

Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see

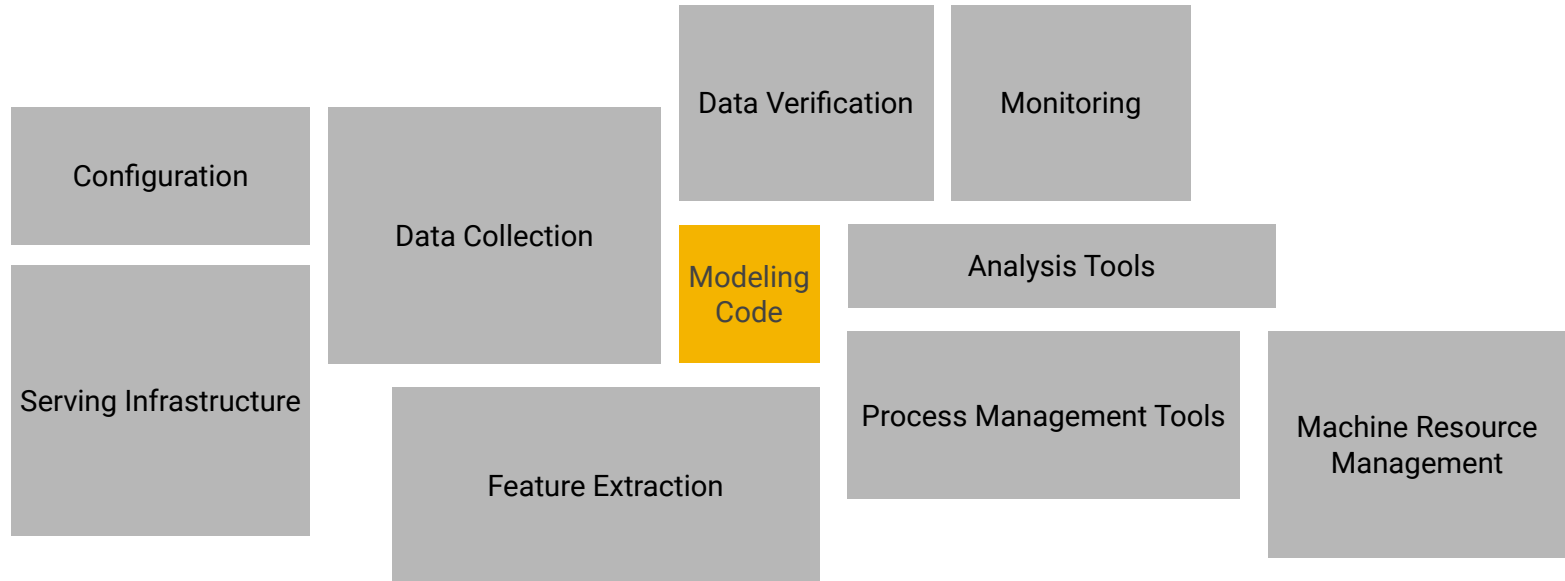
<https://creativecommons.org/licenses/by-sa/2.0/legalcode>

Building Models is just a small part of ML...

