# Data Science from Whiteboard to Production: End-to-End Image Captioning

TUM Data Innovation Lab

February 18, 2020

Co-Mentor: Michael Rauchensteiner
Project Lead: Dr. Ricardo Acevedo Cabra
Supervisor: Prof. Dr. Massimo Fornasier

# Our Team

Dr. Robert Pesch

Sebastian Blank

Julia Kronburger

Oliver Borchert

Nika Dogonadze
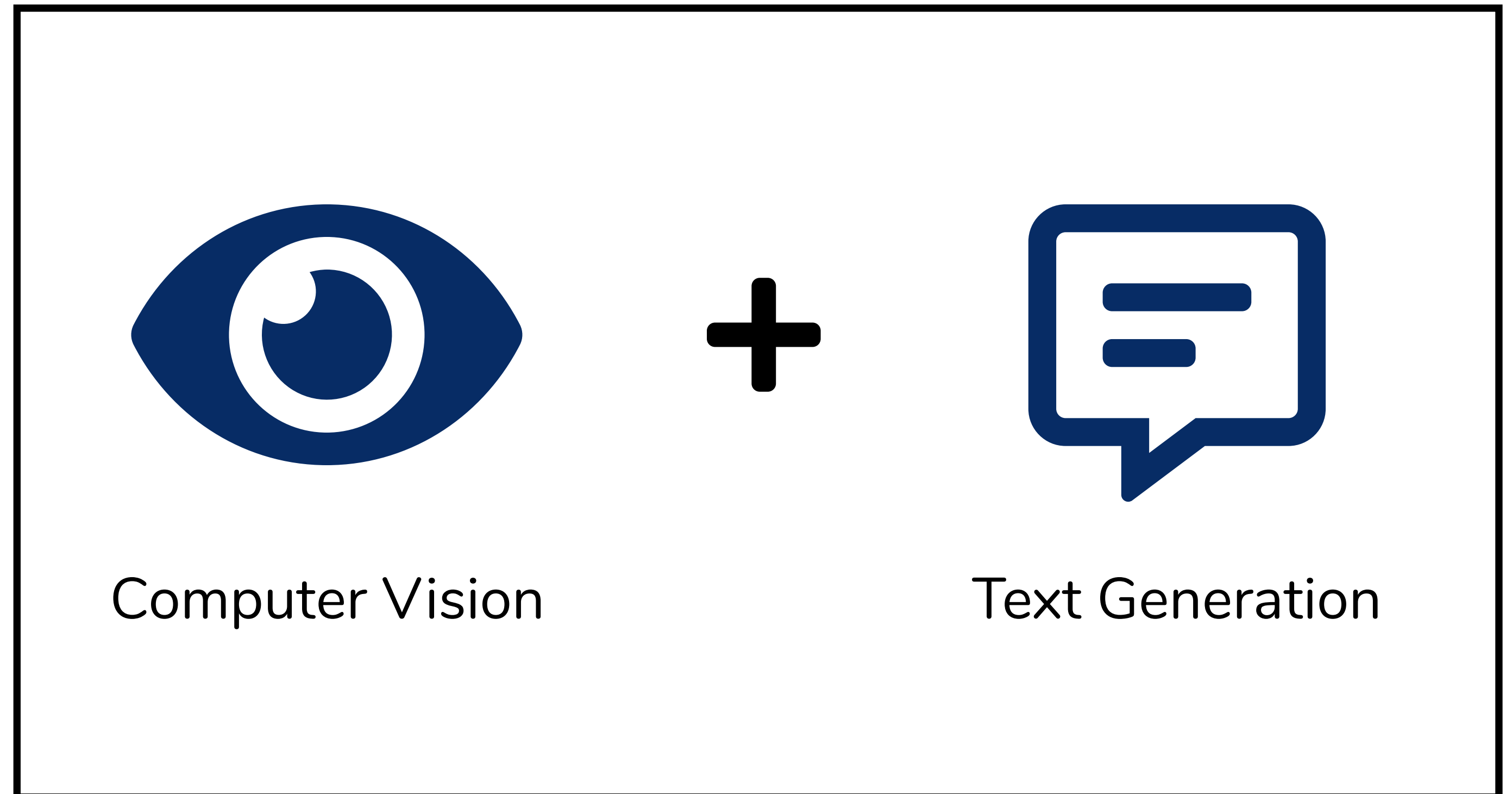
Mürüvvet Hasanbasoglu

Murtaza Raza

Anastasia Stamatouli

# What is Image Captioning?



A dog is playing with
a ball on the grass

Computer Vision + Text Generation

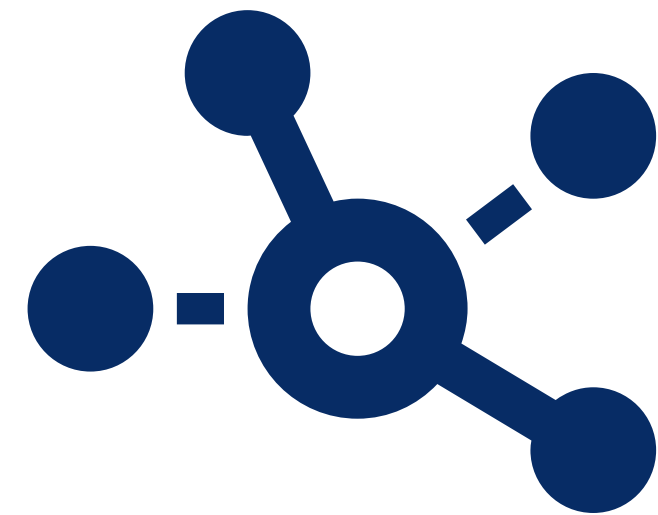# From Whiteboard to Production

