

ASR Automate Speech Recognition

- STT (Speech to Text)
- First, speech recognition that allows the machine to catch the words, phrases and sentences we speak
- Second, natural language processing to allow the machine to understand what we speak, and
- Third, speech synthesis to allow the machine to speak.

Difficulties in developing a speech recognition system

- **Size of the vocabulary** – Size of the vocabulary impacts the ease of developing an ASR. Consider the following sizes of vocabulary for a better understanding.
 - A small size vocabulary consists of 2-100 words, for example, as in a voice-menu system
 - A medium size vocabulary consists of several 100s to 1,000s of words, for example, as in a database-retrieval task
 - A large size vocabulary consists of several 10,000s of words, as in a general dictation task.
- **Channel characteristics** – Channel quality is also an important dimension. For example, human speech contains high bandwidth with full frequency range, while a telephone speech consists of low bandwidth with limited frequency range. Note that it is harder in the latter.
- **Speaking mode**– Ease of developing an ASR also depends on the speaking mode, that is whether the speech is in isolated word mode, or connected word mode, or in a continuous speech mode. Note that a continuous speech is harder to recognize.
- **Speaking style** – A read speech may be in a formal style, or spontaneous and conversational with casual style. The latter is harder to recognize.
- **Speaker dependency**– Speech can be speaker dependent, speaker adaptive, or speaker independent. A speaker independent is the hardest to build.
- **Type of noise** bold text – Noise is another factor to consider while developing an ASR. Signal to noise ratio may be in various ranges, depending on the acoustic environment that observes less versus more background noise
 - If the signal to noise ratio is greater than 30dB, it is considered as high range
 - If the signal to noise ratio lies between 30dB to 10db, it is considered as medium SNR
 - If the signal to noise ratio is lesser than 10dB, it is considered as low range
- **Microphone characteristics** – The quality of microphone may be good, average, or below average. Also, the distance between mouth and micro-phone can vary. These factors also should be considered for recognition systems.

Simple audio recognition: Recognizing keywords

DataSet:

speech_commands

An audio dataset of spoken words designed to help train and evaluate keyword spotting systems. Its primary goal is to provide a way to build and test small models that detect when a single word is spoken, from a set of ten target words, with as few false positives as possible from background noise or unrelated speech. Note that in the train and validation set, the label "unknown" is much more prevalent than the labels of the target words or background noise. One difference from the release version is the handling of silent segments. While in the test set the silence segments are regular 1 second files, in the training they are provided as long segments under "background_noise" folder. Here we split these background noise into 1 second clips, and also keep one of the files for the validation set. [link \(https://www.tensorflow.org/datasets/catalog/speech_commands\)](https://www.tensorflow.org/datasets/catalog/speech_commands)

```
In [ ]: 1 !pip install -U -q tensorflow tensorflow_datasets
2 !apt install --allow-change-held-packages libcudnn8=8.1.0.77-1+cuda11.2
```

5.4/5.4 MB 28.2 MB/s eta 0:00:00
3.0/3.0 MB 52.5 MB/s eta 0:00:00

Reading package lists... Done
Building dependency tree
Reading state information... Done
The following packages will be REMOVED:
libcudnn8-dev
The following held packages will be changed:
libcudnn8
The following packages will be DOWNGRADED:
libcudnn8
0 upgraded, 0 newly installed, 1 downgraded, 1 to remove and 22 not upgraded.
Need to get 430 MB of archives.
After this operation, 1,153 MB disk space will be freed.
Get:1 https://developer.download.nvidia.com/compute/cuda/repos/ubuntu2004/x86_64 libcudnn8 8.1.0.77-1+cuda11.2 [430 MB]
Fetched 430 MB in 6s (75.1 MB/s)
(Reading database ... 122518 files and directories currently installed.)
Removing libcudnn8-dev (8.7.0.84-1+cuda11.8) ...
update-alternatives: removing manually selected alternative - switching libcudnn to auto mode
dpkg: **warning:** downgrading libcudnn8 from 8.7.0.84-1+cuda11.8 to 8.1.0.77-1+cuda11.2
(Reading database ... 122485 files and directories currently installed.)
Preparing to unpack .../libcudnn8_8.1.0.77-1+cuda11.2_amd64.deb ...
Unpacking libcudnn8 (8.1.0.77-1+cuda11.2) over (8.7.0.84-1+cuda11.8) ...
Setting up libcudnn8 (8.1.0.77-1+cuda11.2) ...

```
In [ ]: 1 import os
2 import pathlib
3
4 import matplotlib.pyplot as plt
5 %matplotlib inline
6 import numpy as np
7 import seaborn as sns
8 import tensorflow as tf
9
10 from tensorflow.keras import layers
11 from tensorflow.keras import models
12 from IPython import display
13
14 # Set the seed value for experiment reproducibility.
15 seed = 42
16 tf.random.set_seed(seed)
17 np.random.seed(seed)
```

```
In [ ]: 1 DATASET_PATH = 'data/mini_speech_commands'
2
3 data_dir = pathlib.Path(DATASET_PATH)
4 if not data_dir.exists():
5     tf.keras.utils.get_file(
6         'mini_speech_commands.zip',
7         origin='http://storage.googleapis.com/download.tensorflow.org/data/mini_speech_commands.zip',
8         extract=True,
9         cache_dir='.', cache_subdir='data')
```

