

The Team & Agenda



Camus Ma

AI/ML Lead (Google Cloud)



Ayan Kar

EMEA Head of Data & AI (Capgemini)



Luc Schamhart

Account Manager Nationale-Nederlanden



Turan Bulmus

AI/ML Practice Lead Benelux



Marijn Markus

AI/ML Consultant (Capgemini)



Joost Carpaij

Insight/Data Consultant (Capgemini)

1

10:00 - 10:25 - Google Cloud AI
Adoption Framework - Camus

2

10:25 - 10:50 - Applying AI
governance & ML Opps.&
Glassbox Demo - Ayan & Bikas

3

10:50 - 11:00 - Stretch legs

4

11:00 - 11:45 - Google Cloud Vertex
Framework -Turan (incl. demo)

5

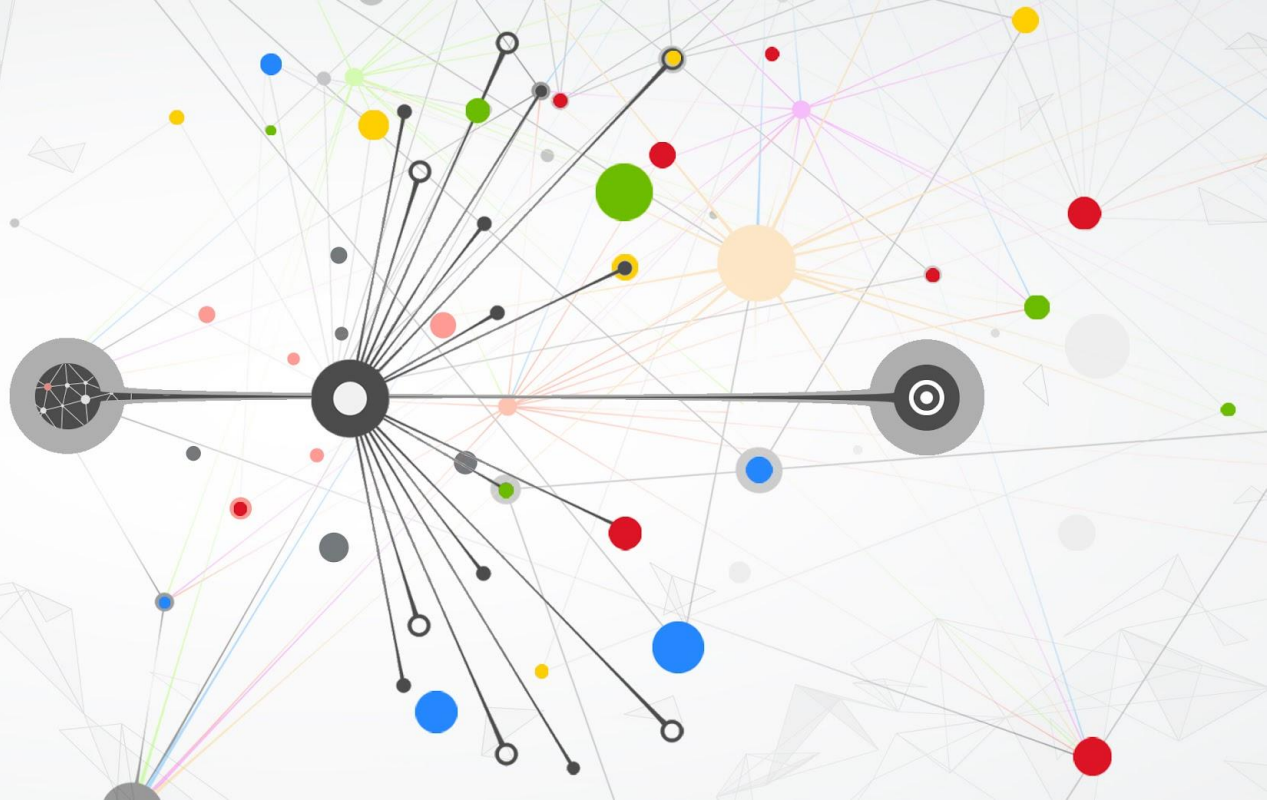
11:45 - 11:55 - Experience @ NN
(Life & Pensions) Joost

6

11:55 - 12:20 - Cool Demo's Marijn

12:20 - 13:00 - Lunch

AI



Google Cloud's AI Adoption Framework



Gain a competitive advantage through AI

Enterprises that invest in building industry-specific AI solutions are proven to be better positioned as future global economic leaders

Companies that fully absorb AI could double their cash flow

2X more
data-driven
decisions

5X faster
decisions
than others

3X faster
execution

Sources: [McKinsey](#) 2018 and [MIT Tech Review](#) 2017

The Cloud AI Adoption Framework

- A **guiding framework** for leaders who want to leverage the power of AI to transform their business
- A **tool to assess** where you are in your journey and **define where you want to be**



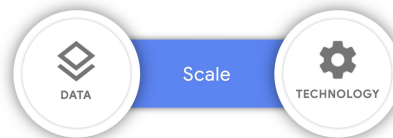
AI Maturity Themes



How are the teams structured?

What is the level of executive sponsorship?

How is budget allocated for AI/ML projects?



How are cloud-based services provisioned?

How is capacity for workloads allocated?

Does an organization use accelerators?



What the data and ML skill sets are required in the organisation?

How does an organization develop ML talent?



What controls are in place?

How does an organisation establish trust in it's AI capability?

Can you explain the decisions made by your AI systems?



How are datasets created, curated, and annotated?

Are they discoverable and reusable?

How are data and ML assets managed?

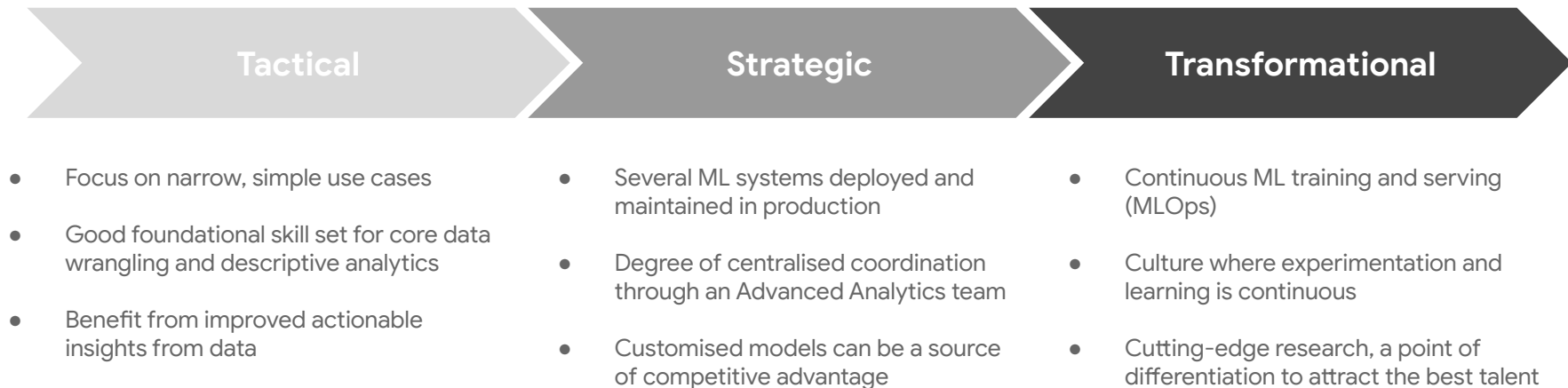


How models are continuously trained and deployed for serving?

How are model updates managed?

What ML quality control are in place?

AI Maturity Phases



The Cloud AI Maturity Scale

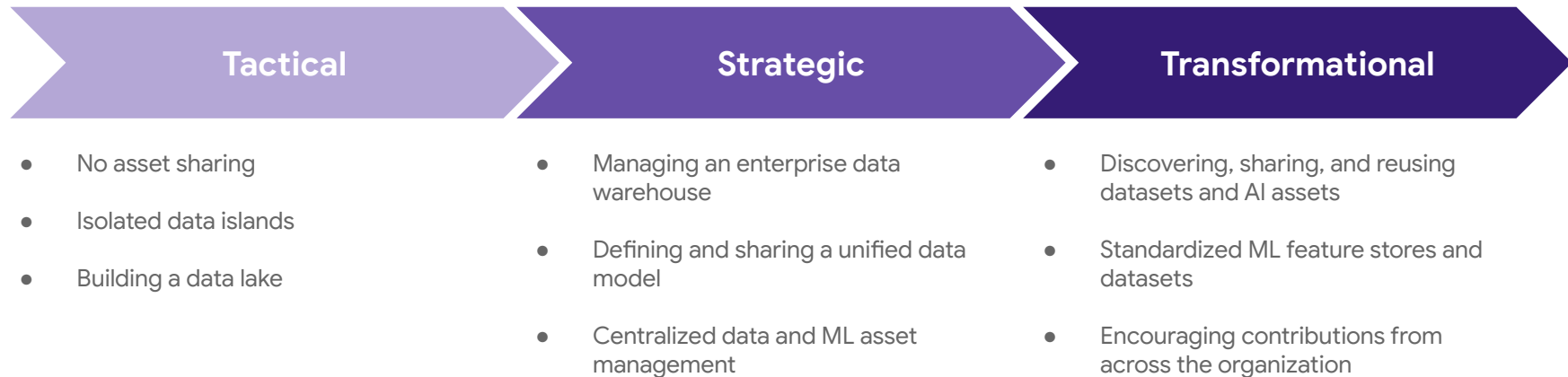
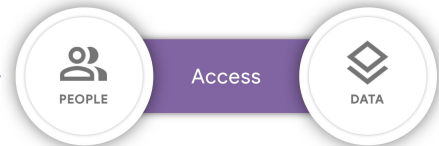
Lead

The theme concerns the extent to which your data scientists are supported by a mandate from leadership to apply ML to business use cases, and the degree to which the data scientists are cross-functional, collaborative, and self-motivated.



