

Data Platform Architectures & Machine Learning Operations (MLOps)

Dr. Christoph Gröger
Arnold Lutsch

Robert Bosch GmbH
IoT & Digitalization - Data Strategy

Prof. Dr. Jan Kirenz
HdM Stuttgart

Data platform architectures & MLOps

Modul-Nr: 338025-338027

SWS/ECTS: 5/10

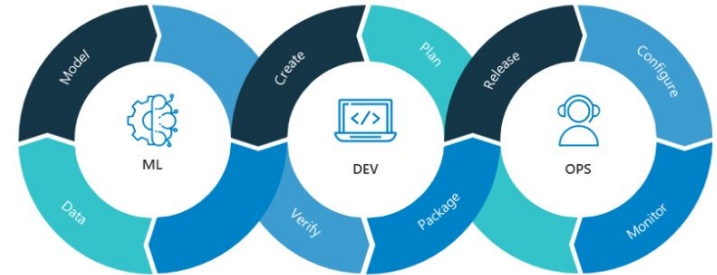
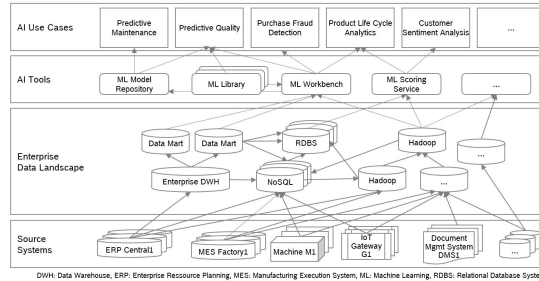
Prüfungsform: PP

Projektpartner:

Dr. Christoph Gröger
Arnold Lutsch
Robert Bosch GmbH

Themenschwerpunkte

- Datenplattform-architektur
- Open Source Technologien (bspw. Delta Lake, Kubeflow, TFX)
- Machine Learning Operations
- Entwicklung in Python



Im Rahmen des Projekts soll mit Hilfe von Open-Source-Software eine prototypische state-of-the-art **Data Lake bzw. Datenplattformarchitektur** (bspw. Lakehouse) für die Realisierung unterschiedlicher Machine Learning Anwendungsfälle konzipiert und implementiert werden.

In der Architektur sollen für eine möglichst umfassende Automatisierung des **Machine Learning Lifecycles** Komponenten der Disziplin **Machine Learning Operations (MLOps)** berücksichtigt werden (bspw. feature store, model registry, data pipelines).

“... [we need to] help companies progress on their AI journey, from one-off AI experimentation to gaining a robust organization-wide capability that acts as a source of competitive agility and growth.”

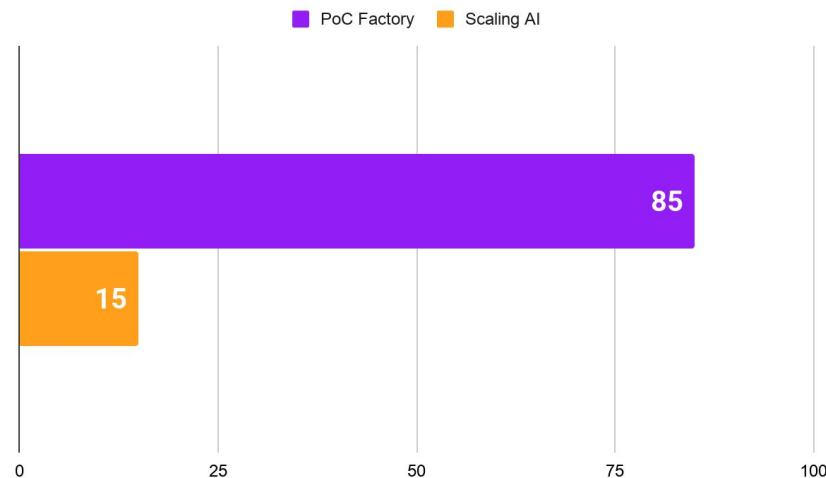
Accenture (2019)

The Proof of Concept Factory

80-85% PoC Factory

Most companies...

- ... conduct AI experiments and pilots but achieve a low scaling success rate
- ... have significant AI under investments, yielding low returns



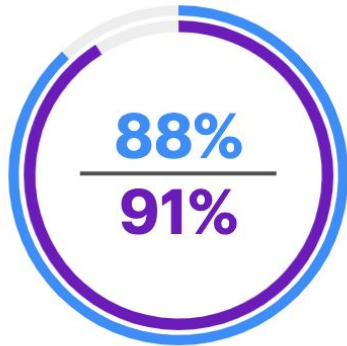
How crucial is scaling AI to your business?

