

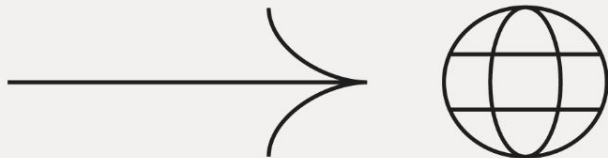


# Observability for GenAI Apps: Logging, Metrics, & Tracing on GCP

Ananda Dwi Rahmawati  
Google Developer Expert - Cloud

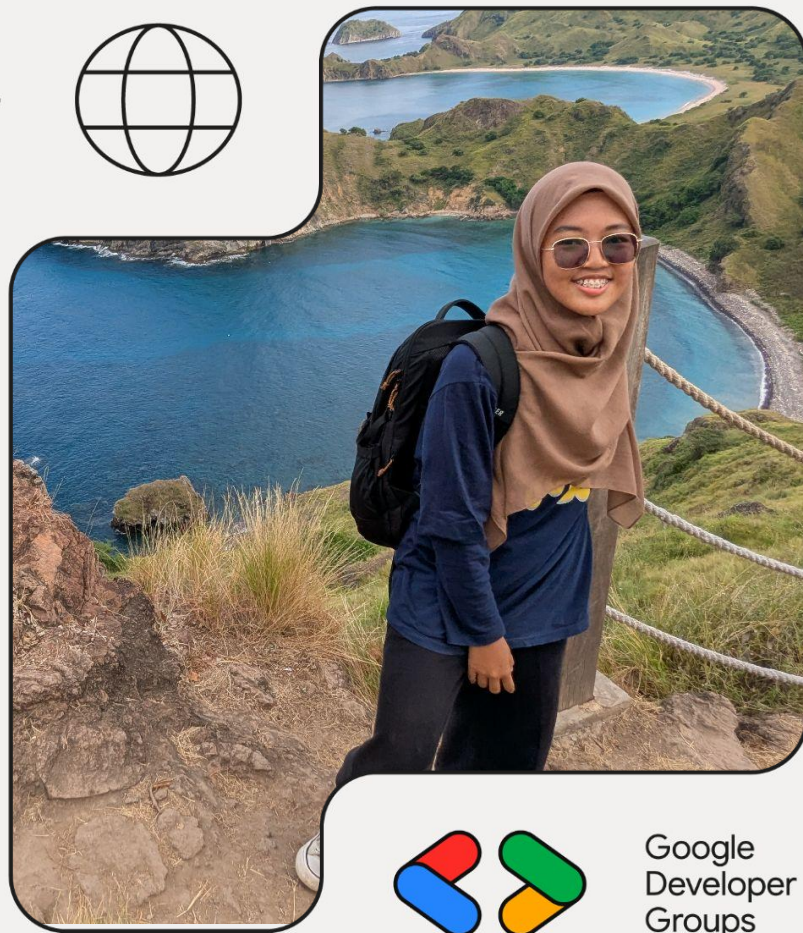


Google  
Developer  
Groups



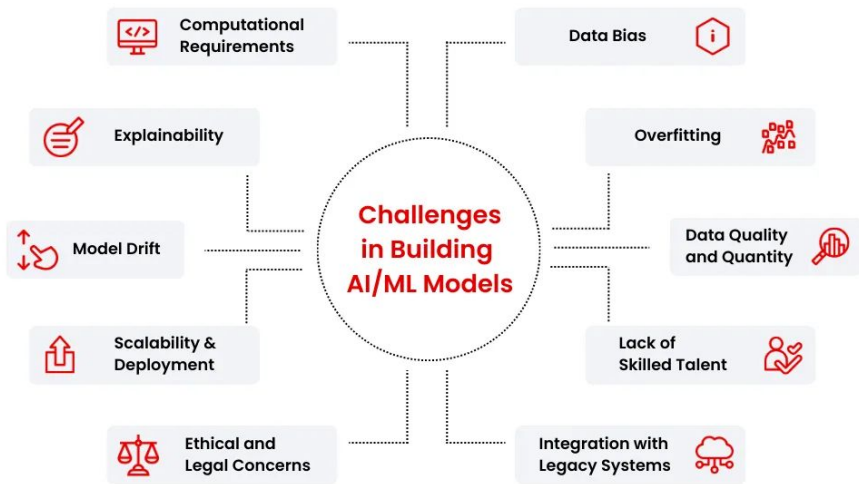
## Ananda Dwi Rahmawati

- ❑ Cloud & DevOps Engineer, Singapore
- ❑ Google Developer Expert Cloud - Modern Architecture
- ❑ Master of Computer Science - University of Texas at Austin
- ❑ <https://linktr.ee/misskecupbung>