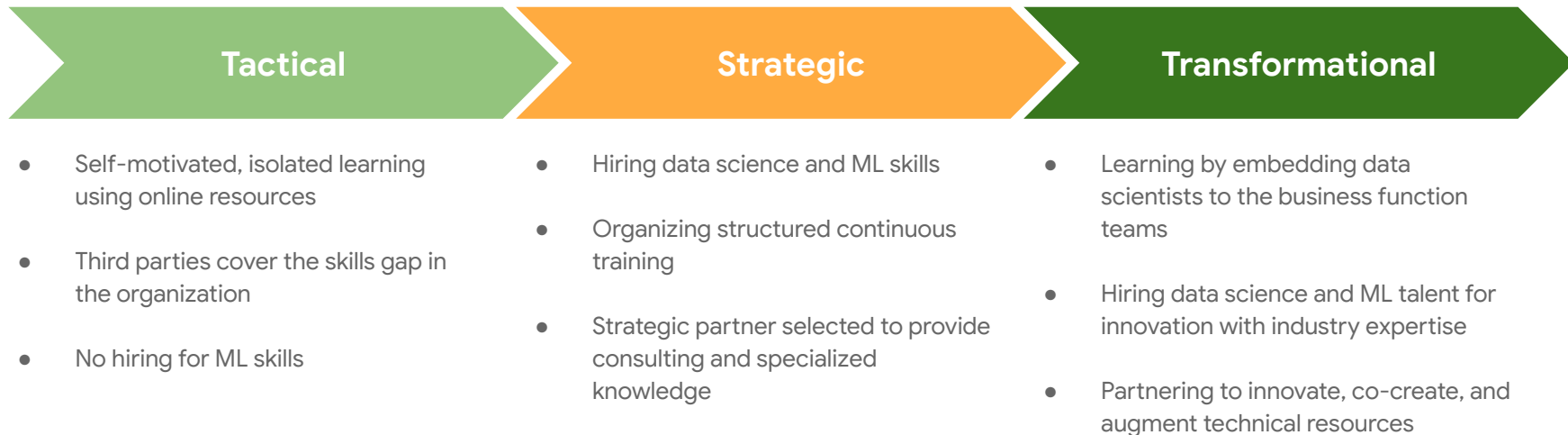
Access

The theme concerns the extent to which your organization recognizes data management as a key element to enable AI and the degree to which data scientists can share, discover, and reuse data and other ML assets.



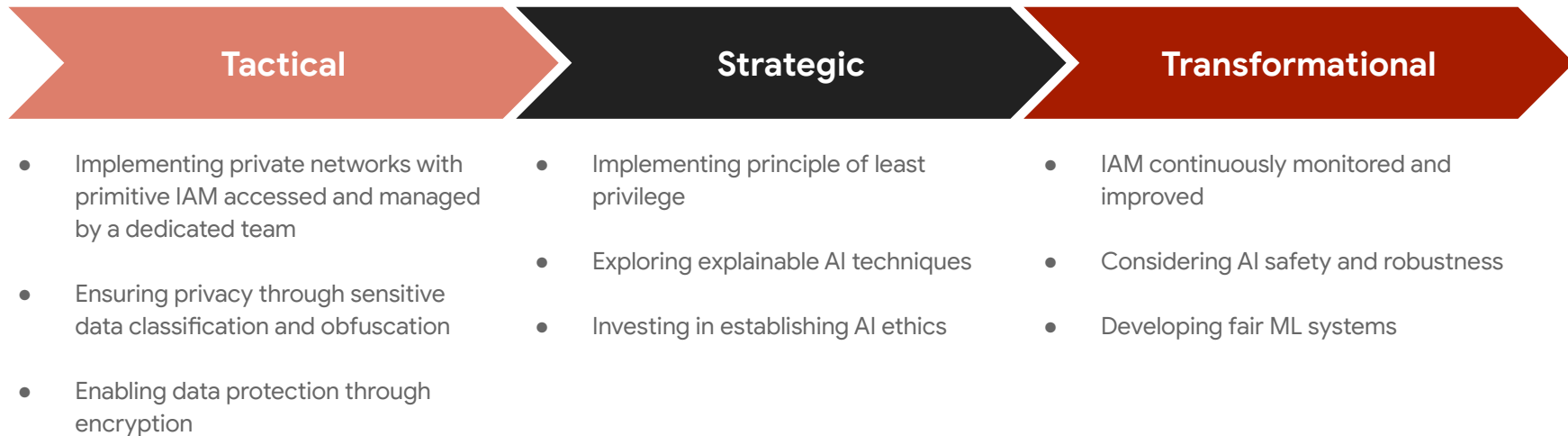
Learn

The theme concerns the quality and scale of learning programs to upskill your staff, hire outside talent, and augment your data science and ML engineering staff with experienced partners.



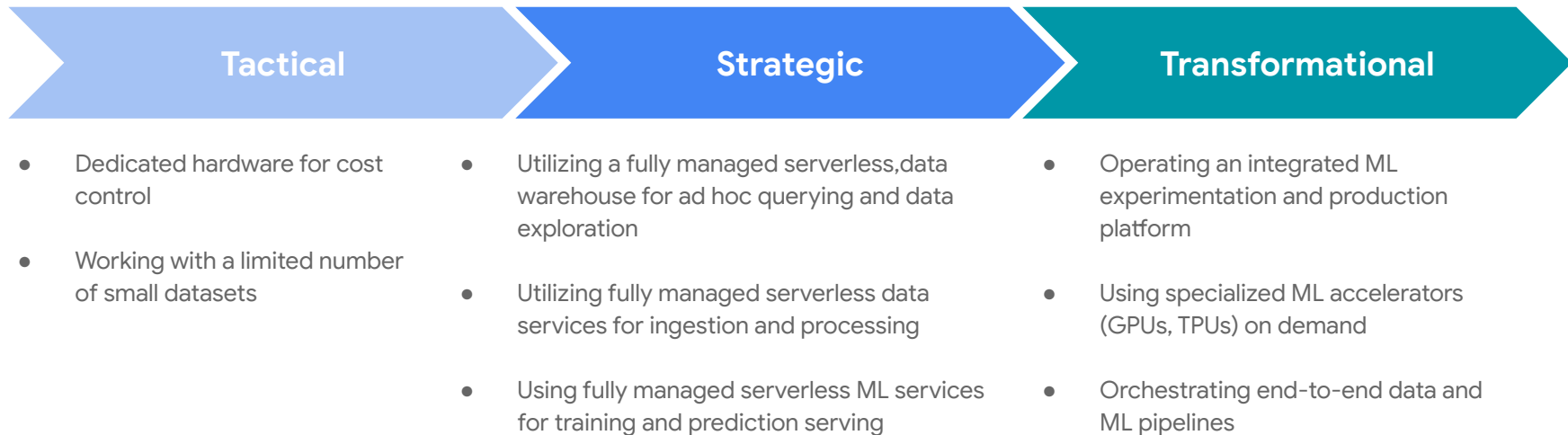
Secure

The theme concerns the extent to which you understand and protect your data and ML services from unauthorized and inappropriate access, in addition to ensuring responsible and explainable AI.



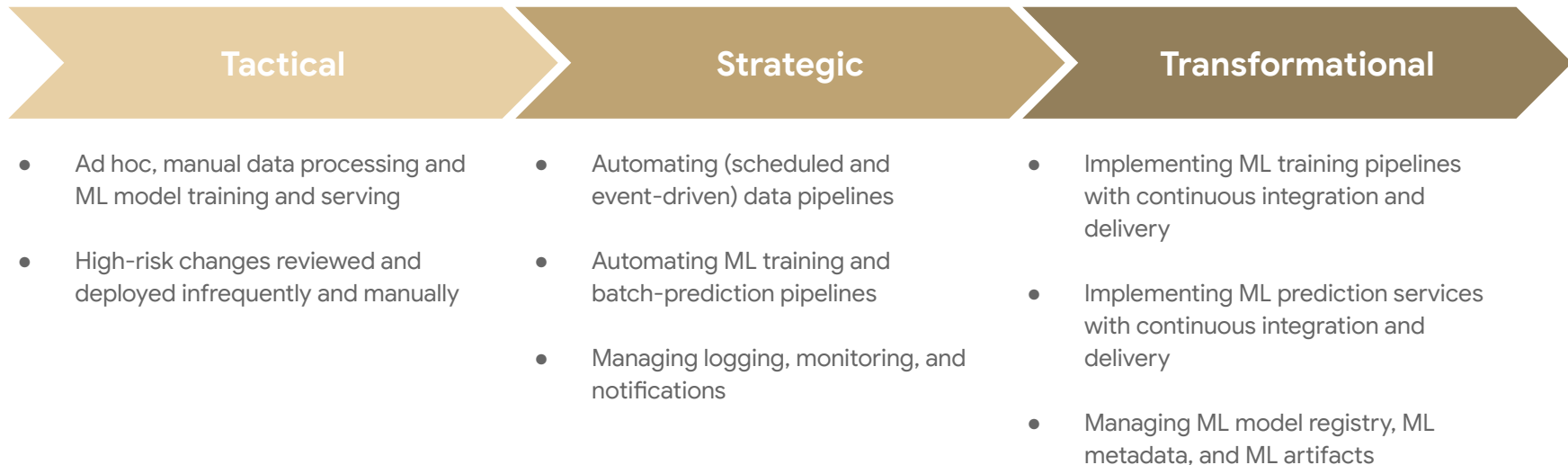
Scale

The theme concerns the extent to which you use cloud-native ML services that scale with large amounts of data and large numbers of data processing and ML jobs, with reduced operational overhead.



Automate

The theme concerns the extent to which you are able to deploy, execute, and operate technology for data processing and ML pipelines in production efficiently, frequently, and reliably.



