

# **Queensland Prosidy - Exploring the sounds of Queensland voices through machine learning**

**SLQ Wiki Fabrication Lab 2024/05/30 01:01**

This workshop is a bit of a provocation, intended to lead into deeper engagement with State Library's original oral history collection with the potential for a community developed AI Queensland Voice(s). It cover the existing state of the art, and how even extremely large models

## Queensland Prosidy - Exploring the sounds of Queensland voices through machine learning

In this workshop we'll look at how the unique accents, rhythms and tones of Queensland voices from State Library's collection fit into the burgeoning world of Text To Speech (TTS). We'll search the collection for some original audio sources, explore the resources available at The Edge for audio restoration and finally learn how to use free and open source machine learning software to create a realistic voice. We'll cover briefly the legal and ethical background to 'voice fakes' - and discover how hard it really is to make a machine speak [Strine](#).

### Voices in the Collection

How to find them? Digitised oral histories are a good bet, meaning a transcript should be available.

- “oral history digital”
- refine by “original materials”
- sort by date - “oldest to newest”

Lets go with [164 oral history interviews regarding the history of the Cape York Peninsula by interviewer Duncan Jackson](#)

### Finding Digital Audio and transcripts

Lets have a listen to Duncan Jackson's interview with Kathleen Jackson.

Click on [online access](#) to access the digitool viewer, then open:

- [mp3](#)
- [PDF transcript](#)

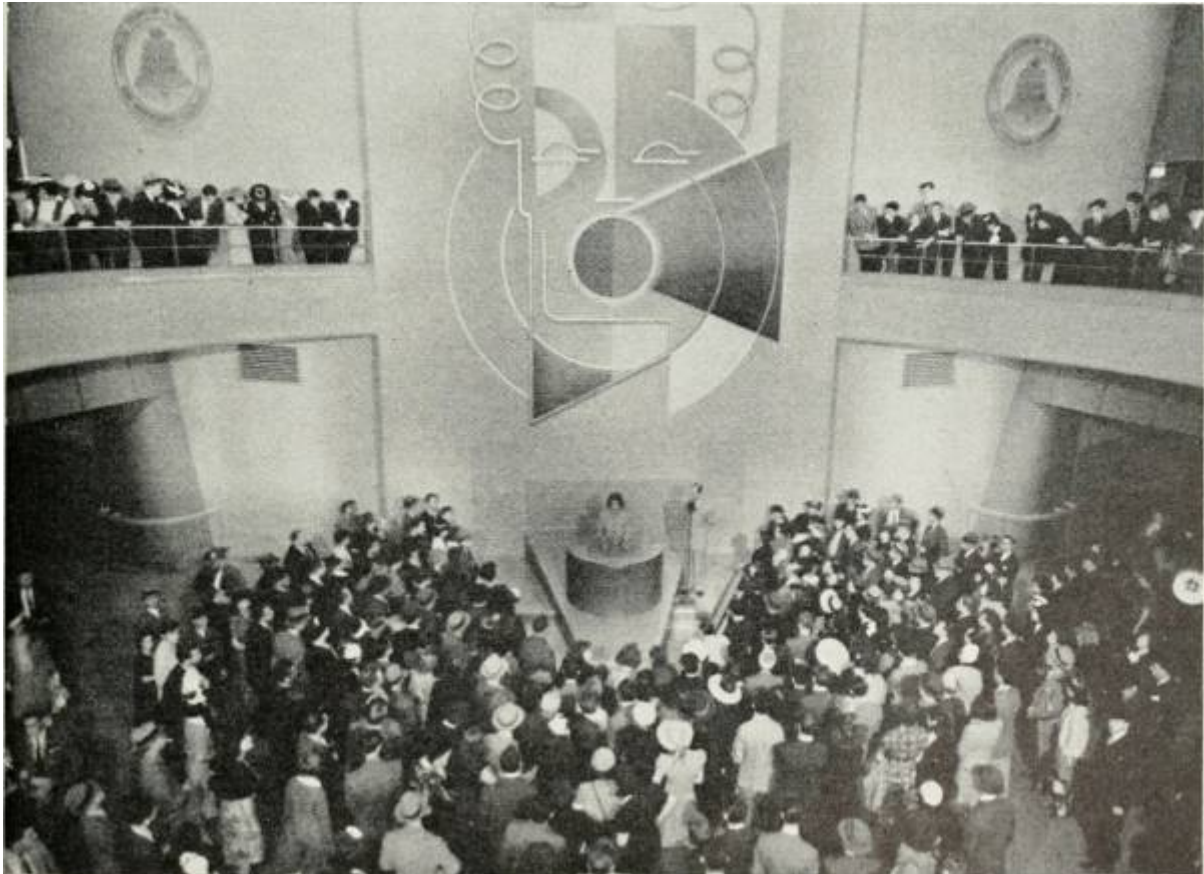
### Copyright and ethical considerations

