



► [Code examples](#) / [Audio Data](#) / Automatic Speech Recognition with Transformer

# Automatic Speech Recognition with Transformer

**Author:** [Apoorv Nandan](#)  
**Date created:** 2021/01/13  
**Last modified:** 2021/01/13  
**Description:** Training a sequence-to-sequence Transformer for automatic speech recognition.

This example uses Keras 3

[View in Colab](#) · [GitHub source](#)

## Introduction

Automatic speech recognition (ASR) consists of transcribing audio speech segments into text. ASR can be treated as a sequence-to-sequence problem, where the audio can be represented as a sequence of feature vectors and the text as a sequence of characters, words, or subword tokens.

For this demonstration, we will use the LJSpeech dataset from the [LibriVox](#) project. It consists of short audio clips of a single speaker reading passages from 7 non-fiction books. Our model will be similar to the original Transformer (both encoder and decoder) as proposed in the paper, "Attention is All You Need".

### References:

- [Attention is All You Need](#)
- [Very Deep Self-Attention Networks for End-to-End Speech Recognition](#)
- [Speech Transformers](#)
- [LJSpeech Dataset](#)

```
import os

os.environ["KERAS_BACKEND"] = "tensorflow"

from glob import glob
import tensorflow as tf
import keras
from keras import layers
```

## Define the Transformer Input Layer

When processing past target tokens for the decoder, we compute the sum of position embeddings and token embeddings.

When processing audio features, we apply convolutional layers to downsample them (via convolution strides) and process local relationships.

### Automatic Speech Recognition with Transformer

- ◆ [Introduction](#)
- ◆ [Define the Transformer Input Layer](#)
- ◆ [Transformer Encoder Layer](#)
- ◆ [Transformer Decoder Layer](#)
- ◆ [Complete the Transformer model](#)
- ◆ [Download the dataset](#)
- ◆ [Preprocess the dataset](#)
- ◆ [Callbacks to display predictions](#)
- ◆ [Learning rate schedule](#)
- ◆ [Create & train the end-to-end model](#)

```
class TokenEmbedding(layers.Layer):
    def __init__(self, num_vocab=1000, maxlen=100, num_hid=64):
        super().__init__()
        self.emb = keras.layers.Embedding(num_vocab, num_hid)
        self.pos_emb = layers.Embedding(input_dim=maxlen, output_dim=num_hid)

    def call(self, x):
        maxlen = tf.shape(x)[-1]
        x = self.emb(x)
        positions = tf.range(start=0, limit=maxlen, delta=1)
        positions = self.pos_emb(positions)
        return x + positions

class SpeechFeatureEmbedding(layers.Layer):
    def __init__(self, num_hid=64, maxlen=100):
        super().__init__()
        self.conv1 = keras.layers.Conv1D(
            num_hid, 11, strides=2, padding="same", activation="relu"
        )
        self.conv2 = keras.layers.Conv1D(
            num_hid, 11, strides=2, padding="same", activation="relu"
        )
        self.conv3 = keras.layers.Conv1D(
            num_hid, 11, strides=2, padding="same", activation="relu"
        )

    def call(self, x):
        x = self.conv1(x)
        x = self.conv2(x)
        return self.conv3(x)
```

Automatic Speech Recognition with Transformer

- ◆ [Introduction](#)
- ◆ [Define the Transformer Input Layer](#)
- ◆ [Transformer Encoder Layer](#)
- ◆ [Transformer Decoder Layer](#)
- ◆ [Complete the Transformer model](#)
- ◆ [Download the dataset](#)
- ◆ [Preprocess the dataset](#)
- ◆ [Callbacks to display predictions](#)
- ◆ [Learning rate schedule](#)
- ◆ [Create & train the end-to-end model](#)

