

Keras

Keras is a high-level API to build and train deep learning models. It's used for fast prototyping, advanced research, and production, with three key advantages:

- *User friendly*
Keras has a simple, consistent interface optimized for common use cases. It provides clear and actionable feedback for user errors.
- *Modular and composable*
Keras models are made by connecting configurable building blocks together, with few restrictions.
- *Easy to extend*
Write custom building blocks to express new ideas for research. Create new layers, loss functions, and develop state-of-the-art models.

Import tf.keras

tf.keras (https://www.tensorflow.org/api_docs/python/tf/keras) is TensorFlow's implementation of the Keras API specification (<https://keras.io>). This is a high-level API to build and train models that includes first-class support for TensorFlow-specific functionality, such as eager execution (`#eager_execution`), **tf.data** (https://www.tensorflow.org/api_docs/python/tf/data) pipelines, and Estimators (<https://www.tensorflow.org/guide/estimators>). **tf.keras** (https://www.tensorflow.org/api_docs/python/tf/keras) makes TensorFlow easier to use without sacrificing flexibility and performance.

To get started, import **tf.keras** (https://www.tensorflow.org/api_docs/python/tf/keras) as part of your TensorFlow program setup:

```
import tensorflow as tf
from tensorflow import keras
```



tf.keras (https://www.tensorflow.org/api_docs/python/tf/keras) can run any Keras-compatible code, but keep in mind:

- The **tf.keras** (https://www.tensorflow.org/api_docs/python/tf/keras) version in the latest TensorFlow release might not be the same as the latest **keras** version from PyPI.

Check `tf.keras.version`

(https://www.tensorflow.org/api_docs/python/tf/keras/__version__).

- When saving a model's weights (#weights_only), `tf.keras` (https://www.tensorflow.org/api_docs/python/tf/keras) defaults to the checkpoint format (<https://www.tensorflow.org/guide/checkpoints>). Pass `save_format='h5'` to use HDF5.

Build a simple model

Sequential model

In Keras, you assemble *layers* to build *models*. A model is (usually) a graph of layers. The most common type of model is a stack of layers: the `tf.keras.Sequential` (https://www.tensorflow.org/api_docs/python/tf/keras/Sequential) model.

To build a simple, fully-connected network (i.e. multi-layer perceptron):

```
model = keras.Sequential()  
# Adds a densely-connected layer with 64 units to the model:  
model.add(keras.layers.Dense(64, activation='relu'))  
# Add another:  
model.add(keras.layers.Dense(64, activation='relu'))  
# Add a softmax layer with 10 output units:  
model.add(keras.layers.Dense(10, activation='softmax'))
```



Configure the layers

There are many `tf.keras.layers` (https://www.tensorflow.org/api_docs/python/tf/keras/layers) available with some common constructor parameters:

- **activation**: Set the activation function for the layer. This parameter is specified by the name of a built-in function or as a callable object. By default, no activation is applied.
- **kernel_initializer** and **bias_initializer**: The initialization schemes that create the layer's weights (kernel and bias). This parameter is a name or a callable object. This defaults to the "Glorot uniform" initializer.
- **kernel_regularizer** and **bias_regularizer**: The regularization schemes that apply the layer's weights (kernel and bias), such as L1 or L2 regularization. By default, no regularization is applied.

The following instantiates `tf.keras.layers.Dense`

(https://www.tensorflow.org/api_docs/python/tf/keras/layers/Dense) layers using constructor arguments:



```
# Create a sigmoid layer:
layers.Dense(64, activation='sigmoid')
# Or:
layers.Dense(64, activation=tf.sigmoid)

# A linear layer with L1 regularization of factor 0.01 applied to the kernel
layers.Dense(64, kernel_regularizer=keras.regularizers.l1(0.01))
# A linear layer with L2 regularization of factor 0.01 applied to the bias vector
layers.Dense(64, bias_regularizer=keras.regularizers.l2(0.01))

# A linear layer with a kernel initialized to a random orthogonal matrix:
layers.Dense(64, kernel_initializer='orthogonal')
# A linear layer with a bias vector initialized to 2.0s:
layers.Dense(64, bias_initializer=keras.initializers.constant(2.0))
```

