

FIAP

MBA



Human-Centered Data & AI



Vinicius Caridá, Ph.D.

- Head of Digital Customer Service Platforms, PCP, WFM, Data and AI - Itaú Unibanco
- MBA Professor - FIAP



Machine Learning

Data science and ML are **becoming core capabilities** for solving complex real-world problems, transforming industries, and delivering value in all domains. Therefore, many businesses are **investing in their data science teams and ML capabilities** to develop predictive models that can deliver business value to their users. While some companies have figured out how to succeed, **many companies still have difficulty to generate value with AI.**

“

As ML matures from research to applied
business solutions, so do we need to improve
the maturity of its operation processes

1

What is "ML Ops"?

DevOps for ML

What is **DevOps**

“DevOps is a software engineering culture and practice that aims at **unifying** software development (Dev) and software operation (Ops).”

“(DevOps is to) strongly advocate **automation and monitoring** at all steps of software construction, from integration, testing, releasing to deployment and infrastructure management.”

- Wikipedia



Code



Testing



Deploy



Monitor



Dev team

Ops team



**People respond
to incentives**



Machine Learning

Além de treinar um modelo incrível...

Código ML

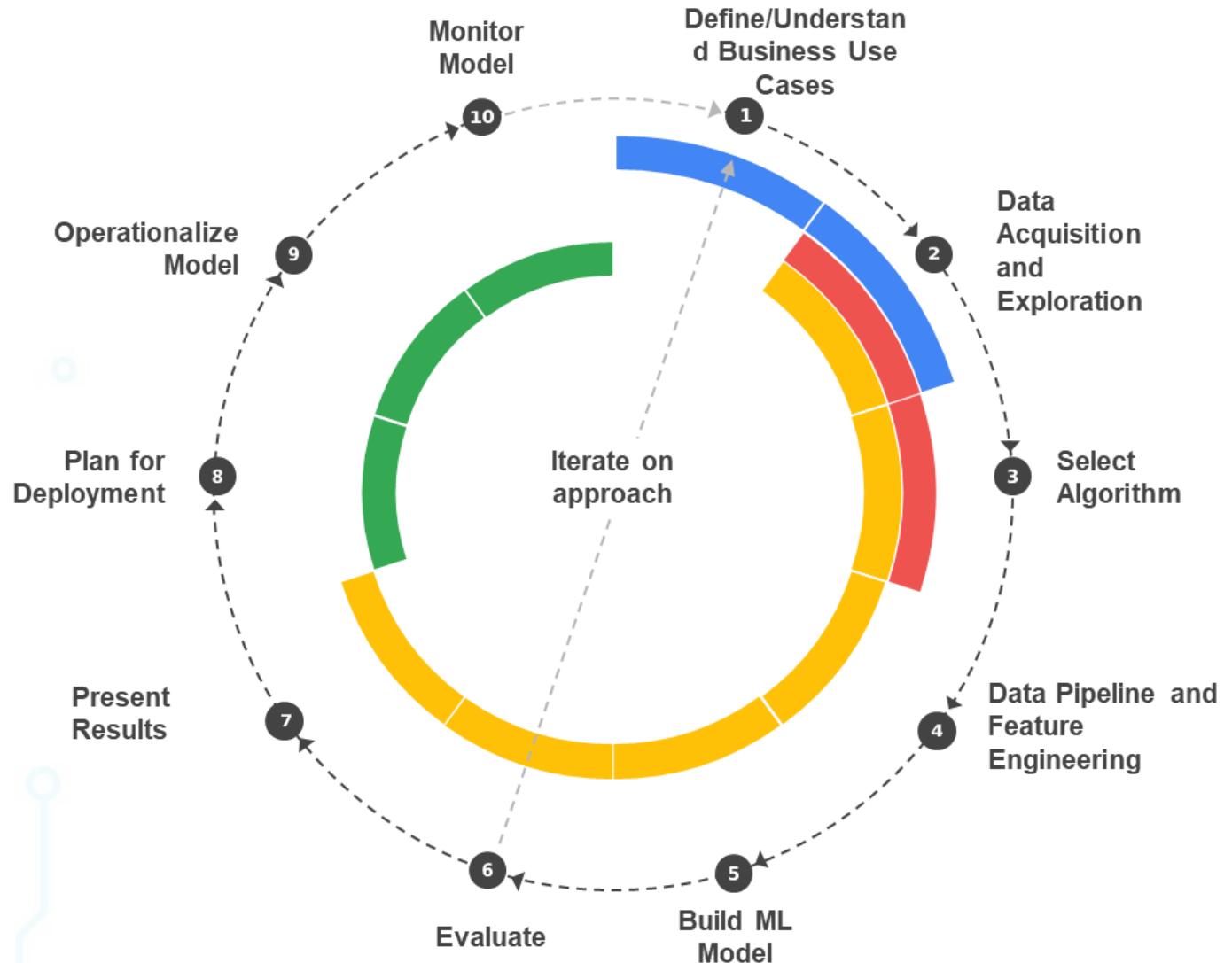
Machine Learning

Realidade: ML requer DevOps









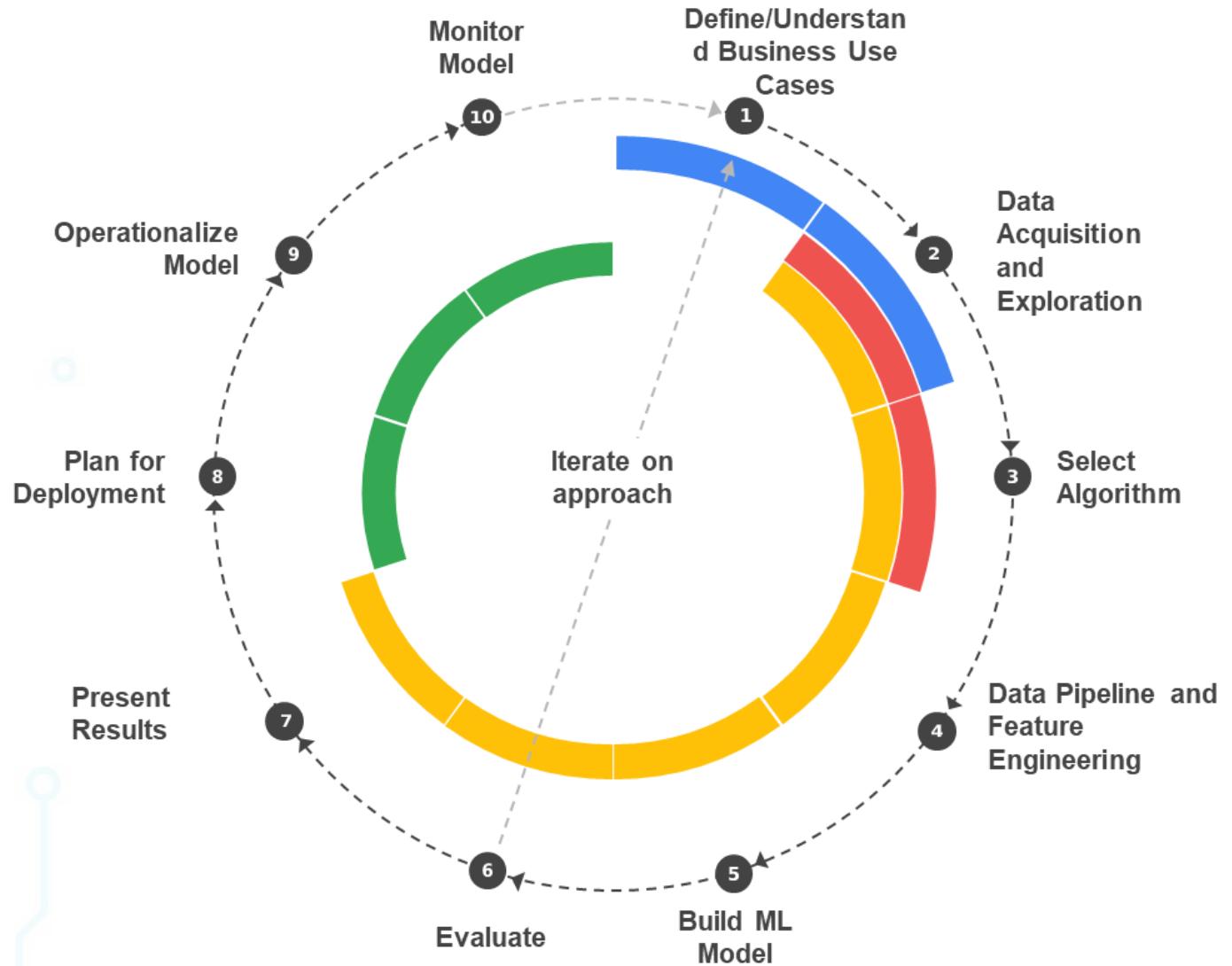
If ML is a rocket engine,
data is the fuel



**Launching is easy,
Operating is hard.**

**"The real problems with a
ML system will be found
while you are continuously
operating it for the long term"**





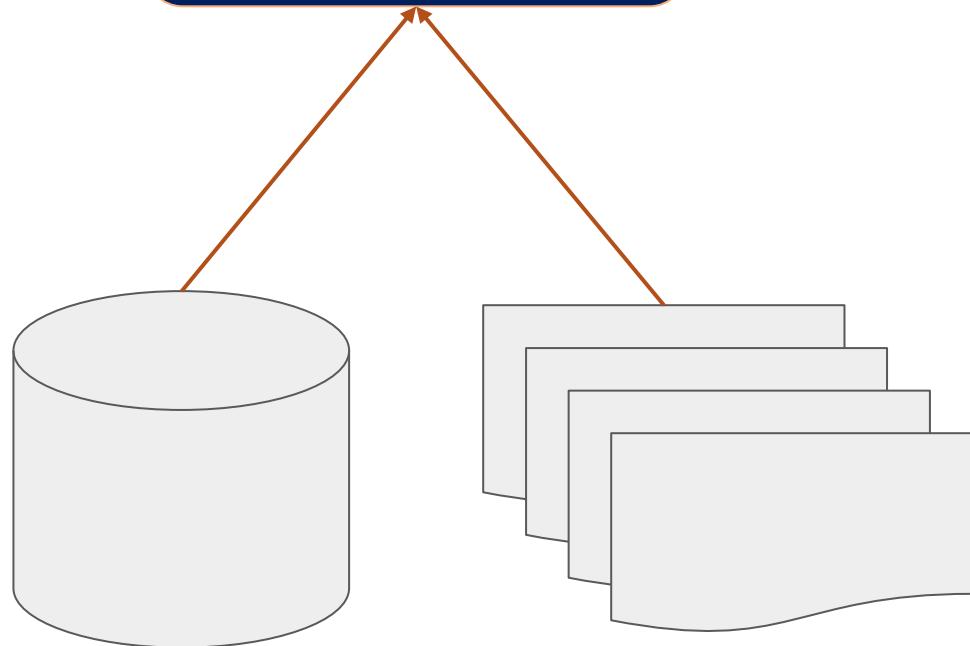
“More than 87% of data science projects never make it into production”

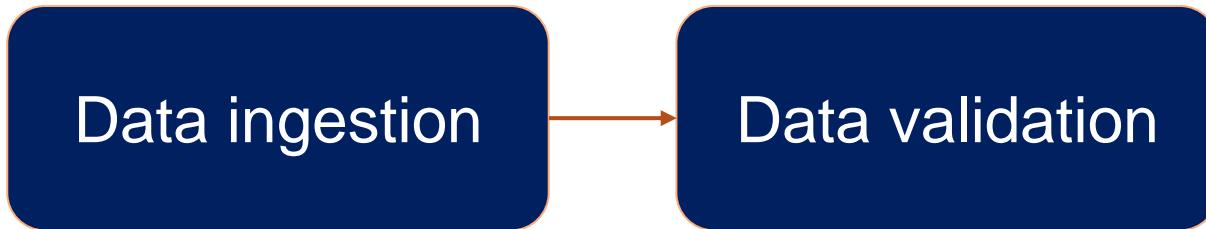
- Multiple studies and surveys -



Source: <https://venturebeat.com/2019/07/19/why-do-87-of-data-science-projects-never-make-it-into-production/>

Data ingestion



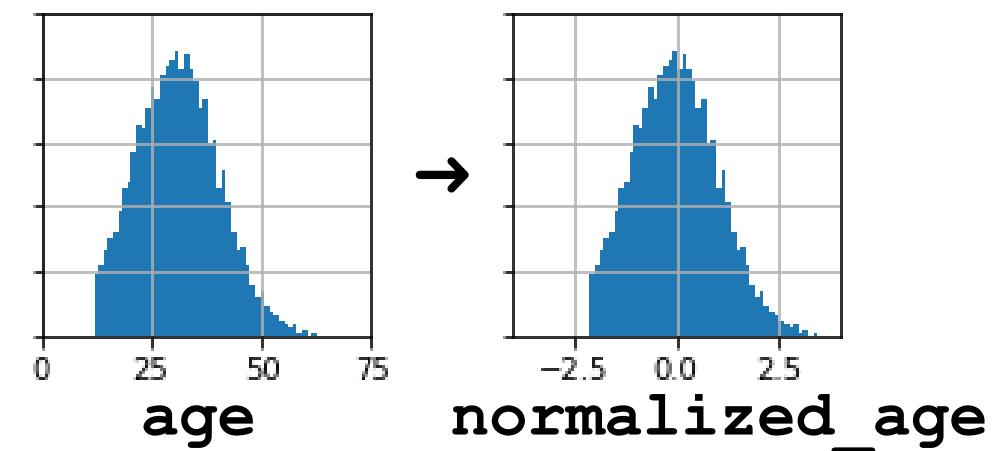


age is missing

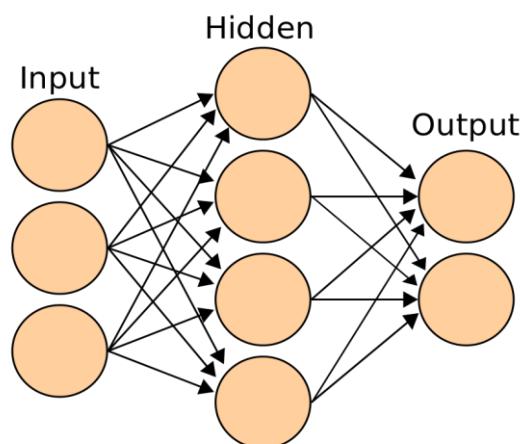
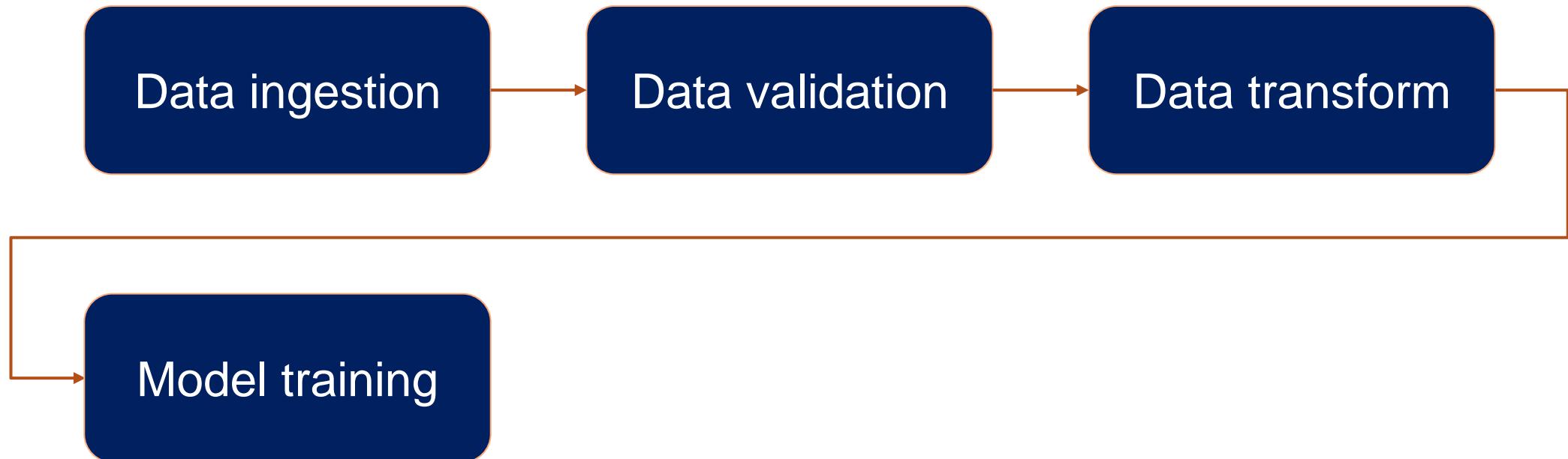


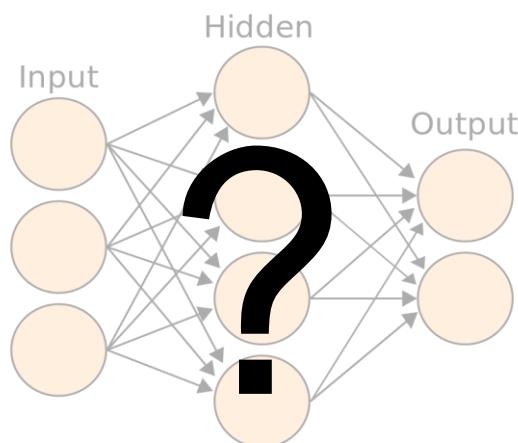
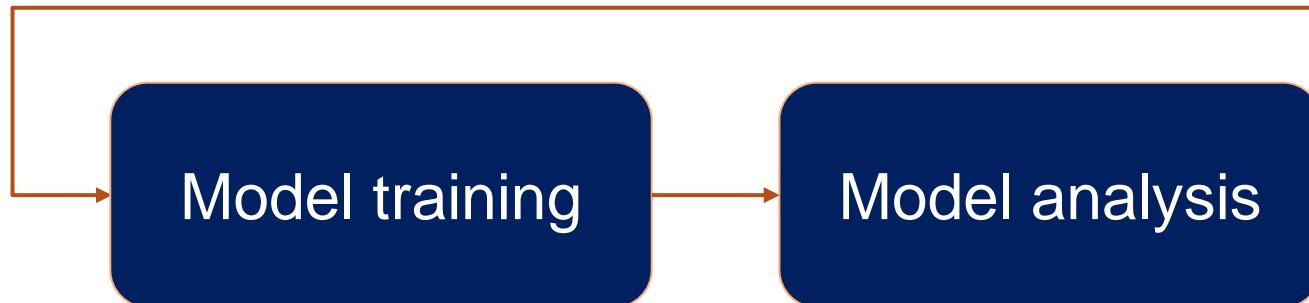
country not in:

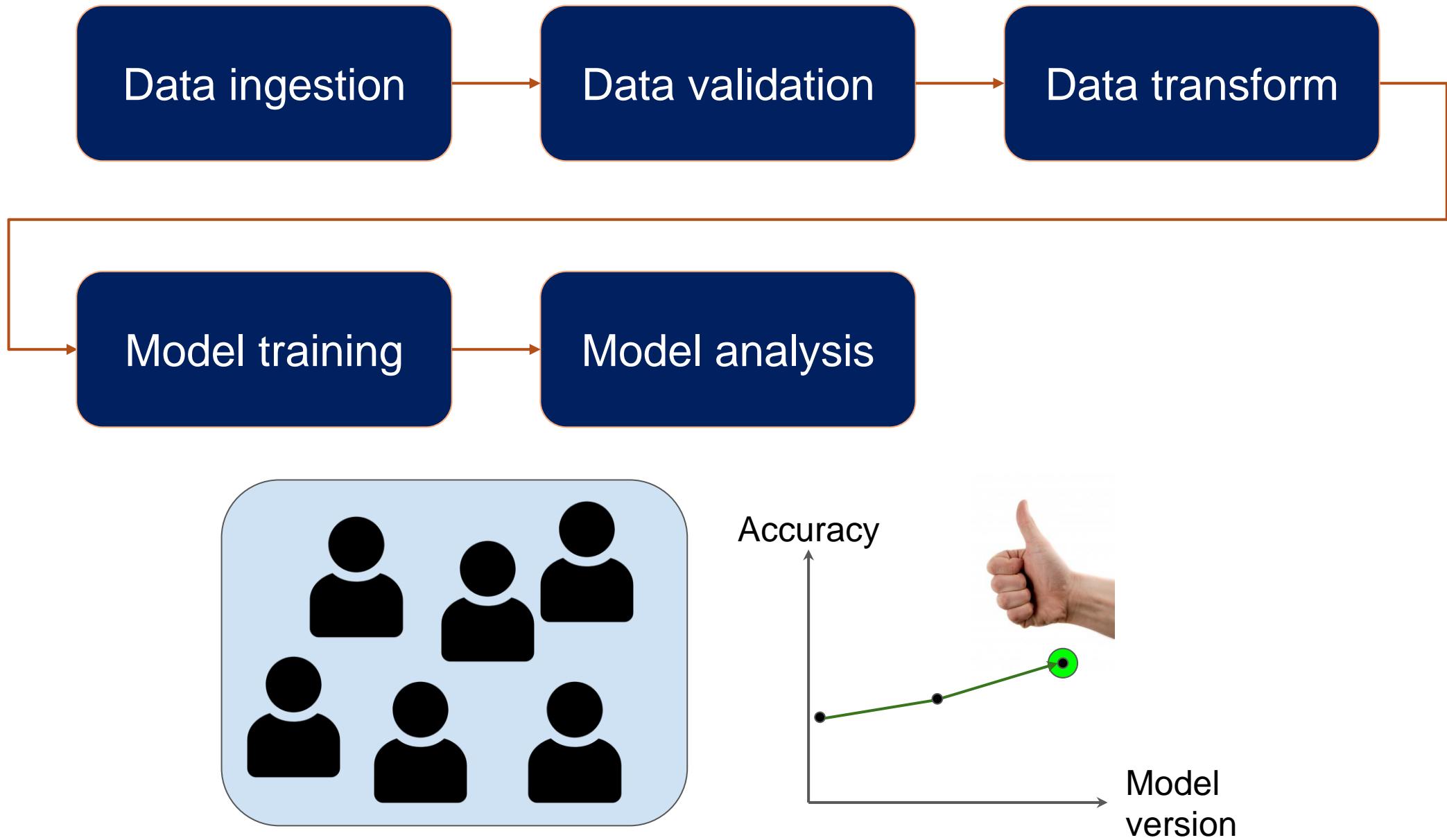
- China
- India
- USA

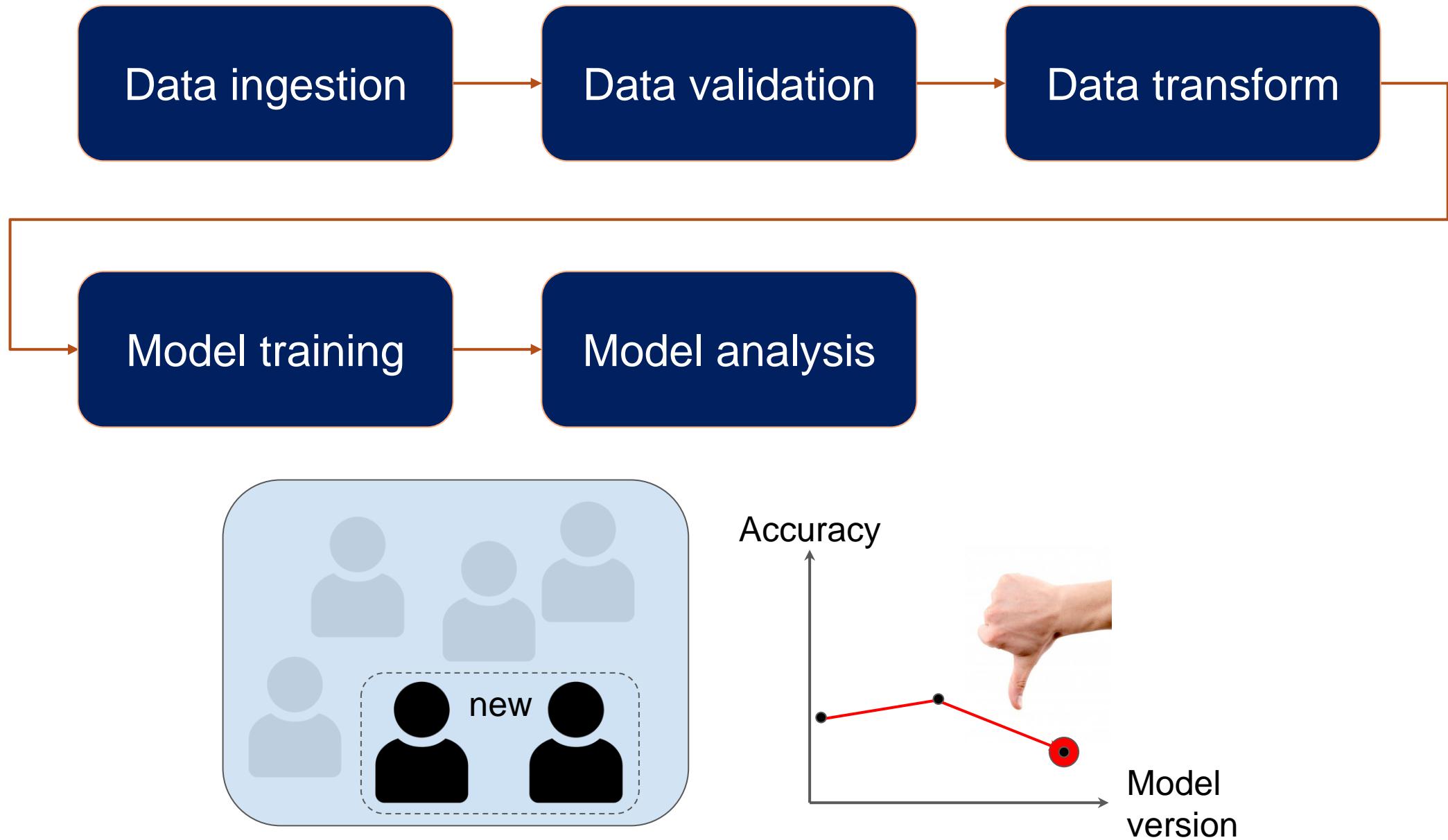


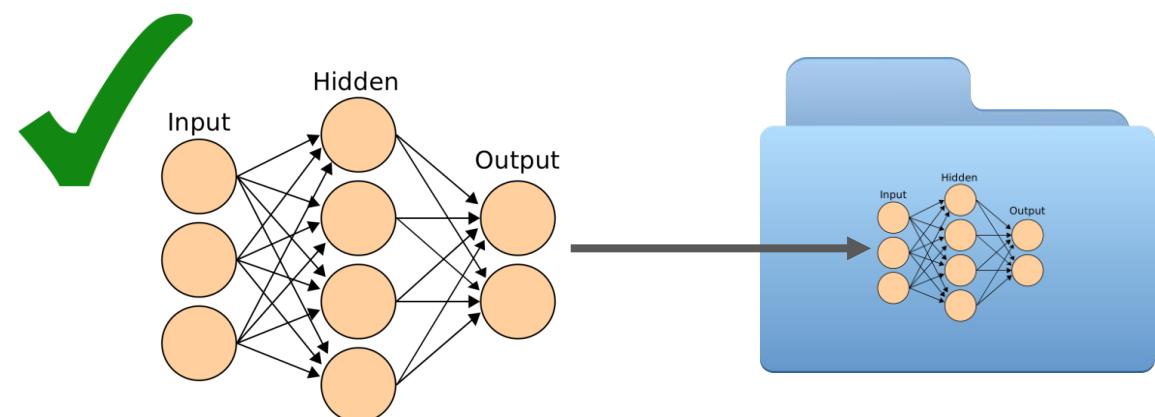
China → [1, 0, 0]
India → [0, 1, 0]
USA → [0, 0, 1]

