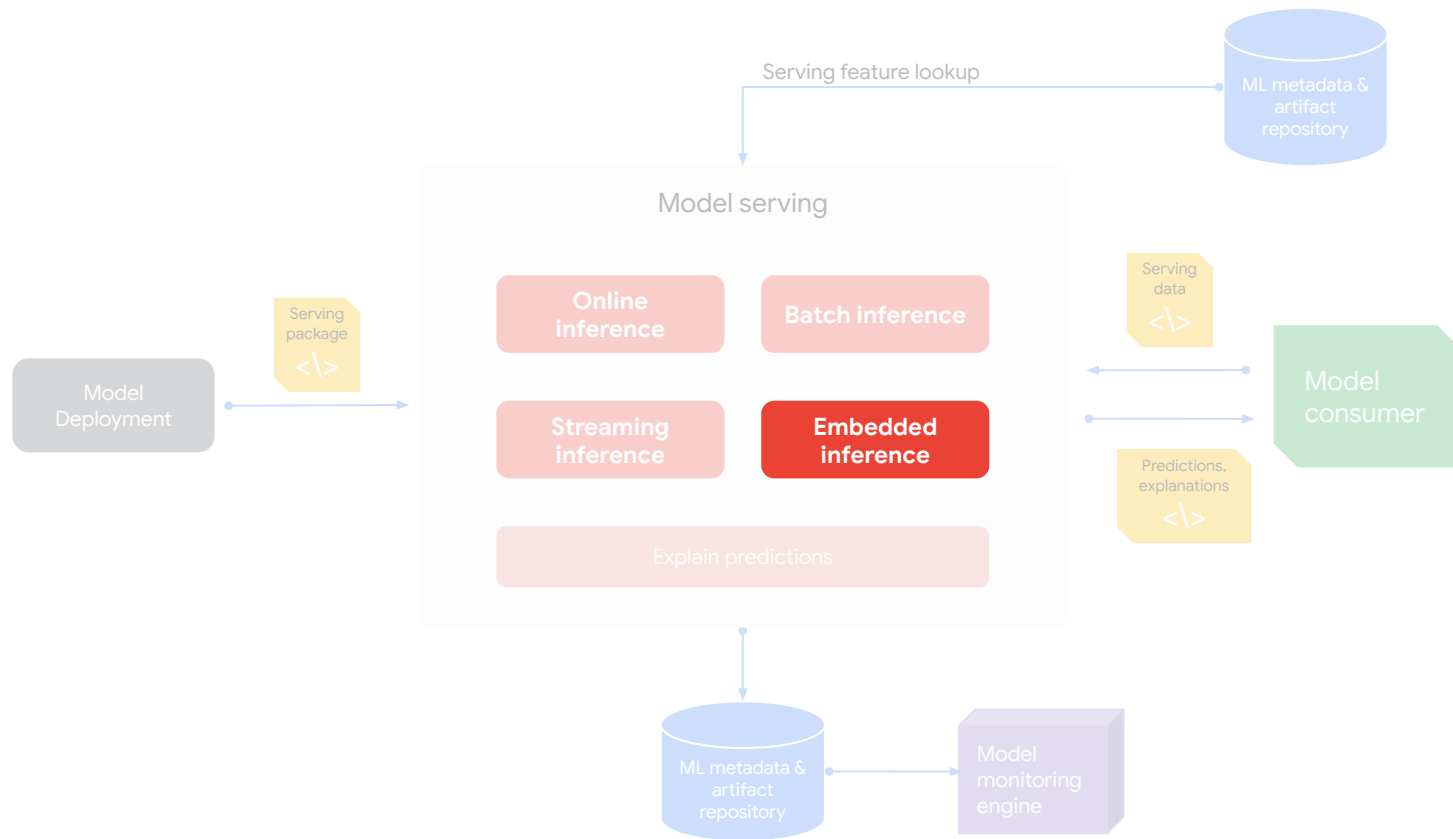


# Prediction Serving Scenarios: Embedded



# MLOps: Prediction Serving



# The MLOps Personas



ML  
Engineer



ML  
Researcher



Data  
Scientist



Data  
Engineer



Software  
Engineer



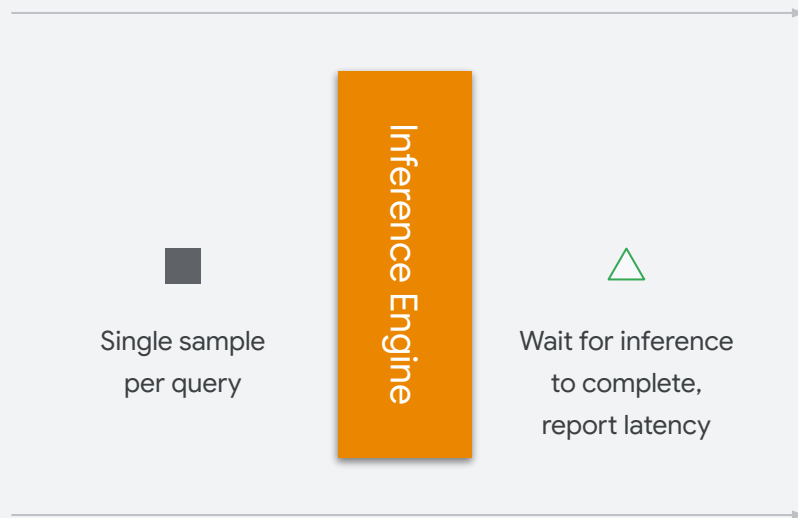
DevOps



Business  
Analyst

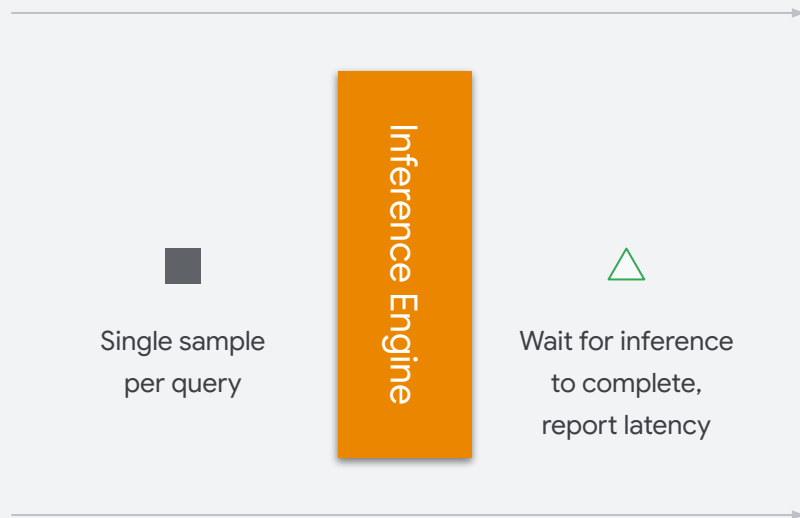
# Embedded Inference: What is it?

- Predict **on-demand**



# Embedded Inference: What is it?

- Predict **on-demand**
- Online inference in near real time for *low-frequency* **singleton** requests



# Embedded Inference:

## When is it useful?

- Smartphone camera
- TinyML use cases
- ...



# Embedded Inference:

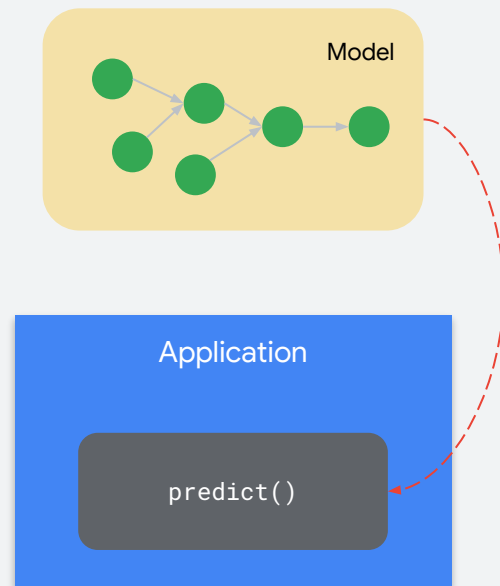
## When is it useful?

- Smartphone camera
- TinyML use cases
- ...



# Embedded Inference: How it works?

Model is packaged into the application for easy deployment at the endpoint device





# Embedded Inference:

## What metrics?

- Single-stream
- Latency metric



Latency

# Embedded Inference:

## Pros & Cons

### Pros

- + Can make a **on-demand** predictions on items

# Embedded Inference:

## Pros & Cons

### Pros

- + Can make a **on-demand** predictions on items
- + Great for
  - + **Bandwidth**
  - + **Latency**
  - + **Energy-efficiency**
  - + **Reliability**
  - + **Privacy**

# Embedded Inference:

## Pros & Cons

### Pros

- + Can make a **on-demand** predictions on items
- + Great for
  - + **Bandwidth**
  - + **Latency**
  - + **Energy-efficiency**
  - + **Reliability**
  - + **Privacy**

