# Putting Deep Learning Models in Production

## Sahil Dua

@sahildua2305

Booking.com

# Let's imagine!

But ...

# whoami

➔ Software Developer @ Booking.com

➔ Previously - Deep Learning Infrastructure

➔ Open Source Contributor (Git, Pandas, Kinto, go-github, etc.)

➔ Tech Speaker

Booking.com

# Agenda

➔ Deep Learning at Booking.com
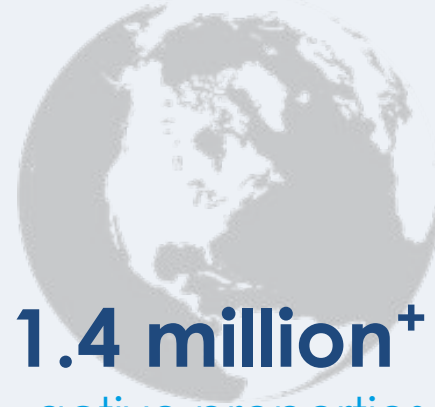
➔ Life-cycle of a model

➔ Training Models

➔ Serving Predictions

Booking.com

# Scale
highlights.

1,500,000+
room nights
booked
every 24 hours

1.4 million+
active properties
in 220+ countries

@sahildua2305

Booking.com

# Deep Learning

➔ Image understanding

➔ Translations

➔ Ads bidding

➔ ...

# Image Tagging

# Image Tagging



| Classes | Score | | |
|---|---|---|---|
| oceanfront | **0.79** | 0 | 1 |
| nature | **0.79** | 0 | 1 |
| beach house | **0.62** | 0 | 1 |
| building | **0.62** | 0 | 1 |
| penthouse | **0.61** | 0 | 1 |
| apartment | **0.61** | 0 | 1 |
| housing | **0.61** | 0 | 1 |

# Image Tagging



| Classes | Score | | |
|---|---|---|---|
| oceanfront | **0.79** | 0 | 1 |
| nature | **0.79** | 0 | 1 |
| beach house | **0.62** | 0 | 1 |
| building | **0.62** | 0 | 1 |
| penthouse | **0.61** | 0 | 1 |
| apartment | **0.61** | 0 | 1 |
| housing | **0.61** | 0 | 1 |

**Sea view:** 6.38
**Balcony/Terrace:** 4.82
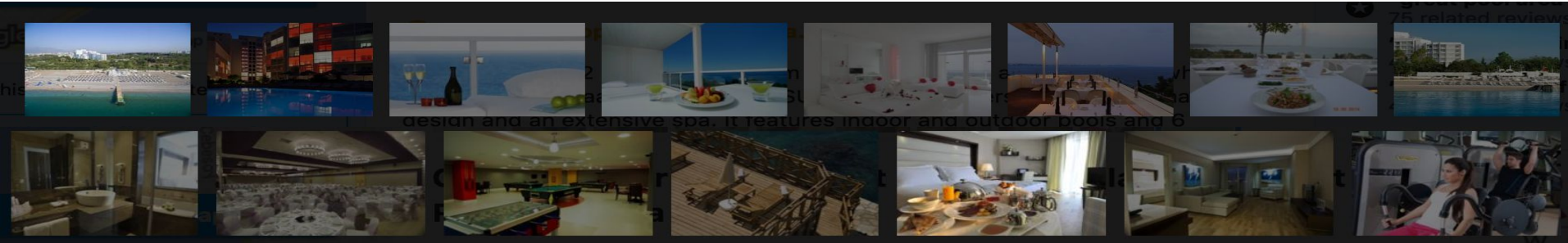**Photo of the whole room:** 4.21
**Bed:** 3.47
**Decorative details:** 3.15
**Seating area:** 2.70