**Data Acquisition**

MSCoco

Flickr30k

Flickr8k

**Model Development**
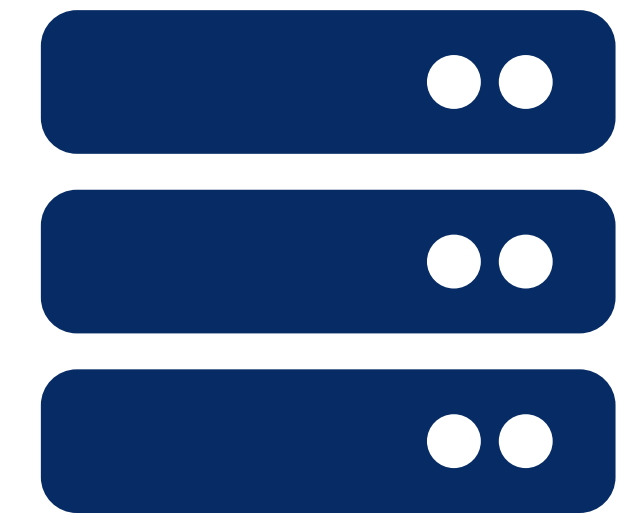
Vanilla Encoder-Decoder

Visual Attention

Beam Search

**Model Evaluation**

BLEU

METEOR

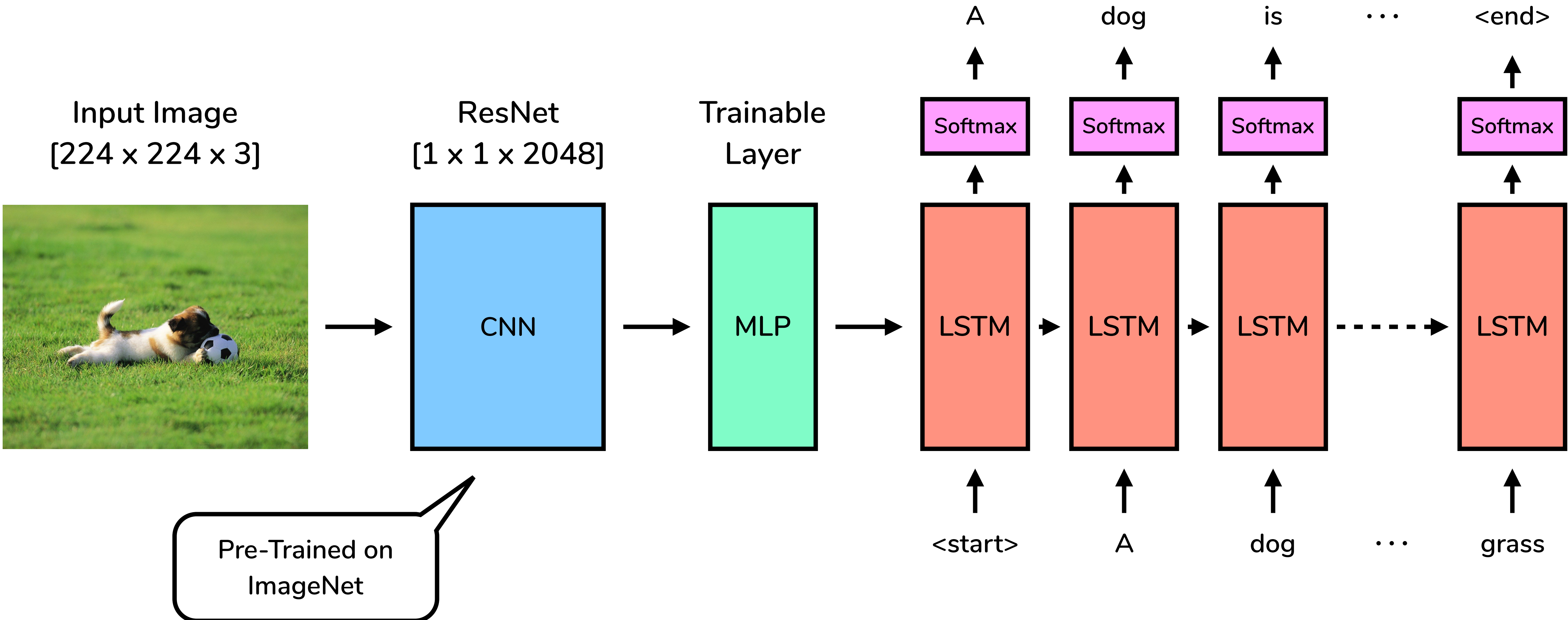ROUGE

CIDEr

**Deployment**

Architecture

Continuous Integration
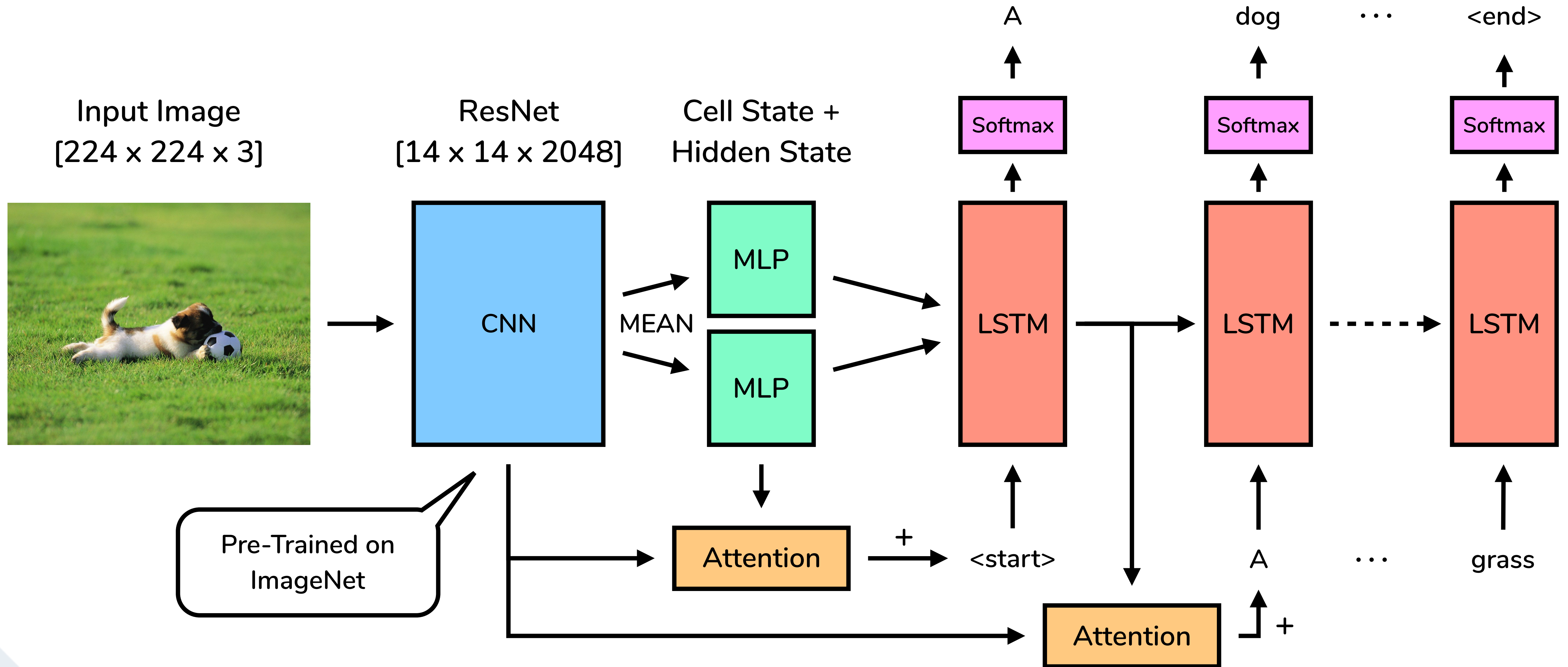
Monitoring

# Model Development and Evaluation
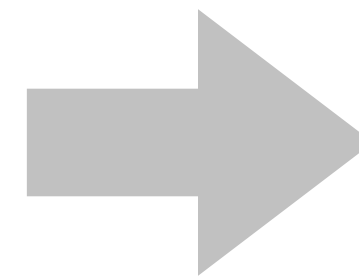
# Vanilla Encoder-Decoder

# Visual Attention
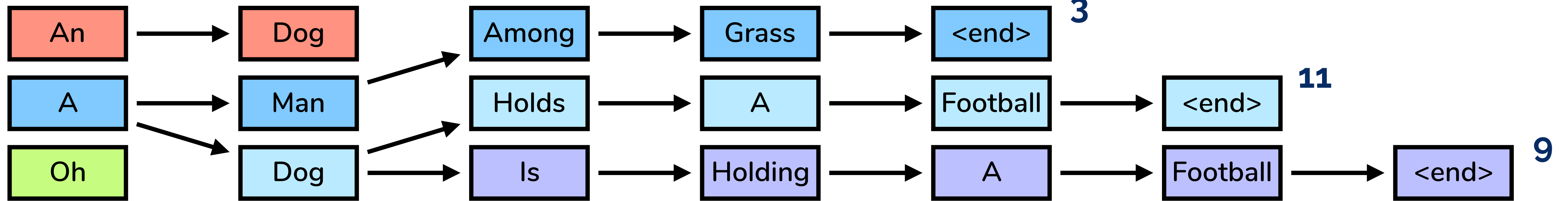
# Attentive Encoder-Decoder

# Beam Search



**Greedy Search**

Sequence of Most Likely States