Downloading data from http://storage.googleapis.com/download.tensorflow.org/data/mini_speech_commands.zip (http://storage.googleapis.com/download.tensorflow.org/data/mini_speech_commands.zip)
182082353/182082353 [=====] - 2s 0us/step

```
In [ ]: 1 from IPython.display import Audio, display
2
3 display(Audio('data/mini_speech_commands/download/004ae714_nohash_0.wav', autoplay=True))
```

0:00 / 0:00

```
In [ ]: 1 commands = np.array(tf.io.gfile.listdir(str(data_dir)))
2 print('Commands:', commands)
```

Commands: ['left' 'no' 'right' 'stop' 'yes' 'down' 'README.md' 'go' 'up']

```
In [ ]: 1 type(commands)
```

```
Out[6]: numpy.ndarray
```

```
In [ ]: 1 commands = commands[(commands != 'README.md') & (commands != '.DS_Store')]
        2 print('Commands:', commands)
```

Commands: ['left' 'no' 'right' 'stop' 'yes' 'down' 'go' 'up']

```
In [ ]: 1 train_ds, val_ds = tf.keras.utils.audio_dataset_from_directory(
        2     directory=data_dir,
        3     batch_size=64,
        4     validation_split=0.2,
        5     seed=0,
        6     output_sequence_length=16000,
        7     subset='both')
        8
        9 label_names = np.array(train_ds.class_names)
        10 print()
        11 print("label names:", label_names)
```

Found 8000 files belonging to 8 classes.
Using 6400 files for training.
Using 1600 files for validation.

label names: ['down' 'go' 'left' 'no' 'right' 'stop' 'up' 'yes']

```
In [ ]: 1 train_ds
```

```
Out[9]: <BatchDataset element_spec=(TensorSpec(shape=(None, 16000, None), dtype=tf.float32, name=None), TensorSpec(shape=(None,), dtype=tf.int32, name=None))>
```

```
In [ ]: 1 val_ds
```

```
Out[10]: <BatchDataset element_spec=(TensorSpec(shape=(None, 16000, None), dtype=tf.float32, name=None), TensorSpec(shape=(None,), dtype=tf.int32, name=None))>
```

```
In [ ]: 1 train_ds.element_spec
```

```
Out[11]: (TensorSpec(shape=(None, 16000, None), dtype=tf.float32, name=None),
TensorSpec(shape=(None,), dtype=tf.int32, name=None))
```

This dataset only contains single channel audio, so use the `tf.squeeze` function to drop the extra axis:

```
In [ ]: 1 def squeeze(audio, labels):
        2     audio = tf.squeeze(audio, axis=-1)
        3     return audio, labels
```

```
In [ ]: 1 train_ds = train_ds.map(squeeze, tf.data.AUTOTUNE)
        2 val_ds = val_ds.map(squeeze, tf.data.AUTOTUNE)
```

The `utils.audio_dataset_from_directory` function only returns up to two splits. It's a good idea to keep a test set separate from your validation set. Ideally you'd keep it in a separate directory, but in this case you can use `Dataset.shard` to split the validation set into two halves. Note that iterating over any shard will load all the data, and only keep its fraction.

```
In [ ]: 1 test_ds = val_ds.shard(num_shards=2, index=0)
        2 val_ds = val_ds.shard(num_shards=2, index=1)
```

```
In [ ]: 1 for example_audio, example_labels in train_ds.take(1):
        2     print(example_audio.shape)
        3     print(example_labels.shape)
```

(64, 16000)
(64,)

