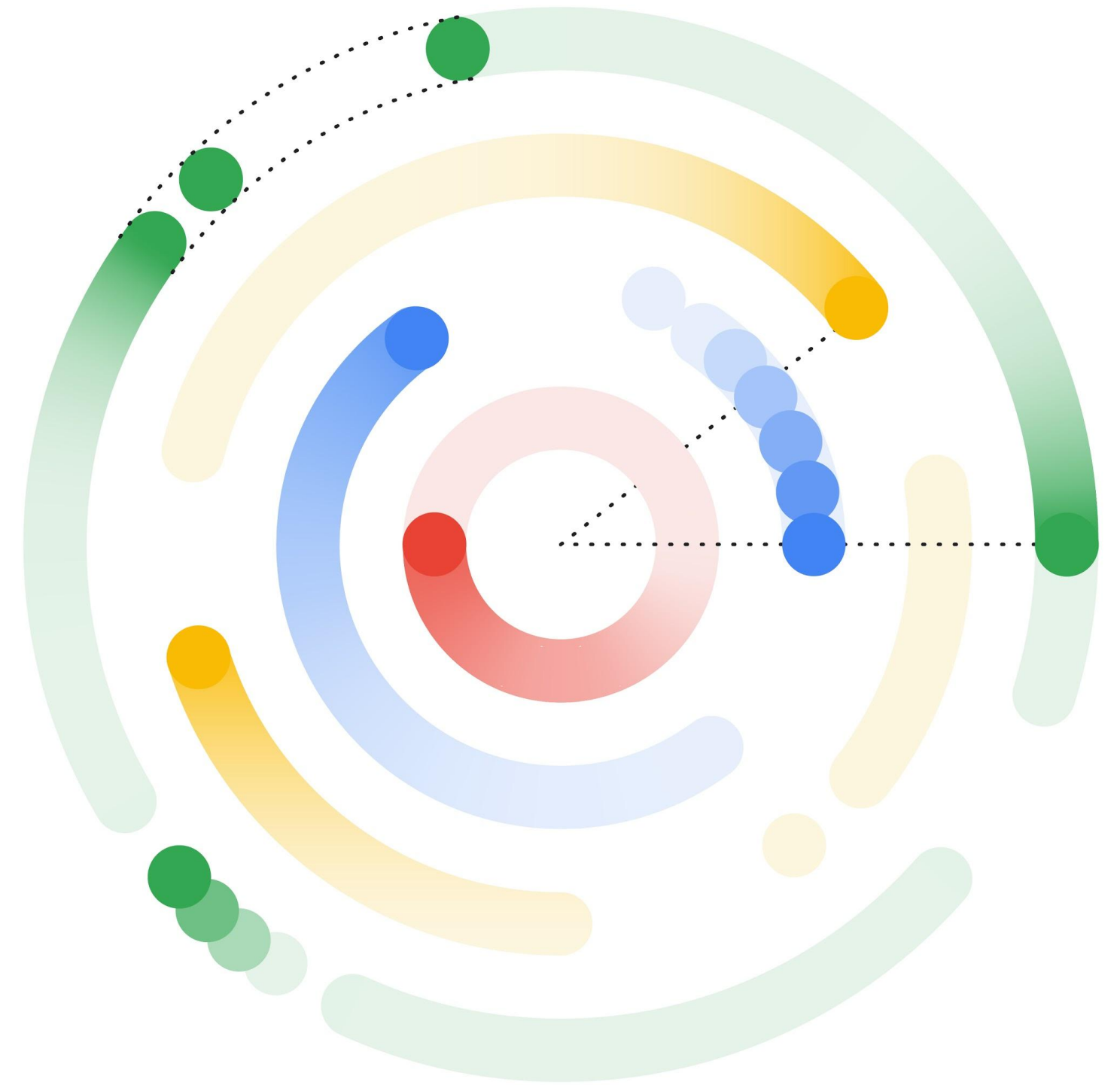Google Cloud

# Scalable ML Deployment Using PyTorch and Kubeflow Pipelines

**Google Cloud Applied ML Summit**
Solving for the future.

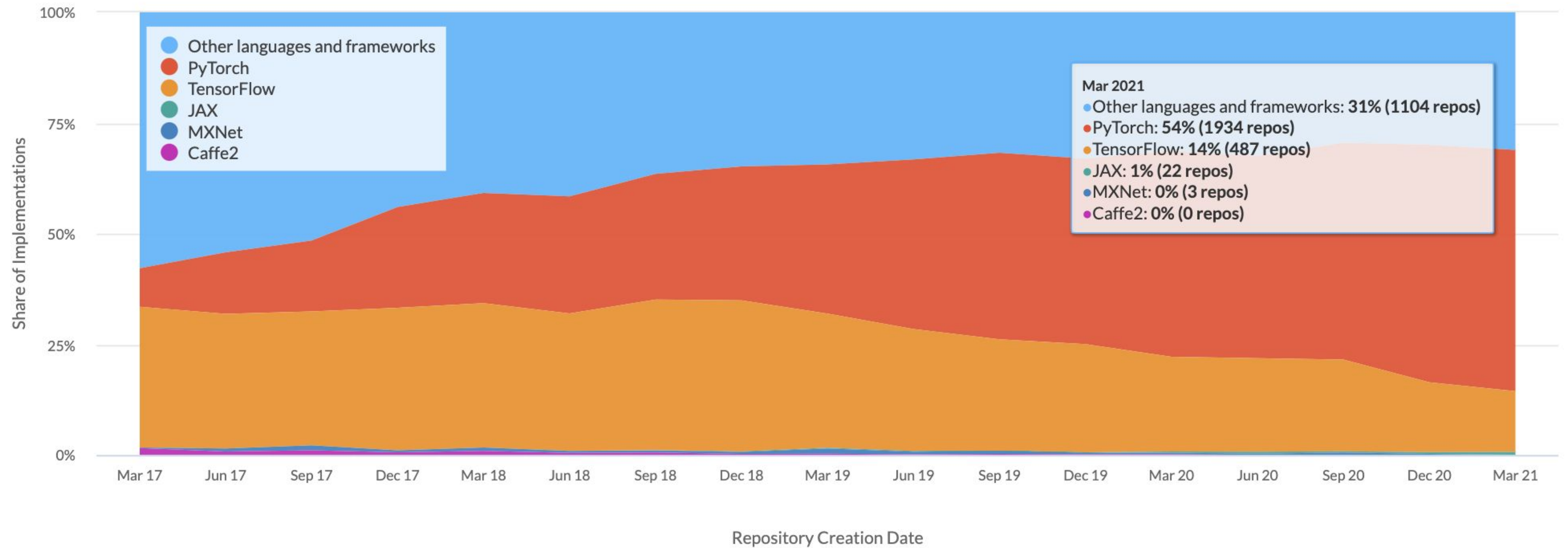06/10/21

**Geeta Chauhan**

Head of AI Partner
Engineering
Facebook AI

Google Cloud

# Table of Contents

Google Cloud

# PyTorch Community Growth

Paper Implementations grouped by framework



Legend:
- Other languages and frameworks
- PyTorch
- TensorFlow
- JAX
- MXNet
- Caffe2

**Mar 2021**
- Other languages and frameworks: **31% (1104 repos)**
- PyTorch: **54% (1934 repos)**
- TensorFlow: **14% (487 repos)**
- JAX: **1% (22 repos)**
- MXNet: **0% (3 repos)**
- Caffe2: **0% (0 repos)**

Y-axis: Share of Implementations (0%, 25%, 50%, 75%, 100%)