We can always find the conditions of access and use on onesearch and the digitool viewer. In this case we have unrestricted access, and the material is in copyright.

You are free to use for personal research and study. For other uses contact [copyright@slq.qld.gov.au](mailto:copyright@slq.qld.gov.au)

# Speech Synthesis

Like many of the 20th century's technological innovations, the first modern speech synthesiser can be traced back to the invention of the [vocoder](#) at [Bell Labs](#). Derived from this, the [Voder](#) was demonstrated at the 1939 World Fair.



1)

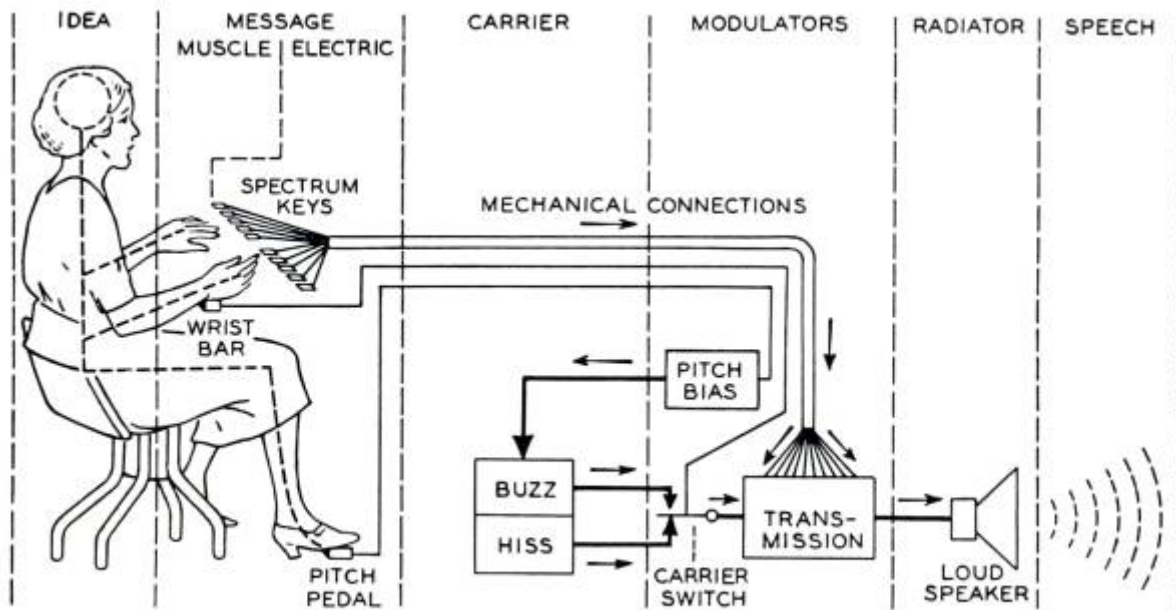


Fig. 8—Schematic circuit of the voder.

## Historical Audio Examples

Here is a playlist of various historical TTS methods.

<https://soundcloud.com/user-552764043>

## Modern State of the Art TTS

Now - it time to have some fun with TTS - check out the man holding the frog below...

<https://vo.codes/#speak>

And have a listen to some interesting examples from pop/meme culture.