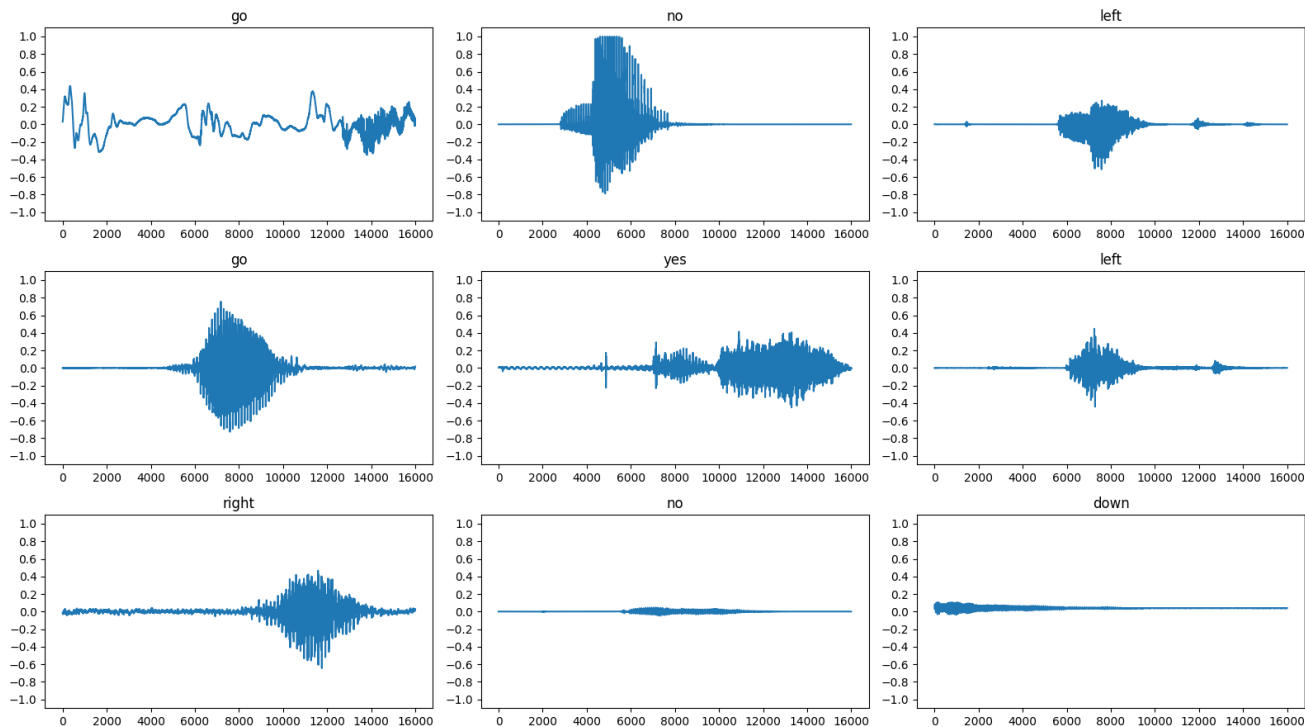
```
In [ ]: 1 label_names[[1,1,3,0]]
```

```
Out[16]: array(['go', 'go', 'no', 'down'], dtype='<U5')
```

```

In [ ]: 1 rows = 3
        2 cols = 3
        3 n = rows * cols
        4 fig, axes = plt.subplots(rows, cols, figsize=(16, 9))
        5
        6 for i in range(n):
        7     if i >= n:
        8         break
        9     r = i // cols
        10    c = i % cols
        11    ax = axes[r][c]
        12    ax.plot(example_audio[i].numpy())
        13    ax.set_yticks(np.arange(-1.2, 1.2, 0.2))
        14    label = label_names[example_labels[i]]
        15    ax.set_title(label)
        16    ax.set_ylim([-1.1, 1.1])
        17 plt.tight_layout()
        18 plt.show()

```



Convert waveforms to spectrograms

The waveforms in the dataset are represented in the time domain. Next, you'll transform the waveforms from the time-domain signals into the time-frequency-domain signals by computing the short-time Fourier transform (STFT) to convert the waveforms to as spectrograms, which show frequency changes over time and can be represented as 2D images. You will feed the spectrogram images into your neural network to train the model

```

In [ ]: 1 def get_spectrogram(waveform):
        2     # Convert the waveform to a spectrogram via a STFT.
        3     spectrogram = tf.signal.stft(
        4         waveform, frame_length=255, frame_step=128)
        5     # Obtain the magnitude of the STFT.
        6     spectrogram = tf.abs(spectrogram)
        7     # Add a `channels` dimension, so that the spectrogram can be used
        8     # as image-like input data with convolution layers (which expect
        9     # shape (`batch_size`, `height`, `width`, `channels`).
        10    spectrogram = spectrogram[..., tf.newaxis]
        11    return spectrogram

```

```
In [ ]: 1 for i in range(3):
2         label = label_names[example_labels[i]]
3         waveform = example_audio[i]
4         spectrogram = get_spectrogram(waveform)
5
6         print('Label:', label)
7         print('Waveform shape:', waveform.shape)
8         print('Spectrogram shape:', spectrogram.shape)
9         print('Audio playback')
10        display.display(display.Audio(waveform, rate=16000))
```

Label: go
Waveform shape: (16000,)
Spectrogram shape: (124, 129, 1)
Audio playback

0:00 / 0:00

Label: no
Waveform shape: (16000,)
Spectrogram shape: (124, 129, 1)
Audio playback

