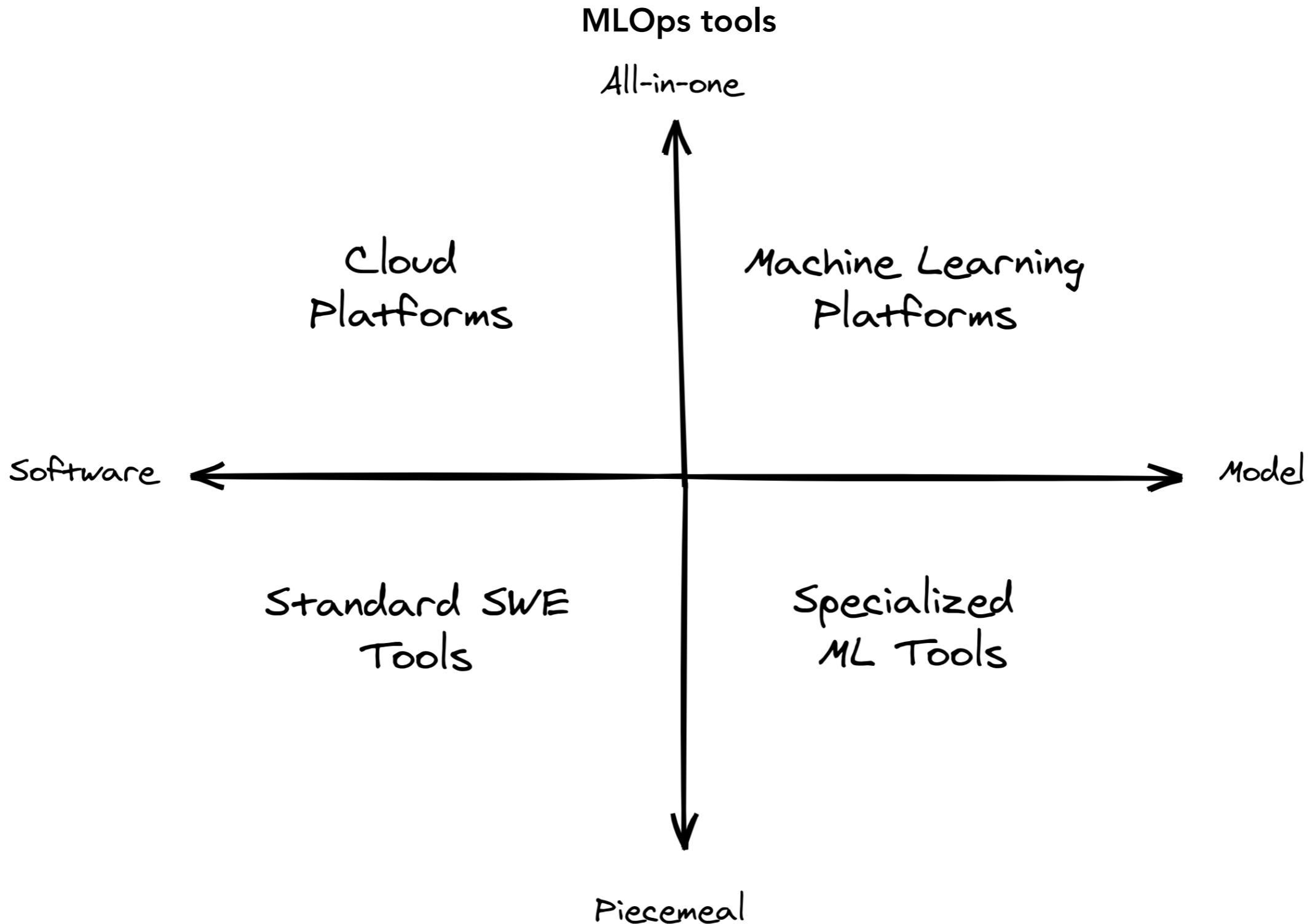


MLOps tools

MLOps Principles	Data	ML Model	Code
Versioning	1) Data preparation pipelines 2) Features store 3) Datasets 4) Metadata	1) ML model training pipeline 2) ML model (object) 3) Hyperparameters 4) Experiment tracking	1) Application code 2) Configurations
Testing	1) Data Validation (error detection) 2) Feature creation unit testing	1) Model specification is unit tested 2) ML model training pipeline is integration tested 3) ML model is validated before being operationalized 4) ML model staleness test (in production) 5) Testing ML model relevance and correctness 6) Testing non-functional requirements (security, fairness, interpretability)	1) Unit testing 2) Integration testing for the end-to-end pipeline
Automation	1) Data transformation 2) Feature creation and manipulation	1) Data engineering pipeline 2) ML model training pipeline 3) Hyperparameter/Parameter selection	1) ML model deployment with CI/CD 2) Application build
Reproducibility	1) Backup data 2) Data versioning 3) Extract metadata 4) Versioning of feature engineering	1) Hyperparameter tuning is identical between dev and prod 2) The order of features is the same 3) Ensemble learning: the combination of ML models is same 4) The model pseudo-code is documented	1) Versions of all dependencies in dev and prod are identical 2) Same technical stack for dev and production environments 3) Reproducing results by providing container images or virtual machines
Deployment	1) Feature store is used in dev and prod environments	1) Containerization of the ML stack 2) REST API 3) On-premise, cloud, or edge	1) On-premise, cloud, or edge
Monitoring	1) Data distribution changes (training vs. serving data) 2) Training vs serving features	1) ML model decay 2) Numerical stability 3) Computational performance of the ML model	1) Predictive quality of the application on serving data