UNITED STATES

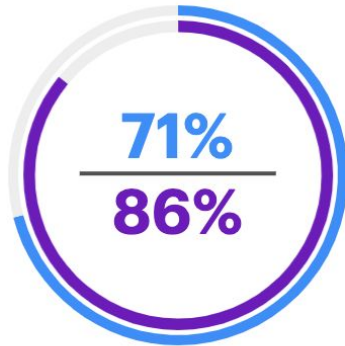


vs.

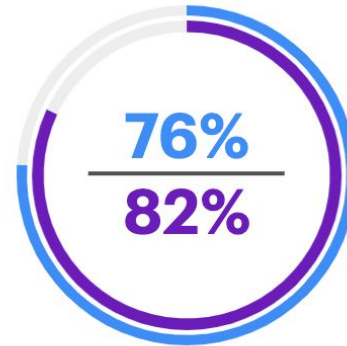
GERMANY



of executives say they won't achieve their growth objectives without scaling AI.



of executives believe they risk going out of business in 5 years if they don't scale AI.

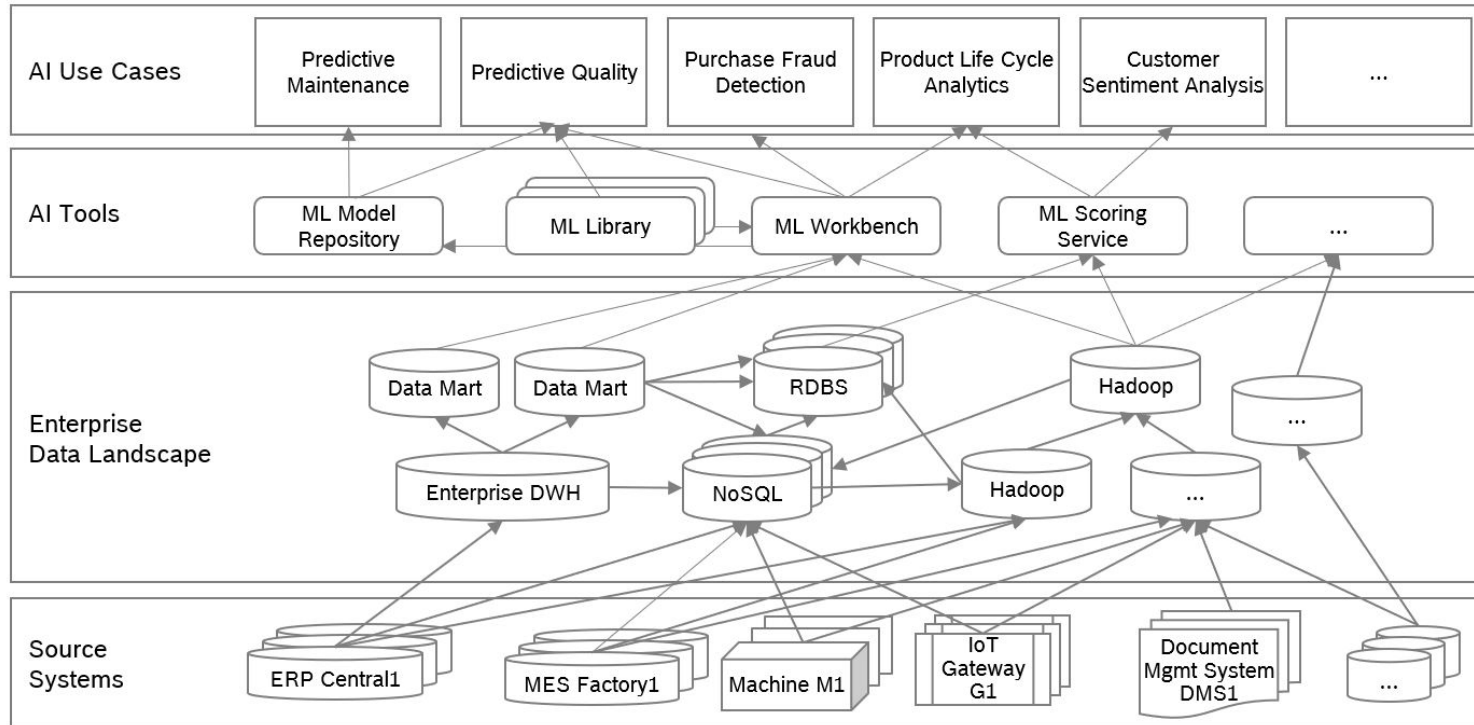


of executives acknowledge they know how to pilot, but struggle to scale AI across the business.

“[...] AI is currently done in an insular fashion leading to a polyglot and heterogeneous enterprise data landscape. This makes systematic data management, comprehensive data democratization and an overall data governance considerably challenging and prevents the wide-spread use of AI in industrial enterprises.”