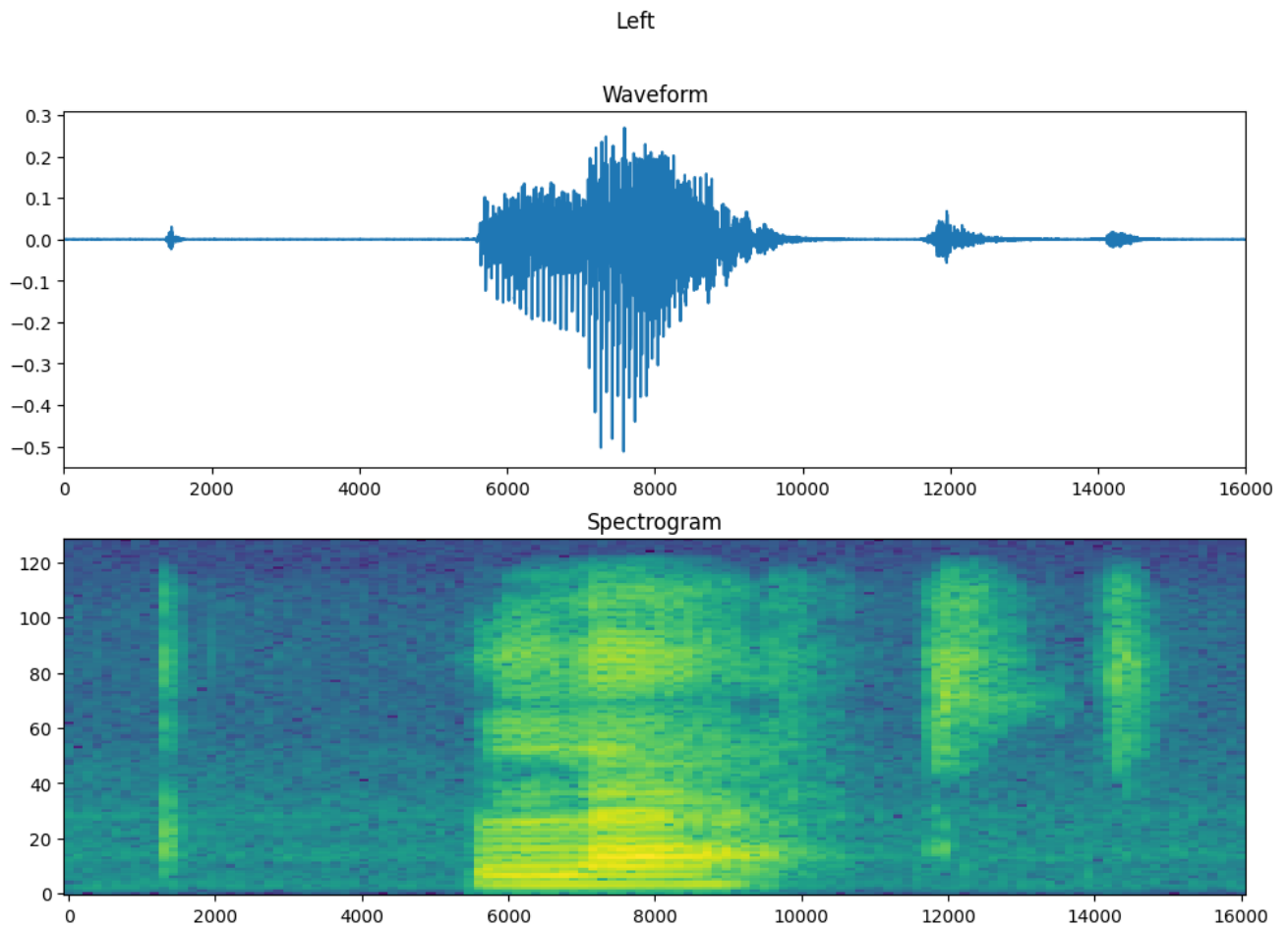
0:00 / 0:00

Label: left
Waveform shape: (16000,)
Spectrogram shape: (124, 129, 1)
Audio playback

0:00 / 0:00

```
In [ ]: 1 def plot_spectrogram(spectrogram, ax):
2         if len(spectrogram.shape) > 2:
3             assert len(spectrogram.shape) == 3
4             spectrogram = np.squeeze(spectrogram, axis=-1)
5         # Convert the frequencies to log scale and transpose, so that the time is
6         # represented on the x-axis (columns).
7         # Add an epsilon to avoid taking a log of zero.
8         log_spec = np.log(spectrogram.T + np.finfo(float).eps)
9         height = log_spec.shape[0]
10        width = log_spec.shape[1]
11        X = np.linspace(0, np.size(spectrogram), num=width, dtype=int)
12        Y = range(height)
13        ax.pcolormesh(X, Y, log_spec)
```

```
In [ ]: 1 fig, axes = plt.subplots(2, figsize=(12, 8))
2 timescale = np.arange(waveform.shape[0])
3 axes[0].plot(timescale, waveform.numpy())
4 axes[0].set_title('Waveform')
5 axes[0].set_xlim([0, 16000])
6
7 plot_spectrogram(spectrogram.numpy(), axes[1])
8 axes[1].set_title('Spectrogram')
9 plt.suptitle(label.title())
10 plt.show()
```



Now, create spectrogram datasets from the audio datasets:

```
In [ ]: 1 def make_spec_ds(ds):
2         return ds.map(
3             map_func=lambda audio, label: (get_spectrogram(audio), label),
4             num_parallel_calls=tf.data.AUTOTUNE)
```