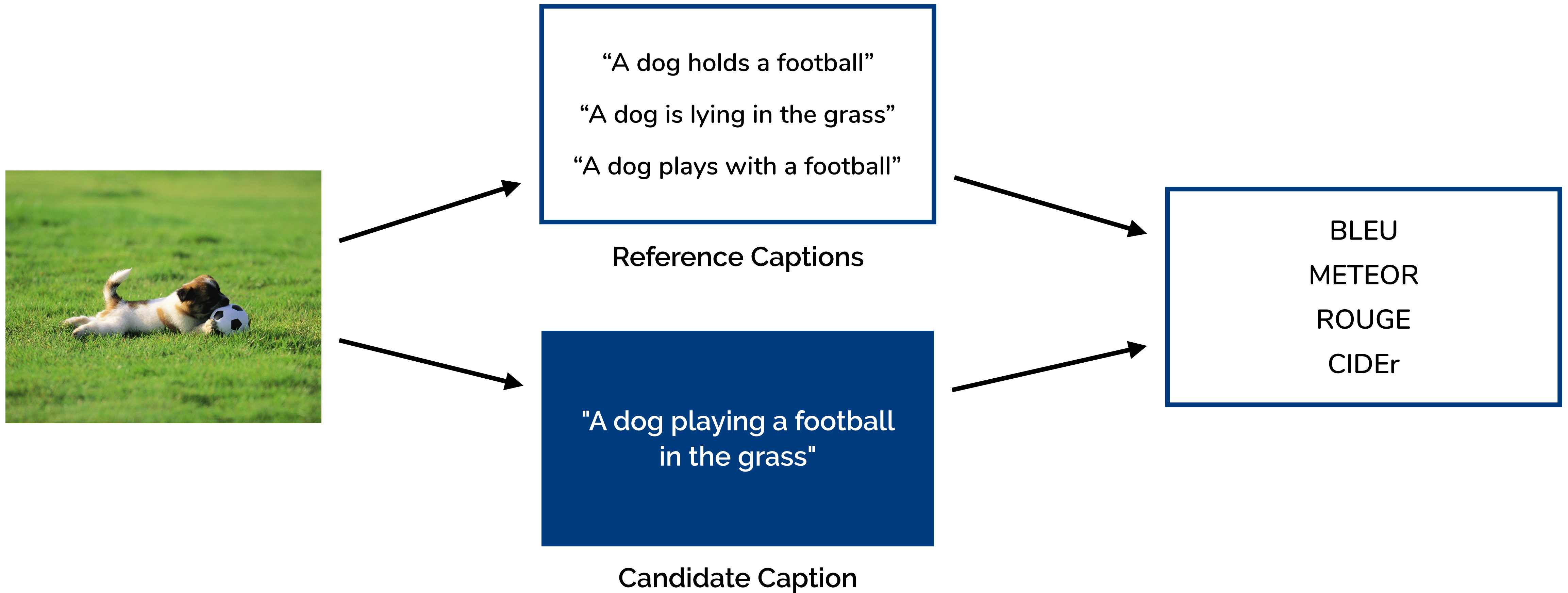**Beam Search**

Most Likely Sequence of States

# Evaluation Metrics



"A dog holds a football"

"A dog is lying in the grass"

"A dog plays with a football"

**Reference Captions**

"A dog playing a football in the grass"

**Candidate Caption**

BLEU

METEOR

ROUGE

CIDEr

# Evaluation Metrics

"A dog holds a football"

"A dog is lying in the grass"

"A dog plays with a football"

Reference Captions

## BLEU

Exact n-gram matches

## METEOR

Exact n-gram + word stem + paraphrase + synonym matches

## ROUGE

Longest common subsequence

## CIDEr

TF-IDF weighted n-gram matches