Gröger (2021)

Insular AI and Enterprise Data Landscape

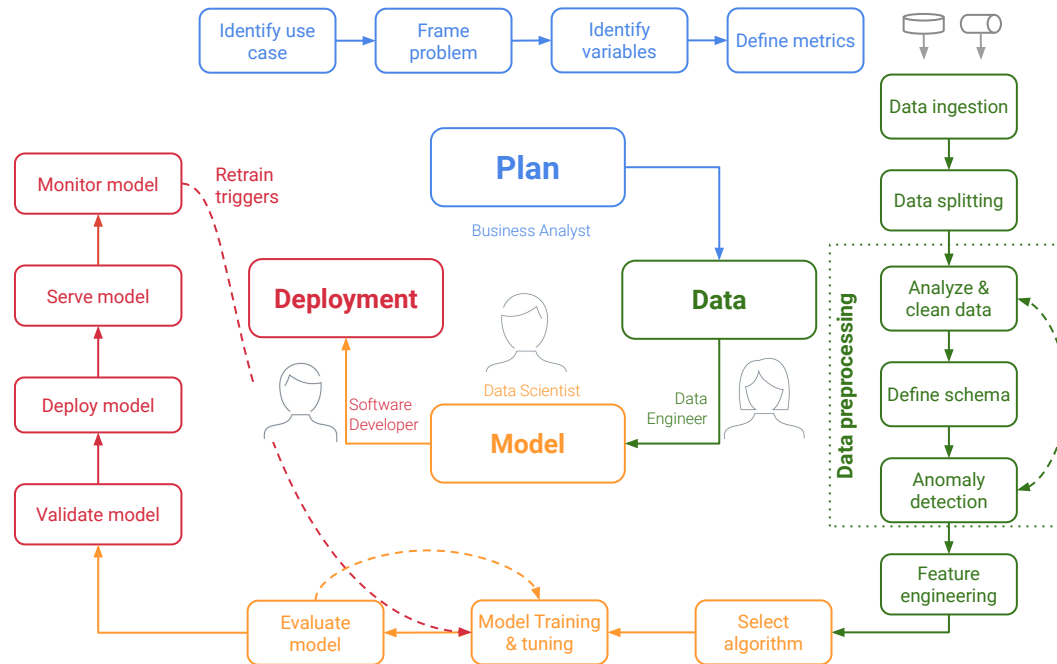


DWH: Data Warehouse, ERP: Enterprise Ressource Planning, MES: Manufacturing Execution System, ML: Machine Learning, RDBS: Relational Database System

Machine Learning Lifecycle, Data Platform Architectures & Machine Learning Operations (MLOps)

Lifecycle of an ML System

Plan | Data | Model | Deployment



Data platform architectures

Evolution

Lakehouse: A New Generation of Open Platforms that Unify Data Warehousing and Advanced Analytics

Michael Amme¹, Ali Ghodsi^{1,2}, Reynold Xin¹, Matei Zaharia¹
¹Databricks, ²UC Berkeley, ³Stanford University

Abstract

This paper argues that the data warehouse architecture on its own is not the answer. It will either be replaced by or be replaced by an architectural pattern, the Labhouse, which will (3) be based on a direct-access data format, such as Apache Parquet; (4) have full client support for machine learning and data science; and (5) fit the rest of the art performance. Labhouses can help address some major challenges with data warehouses, including data siloing, scalability, total cost of ownership, data lock in, and limited user support. We discuss how the industry is already moving toward Labhouses and how this shift may affect work in data management. We also report results from a Labhouse system using Parquet that is competitive with popular cloud data warehouses on TPC-DS.

quality and governance characteristics. In this architecture, a small number of data in the lake would likely be tiled to a distributed storage system (such as Hadoop) for the large majority of data requests, and the top 100 applications would be able to request data directly accessible to a wide range of other analytics engines, such as machine learning systems [30, 33, 42].

From HDFS onwards, cloud data lakes, such as S3, ADLS and GCS started replacing Hadoop. They have superior elasticity (allowing to store) georeplication, and most importantly, extremely low cost along with the possibility of automatic, even cheaper, archival storage (e.g., S3 Glacier). The rest of the architecture is largely the same as the cloud is in the second generations systems, with a distributed data warehouse such as Redshift or Snowflake. The two data lakes + warehouse architecture is now dominant in the industry (in our recent assessment based at virtually all Fortune 500 companies).

2. Reflection

This paper argues that the data warehouse architecture as we know it today will waste in the coming years and be replaced by a new architectural pattern, which we refer to as the *Lakehouse*, characterized by (1) open direct-access data formats, such as *Apache*

The history of data warehousing started with helping business leaders get analytical insights by collecting data from operations.

databases into centralized warehouses, which then could be used for decision support and business intelligence (BI). Data in the warehouses would be written with schema-on-write, which means that the data model was optimized for downstream BI consumption. We refer to this as the first generation data analytics platform.

A decade ago, the first generation systems started to face new challenges. First, they typically coupled compute and storage into monolithic appliances. This forced enterprises to provision and pay for the peak of user load and data under management, which became costly as datasets grew. Second, not only were datasets grown

equally, but more and more datasets were completely unstructured
e.g., video, audio, and text documents, which data warehouses could
not store and query at all.