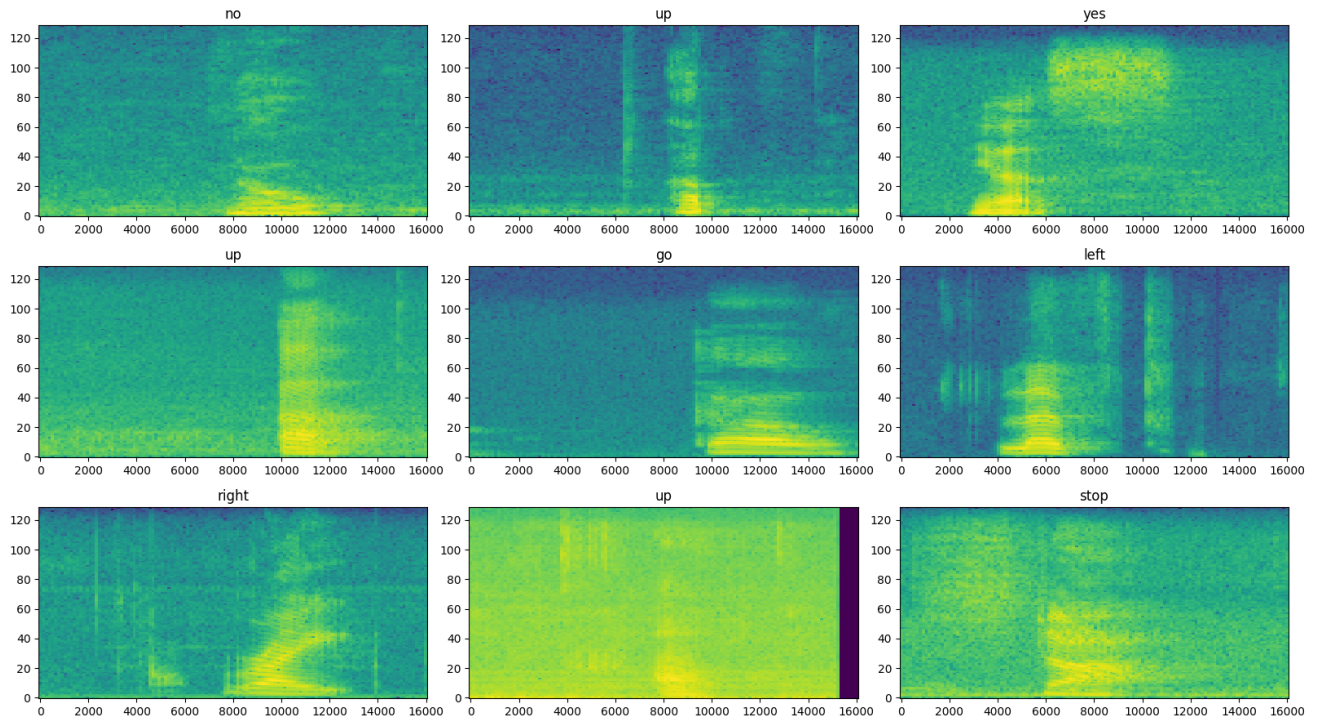
```
In [ ]: 1 train_spectrogram_ds = make_spec_ds(train_ds)
2 val_spectrogram_ds = make_spec_ds(val_ds)
3 test_spectrogram_ds = make_spec_ds(test_ds)
```

```
In [ ]: 1 for example_spectrograms, example_spect_labels in train_spectrogram_ds.take(1):
2         break
```

```

In [ ]: 1 rows = 3
        2 cols = 3
        3 n = rows*cols
        4 fig, axes = plt.subplots(rows, cols, figsize=(16, 9))
        5
        6 for i in range(n):
        7     r = i // cols
        8     c = i % cols
        9     ax = axes[r][c]
       10     plot_spectrogram(example_spectrograms[i].numpy(), ax)
       11     ax.set_title(label_names[example_spect_labels[i].numpy()])
       12 plt.tight_layout()
       13 plt.show()

```



Build and train the model

```

In [ ]: 1 train_spectrogram_ds = train_spectrogram_ds.cache().shuffle(10000).prefetch(tf.data.AUTOTUNE)
        2 val_spectrogram_ds = val_spectrogram_ds.cache().prefetch(tf.data.AUTOTUNE)
        3 test_spectrogram_ds = test_spectrogram_ds.cache().prefetch(tf.data.AUTOTUNE)

```

```

In [ ]: 1 input_shape = example_spectrograms.shape[1:]
2 print('Input shape:', input_shape)
3 num_labels = len(label_names)
4
5 # Instantiate the `tf.keras.layers.Normalization` layer.
6 norm_layer = layers.Normalization()
7 # Fit the state of the layer to the spectrograms
8 # with `Normalization.adapt`.
9 norm_layer.adapt(data=train_spectrogram_ds.map(map_func=lambda spec, label: spec))
10
11 model = models.Sequential([
12     layers.Input(shape=input_shape),
13     # Downsample the input.
14     layers.Resizing(32, 32),
15     # Normalize.
16     norm_layer,
17     layers.Conv2D(32, 3, activation='relu'),
18     layers.Conv2D(64, 3, activation='relu'),
19     layers.MaxPooling2D(),
20     layers.Dropout(0.25),
21     layers.Flatten(),
22     layers.Dense(128, activation='relu'),
23     layers.Dropout(0.5),
24     layers.Dense(num_labels),
25 ])
26
27 model.summary()

```

Input shape: (124, 129, 1)
Model: "sequential"