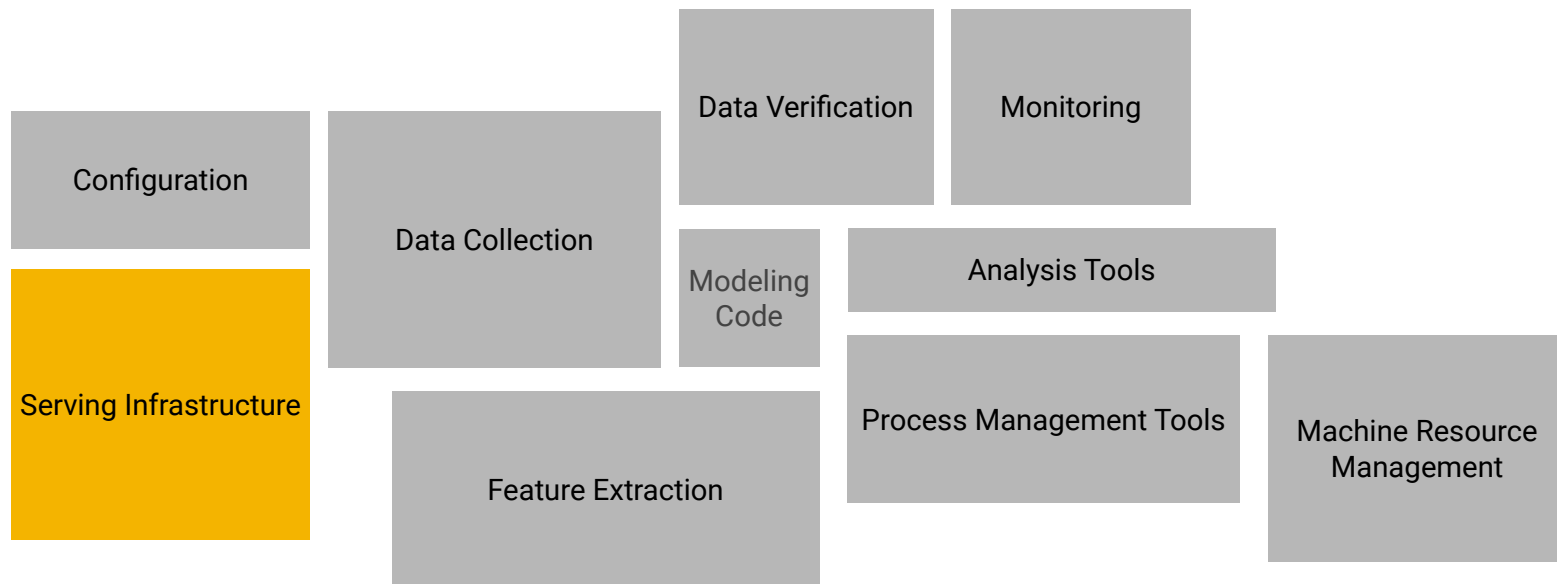


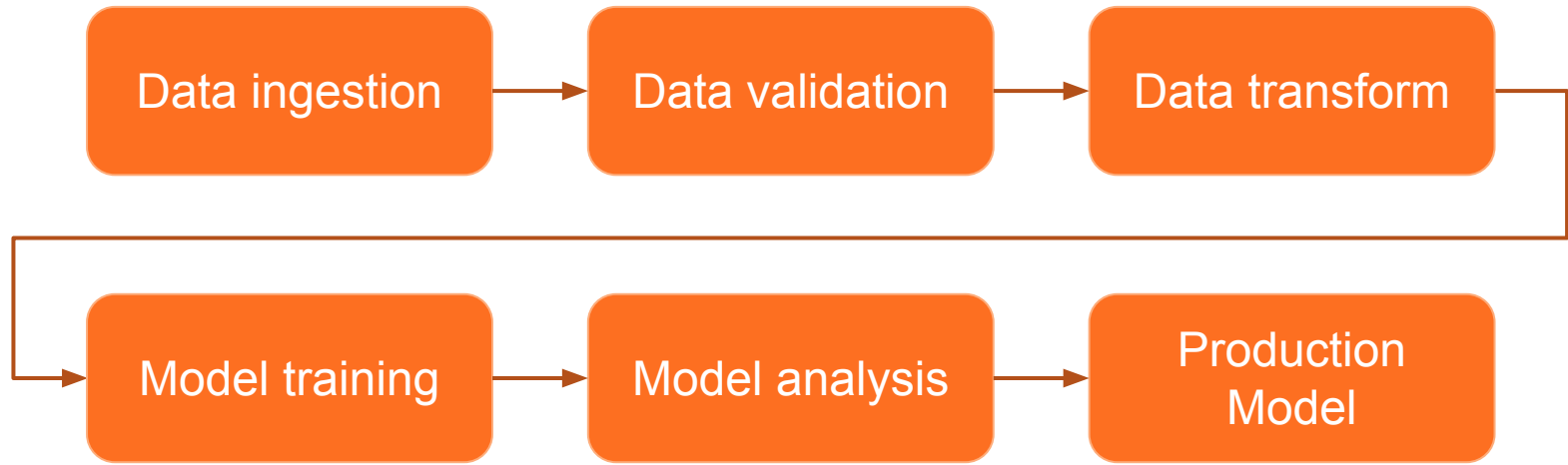
Modeling Code

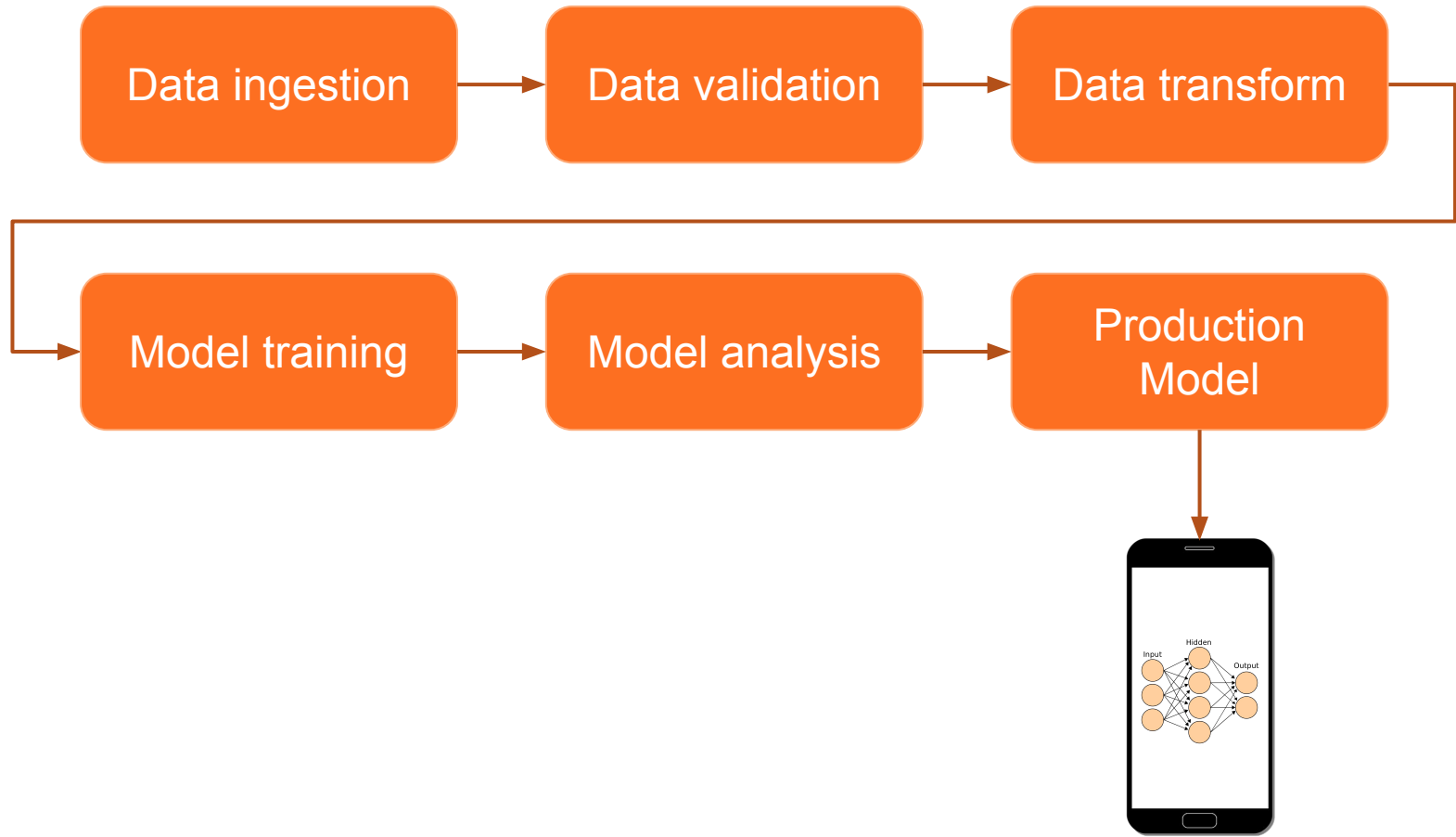
... a production solution requires so much more