## Operating ML models, presents several challenges:

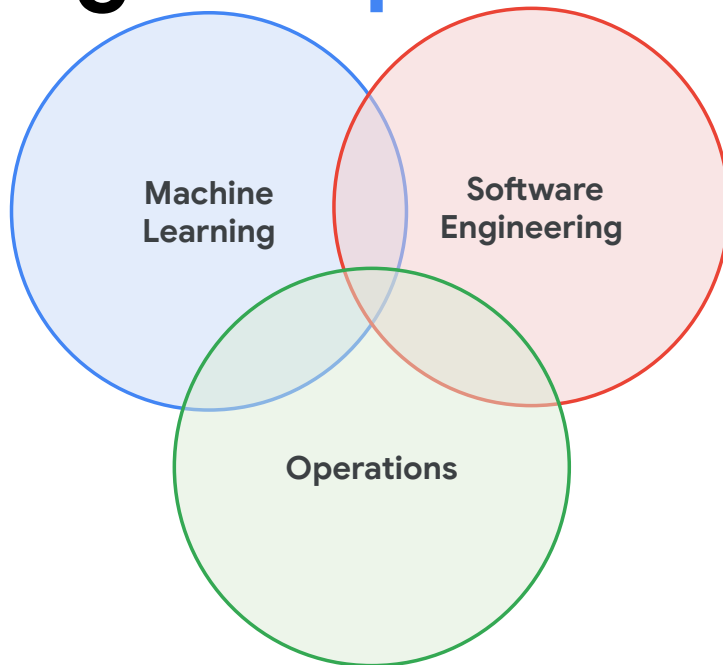
- **Model drift:** As real-world data changes, models become less accurate, requiring frequent retraining.
- **Resource management:** ML workloads have varying demands, making efficient allocation crucial.
- **Data quality:** Consistent, reliable input data is essential for model performance.
- **Compliance:** Meeting governance and regulatory requirements is challenging.
- **Versioning:** Tracking models, datasets, and experiments is difficult at scale.



