

07 Milestone Project 1: 🍷 Food Vision Big™

In the previous notebook ([transfer learning part 3: scaling up](https://github.com/mrdbourke/tensorflow-deep-learning/blob/main/06_transfer_learning_in_tensorflow_part_3_scaling_up) (https://github.com/mrdbourke/tensorflow-deep-learning/blob/main/06_transfer_learning_in_tensorflow_part_3_scaling_up) we built Food Vision mini: a transfer learning model which beat the original results of the [Food101 paper](https://data.vision.ee.ethz.ch/cvl/datasets_extra/food-101/) (https://data.vision.ee.ethz.ch/cvl/datasets_extra/food-101/) with only 10% of the data.

But you might be wondering, what would happen if we used all the data?

Well, that's what we're going to find out in this notebook!

We're going to be building Food Vision Big™, using all of the data from the Food101 dataset.

Yep. All 75,750 training images and 25,250 testing images.

And guess what...

This time we've got the goal of beating [DeepFood](https://www.researchgate.net/publication/304163308_DeepFood_Deep_Learn_Based_Food_Image_Recognition_for_Computer-Aided_Dietary_Assessment) (https://www.researchgate.net/publication/304163308_DeepFood_Deep_Learn_Based_Food_Image_Recognition_for_Computer-Aided_Dietary_Assessment), a 2016 paper which used a Convolutional Neural Network trained for 2-3 days to achieve 77.4% top-1 accuracy.

📌 **Note: Top-1 accuracy** means "accuracy for the top softmax activation value output by the model" (because softmax outputs a value for every class, but top-1 means only the highest one is evaluated). **Top-5 accuracy** means "accuracy for the top 5 softmax activation values output by the model", in other words, did the true label appear in the top 5 activation values? Top-5 accuracy scores are usually noticeably higher than top-1.

🍷 Food Vision Big™

Dataset source	TensorFlow Datasets	Preprocessed
Train data	75,750 images	
Test data	25,250 images	

Mixed precision	Yes		
Data loading	Permanant tf.data API	TensorFlow	
Target results	77.4% top-1 accuracy (beat DeepFood paper (https://arxiv.org/abs/1606.05675))	50.76% top-1 accuracy	(https://data.vision.ee.ethz.ch/cv101/static/bossard)

Table comparing difference between Food Vision Big (this notebook) versus Food Vision mini (previous notebook).

Alongside attempting to beat the DeepFood paper, we're going to learn about two methods to significantly improve the speed of our model training:

1. Prefetching
2. Mixed precision training

But more on these later.

What we're going to cover

- Using TensorFlow Datasets to download and explore data
- Creating preprocessing function for our data
- Batching & preparing datasets for modelling (**making our datasets run fast**)
- Creating modelling callbacks
- Setting up **mixed precision training**
- Building a feature extraction model (see [transfer learning part 1: feature extraction](https://github.com/mrdbourke/tensorflow-deep-learning/blob/main/04_transfer_learning_in_tensorflow_part_1_feature_extraction) (https://github.com/mrdbourke/tensorflow-deep-learning/blob/main/04_transfer_learning_in_tensorflow_part_1_feature_extraction))
- Fine-tuning the feature extraction model (see [transfer learning part 2: fine-tuning](https://github.com/mrdbourke/tensorflow-deep-learning/blob/main/05_transfer_learning_in_tensorflow_part_2_fine_tuning) (https://github.com/mrdbourke/tensorflow-deep-learning/blob/main/05_transfer_learning_in_tensorflow_part_2_fine_tuning))
- Viewing training results on TensorBoard

How you should approach this notebook

You can read through the descriptions and the code (it should all run, except for the cells which error on purpose), but there's a better option.

Write all of the code yourself.

Yes. I'm serious. Create a new notebook, and rewrite each line by yourself. Investigate it, see if you can break it, why does it break?

You don't have to write the text descriptions but writing the code yourself is a great way to get hands-on experience.

Don't worry if you make mistakes, we all do. The way to get better and make less mistakes is to write more code.

□ **Resource:** See the full set of course materials on

Check GPU

For this notebook, we're going to be doing something different.

We're going to be using mixed precision training.

Mixed precision training was introduced in [TensorFlow 2.4.0](https://blog.tensorflow.org/2020/12/whats-new-in-tensorflow-24.html) (<https://blog.tensorflow.org/2020/12/whats-new-in-tensorflow-24.html>) (a very new feature at the time of writing).

What does **mixed precision training** do?

Mixed precision training uses a combination of single precision (float32) and half-precision (float16) data types to speed up model training (up 3x on modern GPUs).

We'll talk about this more later on but in the meantime you can read the [TensorFlow documentation on mixed precision](https://www.tensorflow.org/guide/mixed_precision) (https://www.tensorflow.org/guide/mixed_precision) for more details.

For now, before we can move forward if we want to use mixed precision training, we need to make sure the GPU powering our Google Colab instance (if you're using Google Colab) is compatible.

For mixed precision training to work, you need access to a GPU with a compute capability score of 7.0+.

Google Colab offers P100, K80 and T4 GPUs, however, **the P100 and K80 aren't compatible with mixed precision training.**

Therefore before we proceed we need to make sure we have **access to a Tesla T4 GPU in our Google Colab instance.**

If you're not using Google Colab, you can find a list of various [Nvidia GPU compute capabilities on Nvidia's developer website](https://developer.nvidia.com/cuda-gpus#compute) (<https://developer.nvidia.com/cuda-gpus#compute>).

□ **Note:** If you run the cell below and see a P100 or K80, try going to Runtime -> Factory Reset Runtime (note: this will remove any saved variables and data from your Colab instance) and then retry to get a T4.

```
In [1]: # If using Google Colab, this should output "Tesla T4" otherwise,
# you won't be able to use mixed precision training
!nvidia-smi -L
```

GPU 0: NVIDIA TITAN RTX (UUID: GPU-64b1678c-cec3-56bb-af0c-8ae69de44cbd)

Since mixed precision training was introduced in TensorFlow 2.4.0, make sure you've got at least TensorFlow 2.4.0+.

```
In [2]: # Check TensorFlow version (should be 2.4.0+)
import tensorflow as tf
print(tf.__version__)
```

2.6.2

Get helper functions

We've created a series of helper functions throughout the previous notebooks in the course. Instead of rewriting them (tedious), we'll import the `helper_functions.py` (https://github.com/mrdbourke/tensorflow-deep-learning/blob/main/extras/helper_functions.py) file from the GitHub repo.

```
In [4]: # Get helper functions file
import os

if not os.path.exists("helper_functions.py"):
    !wget https://raw.githubusercontent.com/mrdbourke/tensorflow-deep-learning/main/extras/helper_functions.py
else:
    print("[INFO] 'helper_functions.py' already exists, skipping download")

[INFO] 'helper_functions.py' already exists, skipping download.
```

```
In [5]: # Import series of helper functions for the notebook (we've created/use
from helper_functions import create_tensorboard_callback, plot_loss_cur
```

Use TensorFlow Datasets to Download Data

In previous notebooks, we've downloaded our food images (from the [Food101 dataset \(https://www.kaggle.com/dansbecker/food-101/home\)](https://www.kaggle.com/dansbecker/food-101/home)) from Google Storage.

And this is a typical workflow you'd use if you're working on your own datasets.

However, there's another way to get datasets ready to use with TensorFlow.

For many of the most popular datasets in the machine learning world (often referred to and used as benchmarks), you can access them through [TensorFlow Datasets \(TFDS\)](https://www.tensorflow.org/datasets/overview).
(<https://www.tensorflow.org/datasets/overview>).

What is **TensorFlow Datasets**?

A place for prepared and ready-to-use machine learning datasets.

Why use TensorFlow Datasets?

- Load data already in Tensors
- Practice on well established datasets
- Experiment with different data loading techniques (like we're going to use in this notebook)
- Experiment with new TensorFlow features quickly (such as mixed precision training)

Why *not* use TensorFlow Datasets?

- The datasets are static (they don't change, like your real-world datasets would)
- Might not be suited for your particular problem (but great for experimenting)

To begin using TensorFlow Datasets we can import it under the alias `tfds`.

```
In [6]: # Get TensorFlow Datasets  
import tensorflow_datasets as tfds
```

To find all of the available datasets in TensorFlow Datasets, you can use the `list_builders()` method.

After doing so, we can check to see if the one we're after (`"food101"`) is present.

```
In [7]: # List available datasets  
datasets_list = tfds.list_builders() # get all available datasets in TF  
print("food101" in datasets_list) # is the dataset we're after available
```

True

Beautiful! It looks like the dataset we're after is available (note there are plenty more available but we're on Food101).

To get access to the Food101 dataset from the TFDS, we can use the `tfds.load()`
(https://www.tensorflow.org/datasets/api_docs/python/tfds/load) method.

In particular, we'll have to pass it a few parameters to let it know what we're after:

- `name` (str) : the target dataset (e.g. "food101")
- `split` (list, optional) : what splits of the dataset we're after (e.g. ["train", "validation"])
 - the `split` parameter is quite tricky. See [the documentation for more](https://github.com/tensorflow/datasets/blob/master/docs/splits.md) (<https://github.com/tensorflow/datasets/blob/master/docs/splits.md>)
- `shuffle_files` (bool) : whether or not to shuffle the files on download, defaults to `False`
- `as_supervised` (bool) : `True` to download data samples in tuple format `((data, label))` or `False` for dictionary format
- `with_info` (bool) : `True` to download dataset metadata (labels, number of samples, etc)

□ **Note:** Calling the `tfds.load()` method will start to download a target dataset to disk if the `download=True` parameter is set (default). This dataset could be 100GB+, so make sure you have space.

```
In [9]: # Load in the data (takes about 5-6 minutes in Google Colab)
(train_data, test_data), ds_info = tfds.load(name="food101", # target dataset
                                             split=["train", "validation"],
                                             shuffle_files=True, # shuffle files
                                             as_supervised=True, # download as supervised
                                             with_info=True) # include dataset info
```

Wonderful! After a few minutes of downloading, we've now got access to entire Food101 dataset (in tensor format) ready for modelling.

Now let's get a little information from our dataset, starting with the class names.

Getting class names from a TensorFlow Datasets dataset requires downloading the "dataset_info" variable (by using the `as_supervised=True` parameter in the `tfds.load()` method, **note:** this will only work for supervised datasets in TFDS).

We can access the class names of a particular dataset using the `dataset_info.features` attribute and accessing `names` attribute of the the "label" key.

```
In [10]: # Features of Food101 TFDS
ds_info.features
```

```
Out[10]: FeaturesDict({
    'image': Image(shape=(None, None, 3), dtype=tf.uint8),
    'label': ClassLabel(shape=(), dtype=tf.int64, num_classes=101),
})
```

```
In [11]: # Get class names
class_names = ds_info.features["label"].names
class_names[:10]
```

```
Out[11]: ['apple_pie',
    'baby_back_ribs',
    'baklava',
    'beef_carpaccio',
    'beef_tartare',
    'beet_salad',
    'beignets',
    'bibimbap',
    'bread_pudding',
    'breakfast_burrito']
```

Exploring the Food101 data from TensorFlow Datasets

Now we've downloaded the Food101 dataset from TensorFlow Datasets, how about we do what any good data explorer should?