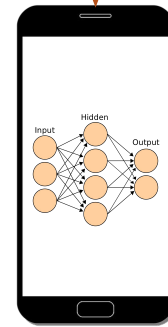
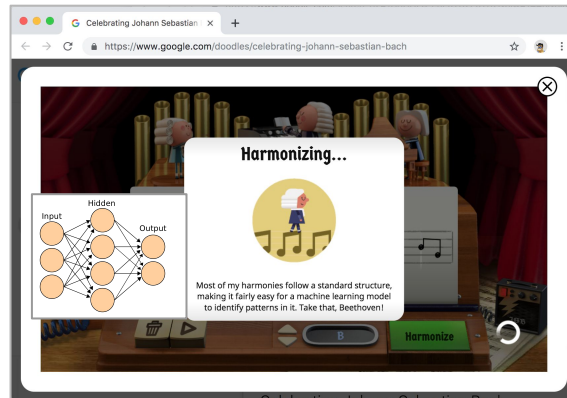
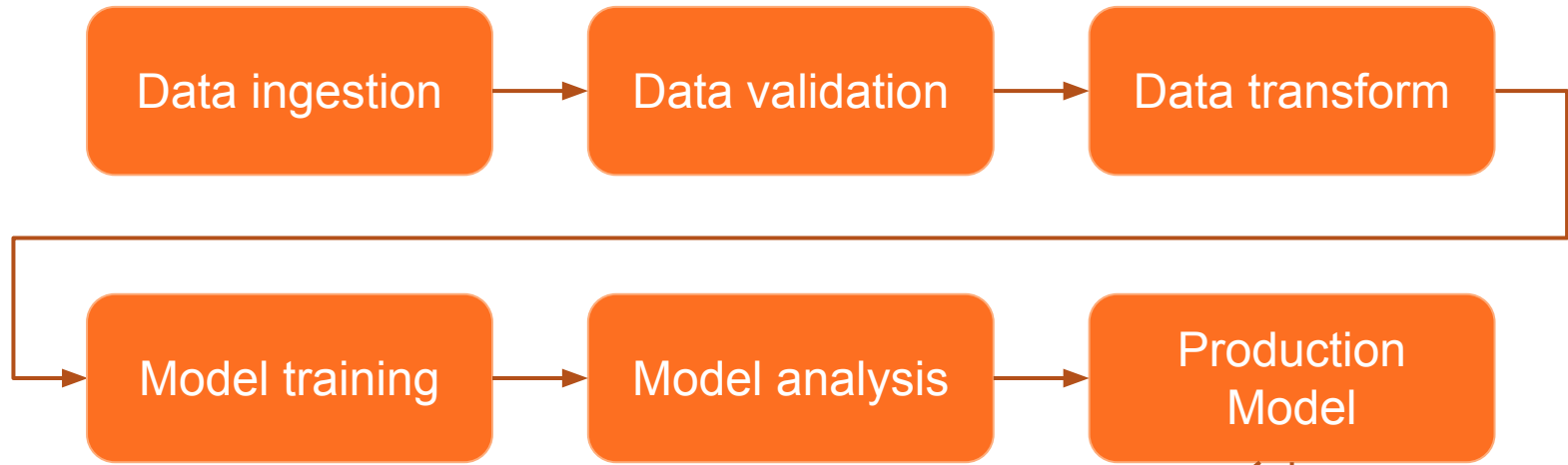


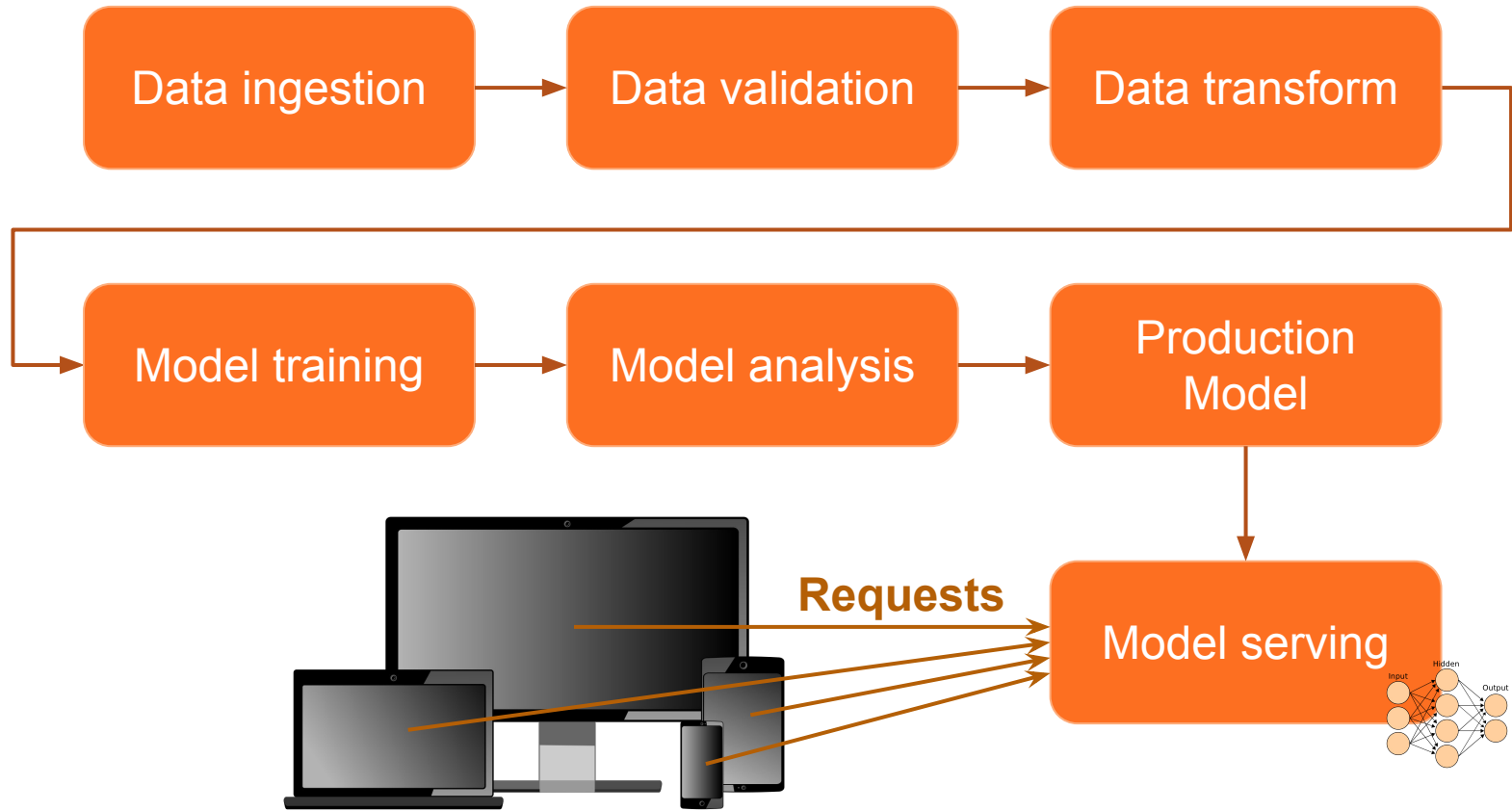
So this week we will focus on Serving...