To reduce storage problems, the several generations of data analysis platforms started offloading all the raw data into distributed low-cost storage systems with a flat API that hold data in generic and common open file formats, such as Apache Parquet and ORC [8, 9]. The approach started with the Apache Hadoop movement [13], using flat

This article is published under the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>). All journal components are available under the Creative Commons Attribution License.

Decision Systems Research, Vol. 23, January 20-23, 1991, Orlando.

<https://databricks.com>

[content/uploads](#)
[unknown.pdf](#)

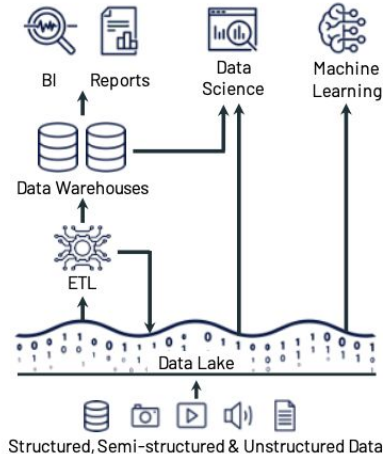
[kenhouse.pdf](#)

Source: Armbrust

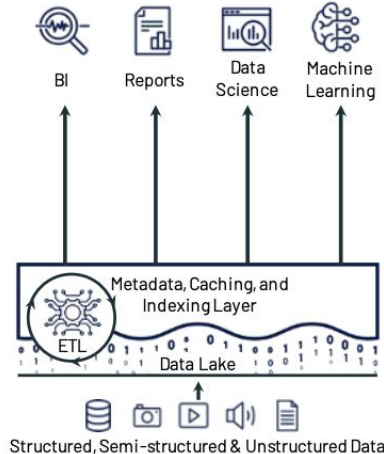
First generation platforms



Current two-tiers architectures



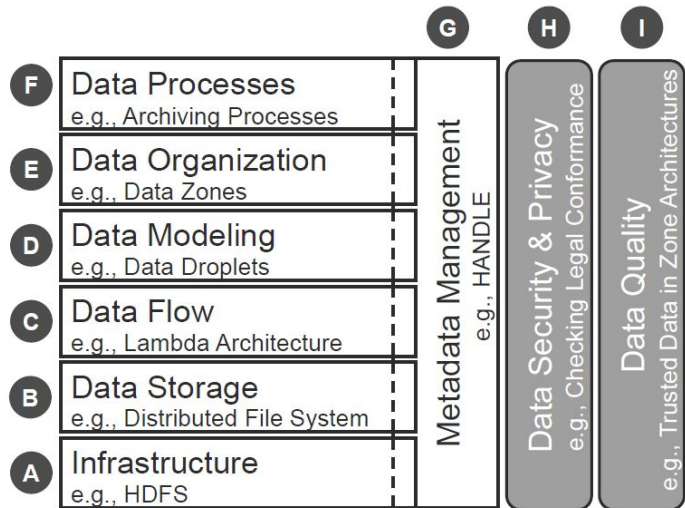
Lakehouse platform



Data platform architectures

Architecture aspects

Aspects



Conceptual and physical