X-axis: Repository Creation Date (Mar 17, Jun 17, Sep 17, Dec 17, Mar 18, Jun 18, Sep 18, Dec 18, Mar 19, Jun 19, Sep 19, Dec 19, Mar 20, Jun 20, Sep 20, Dec 20, Mar 21)
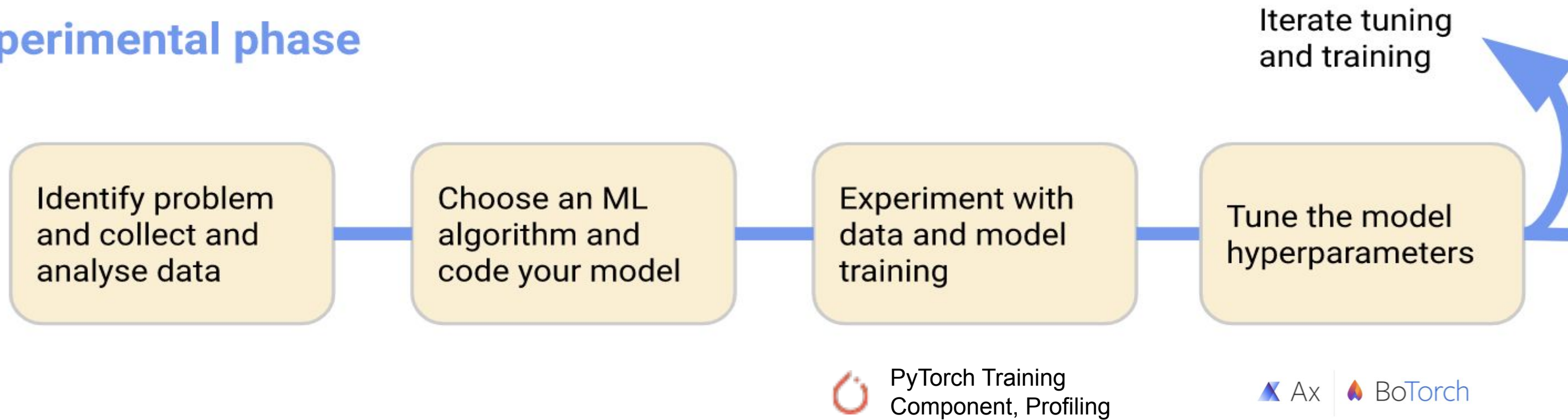
Source: https://paperswithcode.com/trends