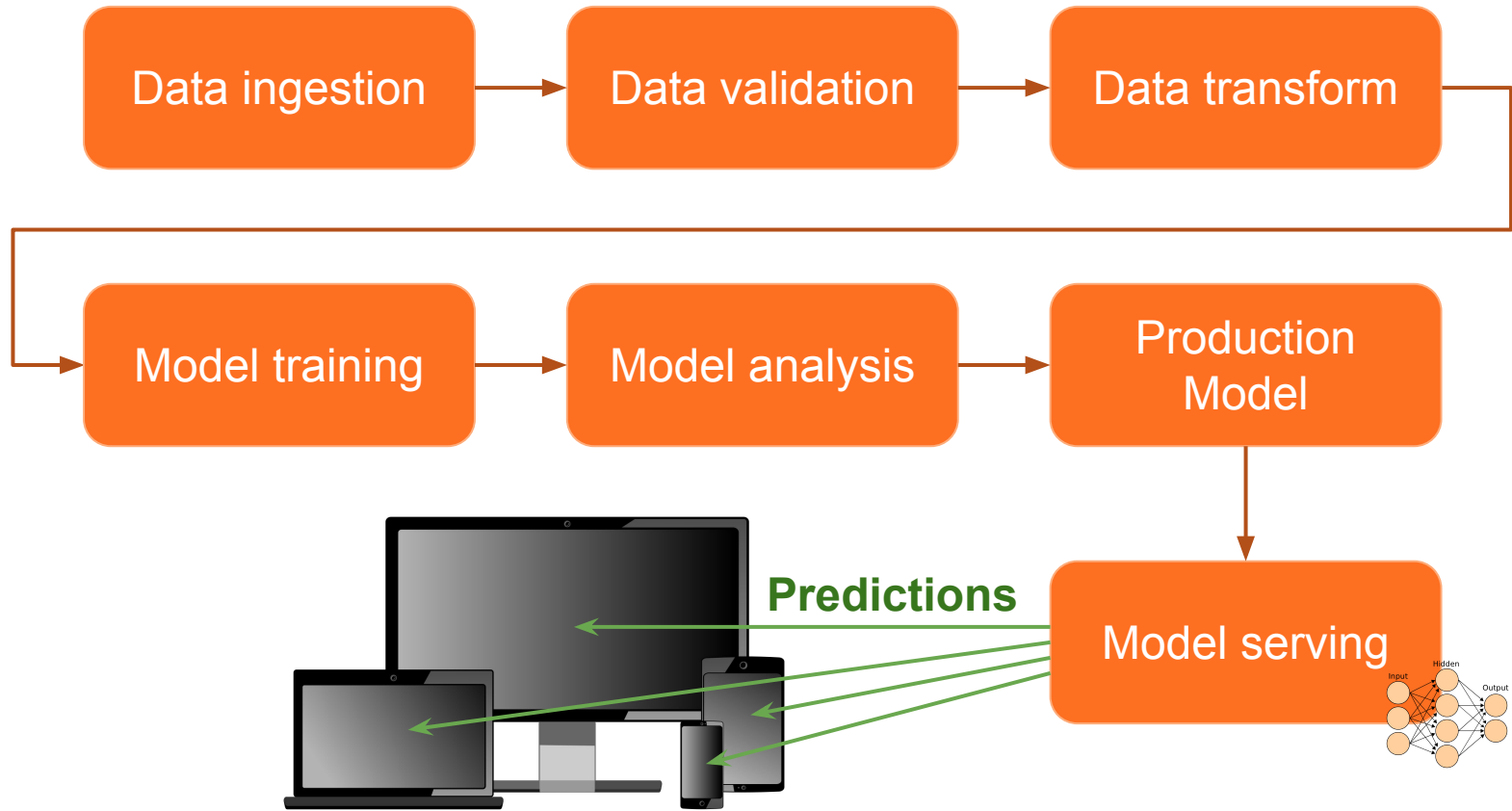














Client 1

Model v1

Client 2

Model v1

Client 3

Model v1



Client 1

Model v1

Client 2

Model v2

Client 3

Model v1



Client 1

Model v2

Client 2

Model v2

Client 3

Model v1



Client 1

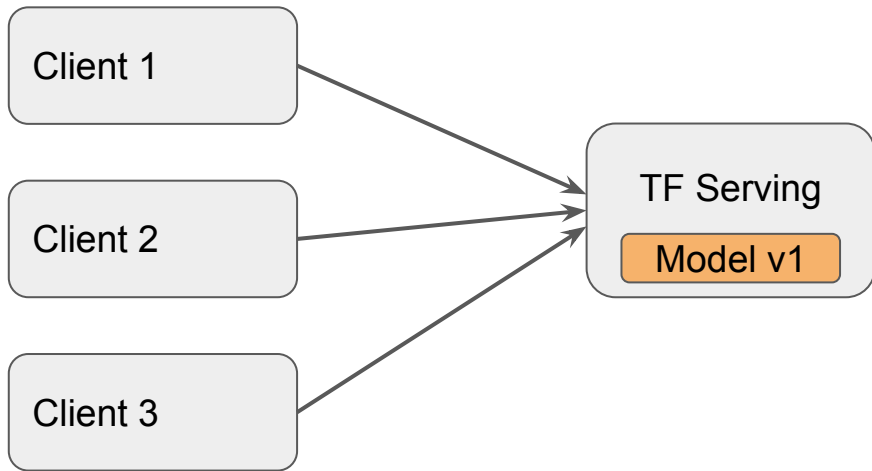