Train and evaluate

Set up training

After the model is constructed, configure its learning process by calling the `compile` method:



```
model.compile(optimizer=tf.train.AdamOptimizer(0.001),
              loss='categorical_crossentropy',
              metrics=['accuracy'])
```

`tf.keras.Model.compile` (https://www.tensorflow.org/api_docs/python/tf/keras/Model#compile) takes three important arguments:

- **optimizer**: This object specifies the training procedure. Pass it optimizer instances from the `tf.train` (https://www.tensorflow.org/api_docs/python/tf/train) module, such as `AdamOptimizer` (https://www.tensorflow.org/api_docs/python/tf/train/AdamOptimizer), `RMSPropOptimizer` (https://www.tensorflow.org/api_docs/python/tf/train/RMSPropOptimizer), or `GradientDescentOptimizer` (https://www.tensorflow.org/api_docs/python/tf/train/GradientDescentOptimizer).
- **loss**: The function to minimize during optimization. Common choices include mean square error (mse), categorical_crossentropy, and binary_crossentropy. Loss

functions are specified by name or by passing a callable object from the `tf.keras.losses` (https://www.tensorflow.org/api_docs/python/tf/keras/losses) module.

- **metrics:** Used to monitor training. These are string names or callables from the `tf.keras.metrics` (https://www.tensorflow.org/api_docs/python/tf/keras/metrics) module.

The following shows a few examples of configuring a model for training:

```
# Configure a model for mean-squared error regression.
model.compile(optimizer=tf.train.AdamOptimizer(0.01),
              loss='mse',          # mean squared error
              metrics=['mae'])    # mean absolute error

# Configure a model for categorical classification.
model.compile(optimizer=tf.train.RMSPropOptimizer(0.01),
              loss=keras.losses.categorical_crossentropy,
              metrics=[keras.metrics.categorical_accuracy])
```



Input NumPy data

For small datasets, use in-memory NumPy (<https://www.numpy.org/>) arrays to train and evaluate a model. The model is "fit" to the training data using the `fit` method:

```
import numpy as np

data = np.random.random((1000, 32))
labels = np.random.random((1000, 10))

model.fit(data, labels, epochs=10, batch_size=32)
```



`tf.keras.Model.fit` (https://www.tensorflow.org/api_docs/python/tf/keras/Model#fit) takes three important arguments:

- **epochs:** Training is structured into *epochs*. An epoch is one iteration over the entire input data (this is done in smaller batches).
- **batch_size:** When passed NumPy data, the model slices the data into smaller batches and iterates over these batches during training. This integer specifies the size of each batch. Be aware that the last batch may be smaller if the total number of samples is not divisible by the batch size.
- **validation_data:** When prototyping a model, you want to easily monitor its performance on some validation data. Passing this argument—a tuple of inputs and labels—allows the model to display the loss and metrics in inference mode for the passed data, at the end of each epoch.

Here's an example using `validation_data`:

```
import numpy as np

data = np.random.random((1000, 32))
labels = np.random.random((1000, 10))

val_data = np.random.random((100, 32))
val_labels = np.random.random((100, 10))

model.fit(data, labels, epochs=10, batch_size=32,
          validation_data=(val_data, val_labels))
```



Input tf.data datasets

Use the [Datasets API](https://www.tensorflow.org/guide/datasets) (<https://www.tensorflow.org/guide/datasets>) to scale to large datasets or multi-device training. Pass a `tf.data.Dataset` (https://www.tensorflow.org/api_docs/python/tf/data/Dataset) instance to the `fit` method:

```
# Instantiates a toy dataset instance:
dataset = tf.data.Dataset.from_tensor_slices((data, labels))
dataset = dataset.batch(32)
dataset = dataset.repeat()

# Don't forget to specify `steps_per_epoch` when calling `fit` on a dataset.
model.fit(dataset, epochs=10, steps_per_epoch=30)
```



Here, the `fit` method uses the steps_per_epoch argument—this is the number of training steps the model runs before it moves to the next epoch. Since the `Dataset` yields batches of data, this snippet does not require a `batch_size`.