Transformer Encoder Layer

```
class TransformerEncoder(layers.Layer):
    def __init__(self, embed_dim, num_heads, feed_forward_dim, rate=0.1):
        super().__init__()
        self.att = layers.MultiHeadAttention(num_heads=num_heads, key_dim=embed_dim)
        self.ffn = keras.Sequential(
            [
                layers.Dense(feed_forward_dim, activation="relu"),
                layers.Dense(embed_dim),
            ]
        )
        self.layernorm1 = layers.LayerNormalization(epsilon=1e-6)
        self.layernorm2 = layers.LayerNormalization(epsilon=1e-6)
        self.dropout1 = layers.Dropout(rate)
        self.dropout2 = layers.Dropout(rate)

    def call(self, inputs, training=False):
        attn_output = self.att(inputs, inputs)
        attn_output = self.dropout1(attn_output, training=training)
        out1 = self.layernorm1(inputs + attn_output)
        ffn_output = self.ffn(out1)
        ffn_output = self.dropout2(ffn_output, training=training)
        return self.layernorm2(out1 + ffn_output)
```

# Transformer Decoder Layer

```
class TransformerDecoder(layers.Layer):
    def __init__(self, embed_dim, num_heads, feed_forward_dim, dropout_rate=0.1):
        super().__init__()
        self.layernorm1 = layers.LayerNormalization(epsilon=1e-6)
        self.layernorm2 = layers.LayerNormalization(epsilon=1e-6)
        self.layernorm3 = layers.LayerNormalization(epsilon=1e-6)
        self.self_att = layers.MultiHeadAttention(
            num_heads=num_heads, key_dim=embed_dim
        )
        self.enc_att = layers.MultiHeadAttention(num_heads=num_heads, key_dim=embed_dim)
        self.self_dropout = layers.Dropout(0.5)
        self.enc_dropout = layers.Dropout(0.1)
        self.ffn_dropout = layers.Dropout(0.1)
        self.ffn = keras.Sequential(
            [
                layers.Dense(feed_forward_dim, activation="relu"),
                layers.Dense(embed_dim),
            ]
        )

    def causal_attention_mask(self, batch_size, n_dest, n_src, dtype):
        """Masks the upper half of the dot product matrix in self attention.

        This prevents flow of information from future tokens to current token.
        1's in the lower triangle, counting from the lower right corner.
        """
        i = tf.range(n_dest)[: , None]
        j = tf.range(n_src)
        m = i >= j - n_src + n_dest
        mask = tf.cast(m, dtype)
        mask = tf.reshape(mask, [1, n_dest, n_src])
        mult = tf.concat(
            [tf.expand_dims(batch_size, -1), tf.constant([1, 1], dtype=tf.int32)], 0
        )
        return tf.tile(mask, mult)

    def call(self, enc_out, target):
        input_shape = tf.shape(target)
        batch_size = input_shape[0]
        seq_len = input_shape[1]
        causal_mask = self.causal_attention_mask(batch_size, seq_len, seq_len, tf.bool)
        target_att = self.self_att(target, target, attention_mask=causal_mask)
        target_norm = self.layernorm1(target + self.self_dropout(target_att))
        enc_out = self.enc_att(target_norm, enc_out)
        enc_out_norm = self.layernorm2(self.enc_dropout(enc_out) + target_norm)
        ffn_out = self.ffn(enc_out_norm)
        ffn_out_norm = self.layernorm3(enc_out_norm + self.ffn_dropout(ffn_out))
        return ffn_out_norm
```

## Automatic Speech Recognition with Transformer

- ◆ [Introduction](#)
- ◆ [Define the Transformer Input Layer](#)
- ◆ [Transformer Encoder Layer](#)
- ◆ [Transformer Decoder Layer](#)
- ◆ [Complete the Transformer model](#)
- ◆ [Download the dataset](#)
- ◆ [Preprocess the dataset](#)
- ◆ [Callbacks to display predictions](#)
- ◆ [Learning rate schedule](#)
- ◆ [Create & train the end-to-end model](#)