<https://fifteen.ai/examples>

<https://www.youtube.com/watch?v=drirw-XvzzQ>

## Wavenet

Modern deep learning based synthesis started with the release of [Wavenet](#) in 2016 by Google's Deepmind.

WaveNet changes this paradigm by directly modelling the raw waveform of the audio signal, one sample at a time. As well as yielding more natural-sounding speech, using raw waveforms means that WaveNet can model any kind of audio, including music.<sup>2)</sup>

## Tacotron and Tacotron2

Wavenet was followed by Tacotron (also from Google) in 2017.

<https://google.github.io/tacotron/publications/tacotron/index.html>

Then Tacotron2

<https://ai.googleblog.com/2017/12/tacotron-2-generating-human-like-speech.html>

## The next wave - Diffusion

In April 2022 open.ai dropped [DALL-E 2](#), which uses diffusion models.

“Diffusion Models are generative models, meaning that they are used to generate data similar to the data on which they are trained. Fundamentally, Diffusion Models work by destroying training data through the successive addition of Gaussian noise, and then learning to recover the data by reversing this noising process. After training, we can use the Diffusion Model to generate data by simply passing randomly sampled noise through the learned denoising process.”

These models can be applied to TTS, with [tortoise-tts](#) producing some [excellent examples](#) of generated speech.

## Getting Started

## Google Colab

Google's Colaboratory<sup>3)</sup>, or “Colab” for short, allows you to write and execute Python in your browser, with

- Zero configuration required
- Free access to GPUs
- Easy sharing

## Python

Python is an open source programming language that was made to be easy-to-read and powerful<sup>(4)</sup>. Python is:

- a high-level language, (Meaning programmer can focus on what to do instead of how to do it.)
- an interpreted language (Interpreted languages do not need to be compiled to run.)
- is often described as a “batteries included” language due to its comprehensive standard library.

A program called an interpreter runs Python code on almost any kind of computer. In our case python will be interpreted by google colab, which is based on Jupyter notebooks.

## Jupyter Notebooks

Jupyter Notebook is an open-source web application that allows you to create and share documents that contain live code, equations, visualizations and narrative text<sup>(5)</sup>. Usually Jupyter notebooks require set-up for a specific purpose, but Colab takes care of all this for us.