Booking.com

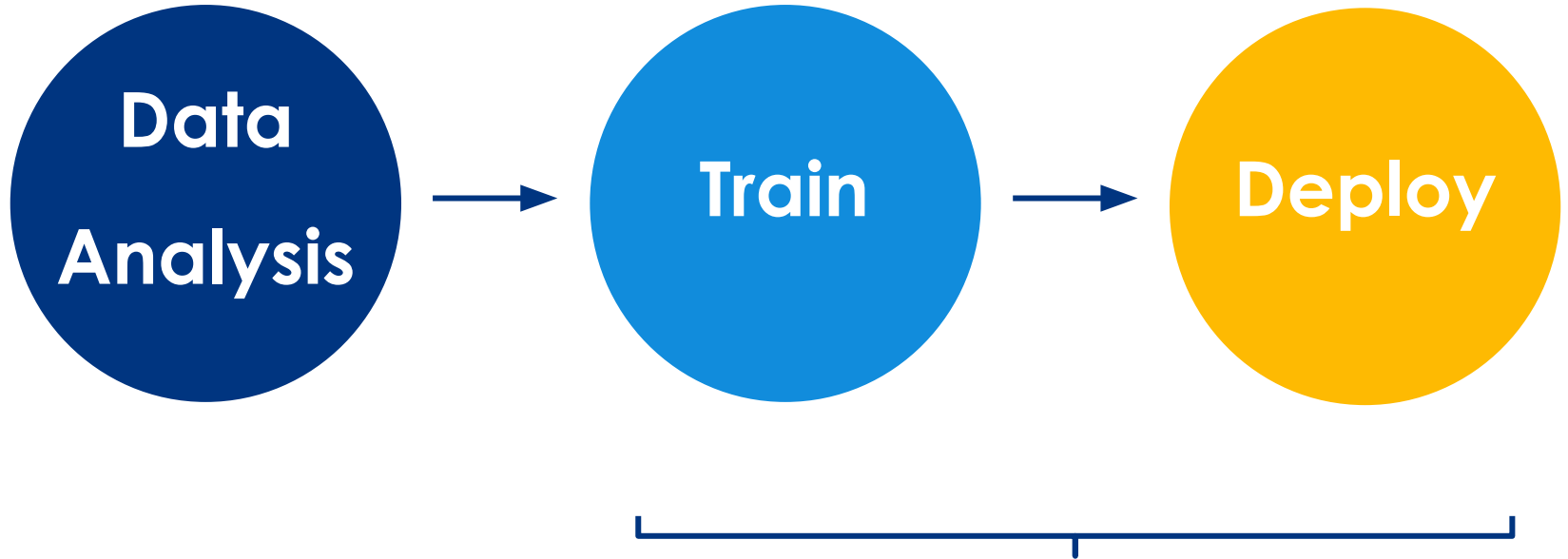# Image Tagging

**Using the image tag information in the right context**
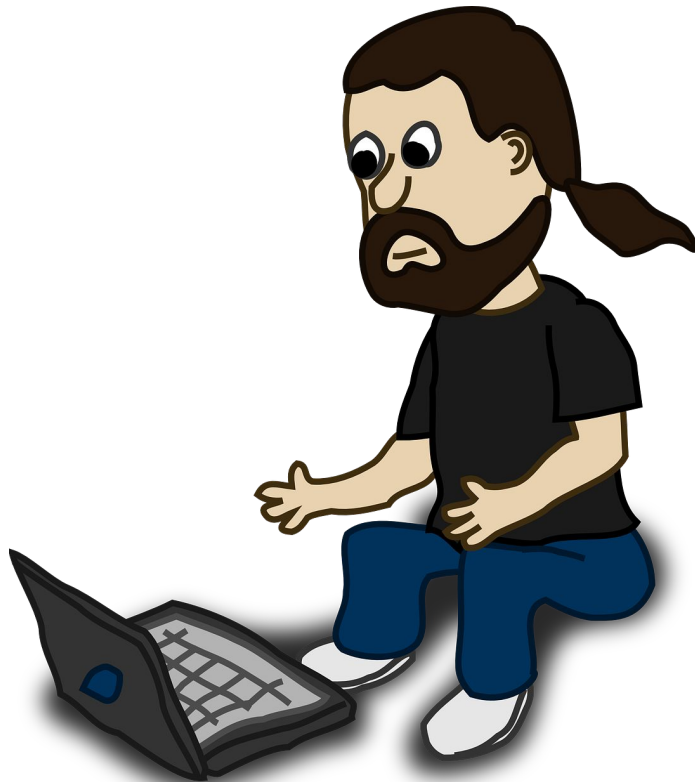Swimming pool, Breakfast Buffet, etc.