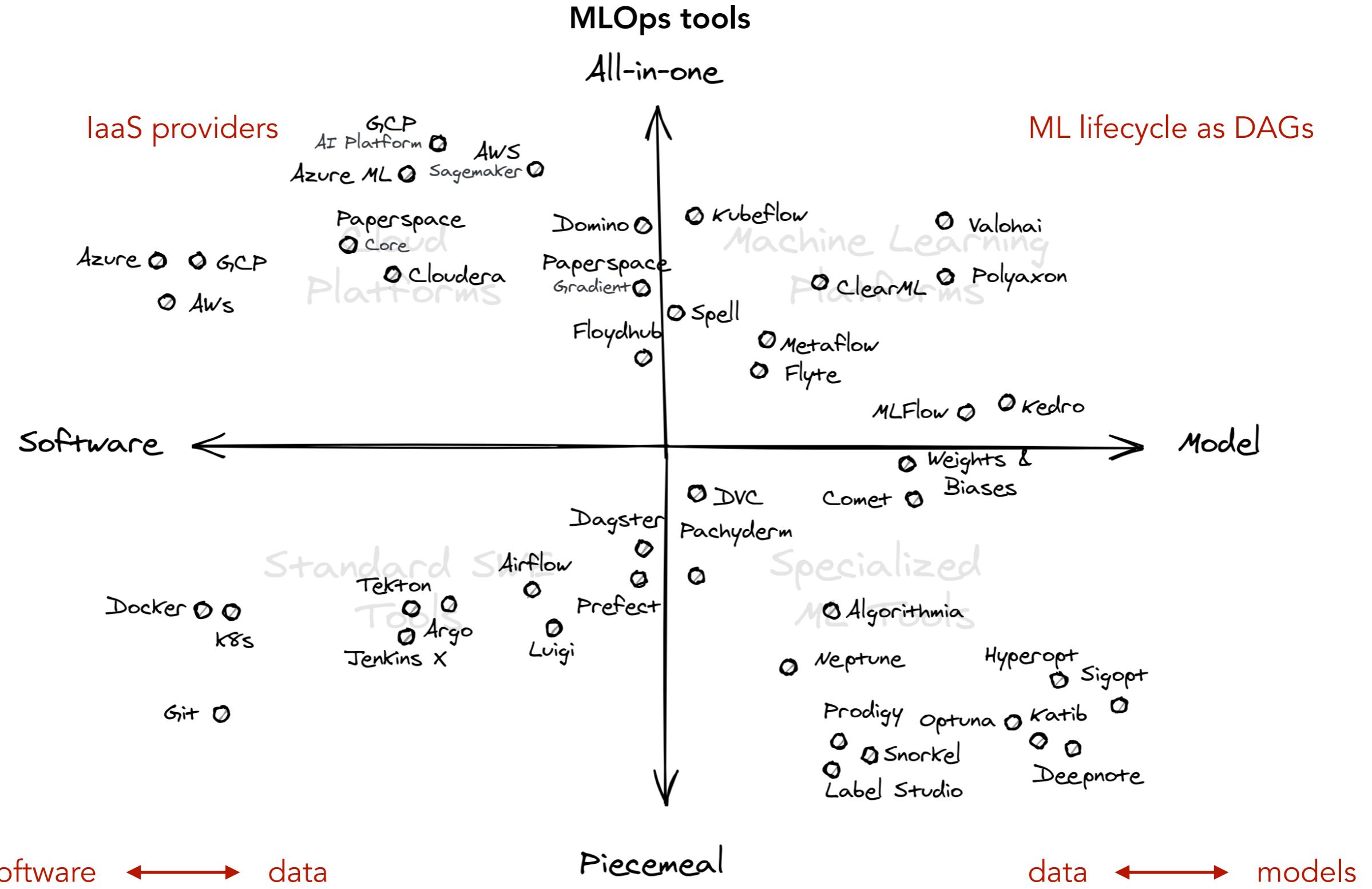
MLOps tools





x-axis: artifact

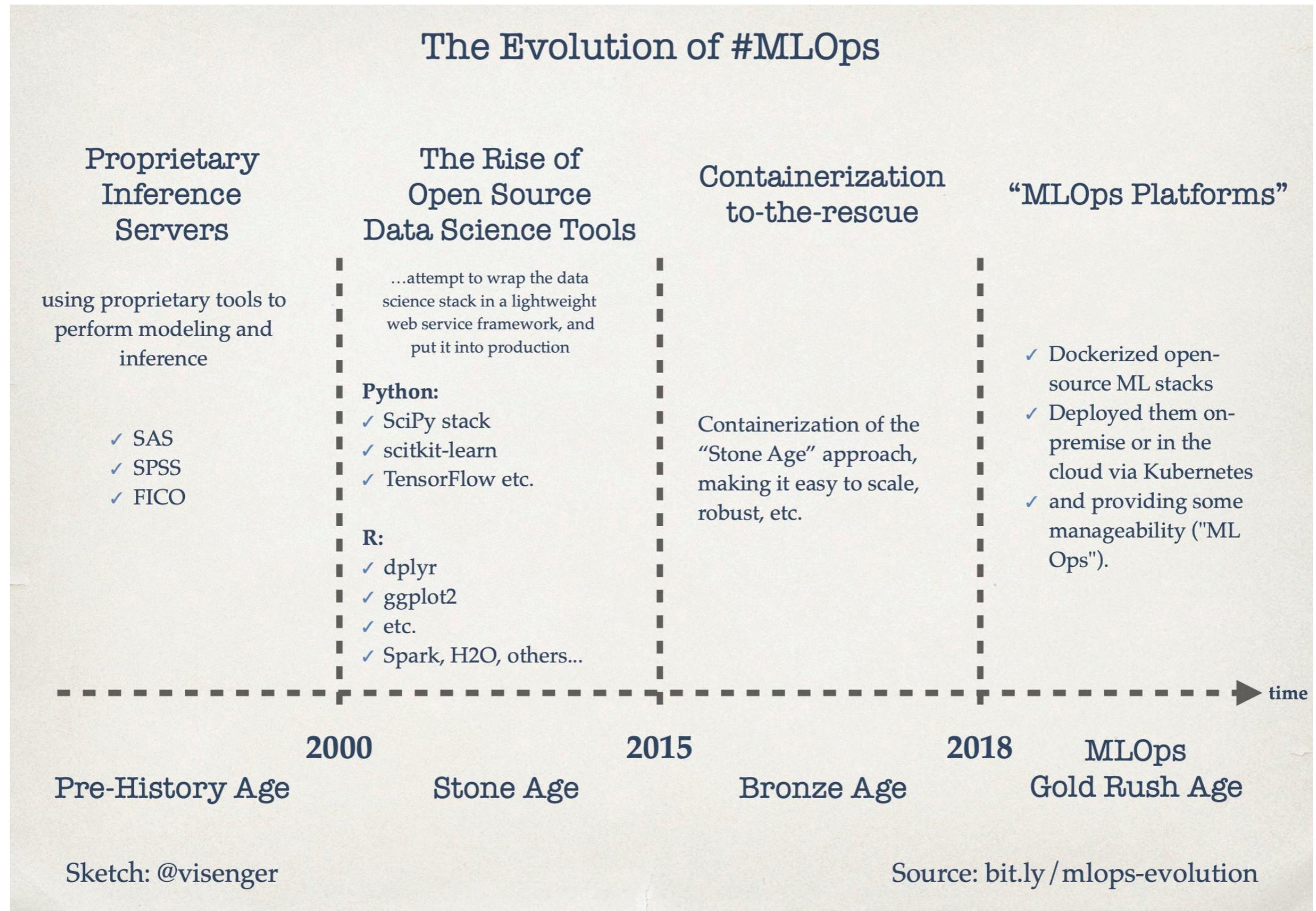
y-axis: scope



x-axis: artifact

y-axis: scope

Predictions



Signals from Big AI Companies

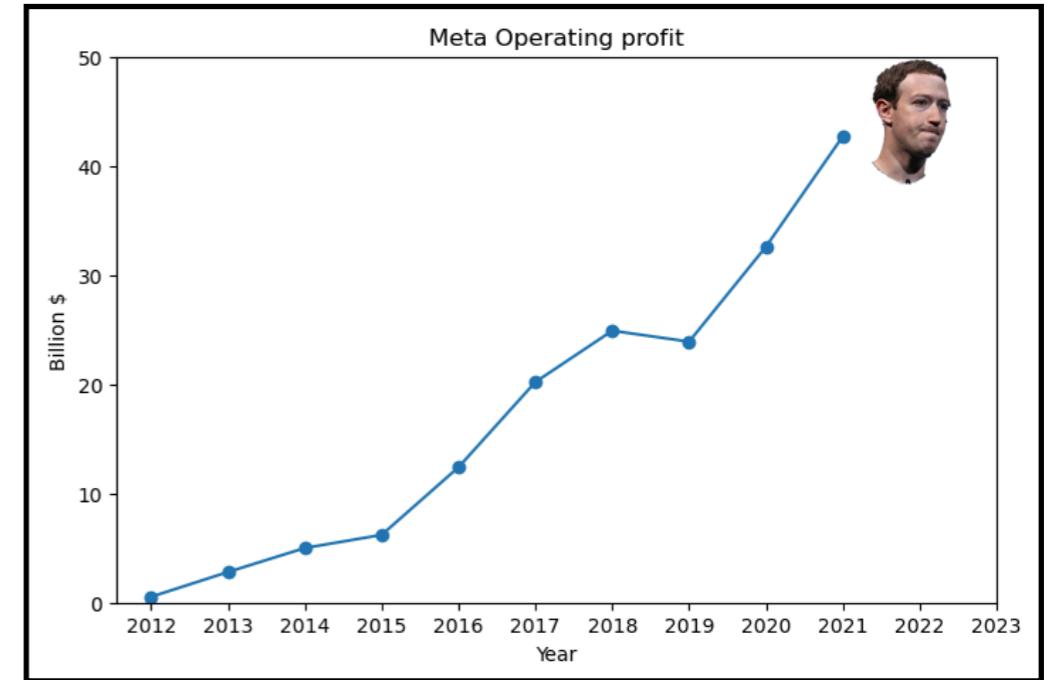
Signals from Big AI Companies

Layoffs



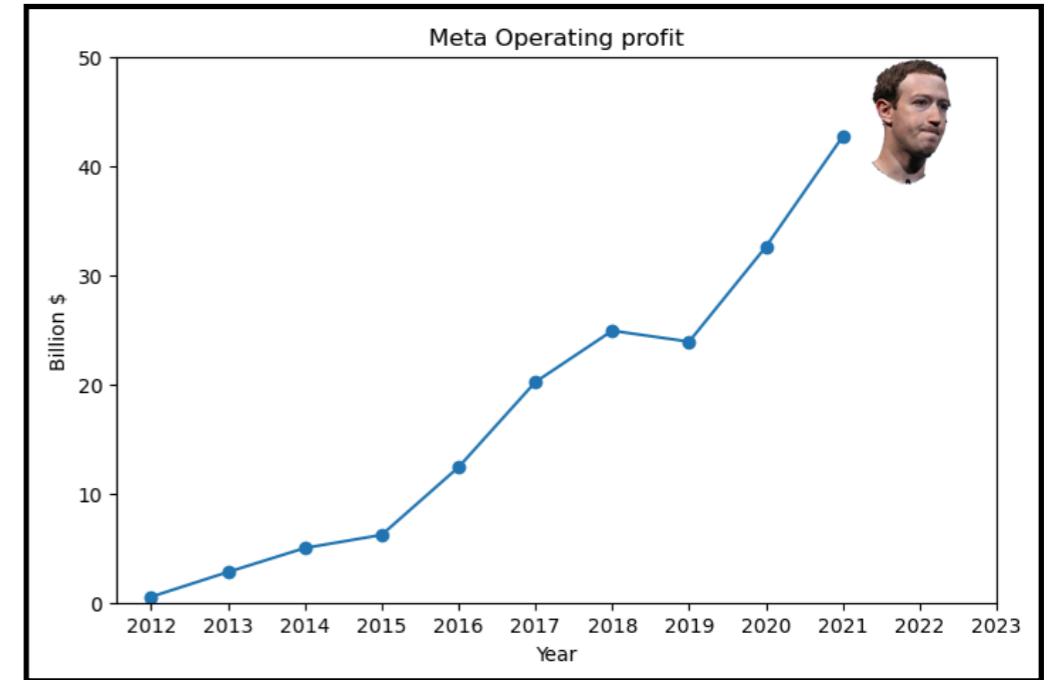
Signals from Big AI Companies

Layoffs



Signals from Big AI Companies

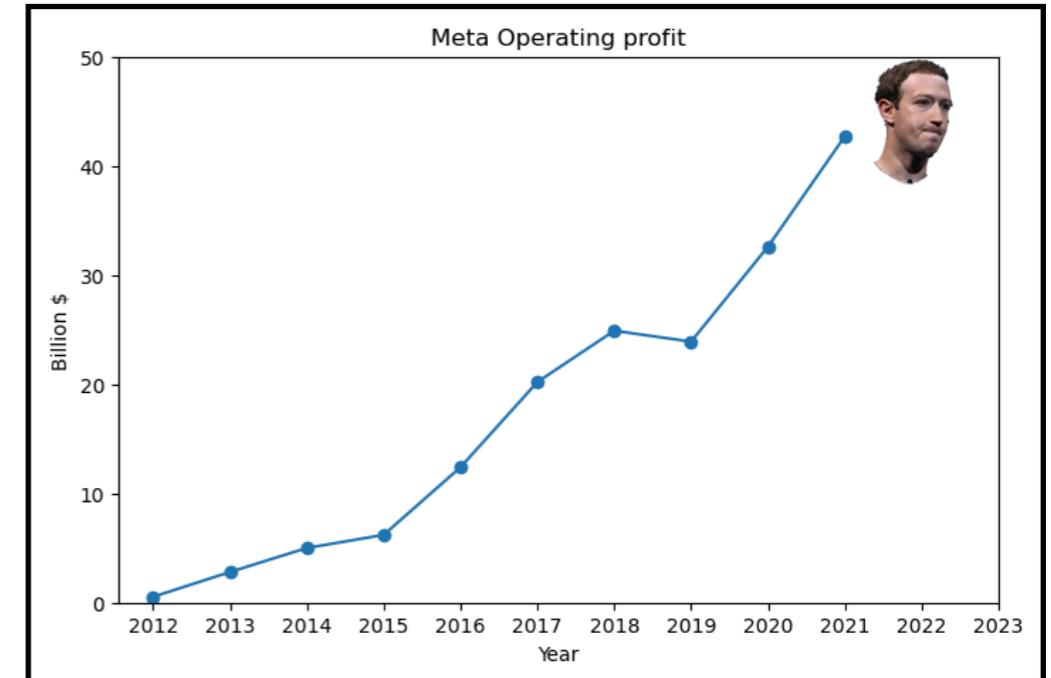
Layoffs



Physical world's problems

Signals from Big AI Companies

Layoffs