Next Steps

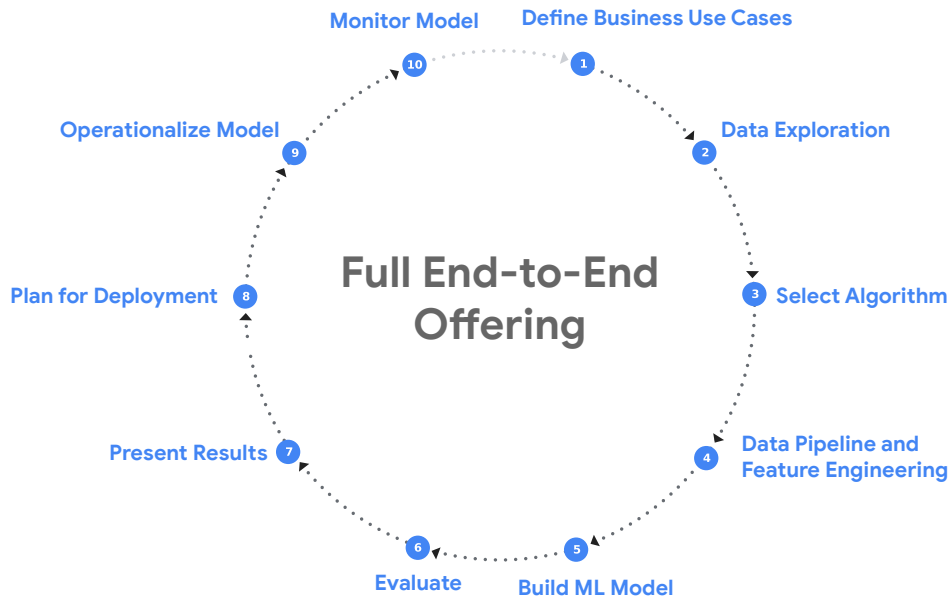
- **Understand where you are** - complete the maturity assessment
- **Set your goal** - where you do you want to go?
- **Create a group of leaders** - who are responsible for building your AI capability?
- **Devise a strategy** - based on the gaps, establish a plan for evolving your AI capability
- **We are here to help you every step of the way** - speak to your account representative about how you can engage Google

Questions?

Appendix

The transformational AI journey offering

AI Services provides a **complete end-to-end path** to accelerate your AI transformation



Address your business challenge
from start to finish



Build for production from the beginning

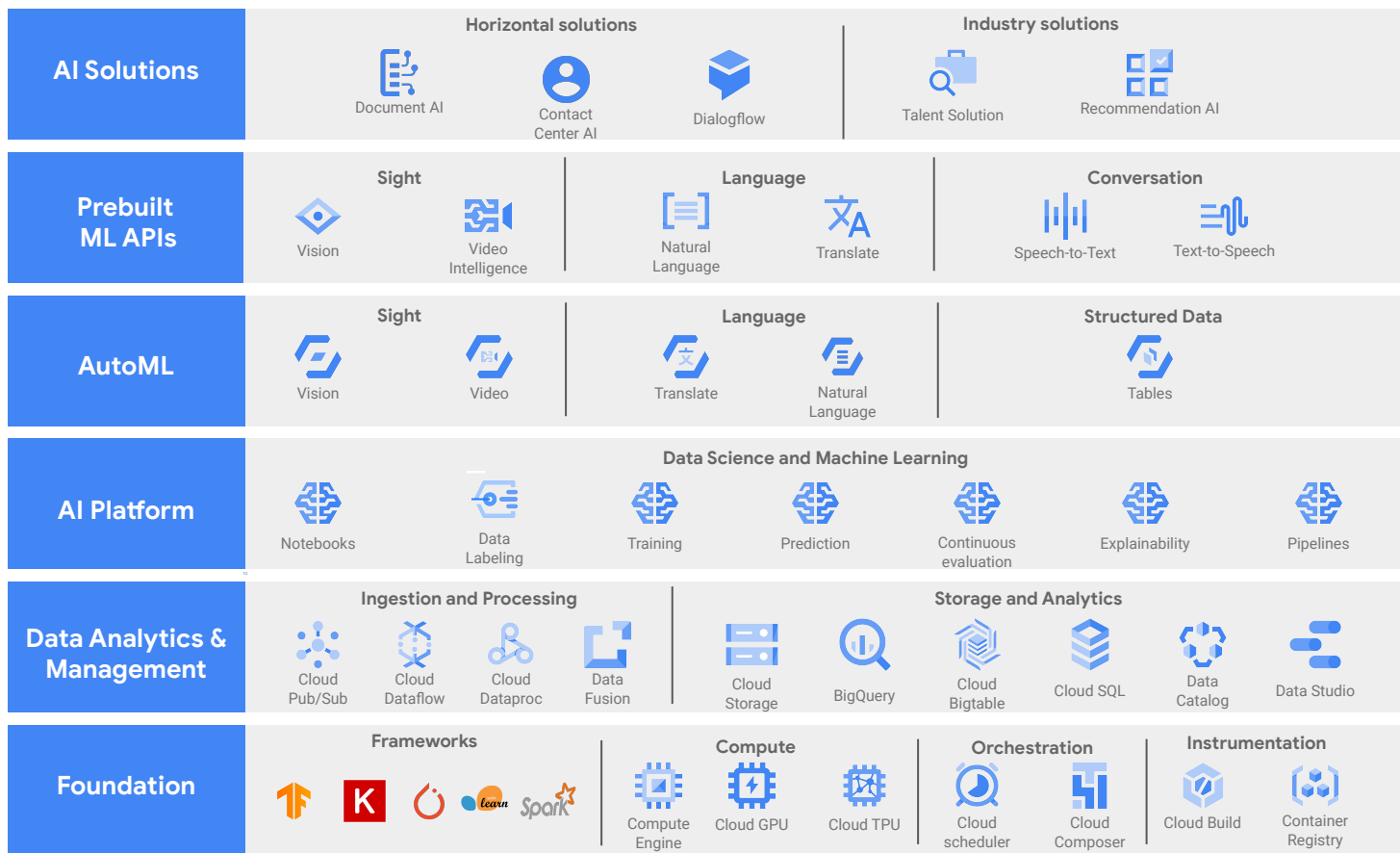


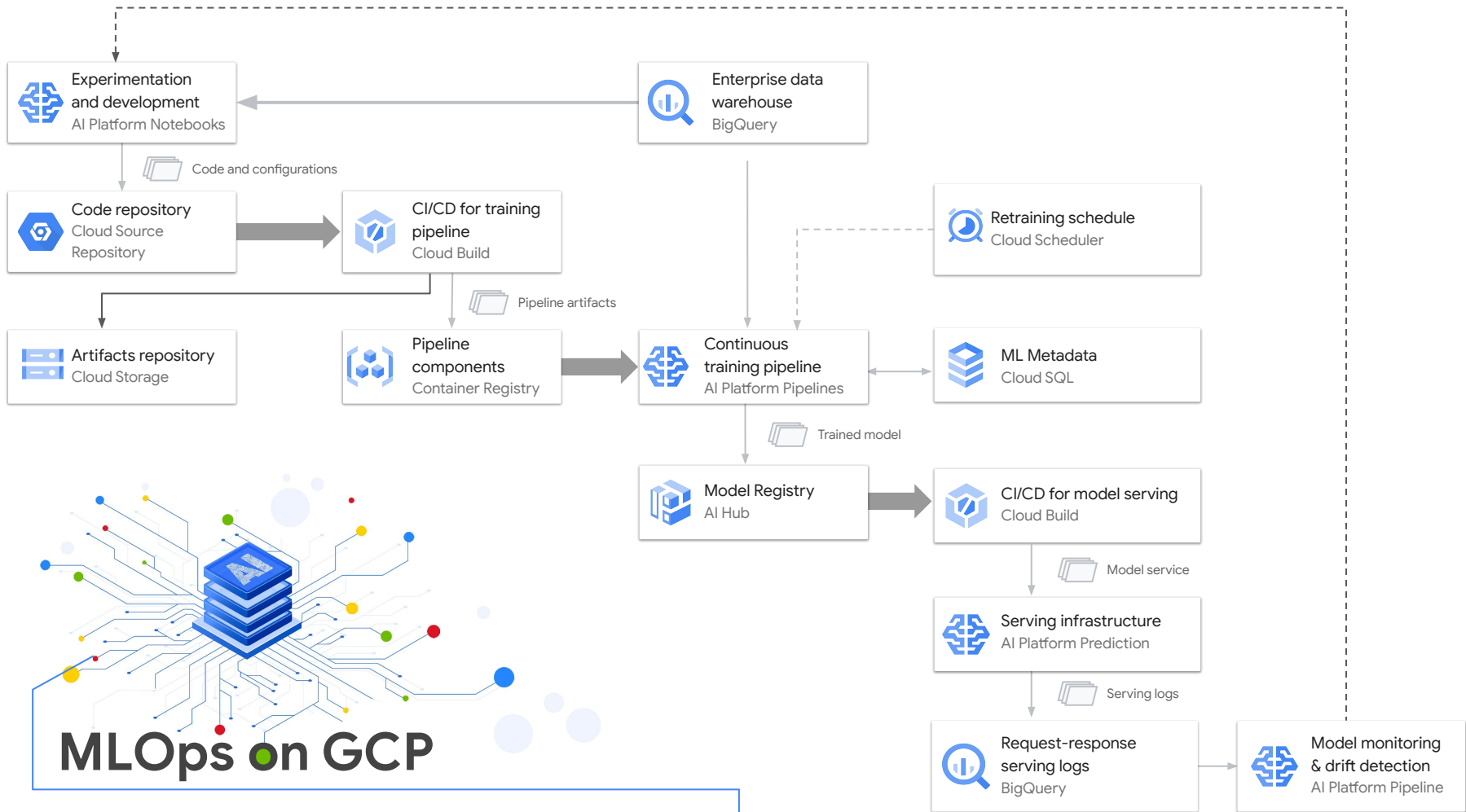
Deep understanding of the problem and solution