Model v2

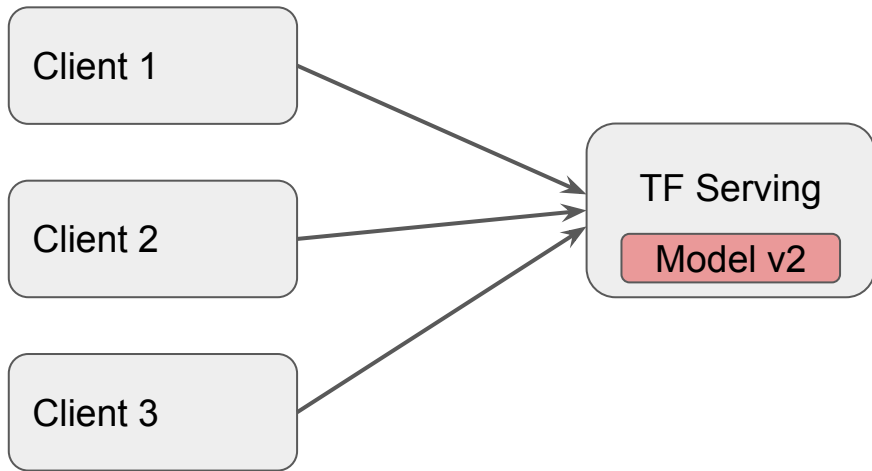
Client 2

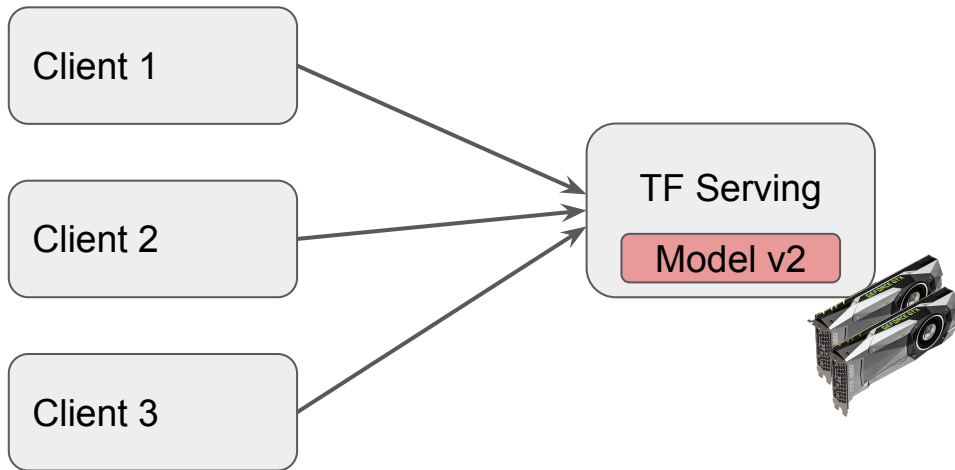
Model v2

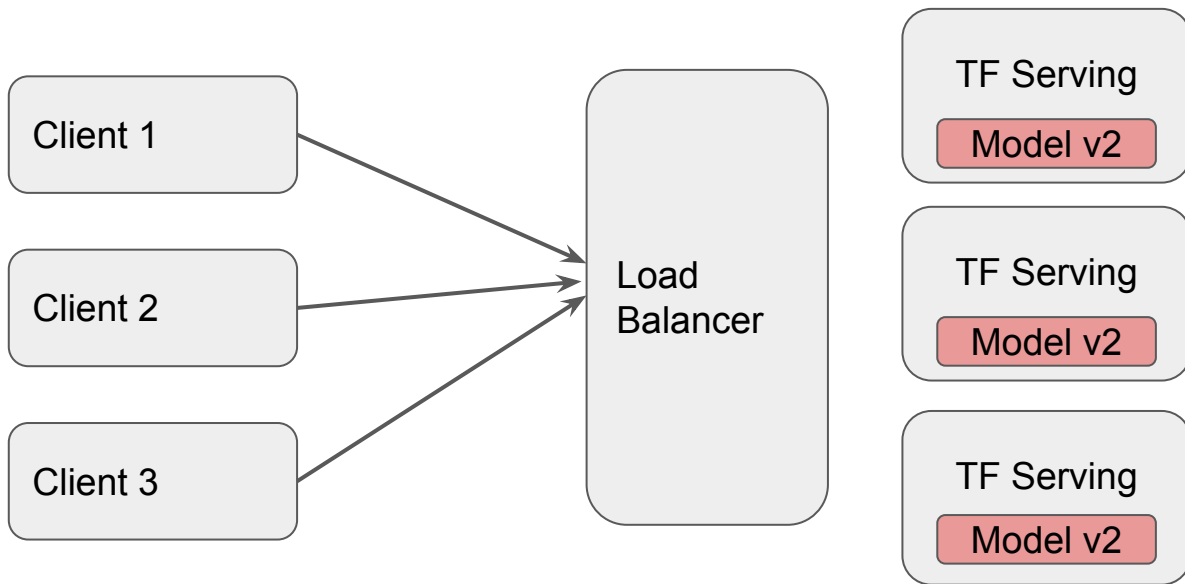
Client 3

Model v2

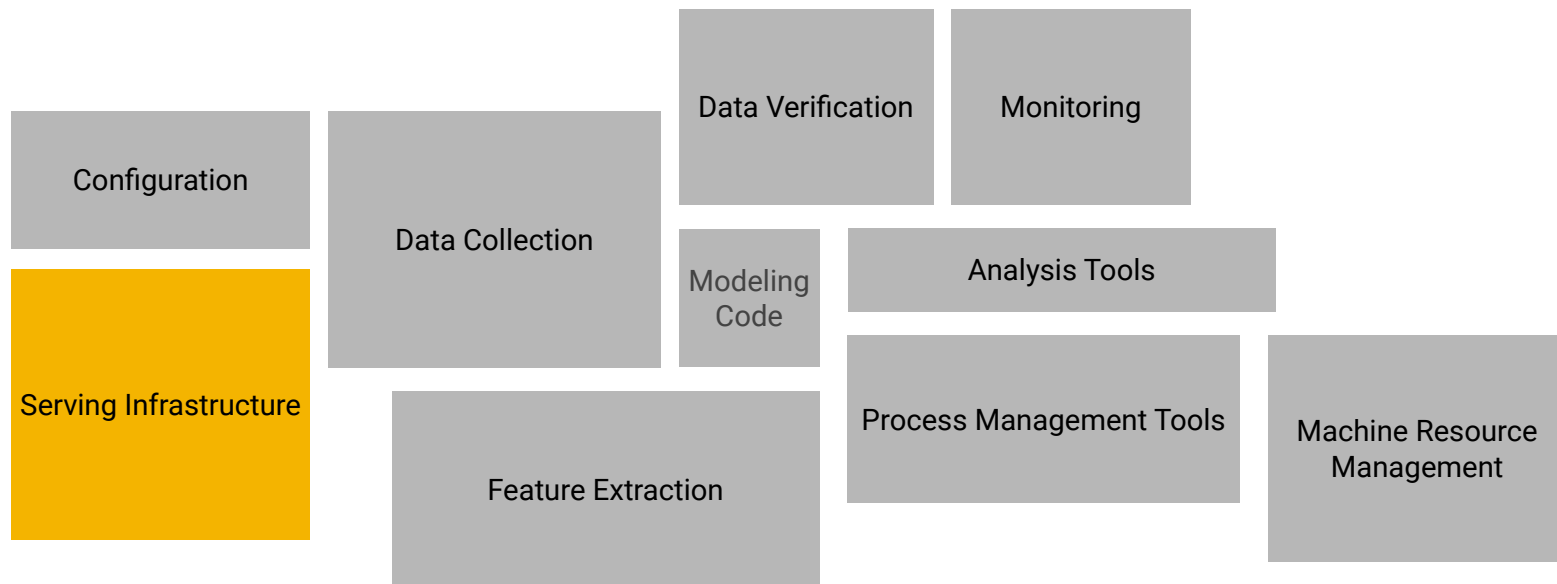








TensorFlow Serving is part of TFX





Install TensorFlow Serving...