# Complete the Transformer model

Our model takes audio spectrograms as inputs and predicts a sequence of characters. During training, we give the decoder the target character sequence shifted to the left as input. During inference, the decoder uses its own past predictions to predict the next token.

```

class Transformer(keras.Model):
    def __init__(
        self,
        num_hid=64,
        num_head=2,
        num_feed_forward=128,
        source_maxlen=100,
        target_maxlen=100,
        num_layers_enc=4,
        num_layers_dec=1,
        num_classes=10,
    ):
        super().__init__()
        self.loss_metric = keras.metrics.Mean(name="loss")
        self.num_layers_enc = num_layers_enc
        self.num_layers_dec = num_layers_dec
        self.target_maxlen = target_maxlen
        self.num_classes = num_classes

        self.enc_input = SpeechFeatureEmbedding(num_hid=num_hid, maxlen=source_maxlen)
        self.dec_input = TokenEmbedding(
            num_vocab=num_classes, maxlen=target_maxlen, num_hid=num_hid
        )

        self.encoder = keras.Sequential(
            [self.enc_input]
            + [
                TransformerEncoder(num_hid, num_head, num_feed_forward)
                for _ in range(num_layers_enc)
            ]
        )

        for i in range(num_layers_dec):
            setattr(
                self,
                f"dec_layer_{i}",
                TransformerDecoder(num_hid, num_head, num_feed_forward),
            )

        self.classifier = layers.Dense(num_classes)

    def decode(self, enc_out, target):
        y = self.dec_input(target)
        for i in range(self.num_layers_dec):
            y = getattr(self, f"dec_layer_{i}")(enc_out, y)
        return y

    def call(self, inputs):
        source = inputs[0]
        target = inputs[1]
        x = self.encoder(source)
        y = self.decode(x, target)
        return self.classifier(y)

    @property
    def metrics(self):
        return [self.loss_metric]

    def train_step(self, batch):
        """Processes one batch inside model.fit()."""
        source = batch["source"]
        target = batch["target"]
        dec_input = target[:, :-1]
        dec_target = target[:, 1:]
        with tf.GradientTape() as tape:
            preds = self([source, dec_input])
            one_hot = tf.one_hot(dec_target, depth=self.num_classes)
            mask = tf.math.logical_not(tf.math.equal(dec_target, 0))
            loss = model.compute_loss(None, one_hot, preds, sample_weight=mask)
        trainable_vars = self.trainable_variables
        gradients = tape.gradient(loss, trainable_vars)
        self.optimizer.apply_gradients(zip(gradients, trainable_vars))
        self.loss_metric.update_state(loss)
        return {"loss": self.loss_metric.result()}

    def test_step(self, batch):
        source = batch["source"]
        target = batch["target"]

```

## Automatic Speech Recognition with Transformer

- ◆ [Introduction](#)
- ◆ [Define the Transformer Input Layer](#)
- ◆ [Transformer Encoder Layer](#)
- ◆ [Transformer Decoder Layer](#)
- ◆ [Complete the Transformer model](#)
- ◆ [Download the dataset](#)
- ◆ [Preprocess the dataset](#)
- ◆ [Callbacks to display predictions](#)
- ◆ [Learning rate schedule](#)
- ◆ [Create & train the end-to-end model](#)

```
dec_input = target[:, :-1]
dec_target = target[:, 1:]
preds = self([source, dec_input])
one_hot = tf.one_hot(dec_target, depth=self.num_classes)
mask = tf.math.logical_not(tf.math.equal(dec_target, 0))
loss = model.compute_loss(None, one_hot, preds, sample_weight=mask)
self.loss_metric.update_state(loss)
return {"loss": self.loss_metric.result()}

def generate(self, source, target_start_token_idx):
    """Performs inference over one batch of inputs using greedy decoding."""
    bs = tf.shape(source)[0]
    enc = self.encoder(source)
    dec_input = tf.ones((bs, 1), dtype=tf.int32) * target_start_token_idx
    dec_logits = []
    for i in range(self.target_maxlen - 1):
        dec_out = self.decode(enc, dec_input)
        logits = self.classifier(dec_out)
        logits = tf.argmax(logits, axis=-1, output_type=tf.int32)
        last_logit = tf.expand_dims(logits[:, -1], axis=-1)
        dec_logits.append(last_logit)
        dec_input = tf.concat([dec_input, last_logit], axis=-1)
    return dec_input
```

Automatic Speech Recognition with Transformer

- ◆ [Introduction](#)
- ◆ [Define the Transformer Input Layer](#)
- ◆ [Transformer Encoder Layer](#)
- ◆ [Transformer Decoder Layer](#)
- ◆ [Complete the Transformer model](#)
- ◆ [Download the dataset](#)
- ◆ [Preprocess the dataset](#)
- ◆ [Callbacks to display predictions](#)
- ◆ [Learning rate schedule](#)
- ◆ [Create & train the end-to-end model](#)