Trusted partner throughout the journey

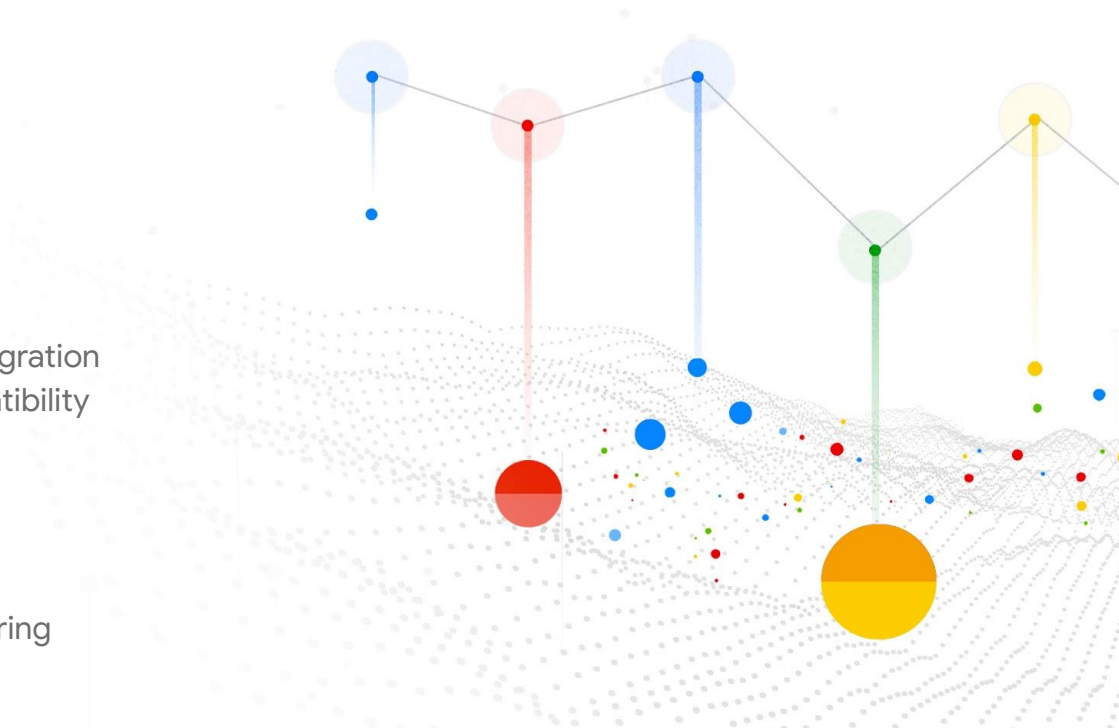
Google Cloud Smart Analytics & AI





ML Quality Control

- **Testing in development**
 - ✓ Data and feature testing
 - ✓ Model testing and debugging
 - ✓ Model evaluation
- **Testing in deployment**
 - ✓ Testing ML pipeline components integration
 - ✓ Validate model-infrastructure compatibility
 - ✓ Test model API
- **Testing in production**
 - ✓ Pipeline data and model validation
 - ✓ A/B testing and performance monitoring
 - ✓ Data drift and shift detection





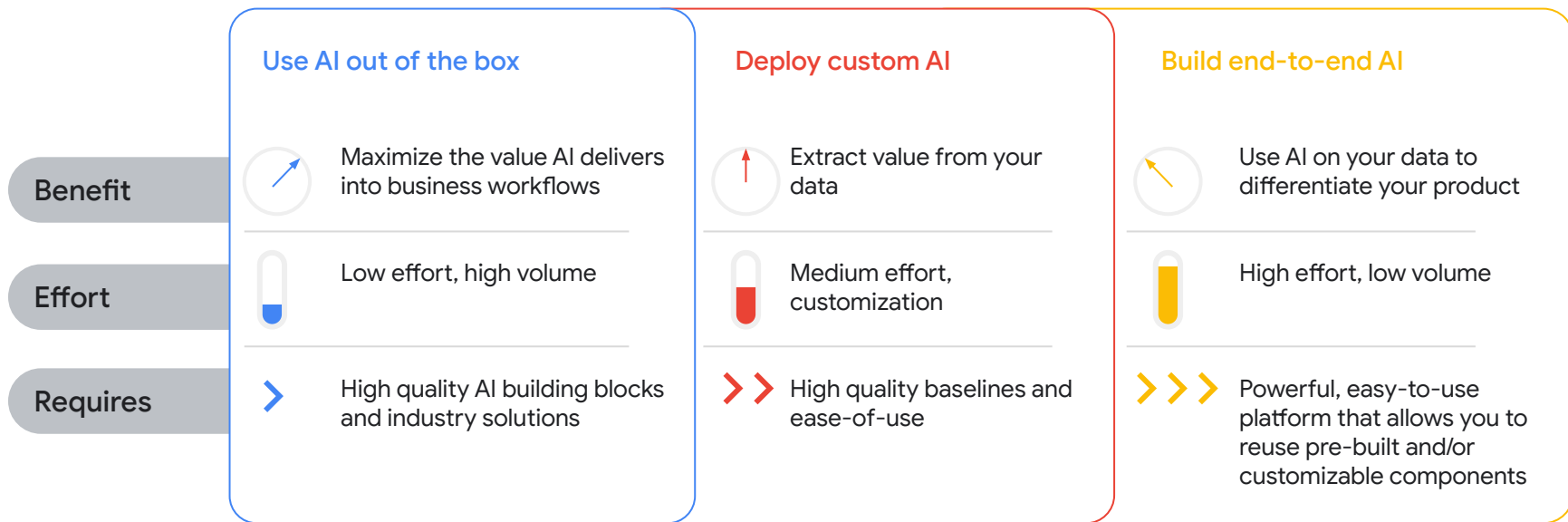
Machine Learning and MLOps on Google Cloud

Turan Bulmus
AI/ML Practice Lead Benelux
April 2022



Enable every company to be an AI company by reducing the challenges of AI model creation down to only the steps that require human judgement or creativity.

Build a portfolio of AI use cases



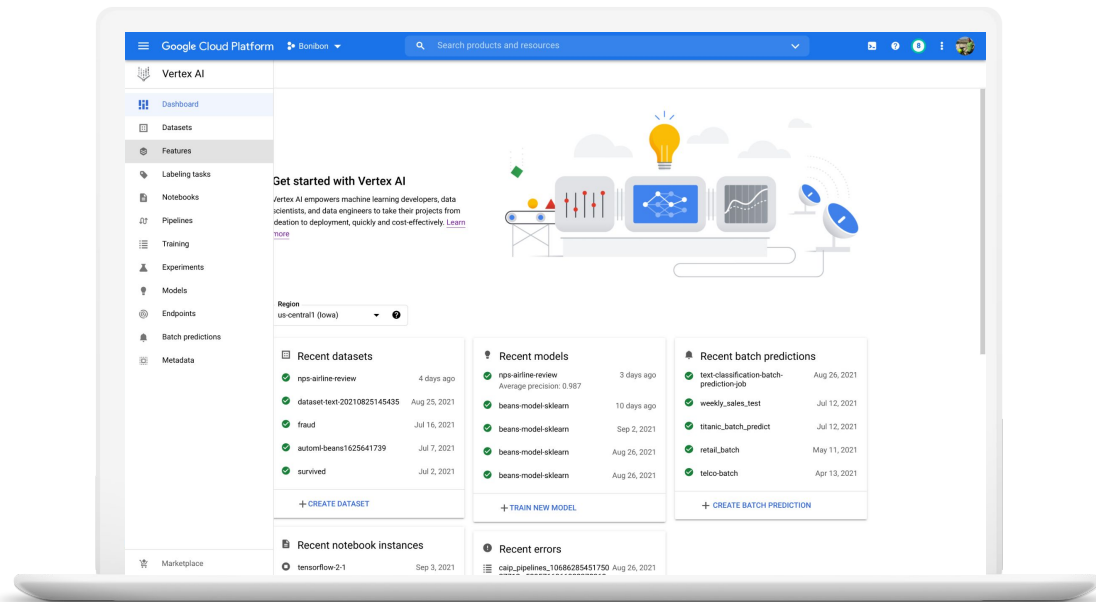
“Business value” generated from all three buckets
Need a unified platform that supports all three buckets

A Unified ML Platform for Solving All Business Problems

Processing all sources of data including images, documents, tables, video

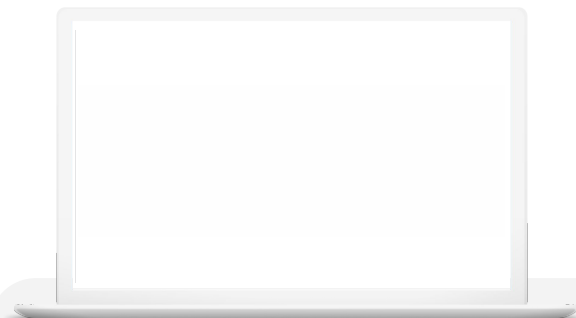


- One unified experience to create, deploy, and manage models over time, at scale
- Tools for all levels of expertise and for all types of data
- Accuracy and fairness of predictions and resulting decisions
- Flexible and secure



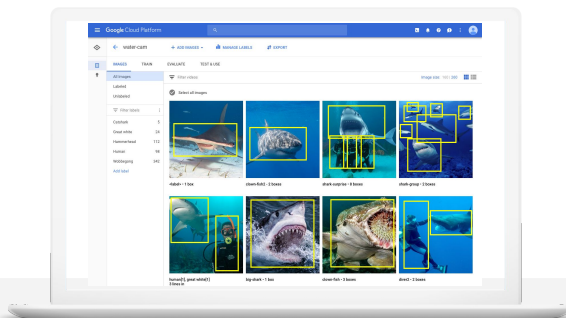
How to design your ML workflow?