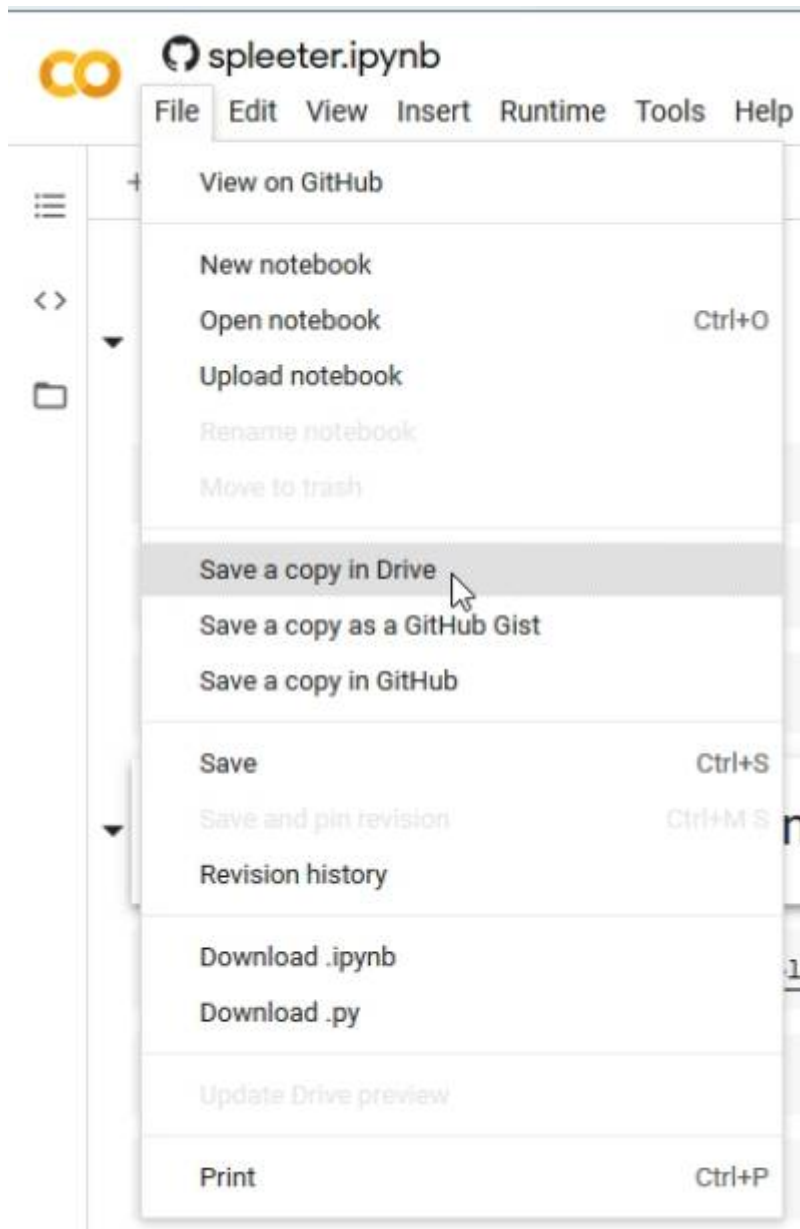
# Getting Started with Colab

The only requirment for using Colab is (unsurprisingly) a Google account. Once you have a google account, lets jump into our first ML example - [Spleeter](#) - that we mentioned earlier. Go to the Colab here:

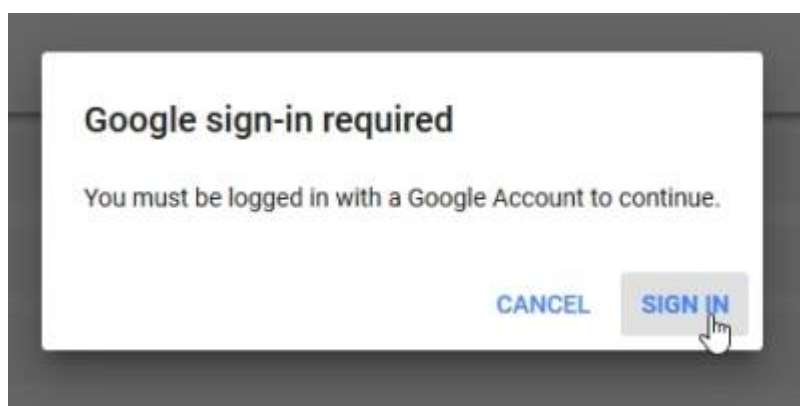
<https://colab.research.google.com/github/deezer/spleeter/blob/master/spleeter.ipynb>

## Making a Colab Copy

The first step is make a copy of the notebook to our Google drive - this means we can save any changes we like.