Booking.com

# Lifecycle of a model

# Lifecycle of a model

**Data Analysis** → **Train** → **Deploy**

# Training a Model - on laptop

# Training a Model - on laptop



@sahildua2305

Booking.com

# Machine Learning workload

➔ Computationally intensive workload

➔ Often not highly parallelizable algorithms

➔ 10 to 100 GBs of data

Booking.com

# Why Kubernetes (k8s)?

➔ Isolation

➔ Elasticity
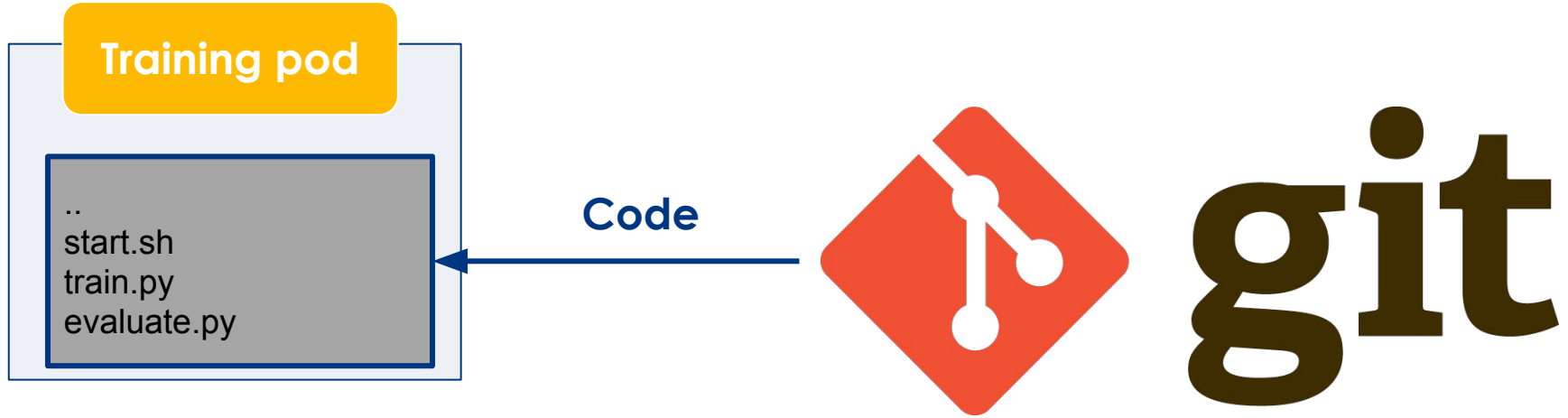
➔ Flexibility

# Why k8s – GPUs?

➔ In alpha since 1.3

➔ Speed up 20X-50X

```
resources:
  limits:
    alpha.kubernetes.io/nvidia-gpu: 1
```
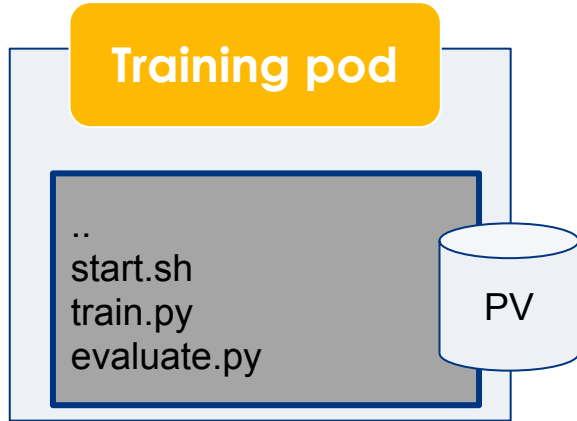
Booking.com

# Training with k8s

➜ Base images with ML frameworks

◆ TensorFlow, Torch, VowpalWabbit, etc.