# Training Results

| | BLEU-1 | BLEU-4 | METEOR | ROUGE | CIDEr |
|---|---|---|---|---|---|
| **Vanilla Encoder-Decoder** | 41.4 | 11.0 | 19.6 | 43.7 | 35.3 |
| **Attentive Encoder-Decoder** | 70.8 | 24.0 | 22.7 | 52.9 | 48.3 |
| **+ Beam Search** | **72.7** | **30.3** | **24.1** | **54.4** | **87.6** |
| **Show, Attend & Tell [1]** | 70.7 | 24.3 | 23.9 | — | — |

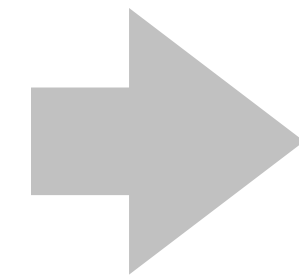[1] Xu, K. et. al. "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention" (ICML 2015)

Models trained and evaluated on MSCOCO
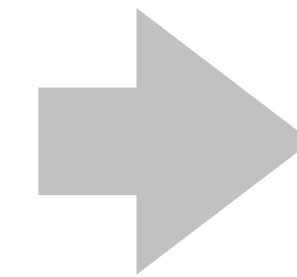
# Experiment and Model Management