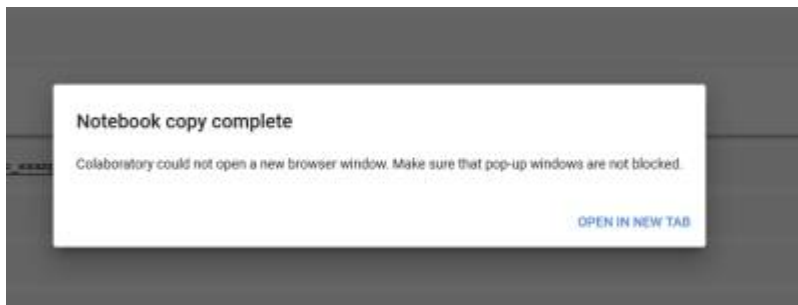


This will trigger a google sign-in



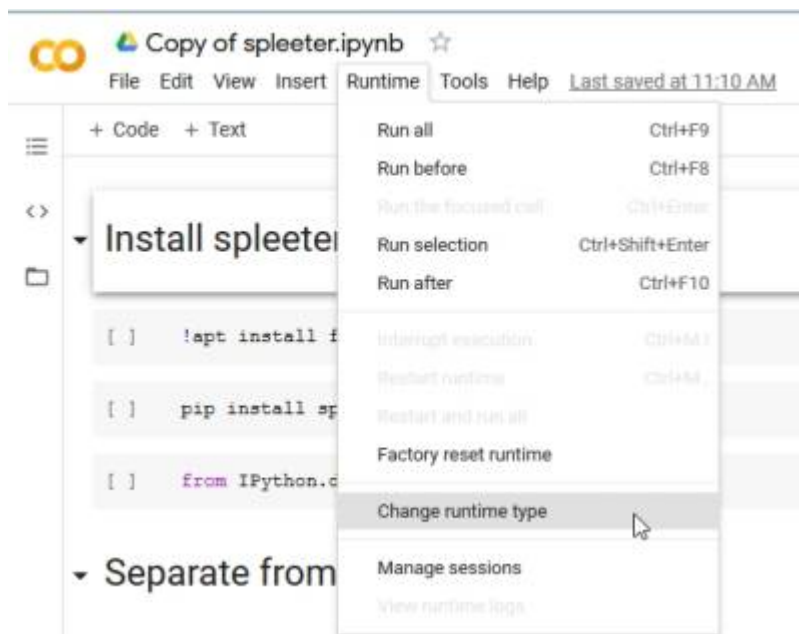


and the your copy will open in a new tab.



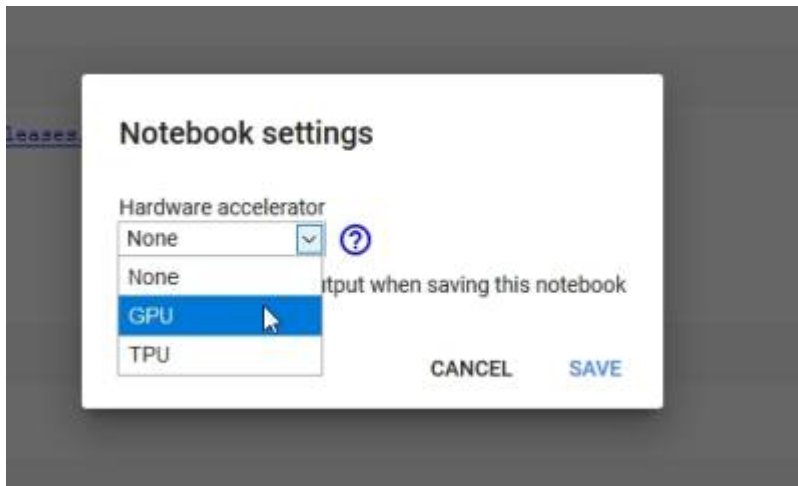
## Select a Runtime

Next we change our runtime (the kind or processor we use)

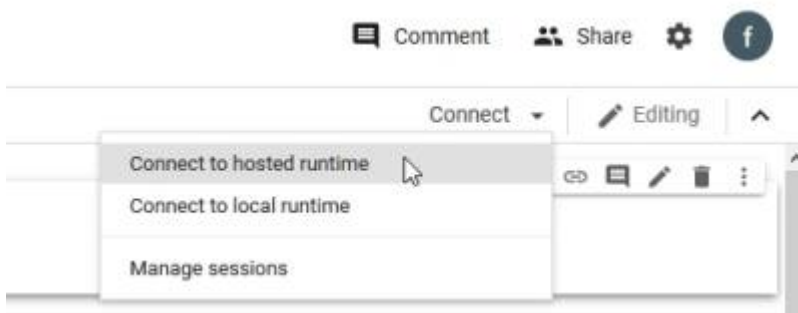


to a GPU to take advantage of Googles free GPU offer.





Now lets connect to our hosted runtime



and check the specs...



## Step Through the Notebook

Now its time to actually use the notebook! Before we start, lets go over how the notebooks work:

- The notebook is divided into sections, with each section made up of cells.
- These cells have code pre-entered into them,

- A play button on the runs (executes) the code in the cell.
- The output of the cell is printed (or displayed) directly below each cell.
- The output could be text, pictures, audio or video.