Layer (type)	Output Shape	Param #
=====		
resizing (Resizing)	(None, 32, 32, 1)	0
normalization (Normalization)	(None, 32, 32, 1)	3
conv2d (Conv2D)	(None, 30, 30, 32)	320
conv2d_1 (Conv2D)	(None, 28, 28, 64)	18496
max_pooling2d (MaxPooling2D)	(None, 14, 14, 64)	0
dropout (Dropout)	(None, 14, 14, 64)	0
flatten (Flatten)	(None, 12544)	0
dense (Dense)	(None, 128)	1605760
dropout_1 (Dropout)	(None, 128)	0
dense_1 (Dense)	(None, 8)	1032
=====		
Total params: 1,625,611		
Trainable params: 1,625,608		
Non-trainable params: 3		

```

In [ ]: 1 model.compile(
2     optimizer=tf.keras.optimizers.Adam(),
3     loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True),
4     metrics=['accuracy'],
5 )

```

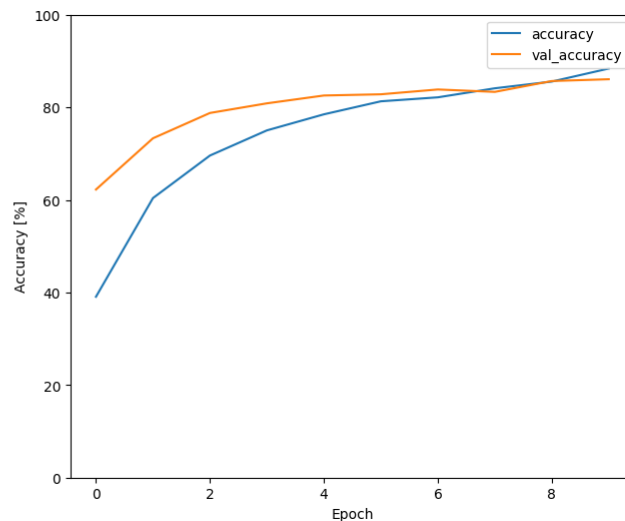
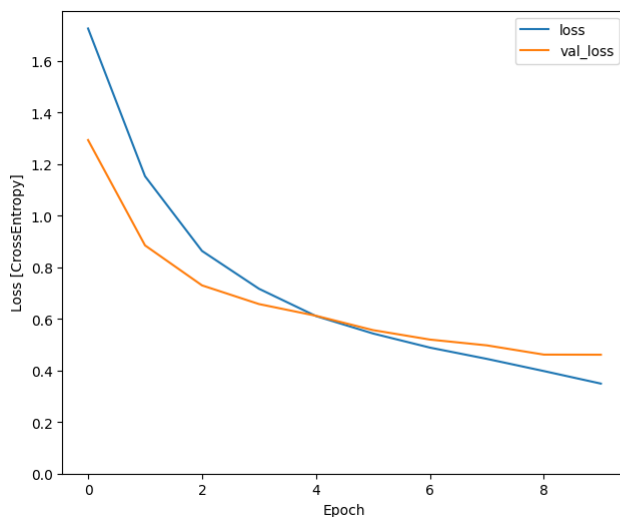


```
In [ ]: 1 EPOCHS = 10
2 history = model.fit(
3     train_spectrogram_ds,
4     validation_data=val_spectrogram_ds,
5     epochs=EPOCHS,
6     callbacks=tf.keras.callbacks.EarlyStopping(verbose=1, patience=2),
7 )
```

```
Epoch 1/10
100/100 [=====] - 43s 393ms/step - loss: 1.7266 - accuracy: 0.3908 - val_loss: 1.2939 - val_accuracy:
0.6224
Epoch 2/10
100/100 [=====] - 40s 404ms/step - loss: 1.1535 - accuracy: 0.6039 - val_loss: 0.8852 - val_accuracy:
0.7331
Epoch 3/10
100/100 [=====] - 39s 391ms/step - loss: 0.8637 - accuracy: 0.6958 - val_loss: 0.7302 - val_accuracy:
0.7878
Epoch 4/10
100/100 [=====] - 36s 364ms/step - loss: 0.7170 - accuracy: 0.7502 - val_loss: 0.6577 - val_accuracy:
0.8086
Epoch 5/10
100/100 [=====] - 31s 313ms/step - loss: 0.6111 - accuracy: 0.7848 - val_loss: 0.6123 - val_accuracy:
0.8255
Epoch 6/10
100/100 [=====] - 29s 288ms/step - loss: 0.5436 - accuracy: 0.8130 - val_loss: 0.5567 - val_accuracy:
0.8281
Epoch 7/10
100/100 [=====] - 29s 290ms/step - loss: 0.4888 - accuracy: 0.8216 - val_loss: 0.5200 - val_accuracy:
0.8385
Epoch 8/10
100/100 [=====] - 28s 280ms/step - loss: 0.4451 - accuracy: 0.8411 - val_loss: 0.4975 - val_accuracy:
0.8333
Epoch 9/10
100/100 [=====] - 32s 312ms/step - loss: 0.3978 - accuracy: 0.8556 - val_loss: 0.4617 - val_accuracy:
0.8568
Epoch 10/10
100/100 [=====] - 28s 277ms/step - loss: 0.3492 - accuracy: 0.8838 - val_loss: 0.4613 - val_accuracy:
0.8607
```

```
In [ ]: 1 metrics = history.history
2 plt.figure(figsize=(16,6))
3 plt.subplot(1,2,1)
4 plt.plot(history.epoch, metrics['loss'], metrics['val_loss'])
5 plt.legend(['loss', 'val_loss'])
6 plt.ylim([0, max(plt.ylim())])
7 plt.xlabel('Epoch')
8 plt.ylabel('Loss [CrossEntropy]')
9
10 plt.subplot(1,2,2)
11 plt.plot(history.epoch, 100*np.array(metrics['accuracy']), 100*np.array(metrics['val_accuracy']))
12 plt.legend(['accuracy', 'val_accuracy'])
13 plt.ylim([0, 100])
14 plt.xlabel('Epoch')
15 plt.ylabel('Accuracy [%]')
```