Download the dataset

Note: This requires ~3.6 GB of disk space and takes ~5 minutes for the extraction of files.

```
keras.utils.get_file(
    os.path.join(os.getcwd(), "data.tar.gz"),
    "https://data.keithito.com/data/speech/LJSpeech-1.1.tar.bz2",
    extract=True,
    archive_format="tar",
    cache_dir=".",
)

saveto = "./datasets/LJSpeech-1.1"
wavs = glob("{}/**/*.*wav".format(saveto), recursive=True)

id_to_text = {}
with open(os.path.join(saveto, "metadata.csv"), encoding="utf-8") as f:
    for line in f:
        id = line.strip().split("|")[0]
        text = line.strip().split("|")[2]
        id_to_text[id] = text

def get_data(wavs, id_to_text, maxlen=50):
    """returns mapping of audio paths and transcription texts"""
    data = []
    for w in wavs:
        id = w.split("/")[-1].split(".")[0]
        if len(id_to_text[id]) < maxlen:
            data.append({"audio": w, "text": id_to_text[id]})
    return data
```

Downloading data from https://data.keithito.com/data/speech/LJSpeech-1.1.tar.bz2  
2748572632/2748572632 18s 0us/step

# Preprocess the dataset

```
class VectorizeChar:
    def __init__(self, max_len=50):
        self.vocab = (
            ["-", "#", "<", ">"]
            + [chr(i + 96) for i in range(1, 27)]
            + [" ", ".", ",", "?"]
        )
        self.max_len = max_len
        self.char_to_idx = {}
        for i, ch in enumerate(self.vocab):
            self.char_to_idx[ch] = i

    def __call__(self, text):
        text = text.lower()
        text = text[: self.max_len - 2]
        text = "<" + text + ">"
        pad_len = self.max_len - len(text)
        return [self.char_to_idx.get(ch, 1) for ch in text] + [0] * pad_len

    def get_vocabulary(self):
        return self.vocab

max_target_len = 200 # all transcripts in out data are < 200 characters
data = get_data(wavs, id_to_text, max_target_len)
vectorizer = VectorizeChar(max_target_len)
print("vocab size", len(vectorizer.get_vocabulary()))

def create_text_ds(data):
    texts = [_["text"] for _ in data]
    text_ds = [vectorizer(t) for t in texts]
    text_ds = tf.data.Dataset.from_tensor_slices(text_ds)
    return text_ds

def path_to_audio(path):
    # spectrogram using stft
    audio = tf.io.read_file(path)
    audio, _ = tf.audio.decode_wav(audio, 1)
    audio = tf.squeeze(audio, axis=-1)
    stfts = tf.signal.stft(audio, frame_length=200, frame_step=80, fft_length=256)
    x = tf.math.pow(tf.abs(stfts), 0.5)
    # normalisation
    means = tf.math.reduce_mean(x, 1, keepdims=True)
    stddevs = tf.math.reduce_std(x, 1, keepdims=True)
    x = (x - means) / stddevs
    audio_len = tf.shape(x)[0]
    # padding to 10 seconds
    pad_len = 2754
    paddings = tf.constant([[0, pad_len], [0, 0]])
    x = tf.pad(x, paddings, "CONSTANT")[:pad_len, :]
    return x

def create_audio_ds(data):
    flist = [_["audio"] for _ in data]
    audio_ds = tf.data.Dataset.from_tensor_slices(flist)
    audio_ds = audio_ds.map(path_to_audio, num_parallel_calls=tf.data.AUTOTUNE)
    return audio_ds

def create_tf_dataset(data, bs=4):
    audio_ds = create_audio_ds(data)
    text_ds = create_text_ds(data)
    ds = tf.data.Dataset.zip((audio_ds, text_ds))
    ds = ds.map(lambda x, y: {"source": x, "target": y})
    ds = ds.batch(bs)
    ds = ds.prefetch(tf.data.AUTOTUNE)
    return ds

split = int(len(data) * 0.99)
train_data = data[:split]
test_data = data[split:]
```