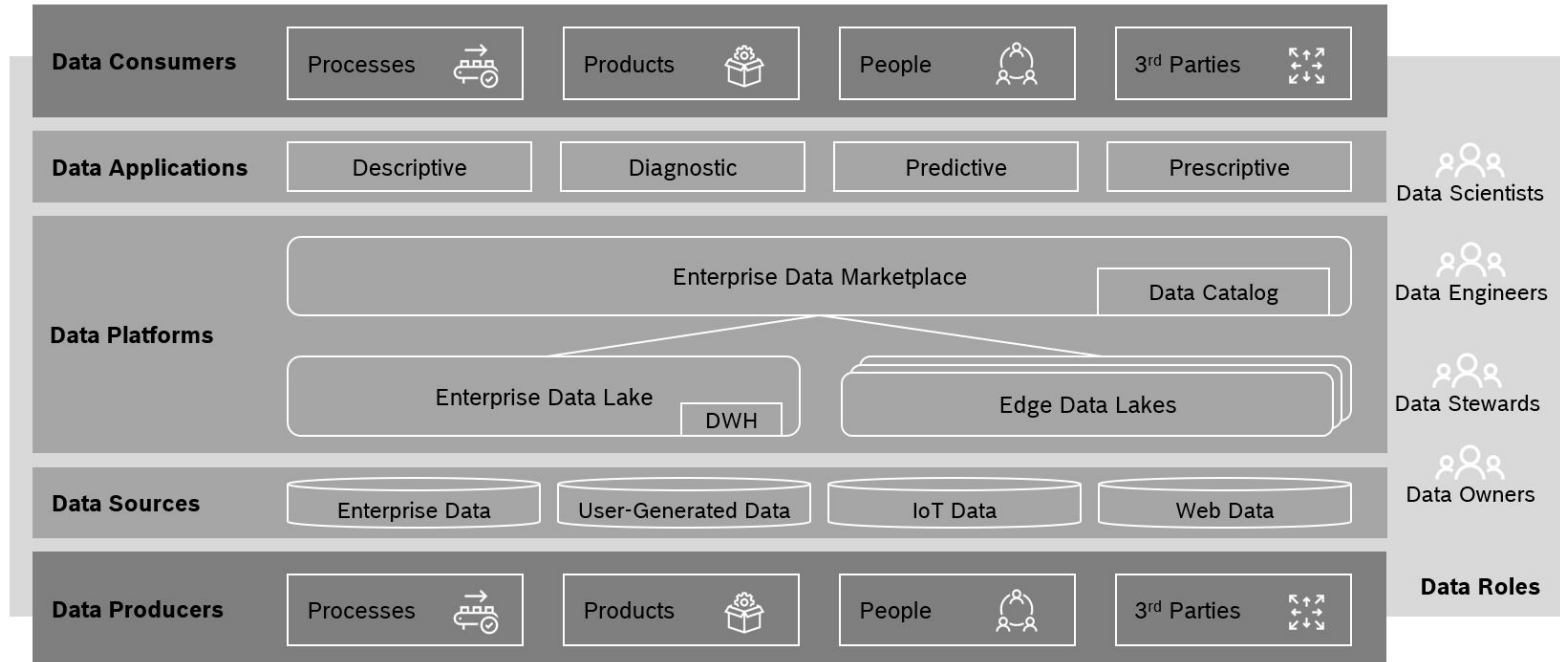
Only conceptual, implementation through individual layers

Example

DLAF Layer	AIRPORTS DL [Ma17a]
A. Infrastructure	Hadoop (HDFS, MapReduce), Apache Flume, Apache Spark, Apache Oozie, Apache Pig, Apache Atlas, R Studio, Shiny, Apache Sqoop
B. Data Storage	Single File System
C. Data Flow	Data are ingested as streams, but processed as batches
D. Data Modeling	Raw Messages, AIRPORTS Data Model
E. Data Organization	Four Zone Architecture
F. Data Processes	Processing Pipeline for Messages (ETL Processes), Processes for Ingestion and Use
G. Metadata Management	Managed by Apache Atlas
H. Data Security & Privacy	Tracking manipulation of data
I. Data Quality	Tracking manipulation of data, Quality through Zones

Data platforms

for industrial enterprises



Core elements of a data ecosystem for industrial enterprises

ML platform @ Spotify

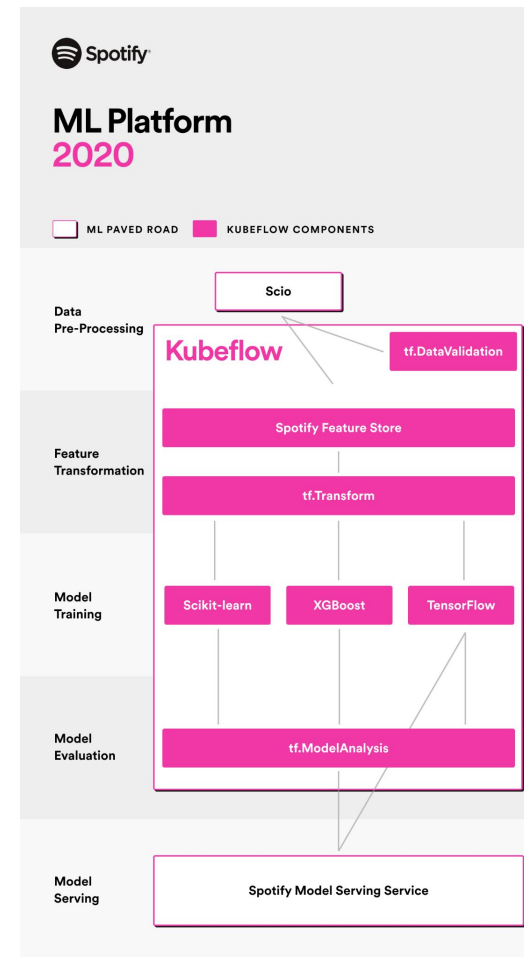
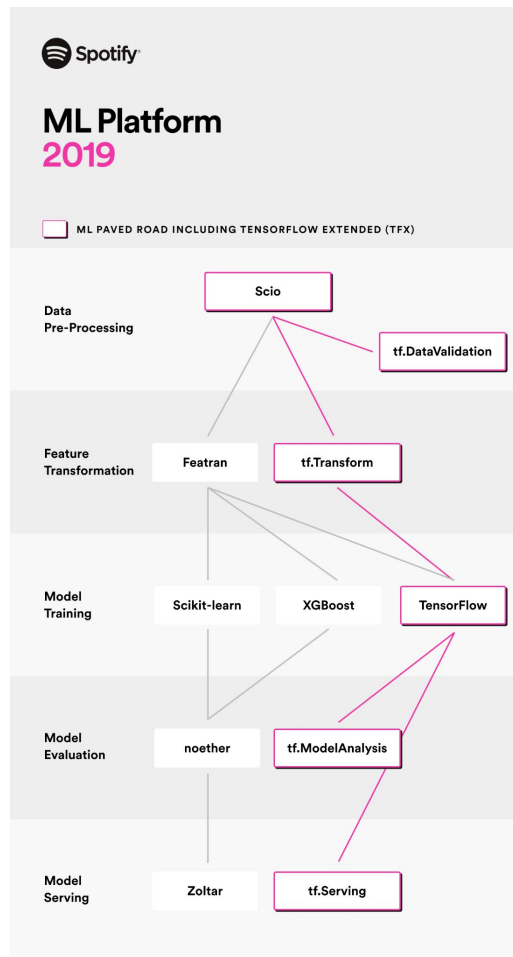
Reference architecture