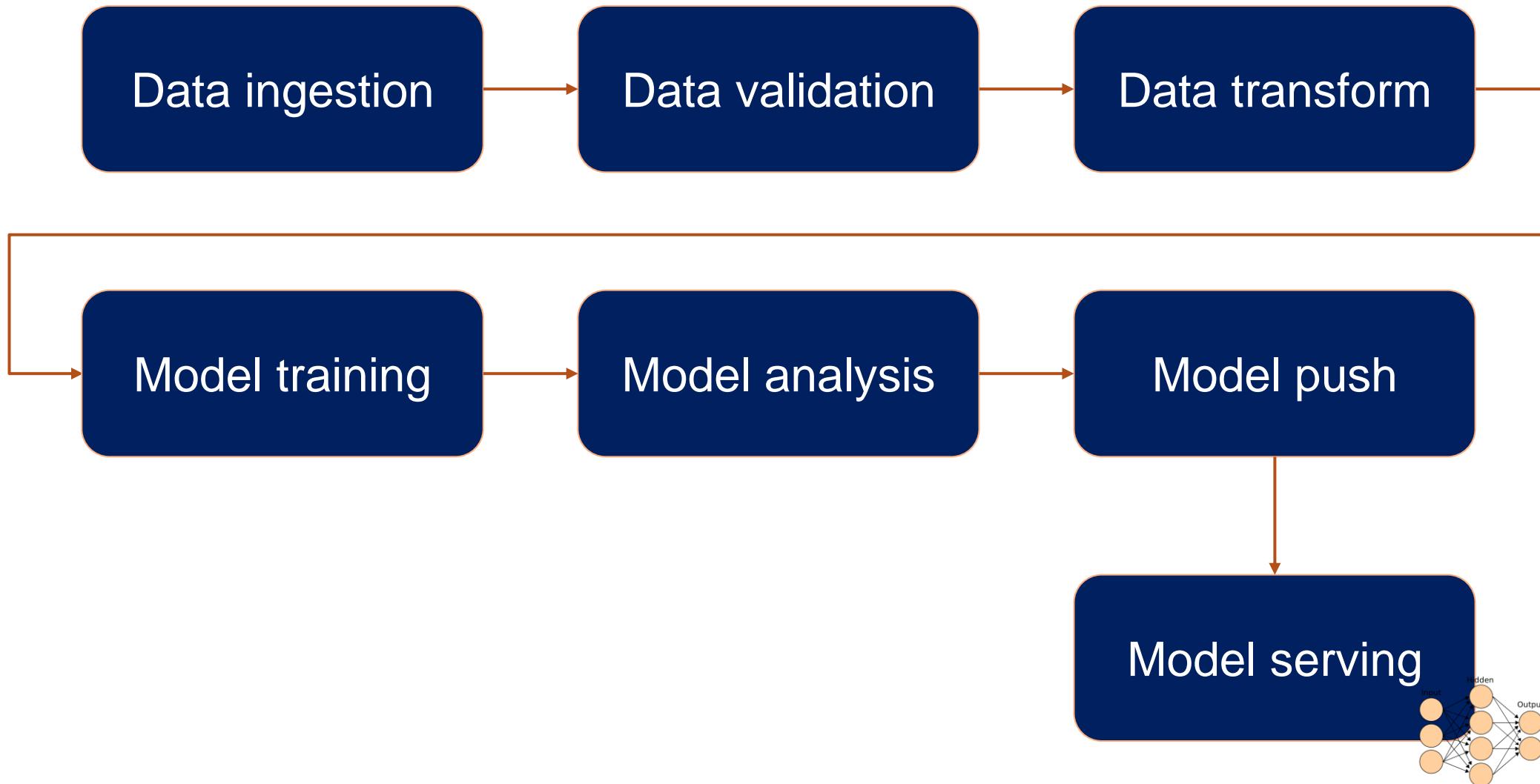


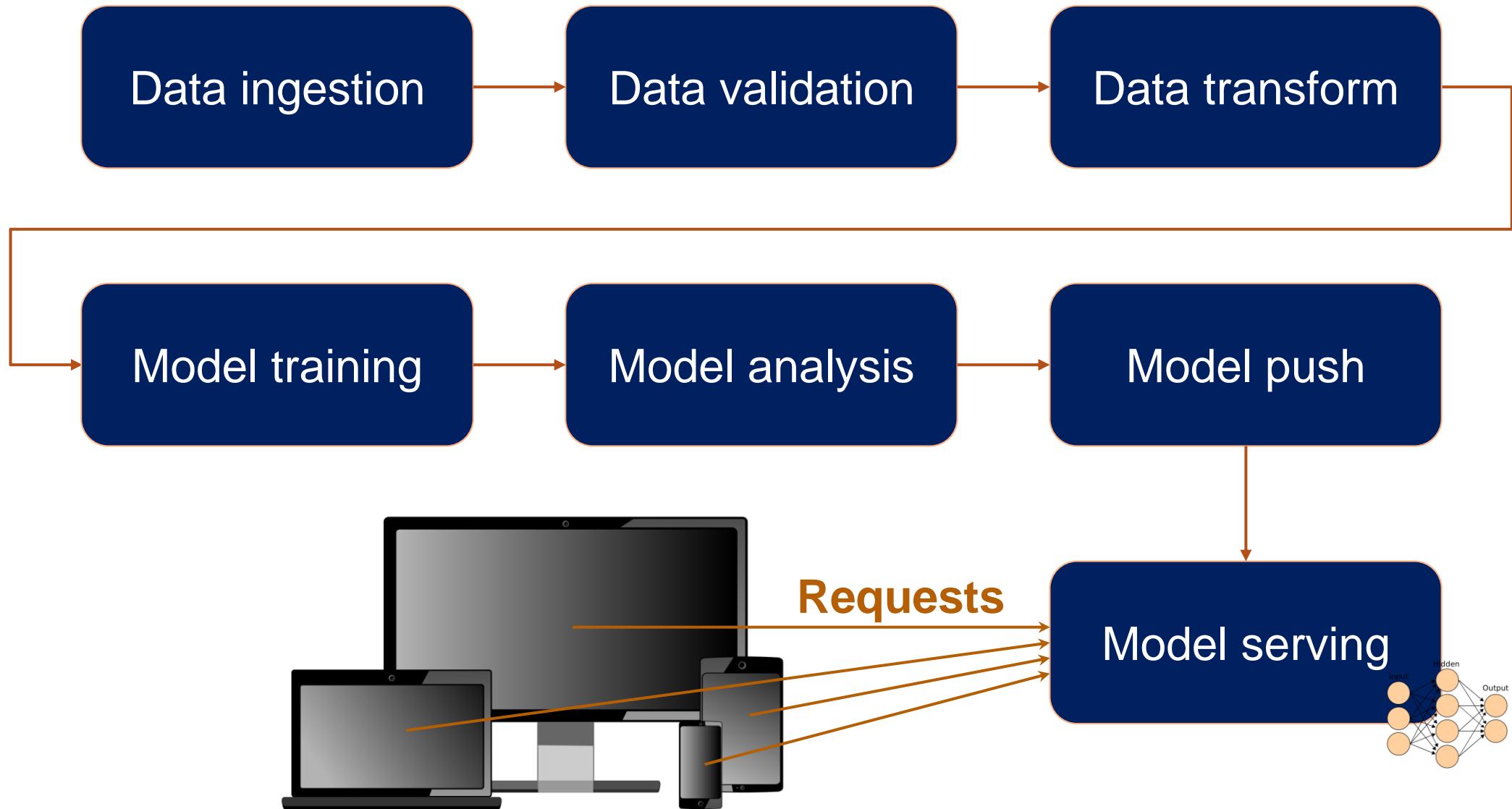


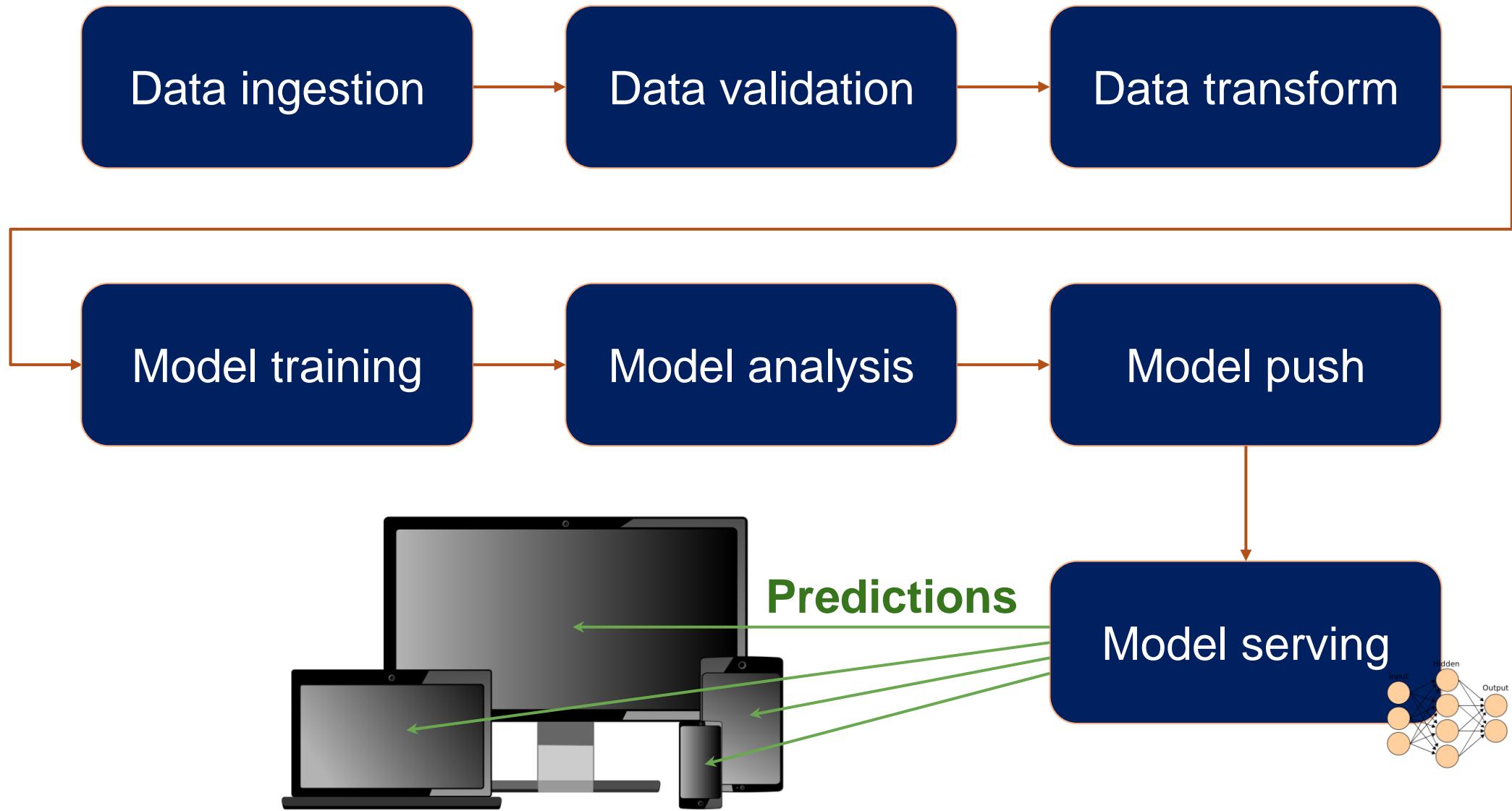


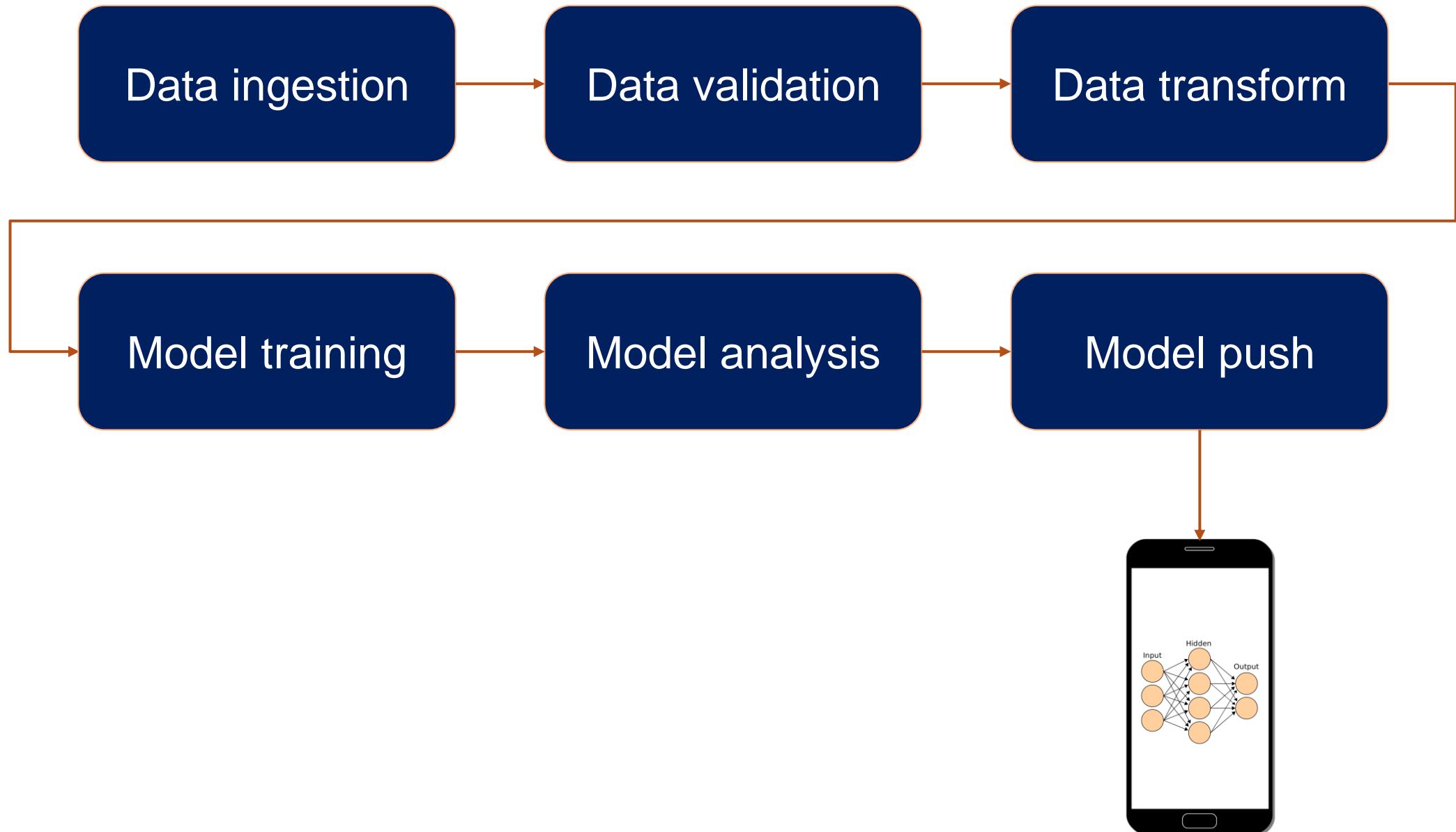


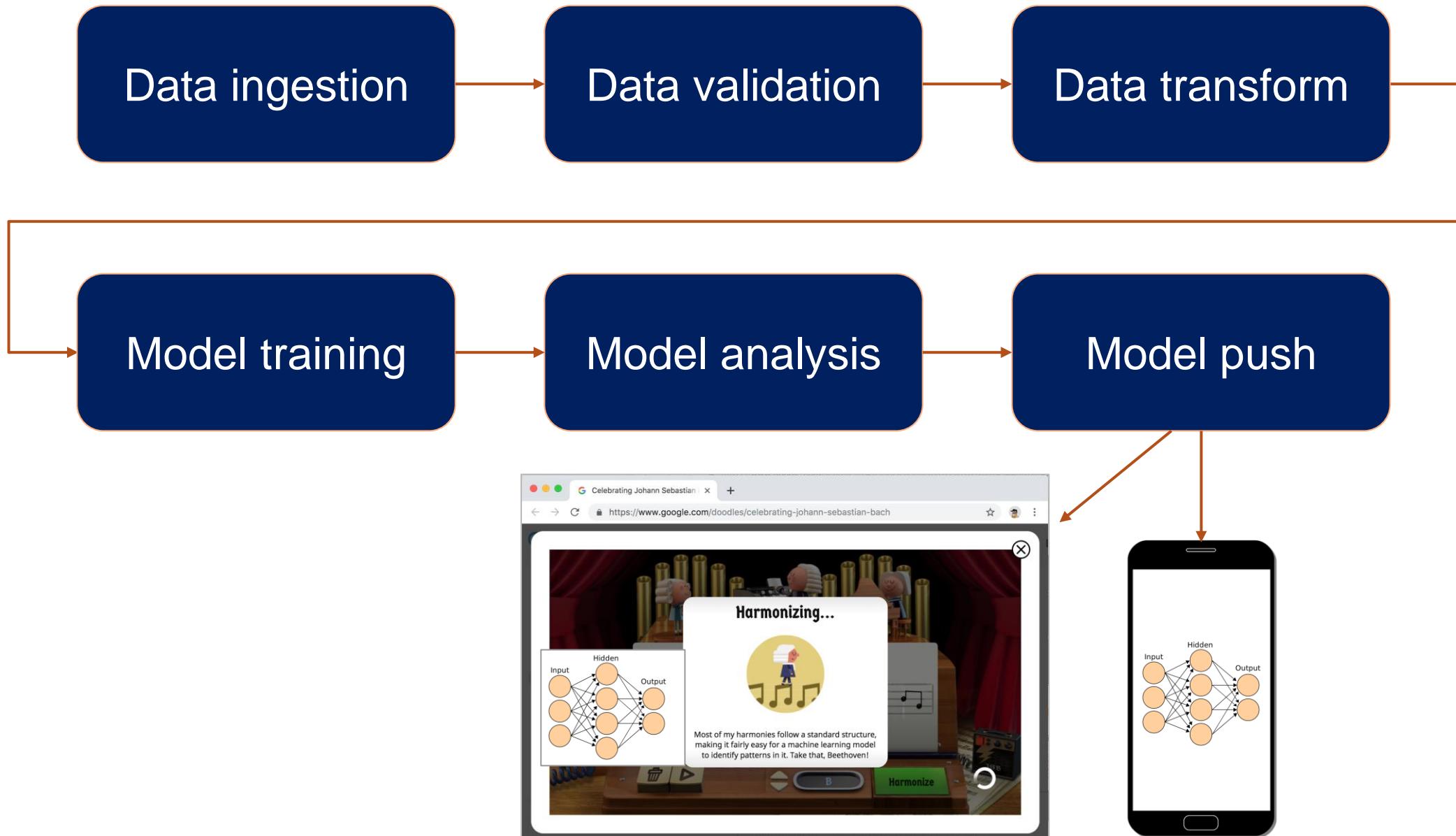














Production Machine Learning

Machine Learning Development

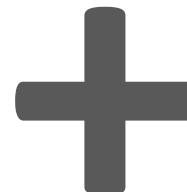
- Labeled data
- Feature space coverage
- Minimal dimensionality
- Maximum predictive data
- Fairness
- Rare conditions
- Data lifecycle management



Production Machine Learning

Machine Learning Development

- Labeled data
- Feature space coverage
- Minimal dimensionality
- Maximum predictive data
- Fairness
- Rare conditions
- Data lifecycle management



Modern Software Development

- Scalability
- Extensibility
- Configuration
- Consistency & Reproducibility
- Modularity
- Best Practices
- Testability
- Monitoring
- Safety & Security

Writing Software (Programming)

Programming in the small (Coding)

Monolithic code

Non-reusable code

Undocumented code

Untested code

Unbenchmarked or hack-optimized once code

Unverified code

Undebuggable code or adhoc tooling

Uninstrumented code

...

Programming in the large (Engineering)

Modular design and implementation

Libraries for reuse (ideally across languages)

Well documented contracts and abstractions

Well tested code (exhaustively and at scale)

Continuously benchmarked and optimized code

Reviewed and peer verified code

Debuggable code and debug tooling

Instrumentable and instrumented code

...

Writing ML Software (The “Code” view)

ML Programming in the small (Coding)

- Monolithic code
- Non-reusable code
- Undocumented code
- Untested code
- Unbenchmarked or hack-optimized once code
- Unverified code
- Undebuggable code or adhoc tooling
- Uninstrumented code

ML Programming in the large (Engineering)

- Modular design and implementation
- Libraries for reuse (ideally across languages)
- Well documented contracts and abstractions
- Well tested code (exhaustively and at scale)
- Continuously benchmarked and optimized code
- Reviewed and peer verified code
- Debuggable code and debug tooling
- Instrumentable and instrumented code

This slide is, not surprisingly, the same as the previous one however it is **only half** the story :)

Writing ML Software (The “Data and other Artifacts” view)

ML Programming in the small (Coding)

Monolithic code Fixed Datasets

Non-reusable code Unmergeable Artifacts

Undocumented code No Problem Statements

Untested code Non-validated Datasets, Models

Unbenchmarked or hack-optimized once code Models

Unverified code Biased Datasets / Artifacts

Undebuggable code or adhoc tooling

Uninstrumented code

ML Programming in the large (Engineering)

Evolving Datasets (Data, Features, ...) and Objectives

Reusable Models aka Modules, Mergeable Statistics ...

Problem Statements, Discoverable Artifacts

Expectations, Data Validation, Model Validation ...

Quality and Performance Benchmarked Models ...

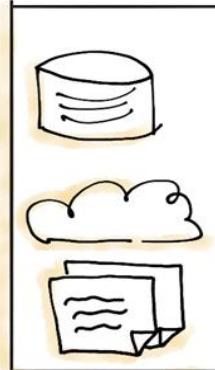
{Data, Model} x {Understanding, Fairness}

Visualizations, Summarizations, Understanding ...

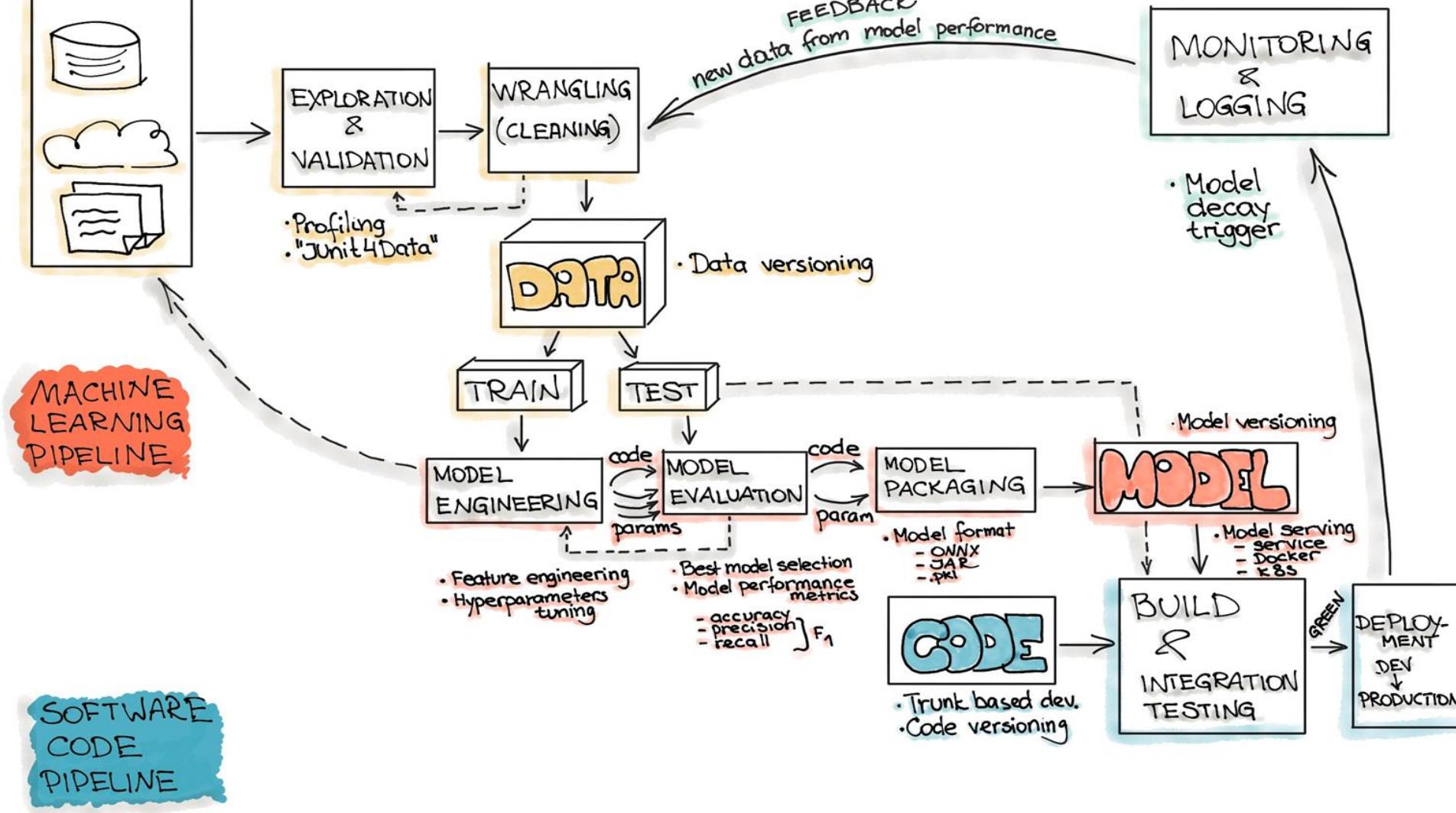
Full Artifact Lineage

This is the **remaining half!**

DATA PIPELINE



MACHINE LEARNING ENGINEERING.



SOFTWARE CODE PIPELINE

Intrinsic Complexities with Machine Learning

- Understanding the business domain
- Selecting the appropriate Model
- Selecting the appropriate Features
- Fine tuning

Incidental Complexities with Machine Learning

- Integrating with Data Warehouse
- Scaling model training & serving
- Keeping consistency between: Prototyping vs Production,
Training vs Inference
- Keeping track of multiple models, versions, experiments
- Supporting iteration on ML models

→ ML models take on average 8 to 12 weeks to build
→ ML workflows tended to be slow, fragmented, and brittle

ML Deploy Platforms

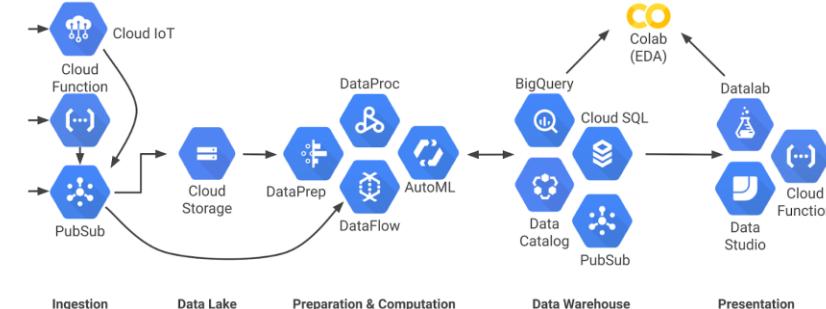
- Uber - [Michelangelo](#)
- AirBnB - [Bighead](#)
- Facebook - [FB Learner](#)
- Lyft - [Lyft Learn](#)
- Data Robot - [ParallelM](#)

Deploying big-data analytics, data science, and machine learning (ML) applications in real-world

Google

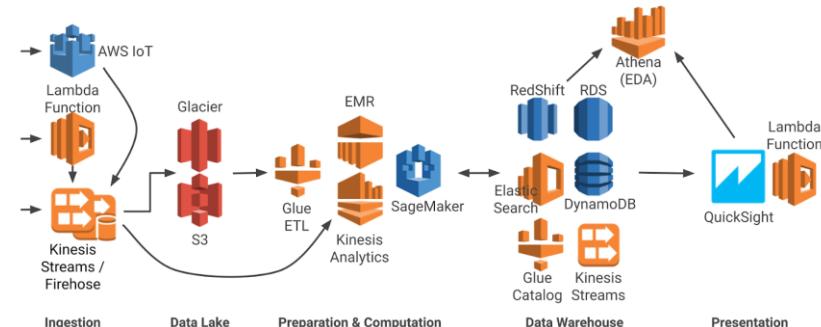


AI Platform



aws

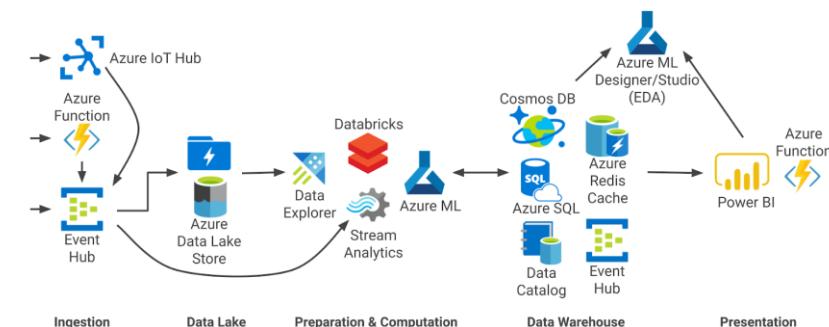
Amazon SageMaker



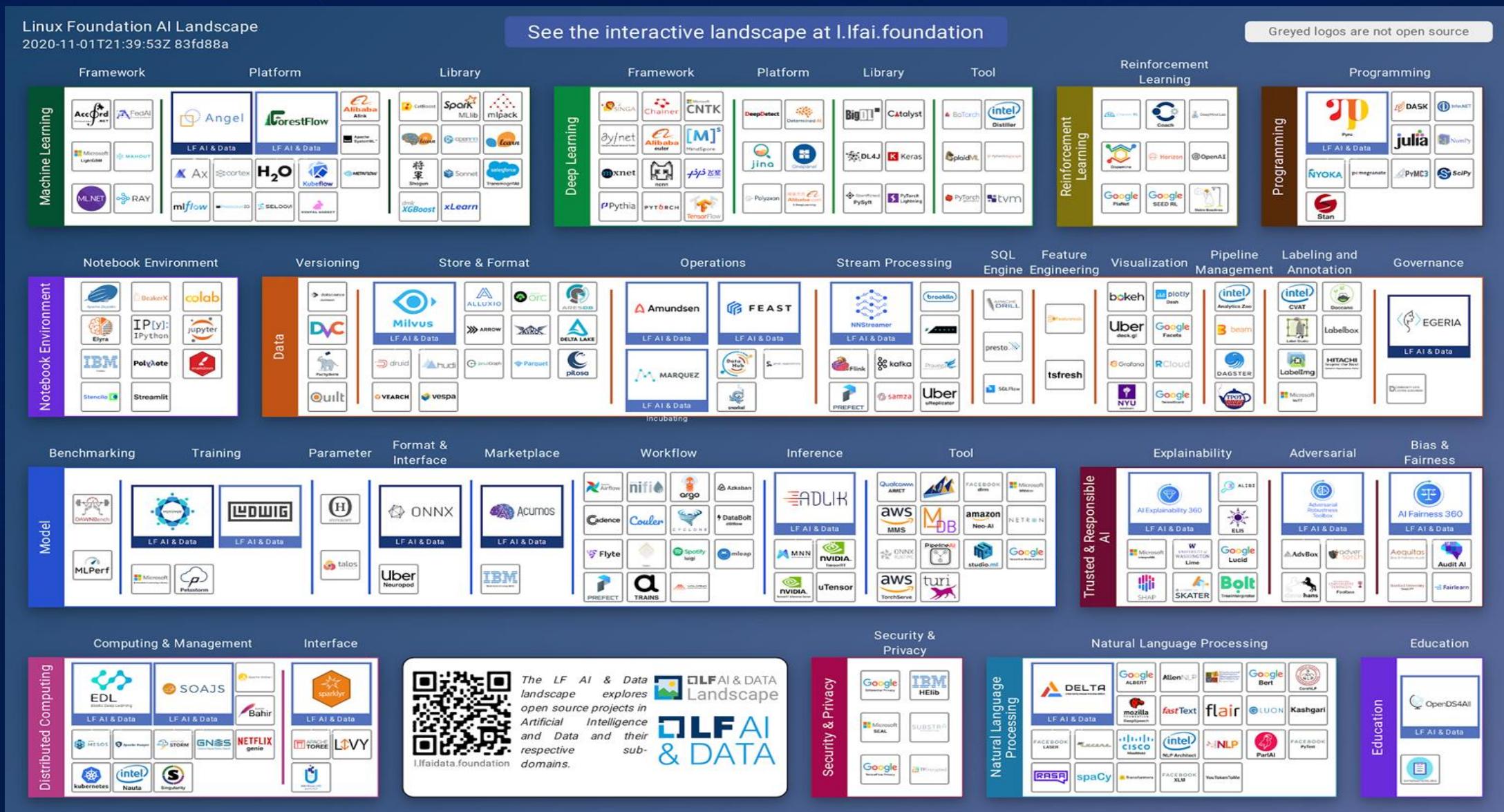
Microsoft Azure



Azure Machine Learning



ML Landscape



Flexible End-to-End ML Platform

| Best practices | Libraries | Binaries | Components | Pipelines |
|--|--|--|---|---|
| ML is hard; doing it well is even harder | "Vanilla" as well as Data Parallel libraries | Optimized packaging of libraries with APIs | Pluggable binaries / UIs conforming to Artifacts APIs | Task-driven and data-driven pipelines of components |