- Docker
- APT
- Build From Source
- PIP Packages

<https://www.tensorflow.org/tfx/serving/setup>

```
!echo "deb http://storage.googleapis.com/tensorflow-serving-apt stable tensorflow-model-server
tensorflow-model-server-universal"
    | tee /etc/apt/sources.list.d/tensorflow-serving.list && \

curl https://storage.googleapis.com/tensorflow-serving-apt/tensorflow-serving.release.pub.gpg
    | apt-key add -

!apt update

!apt-get install tensorflow-model-server
```

```
import tensorflow as tf
import numpy as np
from tensorflow import keras

model = tf.keras.Sequential([keras.layers.Dense(units=1, input_shape=[1])])
model.compile(optimizer='sgd', loss='mean_squared_error')

xs = np.array([-1.0, 0.0, 1.0, 2.0, 3.0, 4.0], dtype=float)
ys = np.array([-3.0, -1.0, 1.0, 3.0, 5.0, 7.0], dtype=float)

model.fit(xs, ys, epochs=500, verbose=2)

print(model.predict([10.0]))
```

```
tf.saved_model.simple_save(  
    keras.backend.get_session(),  
    export_path,  
    inputs={'input_image': model.input},  
    outputs={t.name:t for t in model.outputs})
```

```
tf.saved_model.simple_save(  
    keras.backend.get_session(),  
    export_path,  
    inputs={'input_image': model.input},  
    outputs={t.name:t for t in model.outputs})
```

```
tf.saved_model.simple_save(  
    keras.backend.get_session(),  
    export_path,  
    inputs={'input_image': model.input},  
    outputs={t.name:t for t in model.outputs})
```



```
tf.saved_model.simple_save(  
    keras.backend.get_session(),  
    export_path,  
    inputs={'input_image': model.input},  
    outputs={t.name:t for t in model.outputs})
```

```
tf.saved_model.simple_save(  
    keras.backend.get_session(),  
    export_path,  
    inputs={'input_image': model.input},  
    outputs={t.name:t for t in model.outputs})
```

```
# Fetch the Keras session and save the model
# The signature definition is defined by the input and output tensors,
# and stored with the default serving key
import tempfile

MODEL_DIR = tempfile.gettempdir()
version = 1
export_path = os.path.join(MODEL_DIR, str(version))
print('export_path = {}'.format(export_path))
if os.path.isdir(export_path):
    print('\nAlready saved a model, cleaning up\n')
    !rm -r {export_path}

tf.saved_model.simple_save(
    keras.backend.get_session(),
    export_path,
    inputs={'input_image': model.input},
    outputs={t.name:t for t in model.outputs})

print('\nSaved model:')
!ls -l {export_path}
```

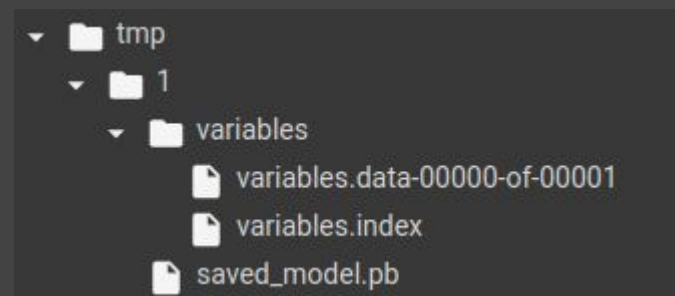
```
export_path = /tmp/1
```

Saved model:

total 44

```
-rw-r--r-- 1 root root 39532 Nov  3 23:35 saved_model.pb
```

```
drwxr-xr-x 2 root root  4096 Nov  3 23:35 variables
```



```
!saved_model_cli show --dir {export_path} --all
```

MetaGraphDef with tag-set: 'serve' contains the following SignatureDefs:

signature_def['serving_default']:

The given SavedModel SignatureDef contains the following input(s):

inputs['input_image'] tensor_info:

dtype: DT_FLOAT

shape: (-1, 1)

name: dense_input:0

The given SavedModel SignatureDef contains the following output(s):

outputs['dense/BiasAdd:0'] tensor_info:

dtype: DT_FLOAT

shape: (-1, 1)

name: dense/BiasAdd:0

Method name is: tensorflow/serving/predict

MetaGraphDef with tag-set: 'serve' contains the following SignatureDefs:

signature_def['serving_default']:

The given SavedModel SignatureDef contains the following input(s):

inputs['input_image'] tensor_info:

dtype: DT_FLOAT

shape: (-1, 1)

name: dense_input:0

The given SavedModel SignatureDef contains the following output(s):

outputs['dense/BiasAdd:0'] tensor_info:

dtype: DT_FLOAT

shape: (-1, 1)

name: dense/BiasAdd:0

Method name is: tensorflow/serving/predict

MetaGraphDef with tag-set: 'serve' contains the following SignatureDefs:

signature_def['serving_default']:

The given SavedModel SignatureDef contains the following input(s):

inputs['input_image'] tensor_info:

dtype: DT_FLOAT

shape: (-1, 1)

name: dense_input:0

The given SavedModel SignatureDef contains the following output(s):

outputs['dense/BiasAdd:0'] tensor_info:

dtype: DT_FLOAT

shape: (-1, 1)

name: dense/BiasAdd:0

Method name is: tensorflow/serving/predict

MetaGraphDef with tag-set: 'serve' contains the following SignatureDefs:

signature_def['serving_default']:

The given SavedModel SignatureDef contains the following input(s):

inputs['input_image'] tensor_info:

dtype: DT_FLOAT

shape: (-1, 1)

name: dense_input:0

The given SavedModel SignatureDef contains the following output(s):

outputs['dense/BiasAdd:0'] tensor_info:

dtype: DT_FLOAT

shape: (-1, 1)

name: dense/BiasAdd:0

Method name is: tensorflow/serving/predict

```
os.environ["MODEL_DIR"] = MODEL_DIR
```

```
%%bash --bg
```

```
nohup tensorflow_model_server \
```

```
--rest_api_port=8501 \
```

```
--model_name=helloworld \
```

```
--model_base_path="${MODEL_DIR}" >server.log 2>&1
```

```
os.environ["MODEL_DIR"] = MODEL_DIR
```

```
%%bash --bg  
nohup tensorflow_model_server \  
  --rest_api_port=8501 \  
  --model_name=helloworld \  
  --model_base_path="${MODEL_DIR}" >server.log 2>&1
```

```
os.environ["MODEL_DIR"] = MODEL_DIR
```

```
%%bash --bg  
nohup tensorflow_model_server \  
  --rest_api_port=8501 \  
  --model_name=helloworld \  
  --model_base_path="${MODEL_DIR}" >server.log 2>&1
```

```
os.environ["MODEL_DIR"] = MODEL_DIR
```

```
%%bash --bg
```

```
nohup tensorflow_model_server \
```

```
--rest_api_port=8501 \
```

```
--model_name=helloworld \
```

```
--model_base_path="${MODEL_DIR}" >server.log 2>&1
```

```
os.environ["MODEL_DIR"] = MODEL_DIR
```

```
%%bash --bg
```

```
nohup tensorflow_model_server \
```

```
--rest_api_port=8501 \
```

```
--model_name=helloworld \
```

```
--model_base_path="${MODEL_DIR}" >server.log 2>&1
```

```
os.environ["MODEL_DIR"] = MODEL_DIR
```

```
%%bash --bg
```

```
nohup tensorflow_model_server \
```

```
--rest_api_port=8501 \
```

```
--model_name=helloworld \
```

```
--model_base_path="${MODEL_DIR}" >server.log 2>&1
```

```
os.environ["MODEL_DIR"] = MODEL_DIR
```

```
%%bash --bg
```

```
nohup tensorflow_model_server \
```

```
--rest_api_port=8501 \
```

```
--model_name=helloworld \
```

```
--model_base_path="${MODEL_DIR}" >server.log 2>&1
```