The Winding Road to Better Machine Learning Infrastructure Through Tensorflow Extended and Kubeflow



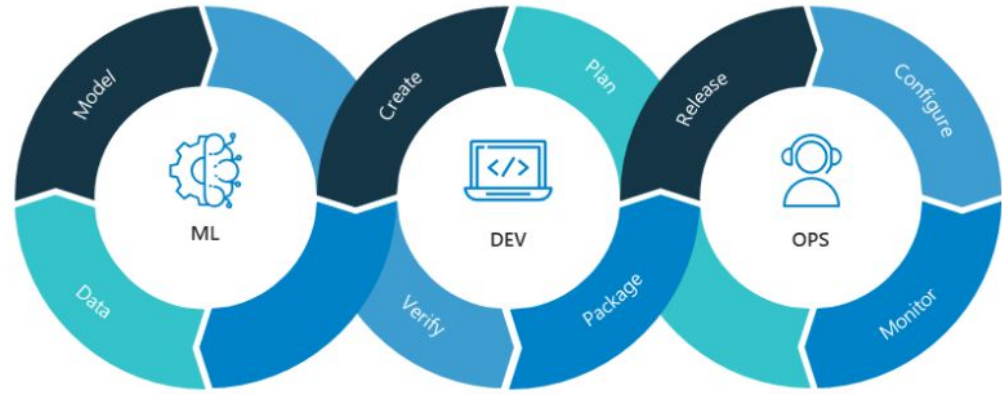
December 13, 2019
Published by Josh Baer, Samuel Ngahane

<https://engineering.spotify.com/2019/12/13/the-winding-road-to-better-machine-learning-infrastructure-through-tensorflow-extended-and-kubeflow/>



Machine learning operations (MLOps)

- ML Engineering culture and practice that aims at **unifying** ML System **development** (Dev) and ML system **operations** (Ops)
- Tools and principles to support workflow **standardization** and **automation** through the ML system lifecycle.