Platform

Seamlessly Integrated
Interoperable

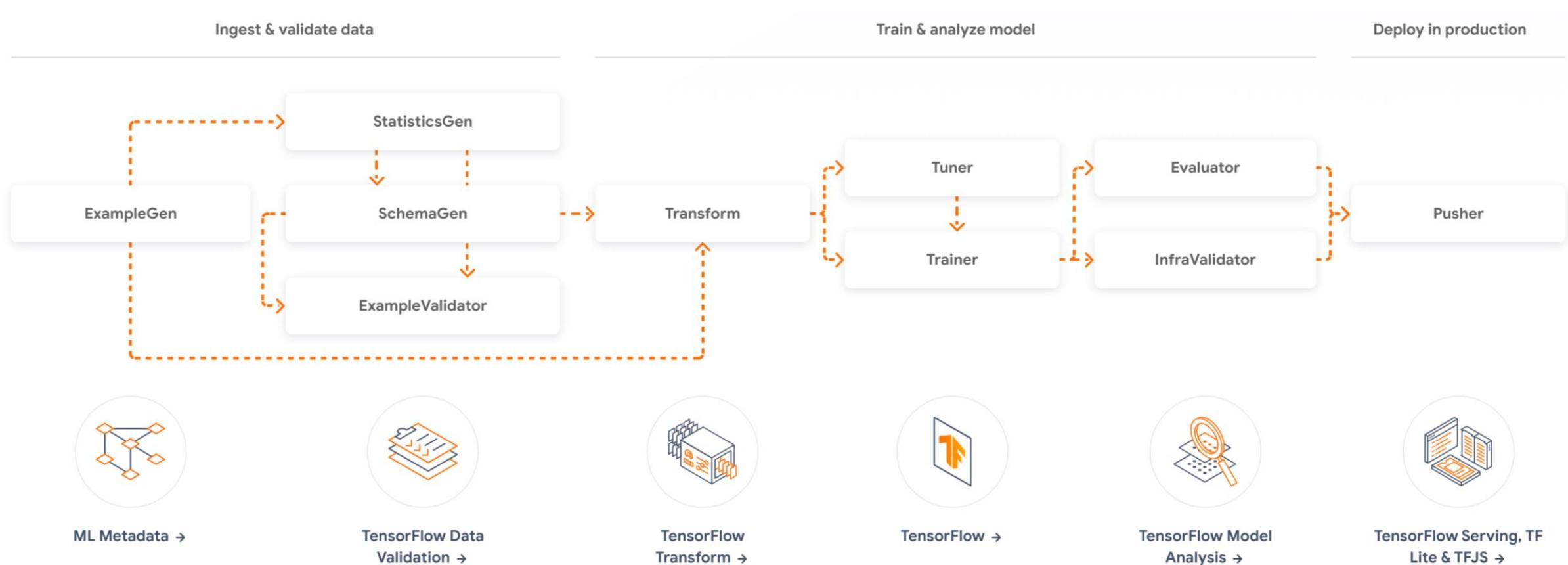
Extensible
Evolving

Supported
Service-able and Service-ed



TensorFlow Extended (TFX)

Flexible End-to-End ML Platform





Define problem

Construct and prepare data

Build and train model

Evaluate model

Deploy and monitor

Who is my ML system for?

The way actual users experience your system is essential to assessing the true impact of its predictions, recommendations, and decisions. Make sure to get input from a diverse set of users early on in your development process.





Step 1

Define problem

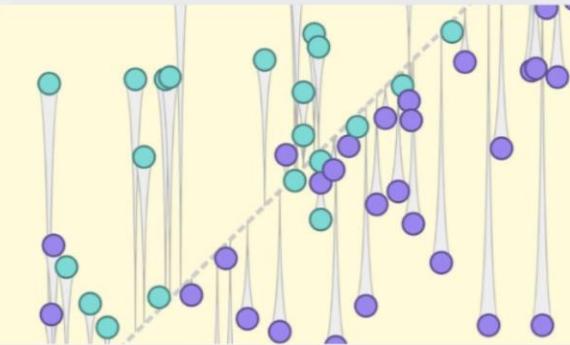
Use the following resources to design models with Responsible AI in mind.

People + AI Guidebook

People + AI Research (PAIR) Guidebook

Learn more about the AI development process and key considerations.

[Learn more ↗](#)



PAIR Explorables

Explore, via interactive visualizations, key questions and concepts in the realm of Responsible AI.

[Learn more ↗](#)

Define problem

Construct and prepare data

Build and train model

Evaluate model

Deploy and monitor

Am I using a representative dataset?

Is your data sampled in a way that represents your users (e.g. will be used for all ages, but you only have training data from senior citizens) and the real-world setting (e.g. will be used year-round, but you only have training data from the summer)?

Is there real-world/human bias in my data?

Underlying biases in data can contribute to complex feedback loops that reinforce existing stereotypes.

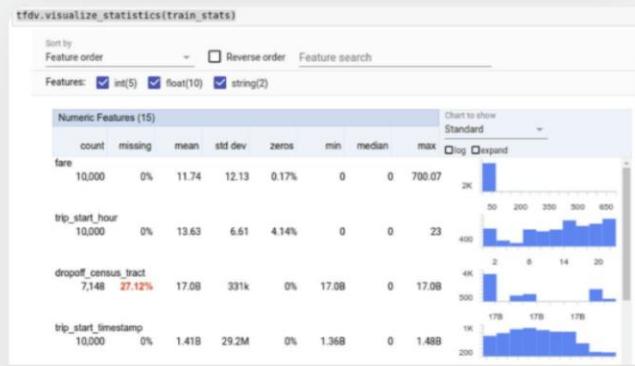




Step 2

Construct and prepare data

Use the following tools to examine data for potential biases.



TF Data Validation

Analyze and transform data to detect problems and engineer more effective feature sets.

[Learn more →](#)

Crowdsourced high-quality Colombian Spanish [es-co] multi-speaker speech dataset

research.google/tf2ai/datasets/colombian-spanish.tflite

PUBLISHED

Google LLC

INDUSTRY TYPE

Corporate - Tech

KEY APPLICATIONS

Machine Learning, Speech Technology

This dataset was created for speech research purposes and contains about 4,700 recordings of participants reading a script in Spanish as spoken in Colombia, one sentence at a time. Each example contains the audio files and the associated text. The audio is in raw format, stereo, 16 bit, 44.1 kHz, recorded with a Rode NT1 condenser microphone. The dataset is multi-speaker, containing recordings from 33 volunteers (male and female), where each volunteer contributed up to 100 recordings. The recordings took place in Bogota, Colombia in 2016.

PRIMARY DATA TYPE

Speech data

DATASET FUNCTION(S)

Training, Testing

INTENDED USE CASES
Multi-speaker and multi-lingual model speech synthesis
models building
Evaluating dialects effects on speech recognition models
Linguistic research

NATURE OF CONTENT

The dataset contains recordings of Spanish as spoken in Colombia in 2016. The participants read a script, approximately 100 sentences per participant. The script lines are delivered in audio files and the associated transcription of the audio. All the script lines are listed with the corresponding audio files in the file "scripts.txt". The file has two columns. The first column contains the FileID of the file, and second the column contains the text read in the corresponding audio file. The columns are tab separated.

EXAMPLE COMPONENTS

The file line, index 10, gives a transcription of each audio file in the format:

DESCRIPTIONS OF EXAMPLE COMPONENTS
.cof_3248_024 is the FileID of the file containing the text in the line. The FileID is composed of three parts, delimited by an underscore ".": The first part is unique for the dataset and

Data Cards

Create a transparency report for your dataset.

[Learn more →](#)



Define problem

Construct and prepare data

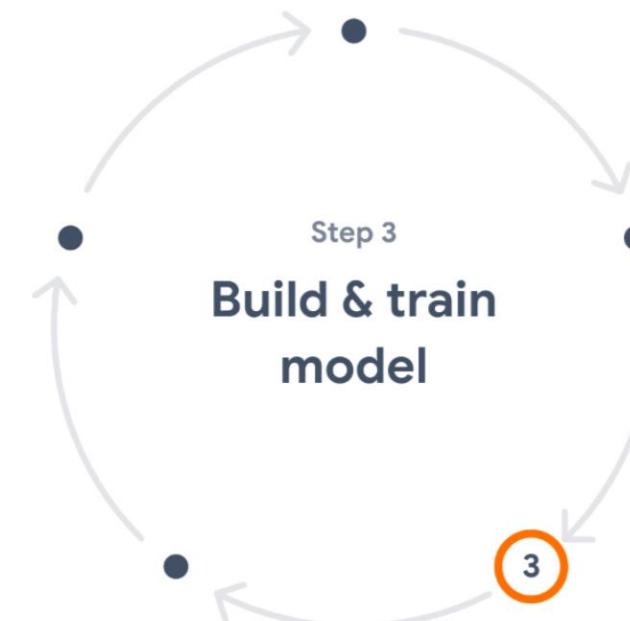
Build and train model

Evaluate model

Deploy and monitor

What methods should I use to train my model?

Use training methods that build fairness, interpretability, privacy, and security into the model.

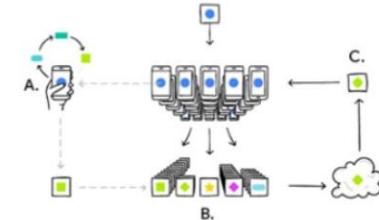




Step 3

Build and train model

Use the following tools to train models using privacy-preserving, interpretable techniques, and more.



TF Model Remediation

Train machine learning models to promote more equitable outcomes.

[Learn more →](#)

TF Privacy

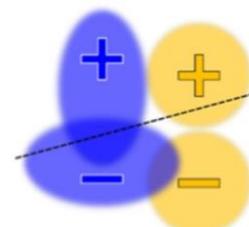
Train machine learning models with privacy.

[Learn more ↗](#)

TF Federated

Train machine learning models using federated learning techniques.

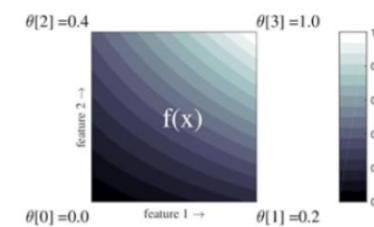
[Learn more →](#)



TF Constrained Optimization

Optimize inequality-constrained problems.

[Learn more ↗](#)



TF Lattice

Implement flexible, controlled, and interpretable lattice-based models.

[Learn more →](#)



Análise facial



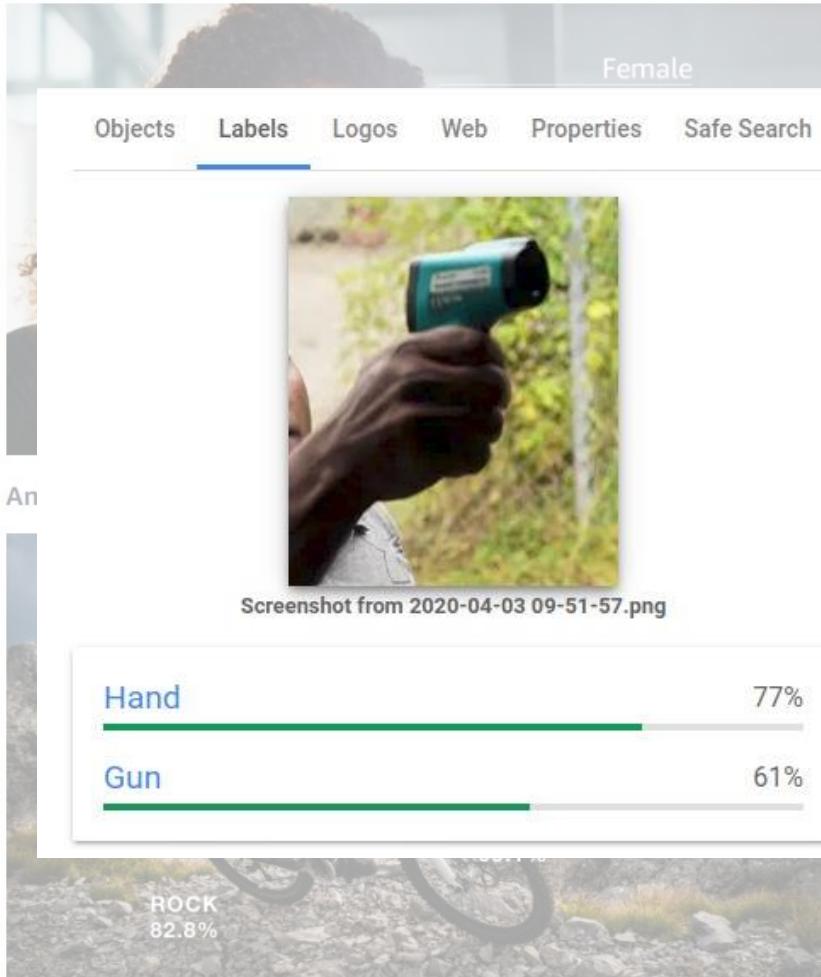
Determinação de caminhos



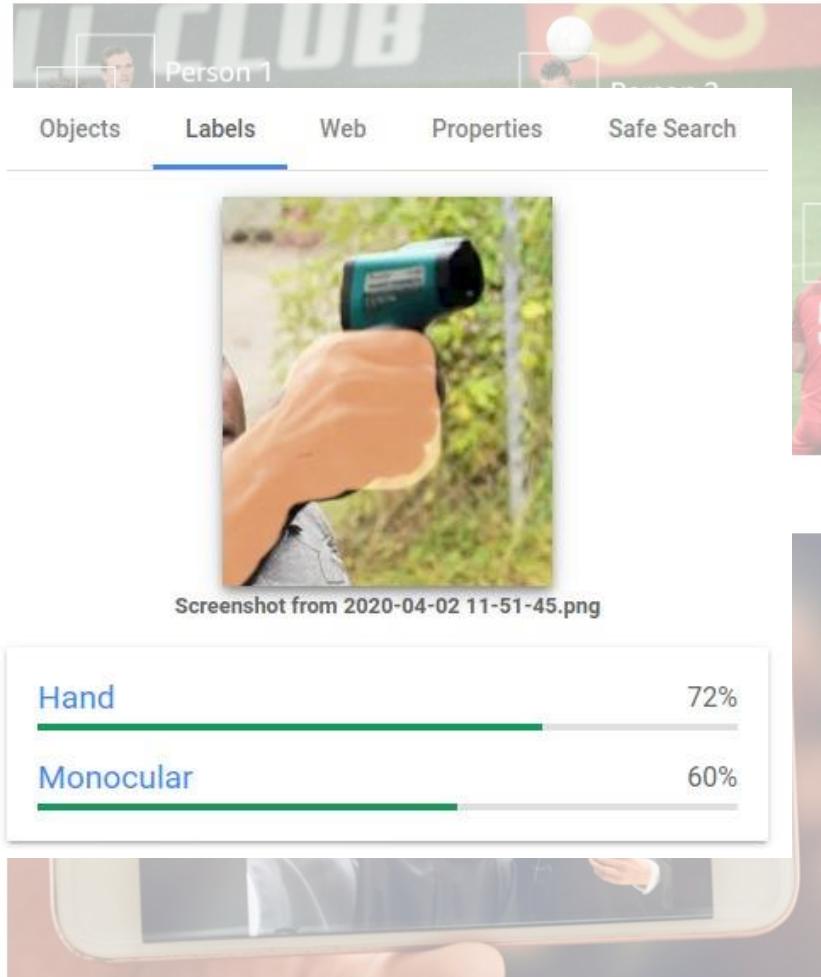
Detecção de objetos, cenas e atividades



Reconhecimento facial



Detectão de objetos, cenas e atividades



Reconhecimento facial

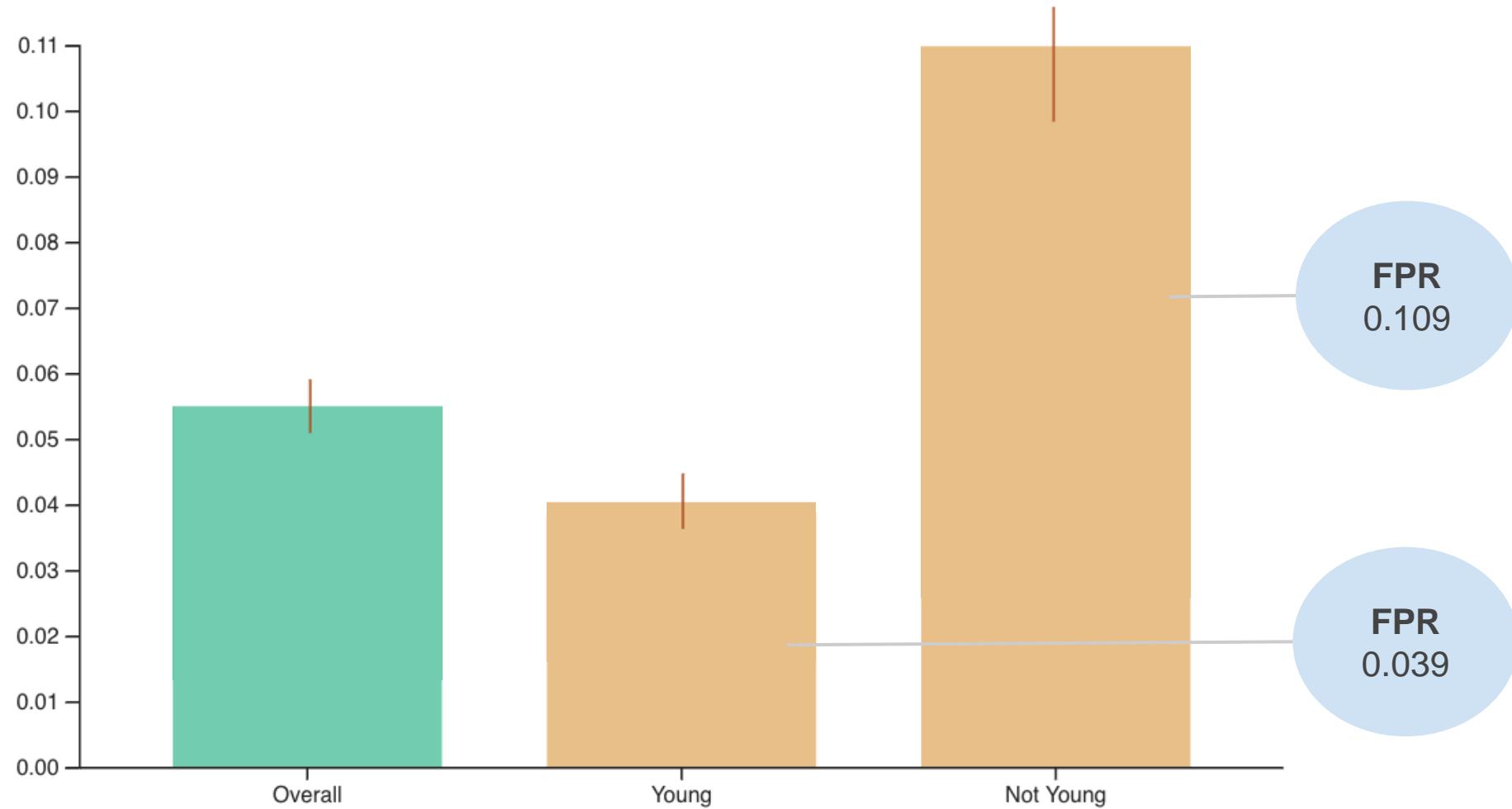
Remediation:

Smile Detection on CelebA using TF

Constrained Optimization



Results: Unconstrained `tf.keras.Sequential` model





Remediation: CelebA example using TFCO

1. Define subsets of interest

```
context = tfco.rate_context(predictions, labels=lambda:labels_tensor)
context_subset = context.subset(lambda:groups_tensor < 1)
```



Remediation: CelebA example using TFCO

1. Define subsets of interest

```
context = tfco.rate_context(predictions, labels=lambda:labels_tensor)
context_subset = context.subset(lambda:groups_tensor < 1)
```

2. Set constraints on subset using rate helpers

```
constraints = [tfco.false_positive_rate(context_subset) <= 0.05]
```



Remediation: CelebA example using TFCO*

1. Define subsets of interest

```
context = tfco.rate_context(predictions, labels=lambda:labels_tensor)
context_subset = context.subset(lambda:groups_tensor < 1)
```

2. Set constraints on subset using rate helpers

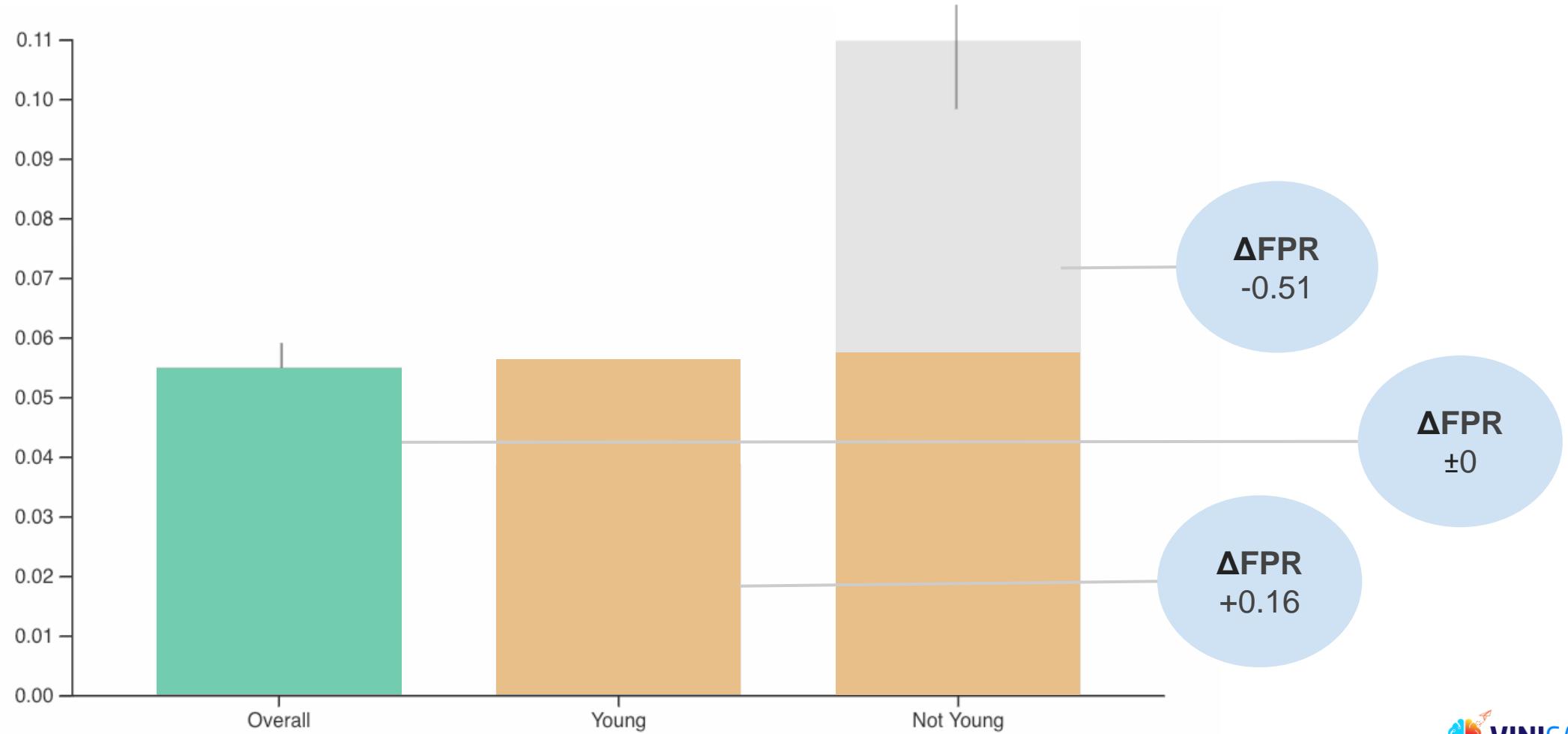
```
constraints = [tfco.false_positive_rate(context_subset) <= 0.05]
```

3. Define optimizer and train

```
problem = tfco.RateMinimizationProblem(tfco.error_rate(context), constraints)
optimizer = tfco.ProxyLagrangianOptimizerV2(
    optimizer=tf.keras.optimizers.Adam(learning_rate=0.001),
    constraint_optimizer=tf.keras.optimizers.Adam(learning_rate=0.001),
    num_constraints=problem.num_constraints)
```

https://github.com/google-research/tensorflow_constrained_optimization

Results: Constrained `tf.keras.Sequential` model





Define problem

Construct and prepare data

Build and train model

Evaluate model

Deploy and monitor

How is my model performing?

Evaluate user experience in real-world scenarios across a broad spectrum of users, use cases, and contexts of use. Test and iterate in dogfood first, followed by continued testing after launch.

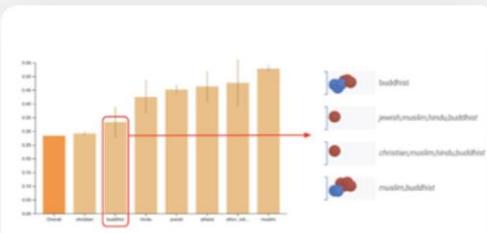




Step 4

Evaluate model

Debug, evaluate, and visualize model performance using the following tools.



Fairness Indicators

Evaluate commonly-identified fairness metrics for binary and multi-class classifiers.

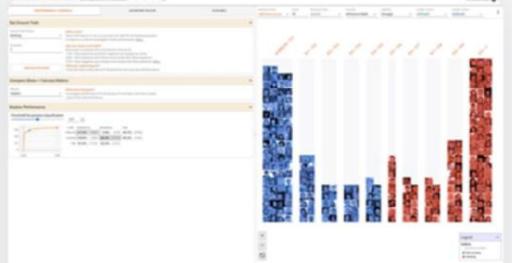
[Learn more →](#)



TF Model Analysis

Evaluate models in a distributed manner and compute over different slices of data.

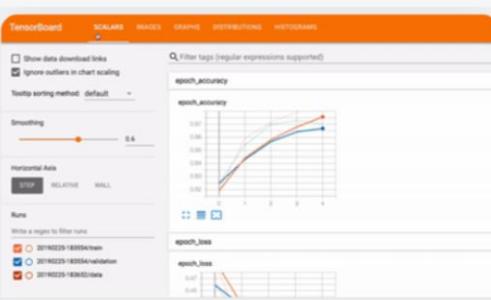
[Learn more →](#)



What-If Tool

Examine, evaluate, and compare machine learning models.

[Learn more →](#)



TensorBoard

Measure and visualize the machine learning workflow.