Freedom of choice from no/low code to custom code



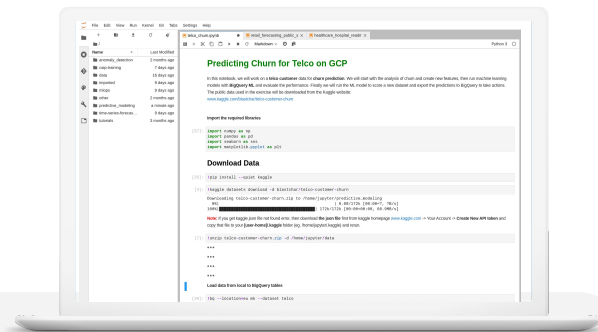
BigQuery ML

- Descriptive and predictive modeling on structured data
- Hyper-parameter tuning
- Feature engineering
- Explainability
- Simple SQL code



AutoML Models in Vertex AI

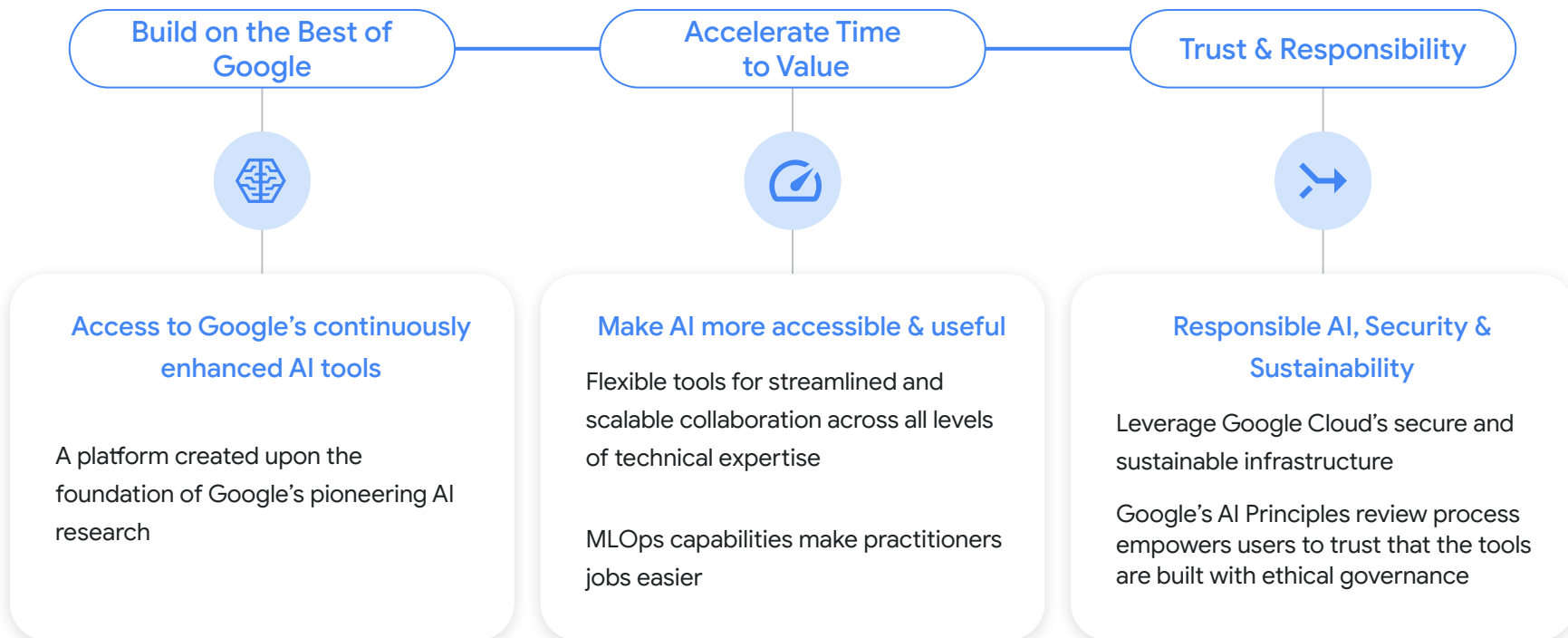
- Predictive modeling on structured & unstructured data
- Hyper-parameter tuning
- Feature engineering
- Explainability
- No code



End-to-end AI with Vertex AI

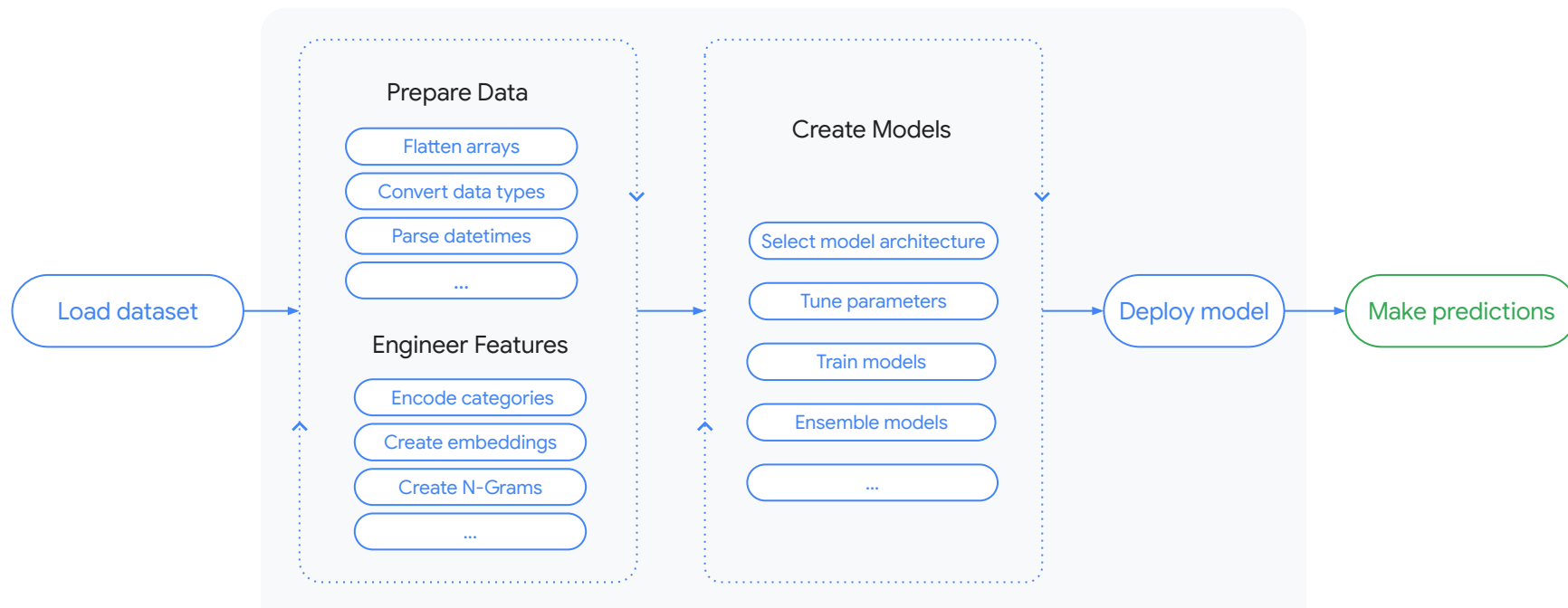
- Custom models on pre-built frameworks
- Noops, serverless training with hyperparameter tuning
- Explainability
- Custom code

Why organizations choose Vertex AI



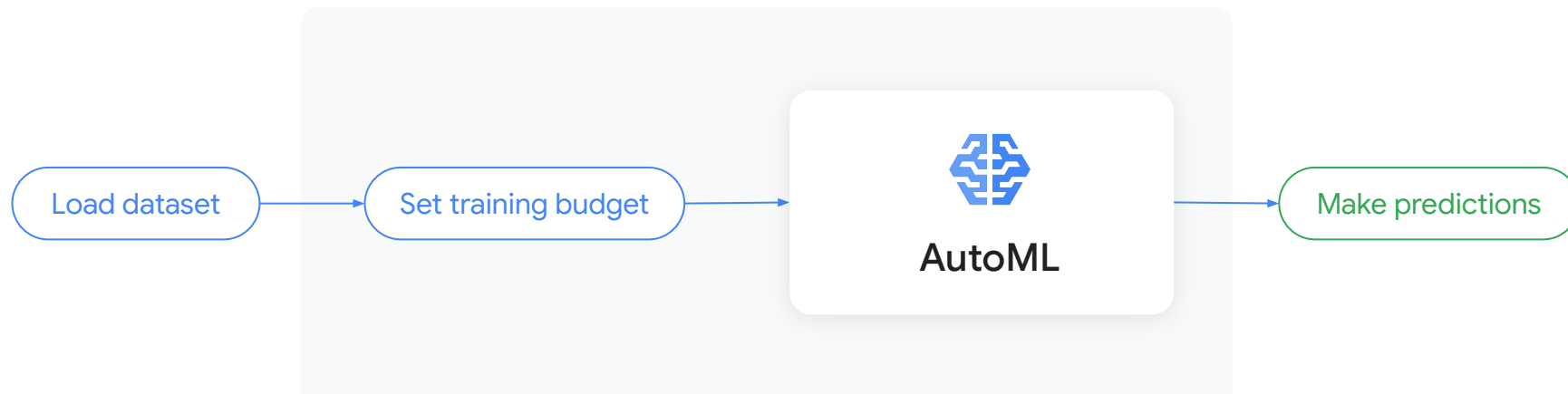
AutoML - Fastest path from data to value