In other words, "visualize, visualize, visualize".

Let's find out a few details about our dataset:

- The shape of our input data (image tensors)
- The datatype of our input data
- What the labels of our input data look like (e.g. one-hot encoded versus label-encoded)
- Do the labels match up with the class names?

To do, let's take one sample off the training data (using the [.take\(\) method](https://www.tensorflow.org/api_docs/python/tf/data/Dataset#take) (https://www.tensorflow.org/api_docs/python/tf/data/Dataset#take)) and explore it.

```
In [12]: # Take one sample off the training data
train_one_sample = train_data.take(1) # samples are in format (image_tensor, label)
```

Because we used the `as_supervised=True` parameter in our `tfds.load()` method above, data samples come in the tuple format structure `(data, label)` or in our case `(image_tensor, label)`.

```
In [13]: # What does one sample of our training data look like?  
train_one_sample
```

```
Out[13]: <TakeDataset shapes: ((None, None, 3), ()), types: (tf.uint8, tf.int64)>
```

Let's loop through our single training sample and get some info from the `image_tensor` and `label`.

```
In [15]: # Output info about our training sample  
for image, label in train_one_sample:  
    print(f"""  
        Image shape: {image.shape}  
        Image dtype: {image.dtype}  
        Target class from Food101 (tensor form): {label}  
        Class name (str form): {class_names[label.numpy()]}  
        """)
```

```
Image shape: (512, 512, 3)  
Image dtype: <dtype: 'uint8'>  
Target class from Food101 (tensor form): 8  
Class name (str form): bread_pudding
```

Because we set the `shuffle_files=True` parameter in our `tfds.load()` method above, running the cell above a few times will give a different result each time.

Checking these you might notice some of the images have different shapes, for example (512, 342, 3) and (512, 512, 3) (height, width, color_channels).

Let's see what one of the image tensors from TFDS's Food101 dataset looks like.


```
In [16]: # What does an image tensor from TFDS's Food101 look like?
image
```

```
Out[16]: <tf.Tensor: shape=(512, 512, 3), dtype=uint8, numpy=
array([[18,  6,  8],
       [18,  6,  8],
       [18,  6,  8],
       ...,
       [30, 15, 22],
       [29, 14, 21],
       [26, 11, 18]],

       [[22, 10, 12],
        [21,  9, 11],
        [20,  8, 10],
        ...,
        [35, 20, 27],
        [31, 16, 23],
        [26, 11, 18]],

       [[23, 13, 14],
        [21, 11, 12],
        [19,  9, 10],
        ...,
        [39, 26, 33],
        [36, 21, 28],
        [30, 15, 22]],

       ...,

       [[15,  4,  8],
        [15,  4,  8],
        [14,  5, 10],
        ...,
        [41,  9, 10],
        [39,  7,  8],
        [36,  4,  5]],

       [[16,  5,  9],
        [16,  5,  9],
        [16,  5, 11],
        ...,
        [42, 12, 12],
        [39,  9,  9],
        [35,  5,  5]],

       [[15,  4,  8],
        [15,  4,  8],
        [16,  5, 11],
        ...,
        [41, 11, 11],
        [39,  9,  9],
        [35,  5,  5]]], dtype=uint8)>
```

```
In [17]: # What are the min and max values?  
tf.reduce_min(image), tf.reduce_max(image)
```

```
Out[17]: (<tf.Tensor: shape=(), dtype=uint8, numpy=0>,  
         <tf.Tensor: shape=(), dtype=uint8, numpy=255>)
```

Alright looks like our image tensors have values of between 0 & 255 (standard red, green, blue colour values) and the values are of data type uint8 .

We might have to preprocess these before passing them to a neural network. But we'll handle this later.

In the meantime, let's see if we can plot an image sample.

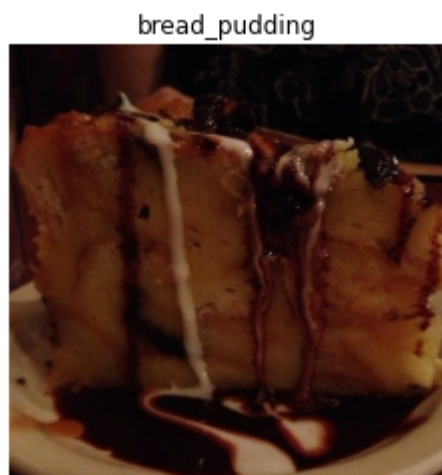
Plot an image from TensorFlow Datasets

We've seen our image tensors in tensor format, now let's really adhere to our motto.

"Visualize, visualize, visualize!"

Let's plot one of the image samples using `matplotlib.pyplot.imshow()` (https://matplotlib.org/stable/api/as_gen/matplotlib.pyplot.imshow.html) and set the title to target class name.

```
In [18]: # Plot an image tensor  
import matplotlib.pyplot as plt  
plt.imshow(image)  
plt.title(class_names[label.numpy()]) # add title to image by indexing  
plt.axis(False);
```



Delicious!

Okay, looks like the Food101 data we've got from TFDS is similar to the datasets we've been using in previous notebooks.

Now let's preprocess it and get it ready for use with a neural network.

Create preprocessing functions for our data

In previous notebooks, when our images were in folder format we used the method `tf.keras.utils.image_dataset_from_directory()`. (https://www.tensorflow.org/api_docs/python/tf/keras/utils/image_dataset) to load them in.

Doing this meant our data was loaded into a format ready to be used with our models.

However, since we've downloaded the data from TensorFlow Datasets, there are a couple of preprocessing steps we have to take before it's ready to model.

More specifically, our data is currently:

- In `uint8` data type
- Comprised of all different sized tensors (different sized images)
- Not scaled (the pixel values are between 0 & 255)

Whereas, models like data to be:

- In `float32` data type
- Have all of the same size tensors (batches require all tensors have the same shape, e.g. (224, 224, 3))
- Scaled (values between 0 & 1), also called normalized

To take care of these, we'll create a `preprocess_img()` function which:

- Resizes an input image tensor to a specified size using `tf.image.resize()`. (https://www.tensorflow.org/api_docs/python/tf/image/resize)
- Converts an input image tensor's current datatype to `tf.float32` using `tf.cast()`. (https://www.tensorflow.org/api_docs/python/tf/cast)

□ **Note:** Pretrained EfficientNetBX models in `tf.keras.applications.efficientnet`
(https://www.tensorflow.org/api_docs/python/tf/keras/applications)

```
In [19]: # Make a function for preprocessing images
def preprocess_img(image, label, img_shape=224):
    """
    Converts image datatype from 'uint8' -> 'float32' and reshapes image
    [img_shape, img_shape, color_channels]
    """
    image = tf.image.resize(image, [img_shape, img_shape]) # reshape to
    return tf.cast(image, tf.float32), label # return (float32_image, label)
```

Our `preprocess_img()` function above takes image and label as input (even though it does nothing to the label) because our dataset is currently in the tuple structure (image, label) .

Let's try our function out on a target image.

```
In [20]: # Preprocess a single sample image and check the outputs
preprocessed_img = preprocess_img(image, label)[0]
print(f"Image before preprocessing:\n {image[:2]}...\nShape: {image.sh
print(f"Image after preprocessing:\n {preprocessed_img[:2]}...\nShape:
```

Image before preprocessing:

```
[[[18  6  8]
   [18  6  8]
   [18  6  8]
   ...
   [30 15 22]
   [29 14 21]
   [26 11 18]]
```

```
[[[22 10 12]
   [21  9 11]
   [20  8 10]
```

```
...
   [35 20 27]
   [31 16 23]
   [26 11 18]]]...,
```

Shape: (512, 512, 3),

Datatype: <dtype: 'uint8'>

Image after preprocessing:

```
[[[20.158163  8.158163 10.158163 ]
   [18.42347  7.6173472 9.020408 ]
   [15.010203  6.423469  9.285714 ]
   ...
   [26.285824 15.714351 23.07156  ]
   [31.091867 17.285728 24.285728 ]
   [28.754953 13.754952 20.754953 ]]
```

```
[[[18.92857  8.928571  9.928571 ]
   [16.214285  7.0765305 8.07653  ]
   [14.739796  8.571429 10.627552 ]
```

```
...
   [26.444029 15.872557 21.658293 ]
   [39.86226  26.862259 33.86226  ]
   [39.49479  24.494787 31.494787 ]]]...,
```

Shape: (224, 224, 3),

Datatype: <dtype: 'float32'>

Excellent! Looks like our `preprocess_img()` function is working as expected.

The input image gets converted from `uint8` to `float32` and gets reshaped from its current shape to `(224, 224, 3)`.

How does it look?

```
In [21]: # We can still plot our preprocessed image as long as we
# divide by 255 (for matplotlib capatibility)
plt.imshow(preprocessed_img/255.)
plt.title(class_names[label])
plt.axis(False);
```

bread_pudding



All this food visualization is making me hungry. How about we start preparing to model it?

Batch & prepare datasets

Before we can model our data, we have to turn it into batches.

Why?

Because computing on batches is memory efficient.

We turn our data from 101,000 image tensors and labels (train and test combined) into batches of 32 image and label pairs, thus enabling it to fit into the memory of our GPU.

To do this in effective way, we're going to be leveraging a number of methods from the [tf.data API](https://www.tensorflow.org/api_docs/python/tf/data) (https://www.tensorflow.org/api_docs/python/tf/data).

□ **Resource:** For loading data in the most performant way possible, see the TensorFlow docuemntation on [Better performance with the tf.data API](https://www.tensorflow.org/guide/data_performance) (https://www.tensorflow.org/guide/data_performance).

Specifically, we're going to be using:

- `map()` (https://www.tensorflow.org/api_docs/python/tf/data/Dataset#map) - maps a predefined function to a target dataset (e.g. `preprocess_img()` to our image tensors)

- `shuffle()`
(https://www.tensorflow.org/api_docs/python/tf/data/Dataset#shuffle)
- randomly shuffles the elements of a target dataset up
buffer_size (ideally, the buffer_size is equal to the size of
the dataset, however, this may have implications on memory)
- `batch()`
(https://www.tensorflow.org/api_docs/python/tf/data/Dataset#batch)
- turns elements of a target dataset into batches (size defined by
parameter batch_size)
- `prefetch()`
(https://www.tensorflow.org/api_docs/python/tf/data/Dataset#prefetch)
- prepares subsequent batches of data whilst other batches of data
are being computed on (improves data loading speed but costs
memory)
- Extra: `cache()`
(https://www.tensorflow.org/api_docs/python/tf/data/Dataset#cache)
- caches (saves them for later) elements in a target dataset,
saving loading time (will only work if your dataset is small
enough to fit in memory, standard Colab instances only have 12GB
of memory)

Things to note:

- Can't batch tensors of different shapes (e.g. different image
sizes, need to reshape images first, hence our preprocess_img()
function)
- `shuffle()` keeps a buffer of the number you pass it images
shuffled, ideally this number would be all of the samples in your
training set, however, if your training set is large, this buffer
might not fit in memory (a fairly large number like 1000 or 10000
is usually suffice for shuffling)
- For methods with the num_parallel_calls parameter available (such
as `map()`), setting it to `num_parallel_calls=tf.data.AUTOTUNE` will
parallelize preprocessing and significantly improve speed
- Can't use `cache()` unless your dataset can fit in memory

Woah, the above is alot. But once we've coded below, it'll start to
make sense.