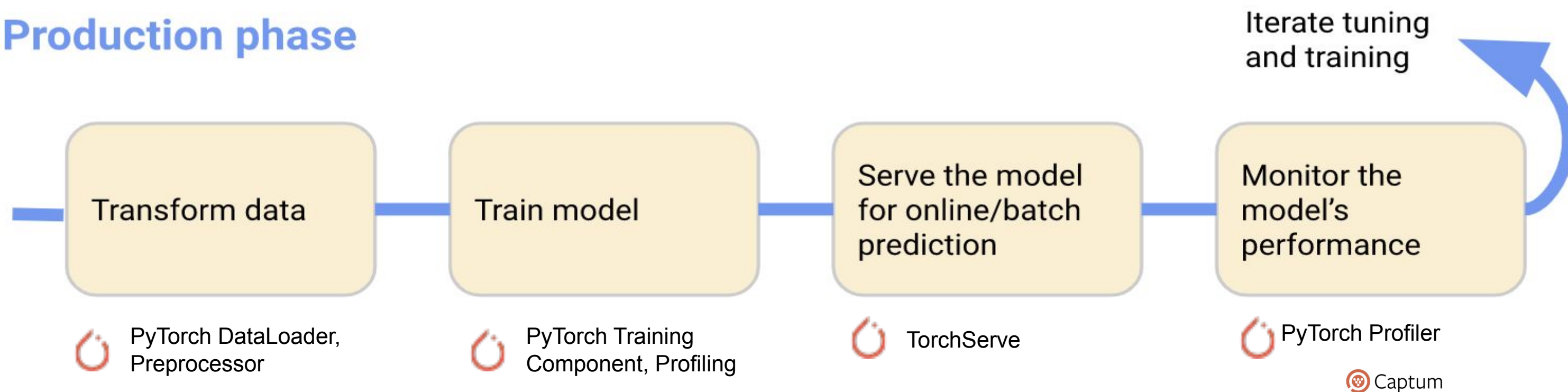Google Cloud

# Challenges for Scalable ML

- Building and Managing pipelines is hard

- Continuous iterative process, optimize for metric

- Over time data changes, model drift

- Experiment tracking is difficult

- Model artifacts get lost

- Diverse deployment environments

- What happens when you hit scale?
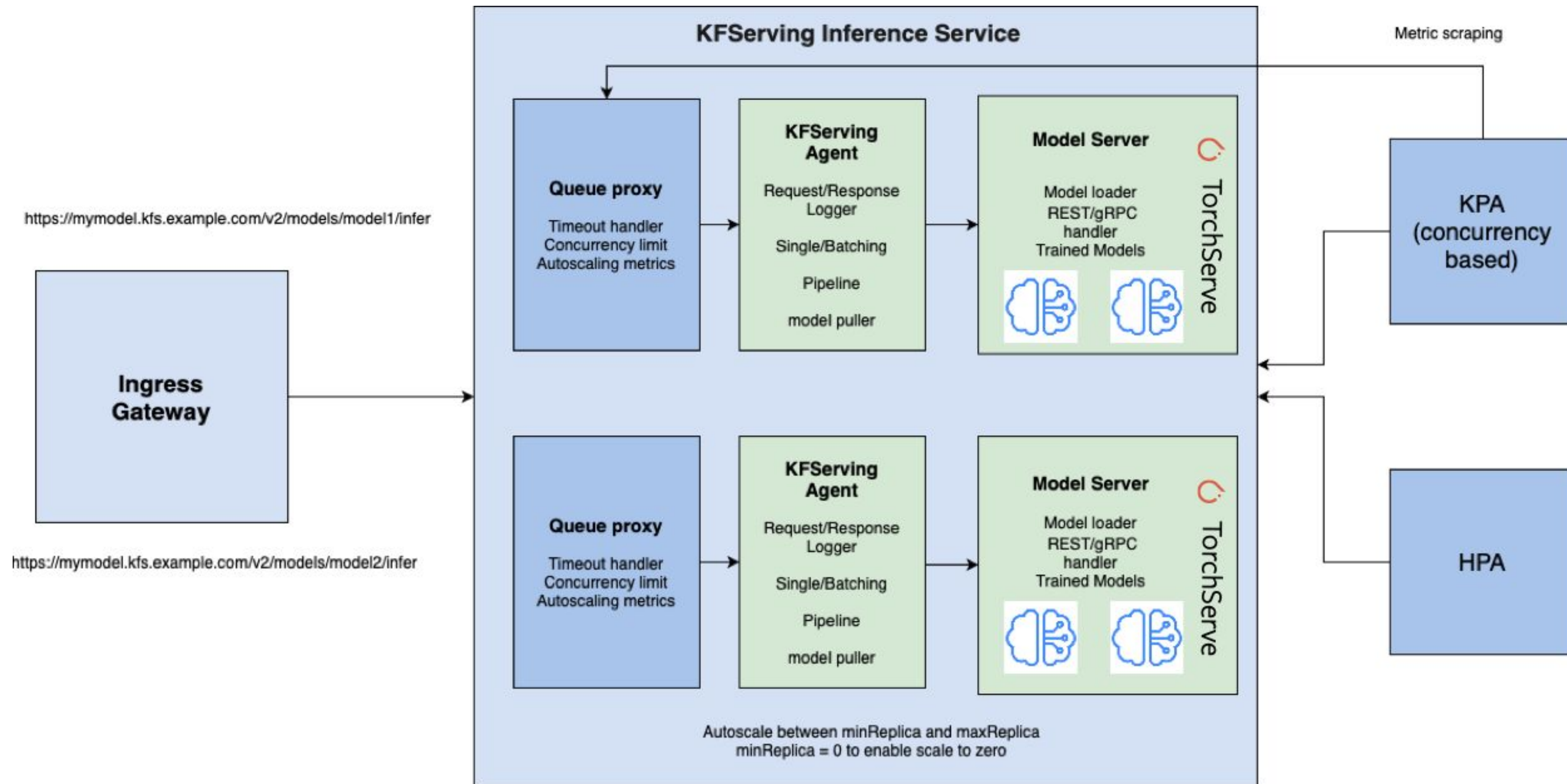
# Kubeflow Pipelines + PyTorch

**Experimental phase**

Iterate tuning and training

| Identify problem and collect and analyse data | Choose an ML algorithm and code your model | Experiment with data and model training | Tune the model hyperparameters |
|---|---|---|---|

PyTorch Training Component, Profiling

Ax | BoTorch

**Production phase**

Iterate tuning and training

| Transform data | Train model | Serve the model for online/batch prediction | Monitor the model's performance |
|---|---|---|---|