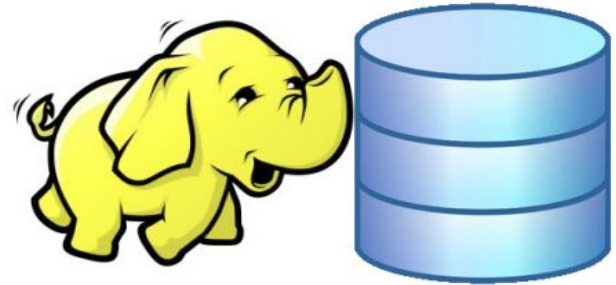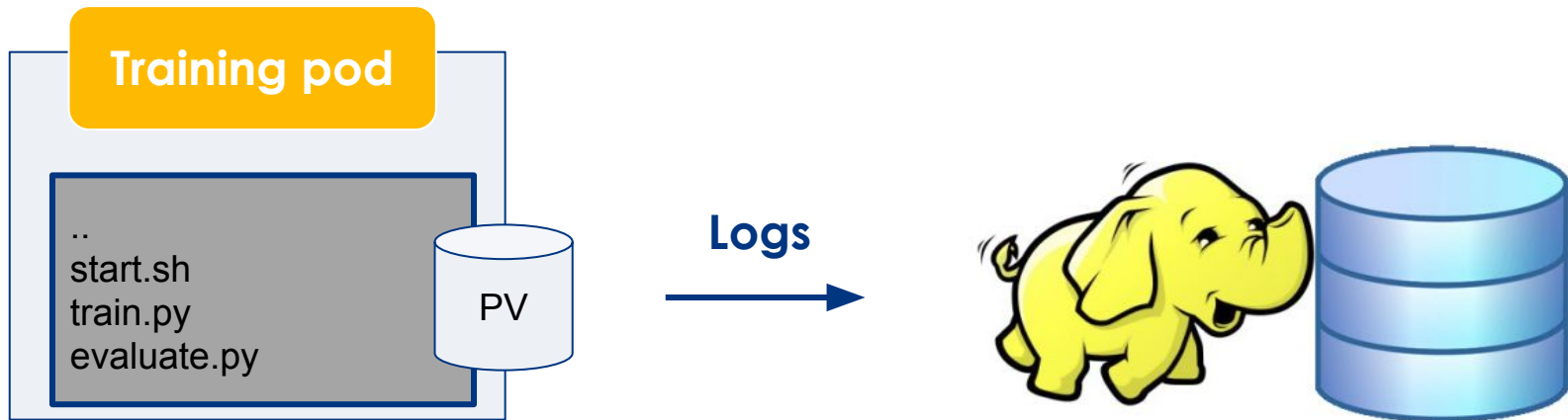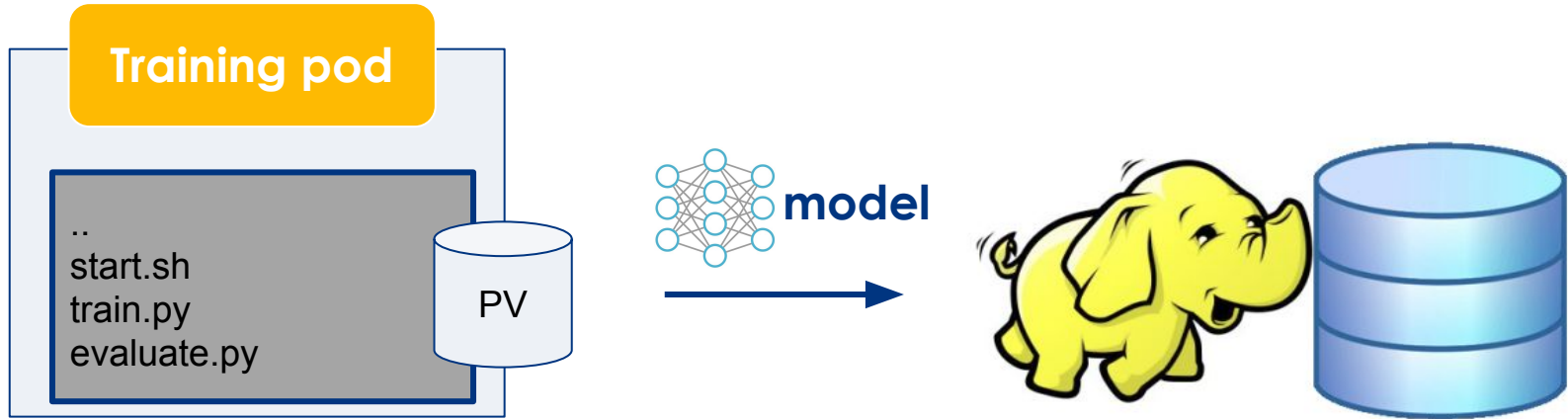➜ Training code is installed at start time

➜ Data access - Hadoop (or PVs)