Traditional Machine Learning Workflow



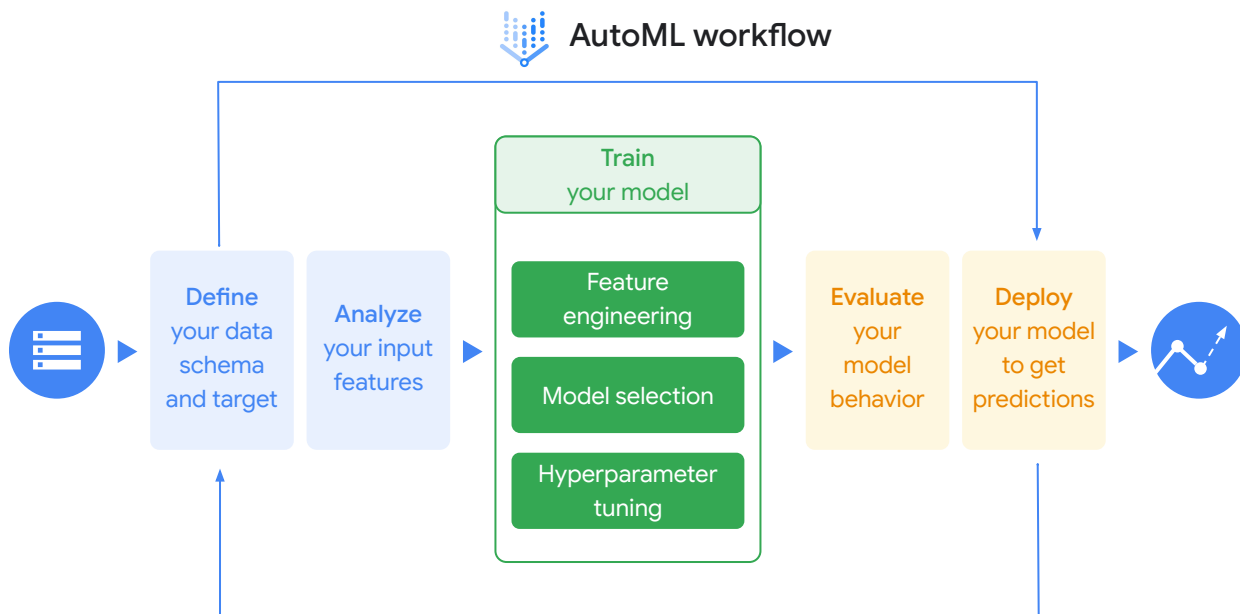
AutoML - Fastest path from data to value

AutoML Workflow



Low/No code

Point and click to build custom, high-quality models using the **AutoML** workflow in **Vertex AI**



Automatically search through Google's whole model zoo...

Linear, logistic

Feedforward DNN

Wide and Deep NN

Gradient Boosted Decision Tree (GBDT)

DNN + GBDT Hybrid

Adanet ensemble




Neural + Tree Architecture Search

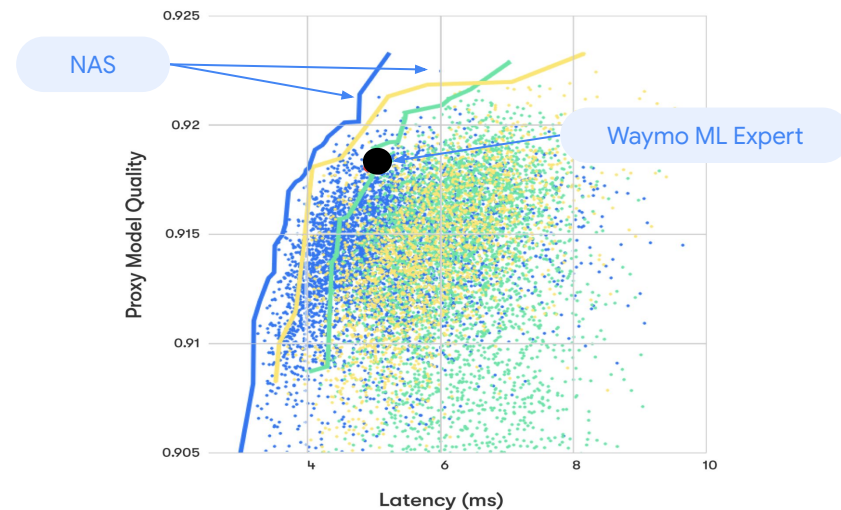
...and more!

Access to Google's best-in-class algorithms like **NAS**

Use of **Neural Architecture Search** (NAS) at Waymo

"Going from months of engineering time to generate and fine tune a architecture manually to "automatically generating" neural nets with NAS"

-  20–30% lower latency/same quality
-  8–10% lower error rate/same latency
-  NAS model in 2 weeks vs months (1 year of GPU time) searching over 10k architectures



MLOps

A set of **standardized** processes and technology capabilities for building, deploying, and operationalizing ML systems **rapidly** and **reliably**



Vertex AI

Applications

Vision and Video

Conversation

Language

Structured Data

Custom machine learning

Workbench

AutoML

NAS

Prediction

ML Metadata

Data Labeling

Training

Explainable AI

Feature Store

Model Monitoring

Experiments

Vizier (Optimization)

Pipelines

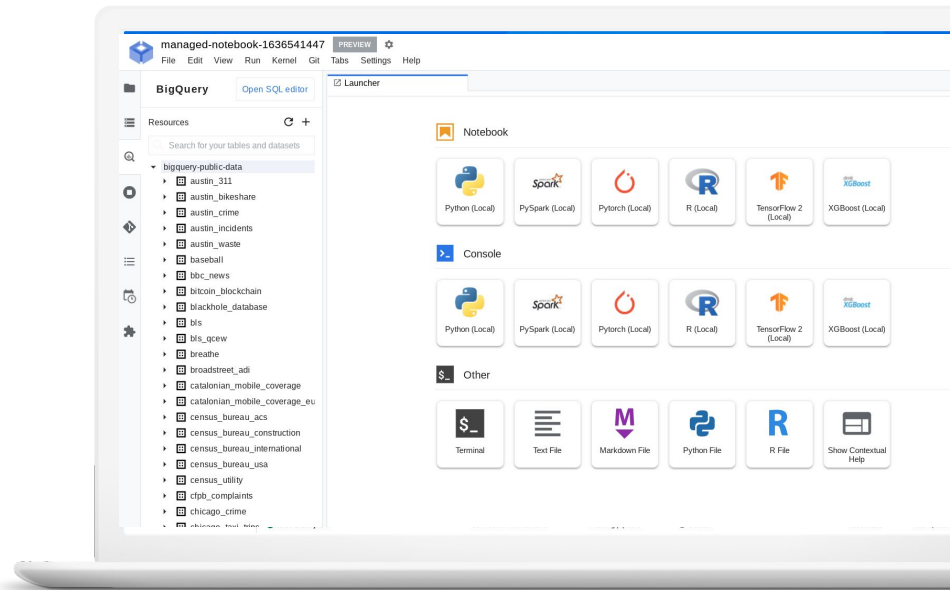
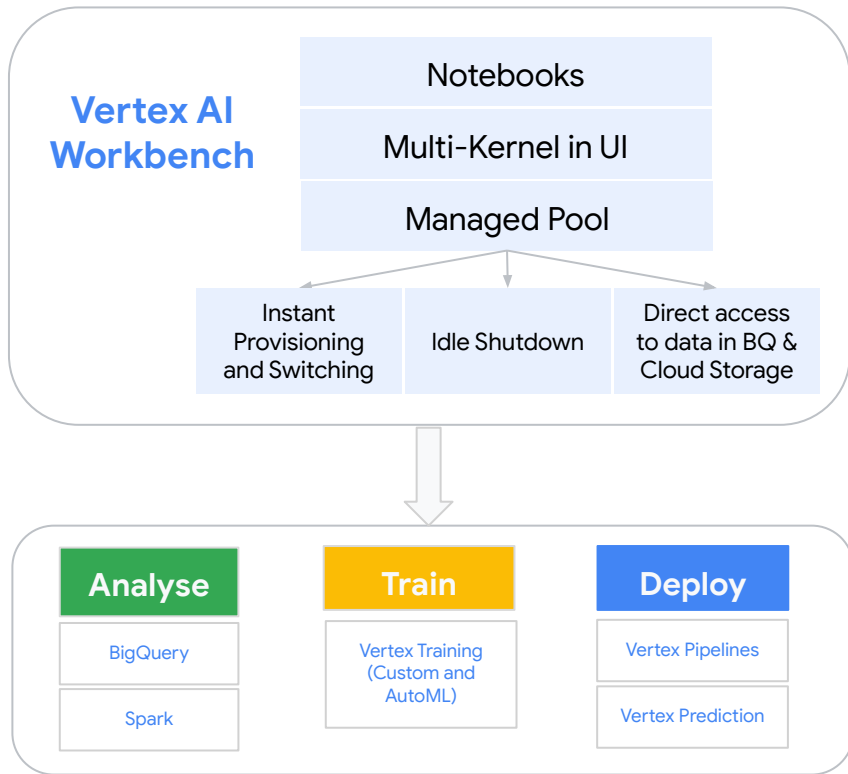
Matching Engine

AI Accelerators

Vertex AI Workbench with Managed Notebooks

Proprietary + Confidential

A one-stop interface for Data Science



Vertex Pipelines

Automate, monitor, and govern your ML systems by orchestrating your ML workflow in a serverless manner, and storing your workflow's artifacts using Vertex ML Metadata



Easy to use Python SDKs: Build your Pipelines using the battle-tested and easy-to-use KFP SDK and TFX SDK



Scalable: Run as many pipelines on as much data as you want without having to worry about compute resources



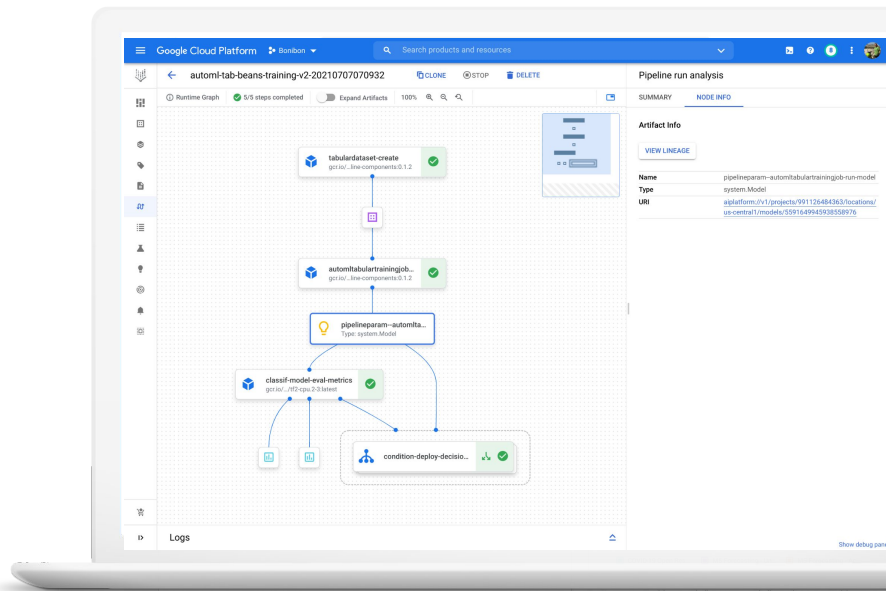
Cost-effective: Pay for the pipelines you run and the resources they use.