[Learn more →](#)



Language Interpretability Tool

Visualize and understand NLP models.

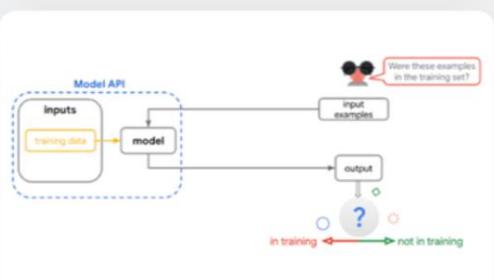
[Learn more →](#)



Explainable AI

Develop interpretable and inclusive machine learning models.

[Learn more →](#)



TF Privacy Tests

Assess the privacy properties of classification models.

[Learn more →](#)

Define problem

Construct and prepare data

Build and train model

Evaluate model

Deploy and monitor

Are there complex feedback loops?

Even if everything in the overall system design is carefully crafted, ML-based models rarely operate with 100% perfection when applied to real, live data. When an issue occurs in a live product, consider whether it aligns with any existing societal disadvantages, and how it will be impacted by both short- and long-term solutions.

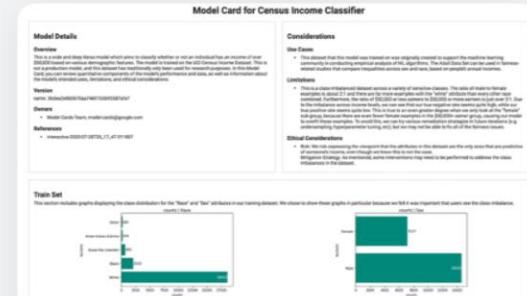




Step 5

Deploy and monitor

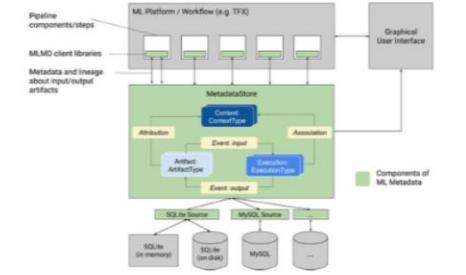
Use the following tools to track and communicate about model context and details.



Model Card Toolkit

Generate model cards with ease using the Model Card toolkit.

[Learn more →](#)



ML Metadata

Record and retrieve metadata associated with ML developer and data scientist workflows.

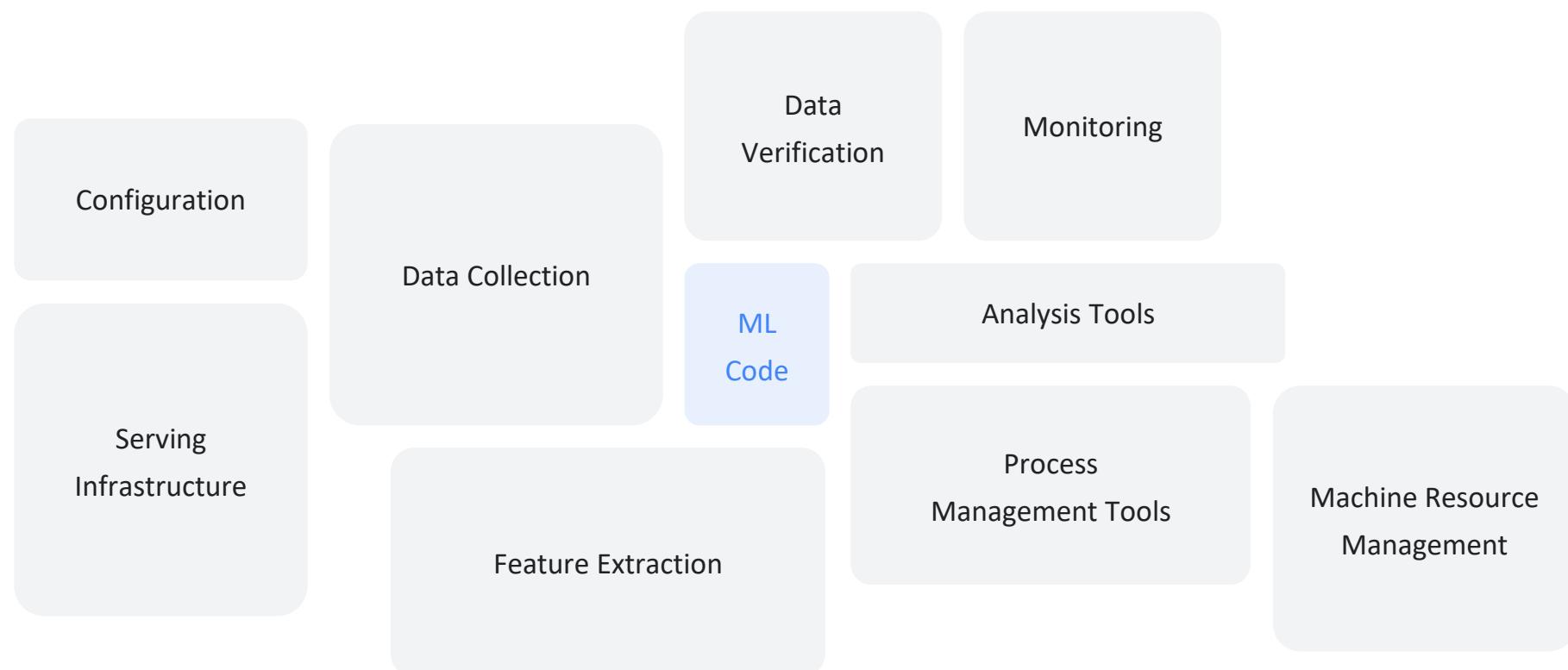
[Learn more →](#)



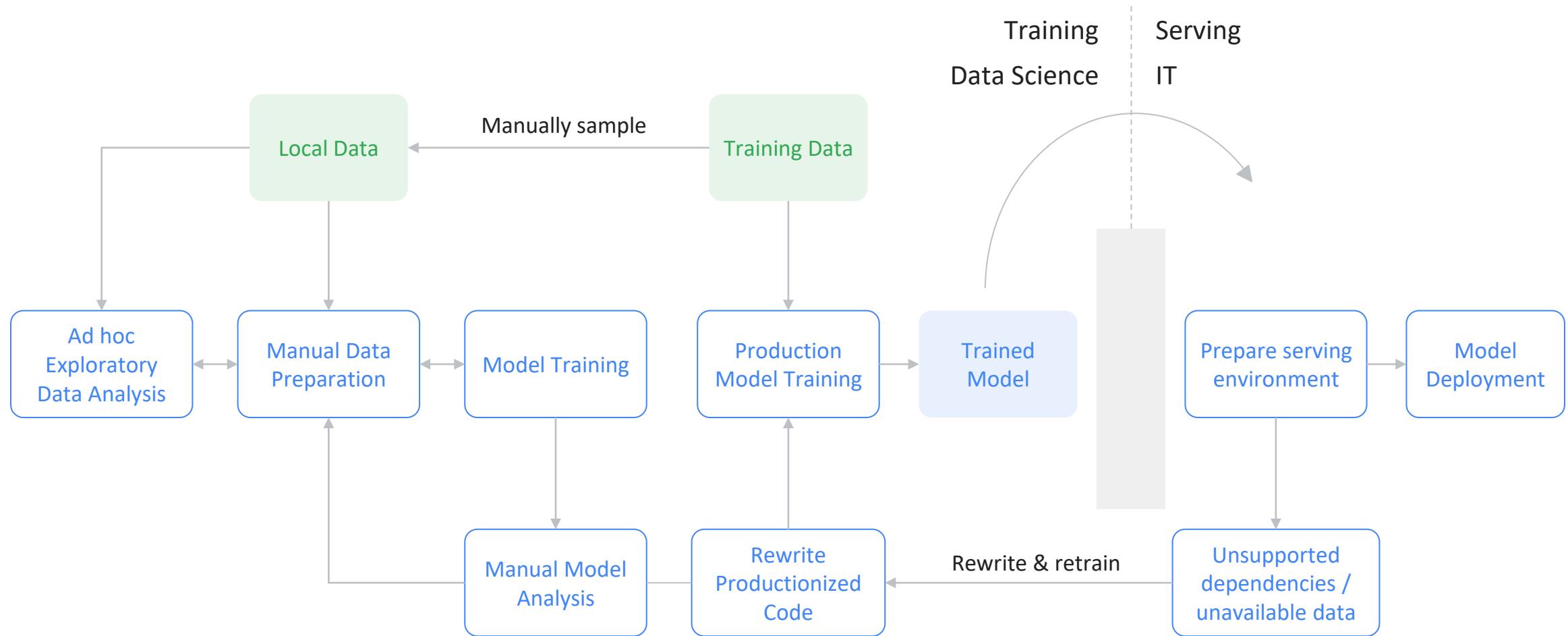
Model Cards

Organize the essential facts of machine learning in a structured way.

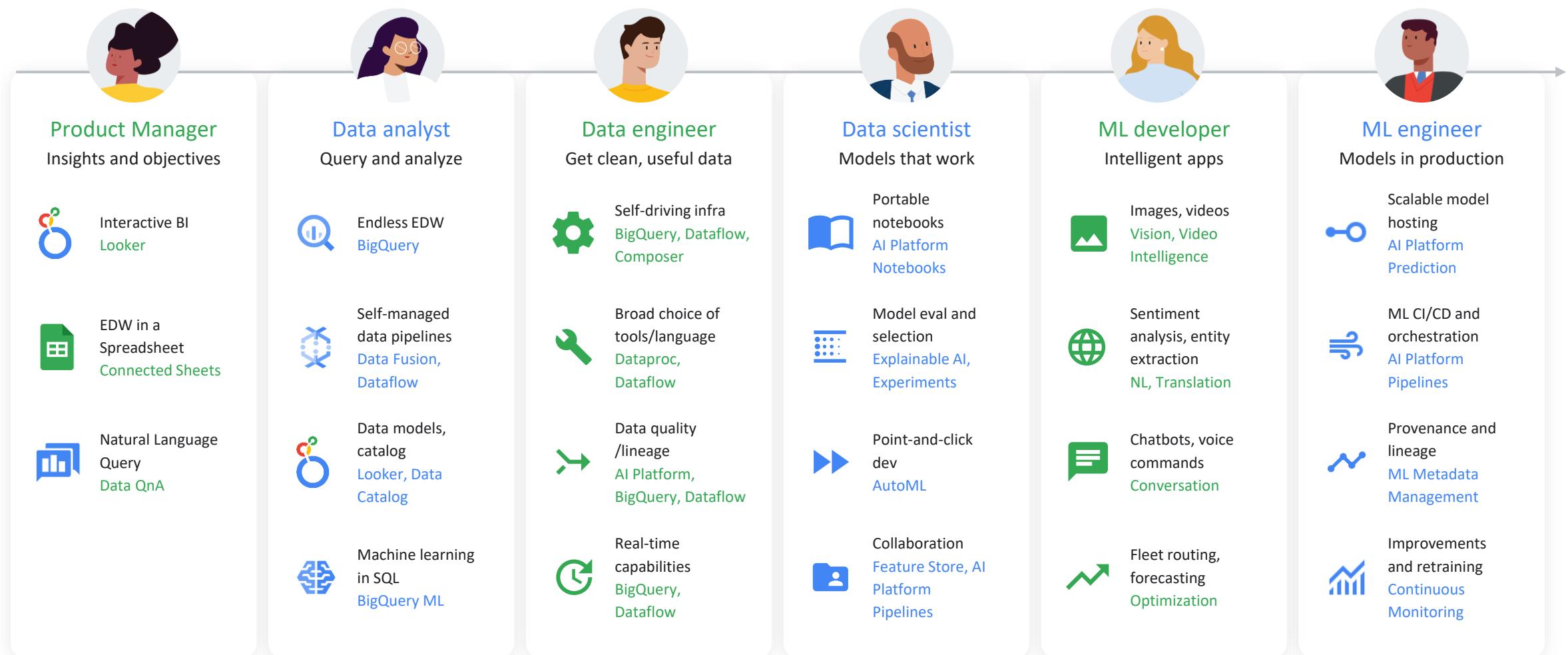
[Learn more →](#)



What's happening today: Data Science and IT (Ops) are isolated



A platform for all users and intents throughout the ML lifecycle



Como conectar usuários para produzir resultados **10x** mais impactantes?



ML Engineer



Data Engineer



Developer



Data Scientist



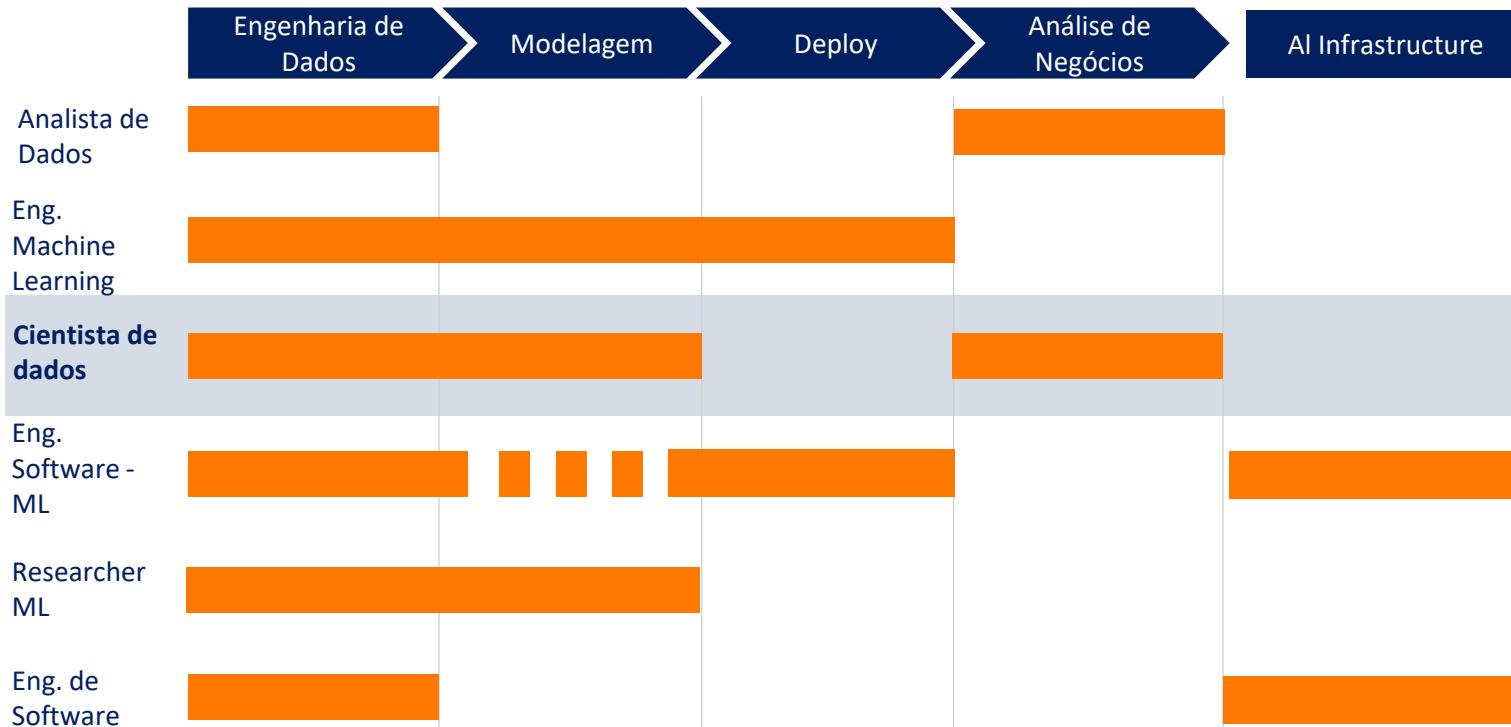
Business Analyst



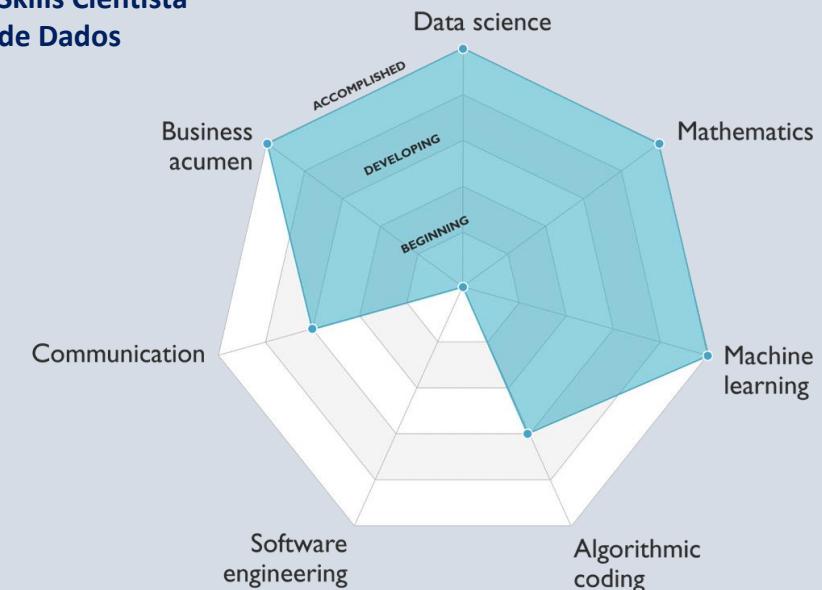
End User

Perfil do Cientista de dados

Frentes de trabalho x profissionais , de acordo com texto publicado no Workera:



Skills Cientista de Dados



Profissionais x Skills

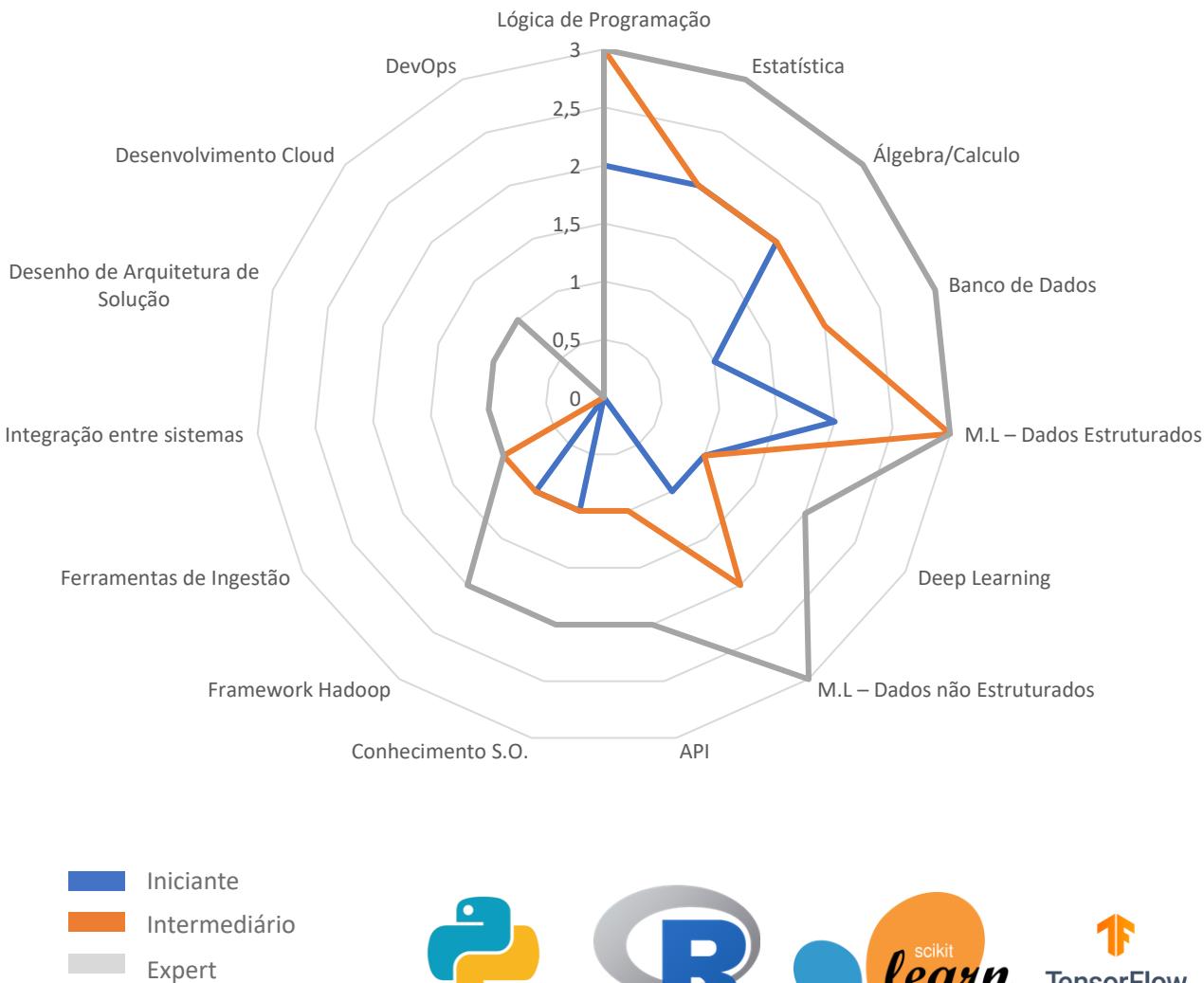
| | Data Science | Mathematics | Machine Learning | Algorithmic coding | Software eng. | Communication | Business acumen |
|-----------------------|--------------|-------------|------------------|--------------------|---------------|---------------|-----------------|
| Analista de Dados | ✓ | | | ✓ | | ✓ | ✓ |
| Eng. Machine Learning | ✓ | | ✓ | ✓ | ✓ | | |
| Eng. Software - ML | | | ✓ | ✓ | ✓ | | |
| Researcher ML | | ✓ | ✓ | ✓ | ✓ | | |
| Eng. de Software | | | | ✓ | ✓ | ✓ | |

O Perfil do Cientista de dados

Para cientista de dados é realizado uma prova com duração de 3 horas e nota de corte 5 e após aprovação na prova existe um entrevista técnica (sabatina), e após aprovação em ambas existe a seleção de perfil por vaga (gestor funcional).

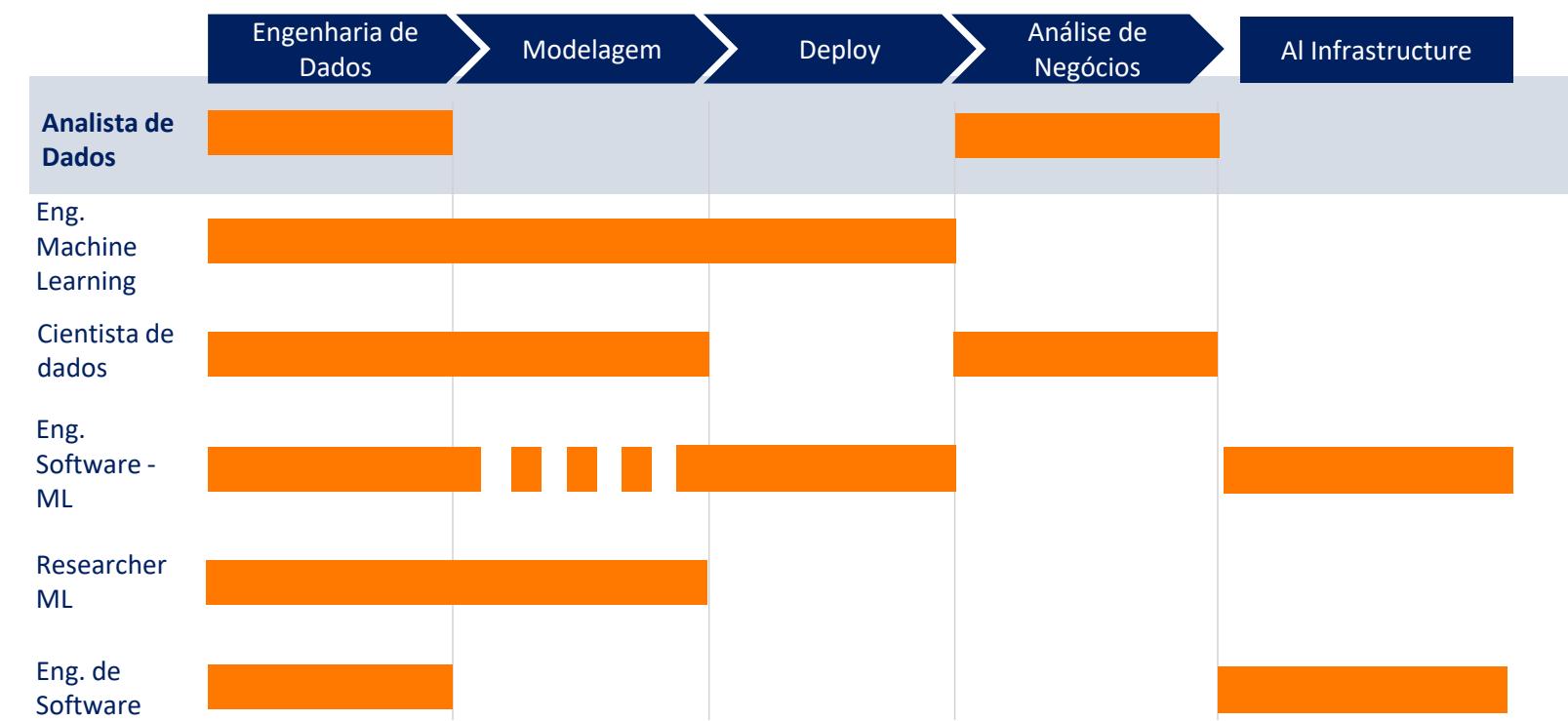
PROCESSOS

| | Iniciante | Intermed. | Expert |
|---|-----------|-----------|--------|
| Conhecimento | | | |
| Conhecimento Técnicos | | | |
| Lógica de Programação | 2 | 3 | 3 |
| Estatística | 2 | 2 | 3 |
| Álgebra/Calculo | 2 | 2 | 3 |
| Banco de Dados | 1 | 2 | 3 |
| Machine Learning – Dados Estruturados | 2 | 3 | 3 |
| Deep Learning | 1 | 1 | 2 |
| Machine Learning – Dados não Estruturados | 1 | 2 | 3 |
| Construção e Funcionamento de API | - | 1 | 2 |
| Conhecimento S.O. | 1 | 1 | 2 |
| Framework Hadoop | 1 | 1 | 2 |
| Ferramentas de Ingestão | - | 1 | 1 |
| Integração entre sistemas | - | - | 1 |
| Desenho de Arquitetura de Solução | - | - | 1 |
| Desenvolvimento Cloud | - | - | 1 |
| DevOps | - | - | - |
| Habilidades | | | |
| Solução de Problemas | 2 | 2 | 3 |
| Comunicação | 1 | 2 | 3 |
| Transformação/Inovação | 1 | 2 | 2 |
| Visão Estratégica e Sistêmica | 1 | 2 | 3 |
| Influência | - | - | 1 |
| Tutoria | - | 2 | 3 |

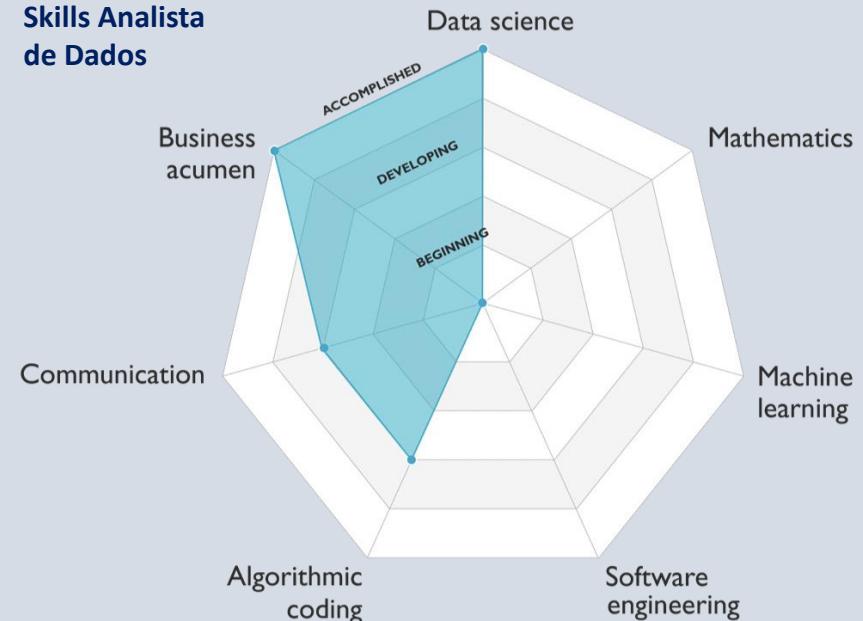


Perfil do Analista de dados

Frentes de trabalho x profissionais , de acordo com texto publicado no Workera:



Skills Analista de Dados



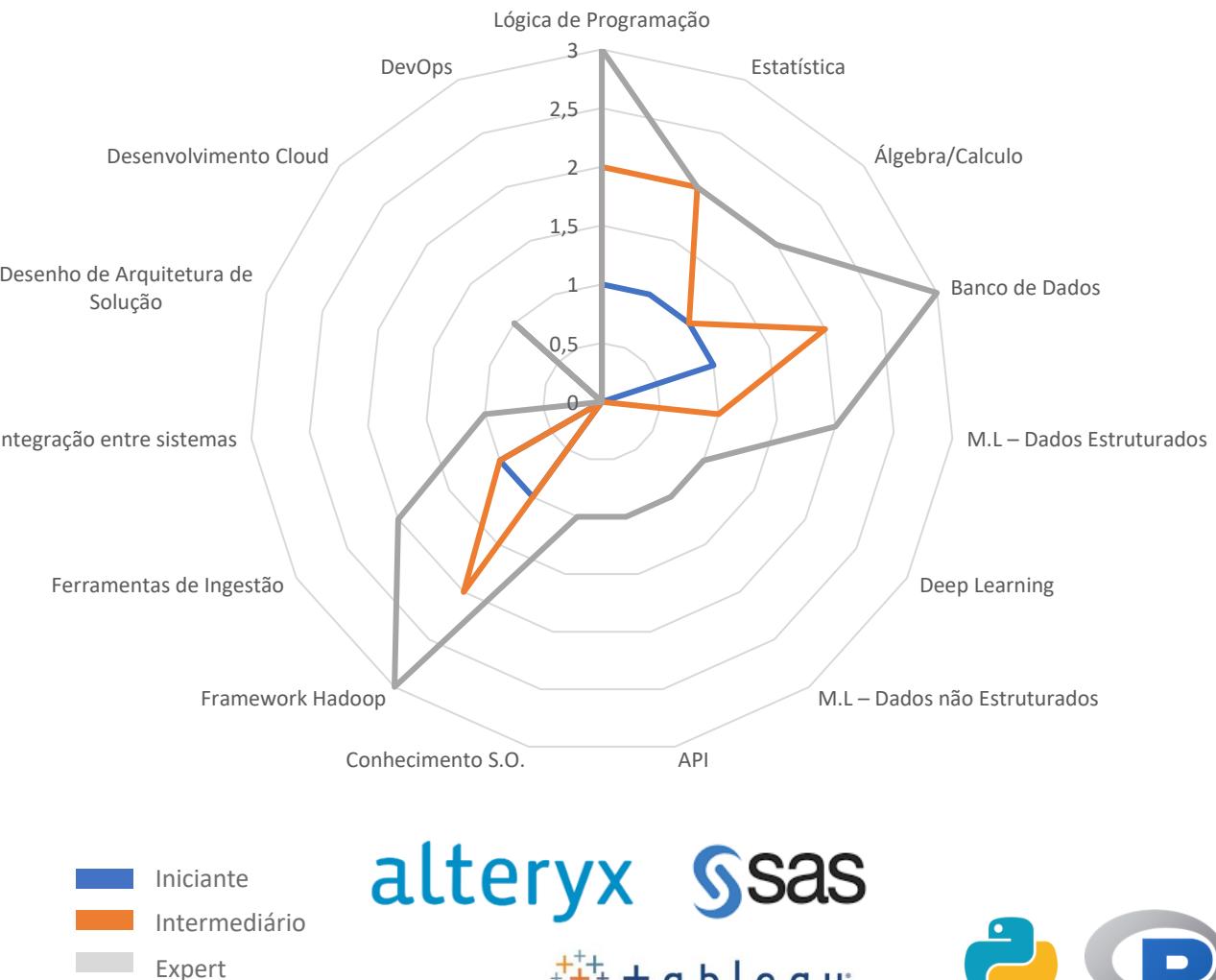
Profissionais x Skills

| | Data Science | Mathematics | Machine Learning | Algorithmic coding | Software eng. | Communication | Business acumen |
|-----------------------|--------------|-------------|------------------|--------------------|---------------|---------------|-----------------|
| Analista de Dados | ✓ | | | ✓ | | ✓ | ✓ |
| Eng. Machine Learning | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | |
| Eng. Software - ML | | | ✓ | ✓ | ✓ | ✓ | |
| Researcher ML | | ✓ | ✓ | ✓ | ✓ | | |
| Eng. de Software | | | | ✓ | ✓ | ✓ | |

O Perfil do Analista de Dados

O perfil é selecionado em entrevista técnica, análise de portfólio (Github), e seleção de perfil (gestor funcional).

| PROCESSOS | Iniciante | Intermed. | Expert |
|---|-----------|-----------|--------|
| Conhecimento | | | |
| Conhecimento Técnicos | | | |
| Lógica de Programação | 1 | 2 | 3 |
| Estatística | 1 | 2 | 2 |
| Álgebra/Calculo | 1 | 1 | 2 |
| Banco de Dados | 1 | 2 | 3 |
| Machine Learning – Dados Estruturados | - | 1 | 2 |
| Deep Learning | - | - | 1 |
| Machine Learning – Dados não Estruturados | - | - | 1 |
| Construção e Funcionamento de API | - | - | 1 |
| Conhecimento S.O. | - | - | 1 |
| Framework Hadoop | 1 | 2 | 3 |
| Ferramentas de Ingestão | 1 | 1 | 2 |
| Integração entre sistemas | - | - | 1 |
| Desenho de Arquitetura de Solução | - | - | - |
| Desenvolvimento Cloud | - | - | 1 |
| DevOps | - | - | - |
| Habilidades | | | |
| Solução de Problemas | 1 | 2 | 3 |
| Comunicação | 1 | 2 | 3 |
| Transformação/Inovação | 1 | 2 | 2 |
| Visão Estratégica e Sistêmica | 1 | 2 | 3 |
| Influência | - | - | 1 |
| Tutoria | - | 2 | 3 |

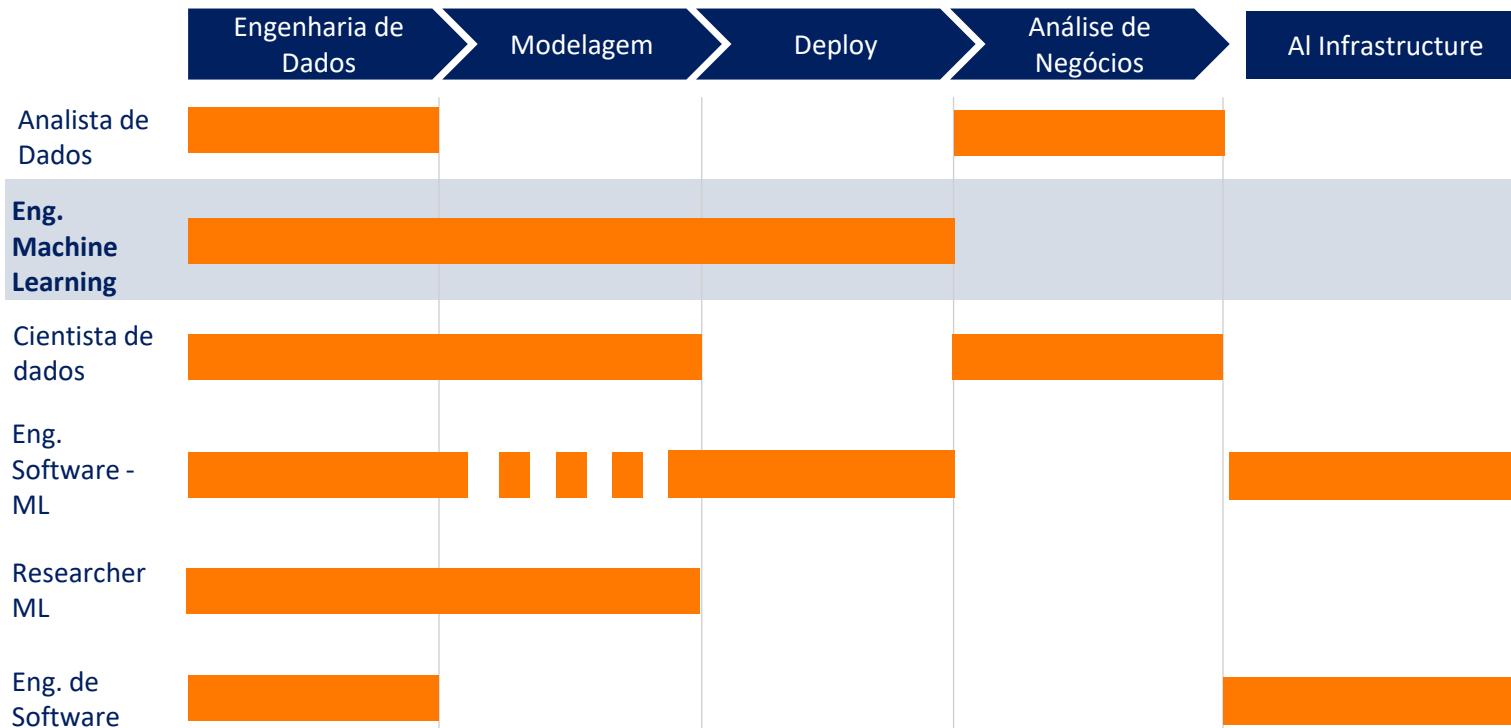


alteryx sas

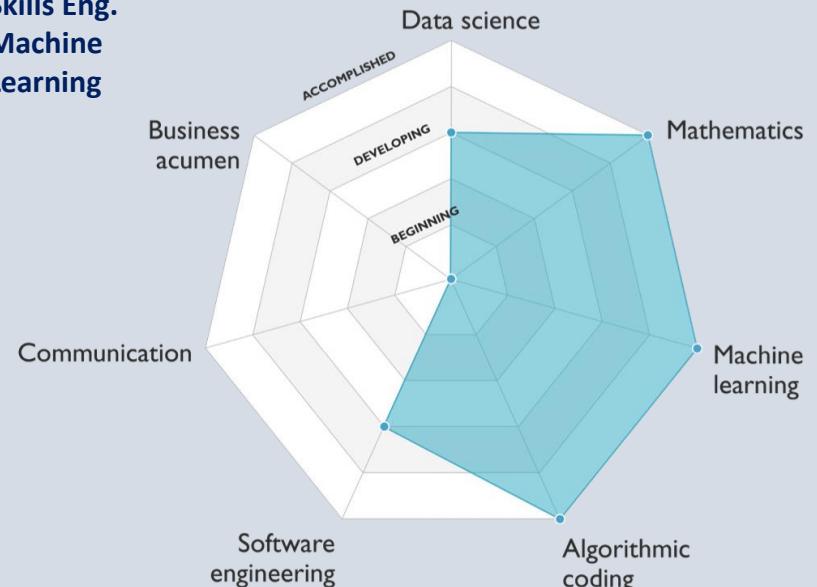


Perfil do Eng. De Machine Learning

Frentes de trabalho x profissionais , de acordo com texto publicado no Workera:



Skills Eng.
Machine
Learning



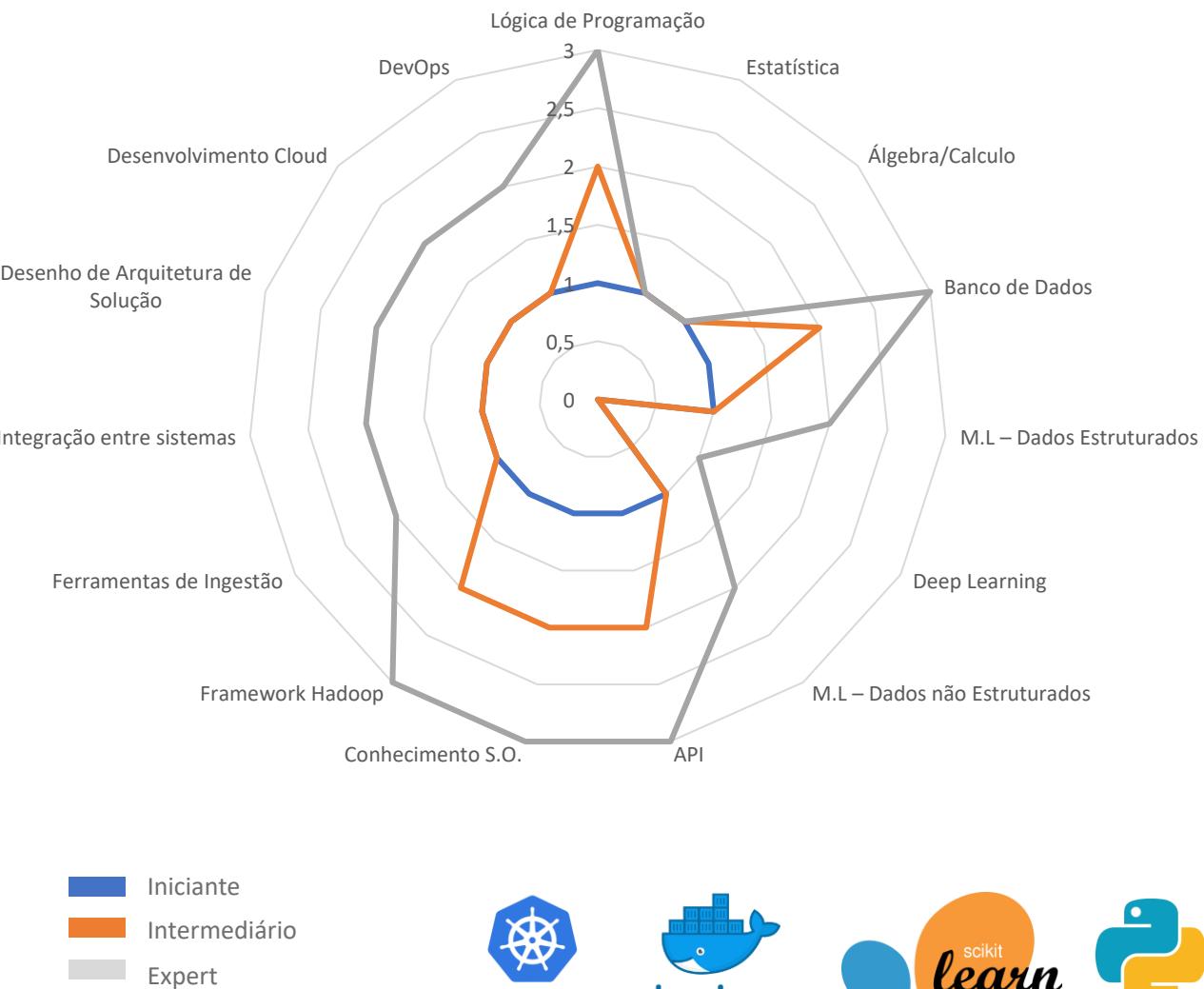
Profissionais x Skills

| | Data Science | Mathematics | Machine Learning | Algorithmic coding | Software eng. | Communication | Business acumen |
|-----------------------|--------------|-------------|------------------|--------------------|---------------|---------------|-----------------|
| Analista de Dados | ✓ | | | ✓ | | ✓ | ✓ |
| Eng. Machine Learning | ✓ | ✓ | ✓ | ✓ | ✓ | | |
| Eng. Software - ML | | | ✓ | ✓ | ✓ | | |
| Researcher ML | | ✓ | ✓ | ✓ | ✓ | | Research |
| Eng. de Software | | | | ✓ | ✓ | | |

O Perfil do Engenheiro de Machine Learning

O perfil é selecionado em entrevista técnica, análise de portfólio (Github), e seleção de perfil (gestor funcional).

| PROCESSOS | Iniciante | Intermed. | Expert |
|---|-----------|-----------|--------|
| Conhecimento | | | |
| Conhecimento Técnicos | | | |
| Lógica de Programação | 1 | 2 | 3 |
| Estatística | 1 | 1 | 1 |
| Álgebra/Calculo | 1 | 1 | 1 |
| Banco de Dados | 1 | 2 | 3 |
| Machine Learning – Dados Estruturados | 1 | 1 | 2 |
| Deep Learning | - | - | 1 |
| Machine Learning – Dados não Estruturados | 1 | 1 | 2 |
| Construção e Funcionamento de API | 1 | 2 | 3 |
| Conhecimento S.O. | 1 | 2 | 3 |
| Framework Hadoop | 1 | 2 | 3 |
| Ferramentas de Ingestão | 1 | 1 | 2 |
| Integração entre sistemas | 1 | 1 | 2 |
| Desenho de Arquitetura de Solução | 1 | 1 | 2 |
| Desenvolvimento Cloud | 1 | 1 | 2 |
| DevOps | 1 | 1 | 2 |
| Habilidades | | | |
| Solução de Problemas | 1 | 2 | 3 |
| Comunicação | 1 | 2 | 3 |
| Transformação/Inovação | 1 | 2 | 2 |
| Visão Estratégica e Sistêmica | 1 | 2 | 3 |
| Influência | - | - | 1 |
| Tutoria | - | 2 | 3 |

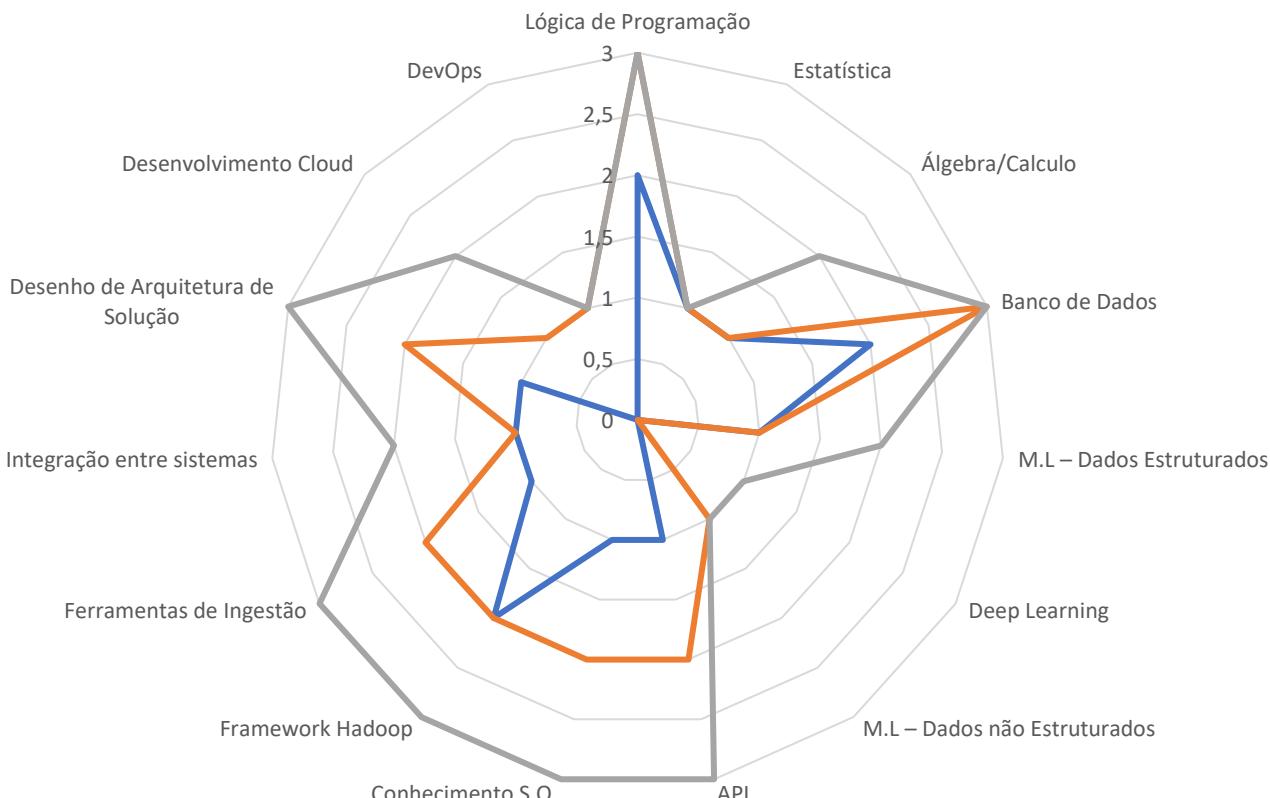


O Perfil do Engenheiro de Dados

O perfil é selecionado em entrevista técnica, análise de portfólio (Github), e seleção de perfil (gestor funcional).

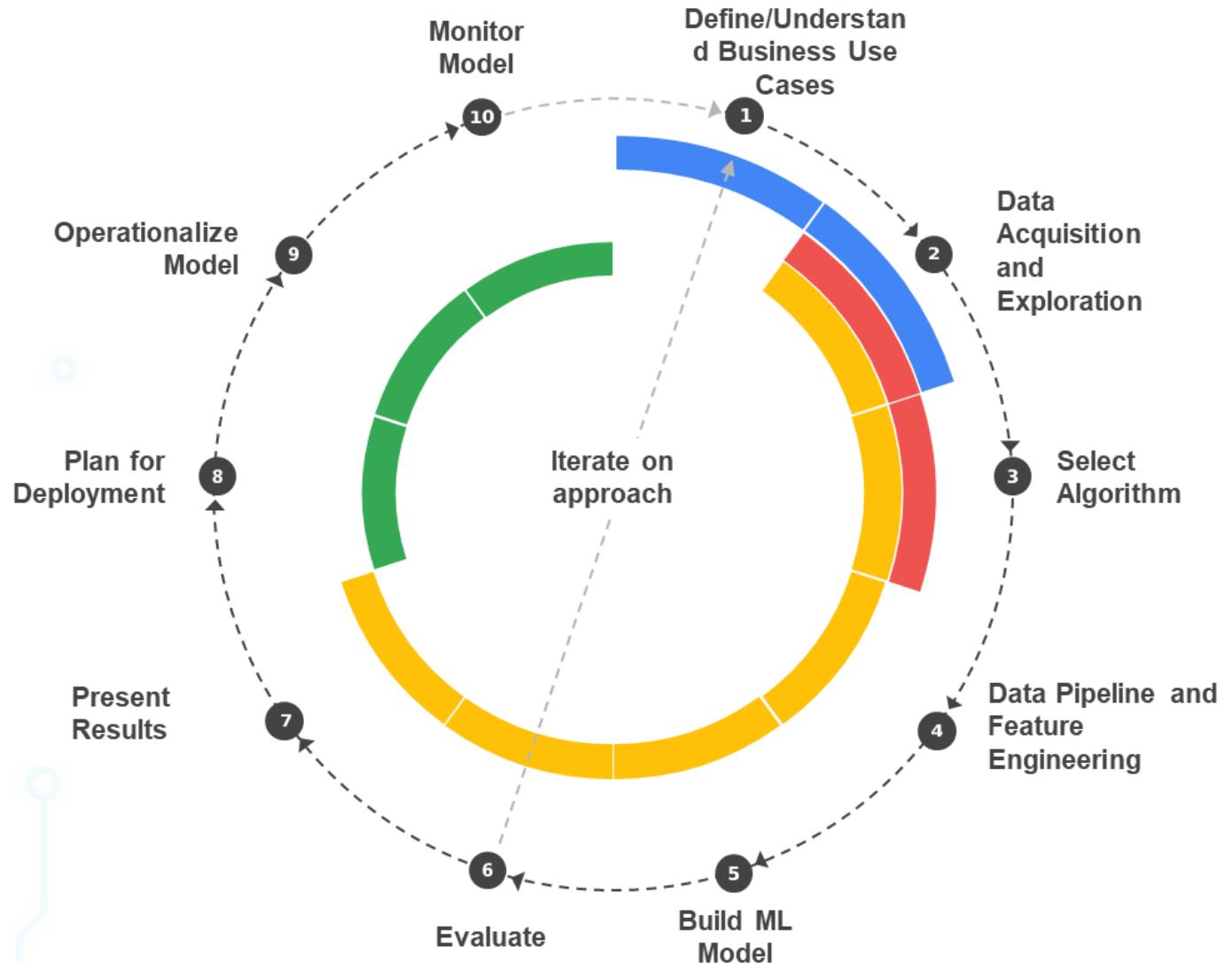
PROCESSOS

| | Iniciante | Intermed. | Expert |
|---|-----------|-----------|--------|
| Conhecimento | | | |
| Conhecimento Técnicos | | | |
| Lógica de Programação | 2 | 3 | 3 |
| Estatística | 1 | 1 | 1 |
| Álgebra/Calculo | 1 | 1 | 2 |
| Banco de Dados | 2 | 3 | 3 |
| Machine Learning – Dados Estruturados | 1 | 1 | 2 |
| Deep Learning | - | - | 1 |
| Machine Learning – Dados não Estruturados | - | 1 | 1 |
| Construção e Funcionamento de API | 1 | 2 | 3 |
| Conhecimento S.O. | 1 | 2 | 3 |
| Framework Hadoop | 2 | 2 | 3 |
| Ferramentas de Ingestão | 1 | 2 | 3 |
| Integração entre sistemas | 1 | 1 | 2 |
| Desenho de Arquitetura de Solução | 1 | 2 | 3 |
| Desenvolvimento Cloud | - | 1 | 2 |
| DevOps | - | 1 | 1 |
| Habilidades | | | |
| Solução de Problemas | 1 | 2 | 3 |
| Comunicação | 2 | 2 | 3 |
| Transformação/Inovação | 1 | 2 | 2 |
| Visão Estratégica e Sistêmica | 1 | 2 | 3 |
| Influência | - | 1 | 2 |
| Tutoria | - | 2 | 3 |