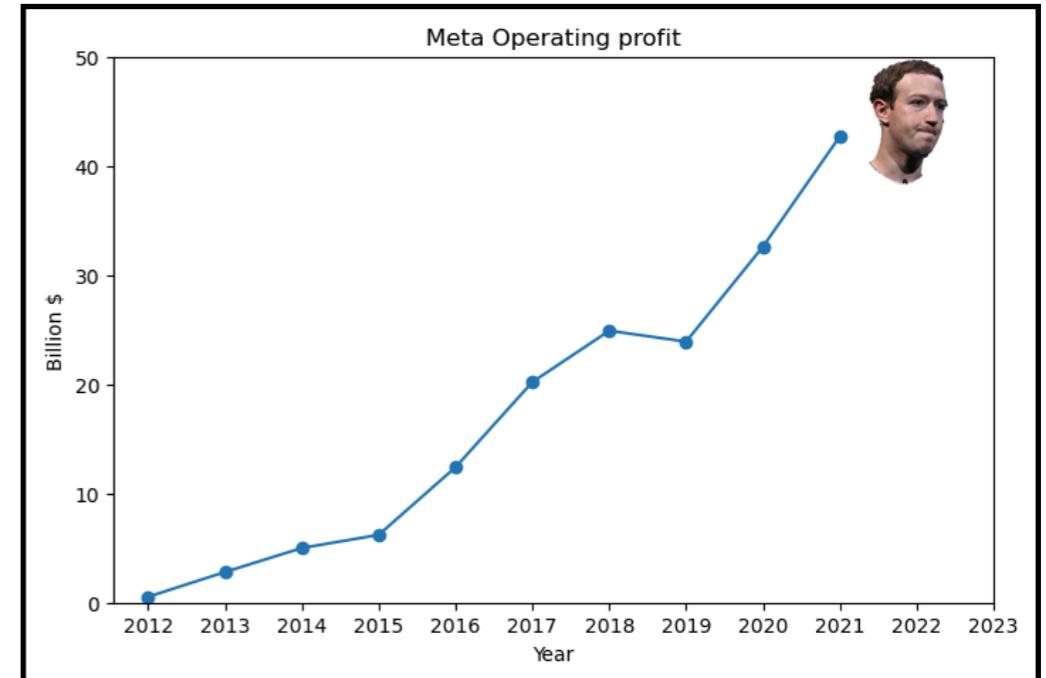


Physical world's problems

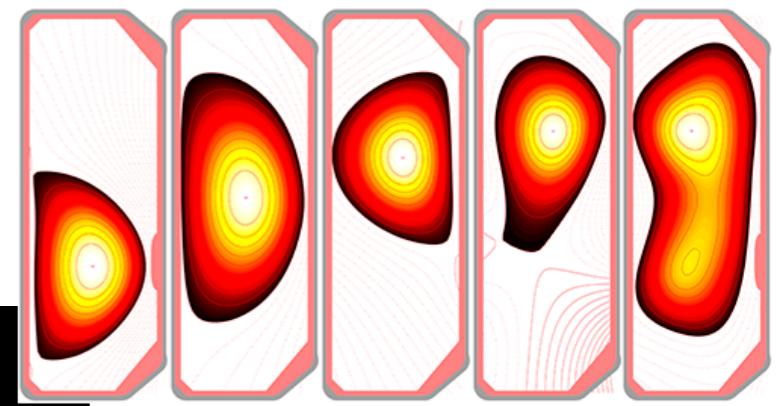


Signals from Big AI Companies

Layoffs



Physical world's problems



Prediction / Hope

Prediction / Hope

MLOps adoption and ease of use might allow different organisations (governments, research institutes, all type of companies) to deploy and use AI to **help solving or mitigating** some of the problems we're inevitably going to face in the next decades (food and energy scarcity, climate change damages)

Prediction / Hope