Datasets can also be used for validation:

```
dataset = tf.data.Dataset.from_tensor_slices((data, labels))
dataset = dataset.batch(32).repeat()

val_dataset = tf.data.Dataset.from_tensor_slices((val_data, val_labels))
val_dataset = val_dataset.batch(32).repeat()

model.fit(dataset, epochs=10, steps_per_epoch=30,
          validation_data=val_dataset,
          validation_steps=3)
```



Evaluate and predict

The `tf.keras.Model.evaluate`

(https://www.tensorflow.org/api_docs/python/tf/keras/Model#evaluate) and

`tf.keras.Model.predict` (https://www.tensorflow.org/api_docs/python/tf/keras/Model#predict)

methods can use NumPy data and a `tf.data.Dataset`

(https://www.tensorflow.org/api_docs/python/tf/data/Dataset).

To *evaluate* the inference-mode loss and metrics for the data provided:

```
model.evaluate(x, y, batch_size=32)
```



```
model.evaluate(dataset, steps=30)
```

And to *predict* the output of the last layer in inference for the data provided, as a NumPy array:

```
model.predict(x, batch_size=32)
```



```
model.predict(dataset, steps=30)
```

Build advanced models

Functional API

The `tf.keras.Sequential` (https://www.tensorflow.org/api_docs/python/tf/keras/Sequential) model is a simple stack of layers that cannot represent arbitrary models. Use the Keras functional API (<https://keras.io/getting-started/functional-api-guide/>) to build complex model topologies such as:

- Multi-input models,
- Multi-output models,
- Models with shared layers (the same layer called several times),
- Models with non-sequential data flows (e.g. residual connections).

Building a model with the functional API works like this:

1. A layer instance is callable and returns a tensor.
2. Input tensors and output tensors are used to define a `tf.keras.Model` (https://www.tensorflow.org/api_docs/python/tf/keras/Model) instance.

3. This model is trained just like the `Sequential` model.

The following example uses the functional API to build a simple, fully-connected network:

```
inputs = keras.Input(shape=(32,)) # Returns a placeholder tensor

# A layer instance is callable on a tensor, and returns a tensor.
x = keras.layers.Dense(64, activation='relu')(inputs)
x = keras.layers.Dense(64, activation='relu')(x)
predictions = keras.layers.Dense(10, activation='softmax')(x)

# Instantiate the model given inputs and outputs.
model = keras.Model(inputs=inputs, outputs=predictions)

# The compile step specifies the training configuration.
model.compile(optimizer=tf.train.RMSPropOptimizer(0.001),
              loss='categorical_crossentropy',
              metrics=['accuracy'])

# Trains for 5 epochs
model.fit(data, labels, batch_size=32, epochs=5)
```

Model subclassing

Build a fully-customizable model by subclassing `tf.keras.Model`

(https://www.tensorflow.org/api_docs/python/tf/keras/Model) and defining your own forward pass. Create layers in the `__init__` method and set them as attributes of the class instance. Define the forward pass in the `call` method.

Model subclassing is particularly useful when eager execution

(<https://www.tensorflow.org/guide/eager>) is enabled since the forward pass can be written imperatively.

Key Point: Use the right API for the job. While model subclassing offers flexibility, it comes at a cost of greater complexity and more opportunities for user errors. If possible, prefer the functional API.

The following example shows a subclassed `tf.keras.Model`

(https://www.tensorflow.org/api_docs/python/tf/keras/Model) using a custom forward pass:

```
class MyModel(keras.Model):

    def __init__(self, num_classes=10):
        super(MyModel, self).__init__(name='my_model')
        self.num_classes = num_classes
```

```

    # Define your layers here.
    self.dense_1 = keras.layers.Dense(32, activation='relu')
    self.dense_2 = keras.layers.Dense(num_classes, activation='sigmoid')

def call(self, inputs):
    # Define your forward pass here,
    # using layers you previously defined (in `__init__`).
    x = self.dense_1(inputs)
    return self.dense_2(x)

def compute_output_shape(self, input_shape):
    # You need to override this function if you want to use the subclassed model
    # as part of a functional-style model.
    # Otherwise, this method is optional.
    shape = tf.TensorShape(input_shape).as_list()
    shape[-1] = self.num_classes
    return tf.TensorShape(shape)

# Instantiates the subclassed model.
model = MyModel(num_classes=10)

# The compile step specifies the training configuration.
model.compile(optimizer=tf.train.RMSPropOptimizer(0.001),
              loss='categorical_crossentropy',
              metrics=['accuracy'])

# Trains for 5 epochs.
model.fit(data, labels, batch_size=32, epochs=5)