- █ Iniciante
- █ Intermediário
- █ Expert





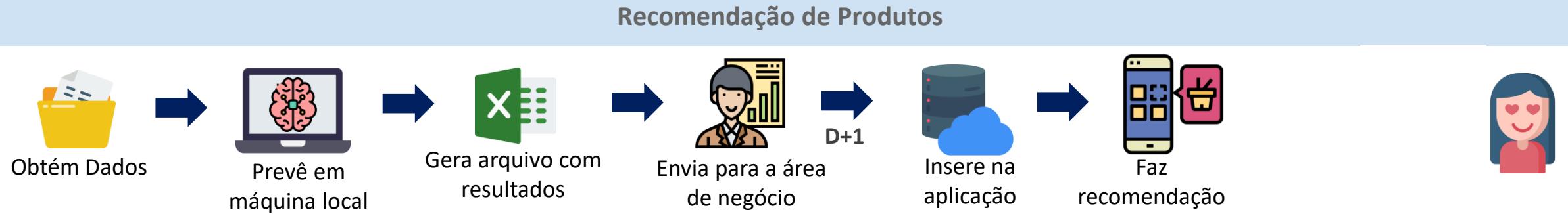
Modelos em Produção

Batch

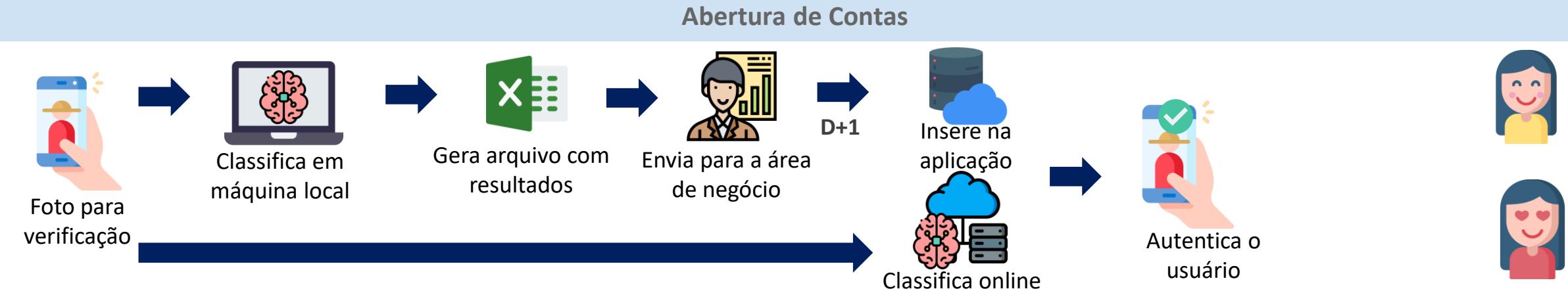


Modelos em Produção

Batch

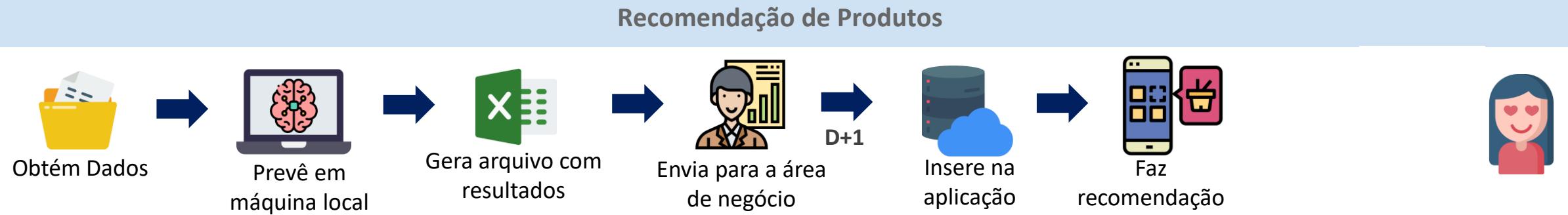


Batch ou Online

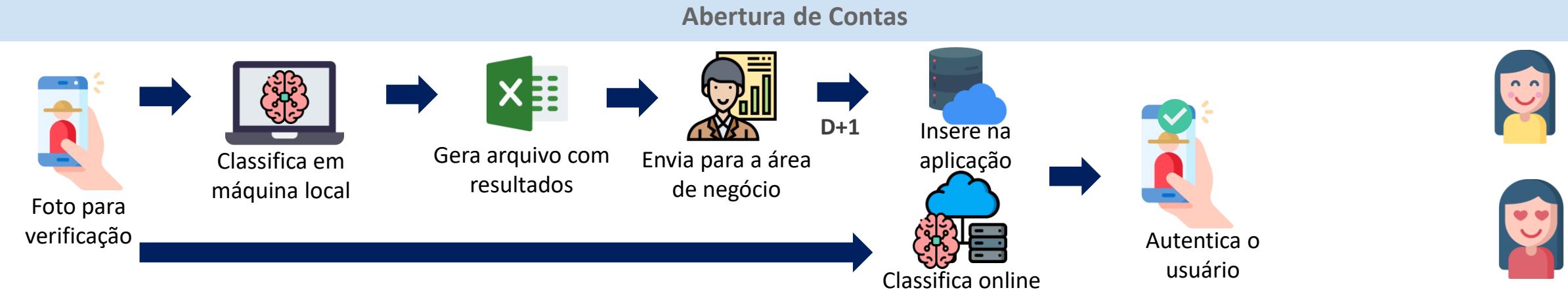


Modelos em Produção

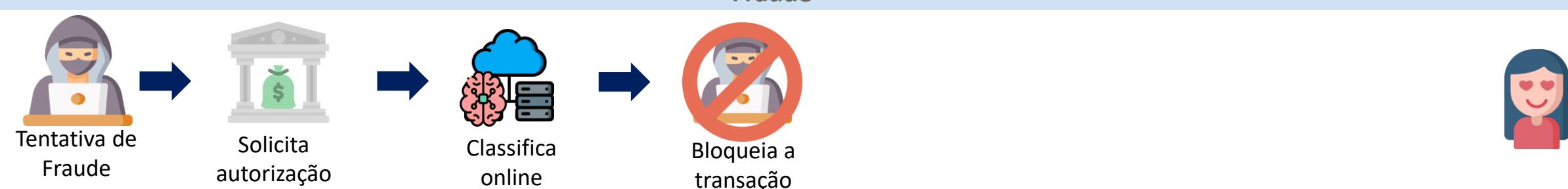
Batch



Batch ou Online



Online



MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

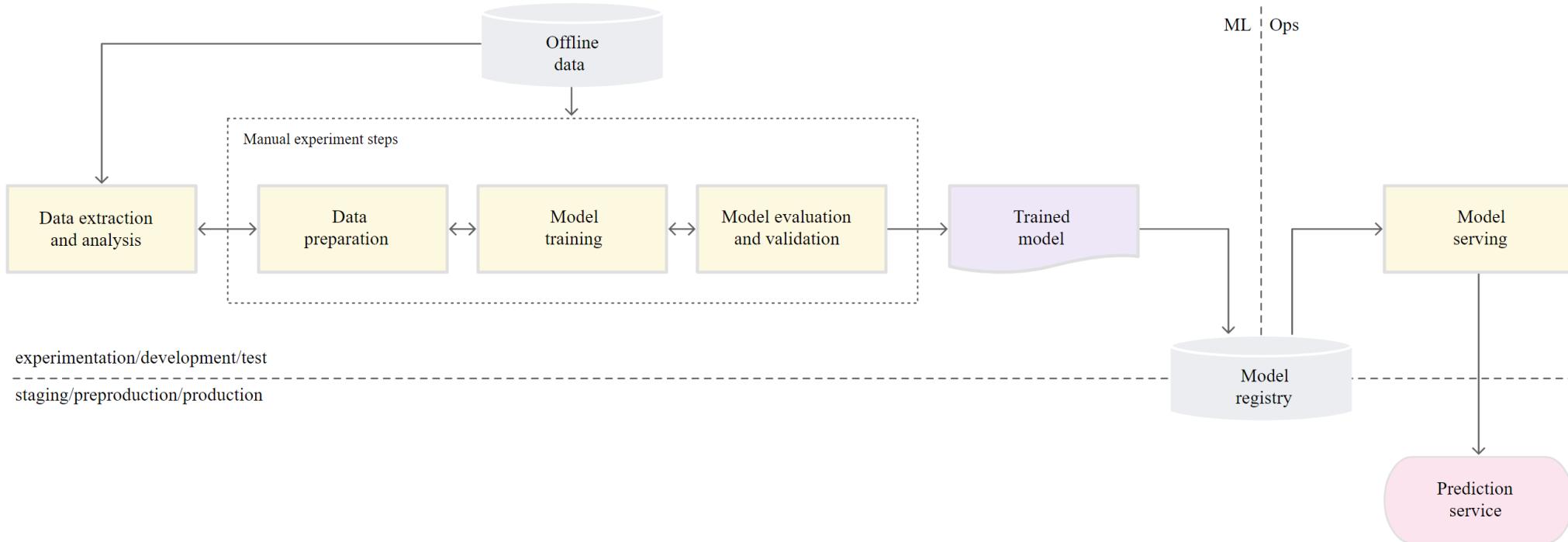


Nível 0 de MLOps: processo manual

Muitas equipes têm cientistas de dados e pesquisadores de ML capazes de criar modelos de última geração, mas o processo de criação e implantação de modelos de ML é totalmente manual. Isso é considerado o nível *básico* de maturidade, ou nível 0

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

Nível 0 de MLOps: processo manual



MLOps: pipelines de entrega contínua e automação no aprendizado de máquina



Nível 0 de MLOps - Desafios

Os MLOps de nível 0 são comuns em muitas empresas que estão começando a aplicar o ML aos casos de uso. Esse processo manual orientado por cientistas de dados pode ser suficiente quando os modelos raramente são alterados ou treinados. Na prática, os modelos costumam falhar quando são implantados no mundo real. Os modelos não se adaptam às mudanças na dinâmica do ambiente ou às alterações nos dados que descrevem o ambiente.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

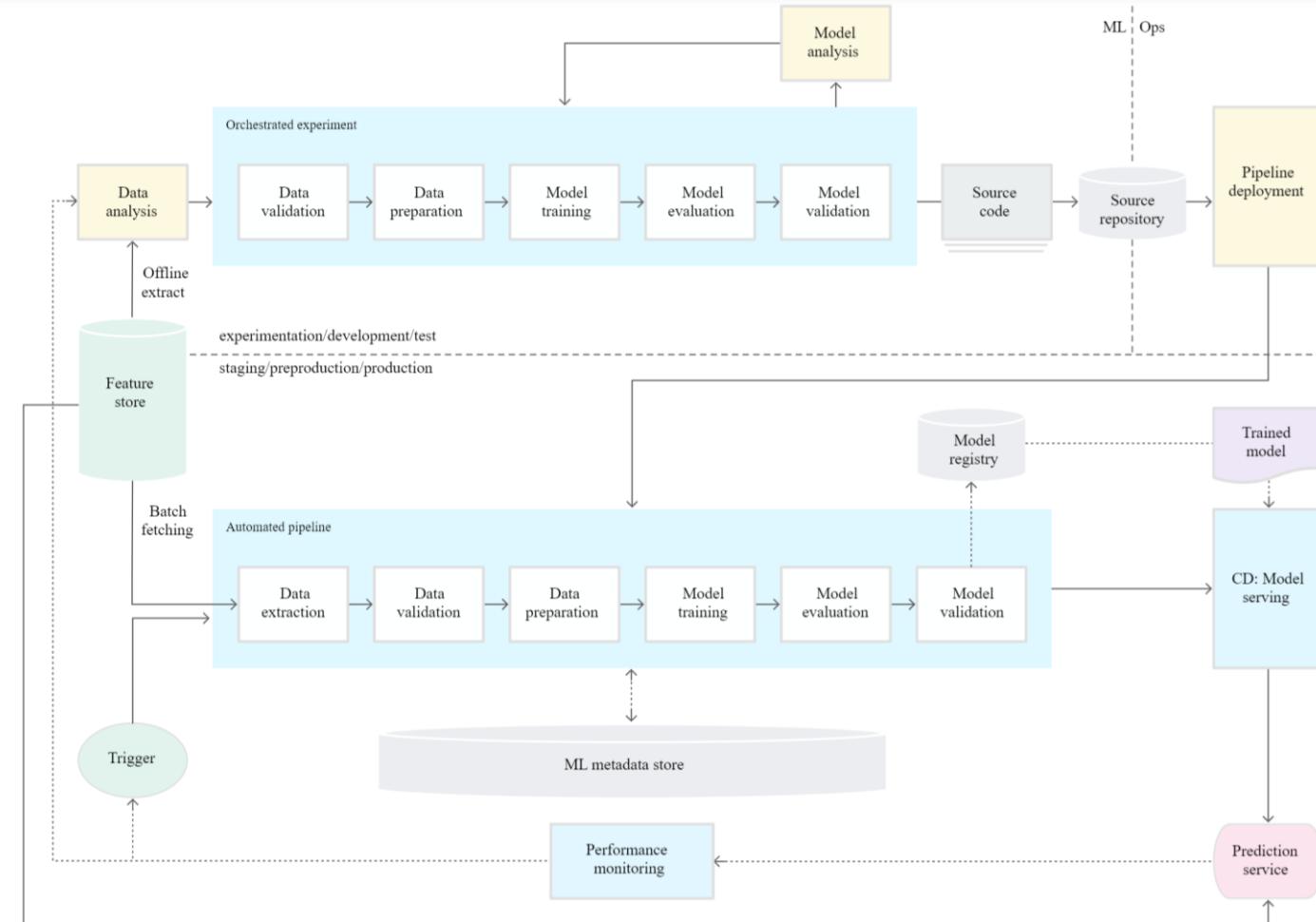


Nível 1 de MLOps: automação de pipeline de ML

O objetivo do nível 1 é realizar o treinamento contínuo do modelo automatizando o pipeline de ML. Isso permite que você alcance a entrega contínua do serviço de predição de modelo. Para automatizar o processo de uso de novos dados para treinar novamente os modelos na produção, é necessário introduzir dados automatizados e passos de validação do modelo, bem como acionadores de pipeline e gerenciamento de metadados.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

Nível 1 de MLOps: automação de pipeline de ML



MLOps: pipelines de entrega contínua e automação no aprendizado de máquina



Nível 1 de MLOps – Desafios

É preciso testar novas ideias de ML e implantar rapidamente novas implementações dos componentes de ML. Ao gerenciar muitos pipelines de ML na produção, é preciso uma configuração de CI/CD para automatizar a criação, o teste e a implantação de pipelines de ML.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

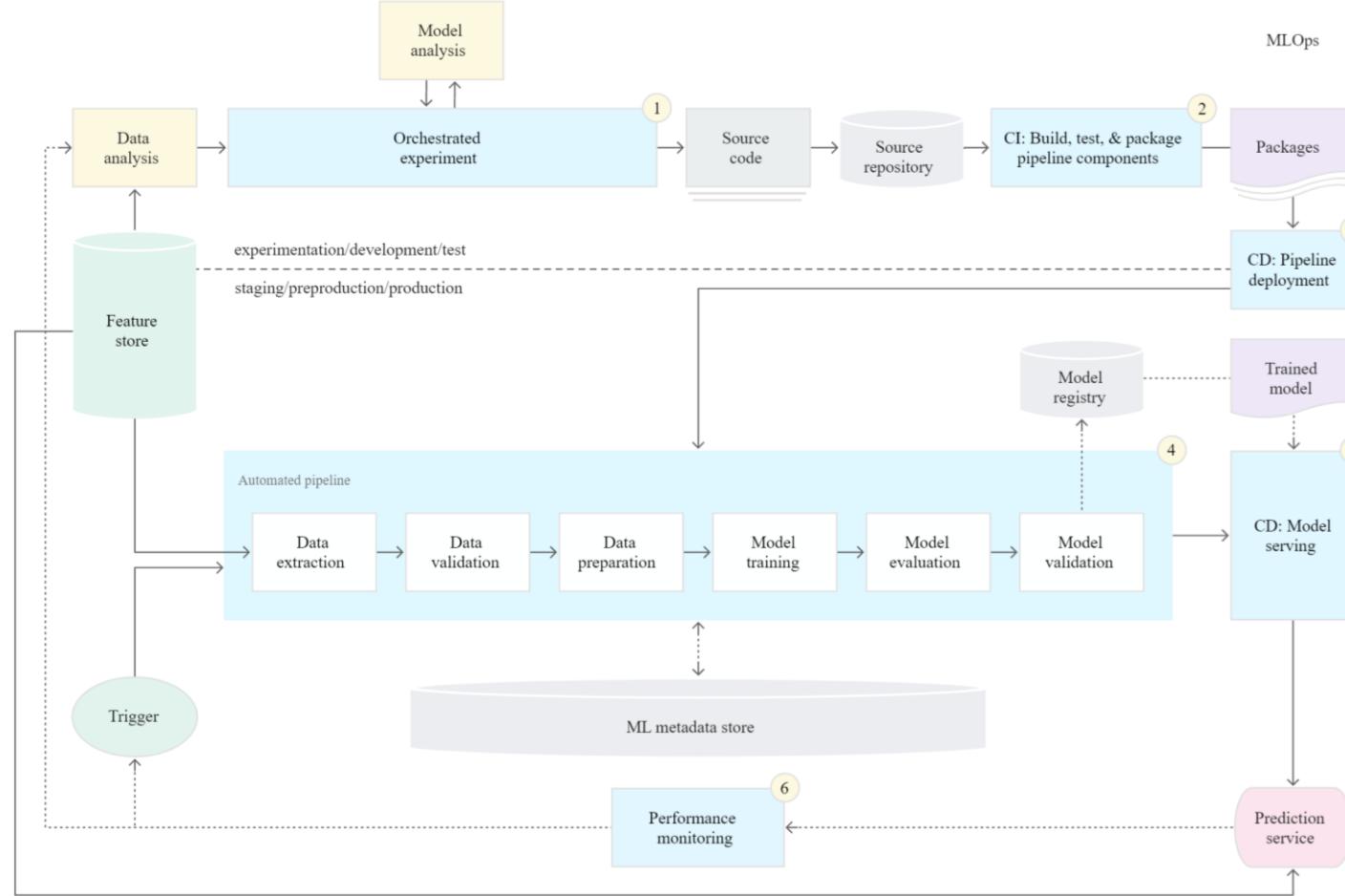


Nível 2 de MLOps: automação de pipeline de CI/CD

Para uma atualização rápida e confiável dos pipelines em produção, é preciso ter um sistema de CI/CD robusto automatizado. Um sistema de CI/CD automatizado permite que os cientistas de dados explorem rapidamente novas ideias sobre engenharia de atributos, arquitetura de modelos e hiperparâmetros. É possível implementar essas ideias e criar, testar e implantar automaticamente os novos componentes do pipeline no ambiente de destino.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

Nível 2 de MLOps: automação de pipeline de CI/CD



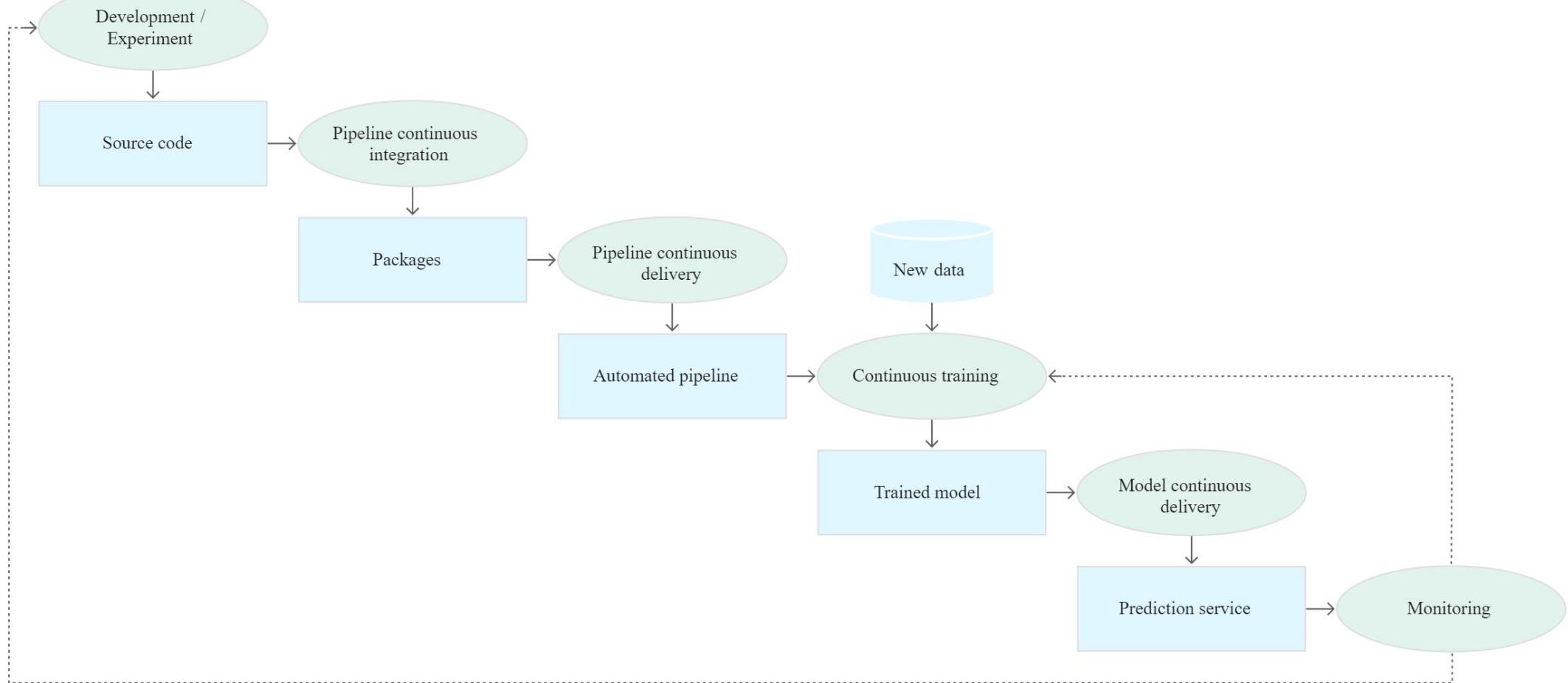
MLOps: pipelines de entrega contínua e automação no aprendizado de máquina



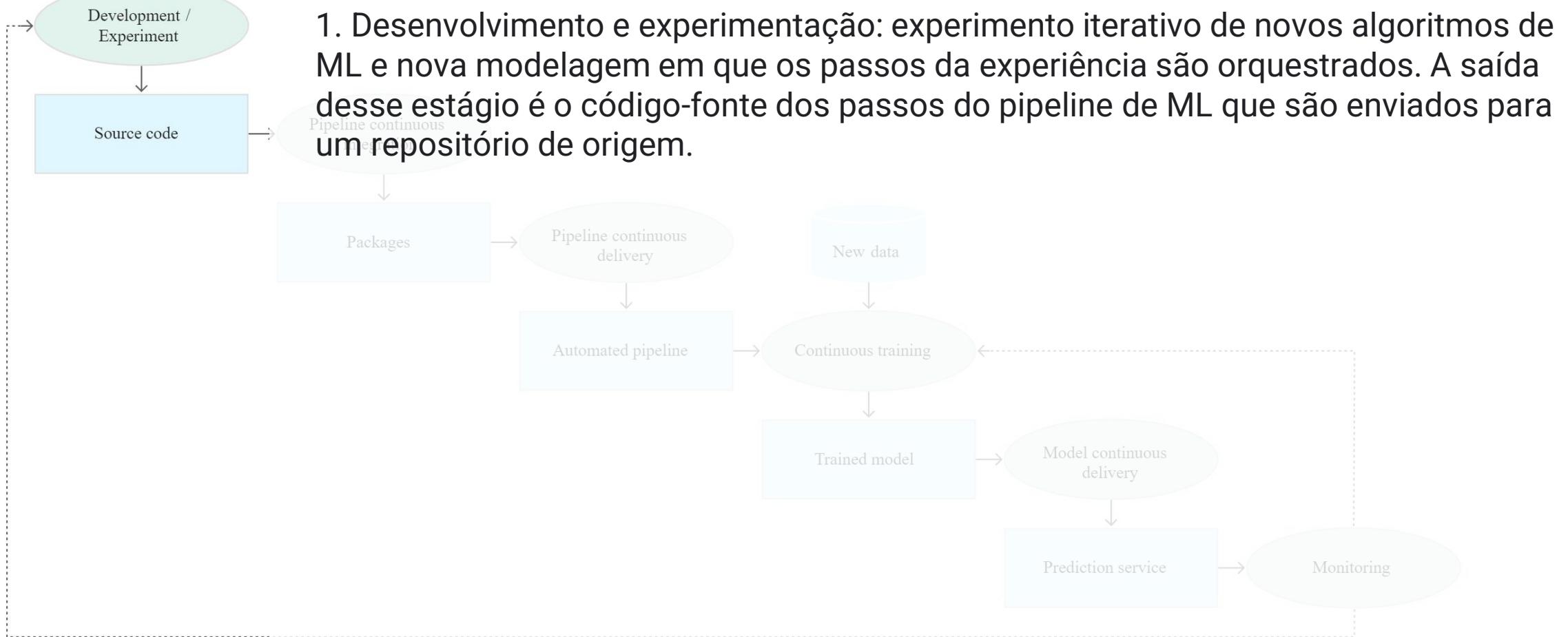
Nível 2 de MLOps: automação de pipeline de CI/CD

- Controle de origem
- Serviços de teste e criação
- Serviços de implantação
- Registro de modelos
- Armazenamento de recursos
- Armazenamento de metadados de ML
- Orquestrador de pipeline de ML

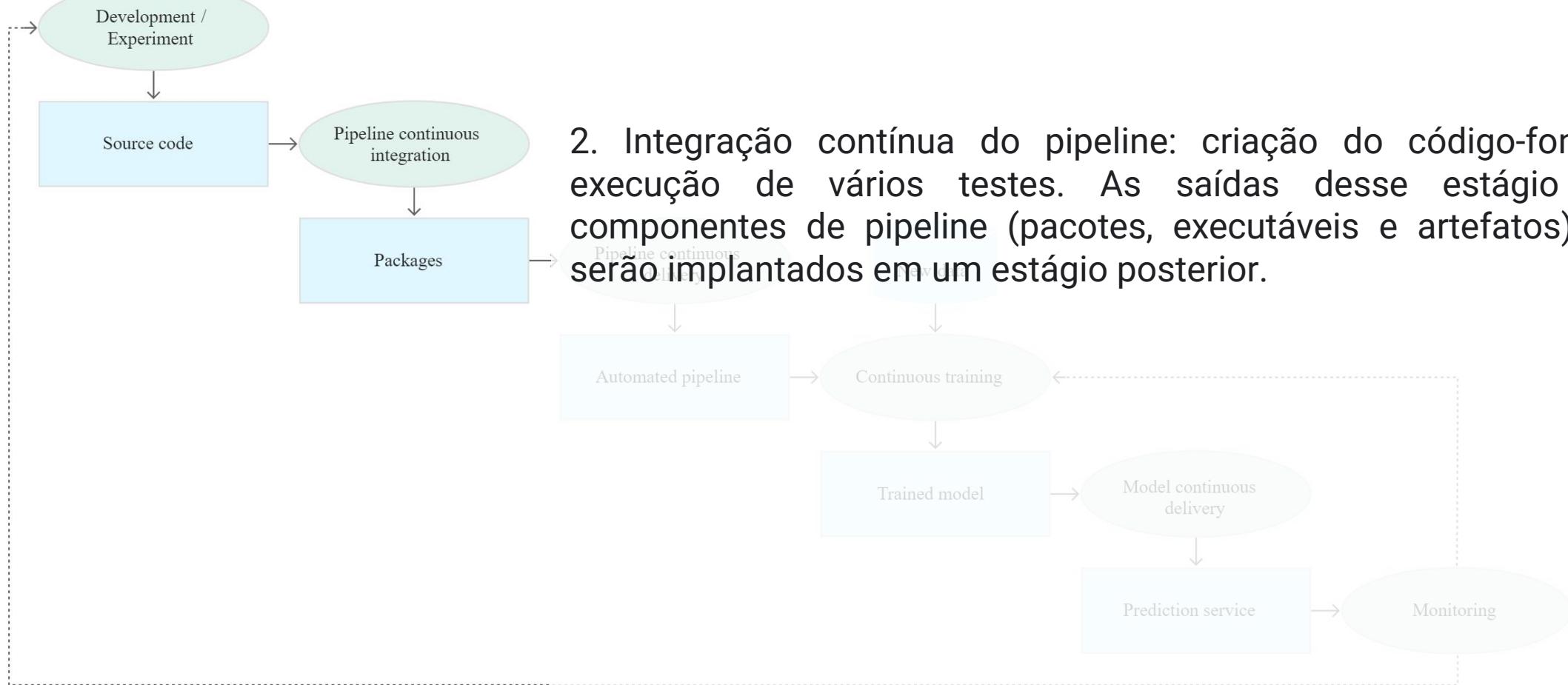
MLOps: pipelines de entrega contínua e automação no aprendizado de máquina



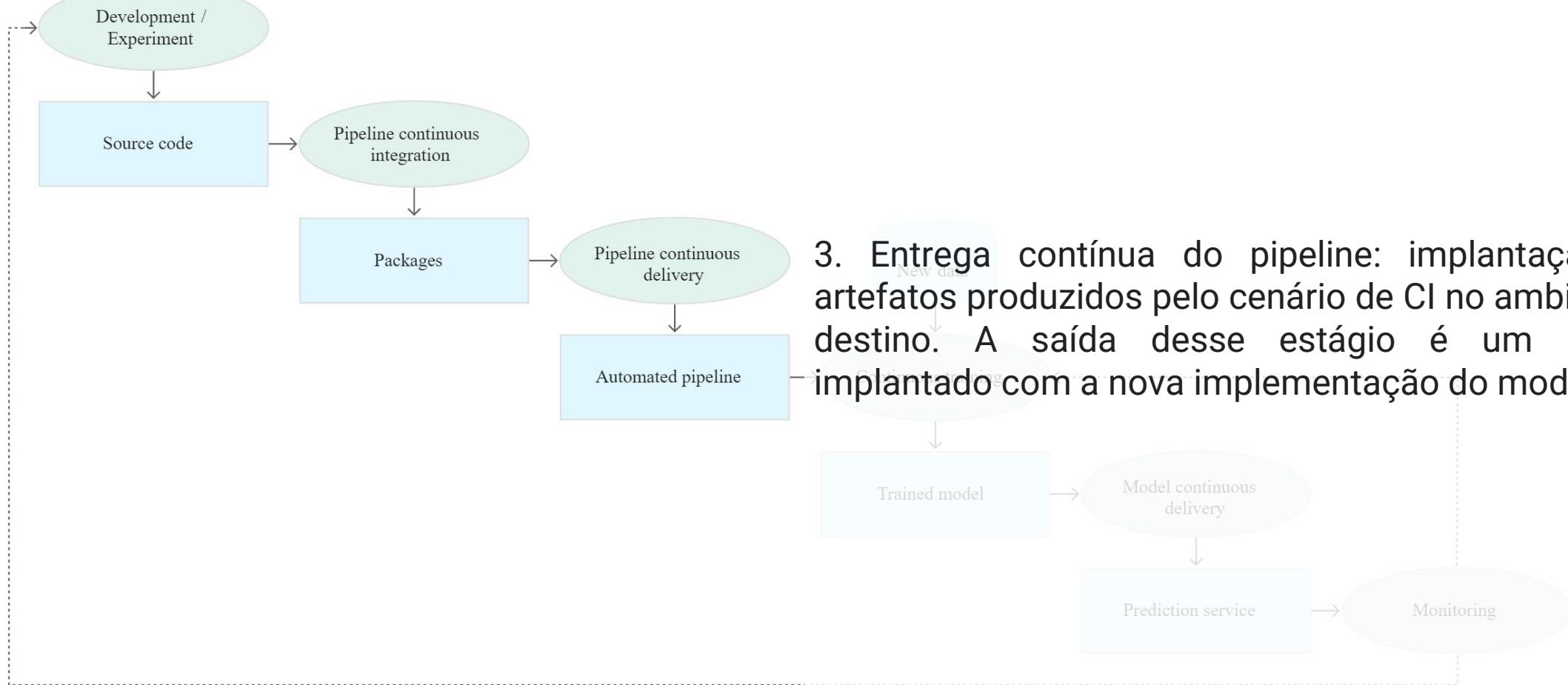
MLOps: pipelines de entrega contínua e automação no aprendizado de máquina