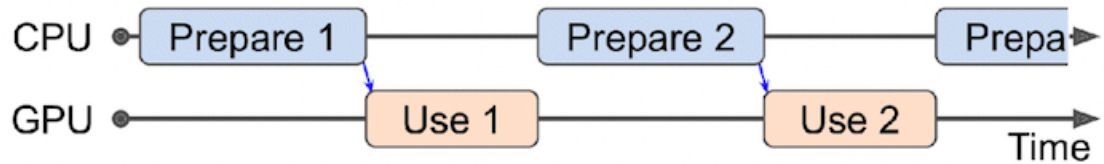
We're going to through things in the following order:

Original dataset (e.g. train_data) -> `map()` -> `shuffle()` -> `batch()` -> `prefetch()` -> `PrefetchDataset`

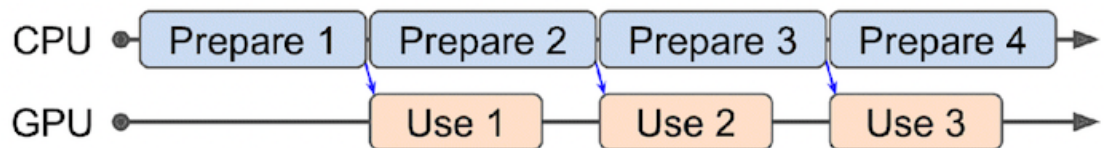
This is like saying,

"Hey, map this preprocessing function across our training dataset, then shuffle a number of elements before batching them together and make sure you prepare new batches (prefetch) whilst the model is looking through the current batch".

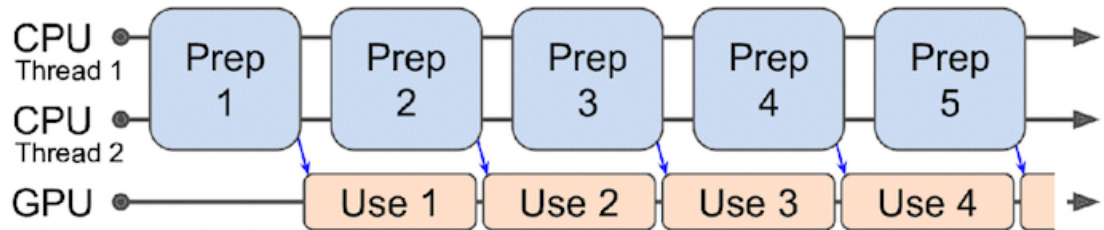
Without prefetching



With prefetching



With prefetching + multithreaded loading & preprocessing



```
In [22]: # Map preprocessing function to training data (and parallelize)
train_data = train_data.map(map_func=preprocess_img, num_parallel_calls=
# Shuffle train_data and turn it into batches and prefetch it (load it
train_data = train_data.shuffle(buffer_size=1000).batch(batch_size=32).

# Map preprocessing function to test data
test_data = test_data.map(preprocess_img, num_parallel_calls=tf.data.AU
# Turn test data into batches (don't need to shuffle)
test_data = test_data.batch(32).prefetch(tf.data.AUTOTUNE)
```

And now let's check out what our prepared datasets look like.

```
In [23]: train_data, test_data
```

```
Out[23]: (<PrefetchDataset shapes: ((None, 224, 224, 3), (None,)), types: (tf.f
loat32, tf.int64)>,
<PrefetchDataset shapes: ((None, 224, 224, 3), (None,)), types: (tf.f
loat32, tf.int64)>)
```


Excellent! Looks like our data is now in tuples of (image, label) with datatypes of (tf.float32, tf.int64) , just what our model is after.

□ **Note:** You can get away without calling the `prefetch()` method on the end of your datasets, however, you'd probably see significantly slower data loading speeds when building a model. So most of your dataset input pipelines should end with a call to `prefetch()`.
https://www.tensorflow.org/api_docs/python/tf/data/Dataset#prefetch

Onward.

Create modelling callbacks

Since we're going to be training on a large amount of data and training could take a long time, it's a good idea to set up some modelling callbacks so we be sure of things like our model's training logs being tracked and our model being checkpointed (saved) after various training milestones.

To do each of these we'll use the following callbacks:

- `tf.keras.callbacks.TensorBoard()`
https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/TensorBoard
- allows us to keep track of our model's training history so we can inspect it later (**note:** we've created this callback before we have imported it from `helper_functions.py` as `create_tensorboard_callback()`)
- `tf.keras.callbacks.ModelCheckpoint()`
https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/ModelCheckpoint
- saves our model's progress at various intervals so we can load it and reuse it later without having to retrain it
 - Checkpointing is also helpful so we can start fine-tuning our model at a particular epoch and revert back to a previous state if fine-tuning offers no benefits

```
In [24]: # Create TensorBoard callback (already have "create_tensorboard_callback"
from helper_functions import create_tensorboard_callback

# Create ModelCheckpoint callback to save model's progress
checkpoint_path = "model_checkpoints/cp.ckpt" # saving weights requires
model_checkpoint = tf.keras.callbacks.ModelCheckpoint(checkpoint_path,
                                                        monitor="val_accu
                                                        save_best_only=True
                                                        save_weights_only
                                                        verbose=0) # don't
```

Setup mixed precision training

We touched on mixed precision training above.

However, we didn't quite explain it.

Normally, tensors in TensorFlow default to the float32 datatype (unless otherwise specified).

In computer science, float32 is also known as [single-precision floating-point format](https://en.wikipedia.org/wiki/Single-precision_floating-point_format) (https://en.wikipedia.org/wiki/Single-precision_floating-point_format). The 32 means it usually occupies 32 bits in computer memory.

Your GPU has a limited memory, therefore it can only handle a number of float32 tensors at the same time.

This is where mixed precision training comes in.

Mixed precision training involves using a mix of float16 and float32 tensors to make better use of your GPU's memory.

Can you guess what float16 means?

Well, if you thought since float32 meant single-precision floating-point, you might've guessed float16 means [half-precision floating-point format](https://en.wikipedia.org/wiki/Half-precision_floating-point_format) (https://en.wikipedia.org/wiki/Half-precision_floating-point_format). And if you did, you're right! And if not, no trouble, now you know.

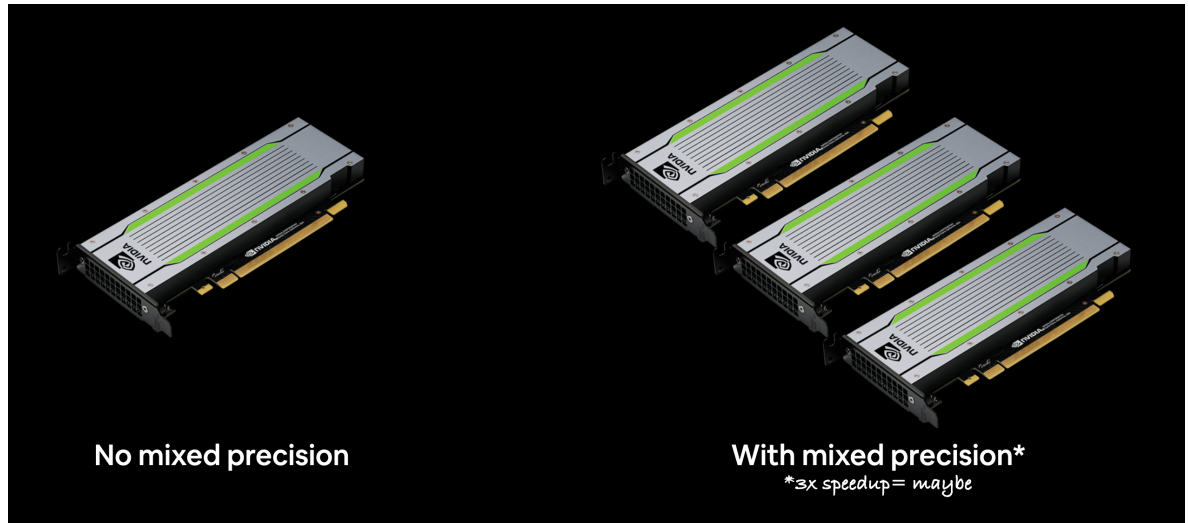
For tensors in float16 format, each element occupies 16 bits in computer memory.

So, where does this leave us?

As mentioned before, when using mixed precision training, your model will make use of float32 and float16 data types to use less memory where possible and in turn run faster (using less memory per tensor means more tensors can be computed on simultaneously).

As a result, using mixed precision training can improve your performance on modern GPUs (those with a compute capability score of 7.0+) by up to 3x.

For a more detailed explanation, I encourage you to read through the [TensorFlow mixed precision guide](https://www.tensorflow.org/guide/mixed_precision) (https://www.tensorflow.org/guide/mixed_precision) (I'd highly recommend at least checking out the summary).



Because mixed precision training uses a combination of float32 and float16 data types, you may see up to a 3x speedup on modern GPUs.

□ **Note:** If your GPU doesn't have a score of over 7.0+ (e.g. P100 in Colab), mixed precision won't work (see: ["Supported Hardware"](https://www.tensorflow.org/guide/mixed_precision#supported_hardware) (https://www.tensorflow.org/guide/mixed_precision#supported_hardware) in the mixed precision guide for more).

□ **Resource:** If you'd like to learn more about precision in computer science (the detail to which a numerical quantity is expressed by a computer), see the [Wikipedia page](https://en.wikipedia.org/wiki/Precision_(computer_science)) ([https://en.wikipedia.org/wiki/Precision_\(computer_science\)](https://en.wikipedia.org/wiki/Precision_(computer_science))) (and accompanying resources).

Okay, enough talk, let's see how we can turn on mixed precision training in TensorFlow.

The beautiful thing is, the `tensorflow.keras.mixed_precision` (https://www.tensorflow.org/api_docs/python/tf/keras/mixed_precision/) API has made it very easy for us to get started.

First, we'll import the API and then use the `set_global_policy()` (https://www.tensorflow.org/api_docs/python/tf/keras/mixed_precision/set_global_policy) method to set the default policy to "mixed_float16".

```
In [26]: # Turn on mixed precision training
from tensorflow.keras import mixed_precision
mixed_precision.set_global_policy(policy="mixed_float16") # set global
```

Nice! As long as the GPU you're using has a compute capability of 7.0+ the cell above should run without error.

Now we can check the global dtype policy (the policy which will be used by layers in our model) using the `mixed_precision.global_policy()` (https://www.tensorflow.org/api_docs/python/tf/keras/mixed_precision/global_policy) method.

```
In [27]: mixed_precision.global_policy() # should output "mixed_float16" (if you
```

```
Out[27]: <Policy "mixed_float16">
```

Great, since the global dtype policy is now "mixed_float16" our model will automatically take advantage of float16 variables where possible and in turn speed up training.

Build feature extraction model

Callbacks: ready to roll.

Mixed precision: turned on.