```

Custom layers

Create a custom layer by subclassing [tf.keras.layers.Layer](https://www.tensorflow.org/api_docs/python/tf/keras/layers/Layer)

(https://www.tensorflow.org/api_docs/python/tf/keras/layers/Layer) and implementing the following methods:

- **build**: Create the weights of the layer. Add weights with the `add_weight` method.
- **call**: Define the forward pass.
- **compute_output_shape**: Specify how to compute the output shape of the layer given the input shape.
- Optionally, a layer can be serialized by implementing the `get_config` method and the `from_config` class method.

Here's an example of a custom layer that implements a `matmul` of an input with a kernel matrix:



```
class MyLayer(keras.layers.Layer):

    def __init__(self, output_dim, **kwargs):
        self.output_dim = output_dim
        super(MyLayer, self).__init__(**kwargs)

    def build(self, input_shape):
        shape = tf.TensorShape((input_shape[1], self.output_dim))
        # Create a trainable weight variable for this layer.
        self.kernel = self.add_weight(name='kernel',
                                       shape=shape,
                                       initializer='uniform',
                                       trainable=True)

        # Be sure to call this at the end
        super(MyLayer, self).build(input_shape)

    def call(self, inputs):
        return tf.matmul(inputs, self.kernel)

    def compute_output_shape(self, input_shape):
        shape = tf.TensorShape(input_shape).as_list()
        shape[-1] = self.output_dim
        return tf.TensorShape(shape)

    def get_config(self):
        base_config = super(MyLayer, self).get_config()
        base_config['output_dim'] = self.output_dim

    @classmethod
    def from_config(cls, config):
        return cls(**config)

# Create a model using the custom layer
model = keras.Sequential([MyLayer(10),
                           keras.layers.Activation('softmax')])

# The compile step specifies the training configuration
model.compile(optimizer=tf.train.RMSPropOptimizer(0.001),
              loss='categorical_crossentropy',
              metrics=['accuracy'])

# Trains for 5 epochs.
model.fit(data, targets, batch_size=32, epochs=5)
```

Callbacks

A callback is an object passed to a model to customize and extend its behavior during training. You can write your own custom callback, or use the built-in [tf.keras.callbacks](https://www.tensorflow.org/api_docs/python/tf/keras/callbacks) (https://www.tensorflow.org/api_docs/python/tf/keras/callbacks) that include:

- [tf.keras.callbacks.ModelCheckpoint](https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/ModelCheckpoint)
(https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/ModelCheckpoint): Save checkpoints of your model at regular intervals.
- [tf.keras.callbacks.LearningRateScheduler](https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/LearningRateScheduler)
(https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/LearningRateScheduler): Dynamically change the learning rate.
- [tf.keras.callbacks.EarlyStopping](https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/EarlyStopping)
(https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/EarlyStopping): Interrupt training when validation performance has stopped improving.
- [tf.keras.callbacks.TensorBoard](https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/TensorBoard)
(https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/TensorBoard): Monitor the model's behavior using [TensorBoard](https://www.tensorflow.org/guide/summaries_and_tensorboard) (https://www.tensorflow.org/guide/summaries_and_tensorboard).