Google  
Developer  
Groups

# Introducing MLOps

- Model development
- Model evaluation
- Parameter tuning



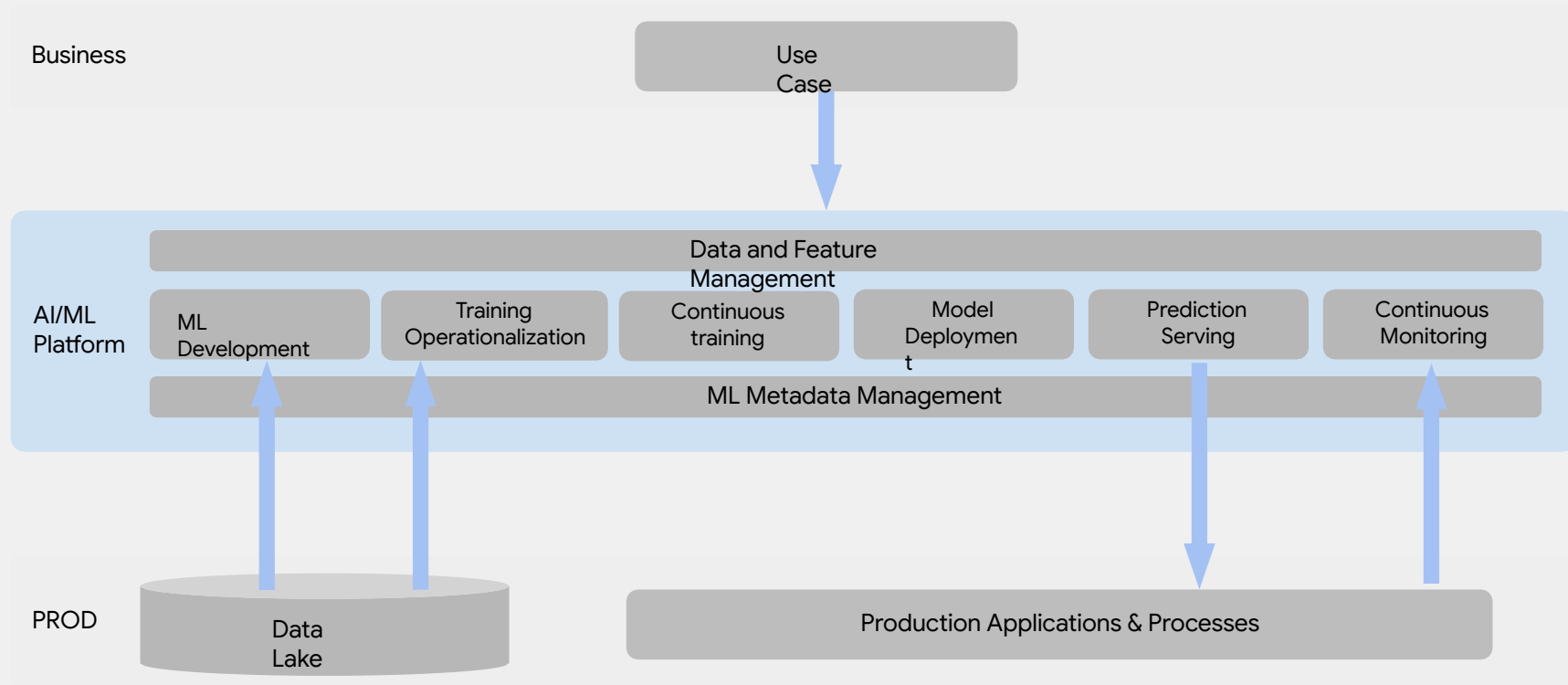
- Data engineering
- Pipeline development
- Integration of model into business application