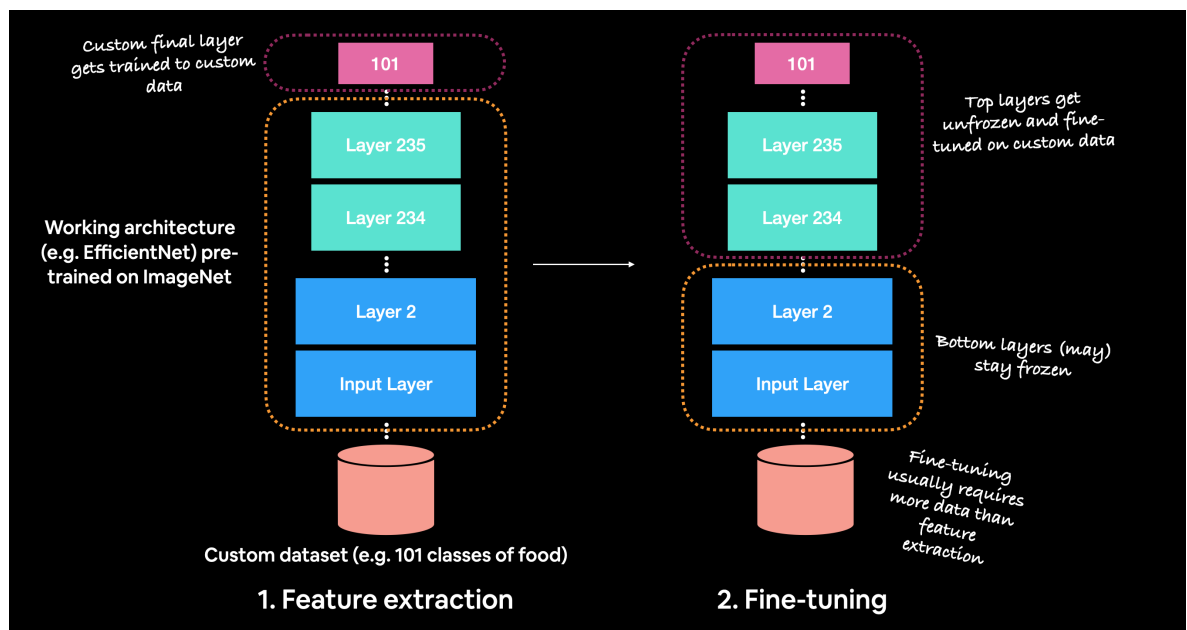
Let's build a model.

Because our dataset is quite large, we're going to move towards fine-tuning an existing pretrained model (EfficientNetB0).

But before we get into fine-tuning, let's set up a feature-extraction model.

Recall, the typical order for using transfer learning is:

1. Build a feature extraction model (replace the top few layers of a pretrained model)
2. Train for a few epochs with lower layers frozen
3. Fine-tune if necessary with multiple layers unfrozen



Before fine-tuning, it's best practice to train a feature extraction model with custom top layers.

To build the feature extraction model (covered in [Transfer Learning in TensorFlow Part 1: Feature extraction](https://github.com/mrdbourke/tensorflow-deep-learning/blob/main/04_transfer_learning_in_tensorflow_part_1_feature_extraction.md) (https://github.com/mrdbourke/tensorflow-deep-learning/blob/main/04_transfer_learning_in_tensorflow_part_1_feature_extraction.md)) we'll:

- Use EfficientNetB0 from `tf.keras.applications` (https://www.tensorflow.org/api_docs/python/tf/keras/applications) pre-trained on ImageNet as our base model
 - We'll download this without the top layers using `include_top=False` parameter so we can create our own output layers
- Freeze the base model layers so we can use the pre-learned patterns the base model has found on ImageNet
- Put together the input, base model, pooling and output layers in a [Functional model](https://keras.io/guides/functional_api/) (https://keras.io/guides/functional_api/)
- Compile the Functional model using the Adam optimizer and [sparse categorical crossentropy](https://www.tensorflow.org/api_docs/python/tf/keras/losses/SparseCategoricalCrossentropy) (https://www.tensorflow.org/api_docs/python/tf/keras/losses/SparseCategoricalCrossentropy) as the loss function (since our labels **aren't** one-hot encoded)
- Fit the model for 3 epochs using the TensorBoard and ModelCheckpoint callbacks

□ **Note:** Since we're using mixed precision training, our model needs a separate output layer with a hard-coded `dtype=float32`, for example, `layers.Activation("softmax", dtype=tf.float32)`. This ensures the outputs of our model are returned back to the `float32` data type which is more numerically stable than the `float16` datatype (important for loss calculations). See the ["Building the model"](#)

```
In [29]: from tensorflow.keras import layers

# Create base model
input_shape = (224, 224, 3)
base_model = tf.keras.applications.EfficientNetB0(include_top=False)
base_model.trainable = False # freeze base model layers

# Create Functional model
inputs = layers.Input(shape=input_shape, name="input_layer")
# Note: EfficientNetBX models have rescaling built-in but if your model
# x = layers.Rescaling(1./255)(x)
x = base_model(inputs, training=False) # set base_model to inference mode
x = layers.GlobalAveragePooling2D(name="pooling_layer")(x)
x = layers.Dense(len(class_names))(x) # want one output neuron per class
# Separate activation of output layer so we can output float32 activations
outputs = layers.Activation("softmax", dtype=tf.float32, name="softmax_float32")(x)
model = tf.keras.Model(inputs, outputs)

# Compile the model
model.compile(loss="sparse_categorical_crossentropy", # Use sparse_categorical_crossentropy
              optimizer=tf.keras.optimizers.Adam(),
              metrics=["accuracy"])
```

```
In [30]: # Check out our model
model.summary()
```

Model: "model_1"

Layer (type)	Output Shape	Param #
input_layer (InputLayer)	[(None, 224, 224, 3)]	0
efficientnetb0 (Functional)	(None, None, None, 1280)	4049571
pooling_layer (GlobalAveragePooling2D)	(None, 1280)	0
dense_1 (Dense)	(None, 101)	129381
softmax_float32 (Activation)	(None, 101)	0
Total params: 4,178,952		
Trainable params: 129,381		
Non-trainable params: 4,049,571		

Checking layer dtype policies (are we using mixed precision?)

Model ready to go!

Before we said the mixed precision API will automatically change our layers' dtype policy's to whatever the global dtype policy is (in our case it's "mixed_float16").

We can check this by iterating through our model's layers and printing layer attributes such as dtype and dtype_policy.

```
In [32]: # Check the dtype_policy attributes of layers in our model
for layer in model.layers:
    print(layer.name, layer.trainable, layer.dtype, layer.dtype_policy)

input_layer True float32 <Policy "float32">
efficientnetb0 False float32 <Policy "mixed_float16">
pooling_layer True float32 <Policy "mixed_float16">
dense_1 True float32 <Policy "mixed_float16">
softmax_float32 True float32 <Policy "float32">
```

Going through the above we see:

- `layer.name` (str) : a layer's human-readable name, can be defined by the name parameter on construction
- `layer.trainable` (bool) : whether or not a layer is trainable (all of our layers are trainable except the `efficientnetb0` layer since we set its trainable attribute to False)
- `layer.dtype` : the data type a layer stores its variables in
- `layer.dtype_policy` : the data type a layer computes in

□ **Note:** A layer can have a dtype of `float32` and a dtype policy of `"mixed_float16"` because it stores its variables (weights & biases) in `float32` (more numerically stable), however it computes in `float16` (faster).

We can also check the same details for our model's base model.

```
In [33]: # Check the layers in the base model and see what dtype policy they're
for layer in model.layers[1].layers[:20]: # only check the first 20 layers
    print(layer.name, layer.trainable, layer.dtype, layer.dtype_policy)
```

```
input_2 False float32 <Policy "float32">
rescaling_1 False float32 <Policy "mixed_float16">
normalization_1 False float32 <Policy "mixed_float16">
stem_conv_pad False float32 <Policy "mixed_float16">
stem_conv False float32 <Policy "mixed_float16">
stem_bn False float32 <Policy "mixed_float16">
stem_activation False float32 <Policy "mixed_float16">
block1a_dwconv False float32 <Policy "mixed_float16">
block1a_bn False float32 <Policy "mixed_float16">
block1a_activation False float32 <Policy "mixed_float16">
block1a_se_squeeze False float32 <Policy "mixed_float16">
block1a_se_reshape False float32 <Policy "mixed_float16">
block1a_se_reduce False float32 <Policy "mixed_float16">
block1a_se_expand False float32 <Policy "mixed_float16">
block1a_se_excite False float32 <Policy "mixed_float16">
block1a_project_conv False float32 <Policy "mixed_float16">
block1a_project_bn False float32 <Policy "mixed_float16">
block2a_expand_conv False float32 <Policy "mixed_float16">
block2a_expand_bn False float32 <Policy "mixed_float16">
block2a_expand_activation False float32 <Policy "mixed_float16">
```

□ **Note:** The mixed precision API automatically causes layers which can benefit from using the "mixed_float16" dtype policy to use it. It also prevents layers which shouldn't use it from using it (e.g. the normalization layer at the start of the base model).

Fit the feature extraction model

Now that's one good looking model. Let's fit it to our data shall we?

Three epochs should be enough for our top layers to adjust their weights enough to our food image data.

To save time per epoch, we'll also only validate on 15% of the test data.


```
In [35]: # Turn off all warnings except for errors
tf.get_logger().setLevel('ERROR')

# Fit the model with callbacks
history_101_food_classes_feature_extract = model.fit(train_data,
                                                    epochs=3,
                                                    steps_per_epoch=1e
                                                    validation_data=te
                                                    validation_steps=i
                                                    callbacks=[create_

                                                    model_c
```

```
2022-09-20 10:06:15.322918: I tensorflow/core/profiler/lib/profiler_session.cc:131] Profiler session initializing.
2022-09-20 10:06:15.322948: I tensorflow/core/profiler/lib/profiler_session.cc:146] Profiler session started.
2022-09-20 10:06:15.427900: I tensorflow/core/profiler/lib/profiler_session.cc:164] Profiler session tear down.
2022-09-20 10:06:15.428018: I tensorflow/core/profiler/internal/gpu/cupti_tracer.cc:1748] CUPTI activity buffer flushed
```

Saving TensorBoard log files to: training_logs/efficientnetb0_101_classes_all_data_feature_extract/20220920-100615