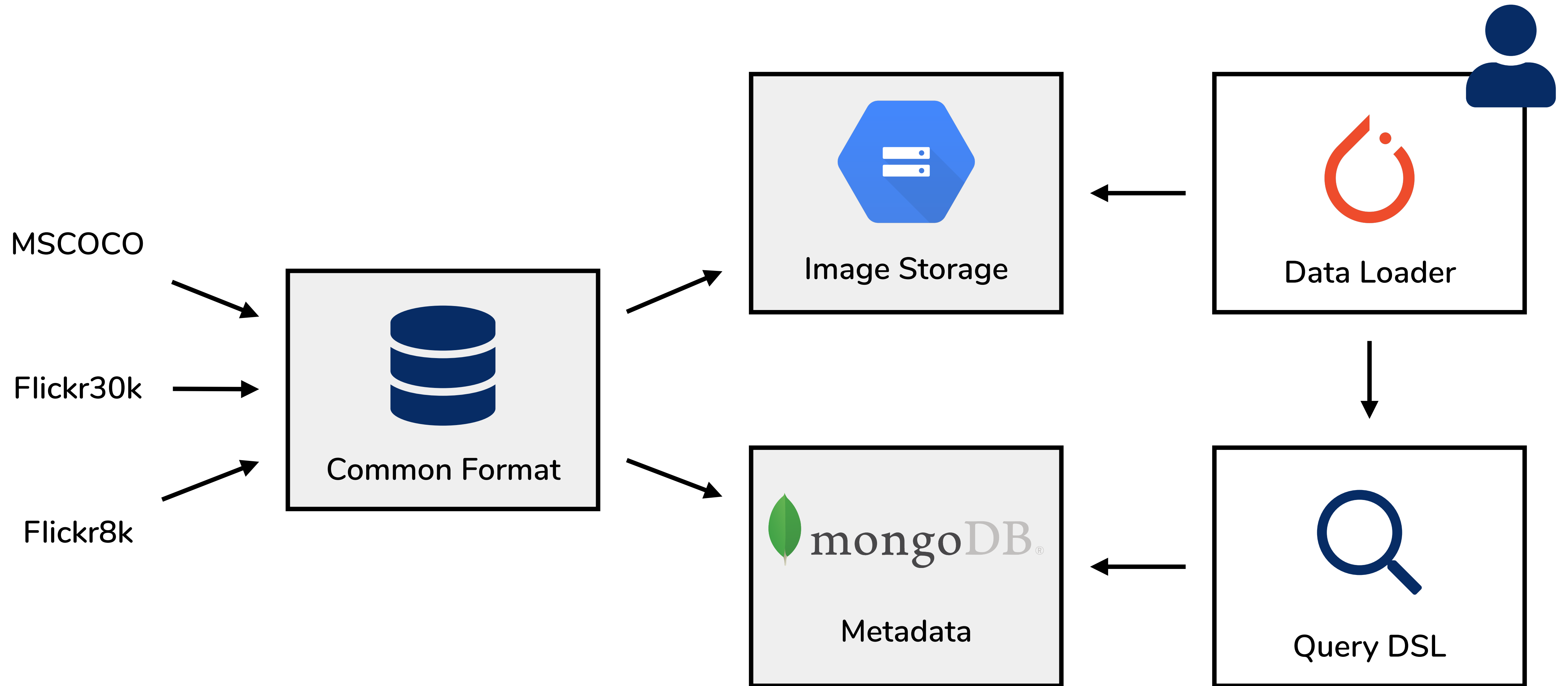# Evaluation Framework

**Data Loading Interface**

**Training Interface**
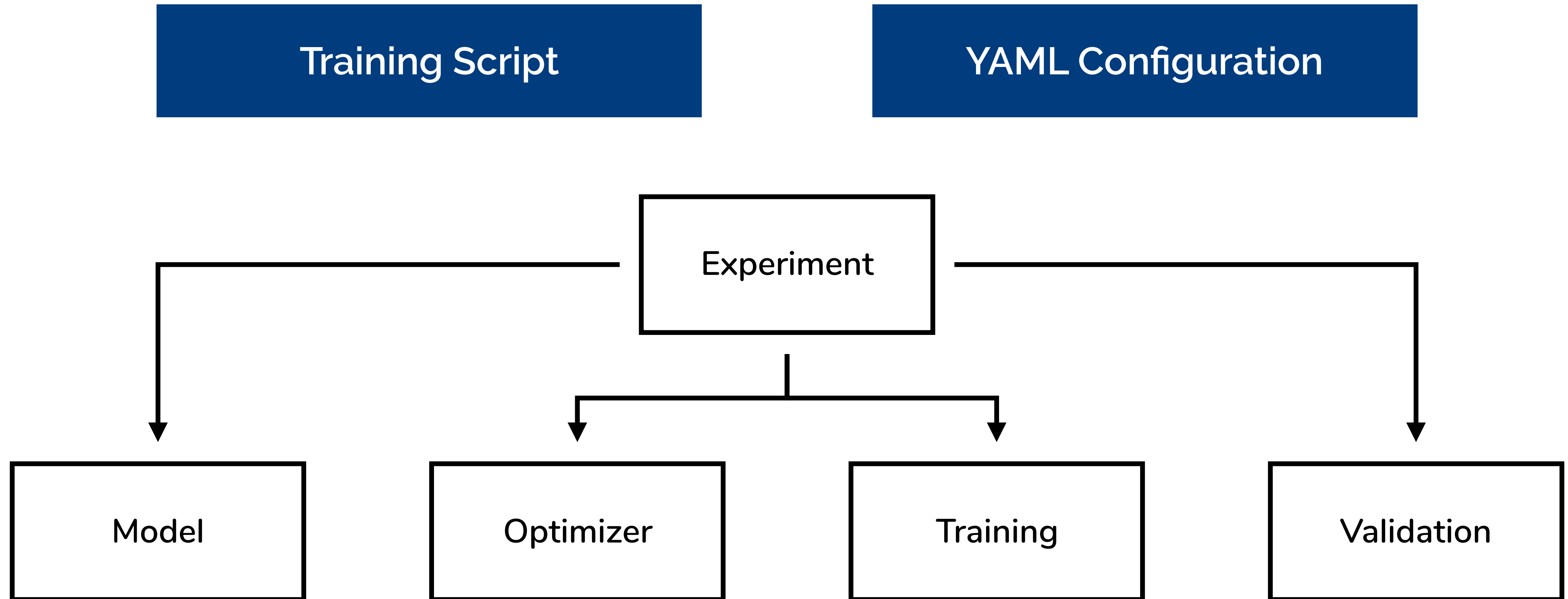
**Experiment Tracking**

# Data Loading Interface



MSCOCO

Flickr30k

Flickr8k

Common Format

Image Storage

Metadata

Data Loader

Query DSL

# Query DSL

# Training Interface

Training Script

YAML Configuration

Experiment

Model

Optimizer