Booking.com

# Startup



**Training pod**

```
..
start.sh
train.py
evaluate.py
```

**Code**

git

# Startup

**Training pod**

```
..
start.sh
train.py
evaluate.py
```

PV

**Data**

# Streaming logs back

**Training pod**

```
..
start.sh
train.py
evaluate.py
```

PV

**Logs** →

Booking.com

# Exports the model



@sahildua2305

Booking.com

# Serving predictions

# Serving Predictions

Booking.com

# Serving Predictions

Client

Input Features →

← Prediction

**Model 1**

Client

Input Features →

← Prediction

**Model X**

Booking.com

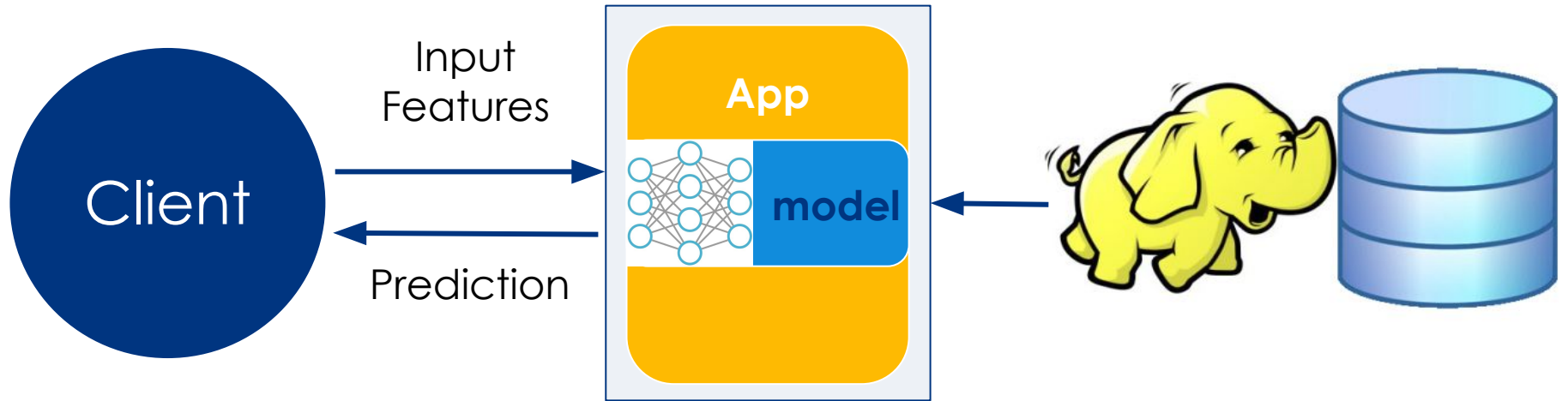# Serving Predictions



@sahildua2305

Booking.com

# Serving Predictions

➔ Stateless app with common code

➔ Containerized

➔ No model in image

➔ REST API for predictions

Booking.com

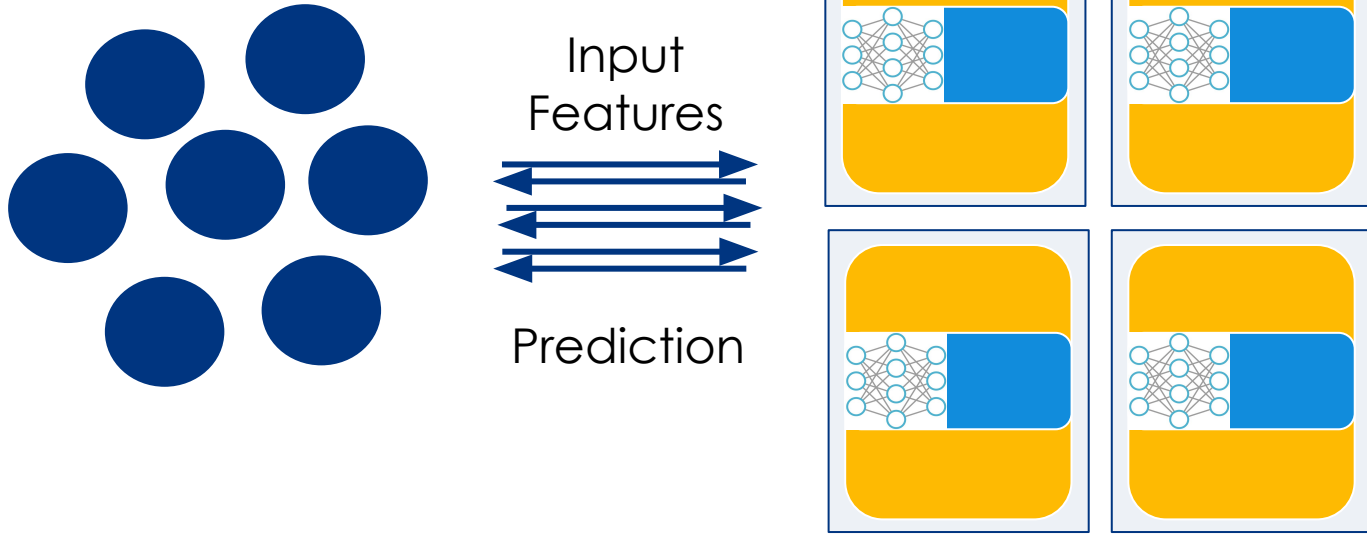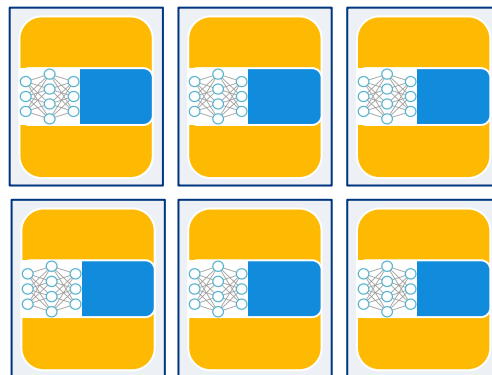# Serving Predictions