- Model deployment
- Metadata management
- Logging and monitoring

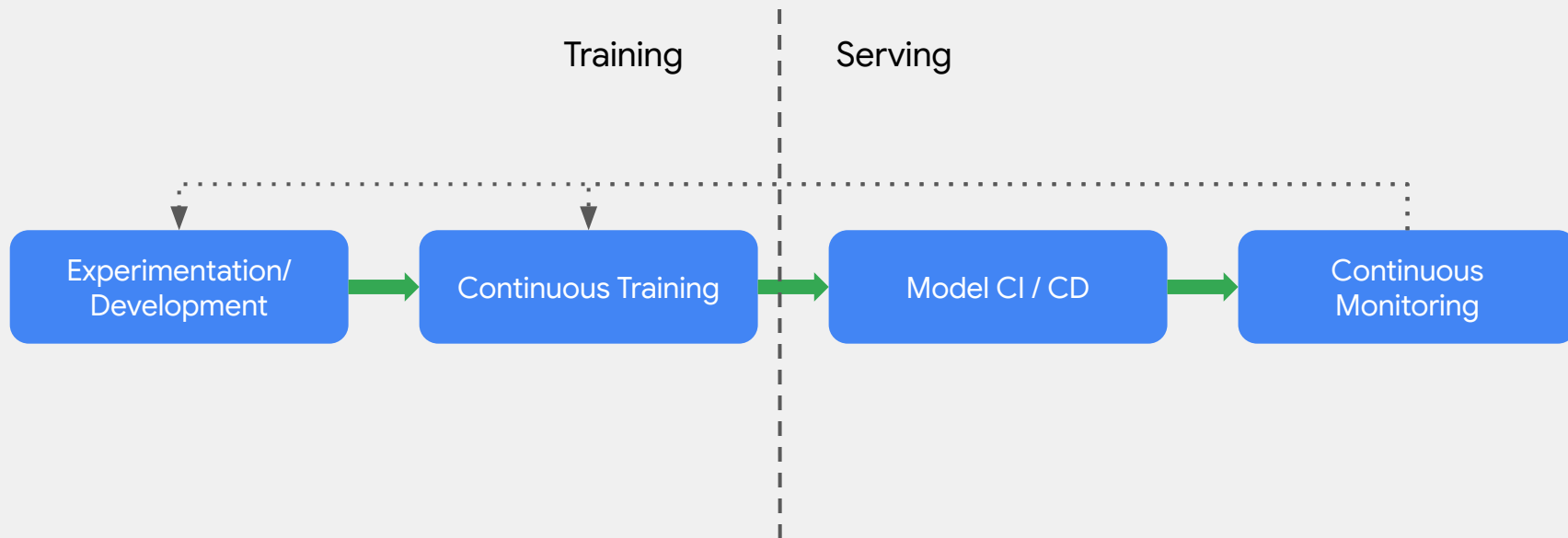


Google  
Developer  
Groups

# MLOps: quick recap

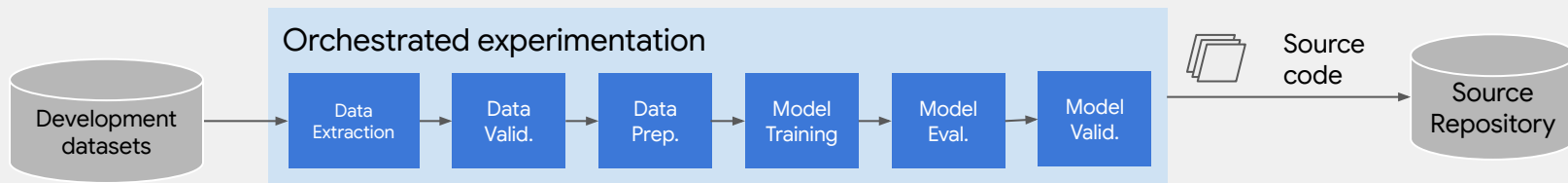


# ML Solution Lifecycle



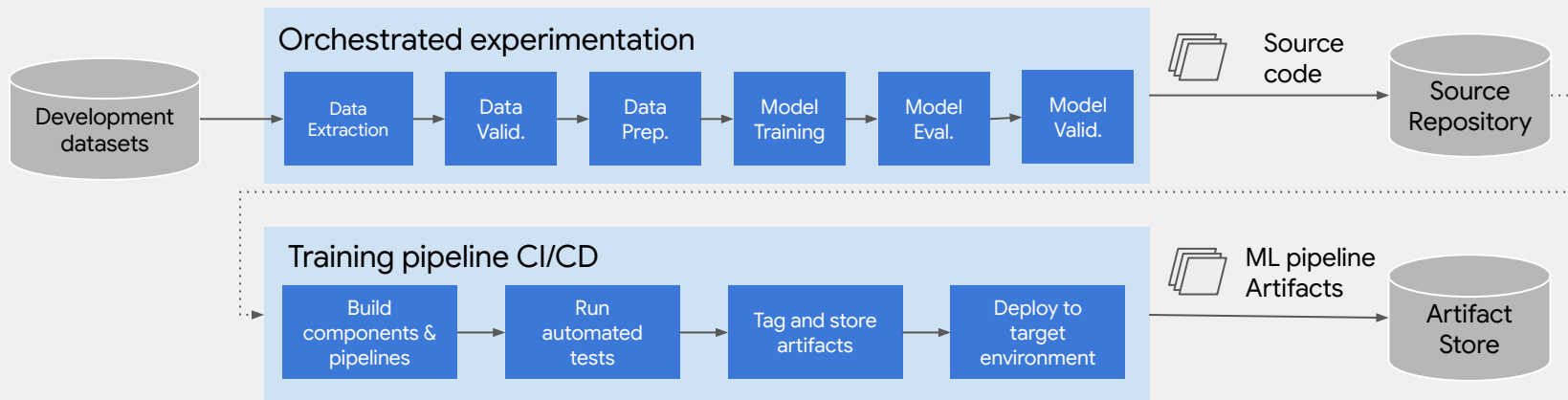