Cells usually contain python code, but can also be coded in bash - the UNIX command line shell. Cells containing bash commands start with an exclamation mark !

Our first section is called "Install Spleeter" and contains the bash command `apt install ffmpeg`. This installs ffmpeg in our runtime, which is used to process audio. Press the go button..

## ▼ Install spleeter



ffmpeg will be downloaded and installed to our runtime.

## Install spleeter



Next we will run a python command `pip` to use the [python package manager](#) to install the spleeter python package.

```

0 upgraded, 0 newly installed, 0 to remove and 35 not upgraded.

pip install spleeter

Collecting spleeter
  Downloading https://files.pythonhosted.org/packages/9a/c4/7d3a4af8bee9843a91582915e0c797c16e860b79c598ee3ec75677405/spleeter-1.5.4.tar.gz
Collecting ffmpeg-python
  Downloading https://files.pythonhosted.org/packages/d7/0c/5b8e52741f7b8ad4dc655991fab2e07b432d333da82c11ad701123888a/ffmpeg-python-0.2.0-py3-none-any.whl
Collecting norbert==0.2.1
  Downloading https://files.pythonhosted.org/packages/22/85/1e4f09c84d2b5541a4a8eece320902c4d2fa264dfe51f779548396f0fea/norbert-0.2.1-py2.py3-none-any.whl
Collecting pandas==0.25.1
  Downloading https://files.pythonhosted.org/packages/73/9b/52e228545d14f14bb2a1422e225f38463c8726645165e1cb7dde95bfe6d4/pandas-0.25.1-cp36-cp36m-manylinux1_x86_64.whl (10.5MB)
Requirement already satisfied: requests in /usr/local/lib/python3.6/dist-packages (from spleeter) (2.23.0)
Requirement already satisfied: setuptools==41.0.0 in /usr/local/lib/python3.6/dist-packages (from spleeter) (49.2.0)
Collecting librosa==0.7.2
  Downloading https://files.pythonhosted.org/packages/77/b5/1817862d64a7c231af15419d8418ee1f000742cac275e85c74b219cbcb/librosa-0.7.2.tar.gz (1.6MB)
Requirement already satisfied: numba==0.48.0 in /usr/local/lib/python3.6/dist-packages (from spleeter) (0.48.0)
Collecting tensorflow==1.15.2
  Downloading https://files.pythonhosted.org/packages/9a/d9/f4234c7bf68638423f8e7f44af7fcfc3bca416b51e6d902391e47ec43/tensorflow-1.15.2-cp36-cp36m-manylinux2010_x86_64.whl (110.5MB)
Collecting importlib_resources
  Downloading https://files.pythonhosted.org/packages/ba/03/0f9595c0a2ef12590877f3c17e5f579759ce5cafb176825645dcb89a1177/importlib_resources-3.0.0-py2.py3-none-any.whl
Requirement already satisfied: future in /usr/local/lib/python3.6/dist-packages (from ffmpeg-python->spleeter) (0.16.0)
Requirement already satisfied: scipy in /usr/local/lib/python3.6/dist-packages (from norbert==0.2.1->spleeter) (1.4.1)
Requirement already satisfied: python-dateutil==2.6.1 in /usr/local/lib/python3.6/dist-packages (from pandas==0.25.1->spleeter) (2.8.1)
Requirement already satisfied: numpy==1.13.3 in /usr/local/lib/python3.6/dist-packages (from pandas==0.25.1->spleeter) (1.16.5)
Requirement already satisfied: pytz==2017.2 in /usr/local/lib/python3.6/dist-packages (from pandas==0.25.1->spleeter) (2018.9)
Requirement already satisfied: idna<3,>=2.5 in /usr/local/lib/python3.6/dist-packages (from requests->spleeter) (2.10)
Requirement already satisfied: urllib3<1.25.0,>=1.25.1,<1.26,>=1.21.1 in /usr/local/lib/python3.6/dist-packages (from requests->spleeter) (1.24.3)
Requirement already satisfied: certifi==2017.4.17 in /usr/local/lib/python3.6/dist-packages (from requests->spleeter) (2020.6.20)
Requirement already satisfied: chardet<4,>=3.0.2 in /usr/local/lib/python3.6/dist-packages (from requests->spleeter) (3.0.4)
Requirement already satisfied: audioread==2.0.0 in /usr/local/lib/python3.6/dist-packages (from librosa==0.7.2->spleeter) (2.1.8)
Requirement already satisfied: scikit-learn==0.19.0,>=0.14.0 in /usr/local/lib/python3.6/dist-packages (from librosa==0.7.2->spleeter) (0.22.2.post1)
Requirement already satisfied: joblib==0.12 in /usr/local/lib/python3.6/dist-packages (from librosa==0.7.2->spleeter) (0.16.0)