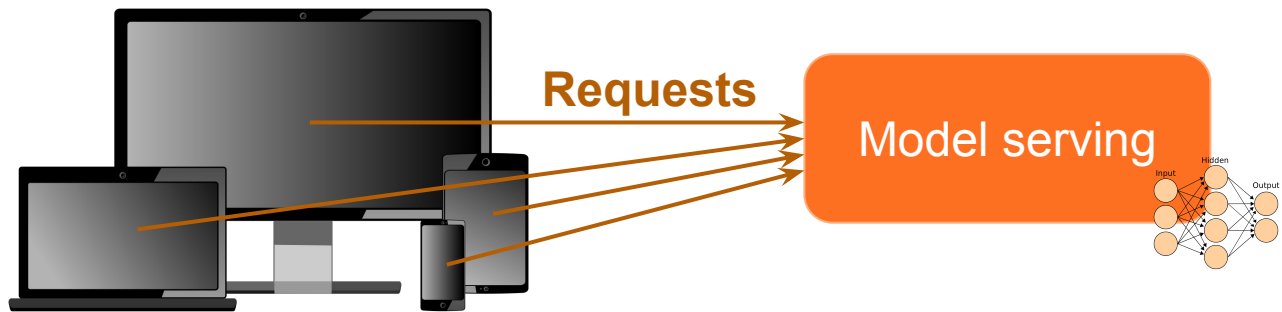


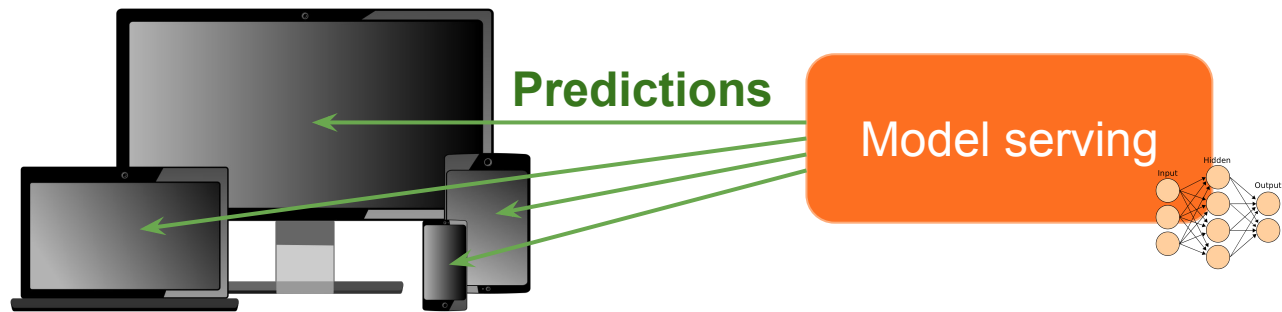
```
!tail server.log
```

```
2019-11-03 23:52:16.178660: I external/org_tensorflow/tensorflow/cc/saved_model/reader.cc:54] Reading meta graph with
tags { serve }
2019-11-03 23:52:16.179306: I external/org_tensorflow/tensorflow/core/platform/cpu_feature_guard.cc:142] Your CPU
supports instructions that this TensorFlow binary was not compiled to use: AVX2 FMA
2019-11-03 23:52:16.192356: I external/org_tensorflow/tensorflow/cc/saved_model/loader.cc:202] Restoring SavedModel
bundle.
2019-11-03 23:52:16.198559: I external/org_tensorflow/tensorflow/cc/saved_model/loader.cc:311] SavedModel load for
tags { serve }; Status: success. Took 20743 microseconds.
2019-11-03 23:52:16.198779: I tensorflow_serving/servables/tensorflow/saved_model_warmup.cc:105] No warmup data file
found at /tmp/1/assets.extra/tf_serving_warmup_requests
2019-11-03 23:52:16.198852: I tensorflow_serving/core/loader_harness.cc:87] Successfully loaded servable version
{name: helloworld version: 1}
2019-11-03 23:52:16.199907: I tensorflow_serving/model_servers/server.cc:353] Running gRPC ModelServer at
0.0.0.0:8500 ...
[warn] getaddrinfo: address family for nodename not supported
2019-11-03 23:52:16.200544: I tensorflow_serving/model_servers/server.cc:373] Exporting HTTP/REST API
at:localhost:8501 ...
[evhttp_server.cc : 238] NET_LOG: Entering the event loop ...
```

```
!tail server.log
```

```
2019-11-03 23:52:16.178660: I external/org_tensorflow/tensorflow/cc/saved_model/reader.cc:54] Reading meta graph with
tags { serve }
2019-11-03 23:52:16.179306: I external/org_tensorflow/tensorflow/core/platform/cpu_feature_guard.cc:142] Your CPU
supports instructions that this TensorFlow binary was not compiled to use: AVX2 FMA
2019-11-03 23:52:16.192356: I external/org_tensorflow/tensorflow/cc/saved_model/loader.cc:202] Restoring SavedModel
bundle.
2019-11-03 23:52:16.198559: I external/org_tensorflow/tensorflow/cc/saved_model/loader.cc:311] SavedModel load for
tags { serve }; Status: success. Took 20743 microseconds.
2019-11-03 23:52:16.198779: I tensorflow_serving/servables/tensorflow/saved_model_warmup.cc:105] No warmup data file
found at /tmp/1/assets.extra/tf_serving_warmup_requests
2019-11-03 23:52:16.198852: I tensorflow_serving/core/loader_harness.cc:87] Successfully loaded servable version
{name: helloworld version: 1}
2019-11-03 23:52:16.199907: I tensorflow_serving/model_servers/server.cc:353] Running gRPC ModelServer at
0.0.0.0:8500 ...
[warn] getaddrinfo: address family for nodename not supported
2019-11-03 23:52:16.200544: I tensorflow_serving/model_servers/server.cc:373] Exporting HTTP/REST API
at:localhost:8501 ...
[evhttp_server.cc : 238] NET_LOG: Entering the event loop ...
```





```
import json
xs = np.array([[9.0], [10.0]])
data = json.dumps({"signature_name": "serving_default", "instances": xs.tolist()})
print(data)
```

```
import json
xs = np.array([[9.0], [10.0]])
data = json.dumps({"signature_name": "serving_default", "instances": xs.tolist()})
print(data)
```

```
import json
xs = np.array([[9.0], [10.0]])
data = json.dumps({"signature_name": "serving_default", "instances": xs.tolist()})
print(data)
```

```
import json
xs = np.array([[9.0], [10.0]])
data = json.dumps({"signature_name": "serving_default", "instances": xs.tolist()})
print(data)
```



```
import json
xs = np.array([[9.0], [10.0]])
data = json.dumps({"signature_name": "serving_default", "instances": xs.tolist()})
print(data)
```

```
import json
xs = np.array([[9.0], [10.0]])
data = json.dumps({"signature_name": "serving_default", "instances": xs.tolist()})
print(data)
```