To use a [tf.keras.callbacks.Callback](https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/Callback)

(https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/Callback), pass it to the model's `fit` method:

```
callbacks = [  
    # Interrupt training if `val_loss` stops improving for over 2 epochs  
    keras.callbacks.EarlyStopping(patience=2, monitor='val_loss'),  
    # Write TensorBoard logs to `./logs` directory  
    keras.callbacks.TensorBoard(log_dir='./logs')  
]  
model.fit(data, labels, batch_size=32, epochs=5, callbacks=callbacks,  
          validation_data=(val_data, val_targets))
```



Save and restore

Weights only

Save and load the weights of a model using [tf.keras.Model.save_weights](https://www.tensorflow.org/api_docs/python/tf/keras/Model#save_weights)

(https://www.tensorflow.org/api_docs/python/tf/keras/Model#save_weights):

```
# Save weights to a TensorFlow Checkpoint file
model.save_weights('./my_model')
```



```
# Restore the model's state,
# this requires a model with the same architecture.
model.load_weights('my_model')
```

By default, this saves the model's weights in the TensorFlow checkpoint (<https://www.tensorflow.org/guide/checkpoints>) file format. Weights can also be saved to the Keras HDF5 format (the default for the multi-backend implementation of Keras):

```
# Save weights to a HDF5 file
model.save_weights('my_model.h5', save_format='h5')

# Restore the model's state
model.load_weights('my_model.h5')
```



Configuration only

A model's configuration can be saved—this serializes the model architecture without any weights. A saved configuration can recreate and initialize the same model, even without the code that defined the original model. Keras supports JSON and YAML serialization formats:

```
# Serialize a model to JSON format
json_string = model.to_json()

# Recreate the model (freshly initialized)
fresh_model = keras.models.from_json(json_string)

# Serializes a model to YAML format
yaml_string = model.to_yaml()

# Recreate the model
fresh_model = keras.models.from_yaml(yaml_string)
```



Caution: Subclassed models are not serializable because their architecture is defined by the Python code in the body of the `call` method.

Entire model

The entire model can be saved to a file that contains the weight values, the model's configuration, and even the optimizer's configuration. This allows you to checkpoint a

model and resume training later—from the exact same state—without access to the original code.



```
# Create a trivial model
model = keras.Sequential([
    keras.layers.Dense(10, activation='softmax', input_shape=(32,)),
    keras.layers.Dense(10, activation='softmax')
])
model.compile(optimizer='rmsprop',
              loss='categorical_crossentropy',
              metrics=['accuracy'])
model.fit(data, targets, batch_size=32, epochs=5)

# Save entire model to a HDF5 file
model.save('my_model.h5')

# Recreate the exact same model, including weights and optimizer.
model = keras.models.load_model('my_model.h5')
```

Eager execution

Eager execution (<https://www.tensorflow.org/guide/eager>) is an imperative programming environment that evaluates operations immediately. This is not required for Keras, but is supported by **`tf.keras`** (https://www.tensorflow.org/api_docs/python/tf/keras) and useful for inspecting your program and debugging.

All of the **`tf.keras`** (https://www.tensorflow.org/api_docs/python/tf/keras) model-building APIs are compatible with eager execution. And while the `Sequential` and functional APIs can be used, eager execution especially benefits *model subclassing* and building *custom layers*—the APIs that require you to write the forward pass as code (instead of the APIs that create models by assembling existing layers).

See the [eager execution guide](https://www.tensorflow.org/guide/eager#build_a_model) (https://www.tensorflow.org/guide/eager#build_a_model) for examples of using Keras models with custom training loops and **`tf.GradientTape`** (https://www.tensorflow.org/api_docs/python/tf/GradientTape).

Distribution

Estimators

The Estimators (<https://www.tensorflow.org/guide/estimators>) API is used for training models for distributed environments. This targets industry use cases such as distributed training on large datasets that can export a model for production.

A `tf.keras.Model` (https://www.tensorflow.org/api_docs/python/tf/keras/Model) can be trained with the `tf.estimator` (https://www.tensorflow.org/api_docs/python/tf/estimator) API by converting the model to an `tf.estimator.Estimator` (https://www.tensorflow.org/api_docs/python/tf/estimator/Estimator) object with `tf.keras.estimator.model_to_estimator` (https://www.tensorflow.org/api_docs/python/tf/keras/estimator/model_to_estimator). See Creating Estimators from Keras models (https://www.tensorflow.org/guide/estimators#creating_estimators_from_keras_models).



```
model = keras.Sequential([layers.Dense(10, activation='softmax'),
                           layers.Dense(10, activation='softmax')])

model.compile(optimizer=tf.train.RMSPropOptimizer(0.001),
              loss='categorical_crossentropy',
              metrics=['accuracy'])

estimator = keras.estimator.model_to_estimator(model)
```

Note: Enable eager execution (<https://www.tensorflow.org/guide/eager>) for debugging Estimator input functions (https://www.tensorflow.org/guide/premade_estimators#create_input_functions) and inspecting data.