```

This will take a while - and at the end we will get a message saying we need to restart our runtime due to some compatibility issues <sup>6)</sup>

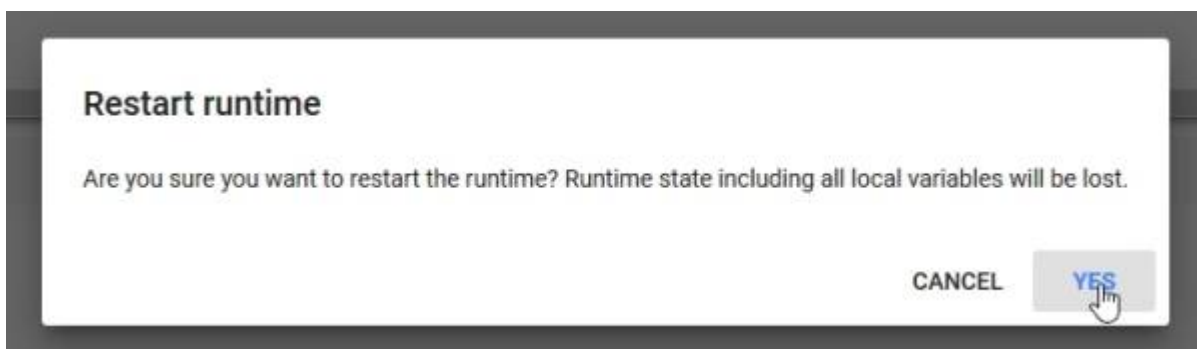
```

Uninstalling tensorboard-2.3.0:
  Successfully uninstalled tensorboard-2.3.0
Found existing installation: gast 0.3.3
Uninstalling gast-0.3.3:
  Successfully uninstalled gast-0.3.3
Found existing installation: tensorflow-estimator 2.3.0
Uninstalling tensorflow-estimator-2.3.0:
  Successfully uninstalled tensorflow-estimator-2.3.0
Found existing installation: tensorflow 2.3.0
Uninstalling tensorflow-2.3.0:
  Successfully uninstalled tensorflow-2.3.0
Successfully installed ffmpeg-python-0.2.0 gast-0.2.2 importlib-resources-3.0.0 keras-applications
WARNING: The following packages were previously imported in this runtime:
[pandas]
You must restart the runtime in order to use newly installed versions.

RESTART RUNTIME

```

Go ahead and restart



Next is another bash command

wget

we use to (web)get our example audio file.



```
!wget https://github.com/deezer/spleeter/raw/master/audio_example.mp3
--2020-08-04 01:17:58-- https://raw.githubusercontent.com/deezer/spleeter/master/audio_example.mp3
Resolving raw.githubusercontent.com (raw.githubusercontent.com)... 151.101.0.133, 151.101.64.133, 151.101.128.133, ...
Connecting to raw.githubusercontent.com (raw.githubusercontent.com)|151.101.0.133|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 262867 (257K) [audio/mpeg]
Saving to: 'audio_example.mp3'

audio_example.mp3 100%[=====>] 256.71K --.-KB/s in 0.02s

2020-08-04 01:17:58 (13.5 MB/s) - 'audio_example.mp3' saved [262867/262867]
```

And the next cell uses the python Audio command to give us a nice little audio player so we can hear our example.