## Automatic Speech Recognition with Transformer

- ◆ [Introduction](#)
- ◆ [Define the Transformer Input Layer](#)
- ◆ [Transformer Encoder Layer](#)
- ◆ [Transformer Decoder Layer](#)
- ◆ [Complete the Transformer model](#)
- ◆ [Download the dataset](#)
- ◆ [Preprocess the dataset](#)
- ◆ [Callbacks to display predictions](#)
- ◆ [Learning rate schedule](#)
- ◆ [Create & train the end-to-end model](#)

```
ds = create_tf_dataset(train_data, bs=64)
val_ds = create_tf_dataset(test_data, bs=4)
```

vocab size 34

## Callbacks to display predictions

```
class DisplayOutputs(keras.callbacks.Callback):
    def __init__(
        self, batch, idx_to_token, target_start_token_idx=27, target_end_token_idx=28
    ):
        """Displays a batch of outputs after every epoch

        Args:
            batch: A test batch containing the keys "source" and "target"
            idx_to_token: A List containing the vocabulary tokens corresponding to their
indices
            target_start_token_idx: A start token index in the target vocabulary
            target_end_token_idx: An end token index in the target vocabulary
        """
        self.batch = batch
        self.target_start_token_idx = target_start_token_idx
        self.target_end_token_idx = target_end_token_idx
        self.idx_to_char = idx_to_token

    def on_epoch_end(self, epoch, logs=None):
        if epoch % 5 != 0:
            return
        source = self.batch["source"]
        target = self.batch["target"].numpy()
        bs = tf.shape(source)[0]
        preds = self.model.generate(source, self.target_start_token_idx)
        preds = preds.numpy()
        for i in range(bs):
            target_text = "".join([self.idx_to_char[_] for _ in target[i, :]])
            prediction = ""
            for idx in preds[i, :]:
                prediction += self.idx_to_char[idx]
                if idx == self.target_end_token_idx:
                    break
            print(f"target:    {target_text.replace('-', '')}")
            print(f"prediction: {prediction}\n")
```

### Automatic Speech Recognition with Transformer

- ◆ [Introduction](#)
- ◆ [Define the Transformer Input Layer](#)
- ◆ [Transformer Encoder Layer](#)
- ◆ [Transformer Decoder Layer](#)
- ◆ [Complete the Transformer model](#)
- ◆ [Download the dataset](#)
- ◆ [Preprocess the dataset](#)
- ◆ [Callbacks to display predictions](#)
- ◆ [Learning rate schedule](#)
- ◆ [Create & train the end-to-end model](#)



# Learning rate schedule

```
class CustomSchedule(keras.optimizers.schedules.LearningRateSchedule):
    def __init__(
        self,
        init_lr=0.00001,
        lr_after_warmup=0.001,
        final_lr=0.00001,
        warmup_epochs=15,
        decay_epochs=85,
        steps_per_epoch=203,
    ):
        super().__init__()
        self.init_lr = init_lr
        self.lr_after_warmup = lr_after_warmup
        self.final_lr = final_lr
        self.warmup_epochs = warmup_epochs
        self.decay_epochs = decay_epochs
        self.steps_per_epoch = steps_per_epoch

    def calculate_lr(self, epoch):
        """linear warm up - linear decay"""
        warmup_lr = (
            self.init_lr
            + ((self.lr_after_warmup - self.init_lr) / (self.warmup_epochs - 1)) * epoch
        )
        decay_lr = tf.math.maximum(
            self.final_lr,
            self.lr_after_warmup
            - (epoch - self.warmup_epochs)
            * (self.lr_after_warmup - self.final_lr)
            / self.decay_epochs,
        )
        return tf.math.minimum(warmup_lr, decay_lr)

    def __call__(self, step):
        epoch = step // self.steps_per_epoch
        epoch = tf.cast(epoch, "float32")
        return self.calculate_lr(epoch)
```

## Automatic Speech Recognition with Transformer

- ◆ [Introduction](#)
- ◆ [Define the Transformer Input Layer](#)
- ◆ [Transformer Encoder Layer](#)
- ◆ [Transformer Decoder Layer](#)
- ◆ [Complete the Transformer model](#)
- ◆ [Download the dataset](#)
- ◆ [Preprocess the dataset](#)
- ◆ [Callbacks to display predictions](#)
- ◆ [Learning rate schedule](#)
- ◆ [Create & train the end-to-end model](#)