Client → Input Features → App [model]
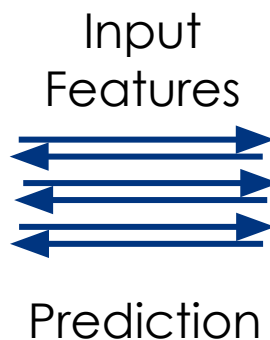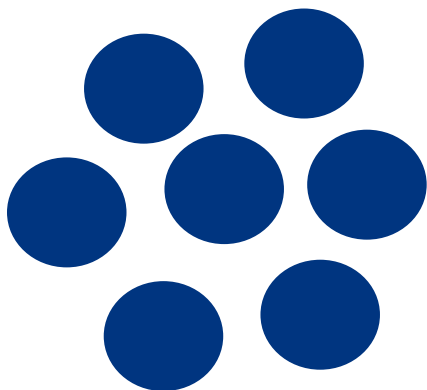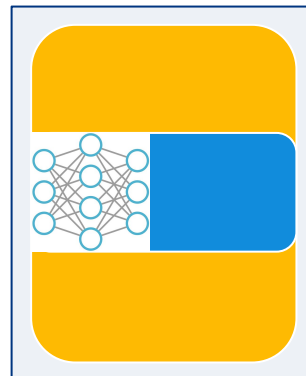
App → Prediction → Client

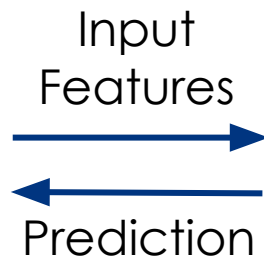# Serving Predictions

➜ Get trained model from Hadoop

➜ Load model in memory

➜ Warm it up

➜ Expose HTTP API

➜ Respond to the probes

Booking.com

# Serving Predictions

Input
Features

Prediction

# Serving Predictions



@sahildua2305

Booking.com

# Deploying a new model

➜ Create new Deployment

➜ Create new HTTP Route

➜ Wait for liveness/readiness probe

Booking.com

# Performance

**PredictionTime = RequestOverhead + N*ComputationTime**

*N is the number of instances to predict on*

Booking.com

# Optimizing for Latency

➔ Do not predict if you can precompute

➔ Reduce Request Overhead

➔ Predict for one instance

➔ Quantization (float 32 => fixed 8)

➔ TensorFlow specific: freeze network & optimize for inference

Booking.com

# Optimizing for Throughput

➔  Do not predict if you can precompute

➔  Batch requests

➔  Parallelize requests

Booking.com

# Summary

➔   Training models in pods

➔   Serving models

➔   Optimizing serving for latency/throughput

# Next steps

➔ Tooling to control hundred deployments

➔ Autoscale prediction service

➔ Hyper parameter tuning for training

Booking.com

# Want to get in touch?

LinkedIn / Twitter / GitHub

**@sahildua2305**

Website

**www.sahildua.com**

Booking.com

**THANK YOU**

@sahildua2305