Multiple GPUs

`tf.keras` (https://www.tensorflow.org/api_docs/python/tf/keras) models can run on multiple GPUs using `tf.contrib.distribute.DistributionStrategy` (https://www.tensorflow.org/api_docs/python/tf/contrib/distribute/DistributionStrategy). This API provides distributed training on multiple GPUs with almost no changes to existing code.

Currently, `tf.contrib.distribute.MirroredStrategy` (https://www.tensorflow.org/api_docs/python/tf/contrib/distribute/MirroredStrategy) is the only supported distribution strategy. `MirroredStrategy` does in-graph replication with synchronous training using all-reduce on a single machine. To use `DistributionStrategy` with Keras, convert the `tf.keras.Model` (https://www.tensorflow.org/api_docs/python/tf/keras/Model) to a `tf.estimator.Estimator` (https://www.tensorflow.org/api_docs/python/tf/estimator/Estimator) with `tf.keras.estimator.model_to_estimator`

(https://www.tensorflow.org/api_docs/python/tf/keras/estimator/model_to_estimator), then train the estimator

The following example distributes a **`tf.keras.Model`**

(https://www.tensorflow.org/api_docs/python/tf/keras/Model) across multiple GPUs on a single machine.

First, define a simple model:

```
model = keras.Sequential()
model.add(keras.layers.Dense(16, activation='relu', input_shape=(10,)))
model.add(keras.layers.Dense(1, activation='sigmoid'))

optimizer = tf.train.GradientDescentOptimizer(0.2)

model.compile(loss='binary_crossentropy', optimizer=optimizer)
model.summary()
```



Define an *input pipeline*. The `input_fn` returns a **`tf.data.Dataset`**

(https://www.tensorflow.org/api_docs/python/tf/data/Dataset) object used to distribute the data across multiple devices—with each device processing a slice of the input batch.

```
def input_fn():
    x = np.random.random((1024, 10))
    y = np.random.randint(2, size=(1024, 1))
    x = tf.cast(x, tf.float32)
    dataset = tf.data.Dataset.from_tensor_slices((x, y))
    dataset = dataset.repeat(10)
    dataset = dataset.batch(32)
    return dataset
```



Next, create a **`tf.estimator.RunConfig`**

(https://www.tensorflow.org/api_docs/python/tf/estimator/RunConfig) and set the

`train_distribute` argument to the **`tf.contrib.distribute.MirroredStrategy`**

(https://www.tensorflow.org/api_docs/python/tf/contrib/distribute/MirroredStrategy) instance. When creating **`MirroredStrategy`**, you can specify a list of devices or set the `num_gpus` argument. The default uses all available GPUs, like the following:

```
strategy = tf.contrib.distribute.MirroredStrategy()
config = tf.estimator.RunConfig(train_distribute=strategy)
```



Convert the Keras model to a **`tf.estimator.Estimator`**

(https://www.tensorflow.org/api_docs/python/tf/estimator/Estimator) instance:

```
keras_estimator = keras.estimator.model_to_estimator(  
    keras_model=model,  
    config=config,  
    model_dir='/tmp/model_dir')
```



Finally, train the `Estimator` instance by providing the `input_fn` and `steps` arguments:

```
keras_estimator.train(input_fn=input_fn, steps=10)
```



Except as otherwise noted, the content of this page is licensed under the [Creative Commons Attribution 3.0 License](https://creativecommons.org/licenses/by/3.0/) (<https://creativecommons.org/licenses/by/3.0/>), and code samples are licensed under the [Apache 2.0 License](https://www.apache.org/licenses/LICENSE-2.0) (<https://www.apache.org/licenses/LICENSE-2.0>). For details, see our [Site Policies](https://developers.google.com/terms/site-policies) (<https://developers.google.com/terms/site-policies>). Java is a registered trademark of Oracle and/or its affiliates.

Last updated August 8, 2018.