Epoch 1/3

```
2/2368 [.....] - ETA: 9:47 - loss: 2.5717
- accuracy: 0.5000
```

```
2022-09-20 10:06:15.858416: I tensorflow/core/profiler/lib/profiler_session.cc:131] Profiler session initializing.
2022-09-20 10:06:15.858442: I tensorflow/core/profiler/lib/profiler_session.cc:146] Profiler session started.
2022-09-20 10:06:16.013564: I tensorflow/core/profiler/lib/profiler_session.cc:66] Profiler session collecting data.
2022-09-20 10:06:16.013961: I tensorflow/core/profiler/internal/gpu/cupti_tracer.cc:1748] CUPTI activity buffer flushed
2022-09-20 10:06:16.022859: I tensorflow/core/profiler/internal/gpu/cupti_collector.cc:673] GpuTracer has collected 551 callback api events and 550 activity events.
2022-09-20 10:06:16.030739: I tensorflow/core/profiler/lib/profiler_session.cc:164] Profiler session tear down.
2022-09-20 10:06:16.042630: I tensorflow/core/profiler/rpc/client/save_profile.cc:136] Creating directory: training_logs/efficientnetb0_101_classes_all_data_feature_extract/20220920-100615/train/plugins/profile/2022_09_20_10_06_16
```

```
2022-09-20 10:06:16.050313: I tensorflow/core/profiler/rpc/client/save_profile.cc:142] Dumped gzipped tool data for trace.json.gz to training_logs/efficientnetb0_101_classes_all_data_feature_extract/20220920-100615/train/plugins/profile/2022_09_20_10_06_16/daniel-2490-UD.trace.json.gz
```

```
9/2368 [.....] - ETA: 2:18 - loss: 2.6338
- accuracy: 0.4653
```

```
2022-09-20 10:06:16.072212: I tensorflow/core/profiler/rpc/client/save_profile.cc:136] Creating directory: training_logs/efficientnetb0_101_classes_all_data_feature_extract/20220920-100615/train/plugins/profile/2022_09_20_10_06_16
```

```
2022-09-20 10:06:16.076049: I tensorflow/core/profiler/rpc/client/save_profile.cc:142] Dumped gzipped tool data for memory_profile.json.gz to training_logs/efficientnetb0_101_classes_all_data_feature_extract/20220920-100615/train/plugins/profile/2022_09_20_10_06_16/daniel-Z490-UD.memory_profile.json.gz
```

```
2022-09-20 10:06:16.076738: I tensorflow/core/profiler/rpc/client/capture_profile.cc:251] Creating directory: training_logs/efficientnetb0_101_classes_all_data_feature_extract/20220920-100615/train/plugins/profile/2022_09_20_10_06_16
```

```
Dumped tool data for xplane.pb to training_logs/efficientnetb0_101_classes_all_data_feature_extract/20220920-100615/train/plugins/profile/2022_09_20_10_06_16/daniel-Z490-UD.xplane.pb
```

```
Dumped tool data for overview_page.pb to training_logs/efficientnetb0_101_classes_all_data_feature_extract/20220920-100615/train/plugins/profile/2022_09_20_10_06_16/daniel-Z490-UD.overview_page.pb
```

```
Dumped tool data for input_pipeline.pb to training_logs/efficientnetb0_101_classes_all_data_feature_extract/20220920-100615/train/plugins/profile/2022_09_20_10_06_16/daniel-Z490-UD.input_pipeline.pb
```

```
Dumped tool data for tensorflow_stats.pb to training_logs/efficientnetb0_101_classes_all_data_feature_extract/20220920-100615/train/plugins/profile/2022_09_20_10_06_16/daniel-Z490-UD.tensorflow_stats.pb
```

```
Dumped tool data for kernel_stats.pb to training_logs/efficientnetb0_101_classes_all_data_feature_extract/20220920-100615/train/plugins/profile/2022_09_20_10_06_16/daniel-Z490-UD.kernel_stats.pb
```

```
2368/2368 [=====] - 52s 22ms/step - loss: 1.6907 - accuracy: 0.5805 - val_loss: 1.2160 - val_accuracy: 0.6748
```

```
Epoch 2/3
```

```
2368/2368 [=====] - 51s 21ms/step - loss: 1.2817 - accuracy: 0.6685 - val_loss: 1.1228 - val_accuracy: 0.6994
```

```
Epoch 3/3
```

```
2368/2368 [=====] - 51s 21ms/step - loss: 1.1366 - accuracy: 0.7028 - val_loss: 1.0876 - val_accuracy: 0.7068
```

Nice, looks like our feature extraction model is performing pretty well. How about we evaluate it on the whole test dataset?

```
In [36]: # Evaluate model (unsaved version) on whole test dataset
results_feature_extract_model = model.evaluate(test_data)
results_feature_extract_model
```

```
790/790 [=====] - 14s 18ms/step - loss: 1.0888 - accuracy: 0.7048
```

```
Out[36]: [1.0887573957443237, 0.704752504825592]
```

And since we used the `ModelCheckpoint` callback, we've got a saved version of our model in the `model_checkpoints` directory.

Let's load it in and make sure it performs just as well

Load and evaluate checkpoint weights

We can load in and evaluate our model's checkpoints by:

1. Cloning our model using `tf.keras.models.clone_model()` (https://www.tensorflow.org/api_docs/python/tf/keras/models/clone_model) to make a copy of our feature extraction model with reset weights.
2. Calling the `load_weights()` method on our cloned model passing it the path to where our checkpointed weights are stored.
3. Calling `evaluate()` on the cloned model with loaded weights.

A reminder, checkpoints are helpful for when you perform an experiment such as fine-tuning your model. In the case you fine-tune your feature extraction model and find it doesn't offer any improvements, you can always revert back to the checkpointed version of your model.

```
In [37]: # Clone the model we created (this resets all weights)
cloned_model = tf.keras.models.clone_model(model)
cloned_model.summary()
```

Model: "model_1"

Layer (type)	Output Shape	Param #
=====		
input_layer (InputLayer)	[(None, 224, 224, 3)]	0
efficientnetb0 (Functional)	(None, None, None, 1280)	4049571
pooling_layer (GlobalAverage)	(None, 1280)	0
dense_1 (Dense)	(None, 101)	129381
softmax_float32 (Activation)	(None, 101)	0
=====		
Total params: 4,178,952		
Trainable params: 129,381		
Non-trainable params: 4,049,571		

```
In [38]: # Where are our checkpoints stored?
checkpoint_path
```

```
Out[38]: 'model_checkpoints/cp.ckpt'
```

```
In [39]: # Load checkpointed weights into cloned_model
cloned_model.load_weights(checkpoint_path)
```

```
Out[39]: <tensorflow.python.training.tracking.util.CheckpointLoadStatus at 0x7fa40c318430>
```

Each time you make a change to your model (including loading weights), you have to recompile.

```
In [41]: # Compile cloned_model (with same parameters as original model)
cloned_model.compile(loss="sparse_categorical_crossentropy",
                    optimizer=tf.keras.optimizers.Adam(),
                    metrics=["accuracy"])
```

```
In [42]: # Evaluate cloned model with loaded weights (should be same score as tr
results_cloned_model_with_loaded_weights = cloned_model.evaluate(test_d

790/790 [=====] - 15s 17ms/step - loss: 1.725
9 - accuracy: 0.5503
```

Our cloned model with loaded weight's results should be very close to the feature extraction model's results (if the cell below errors, something went wrong).

```
In [47]: # Loaded checkpoint weights should return very similar results to check
import numpy as np
assert np.isclose(results_feature_extract_model, results_cloned_model_w
```

```
-----
-----
AssertionError                                Traceback (most recent call
last)
/tmp/ipykernel_1467108/1538537382.py in <module>
      1 # Loaded checkpoint weights should return very similar results
to checkpoint weights prior to saving
      2 import numpy as np
----> 3 assert np.isclose(results_feature_extract_model, results_clone
d_model_with_loaded_weights).all(), "Loaded weights results are not cl
ose to original model." # check if all elements in array are close

AssertionError: Loaded weights results are not close to original mode
l.
```

Cloning the model preserves dtype_policy's of layers (but doesn't preserve weights) so if we wanted to continue fine-tuning with the cloned model, we could and it would still use the mixed precision dtype policy.

```
In [44]: # Check the layers in the base model and see what dtype policy they're
for layer in cloned_model.layers[1].layers[:20]: # check only the first
    print(layer.name, layer.trainable, layer.dtype, layer.dtype_policy)

input_2 True float32 <Policy "float32">
rescaling_1 False float32 <Policy "mixed_float16">
normalization_1 False float32 <Policy "mixed_float16">
stem_conv_pad False float32 <Policy "mixed_float16">
stem_conv False float32 <Policy "mixed_float16">
stem_bn False float32 <Policy "mixed_float16">
stem_activation False float32 <Policy "mixed_float16">
block1a_dwconv False float32 <Policy "mixed_float16">
block1a_bn False float32 <Policy "mixed_float16">
block1a_activation False float32 <Policy "mixed_float16">
block1a_se_squeeze False float32 <Policy "mixed_float16">
block1a_se_reshape False float32 <Policy "mixed_float16">
block1a_se_reduce False float32 <Policy "mixed_float16">
block1a_se_expand False float32 <Policy "mixed_float16">
block1a_se_excite False float32 <Policy "mixed_float16">
block1a_project_conv False float32 <Policy "mixed_float16">
block1a_project_bn False float32 <Policy "mixed_float16">
block2a_expand_conv False float32 <Policy "mixed_float16">
block2a_expand_bn False float32 <Policy "mixed_float16">
block2a_expand_activation False float32 <Policy "mixed_float16">
```

Save the whole model to file

We can also save the whole model using the `save()` (https://www.tensorflow.org/api_docs/python/tf/keras/Model#save) method.

Since our model is quite large, you might want to save it to Google Drive (if you're using Google Colab) so you can load it in for use later.

□ **Note:** Saving to Google Drive requires mounting Google Drive (go to Files -> Mount Drive).

```
In [ ]: # ## Saving model to Google Drive (optional)

# # Create save path to drive
# save_dir = "drive/MyDrive/tensorflow_course/food_vision/07_efficientnet"
# # os.makedirs(save_dir) # Make directory if it doesn't exist

# # Save model
# model.save(save_dir)
```

We can also save it directly to our Google Colab instance.

□ **Note:** Google Colab storage is ephemeral and your model will delete itself (along with any other saved files) when

```
In [48]: # Save model locally (if you're using Google Colab, your saved model will
save_dir = "07_efficientnetb0_feature_extract_model_mixed_precision"
model.save(save_dir)
```

```
2022-09-20 10:10:36.539525: W tensorflow/python/util/util.cc:348] Sets
are not currently considered sequences, but this may change in the future,
so consider avoiding using them.
```

```
/home/daniel/code/tensorflow/env/lib/python3.9/site-packages/keras/uti
ls/generic_utils.py:494: CustomMaskWarning: Custom mask layers require
a config and must override get_config. When loading, the custom mask l
ayer must be passed to the custom_objects argument.
```

```
warnings.warn('Custom mask layers require a config and must override
',
```

And again, we can check whether or not our model saved correctly by loading it in and evaluating it.

```
In [53]: # Load model previously saved above
loaded_saved_model = tf.keras.models.load_model(save_dir)
```

```
WARNING:absl:Importing a function (__inference_block2a_expand_activa
tion_layer_call_and_return_conditional_losses_65022) with ops with u
nsaved custom gradients. Will likely fail if a gradient is requeste
d.
```

```
WARNING:absl:Importing a function (__inference_block2a_se_reduce_lay
er_call_and_return_conditional_losses_65096) with ops with unsaved c
ustom gradients. Will likely fail if a gradient is requested.
```

```
WARNING:absl:Importing a function (__inference_block3a_expand_activa
tion_layer_call_and_return_conditional_losses_65348) with ops with u
nsaved custom gradients. Will likely fail if a gradient is requeste
d.
```

```
WARNING:absl:Importing a function (__inference_block7a_expand_activa
tion_layer_call_and_return_conditional_losses_67313) with ops with u
nsaved custom gradients. Will likely fail if a gradient is requeste
d.
```

```
WARNING:absl:Importing a function (__inference_block6d_expand_activa
tion_layer_call_and_return_conditional_losses_67146) with ops with u
nsaved custom gradients. Will likely fail if a gradient is requeste
d.
```

```
WARNING:absl:Importing a function (__inference_block6d_expand_activa
tion_layer_call_and_return_conditional_losses_67146) with ops with u
nsaved custom gradients. Will likely fail if a gradient is requeste
d.
```