Out[33]: Text(0, 0.5, 'Accuracy [%]')



```
In [ ]: 1 model.evaluate(test_spectrogram_ds, return_dict=True)
```

```
13/13 [=====] - 2s 163ms/step - loss: 0.5173 - accuracy: 0.8341
```

```
Out[34]: {'loss': 0.5173158049583435, 'accuracy': 0.8341346383094788}
```

Display a confusion matrix

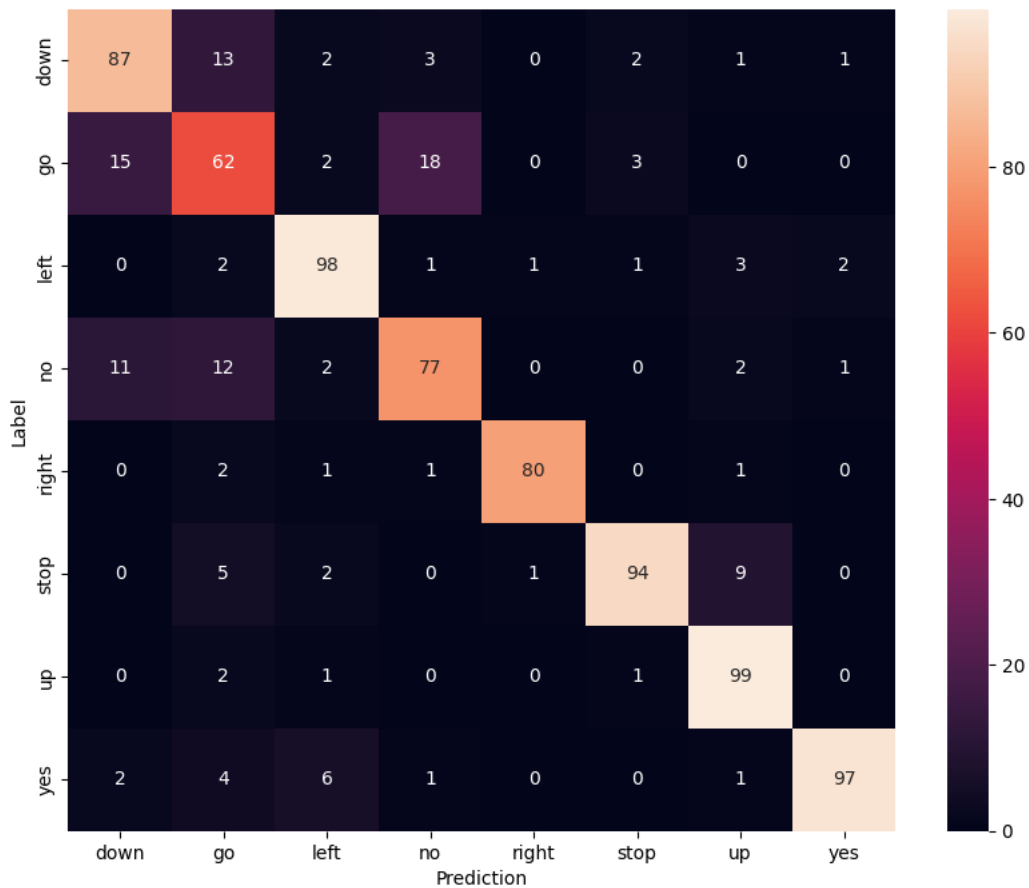
```
In [ ]: 1 y_pred = model.predict(test_spectrogram_ds)
```

```
13/13 [=====] - 1s 53ms/step
```

```
In [ ]: 1 y_pred = tf.argmax(y_pred, axis=1)
```