MetaGraphDef with tag-set: 'serve' contains the following SignatureDefs:

```
signature_def['serving_default']:
```

The given SavedModel SignatureDef contains the following input(s):

inputs['input_image'] tensor_info:

dtype: DT_FLOAT

shape: (-1, 1)

name: dense_input:0

The given SavedModel SignatureDef contains the following output(s):

outputs['dense/BiasAdd:0'] tensor_info:

dtype: DT_FLOAT

shape: (-1, 1)

name: dense/BiasAdd:0

Method name is: tensorflow/serving/predict

```
import json
xs = np.array([[9.0], [10.0]])
data = json.dumps({"signature_name": "serving_default", "instances": xs.tolist()})
print(data)
```

```
import json
xs = np.array([[9.0], [10.0]])
data = json.dumps({"signature_name": "serving_default", "instances": xs.tolist()})
print(data)
```

```
{"signature_name": "serving_default", "instances": [[9.0], [10.0]]}
```

```
!pip install -q requests
```

```
import requests
headers = {"content-type": "application/json"}
json_response = requests.post('http://localhost:8501/v1/models/hello-world:predict',
data=data, headers=headers)

print(json_response.text)
predictions = json.loads(json_response.text)['predictions']
```

```
!pip install -q requests
```

```
import requests
```

```
headers = {"content-type": "application/json"}
```

```
json_response = requests.post('http://localhost:8501/v1/models/helloworld:predict',  
                               data=data,  
                               headers=headers)
```

```
print(json_response.text)
```

```
predictions = json.loads(json_response.text)['predictions']
```

```
!pip install -q requests
```

```
import requests
```

```
headers = {"content-type": "application/json"}
```

```
json_response = requests.post('http://localhost:8501/v1/models/helloworld:predict',  
                               data=data,  
                               headers=headers)
```

```
print(json_response.text)
```

```
predictions = json.loads(json_response.text)['predictions']
```

```
!pip install -q requests
```

```
import requests
```

```
headers = {"content-type": "application/json"}
```

```
json_response = requests.post('http://localhost:8501/v1/models/helloworld:predict',  
                               data=data,  
                               headers=headers)
```

```
print(json_response.text)
```

```
predictions = json.loads(json_response.text)['predictions']
```



```
!pip install -q requests
```

```
import requests
```

```
headers = {"content-type": "application/json"}
```

```
json_response = requests.post('http://localhost:8501/v1/models/helloworld:predict',  
                               data=data,  
                               headers=headers)
```

```
print(json_response.text)
```

```
predictions = json.loads(json_response.text)['predictions']
```

```
!pip install -q requests
```

```
import requests
```

```
headers = {"content-type": "application/json"}
```

```
json_response = requests.post('http://localhost:8501/v1/models/helloworld:predict',  
                               data=data,  
                               headers=headers)
```

```
print(json_response.text)
```

```
predictions = json.loads(json_response.text)['predictions']
```

```
!pip install -q requests
```

```
import requests
```

```
headers = {"content-type": "application/json"}
```

```
json_response = requests.post('http://localhost:8501/v1/models/helloworld:predict',  
                               data=data,  
                               headers=headers)
```

```
print(json_response.text)
```

```
predictions = json.loads(json_response.text)['predictions']
```

```
{  
  "predictions": [[16.9865685], [18.984293]]  
}
```

<http://bit.ly/tfserving-lab1>

Train and serve a TensorFlow model with TensorFlow Serving

Warning: This notebook is designed to be run in a Google Colab only. It installs packages on the system and requires root access. If you want to run it in a local Jupyter notebook, please proceed with caution.

★ **Note:** You can run this example right now in a Jupyter-style notebook, no setup required! Just click "Run in Google Colab"

 Run in Google Colab

 View source on GitHub

This guide trains a neural network model to classify [images of clothing, like sneakers and shirts](#), saves the trained model, and then serves it with [TensorFlow Serving](#). The focus is on TensorFlow Serving, rather than the modeling and training in TensorFlow, so for a complete example which focuses on the modeling and training see the [Basic Classification example](#).

This guide uses [tf.keras](#), a high-level API to build and train models in TensorFlow.

```
# TensorFlow and tf.keras
import tensorflow as tf
from tensorflow import keras

# Helper libraries
import numpy as np
import matplotlib.pyplot as plt
import os
import subprocess

print(tf.__version__)
```

1.12.0

```
import json

data = json.dumps({"signature_name": "serving_default",
                  "instances": test_images[0:3].tolist()})

print('Data: {} ... {}'.format(data[:50], data[len(data)-52:]))

print(data)
```



```
[[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]],
```

```
[[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]],
```

```
[[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]]
```



```
[[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]],
```

```
[[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]],
```

```
[[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]]
```

```
[[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]],
```

```
[[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]],
```

```
[[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]],  
[[0.0], [0.0], [0.0], ... ]]
```

```
!pip install -q requests
```

```
import requests
```

```
headers = {"content-type": "application/json"}
```

```
json_response = requests.post('http://localhost:8501/v1/models/fashion_model:predict',  
                               data=data, headers=headers)
```

```
predictions = json.loads(json_response.text)['predictions']
```

```
[  
[5.77123615e-07, 2.66907847e-08, 4.7217938e-08, 1.97792871e-09, 5.31984341e-08,  
0.00734644197, 3.1462946e-07, 0.0439051725, 0.000500570168, 0.948246837],  
  
[0.00227244, 6.12080342e-09, 0.967876315, 3.0579281e-06, 0.0183339939, 3.18483538e-11,  
0.011510049, 1.38639566e-14, 4.19033222e-06, 4.40264526e-11],  
  
[1.45221502e-05, 0.999841571, 3.96758715e-08, 0.000131023204, 1.22008023e-05,  
1.18227668e-08, 5.97860179e-08, 1.31281848e-08, 5.49047854e-07, 2.97885189e-10]  
]
```

```
[  
[5.77123615e-07, 2.66907847e-08, 4.7217938e-08, 1.97792871e-09, 5.31984341e-08,  
0.00734644197, 3.1462946e-07, 0.0439051725, 0.000500570168, 0.948246837],  
  
[0.00227244, 6.12080342e-09, 0.967876315, 3.0579281e-06, 0.0183339939, 3.18483538e-11,  
0.011510049, 1.38639566e-14, 4.19033222e-06, 4.40264526e-11],  
  
[1.45221502e-05, 0.999841571, 3.96758715e-08, 0.000131023204, 1.22008023e-05,  
1.18227668e-08, 5.97860179e-08, 1.31281848e-08, 5.49047854e-07, 2.97885189e-10]  
]
```

```
show(0,  
'The model saw {} (class {}), and it was actually a {} (class {})' .format(  
    class_names[np.argmax(predictions[0])], test_labels[0],  
    class_names[np.argmax(predictions[0])], test_labels[0]))
```

The model thought this was a Ankle boot (class 9), and it was actually a Ankle boot (class 9)