Secure: Integrated with GCP security features like IAM, VPC-SC, and CMEK.



Metadata Tracking and Lineage: Automatically store metadata about every artifact produced by the Pipelines.



Metadata and Lineage on Vertex AI

Artifact, lineage, and execution tracking for your ML workflow



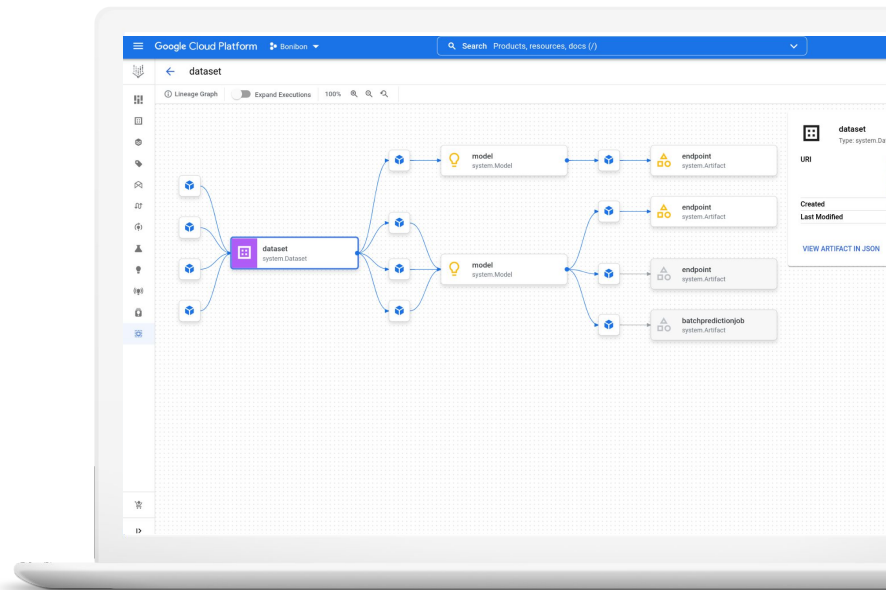
Automatically track inputs and outputs to all components in an ML pipeline, and their lineage.

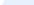



Visualize the workflow for faster debugging with a DAG of all related executions.




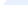
Manage artifacts by projects, group by experiments, and track the usage of datasets and models in your organization



-  **Fully Managed:** Train without provisioning or managing servers. Pay only for the compute you consume. Zero administration.

 **High Performance:** Optimized for Machine Learning.
Scalable distributed orchestration with the most advanced
cloud Accelerators (GPUs and TPUs)

-  **HyperParameter Optimization:** Automatically tune models with Google's Vizier optimizer

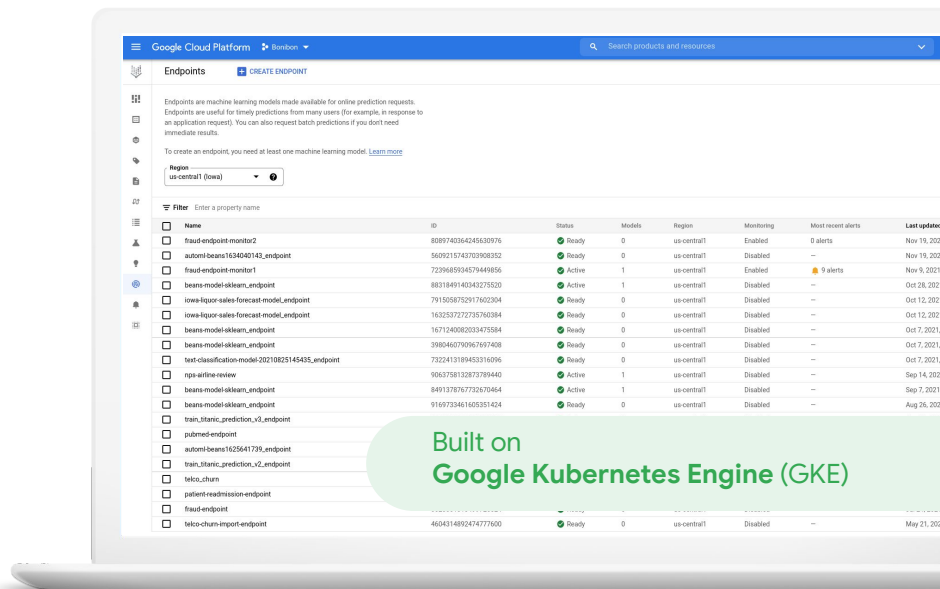
 **Customizable:** Supports predefined (Tensorflow, Sklearn, XGB, Keras) and custom containers with flexible machines

- **Built-in logging and monitoring:** review your execution jobs and monitor resources utilization for your jobs



Vertex Prediction

- Serve **online** endpoints for low-latency predictions, or predictions on massive **batches** of data.
- **Built-in security and compliance:** VPC peering and security perimeter. Custom managed encryption keys. Fine-tuned access control.
- **Low TCO:** Scale automatically based on your traffic, and alleviate operational overhead.
- **Intelligent and assistive:** Built-in Model Explainability and proactive model monitoring.
- Log prediction requests and responses to **BigQuery** for monitoring and debugging
- **Fast inference on GPUs:** Support for a broad range of machine types specialized for ML, such as GPUs.



Feature Store on Vertex AI

A rich feature repository to serve, share and re-use ML features.



Share and reuse ML features across use cases

Centralized feature repository with easy APIs to search & discover features, fetch them for training/serving and manage permissions.



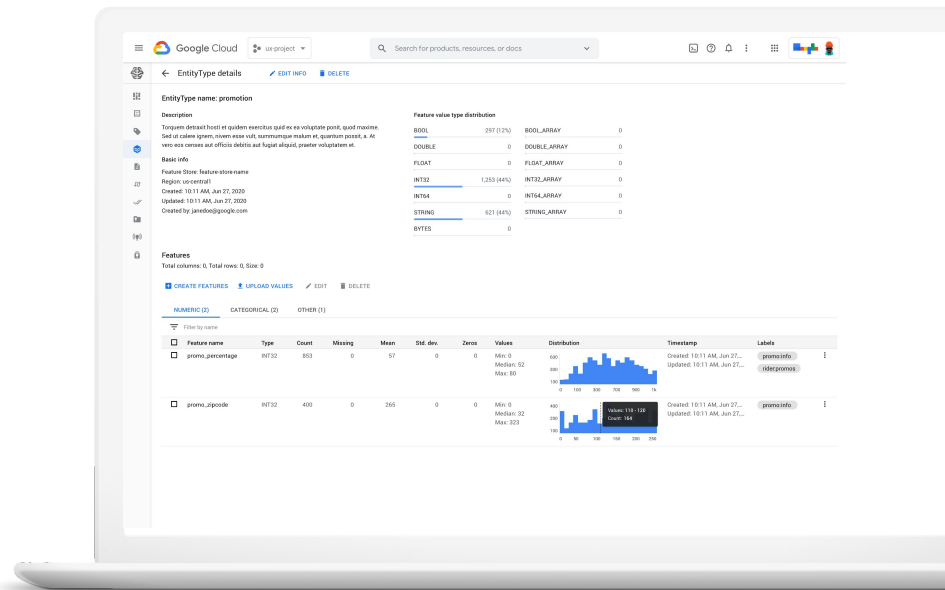
Serve ML Features at scale with low latency

Offload the operational overhead of handling infrastructure for low latency scalable feature serving.



Alleviate training serving skew

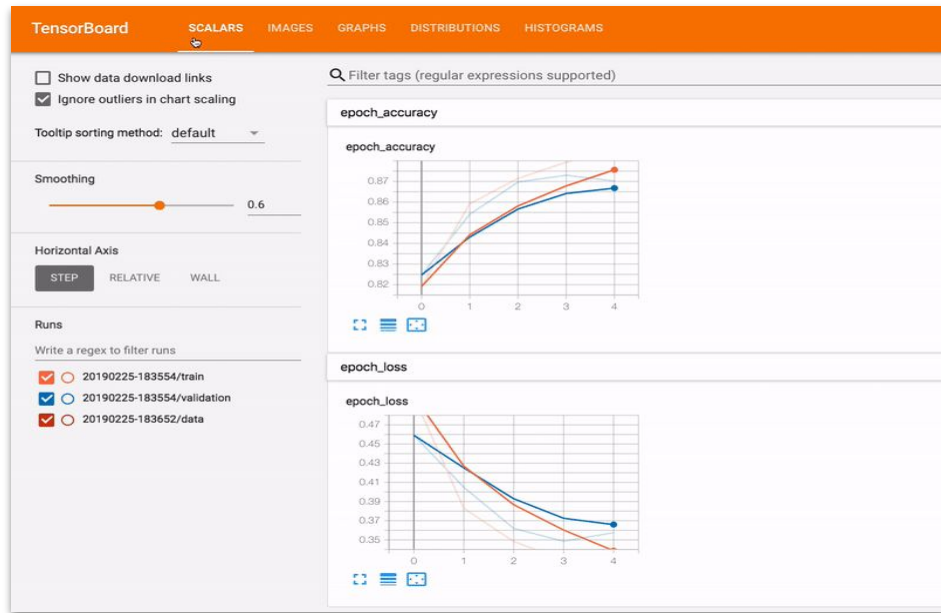
- Compute feature values once, re-use for training and serving
- Track & monitor for drift and other quality issues



TensorBoard: ML visualization toolkit

TensorBoard provides the visualization and tooling needed for ML experimentation

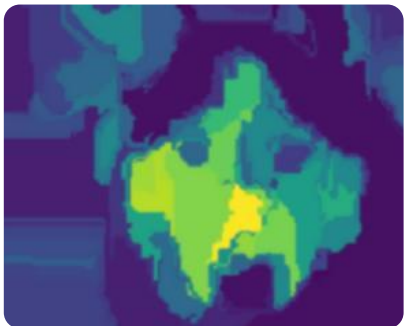
- Tracking and visualizing metrics such as loss and accuracy
- Visualizing the model graph (ops and layers)
- Viewing histograms of weights, biases, or other tensors as they change over time
- Projecting embeddings to a lower dimensional space
- Displaying images, text, and audio data
- Profiling TensorFlow programs



[Get started with TensorBoard Docs](#)

Explainable AI tells you how important each input feature is

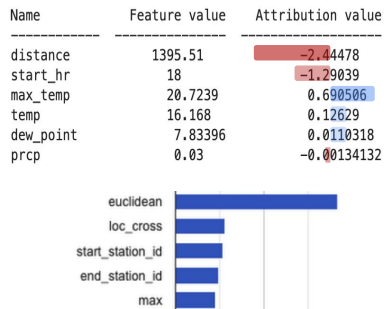
Images



Explanations tell you:

What **image pixels or regions** most contributed to the model's classification?