Training

Validation

17

# Experiment Tracking

# Deployment

# Production-Grade Deployment

Reliability

Security

Logging

Accurate Results

Continuous Model Improvement

Friendly UI

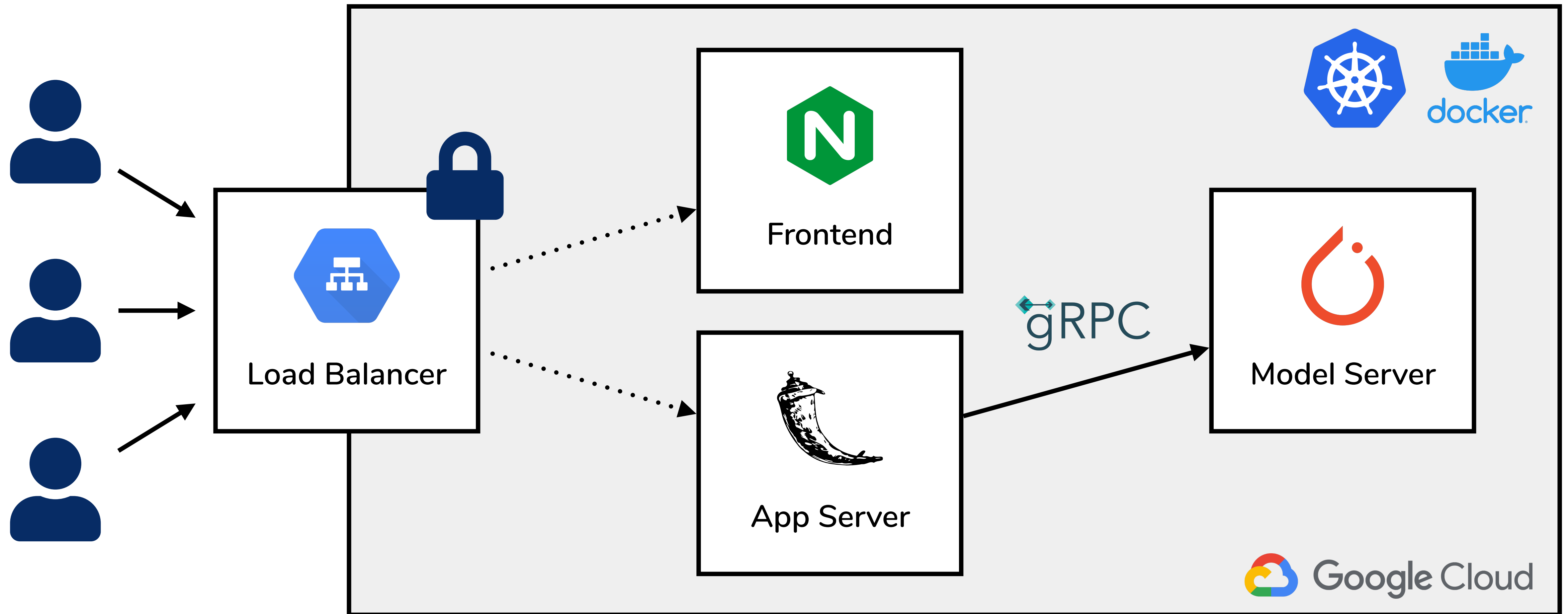Responsive UI

Fault Tolerance
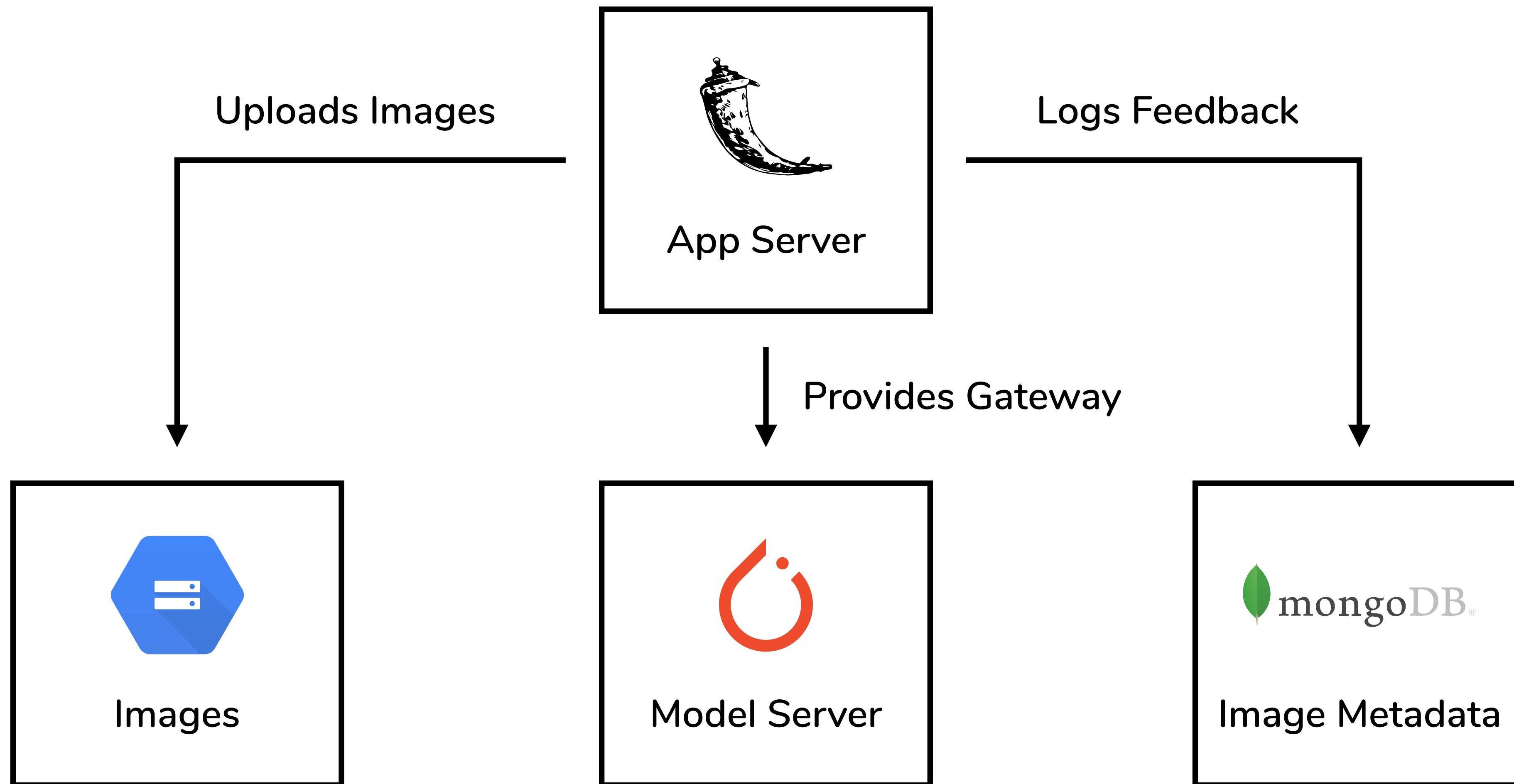
High Availability

Monitoring

Auto-Scalability
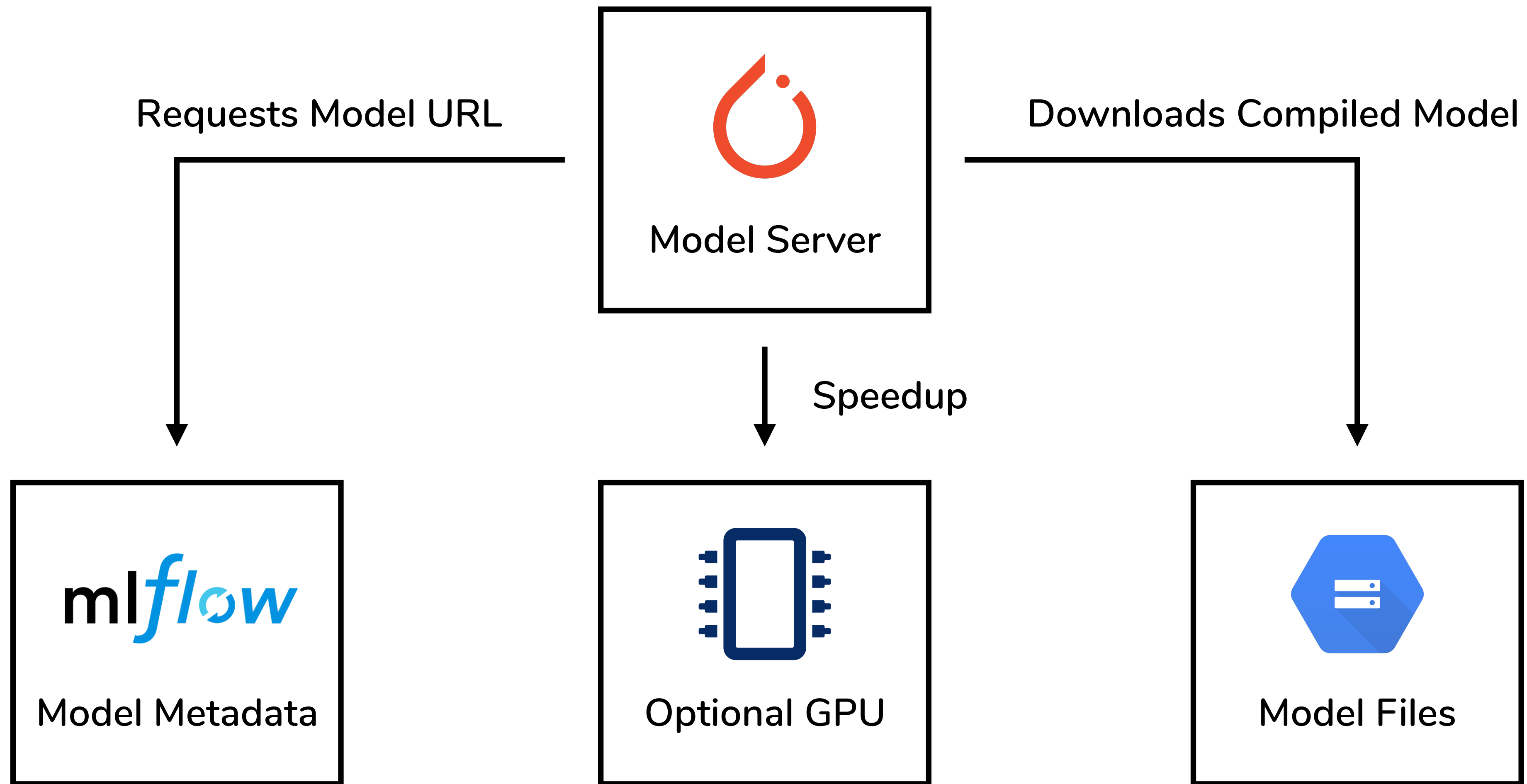
REST API

Fast Response Times