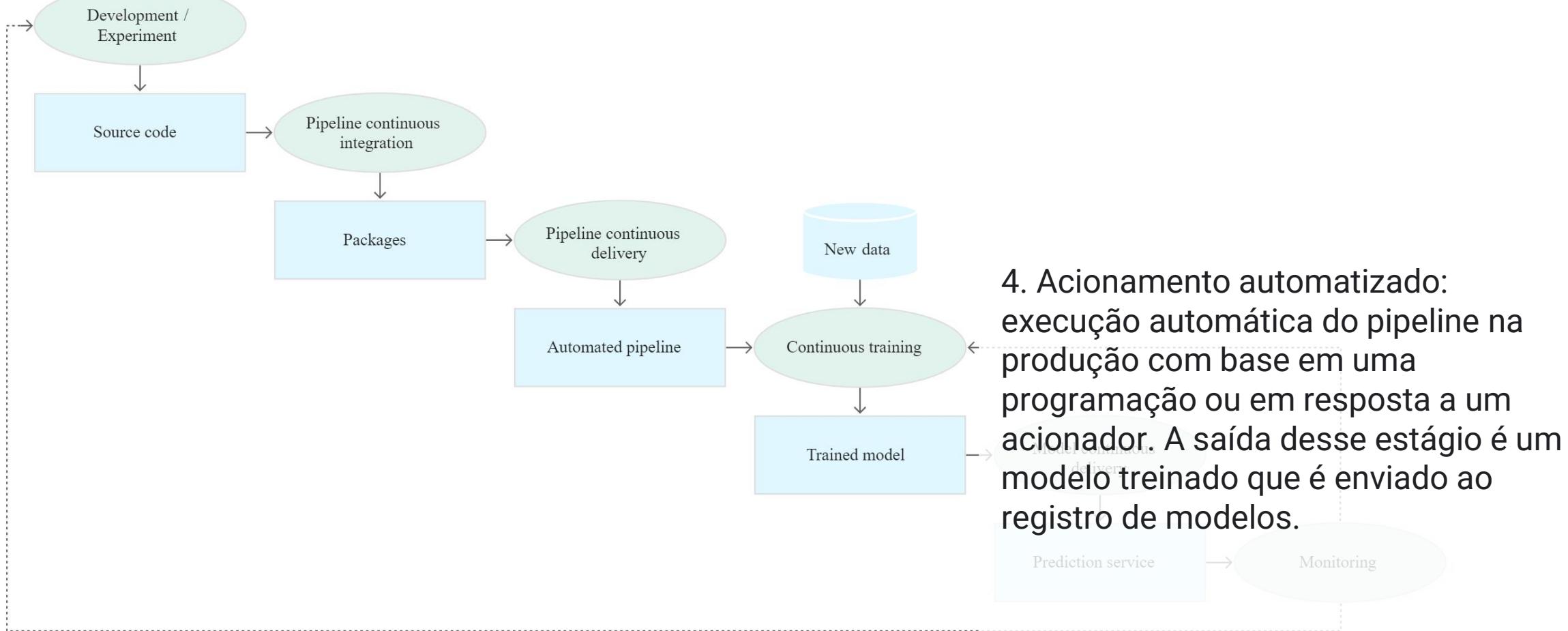
MLOps: pipelines de entrega contínua e automação no aprendizado de máquina



MLOps: pipelines de entrega contínua e automação no aprendizado de máquina



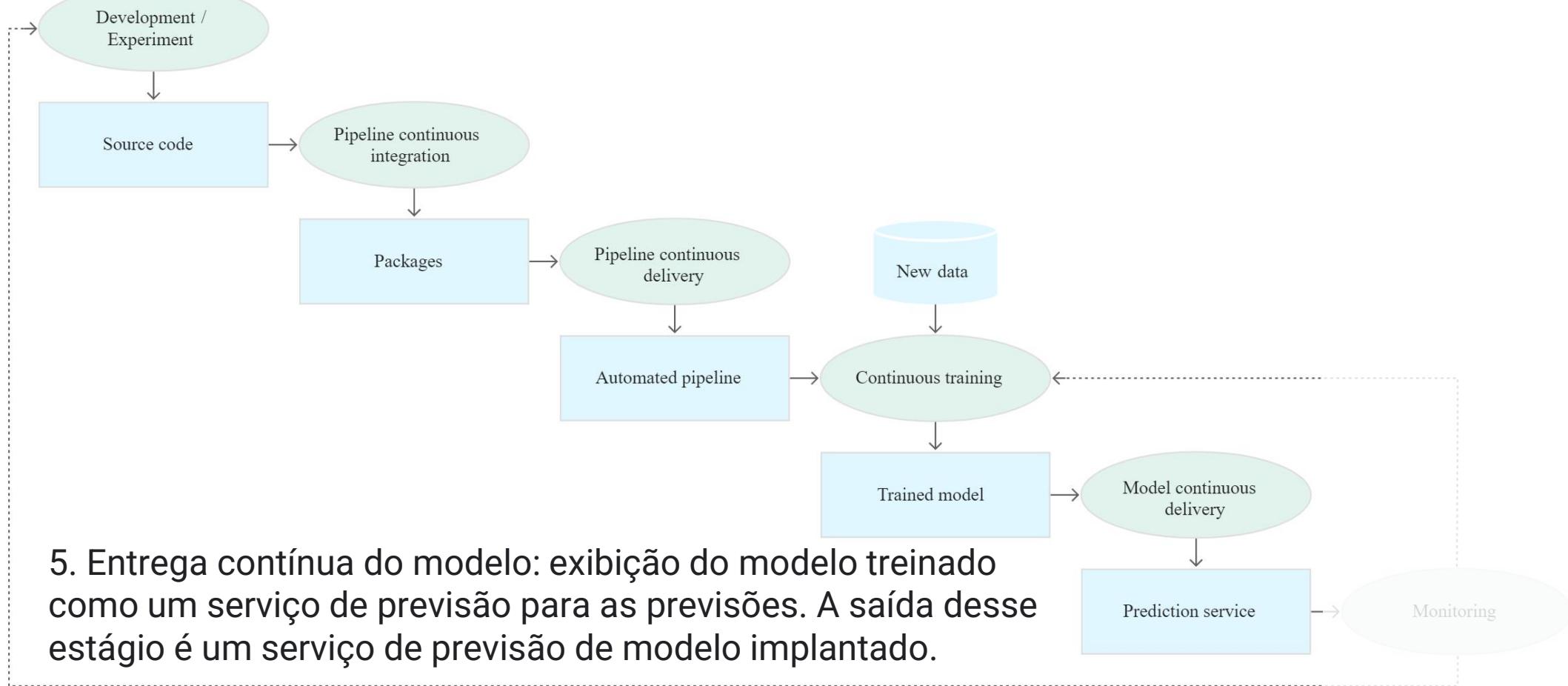
MLOps: pipelines de entrega contínua e automação no aprendizado de máquina



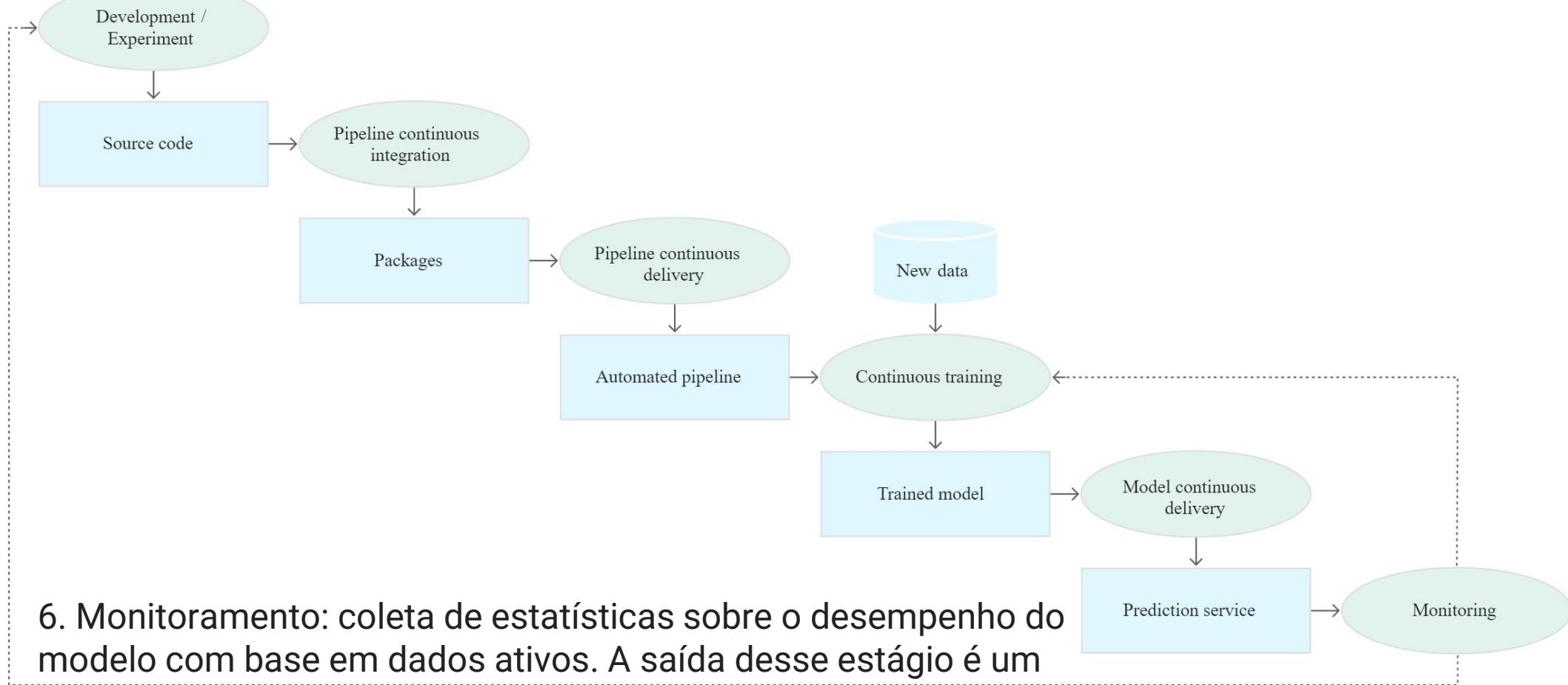
4. Acionamento automatizado: execução automática do pipeline na produção com base em uma programação ou em resposta a um acionador. A saída desse estágio é um modelo treinado que é enviado ao registro de modelos.



MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

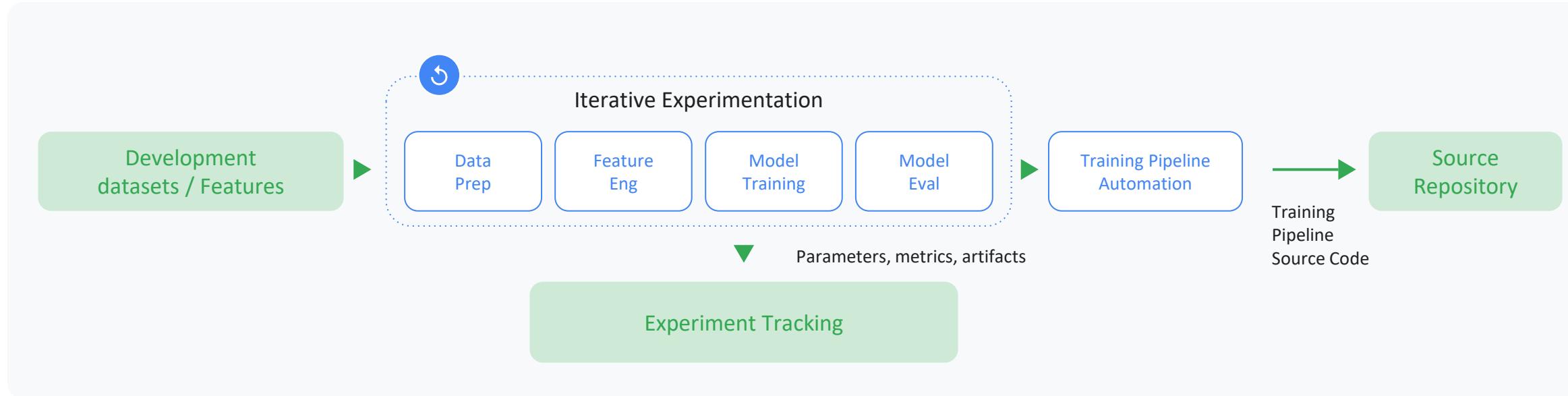


MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

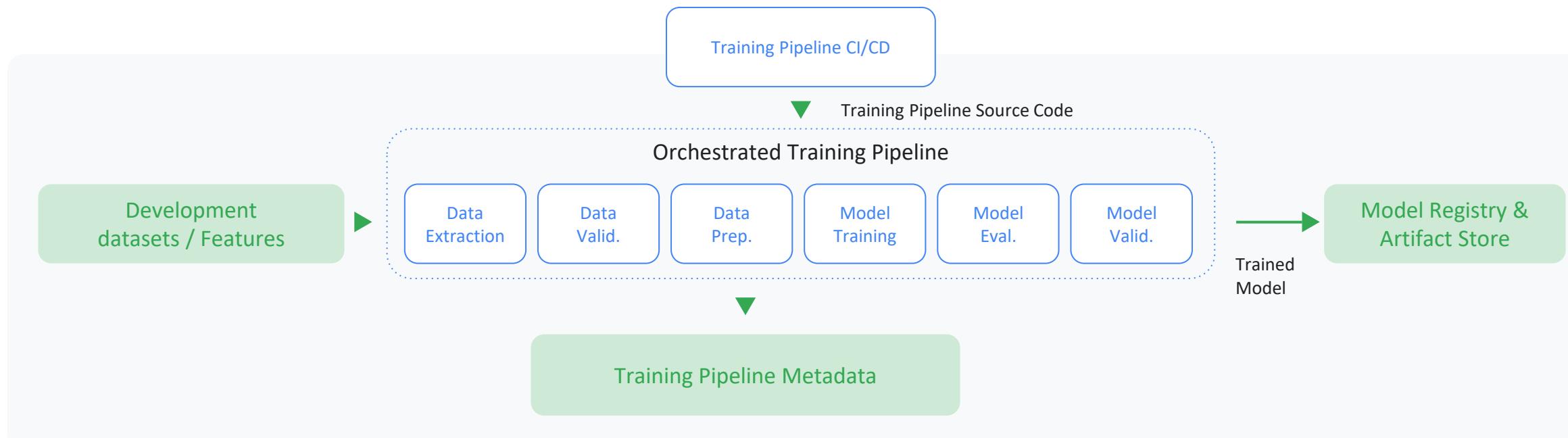


6. Monitoramento: coleta de estatísticas sobre o desempenho do modelo com base em dados ativos. A saída desse estágio é um acionador para executar o pipeline ou executar um novo ciclo de experiência

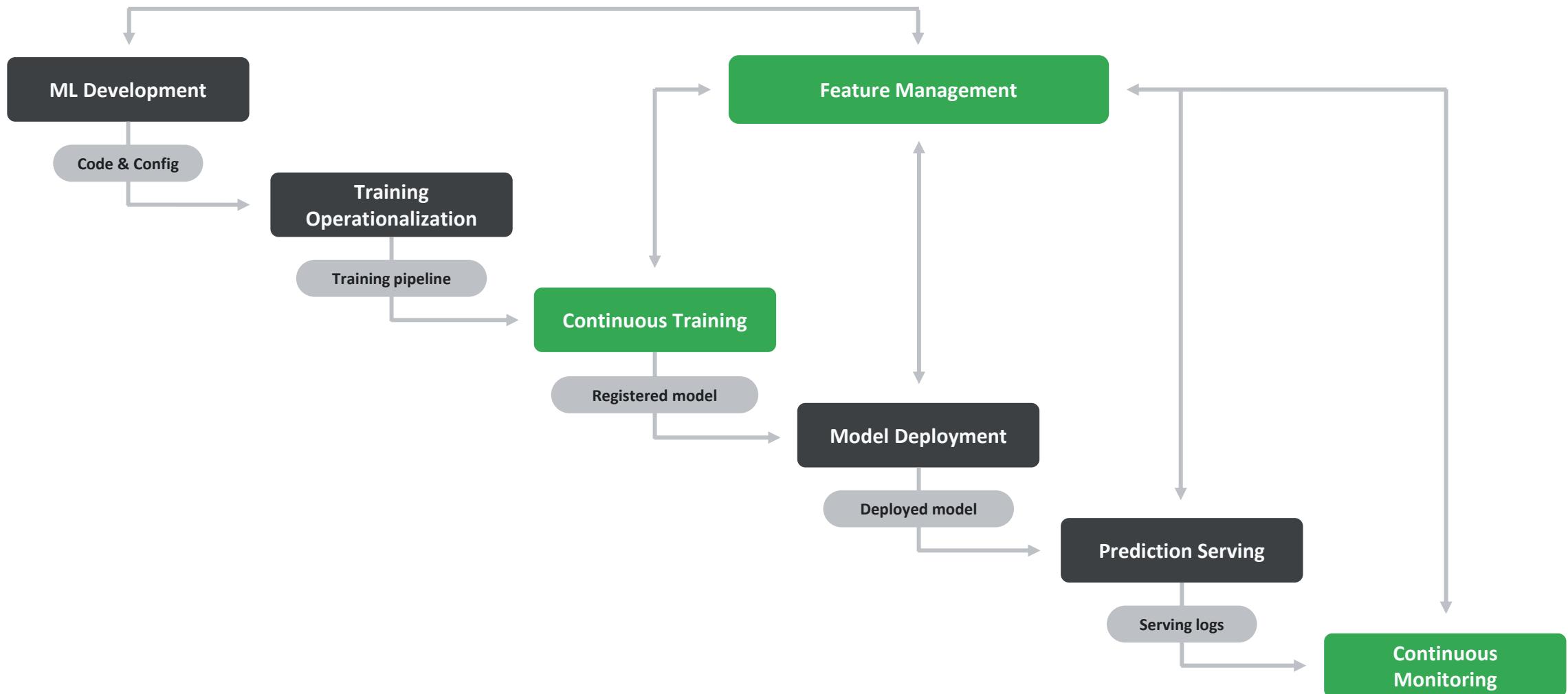
Experimentation management with Vertex Pipelines



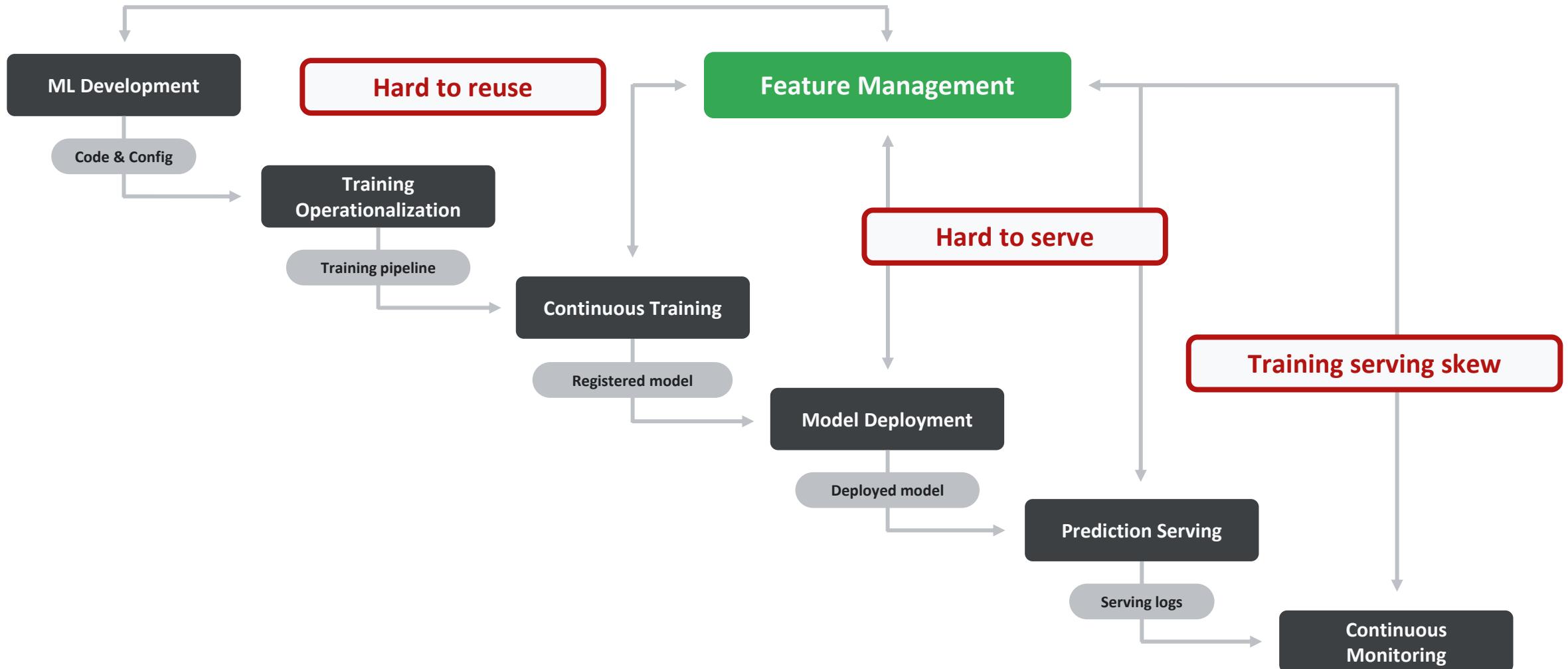
Continuous Training with Vertex Pipelines



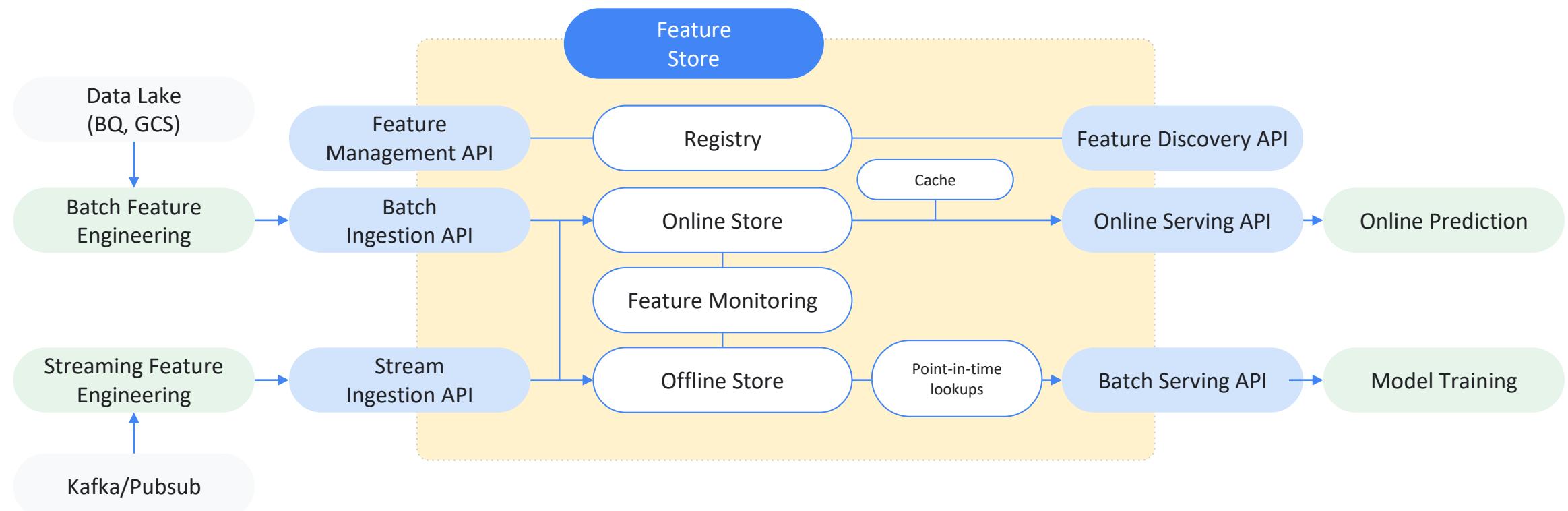
MLOps challenges in ML lifecycle



Feature management pain points

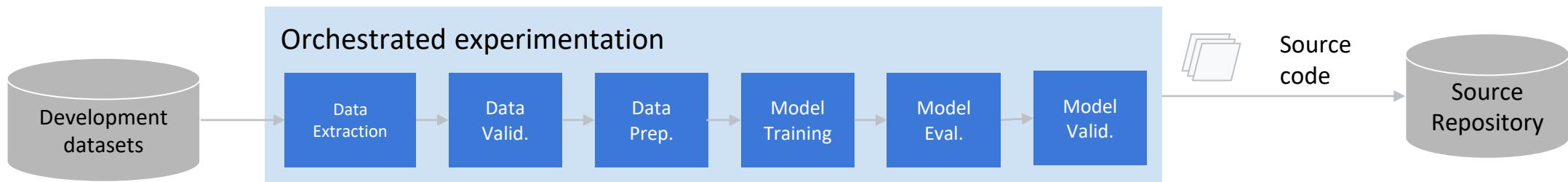


Simple Data Scientist friendly APIs & SDKs abstract away the underlying complexity



