# ▾ Ungraded Lab: Developing Custom TFX Components

[Tensorflow Extended (TFX)](#) provides ready made components for typical steps in a machine learning workflow. Other courses in this specialization focus on specific components and in this lab, you will learn how to make your own. This will be useful in case your project has specific needs that fall outside the standard TFX components. It will make your pipelines more flexible while still leveraging the experiment tracking and orchestration that TFX provides. In particular, you will:

- build a custom component using Python functions
- build a custom component by reusing an existing TFX component
- run a TFX pipeline locally using a built-in pipeline orchestrator

To demonstrate, you will run the pipeline used in this [official tutorial](#) then modify it to have a custom component. Some of the discussions here are also taken from that tutorial to explain the motivation and point to additional resources.

Let's begin!

*Note: If you haven't taken other courses in this specialization and it's the first time you're using TFX, please see [Understanding TFX Pipelines](#) to get an overview of important concepts.*

## ▾ Set Up

We first need to install the TFX Python package and download the dataset which we will use for our model.

### Upgrade Pip

To avoid upgrading Pip in a system when running locally, check to make sure that we are running in Colab. Local systems can of course be upgraded separately.

```
1   try:
2     import colab
3     !pip install --upgrade pip
4   except:
5     pass
```

## ▾ Install TFX

```
1   !pip install -U tfx
```

**Note:** *Please do not proceed to the next steps without restarting the Runtime after installing* `tfx`*. You can do that by either pressing the* `Restart Runtime` *button at the end of the cell output above, or going to the* `Runtime` *button at the Colab toolbar above and selecting* `Restart Runtime`*.*

Check the TensorFlow and TFX versions.

```
1   import tensorflow as tf
2   print('TensorFlow version: {}'.format(tf.__version__))
3   from tfx import v1 as tfx
4   print('TFX version: {}'.format(tfx.__version__))
```

```
TensorFlow version: 2.5.1
TFX version: 1.2.0
```

## ▾ Set up variables

There are some variables used to define a pipeline. You can customize these variables as you want. By default all output from the pipeline will be generated under the current directory.

```
1   import os
2   from absl import logging
3
4   # Pipeline label
5   PIPELINE_NAME = "penguin-simple"
6
7   # Output directory to store artifacts generated from the pipeline.
8   PIPELINE_ROOT = os.path.join('pipelines', PIPELINE_NAME)
9
10  # Path to a SQLite DB file to use as an MLMD storage.
11  METADATA_PATH = os.path.join('metadata', PIPELINE_NAME, 'metadata.db')
12
13  # Output directory where created models from the pipeline will be exported.
14  SERVING_MODEL_DIR = os.path.join('serving_model', PIPELINE_NAME)
15
16  # Set default logging level.
17  logging.set_verbosity(logging.INFO)
```

## ▾ Prepare example data

The dataset you will use is the [Palmer Penguins dataset](#).

There are four numeric features in this dataset:

- culmen_length_mm
- culmen_depth_mm
- flipper_length_mm

- flipper_length_mm
- body_mass_g

All features were already normalized to have range [0,1]. You will build a classification model which predicts the `species` of penguins. The code below creates a directory and copies the dataset to it. After running it, you should see the dataset in the Colab file explorer under the `data` folder.

```
1   import urllib.request
2
3   # Create directory
4   DATA_ROOT = 'data'
5   !mkdir {DATA_ROOT}
6
7   # Copy dataset to directory
8   _data_url = 'https://raw.githubusercontent.com/tensorflow/tfx/master/tfx/examples/penguin/data/labelled/penguins_processed.csv'
9   _data_filepath = os.path.join(DATA_ROOT, "data.csv")
10  urllib.request.urlretrieve(_data_url, _data_filepath)
```

```
('data/data.csv', <http.client.HTTPMessage at 0x7f637304e510>)
```

## Running the pipeline using standard components

The pipeline will consist of three essential TFX components and the graph will look like this:

```
ExampleGen -> Trainer -> Pusher
```

The pipeline includes the most minimal ML workflow which is importing data (ExampleGen), training a model (Trainer) and exporting the trained model (Pusher).

## Model training code

You will first define the trainer module so the `Trainer` component can build the model and train it.

```
1   _trainer_module_file = 'penguin_trainer.py'
```

```
1   %%writefile {_trainer_module_file}
2
3   from typing import List
4   from absl import logging
5   import tensorflow as tf
6   from tensorflow import keras
7   from tensorflow_transform.tf_metadata import schema_utils
8
9   from tfx import v1 as tfx
10  from tfx_bsl.public import tfxio
11  from tensorflow_metadata.proto.v0 import schema_pb2
12
13  _FEATURE_KEYS = [
14      'culmen_length_mm', 'culmen_depth_mm', 'flipper_length_mm', 'body_mass_g'
15  ]
16  _LABEL_KEY = 'species'
17
18  _TRAIN_BATCH_SIZE = 20
19  _EVAL_BATCH_SIZE = 10
```

```
    def _build_keras_model() -> tf.keras.Model:
57
58      """Creates a DNN Keras model for classifying penguin data.
59
60      Returns:
61        A Keras Model.
62      """
63      # The model below is built with Functional API, please refer to
64      # https://www.tensorflow.org/guide/keras/overview for all API options.
65      inputs = [keras.layers.Input(shape=(1,), name=f) for f in _FEATURE_KEYS]
66      d = keras.layers.concatenate(inputs)
67      for _ in range(2):
68        d = keras.layers.Dense(8, activation='relu')(d)
69      outputs = keras.layers.Dense(3)(d)
70
71      model = keras.Model(inputs=inputs, outputs=outputs)
72      model.compile(
73          optimizer=keras.optimizers.Adam(1e-2),
74          loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True),
75          metrics=[keras.metrics.SparseCategoricalAccuracy()])
76
77      model.summary(print_fn=logging.info)
78      return model
79
```