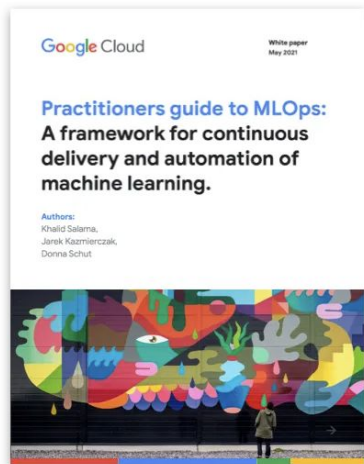


Source: Nvidia (2021) <https://blogs.nvidia.com/blog/2020/09/03/what-is-mlops/>

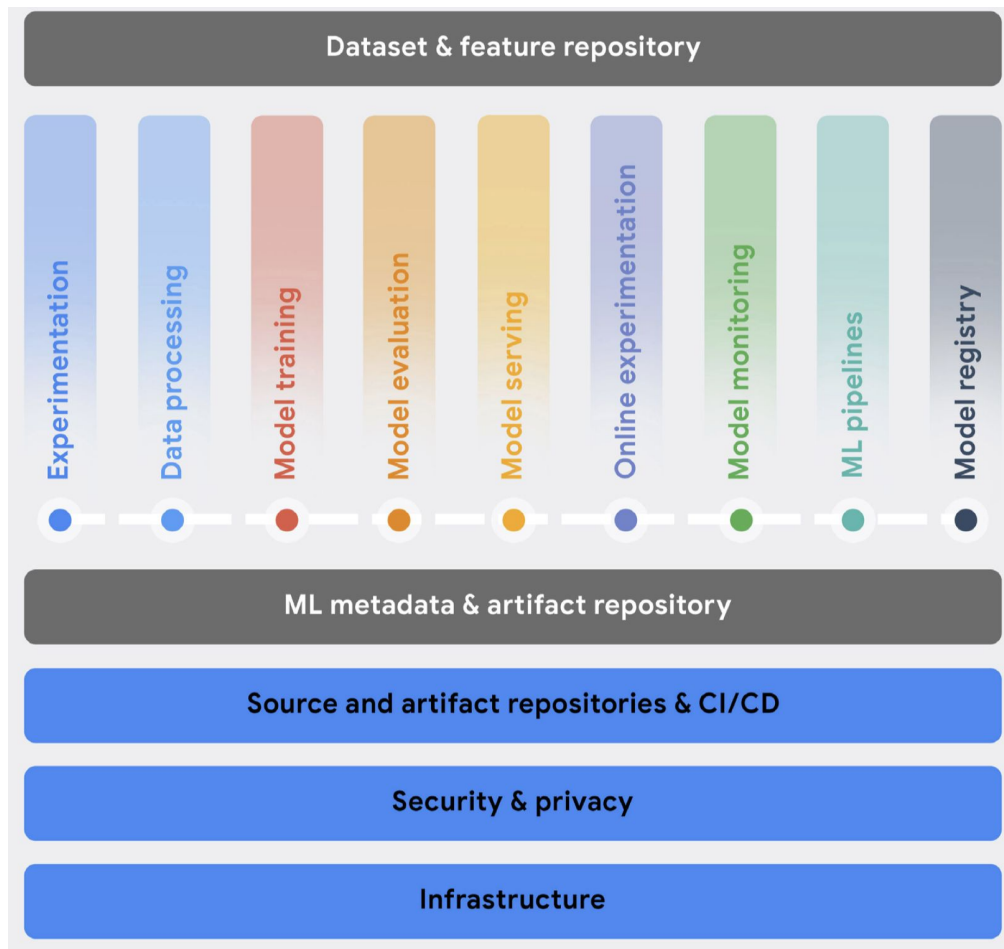
MLOps

Introduction

Learn the basics of MLOps

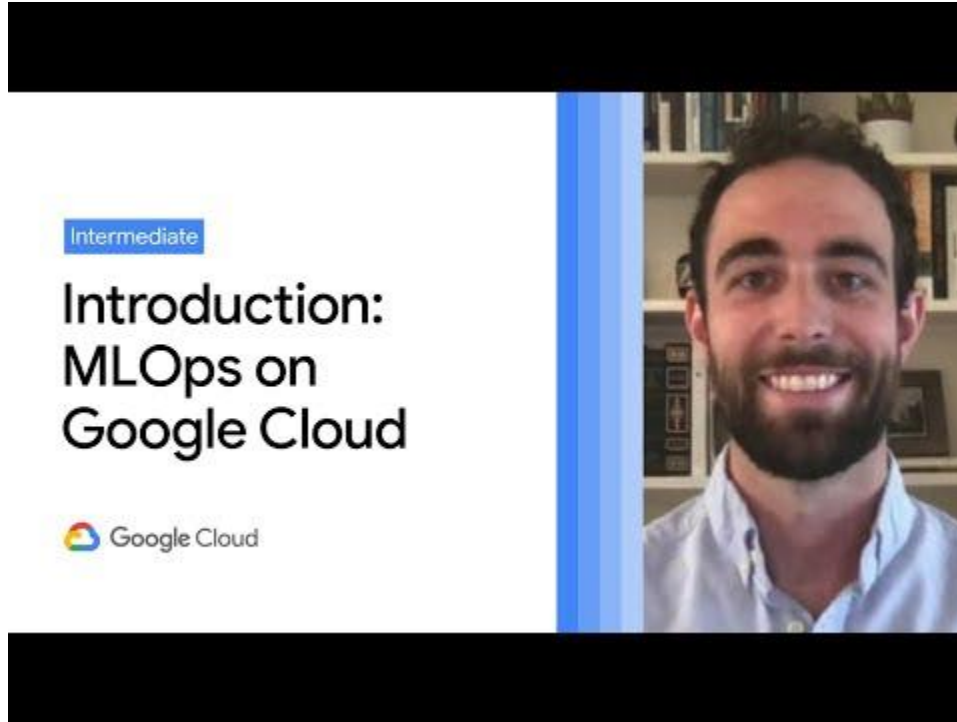


<https://cloud.google.com/resources/mlops-whitepaper>



MLOps Core technical capabilities

Introduction to Machine Learning Operations

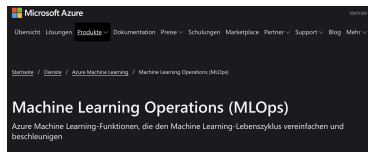


MLOps demo (including Kubeflow)