# Create & train the end-to-end model

```
batch = next(iter(val_ds))

# The vocabulary to convert predicted indices into characters
idx_to_char = vectorizer.get_vocabulary()
display_cb = DisplayOutputs(
    batch, idx_to_char, target_start_token_idx=2, target_end_token_idx=3
) # set the arguments as per vocabulary index for '<' and '>'

model = Transformer(
    num_hid=200,
    num_head=2,
    num_feed_forward=400,
    target_maxlen=max_target_len,
    num_layers_enc=4,
    num_layers_dec=1,
    num_classes=34,
)
loss_fn = keras.losses.CategoricalCrossentropy(
    from_logits=True,
    label_smoothing=0.1,
)

learning_rate = CustomSchedule(
    init_lr=0.00001,
    lr_after_warmup=0.001,
    final_lr=0.00001,
    warmup_epochs=15,
    decay_epochs=85,
    steps_per_epoch=len(ds),
)
optimizer = keras.optimizers.Adam(learning_rate)
model.compile(optimizer=optimizer, loss=loss_fn)

history = model.fit(ds, validation_data=val_ds, callbacks=[display_cb], epochs=1)
```

## Automatic Speech Recognition with Transformer

- ◆ [Introduction](#)
- ◆ [Define the Transformer Input Layer](#)
- ◆ [Transformer Encoder Layer](#)
- ◆ [Transformer Decoder Layer](#)
- ◆ [Complete the Transformer model](#)
- ◆ [Download the dataset](#)
- ◆ [Preprocess the dataset](#)
- ◆ [Callbacks to display predictions](#)
- ◆ [Learning rate schedule](#)
- ◆ [Create & train the end-to-end model](#)

1/203 [37m 9:20:11 166s/step - loss: 2.2387

WARNING: All log messages before absl::InitializeLog() is called are written to STDERR  
I0000 00:00:1700071380.331418 678094 device\_compiler.h:187] Compiled cluster using XLA! This line is logged at most once for the lifetime of the process.

203/203 0s 947ms/step - loss: 1.8285target: <the relations between lee and marina oswald are of great importance in any attempt to understand oswald#s possible motivation.>  
prediction: <the the he at the t the an of t te the ale t he t te ar the in the the s the s tan as t the t as re the te the ast he and t the s s the thee thed the the thes the s te te he t the of in anae o the or

target: <he was in consequence put out of the protection of their internal law, end quote. their code was a subject of some curiosity.>  
prediction: <the the he at the t the an of t te the ale t he t te ar the in the the s the s tan as t the t as re the te the ast he and t the s s the thee thed the the thes the s te te he t the of in anae o the or

target: <that is why i occasionally leave this scene of action for a few days>  
prediction: <the the he at the t the an of t te the ale t he t te ar the in the the s the s tan ase athe t as re the te the ast he and t the s s the thee thed the the thes the s te te he t the of in anse o the or

target: <it probably contributed greatly to the general dissatisfaction which he exhibited with his environment,>  
prediction: <the the he at the t the an of t te the ale t he t te ar the in the the s the s tan as t the t as re the te the ast he and t the s s the thee thed the the thes the s te te he t the of in anae o the or

203/203 428s 1s/step - loss: 1.8276 - val\_loss: 1.5233

In practice, you should train for around 100 epochs or more.

Some of the predicted text at or around epoch 35 may look as follows:

```
target:      <as they sat in the car, frazier asked oswald where his lunch was>
prediction:  <as they sat in the car frazier his lunch ware mis lunch was>

target:      <under the entry for may one, nineteen sixty,>
prediction:  <under the introus for may monee, nin the sixty,>
```

[Terms](#) | [Privacy](#)

[Automatic Speech Recognition with Transformer](#)

- ◆ [Introduction](#)
- ◆ [Define the Transformer Input Layer](#)
- ◆ [Transformer Encoder Layer](#)
- ◆ [Transformer Decoder Layer](#)
- ◆ [Complete the Transformer model](#)
- ◆ [Download the dataset](#)
- ◆ [Preprocess the dataset](#)
- ◆ [Callbacks to display predictions](#)
- ◆ [Learning rate schedule](#)
- ◆ [Create & train the end-to-end model](#)