Loading a SavedModel also retains all of the underlying layers dtype_policy (we want them to be "mixed_float16").

```
In [54]: # Check the layers in the base model and see what dtype policy they're
for layer in loaded_saved_model.layers[1].layers[:20]: # check only the
    print(layer.name, layer.trainable, layer.dtype, layer.dtype_policy)
```

```
input_2 True float32 <Policy "float32">
rescaling_1 False float32 <Policy "mixed_float16">
normalization_1 False float32 <Policy "mixed_float16">
stem_conv_pad False float32 <Policy "mixed_float16">
stem_conv False float32 <Policy "mixed_float16">
stem_bn False float32 <Policy "mixed_float16">
stem_activation False float32 <Policy "mixed_float16">
block1a_dwconv False float32 <Policy "mixed_float16">
block1a_bn False float32 <Policy "mixed_float16">
block1a_activation False float32 <Policy "mixed_float16">
block1a_se_squeeze False float32 <Policy "mixed_float16">
block1a_se_reshape False float32 <Policy "mixed_float16">
block1a_se_reduce False float32 <Policy "mixed_float16">
block1a_se_expand False float32 <Policy "mixed_float16">
block1a_se_excite False float32 <Policy "mixed_float16">
block1a_project_conv False float32 <Policy "mixed_float16">
block1a_project_bn False float32 <Policy "mixed_float16">
block2a_expand_conv False float32 <Policy "mixed_float16">
block2a_expand_bn False float32 <Policy "mixed_float16">
block2a_expand_activation False float32 <Policy "mixed_float16">
```

```
In [55]: # Check loaded model performance (this should be the same as results_fe
results_loaded_saved_model = loaded_saved_model.evaluate(test_data)
results_loaded_saved_model
```

```
790/790 [=====] - 15s 18ms/step - loss: 1.088
8 - accuracy: 0.7048
```

```
Out[55]: [1.0887584686279297, 0.704752504825592]
```

```
In [56]: # The loaded model's results should equal (or at least be very close) t
# Note: this will only work if you've instatiated results variables
import numpy as np
assert np.isclose(results_feature_extract_model, results_loaded_saved_m
```

That's what we want! Our loaded model performing as it should.

□ **Note:** We spent a fair bit of time making sure our model saved correctly because training on a lot of data can be time-consuming, so we want to make sure we don't have to continually train from scratch.

Preparing our model's layers for fine-tuning

Our feature-extraction model is showing some great promise after three epochs. But since we've got so much data, it's probably worthwhile that we see what results we can get with fine-tuning (fine-tuning usually works best when you've got quite a large amount of data).

Remember our goal of beating the [DeepFood paper](https://arxiv.org/pdf/1606.05675.pdf) (<https://arxiv.org/pdf/1606.05675.pdf>)?

They were able to achieve 77.4% top-1 accuracy on Food101 over 2-3 days of training.

Do you think fine-tuning will get us there?

Let's find out.

To start, let's load in our saved model.

□ **Note:** It's worth remembering a traditional workflow for fine-tuning is to freeze a pre-trained base model and then train only the output layers for a few iterations so their weights can be updated inline with your custom data (feature extraction). And then unfreeze a number or all of the layers in the base model and continue training until the model stops improving.

Like all good cooking shows, I've saved a model I prepared earlier (the feature extraction model from above) to Google Storage.

We can download it to make sure we're using the same model going forward.


```
In [57]: # Download the saved model from Google Storage
!wget https://storage.googleapis.com/ztm_tf_course/food_vision/07_efficientnetb0_feature_extract_model_mixed_precision.zip

--2022-09-20 10:11:48-- https://storage.googleapis.com/ztm_tf_course/food_vision/07_efficientnetb0_feature_extract_model_mixed_precision.zip
Resolving storage.googleapis.com (storage.googleapis.com)... 142.250.76.112, 142.250.204.16, 172.217.167.80, ...
Connecting to storage.googleapis.com (storage.googleapis.com)|142.250.76.112|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 16976857 (16M) [application/zip]
Saving to: '07_efficientnetb0_feature_extract_model_mixed_precision.zip.1'

07_efficientnetb0_f 100%[=====>] 16.19M 11.7MB/s in 1.4s

2022-09-20 10:11:51 (11.7 MB/s) - '07_efficientnetb0_feature_extract_model_mixed_precision.zip.1' saved [16976857/16976857]
```

```
In [58]: # Unzip the SavedModel downloaded from Google Stroage
!mkdir downloaded_gs_model # create new dir to store downloaded feature
!unzip 07_efficientnetb0_feature_extract_model_mixed_precision.zip -d d

mkdir: cannot create directory 'downloaded_gs_model': File exists
Archive: 07_efficientnetb0_feature_extract_model_mixed_precision.zip
replace downloaded_gs_model/07_efficientnetb0_feature_extract_model_mixed_precision/variables/variables.data-00000-of-00001? [y]es, [n]o, [A]ll, [N]one, [r]ename: ^C
```

```
In [59]: # Load and evaluate downloaded GS model
loaded_gs_model = tf.keras.models.load_model("downloaded_gs_model/07_efficientnetb0.h5")

WARNING:absl:Importing a function (__inference_block1a_activation_layer_call_and_return_conditional_losses_158253) with ops with unsaved custom gradients. Will likely fail if a gradient is requested.
WARNING:absl:Importing a function (__inference_block2a_activation_layer_call_and_return_conditional_losses_191539) with ops with unsaved custom gradients. Will likely fail if a gradient is requested.
WARNING:absl:Importing a function (__inference_block6d_expand_activation_layer_call_and_return_conditional_losses_196076) with ops with unsaved custom gradients. Will likely fail if a gradient is requested.
WARNING:absl:Importing a function (__inference_block6c_activation_layer_call_and_return_conditional_losses_195780) with ops with unsaved custom gradients. Will likely fail if a gradient is requested.
WARNING:absl:Importing a function (__inference_block6d_activation_layer_call_and_return_conditional_losses_196153) with ops with unsaved custom gradients. Will likely fail if a gradient is requested.
WARNING:absl:Importing a function (__inference_model_layer_call_and_return_conditional_losses_180010) with ops with unsaved custom gradients. Will likely fail if a gradient is requested.
```

```
In [60]: # Get a summary of our downloaded model
loaded_gs_model.summary()
```

Model: "model"

Layer (type)	Output Shape	Param #
=====		
input_layer (InputLayer)	[(None, 224, 224, 3)]	0

efficientnetb0 (Functional)	(None, None, None, 1280)	4049571

pooling_layer (GlobalAverage)	(None, 1280)	0

dense (Dense)	(None, 101)	129381

softmax_float32 (Activation)	(None, 101)	0
=====		
Total params: 4,178,952		
Trainable params: 129,381		
Non-trainable params: 4,049,571		

And now let's make sure our loaded model is performing as expected.

```
In [61]: # How does the loaded model perform?
results_loaded_gs_model = loaded_gs_model.evaluate(test_data)
results_loaded_gs_model

790/790 [=====] - 15s 18ms/step - loss: 1.088
1 - accuracy: 0.7065
```

```
Out[61]: [1.0881085395812988, 0.7064950466156006]
```

Great, our loaded model is performing as expected.

When we first created our model, we froze all of the layers in the base model by setting `base_model.trainable=False` but since we've loaded in our model from file, let's check whether or not the layers are trainable or not.

```
In [63]: # Are any of the layers in our model frozen?
for layer in loaded_gs_model.layers:
    layer.trainable = True # set all layers to trainable
    print(layer.name, layer.trainable, layer.dtype, layer.dtype_policy)

input_layer True float32 <Policy "float32">
efficientnetb0 True float32 <Policy "mixed_float16">
pooling_layer True float32 <Policy "mixed_float16">
dense True float32 <Policy "mixed_float16">
softmax_float32 True float32 <Policy "float32">
```

Alright, it seems like each layer in our loaded model is trainable. But what if we got a little deeper and inspected each of the layers in our base model?

□ **Question:** Which layer in the loaded model is our base model?

Before saving the Functional model to file, we created it with five layers (layers below are 0-indexed): 0. The input layer

1. The pre-trained base model layer
(`tf.keras.applications.EfficientNetB0`)
2. The pooling layer
3. The fully-connected (dense) layer
4. The output softmax activation (with float32 dtype)

Therefore to inspect our base model layer, we can access the `layers` attribute of the layer at index 1 in our model.

```
In [64]: # Check the layers in the base model and see what dtype policy they're
for layer in loaded_gs_model.layers[1].layers[:20]:
    print(layer.name, layer.trainable, layer.dtype, layer.dtype_policy)
```

```
input_1 True float32 <Policy "float32">
rescaling True float32 <Policy "mixed_float16">
normalization True float32 <Policy "float32">
stem_conv_pad True float32 <Policy "mixed_float16">
stem_conv True float32 <Policy "mixed_float16">
stem_bn True float32 <Policy "mixed_float16">
stem_activation True float32 <Policy "mixed_float16">
block1a_dwconv True float32 <Policy "mixed_float16">
block1a_bn True float32 <Policy "mixed_float16">
block1a_activation True float32 <Policy "mixed_float16">
block1a_se_squeeze True float32 <Policy "mixed_float16">
block1a_se_reshape True float32 <Policy "mixed_float16">
block1a_se_reduce True float32 <Policy "mixed_float16">
block1a_se_expand True float32 <Policy "mixed_float16">
block1a_se_excite True float32 <Policy "mixed_float16">
block1a_project_conv True float32 <Policy "mixed_float16">
block1a_project_bn True float32 <Policy "mixed_float16">
block2a_expand_conv True float32 <Policy "mixed_float16">
block2a_expand_bn True float32 <Policy "mixed_float16">
block2a_expand_activation True float32 <Policy "mixed_float16">
```

Wonderful, it looks like each layer in our base model is trainable (unfrozen) and every layer which should be using the dtype policy "mixed_policy16" is using it.

Since we've got so much data (750 images x 101 training classes = 75750 training images), let's keep all of our base model's layers unfrozen.

□ **Note:** If you've got a small amount of data (less than 100 images per class), you may want to only unfreeze and fine-tune a small number of layers in the base model at a time. Otherwise, you risk overfitting.

A couple more callbacks

We're about to start fine-tuning a deep learning model with over 200 layers using over 100,000 (75k+ training, 25K+ testing) images, which means our model's training time is probably going to be much longer than before.

□ **Question:** *How long does training take?*

It could be a couple of hours or in the case of the [DeepFood paper \(https://arxiv.org/pdf/1606.05675.pdf\)](https://arxiv.org/pdf/1606.05675.pdf) (the baseline we're trying to beat), their best performing model took 2-3 days of training time.

You will really only know how long it'll take once you start training.