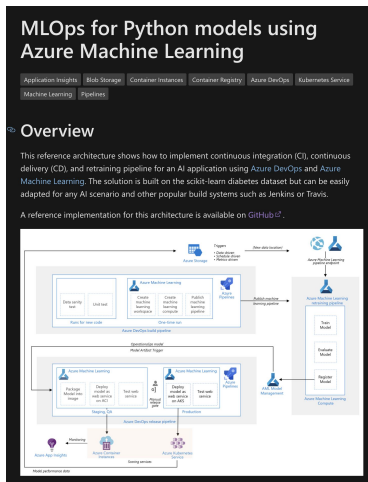


MLOps in Azure

Reference architectures



<https://azure.microsoft.com/de-de/services/machine-learning/mlops/>



<https://docs.microsoft.com/en-us/azure/architecture/reference-architectures/ai/mlops-python>

5 Best Practices to optimize your MLOps lifecycle on Azure:

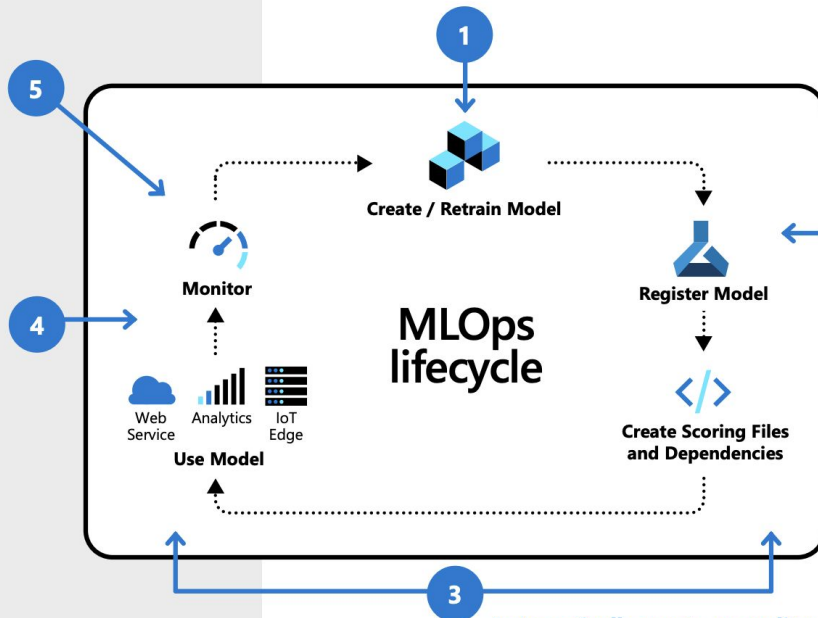
Observe data drift and feed back model information to improve future training.

Deploy and monitor performance so you can release models with confidence and know when to retrain.

Create models with reusable ML pipelines using the Azure Machine Learning extension for Azure DevOps. Store your code in GitHub so it automatically integrates into your MLOps pipeline.

Automate your MLOps rollout using Azure DevOps + Azure Machine Learning for version models with rich metadata and event management.

Automatically create an audit trail for all artifacts in your MLOps pipeline ensure asset integrity and meet regulatory requirements.



Resources

- Armbrust, M., Ghodsi, A., Xin, R. & Zaharia, M. (2021). Lakehouse: A New Generation of Open Platforms that Unify Data Warehousing and Advanced Analytics. 11th Annual Conference on Innovative Data Systems Research (CIDR '21), January 11–15, 2021, Online https://databricks.com/de/wp-content/uploads/2020/12/cidr_lakehouse.pdf
- Baer, J. & Samuel Ngahane, S. (2019). The Winding Road to Better Machine Learning Infrastructure Through Tensorflow Extended and Kubeflow. <https://engineering.atspotify.com/2019/12/13/the-winding-road-to-better-machine-learning-infrastructure-through-tensorflow-extended-and-kubeflow/>
- Giebler, C., Gröger, C., Hoos, E., Schwarz, H., Mitschang, B. The Data Lake Architecture Framework: A Foundation for Building a Comprehensive Data Lake Architecture. In: Proceedings of the 19th Conference for Database Systems for Business, Technology and Web (BTW), pp. 351-370. <https://dx.doi.org/10.18428/btw2021-19>
- Gröger, C. (2021). There Is No AI Without Data. Industry Experiences on the Data Challenges of AI and Call for a Data Ecosystem for Industrial Enterprises. In: Communications of the ACM, 2021, to appear. http://christophgroeger.de/download/Groeger_There_Is_No_AI_Without_Data.pdf
- Kirenz, J. (2021). MLOps with Tensor Flow Extended & Kubeflow. Presentation <https://www.slideshare.net/1/8/mlops-with-tensorflow-extended-and-kubeflow>
- Salama, K., Kazmierczak, J. & Schut, D. (2021). Practitioners guide to MLOps: A framework for continuous delivery and automation of machine learning. Google Whitepaper. https://services.google.com/fh/files/misc/practitioners_guide_to_mlops_whitepaper.pdf