# Reliable and repeatable training

## Automated E2E Pipelines



# Reliable and repeatable training

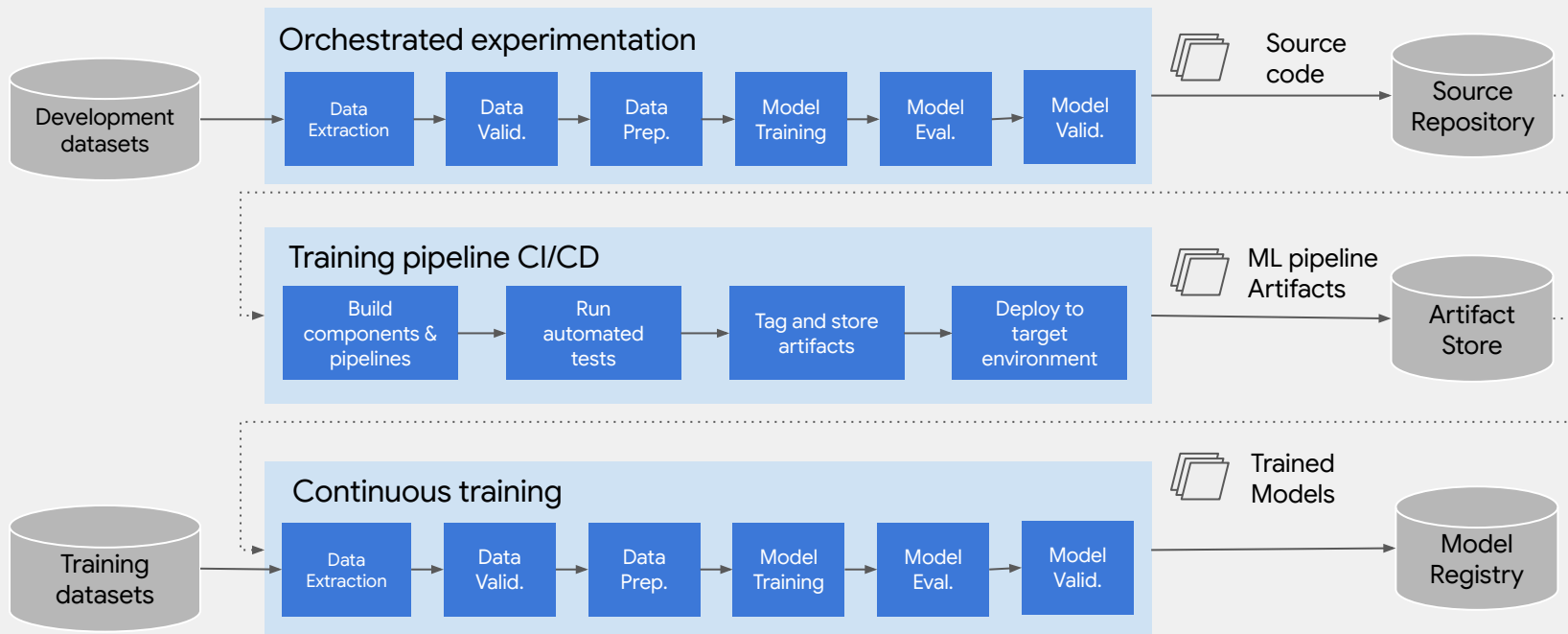
## Automated E2E Pipelines





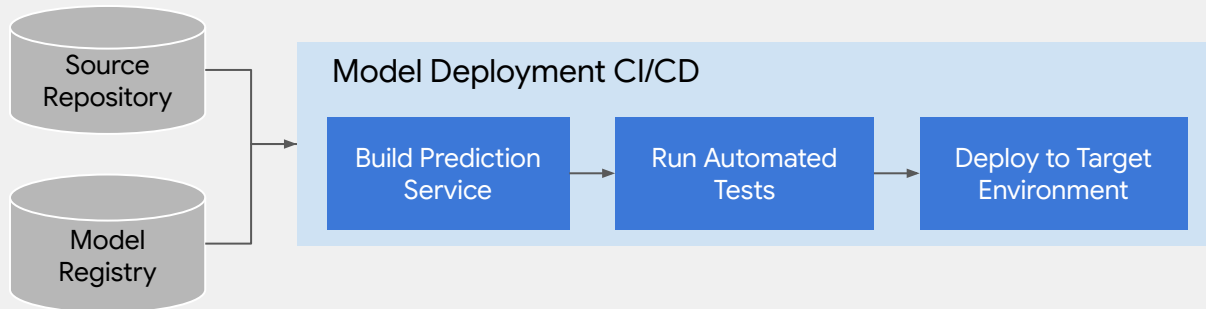