□ **Question:** *When do you stop training?*

Ideally, when your model stops improving. But again, due to the nature of deep learning, it can be hard to know when exactly a model will stop improving.

Luckily, there's a solution: the [EarlyStopping callback \(https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/EarlyStopping\)](https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/EarlyStopping)

The `EarlyStopping` callback monitors a specified model performance metric (e.g. `val_loss`) and when it stops improving for a specified number of epochs, automatically stops training.

Using the `EarlyStopping` callback combined with the `ModelCheckpoint` callback saving the best performing model automatically, we could keep our model training for an unlimited number of epochs until it stops improving.

```
In [65]: # Setup EarlyStopping callback to stop training if model's val_loss does not improve
early_stopping = tf.keras.callbacks.EarlyStopping(monitor="val_loss", # if val loss does not improve
                                                    patience=3) # if val loss does not improve for 3 epochs

# Create ModelCheckpoint callback to save best model during fine-tuning
checkpoint_path = "fine_tune_checkpoints/"
model_checkpoint = tf.keras.callbacks.ModelCheckpoint(checkpoint_path,
                                                        save_best_only=True,
                                                        monitor="val_loss")
```

Woohoo! Fine-tuning callbacks ready.

If you're planning on training large models, the `ModelCheckpoint` and `EarlyStopping` are two callbacks you'll want to become very familiar with.

We're almost ready to start fine-tuning our model but there's one more callback we're going to implement: [ReduceLROnPlateau \(https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/ReduceLROnPlateau\)](https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/ReduceLROnPlateau)

Remember how the learning rate is the most important model hyperparameter you can tune? (if not, treat this as a reminder).

Well, the `ReduceLROnPlateau` callback helps to tune the learning rate for you.

Like the `ModelCheckpoint` and `EarlyStopping` callbacks, the `ReduceLROnPlateau` callback monitors a specified metric and when that metric stops improving, it reduces the learning rate by a specified factor (e.g. divides the learning rate by 10).

□ **Question:** *Why lower the learning rate?*

Imagine having a coin at the back of the couch and you're trying to grab with your fingers.

Now think of the learning rate as the size of the movements your hand makes towards the coin.

The closer you get, the smaller you want your hand movements to be, otherwise the coin will be lost.

Our model's ideal performance is the equivalent of grabbing the coin. So as training goes on and our model gets closer and closer to its ideal performance (also called **convergence**), we want the amount it learns to be less and less.

To do this we'll create an instance of the `ReduceLROnPlateau` callback to monitor the validation loss just like the `EarlyStopping` callback.

Once the validation loss stops improving for two or more epochs, we'll reduce the learning rate by a factor of 5 (e.g. 0.001 to 0.0002).

And to make sure the learning rate doesn't get too low (and potentially result in our model learning nothing), we'll set the minimum learning rate to $1e-7$.

```
In [66]: # Creating learning rate reduction callback
reduce_lr = tf.keras.callbacks.ReduceLROnPlateau(monitor="val_loss",
                                                  factor=0.2, # multiply
                                                  patience=2,
                                                  verbose=1, # print out
                                                  min_lr=1e-7)
```

Learning rate reduction ready to go!

Now before we start training, we've got to recompile our model.

We'll use sparse categorical crossentropy as the loss and since we're fine-tuning, we'll use a 10x lower learning rate than the Adam optimizers default ($1e-4$ instead of $1e-3$).

```
In [67]: # Compile the model
loaded_gs_model.compile(loss="sparse_categorical_crossentropy", # sparse categorical crossentropy
                        optimizer=tf.keras.optimizers.Adam(0.0001), # Adam optimizer with learning rate 0.0001
                        metrics=["accuracy"])
```

Okay, model compiled.

Now let's fit it on all of the data.

We'll set it up to run for up to 100 epochs.

Since we're going to be using the `EarlyStopping` callback, it might stop before reaching 100 epochs.

□ **Note:** Running the cell below will set the model up to fine-tune all of the pre-trained weights in the base model on all of the Food101 data. Doing so with **unoptimized** data pipelines and **without** mixed precision training will take a fairly long time per epoch depending on what type of GPU you're using (about 15-20 minutes on Colab GPUs). But don't worry, **the code we've written above will ensure it runs much faster** (more like 4-5 minutes per epoch).

```
In [68]: # Start to fine-tune (all layers)
history_101_food_classes_all_data_fine_tune = loaded_gs_model.fit(train
                                                                    epochs=100, # t
                                                                    steps_per_epoch
                                                                    validation_data
                                                                    validation_step
                                                                    callbacks=[crea
                                                                    mode
                                                                    earl
                                                                    redu
```

```
2022-09-20 10:14:32.462381: I tensorflow/core/profiler/lib/profiler_session.cc:131] Profiler session initializing.
2022-09-20 10:14:32.462399: I tensorflow/core/profiler/lib/profiler_session.cc:146] Profiler session started.
2022-09-20 10:14:32.640829: I tensorflow/core/profiler/lib/profiler_session.cc:164] Profiler session tear down.
2022-09-20 10:14:32.640979: I tensorflow/core/profiler/internal/gpu/cupti_tracer.cc:1748] CUPTI activity buffer flushed
```

Saving TensorBoard log files to: training_logs/efficientb0_101_classes_all_data_fine_tuning/20220920-101432

Epoch 1/100

```
1/2368 [.....] - ETA: 4:36:13 - loss: 1.1109 - accuracy: 0.7188
```

```
2022-09-20 10:14:39.914227: I tensorflow/core/profiler/lib/profiler_session.cc:131] Profiler session initializing.
```

```
2022-09-20 10:14:39.914257: I tensorflow/core/profiler/lib/profiler_session.cc:146] Profiler session started.
```

```
2/2368 [.....] - ETA: 18:36 - loss: 1.1839 - accuracy: 0.6719
```

```
2022-09-20 10:14:41.065837: I tensorflow/core/profiler/lib/profiler_session.cc:66] Profiler session collecting data.
```

```
2022-09-20 10:14:41.070394: I tensorflow/core/profiler/internal/gpu/cupti_tracer.cc:1748] CUPTI activity buffer flushed
```

```
2022-09-20 10:14:41.114372: I tensorflow/core/profiler/internal/gpu/cupti_collector.cc:673] GpuTracer has collected 2751 callback api events and 2762 activity events.
```

```
2022-09-20 10:14:41.181741: I tensorflow/core/profiler/lib/profiler_session.cc:164] Profiler session tear down.
```

```
2022-09-20 10:14:41.251293: I tensorflow/core/profiler/rpc/client/save_profile.cc:136] Creating directory: training_logs/efficientb0_101_classes_all_data_fine_tuning/20220920-101432/train/plugins/profile/2022_09_20_10_14_41
```

```
3/2368 [.....] - ETA: 33:56 - loss: 1.2313 - accuracy: 0.7083
```


2022-09-20 10:14:41.284385: I tensorflow/core/profiler/rpc/client/save_profile.cc:142] Dumped gzipped tool data for trace.json.gz to training_logs/efficientb0_101_classes_all_data_fine_tuning/20220920-101432/train/plugins/profile/2022_09_20_10_14_41/daniel-Z490-UD.trace.json.gz
2022-09-20 10:14:41.369966: I tensorflow/core/profiler/rpc/client/save_profile.cc:136] Creating directory: training_logs/efficientb0_101_classes_all_data_fine_tuning/20220920-101432/train/plugins/profile/2022_09_20_10_14_41

2022-09-20 10:14:41.375375: I tensorflow/core/profiler/rpc/client/save_profile.cc:142] Dumped gzipped tool data for memory_profile.json.gz to training_logs/efficientb0_101_classes_all_data_fine_tuning/20220920-101432/train/plugins/profile/2022_09_20_10_14_41/daniel-Z490-UD.memory_profile.json.gz

2022-09-20 10:14:41.379940: I tensorflow/core/profiler/rpc/client/capture_profile.cc:251] Creating directory: training_logs/efficientb0_101_classes_all_data_fine_tuning/20220920-101432/train/plugins/profile/2022_09_20_10_14_41

Dumped tool data for xplane.pb to training_logs/efficientb0_101_classes_all_data_fine_tuning/20220920-101432/train/plugins/profile/2022_09_20_10_14_41/daniel-Z490-UD.xplane.pb

Dumped tool data for overview_page.pb to training_logs/efficientb0_101_classes_all_data_fine_tuning/20220920-101432/train/plugins/profile/2022_09_20_10_14_41/daniel-Z490-UD.overview_page.pb

Dumped tool data for input_pipeline.pb to training_logs/efficientb0_101_classes_all_data_fine_tuning/20220920-101432/train/plugins/profile/2022_09_20_10_14_41/daniel-Z490-UD.input_pipeline.pb

Dumped tool data for tensorflow_stats.pb to training_logs/efficientb0_101_classes_all_data_fine_tuning/20220920-101432/train/plugins/profile/2022_09_20_10_14_41/daniel-Z490-UD.tensorflow_stats.pb

Dumped tool data for kernel_stats.pb to training_logs/efficientb0_101_classes_all_data_fine_tuning/20220920-101432/train/plugins/profile/2022_09_20_10_14_41/daniel-Z490-UD.kernel_stats.pb

2368/2368 [=====] - 143s 57ms/step - loss: 0.9223 - accuracy: 0.7527 - val_loss: 0.7910 - val_accuracy: 0.7783

/home/daniel/code/tensorflow/env/lib/python3.9/site-packages/keras/utils/generic_utils.py:494: CustomMaskWarning: Custom mask layers require a config and must override get_config. When loading, the custom mask layer must be passed to the custom_objects argument.

warnings.warn('Custom mask layers require a config and must override',

Epoch 2/100

2368/2368 [=====] - 134s 56ms/step - loss: 0.5792 - accuracy: 0.8375 - val_loss: 0.7917 - val_accuracy: 0.7868

Epoch 3/100

2368/2368 [=====] - 134s 57ms/step - loss: 0.3319 - accuracy: 0.9058 - val_loss: 0.8618 - val_accuracy: 0.7762

Epoch 00003: ReduceLROnPlateau reducing learning rate to 1.9999999494757503e-05.

Epoch 4/100

2368/2368 [=====] - 134s 56ms/step - loss: 0.0854 - accuracy: 0.9793 - val_loss: 0.9299 - val_accuracy: 0.8006

□ **Note:** If you didn't use mixed precision or use techniques such as `prefetch()` (https://www.tensorflow.org/api_docs/python/tf/data/Dataset#prefetch) in the *Batch & prepare datasets* section, your model fine-tuning probably takes up to 2.5-3x longer per epoch (see the output below for an example).