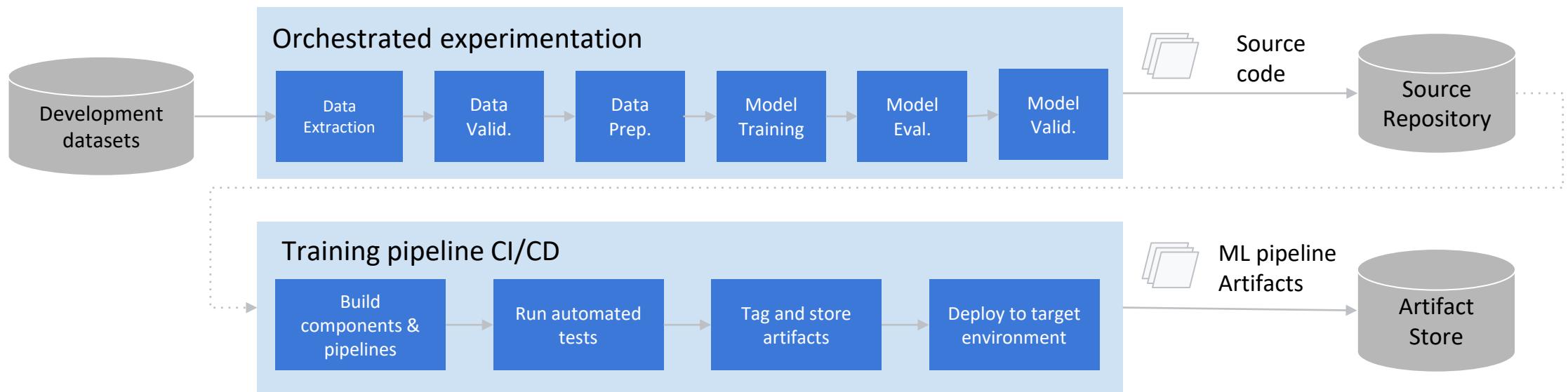
Reliable and repeatable training

Automated E2E Pipelines



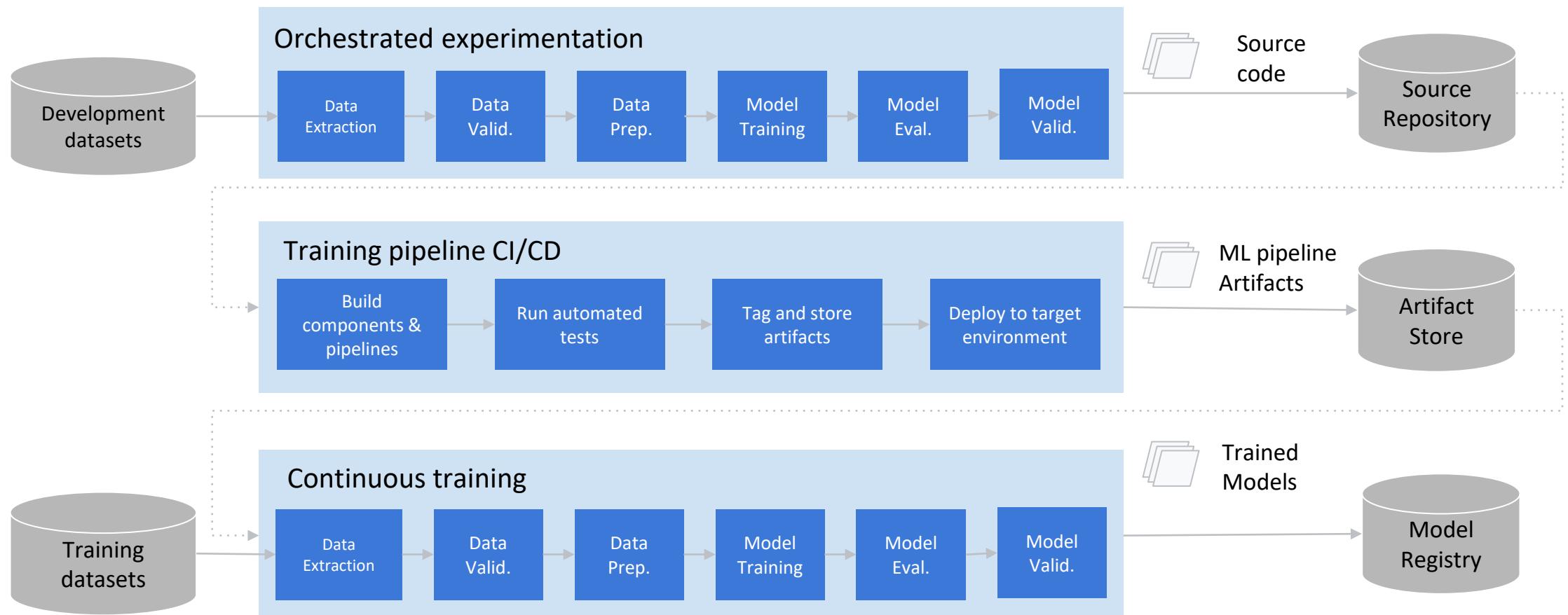
Reliable and repeatable training

Automated E2E Pipelines



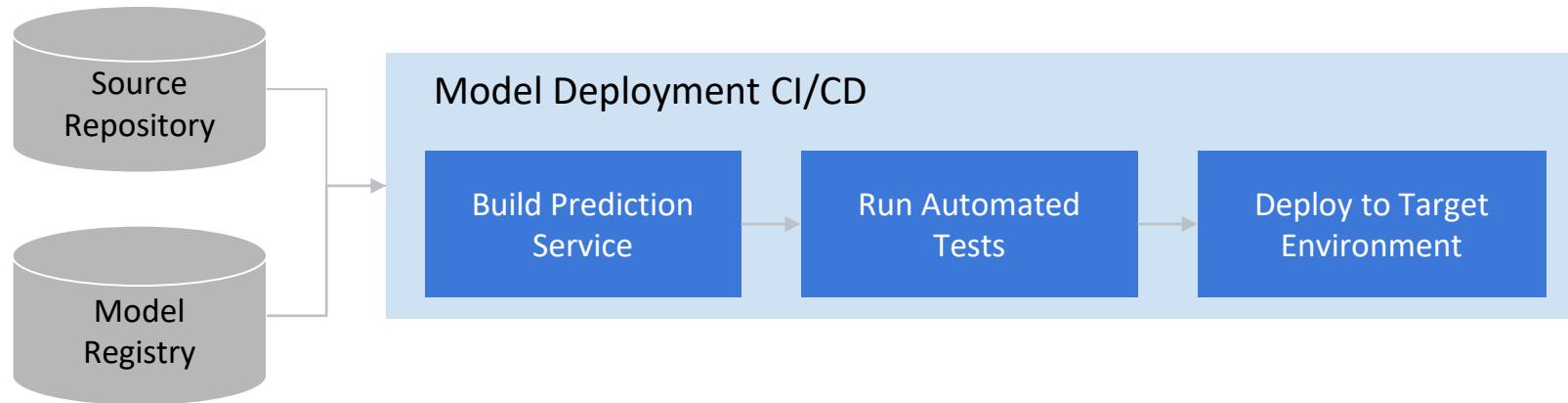
Reliable and repeatable training

Automated E2E Pipelines



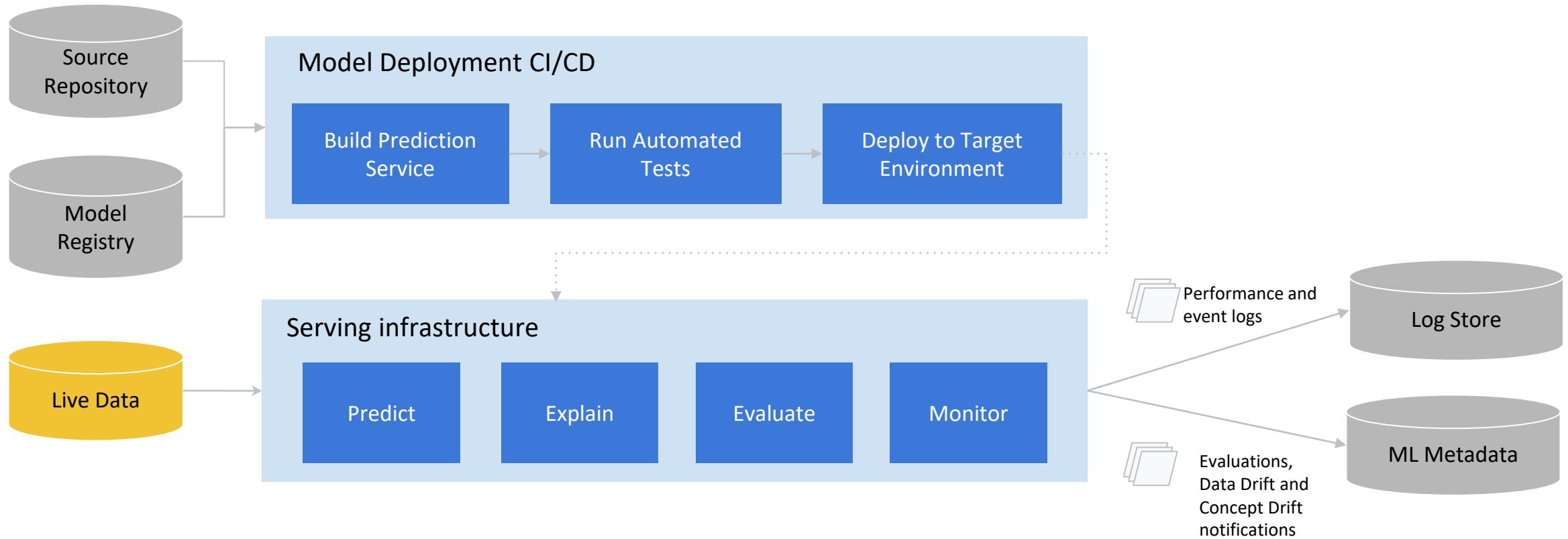
Reliable and monitored serving

Automated E2E Pipelines



Reliable and monitored serving

Automated E2E Pipelines



Leading ML | Best Practices

Continuous Training for Production ML in the TFX Platform. OpML (2019).

Slice Finder: Automated Data Slicing for Model Validation. ICDE (2019).

Data Validation for Machine Learning. SysML (2019).

TFX: A TensorFlow-Based Production-Scale Machine Learning Platform. KDD (2017).

Data Management Challenges in Production Machine Learning. SIGMOD (2017).

Rules of Machine Learning: Best Practices for ML Engineering. Google AI Web (2017).

Machine Learning: The High Interest Credit Card of Technical Debt. NeurIPS (2015).

The High-Interest Credit Card of Technical Debt

D. Sculley, J. Dean, L. Li, Q. V. Le, M. Denn, J. Corrado, R. Barroso, J. Dean, J. Shlens, C. Chaudhuri, J. Zico Kolter, M. Ng, and A. Y. Ng

KDD 2015 Applied Data Science Paper

TFX: A TensorFlow-Based Production-Scale Machine Learning Platform

Denis Baylor, Eric Breck, Heng-Salem Haykal, Mustafa Ispir, Sudip Roy, Clemens Mewald, Akshay Narayan, Steven Euijong Whang, Martin Zinkevich

KDD 2017 Applied Data Science Paper

Data Management Challenges in Production Machine Learning

Neoklis Polyzotis, Sudip Roy, and Martin Zinkevich

SIGMOD 2017

Continuous Training for Production Machine Learning

Denis Baylor, Kevin Haas, Rose Liu, Clemens Mewald, Mitchell Welleck, and Martin Zinkevich

OpML 2019

Automated Data Slicing for Model Validation

Yeounoh Chung, Tim Kraska, Neoklis Polyzotis, Sudip Roy, Steven Euijong Whang, and Martin Zinkevich

ICDE 2019



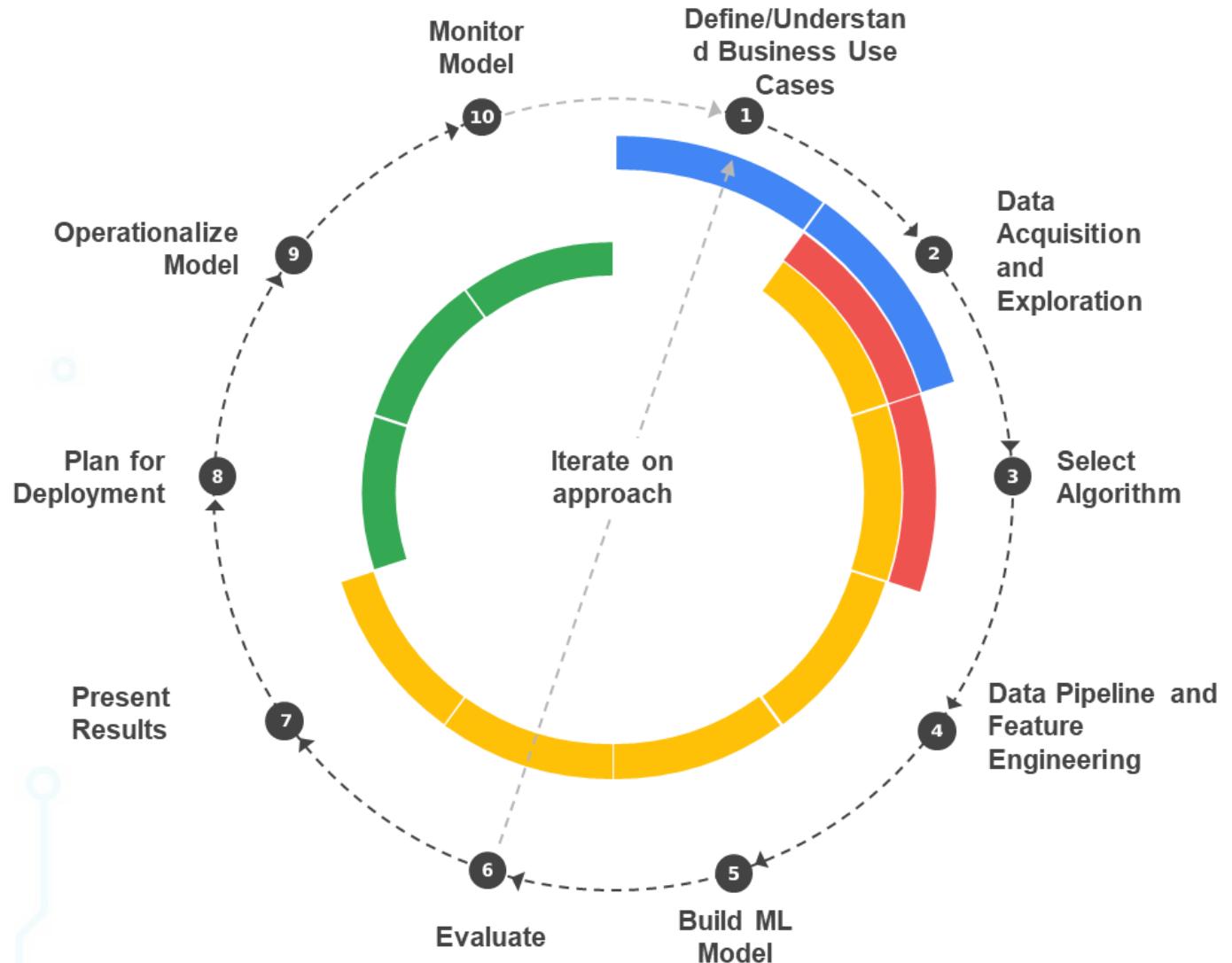
DATA VALIDATION FOR MACHINE LEARNING

Eric Breck¹ Neoklis Polyzotis¹ Sudip Roy¹ Steven Euijong Whang² Martin Zinkevich¹

ABSTRACT

Machine learning is a powerful tool for gleaning knowledge from massive amounts of data. While a great deal of machine learning research has focused on improving the accuracy and efficiency of training and inference algorithms, there is less attention in the equally important problem of monitoring the quality of data fed to machine learning. The importance of this problem is hard to dispute: errors in the input data can nullify any benefits on speed and accuracy for training and inference. This argument points to a data-centric approach to machine learning that treats training and serving data as an important production asset, on par with the algorithm and infrastructure used for learning.

In this paper, we tackle this problem and present a data validation system that is designed to detect anomalies specifically in data fed into machine learning pipelines. This system is deployed in production as an integral part of TFX(Baylor et al., 2017) – an end-to-end machine learning platform at Google. It is used by hundreds of product teams use it to continuously monitor and validate several petabytes of production data per day. We faced several challenges in developing our system, most notably around the ability of ML pipelines to soldier on in the face of unexpected patterns, schema-free data, or training/serving skew. We discuss these challenges, the techniques we used to address them, and the various design choices that we made in implementing the system. Finally, we present evidence from the system’s deployment in production that illustrate the tangible benefits of data validation in the context of ML: early detection of errors, model-quality wins from using better data, savings in engineering hours to debug problems, and a shift towards data-centric workflows in model development.



MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

AutoML

Tools for performing AutoML.

- [AutoGluon](#) - Automates machine learning tasks enabling you to easily achieve strong predictive performance.
- [AutoKeras](#) - AutoKeras goal is to make machine learning accessible for everyone.
- [AutoPyTorch](#) - Automatic architecture search and hyperparameter optimization for PyTorch.
- [AutoSKLearn](#) - Automated machine learning toolkit and a drop-in replacement for a scikit-learn estimator.
- [EvalML](#) - A library that builds, optimizes, and evaluates ML pipelines using domain-specific functions.
- [FLAML](#) - Finds accurate ML models automatically, efficiently and economically.
- [H2O AutoML](#) - Automates ML workflow, which includes automatic training and tuning of models.
- [MindsDB](#) - AI layer for databases that allows you to effortlessly develop, train and deploy ML models.
- [MLBox](#) - MLBox is a powerful Automated Machine Learning python library.
- [Model Search](#) - Framework that implements AutoML algorithms for model architecture search at scale.
- [NNI](#) - An open source AutoML toolkit for automate machine learning lifecycle.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina



CI/CD for Machine Learning

Tools for performing CI/CD for Machine Learning.

- [ClearML](#) - Auto-Magical CI/CD to streamline your ML workflow.
- [CML](#) - Open-source library for implementing CI/CD in machine learning projects.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina



Cron Job Monitoring

Tools for monitoring cron jobs (recurring jobs).

- [Cronitor](#) - Monitor any cron job or scheduled task.
- [HealthchecksIO](#) - Simple and effective cron job monitoring.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina



Data Catalog

Tools for data cataloging.

- [Amundsen](#) - Data discovery and metadata engine for improving the productivity when interacting with data.
- [Apache Atlas](#) - Provides open metadata management and governance capabilities to build a data catalog.
- [CKAN](#) - Open-source DMS (data management system) for powering data hubs and data portals.
- [DataHub](#) - LinkedIn's generalized metadata search & discovery tool.
- [Magda](#) - A federated, open-source data catalog for all your big data and small data.
- [Metacat](#) - Unified metadata exploration API service for Hive, RDS, Teradata, Redshift, S3 and Cassandra.
- [OpenMetadata](#) - A Single place to discover, collaborate and get your data right.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina



Data Enrichment

Tools and libraries for data enrichment.

- [Upgini](#) - Enriches training datasets with features from public and community shared data sources.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina



Data Exploration

Tools for performing data exploration.

- [Apache Zeppelin](#) - Enables data-driven, interactive data analytics and collaborative documents.
- [BambooLib](#) - An intuitive GUI for Pandas DataFrames.
- [Google Colab](#) - Hosted Jupyter notebook service that requires no setup to use.
- [Jupyter Notebook](#) - Web-based notebook environment for interactive computing.
- [JupyterLab](#) - The next-generation user interface for Project Jupyter.
- [Jupytext](#) - Jupyter Notebooks as Markdown Documents, Julia, Python or R scripts.
- [Polynote](#) - The polyglot notebook with first-class Scala support.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

Data Management

Tools for performing data management.

- [Arrikto](#) - Dead simple, ultra fast storage for the hybrid Kubernetes world.
- [BlazingSQL](#) - A lightweight, GPU accelerated, SQL engine for Python. Built on RAPIDS cuDF.
- [Delta Lake](#) - Storage layer that brings scalable, ACID transactions to Apache Spark and other engines.
- [Dolt](#) - SQL database that you can fork, clone, branch, merge, push and pull just like a git repository.
- [Dud](#) - A lightweight CLI tool for versioning data alongside source code and building data pipelines.
- [DVC](#) - Management and versioning of datasets and machine learning models.
- [Git LFS](#) - An open source Git extension for versioning large files.
- [Hub](#) - A dataset format for creating, storing, and collaborating on AI datasets of any size.
- [Intake](#) - A lightweight set of tools for loading and sharing data in data science projects.
- [lakeFS](#) - Repeatable, atomic and versioned data lake on top of object storage.
- [Marquez](#) - Collect, aggregate, and visualize a data ecosystem's metadata.
- [Milvus](#) - An open source embedding vector similarity search engine powered by Faiss, NMSLIB and Annoy.
- [Pinecone](#) - Managed and distributed vector similarity search used with a lightweight SDK.
- [Qdrant](#) - An open source vector similarity search engine with extended filtering support.
- [Quilt](#) - A self-organizing data hub with S3 support.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

Data Processing

Tools related to data processing and data pipelines.

- [Airflow](#) - Platform to programmatically author, schedule, and monitor workflows.
- [Azkaban](#) - Batch workflow job scheduler created at LinkedIn to run Hadoop jobs.
- [Dagster](#) - A data orchestrator for machine learning, analytics, and ETL.
- [Hadoop](#) - Framework that allows for the distributed processing of large data sets across clusters.
- [Spark](#) - Unified analytics engine for large-scale data processing.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina



Data Validation

Tools related to data validation.

- [Cerberus](#) - Lightweight, extensible data validation library for Python.
- [Great Expectations](#) - A Python data validation framework that allows to test your data against datasets.
- [JSON Schema](#) - A vocabulary that allows you to annotate and validate JSON documents.
- [TFDV](#) - An library for exploring and validating machine learning data.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

Data Visualization

Tools for data visualization, reports and dashboards.

- [Count](#) - SQL/drag-and-drop querying and visualisation tool based on notebooks.
- [Dash](#) - Analytical Web Apps for Python, R, Julia, and Jupyter.
- [Data Studio](#) - Reporting solution for power users who want to go beyond the data and dashboards of GA.
- [Facets](#) - Visualizations for understanding and analyzing machine learning datasets.
- [Lux](#) - Fast and easy data exploration by automating the visualization and data analysis process.
- [Metabase](#) - The simplest, fastest way to get business intelligence and analytics to everyone.
- [Redash](#) - Connect to any data source, easily visualize, dashboard and share your data.
- [Superset](#) - Modern, enterprise-ready business intelligence web application.
- [Tableau](#) - Powerful and fastest growing data visualization tool used in the business intelligence industry.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

Feature Engineering

Tools and libraries related to feature engineering.

- [Feature Engine](#) - Feature engineering package with SKlearn like functionality.
- [Featuretools](#) - Python library for automated feature engineering.
- [TSFresh](#) - Python library for automatic extraction of relevant features from time series.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina



Feature Store

Feature store tools for data serving.

- [Butterfree](#) - A tool for building feature stores. Transform your raw data into beautiful features.
- [ByteHub](#) - An easy-to-use feature store. Optimized for time-series data.
- [Feast](#) - End-to-end open source feature store for machine learning.
- [Feathr](#) - An enterprise-grade, high performance feature store.
- [Tecton](#) - A fully-managed feature platform built to orchestrate the complete lifecycle of features.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

Hyperparameter Tuning

Tools and libraries to perform hyperparameter tuning.

- [Advisor](#) - Open-source implementation of Google Vizier for hyper parameters tuning.
- [Hyperas](#) - A very simple wrapper for convenient hyperparameter optimization.
- [Hyperopt](#) - Distributed Asynchronous Hyperparameter Optimization in Python.
- [Katib](#) - Kubernetes-based system for hyperparameter tuning and neural architecture search.
- [KerasTuner](#) - Easy-to-use, scalable hyperparameter optimization framework.
- [Optuna](#) - Open source hyperparameter optimization framework to automate hyperparameter search.
- [Scikit Optimize](#) - Simple and efficient library to minimize expensive and noisy black-box functions.
- [Talos](#) - Hyperparameter Optimization for TensorFlow, Keras and PyTorch.
- [Tune](#) - Python library for experiment execution and hyperparameter tuning at any scale.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

Machine Learning Platform

Complete machine learning platform solutions.

- [aiWARE](#) - aiWARE helps MLOps teams evaluate, deploy, integrate, scale & monitor ML models.
- [Algorithmia](#) - Securely govern your machine learning operations with a healthy ML lifecycle.
- [Allegro AI](#) - Transform ML/DL research into products. Faster.
- [Bodywork](#) - Deploys machine learning projects developed in Python, to Kubernetes.
- [CNVRG](#) - An end-to-end machine learning platform to build and deploy AI models at scale.
- [DAGsHub](#) - A platform built on open source tools for data, model and pipeline management.
- [Dataiku](#) - Platform democratizing access to data and enabling enterprises to build their own path to AI.
- [DataRobot](#) - AI platform that democratizes data science and automates the end-to-end ML at scale.
- [Domino](#) - One place for your data science tools, apps, results, models, and knowledge.
- [Edge Impulse](#) - Platform for creating, optimizing, and deploying AI/ML algorithms for edge devices.
- [envd](#) - Machine learning development environment for data science and AI/ML engineering teams.
- [FedML](#) - Simplifies the workflow of federated learning anywhere at any scale.
- [Gradient](#) - Multicloud CI/CD and MLOps platform for machine learning teams.
- [H2O](#) - Open source leader in AI with a mission to democratize AI for everyone.
- [Hopsworks](#) - Open-source platform for developing and operating machine learning models at scale.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

Model Fairness and Privacy

Tools for performing model fairness and privacy in production.

- [AIF360](#) - A comprehensive set of fairness metrics for datasets and machine learning models.
- [Fairlearn](#) - A Python package to assess and improve fairness of machine learning models.
- [Opacus](#) - A library that enables training PyTorch models with differential privacy.
- [TensorFlow Privacy](#) - Library for training machine learning models with privacy for training data.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

Model Interpretability

Tools for performing model interpretability/explainability.

- [Alibi](#) - Open-source Python library enabling ML model inspection and interpretation.
- [Captum](#) - Model interpretability and understanding library for PyTorch.
- [ELI5](#) - Python package which helps to debug machine learning classifiers and explain their predictions.
- [InterpretML](#) - A toolkit to help understand models and enable responsible machine learning.
- [LIME](#) - Explaining the predictions of any machine learning classifier.
- [Lucid](#) - Collection of infrastructure and tools for research in neural network interpretability.
- [SAGE](#) - For calculating global feature importance using Shapley values.
- [SHAP](#) - A game theoretic approach to explain the output of any machine learning model.
- [Skater](#) - Unified framework to enable Model Interpretation for all forms of model.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

Model Lifecycle

Tools for managing model lifecycle (tracking experiments, parameters and metrics).

- [Aim](#) - A super-easy way to record, search and compare 1000s of ML training runs.
- [Comet](#) - Track your datasets, code changes, experimentation history, and models.
- [Guild AI](#) - Open source experiment tracking, pipeline automation, and hyperparameter tuning.
- [Keepsake](#) - Version control for machine learning with support to Amazon S3 and Google Cloud Storage.
- [Losswise](#) - Makes it easy to track the progress of a machine learning project.
- [MLflow](#) - Open source platform for the machine learning lifecycle.
- [ModelDB](#) - Open source ML model versioning, metadata, and experiment management.
- [Neptune AI](#) - The most lightweight experiment management tool that fits any workflow.
- [Replicate](#) - Library that uploads files and metadata (like hyperparameters) to S3 or GCS.
- [Sacred](#) - A tool to help you configure, organize, log and reproduce experiments.
- [Weights and Biases](#) - A tool for visualizing and tracking your machine learning experiments.

MLOps: pipelines de entrega contínua e automação no aprendizado de máquina

Model Serving

Tools for serving models in production.

- [Banana](#) - Host your ML inference code on serverless GPUs and integrate it into your app with one line of code.
- [BentoML](#) - Open-source platform for high-performance ML model serving.
- [BudgetML](#) - Deploy a ML inference service on a budget in less than 10 lines of code.
- [Cortex](#) - Machine learning model serving infrastructure.
- [Gradio](#) - Create customizable UI components around your models.
- [GraphPipe](#) - Machine learning model deployment made simple.
- [Hydrosphere](#) - Platform for deploying your Machine Learning to production.
- [KFServing](#) - Kubernetes custom resource definition for serving ML models on arbitrary frameworks.
- [Merlin](#) - A platform for deploying and serving machine learning models.
- [MLEM](#) - Version and deploy your ML models following GitOps principles.
- [Opyrator](#) - Turns your ML code into microservices with web API, interactive GUI, and more.
- [Rune](#) - Provides containers to encapsulate and deploy EdgeML pipelines and applications.
- [Seldon](#) - Take your ML projects from POC to production with maximum efficiency and minimal risk.
- [Streamlit](#) - Lets you create apps for your ML projects with deceptively simple Python scripts.
- [TensorFlow Serving](#) - Flexible, high-performance serving system for ML models, designed for production.

Collect
data from internal & external sources

Ingest
data through batch jobs or streams

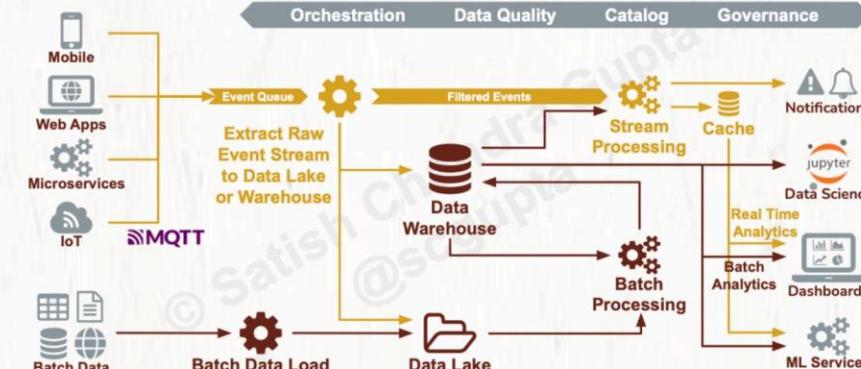
Store
in Data Lake or Data Warehouse

Compute
analytics aggregations, ML features

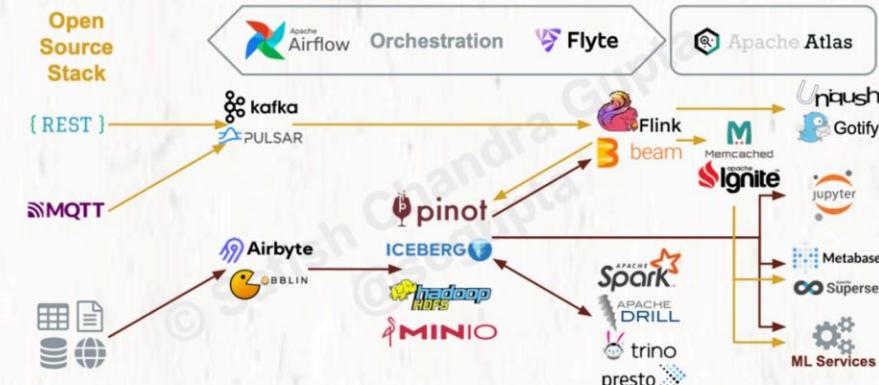
Use
it in dashboards, data science, ML



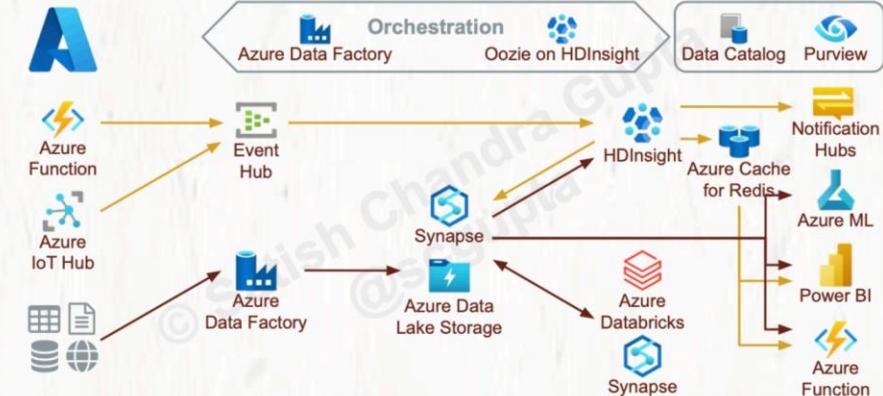
Collect → Ingest → Store → Compute → Use



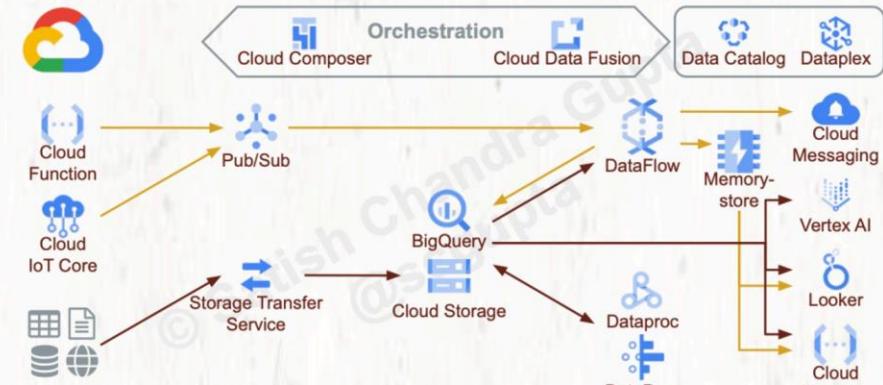
Collect → Ingest → Store → Compute → Use



Collect → Ingest → Store → Compute → Use



Collect → Ingest → Store → Compute → Use



Thanks !



Vinicius Fernandes Caridá
vfcarida@gmail.com



@vinicius caridá



@vfcarida



@vinicius caridá



@vfcarida



@vinicius caridá



@vfcarida

O que você achou da aula de hoje?



Questions and Feedback

MBA⁺

Copyright © Prof. Vinicius Fernandes Caridá
Todos direitos reservados. Reprodução ou divulgação
total ou parcial deste documento é expressamente
proibido sem o consentimento formal, por escrito, do
Professor (autor).

F | A P