# Reliable and repeatable training

## Automated E2E Pipelines



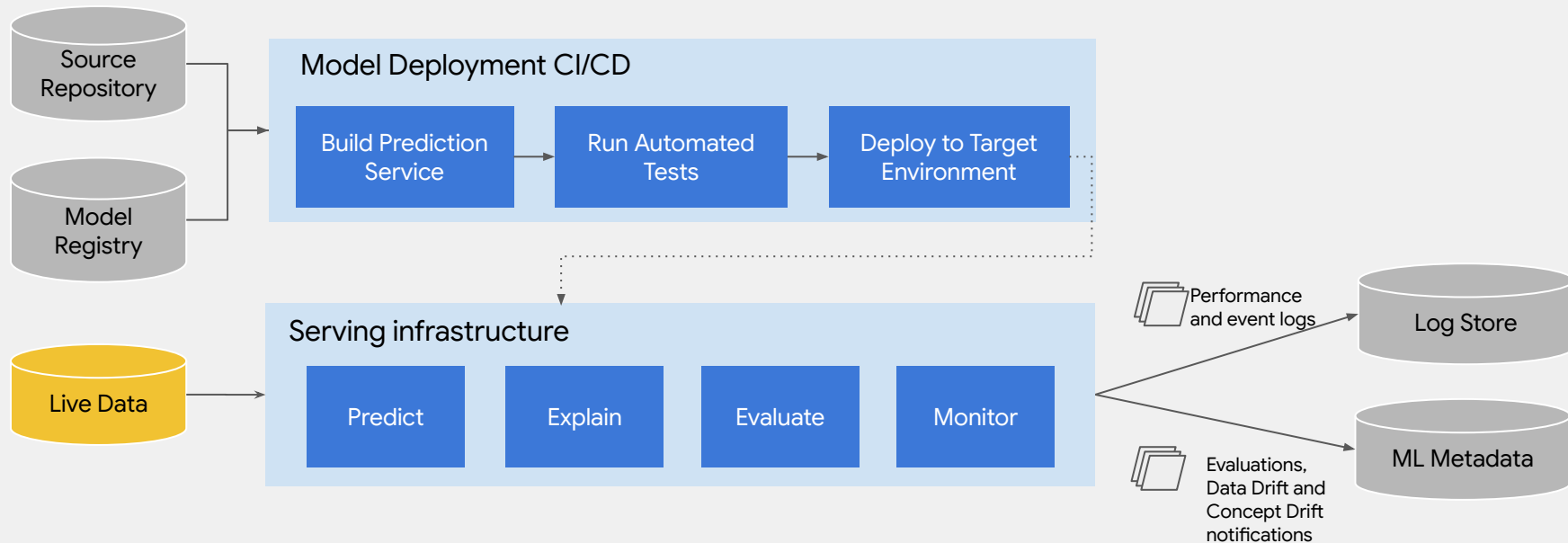
# Reliable and monitored serving

## Automated E2E Pipelines

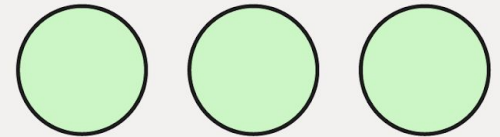
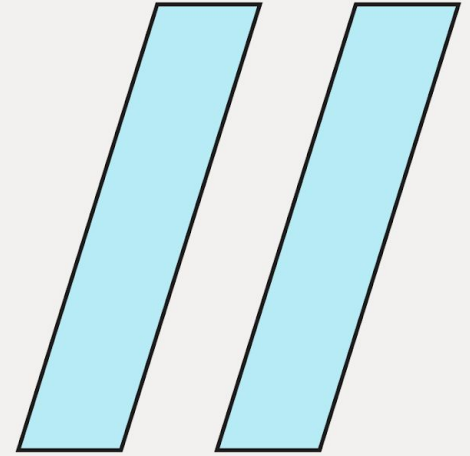