# System Architecture

# Application Server



App Server

Uploads Images

Logs Feedback

Provides Gateway

Images
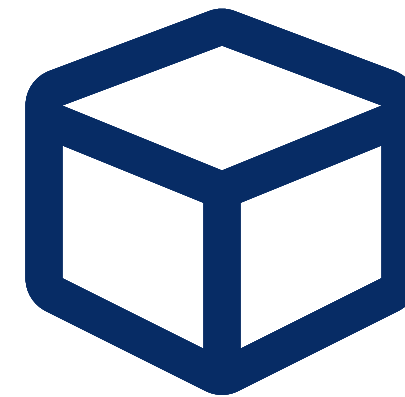
Model Server

Image Metadata
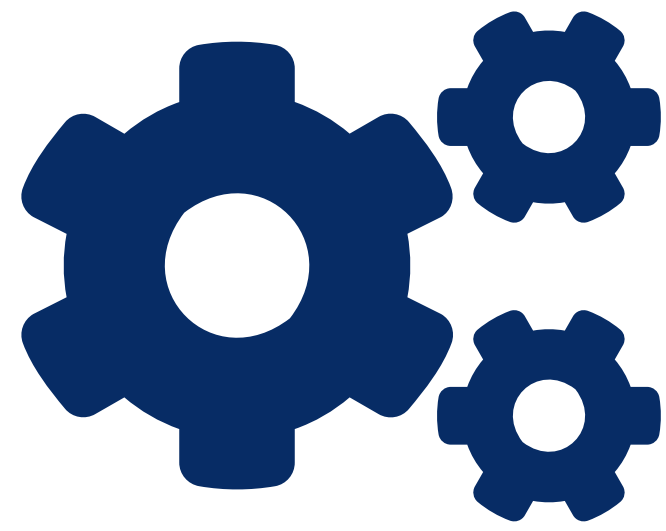
# Model Server

# Autoscaling
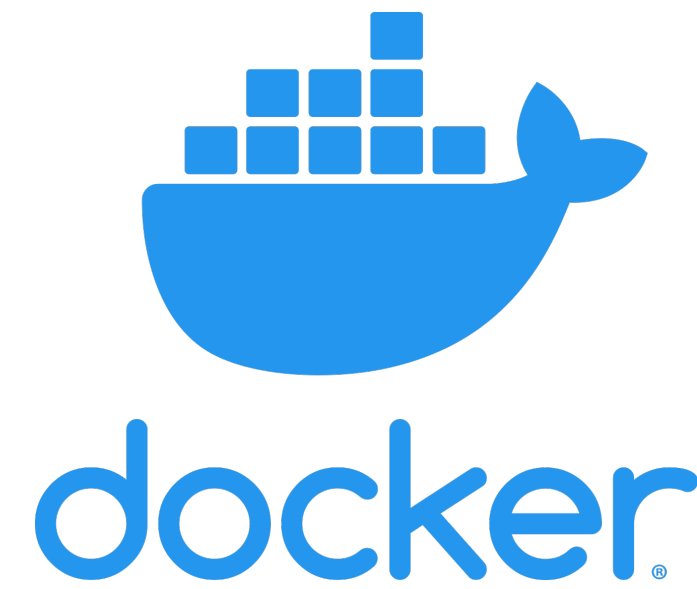
**Demand-Based Number of Replicas**

Independent Scaling

Automatic Load Balancing

# Continuous Integration & Deployment



Run Automated
Tests

Build & Tag
Docker Image

Bundle Required
Resources

Deploy Resources
to Production

# Monitoring


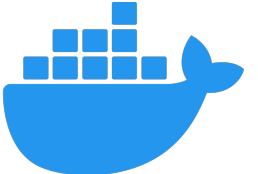
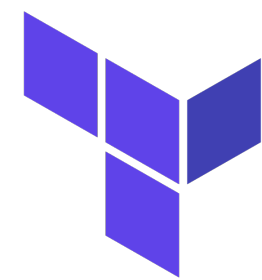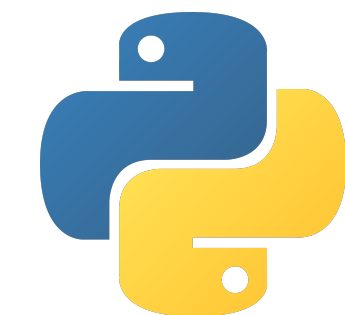Linkerd  +  Prometheus  +  Grafana

# Tools & Technologies

# Thanks for your attention!

Let's see it in action...

https://dilab.inovex.de/