Tabular



How much did each **feature column** contribute to a single prediction or the model overall?

Text

The cake tastes
delicious!

Sentiment score: 0.9

How much did each **word or token** contribute to the text classification?

Monitoring of model performance

For models deployed in the Vertex Prediction service



Monitor and alert

Monitor signals for model's predictive performance, and alert when those signals deviate.



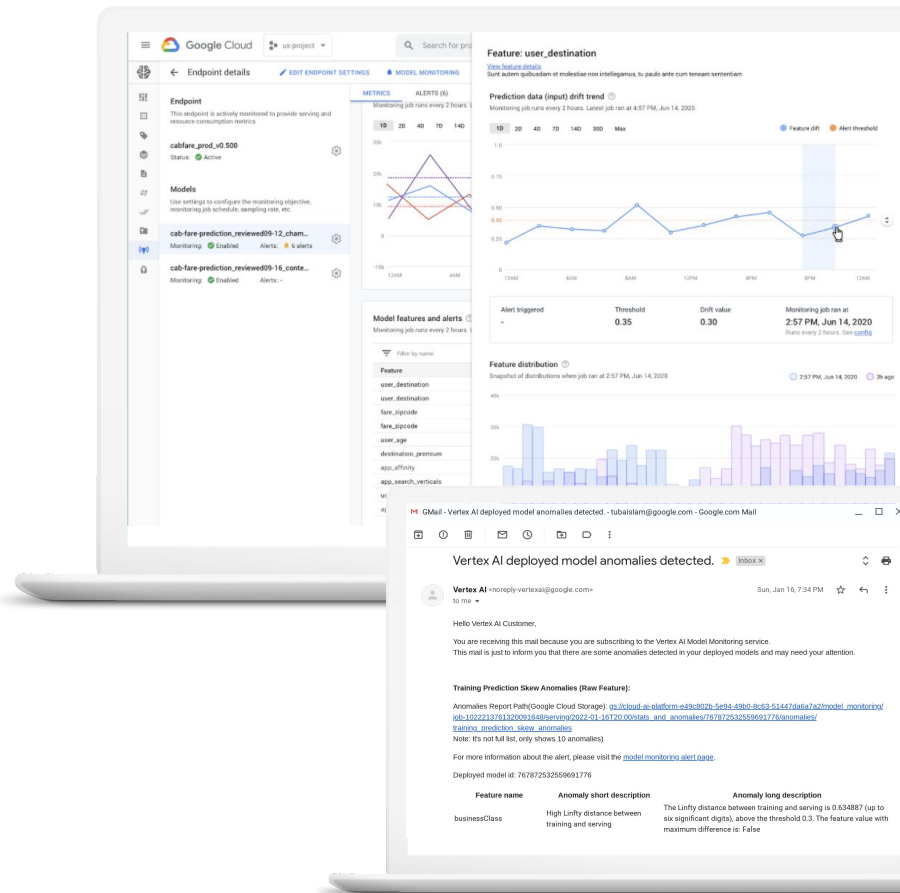
Diagnose

Help identify the cause for the deviation i.e. what changed, how and how much?



Update Model

Trigger model re-training pipeline or collect relevant training data to address performance degradation.



New ML Tools on Vertex AI: Matching Engine

30-50% cheaper than alternatives, while delivering higher scale and lower latencies.



Faster i.e. low latency

Find nearest neighbors in a few millisecond



The most scalable

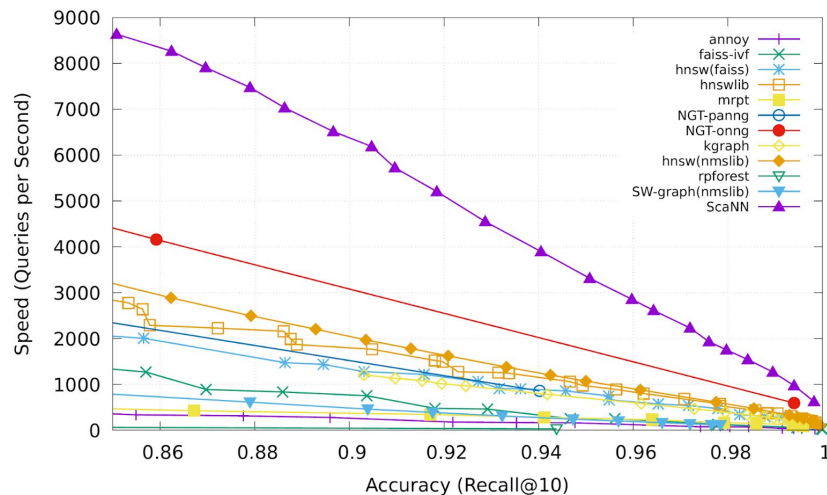
Scales to billions of vectors



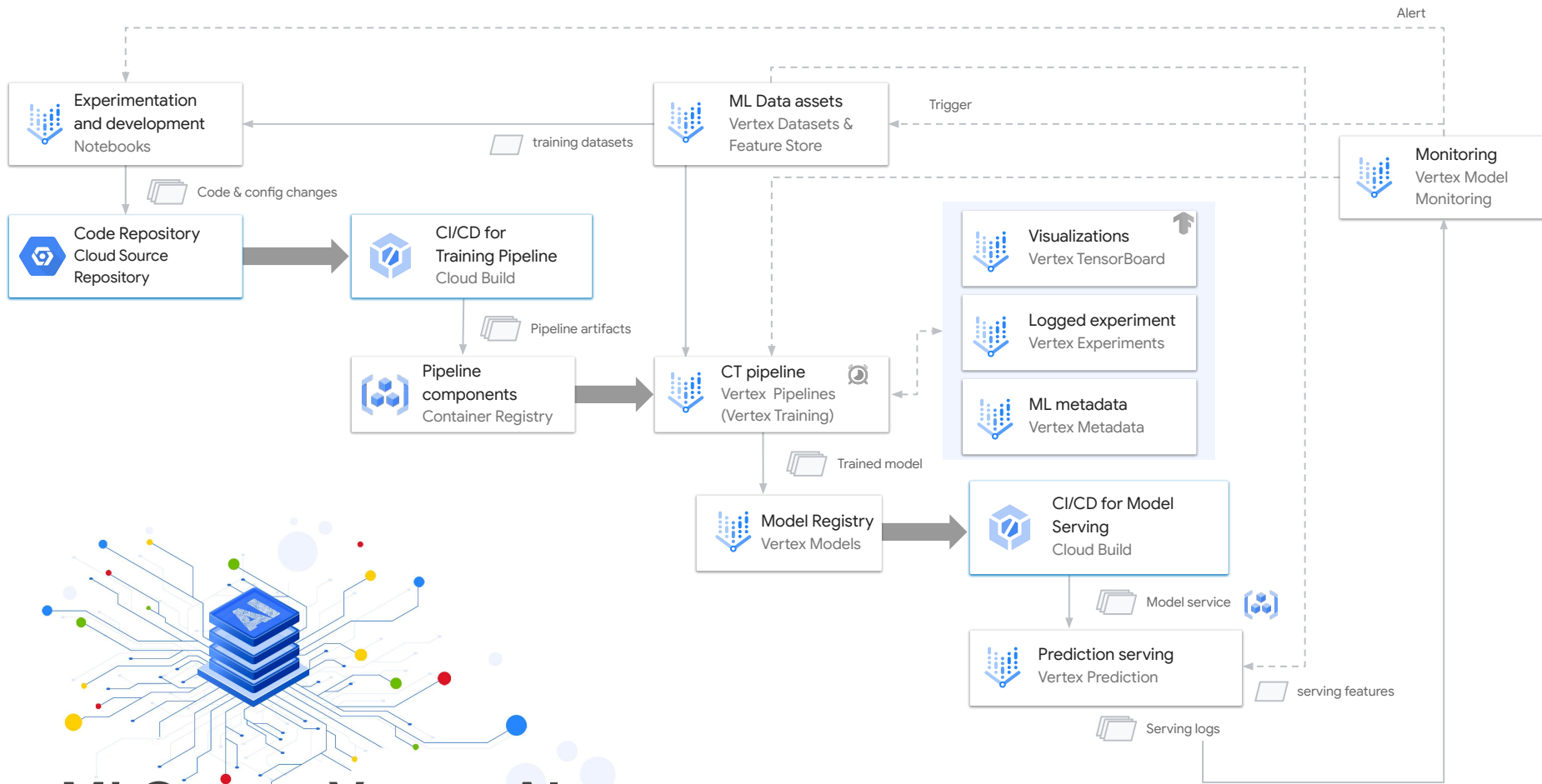
Cheaper

Requires fewer VMs to serve the same workload

- 1/4th the CPU consumption of **faiss**
- 1/3rd the memory consumption of **nmslib**



Google's technology (labelled **ScaNN**) compared with popular ANN services



MLOps on Vertex AI

Data & AI

Demo's Marijn



Advanced analytics



Anomaly detection



Forecasting



Document processing



Image recognition



Conversational AI