# Reliable and monitored serving

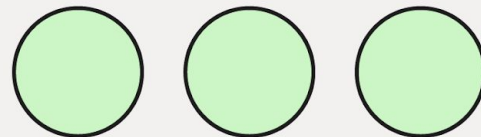
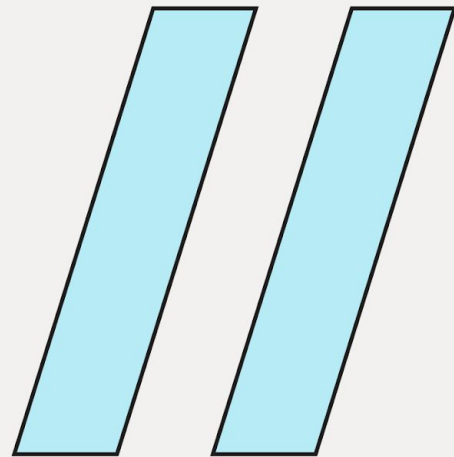
## Automated E2E Pipelines



“In control theory, **observability** is a measure of **how well** internal states of a system can be inferred from knowledge of its external outputs.”



“**Monitoring** tells you whether  
a system is working;  
**Observability** lets you  
understand why isn't  
working.”



# Unique ML Characteristics

## Resource Patterns:

- Sustained high GPU usage during training vs consistent CPU usage in traditional apps
- Specialized GPU node scheduling vs typical short-lived batch jobs
- Variable computational demands requiring dynamic resource allocation

## Monitoring Focus:

- **Model-specific metrics:** Accuracy, F1 scores (irrelevant for standard applications)
- **Data drift monitoring:** Track shifts in user preferences and data patterns
- **Continuous feedback loops:** Analyze interactions for targeted improvements
- **Granular observations:** Sometimes per-prediction monitoring vs standard application metrics



Google  
Developer  
Groups

# Observability

Observability = **gaining insights into ML model behavior & infrastructure.**

## Enables Teams to:

- Quickly identify and diagnose issues
- Optimize resource usage
- Ensure compliance
- Monitor model performance and detect drift
- Track data quality and integrity

## Feedback Loop:

- Continuous monitoring and retraining using real-world data
- Helps models adapt to user behavior, new data patterns, and emerging trends
- Drives better decision-making, user experience, and business value

Your Pods Are  
**Probably Fine...**  
*\*Probably.\**



Google  
Developer  
Groups

# Unique ML Characteristics

## Resource Patterns:

- Sustained high GPU usage during training vs consistent CPU usage in traditional apps
- Specialized GPU node scheduling vs typical short-lived batch jobs
- Variable computational demands requiring dynamic resource allocation

## Monitoring Focus:

- **Model-specific metrics:** Accuracy, F1 scores (irrelevant for standard applications)
- **Data drift monitoring:** Track shifts in user preferences and data patterns
- **Continuous feedback loops:** Analyze interactions for targeted improvements
- **Granular observations:** Sometimes per-prediction monitoring vs standard application metrics