PyTorch DataLoader, Preprocessor

PyTorch Training Component, Profiling

TorchServe

PyTorch Profiler

Captum

Google Cloud

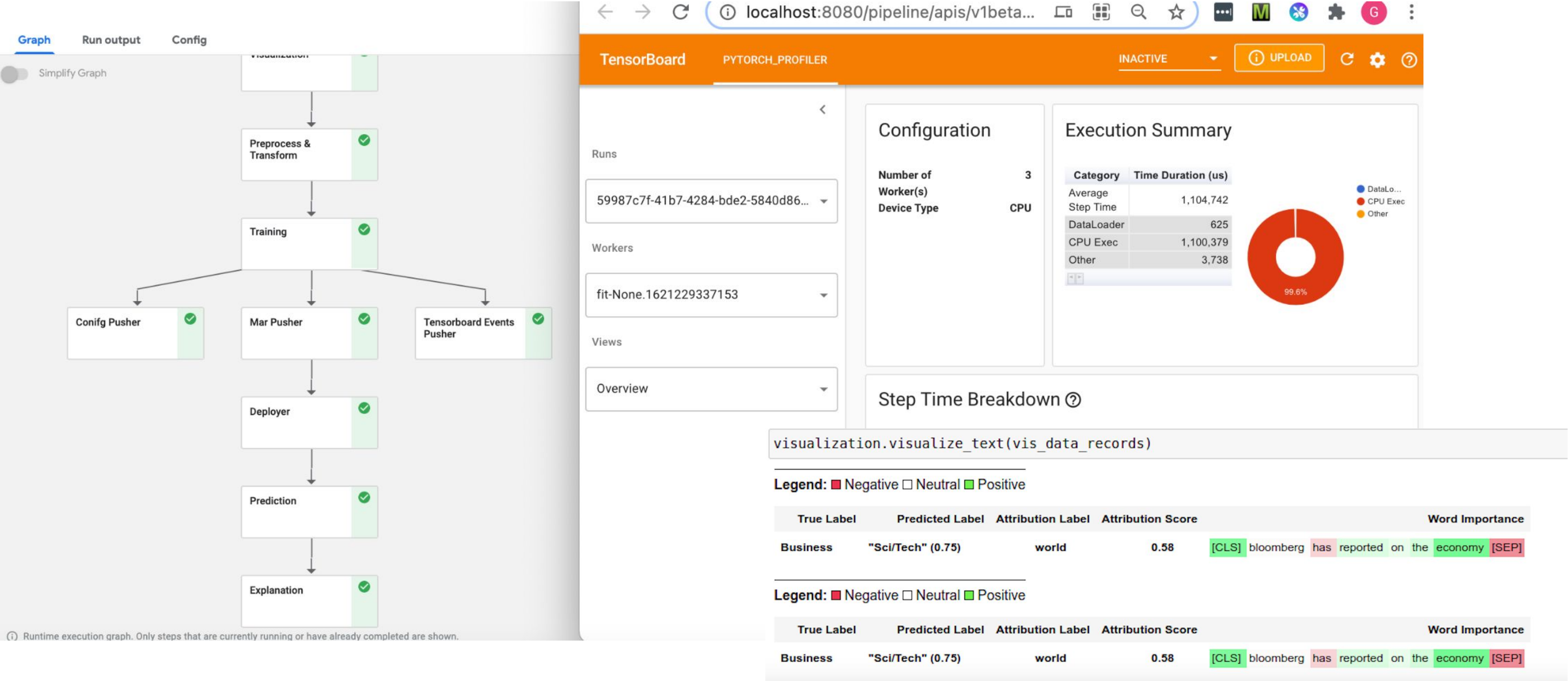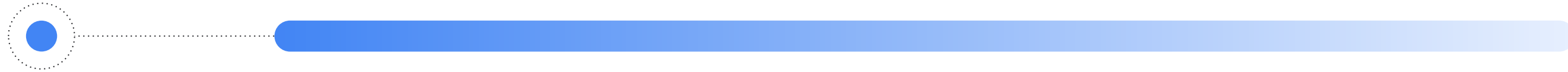# KFServing for Model Serving with TorchServe

# Kubeflow Pipelines + PyTorch

- Dataloading and preprocessing

- Model Training

- Hyper Parameter Optimization using Ax/Botorch

- Model deployment & Serving using TorchServe + KFServing w/ canary rollouts, autoscaling, prometheus monitoring