	Prefetch and mixed precision	No prefetch and no mixed precision
--	------------------------------	------------------------------------

Time per epoch	~280-300s	~1127-1397s
----------------	-----------	-------------

Results from fine-tuning 🍷 Food Vision Big™ on Food101 dataset using an EfficientNetNetB0 backbone using a Google Colab Tesla T4 GPU.

```
Saving TensorBoard log files to: training_logs/efficientB0_101_
classes_all_data_fine_tuning/20200928-013008
```

```
Epoch 1/100
```

```
2368/2368 [=====] - 1397s 590ms/step -
loss: 1.2068 - accuracy: 0.6820 - val_loss: 1.1623 - val_accu
cy: 0.6894
```

```
Epoch 2/100
```

```
2368/2368 [=====] - 1193s 504ms/step -
loss: 0.9459 - accuracy: 0.7444 - val_loss: 1.1549 - val_accu
cy: 0.6872
```

```
Epoch 3/100
```

```
2368/2368 [=====] - 1143s 482ms/step -
loss: 0.7848 - accuracy: 0.7838 - val_loss: 1.0402 - val_accu
cy: 0.7142
```

```
Epoch 4/100
```

```
2368/2368 [=====] - 1127s 476ms/step -
loss: 0.6599 - accuracy: 0.8149 - val_loss: 0.9599 - val_accu
cy: 0.7373
```

Example fine-tuning time for non-prefetched data as well as non-mixed precision training (~2.5-3x longer per epoch).

Let's make sure we save our model before we start evaluating it.

```
In [69]: ## Save model to Google Drive (optional)
# loaded_gs_model.save("/content/drive/MyDrive/tensorflow_course/food_v
```

```
In [70]: # Save model locally (note: if you're using Google Colab and you save y
loaded_gs_model.save("07_efficientnetb0_fine_tuned_101_classes_mixed_pr
```

```
/home/daniel/code/tensorflow/env/lib/python3.9/site-packages/keras/uti
ls/generic_utils.py:494: CustomMaskWarning: Custom mask layers require
a config and must override get_config. When loading, the custom mask l
ayer must be passed to the custom_objects argument.
  warnings.warn('Custom mask layers require a config and must override
,
```

Looks like our model has gained a few performance points from fine-tuning, let's evaluate on the whole test dataset and see if managed to beat the [DeepFood paper's](https://arxiv.org/abs/1606.05675) (<https://arxiv.org/abs/1606.05675>) result of 77.4% accuracy.

Woohoo!!!! It looks like our model beat the results mentioned in the DeepFood paper for Food101 (DeepFood's 77.4% top-1 accuracy versus our ~79% top-1 accuracy).

Download fine-tuned model from Google Storage

As mentioned before, training models can take a significant amount of time.

And again, like any good cooking show, here's something we prepared earlier...

It's a fine-tuned model exactly like the one we trained above but it's saved to Google Storage so it can be accessed, imported and evaluated.

```
In [71]: # Download and evaluate fine-tuned model from Google Storage
!wget https://storage.googleapis.com/ztm_tf_course/food_vision/07_efficientnetb0_fine_tuned_101_classes_mixed_precision.zip

--2022-09-20 10:24:16-- https://storage.googleapis.com/ztm_tf_course/food_vision/07_efficientnetb0_fine_tuned_101_classes_mixed_precision.zip
Resolving storage.googleapis.com (storage.googleapis.com)... 172.217.167.80, 142.250.66.240, 142.250.67.16, ...
Connecting to storage.googleapis.com (storage.googleapis.com)|172.217.167.80|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 46790356 (45M) [application/zip]
Saving to: '07_efficientnetb0_fine_tuned_101_classes_mixed_precision.zip.1'

07_efficientnetb0_f 100%[=====>] 44.62M 13.0MB/s in 3.7s

2022-09-20 10:24:21 (12.2 MB/s) - '07_efficientnetb0_fine_tuned_101_classes_mixed_precision.zip.1' saved [46790356/46790356]
```

The downloaded model comes in zip format (.zip) so we'll unzip it into the Google Colab instance.

```
In [72]: # Unzip fine-tuned model
!mkdir downloaded_fine_tuned_gs_model # create separate directory for t
!unzip 07_efficientnetb0_fine_tuned_101_classes_mixed_precision -d down

mkdir: cannot create directory 'downloaded_fine_tuned_gs_model': File exists
unzip: cannot find or open /content/07_efficientnetb0_fine_tuned_101_classes_mixed_precision, /content/07_efficientnetb0_fine_tuned_101_classes_mixed_precision.zip or /content/07_efficientnetb0_fine_tuned_101_classes_mixed_precision.ZIP.
```

Now we can load it using the `tf.keras.models.load_model()` (https://www.tensorflow.org/tutorials/keras/save_and_load) method and get a summary (it should be the exact same as the model we created above).

```
In [73]: # Load in fine-tuned model from Google Storage and evaluate
loaded_fine_tuned_gs_model = tf.keras.models.load_model("downloaded_fin
```

```
WARNING:absl:Importing a function (__inference_block3b_expand_activation_layer_call_and_return_conditional_losses_443625) with ops with unsaved custom gradients. Will likely fail if a gradient is requested.
WARNING:absl:Importing a function (__inference_block5c_activation_layer_call_and_return_conditional_losses_412189) with ops with unsaved custom gradients. Will likely fail if a gradient is requested.
WARNING:absl:Importing a function (__inference_block1a_se_reduce_layer_call_and_return_conditional_losses_409120) with ops with unsaved custom gradients. Will likely fail if a gradient is requested.
WARNING:absl:Importing a function (__inference_block6c_expand_activation_layer_call_and_return_conditional_losses_446895) with ops with unsaved custom gradients. Will likely fail if a gradient is requested.
WARNING:absl:Importing a function (__inference_block1a_activation_layer_call_and_return_conditional_losses_442329) with ops with unsaved custom gradients. Will likely fail if a gradient is requested.
WARNING:absl:Importing a function (__inference_efficientnetb0_layer_call_and_return_conditional_losses_421687) with ops with unsaved cus
```

```
In [74]: # Get a model summary (same model architecture as above)
loaded_fine_tuned_gs_model.summary()
```

Model: "model"

Layer (type)	Output Shape	Param #
input_layer (InputLayer)	[(None, 224, 224, 3)]	0
efficientnetb0 (Functional)	(None, None, None, 1280)	4049571
pooling_layer (GlobalAverage)	(None, 1280)	0
dense (Dense)	(None, 101)	129381
softmax_float32 (Activation)	(None, 101)	0

Total params: 4,178,952
Trainable params: 4,136,929
Non-trainable params: 42,023

Finally, we can evaluate our model on the test data (this requires the `test_data` variable to be loaded).

```
In [75]: # Note: Even if you're loading in the model from Google Storage, you will
results_downloaded_fine_tuned_gs_model = loaded_fine_tuned_gs_model.evaluate
results_downloaded_fine_tuned_gs_model

790/790 [=====] - 16s 19ms/step - loss: 0.907
2 - accuracy: 0.8016
```

```
Out[75]: [0.9072349071502686, 0.8015841841697693]
```

Excellent! Our saved model is performing as expected (better results than the DeepFood paper!).

Congratulations! You should be excited! You just trained a computer vision model with competitive performance to a research paper and in far less time (our model took ~20 minutes to train versus DeepFood's quoted 2-3 days).

In other words, you brought Food Vision life!

If you really wanted to step things up, you could try using the [EfficientNetB4](https://www.tensorflow.org/api_docs/python/tf/keras/applications/EfficientNetB4) (https://www.tensorflow.org/api_docs/python/tf/keras/applications/EfficientNetB4) model (a larger version of EfficientNetB0). At the time of writing, the EfficientNet family has the [state of the art classification results](https://paperswithcode.com/sota/fine-grained-image-classification-on-food-101) (<https://paperswithcode.com/sota/fine-grained-image-classification-on-food-101>) on the Food101 dataset.

□ **Resource:** To see which models are currently performing the best on a given dataset or problem type as well as the latest trending machine learning research, be sure to check out paperswithcode.com (<http://paperswithcode.com/>) and sotabench.com (<https://sotabench.com/>).

View training results on TensorBoard

Since we tracked our model's fine-tuning training logs using the TensorBoard callback, let's upload them and inspect them on TensorBoard.dev.

```
In [76]: # Upload experiment results to TensorBoard (uncomment to run)
# !tensorboard dev upload --logdir ./training_logs \
#     --name "Fine-tuning EfficientNetB0 on all Food101 Data" \
#     --description "Training results for fine-tuning EfficientNetB0 on Food101" \
#     --one_shot
```

Viewing at our [model's training curves on TensorBoard.dev](https://tensorboard.dev/experiment/2KINdYxgSgW2bUg7dIvevw/) (<https://tensorboard.dev/experiment/2KINdYxgSgW2bUg7dIvevw/>), it looks like our fine-tuning model gains boost in performance but starts to overfit as training goes on.

See the training curves on TensorBoard.dev here:

<https://tensorboard.dev/experiment/2KINdYxgSgW2bUg7dIvevw/>
(<https://tensorboard.dev/experiment/2KINdYxgSgW2bUg7dIvevw/>)

To fix this, in future experiments, we might try things like:

- A different iteration of EfficientNet (e.g. EfficientNetB4 instead of EfficientNetB0).
- Unfreezing less layers of the base model and training them rather than unfreezing the whole base model in one go.

```
In [78]: # View past TensorBoard experiments
# !tensorboard dev list
```

```
In [79]: # Delete past TensorBoard experiments
# !tensorboard dev delete --experiment_id YOUR_EXPERIMENT_ID

# Example
# !tensorboard dev delete --experiment_id OAE6KXizQZKQxDiqI3cnUQ
```

✂ Exercises

1. Use the same evaluation techniques on the large-scale Food Vision model as you did in the previous notebook ([Transfer Learning Part 3: Scaling up](https://github.com/mrdbourke/tensorflow-deep-learning/blob/main/06_transfer_learning_in_tensorflow_part_3_scaling) (https://github.com/mrdbourke/tensorflow-deep-learning/blob/main/06_transfer_learning_in_tensorflow_part_3_scaling)). More specifically, it would be good to see:
 - A confusion matrix between all of the model's predictions and true labels.
 - A graph showing the f1-scores of each class.
 - A visualization of the model making predictions on various images and comparing the predictions to the ground truth.
 - For example, plot a sample image from the test dataset and have the title of the plot show the prediction, the prediction probability and the ground truth label.
2. Take 3 of your own photos of food and use the Food Vision model to make predictions on them. How does it go? Share your images/predictions with the other students.
3. Retrain the model (feature extraction and fine-tuning) we trained in this notebook, except this time use [EfficientNetB4](https://www.tensorflow.org/api_docs/python/tf/keras/applications/EfficientNetB4) (https://www.tensorflow.org/api_docs/python/tf/keras/applications/EfficientNetB4) as the base model instead of EfficientNetB0. Do you notice an

improvement in performance? Does it take longer to train? Are there any tradeoffs to consider?

4. Name one important benefit of mixed precision training, how does this benefit take place?

□ Extra-curriculum

- Read up on learning rate scheduling and the [learning_rate scheduler callback](https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/LearningRateScheduler) (https://www.tensorflow.org/api_docs/python/tf/keras/callbacks/LearningRateScheduler). What is it? And how might it be helpful to this project?
- Read up on TensorFlow data loaders ([improving TensorFlow data loading performance](https://www.tensorflow.org/guide/data_performance) (https://www.tensorflow.org/guide/data_performance)). Is there anything we've missed? What methods you keep in mind whenever loading data in TensorFlow? Hint: check the summary at the bottom of the page for a great round up of ideas.
- Read up on the documentation for [TensorFlow mixed precision training](https://www.tensorflow.org/guide/mixed_precision) (https://www.tensorflow.org/guide/mixed_precision). What are the important things to keep in mind when using mixed precision training?