MLOps adoption and ease of use might allow different organisations (governments, research institutes, all type of companies) to deploy and use AI to **help solving or mitigating** some of the problems we're inevitably going to face in the next decades (food and energy scarcity, climate change damages)

And we might realise that maybe AI was never meant to be used for generating profits or controlling people using recommended ads and customer behaviour profiling **but to solve real problems in the physical world**

References

Books :

- “[Designing Machine Learning Systems](#)” - Chip Huyen, O'Reilly 2022
- “[Introducing MLOps](#)” - Mark Travel & the Dataiku team, O'Reilly 2020

Papers / Blogs / Reports / Videos :

- “[2020 state of enterprise machine learning](#)” - Algorithmia
- “[Hidden Technical Debt in Machine Learning Systems](#)” - Sculley et al. (Google Inc.), NIPS 2015
- “[Practitioners guide to MLOps](#)” - Google cloud white paper, 2021
- “[MLOps: Continuous delivery and automation pipelines in machine learning](#)” - Google Cloud Architecture Center, 2020
- “[Why data scientists shouldn't need to know Kubernetes](#)” - Chip Huyen's blog, 2021
- “[Machine Learning Tools Landscape v2](#)” - Chip Huyen's blog, 2020
- “[Navigating the MLOps tooling landscape \(Part I, II, III\)](#)” - Lj Miranda's blog, 2021
- “[Software 2.0](#)” - Andrej Karpathy Medium blog, 2017 (also on [youtube](#))
- “[Rules of Machine Learning: Best Practices for ML Engineering](#)” - M. Zinkevich (Google developers guide), 2022
- “[Machine Learning Operations \(MLOps\): Overview, Definition, and Architecture](#)” - D. Kreuzberger et al., 2022