### Cons

- Compute **intensive**

# Embedded Inference:

## Pros & Cons

### Pros

- + Can make a **on-demand** predictions on items
- + Great for
  - + **Bandwidth**
  - + **Latency**
  - + **Energy-efficiency**
  - + **Reliability**
  - + **Privacy**

### Cons

- Compute **intensive**
- Latency **sensitive**—may limit model complexity

# Embedded Inference:

## Pros & Cons

### Pros

- + Can make a **on-demand** predictions on items
- + Great for
  - + **Bandwidth**
  - + **Latency**
  - + **Energy-efficiency**
  - + **Reliability**
  - + **Privacy**

### Cons

- Compute **intensive**
- Latency **sensitive**—may limit model complexity
- Monitoring needs are more **important** than for the other types of scenarios

# Embedded Inference:

## Pros & Cons

### Pros

- + Can make a **on-demand** predictions on items
- + Great for
  - + **Bandwidth**
  - + **Latency**
  - + **Energy-efficiency**
  - + **Reliability**
  - + **Privacy**

### Cons

- Compute **intensive**
- Latency **sensitive**—may limit model complexity
- Monitoring needs are more **important** than for the other types of scenarios
- Embedded deployment makes scalability and flexibility **poor**

# Scenario

# Metric



**Batch inference**  
(e.g. photo sorting app)

**Throughput**



**Online inference**  
(e.g. translation app)

**QPS**  
subject to latency bound



**Streaming inference**  
(e.g. multiple camera  
driving assistance)

**Number streams**  
subject to latency bound



**Embedded inference**  
(e.g. cell phone  
augmented vision)

**Latency**