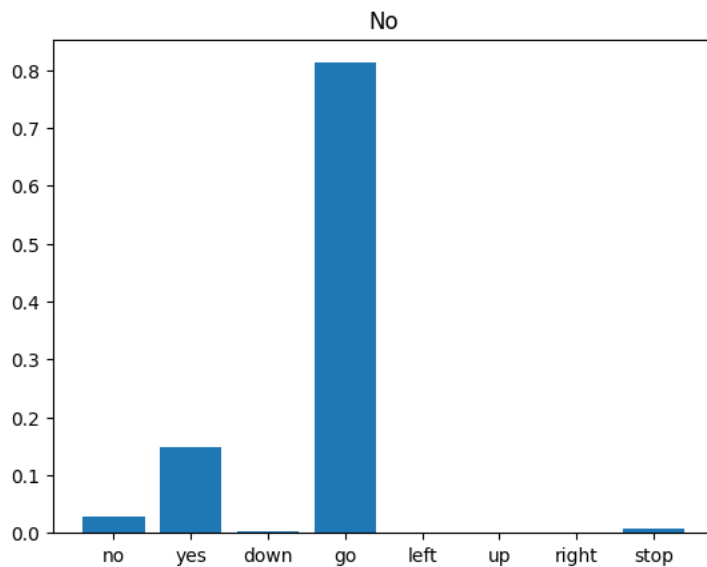
```
In [ ]: 1 y_true = tf.concat(list(test_spectrogram_ds.map(lambda s,lab: lab)), axis=0)
```

```
In [ ]: 1 confusion_mtx = tf.math.confusion_matrix(y_true, y_pred)
2 plt.figure(figsize=(10, 8))
3 sns.heatmap(confusion_mtx,
4             xticklabels=label_names,
5             yticklabels=label_names,
6             annot=True, fmt='g')
7 plt.xlabel('Prediction')
8 plt.ylabel('Label')
9 plt.show()
```



Run inference on an audio file

```
In [ ]: 1 x = data_dir/'no/01bb6a2a_nohash_0.wav'
2 x = tf.io.read_file(str(x))
3 x, sample_rate = tf.audio.decode_wav(x, desired_channels=1, desired_samples=16000,)
4 x = tf.squeeze(x, axis=-1)
5 waveform = x
6 x = get_spectrogram(x)
7 x = x[tf.newaxis,...]
8
9 prediction = model(x)
10 x_labels = ['no', 'yes', 'down', 'go', 'left', 'up', 'right', 'stop']
11 plt.bar(x_labels, tf.nn.softmax(prediction[0]))
12 plt.title('No')
13 plt.show()
14
15 display.display(display.Audio(waveform, rate=16000))
```



0:00 / 0:00

Export the model with preprocessing

```

In [ ]: 1 class ExportModel(tf.Module):
        2     def __init__(self, model):
        3         self.model = model
        4
        5         # Accept either a string-filename or a batch of waveforms.
        6         # You could add additional signatures for a single wave, or a ragged-batch.
        7         self.__call__.get_concrete_function(
        8             x=tf.TensorSpec(shape=(), dtype=tf.string))
        9         self.__call__.get_concrete_function(
        10            x=tf.TensorSpec(shape=[None, 16000], dtype=tf.float32))
        11
        12
        13     @tf.function
        14     def __call__(self, x):
        15         # If they pass a string, load the file and decode it.
        16         if x.dtype == tf.string:
        17             x = tf.io.read_file(x)
        18             x, _ = tf.audio.decode_wav(x, desired_channels=1, desired_samples=16000,)
        19             x = tf.squeeze(x, axis=-1)
        20             x = x[tf.newaxis, :]
        21
        22             x = get_spectrogram(x)
        23             result = self.model(x, training=False)
        24
        25             class_ids = tf.argmax(result, axis=-1)
        26             class_names = tf.gather(label_names, class_ids)
        27             return {'predictions': result,
        28                   'class_ids': class_ids,
        29                   'class_names': class_names}

```

```

In [ ]: 1 export = ExportModel(model)
        2 export(tf.constant(str(data_dir/'no/01bb6a2a_nohash_0.wav')))

```

```

Out[41]: {'predictions': <tf.Tensor: shape=(1, 8), dtype=float32, numpy=
array([[ 0.6647731,  2.3551004, -1.5815372,  4.0603275, -3.4281795,
        -2.5847144, -2.7507148, -0.5820455]], dtype=float32)>,
'class_ids': <tf.Tensor: shape=(1,), dtype=int64, numpy=array([3])>,
'class_names': <tf.Tensor: shape=(1,), dtype=string, numpy=array([b'no'], dtype=object)>}

```

```

In [ ]: 1 tf.saved_model.save(export, "saved")
        2 imported = tf.saved_model.load("saved")
        3 imported(waveform[tf.newaxis, :])

```

WARNING:absl:Found untraced functions such as _jit_compiled_convolution_op, _jit_compiled_convolution_op, _update_step_xla while saving (showing 3 of 3). These functions will not be directly callable after loading.

```

Out[42]: {'class_names': <tf.Tensor: shape=(1,), dtype=string, numpy=array([b'no'], dtype=object)>,
'predictions': <tf.Tensor: shape=(1, 8), dtype=float32, numpy=
array([[ 0.6647731,  2.3551004, -1.5815372,  4.0603275, -3.4281795,
        -2.5847144, -2.7507148, -0.5820455]], dtype=float32)>,
'class_ids': <tf.Tensor: shape=(1,), dtype=int64, numpy=array([3])>}

```

```

In [ ]: 1

```