- Visualizations using Tensorboard PyTorch Profiler

- Model Interpretability using Captum

- Artifact Lineage Tracking

- Open Source KFP: on-prem or any cloud

- GCP Vertex AI Pipelines - Serverless

# NLP BERT workflow - Open source KFP
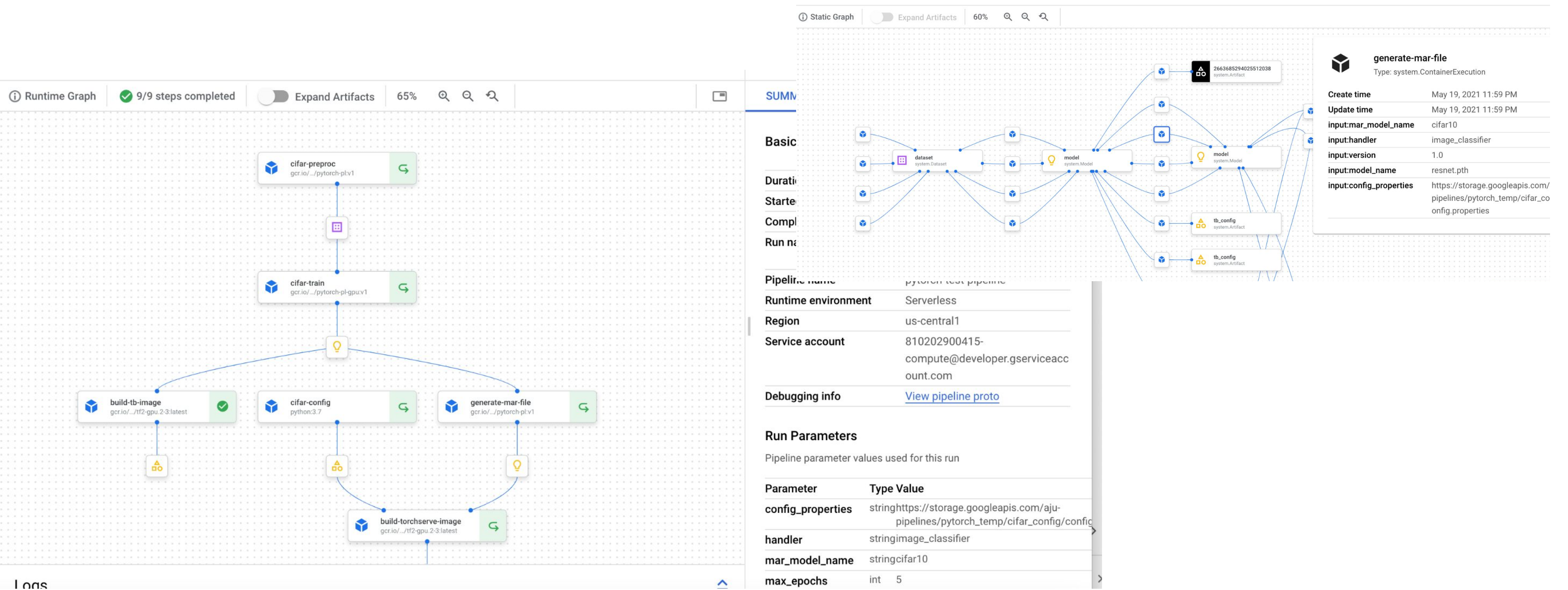


Google Cloud

# Demo

# Scalable Serverless Vertex AI Pipelines

- Automated, Scalable, Serverless, no need to manage cluster

- Cost Effective - Pay only for what you use

- Build using familiar open source KFP SDK

- Metadata and Lineage for all Artifacts, metrics across and execution for the workflow

- Pipeline run analysis, monitoring and debugging

- Security using Cloud IAM, VPC-SC and CMEK
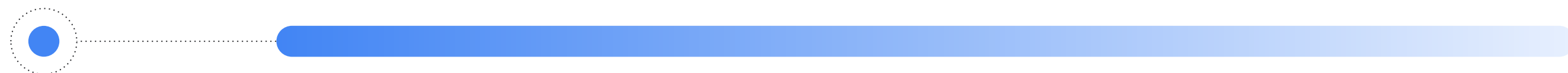
# Computer Vision Workflow - Vertex AI Pipeline



Google Cloud

# Demo

# References

- KFP Examples Github: http://bit.ly/pt-kfp-demos

- Vertex AI Pipeline - Colab Notebook: http://bit.ly/PT-Vertex-AI

- Vertex AI Pipelines: https://cloud.google.com/vertex-ai/docs/pipelines

- Captum.ai: https://captum.ai/tutorials/CIFAR_TorchVision_Captum_Insights

- PyTorch Tensorboard Profiler: https://pytorch.org/blog/introducing-pytorch-profiler-the-new-and-improved-performance-tool/

- Operationalize, Scale & Infuse Trust in AI using KFServing: https://blog.kubeflow.org/release/official/2021/03/08/kfserving-0.5.html

Google Cloud

Google Cloud Summit

# Q&A

Contact:

Email: gchauhan@fb.com

Linkedin:
https://www.linkedin.com/in/geetachauhan/

Google Cloud

Google Cloud Summit
# Thank you for joining