Google  
Developer  
Groups

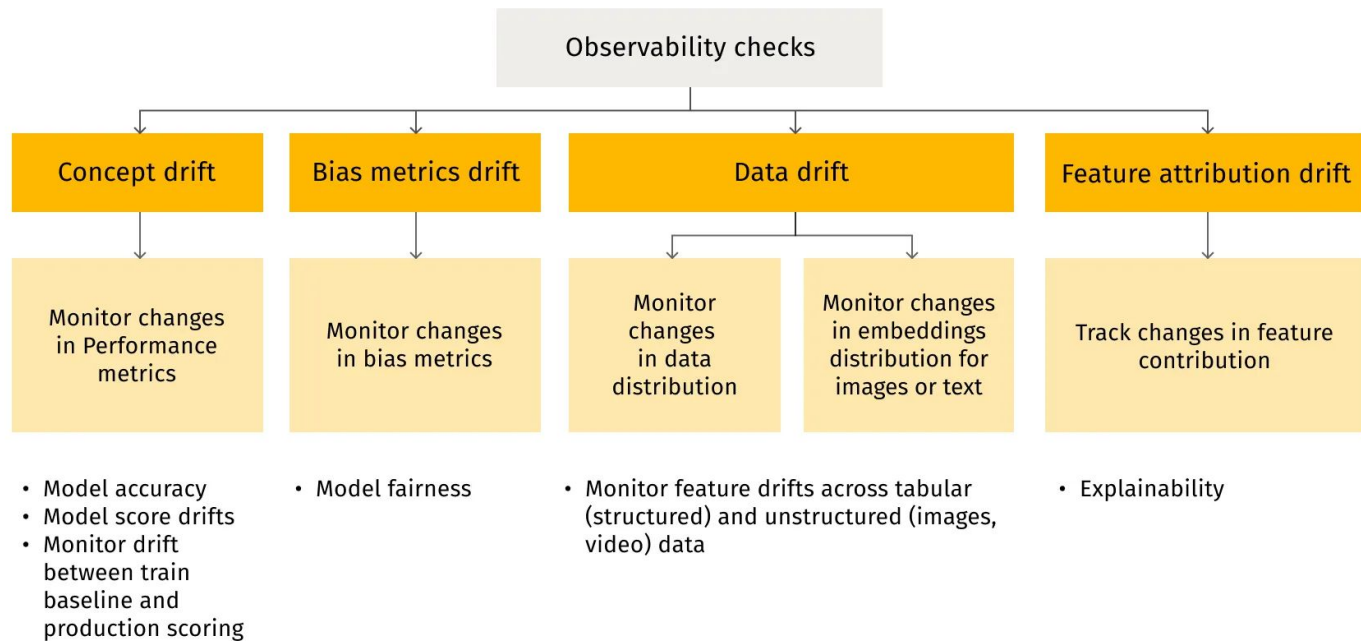


# The Need for ML Observability in MLOps

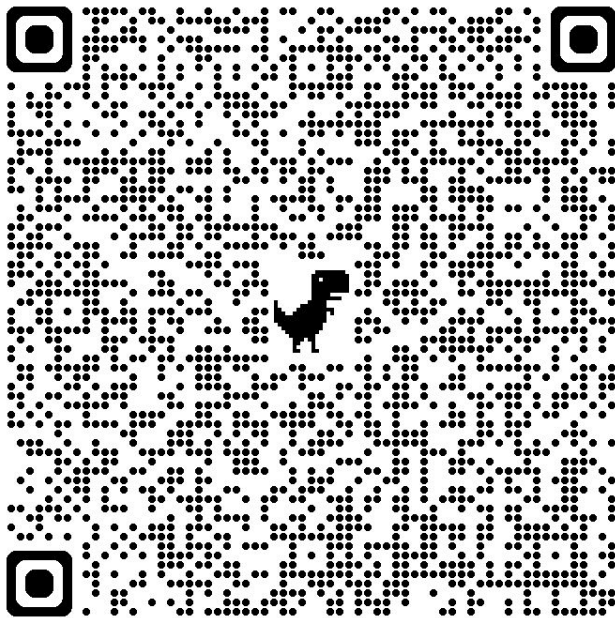
Pillar	What It Covers	GCP Services / Methods
<b>Data Quality</b>	Detect schema mismatches, cardinality shifts, out-of-range values; track distribution drift in features.	<b>Vertex AI Feature Monitoring</b> for feature drift & anomalies; batch & streaming ingestion via <b>BigQuery</b> , <b>Dataflow</b> , or <b>Cloud Storage</b> ; establish baseline datasets with <b>Vertex AI Data Labeling</b> and <b>Data Quality checks</b> .
<b>Fairness / Bias</b>	Pre- and post-training bias detection; monitor predictions' distribution across sensitive groups.	<b>Vertex AI Fairness Indicators</b> to compute fairness metrics across facets; integrate with model evaluation pipelines in <b>Vertex AI Experiments</b> .
<b>Explainability</b>	Understand which features drive predictions (global & local); detect unjustified dependencies.	<b>Vertex AI Explainable AI</b> using SHAP/LIME-like methods; feature attribution for global & local explanations; visualize top features, heatmaps, and partial dependence plots.
<b>Model Performance / Drift</b>	Monitor model accuracy, recall, F1, etc.; detect concept drift; compare predictions vs ground truth.	<b>Vertex AI Model Monitoring</b> for drift detection and performance metrics; optionally use <b>TensorFlow Data Validation (TFDV)</b> or open-source libs like <b>nannyML</b> .



Google  
Developer  
Groups



# Let's Demo



Google  
Developer  
Groups



**Thank you!**