Now its finally time to use the spleeter tool with the separate command <sup>7)</sup> as `!spleeter separate`, and lets pass the `-h` flag <sup>8)</sup> to show us the built in help for the command.

```
!spleeter separate -h
```

usage: spleeter separate [-h] [-a AUDIO\_ADAPTER] [-p CONFIGURATION]
 [--verbose] -i INPUTS [INPUTS ...] [-o OUTPUT\_PATH]
 [-f FILENAME\_FORMAT] [-d DURATION] [-s OFFSET]
 [-c {wav,mp3,ogg,m4a,wma,flac}] [-b BITRATE] [-m]
 [-B {tensorflow,librosa,auto}]

optional arguments:

- h, --help show this help message and exit
- a AUDIO\_ADAPTER, --adapter AUDIO\_ADAPTER
 Name of the audio adapter to use for audio I/O
- p CONFIGURATION, --params\_filename CONFIGURATION
 JSON filename that contains params
- verbose Shows verbose logs
- i INPUTS [INPUTS ...], --inputs INPUTS [INPUTS ...]
 List of input audio filenames
- o OUTPUT\_PATH, --output\_path OUTPUT\_PATH
 Path of the output directory to write audio files in
- f FILENAME\_FORMAT, --filename\_format FILENAME\_FORMAT
 Template string that will be formatted to
 generated output filename. Such template should be
 Python formattable string, and could use {filename},
 {instrument}, and {codec} variables.
- d DURATION, --duration DURATION
 Set a maximum duration for processing audio (only
 separate offset + duration first seconds of the input
 file)
- s OFFSET, --offset OFFSET
 Set the starting offset to separate audio from.
- c {wav,mp3,ogg,m4a,wma,flac}, --codec {wav,mp3,ogg,m4a,wma,flac}
 Audio codec to be used for the separated output
- b BITRATE, --birate BITRATE
 Audio bitrate to be used for the separated output
- m, --mwf
 Whether to use multichannel Wiener filtering for
 separation
- B {tensorflow,librosa,auto}, --stft-backend {tensorflow,librosa,auto}
 Who should be in charge of computing the stfts.
 Librosa is faster than tensorflow on CPU and uses less
 memory. "auto" will use tensorflow when GPU
 acceleration is available and librosa when not.

Now that we know what we are doing - we run the tool for real, and will use the `-i` flag to define the input as our downloaded example, and the `-o` flag to define our output destination as the directory (folder) output. By default spleeter will download and use the [2stems model](#).

```
!spleeter separate -i audio_example.mp3 -o output/
```

INFO:spleeter:Downloading model archive <https://github.com/deezer/spleeter/releases/download/v1.4.0/2stems.tar.gz>

INFO:spleeter:Validating archive checksum

INFO:spleeter:Extracting downloaded 2stems archive

INFO:spleeter:2stems model file(s) extracted

INFO:spleeter:File output/audio\_example/accompaniment.wav written successfully

INFO:spleeter:File output/audio\_example/vocals.wav written successfully

Another bash command `ls` (list) shows us the contents of our output directory



```
!ls output/audio_example  
accompaniment.wav  vocals.wav
```

And finally another couple of audio commands to hear our result!

```
[10] Audio('output/audio_example/vocals.wav')
```



```
Audio('output/audio_example/accompaniment.wav')
```



## Things to try

Check out the [usage instructions](#) for the separate tool on the Github site and try your own 4stem and 5stem separations. Use your own audio files to test the separation.

# Speech to Text with Mozilla DeepSpeech

Our next challenge will be to adapt the latest version of Mozilla's DeepSpeech for use in Google Colab.

We will be using the documentation here:

<https://deepspeech.readthedocs.io/en/v0.8.0/USING.html#getting-the-pre-trained-model>

To adapt this colab notebook to run the latest version of Mozilla DeepSpeech:

<https://colab.research.google.com/github/tugstugi/dl-colab-notebooks/blob/master/notebooks/MozillaDeepSpeech.ipynb#scrollTo=4OAYwPHApuz>

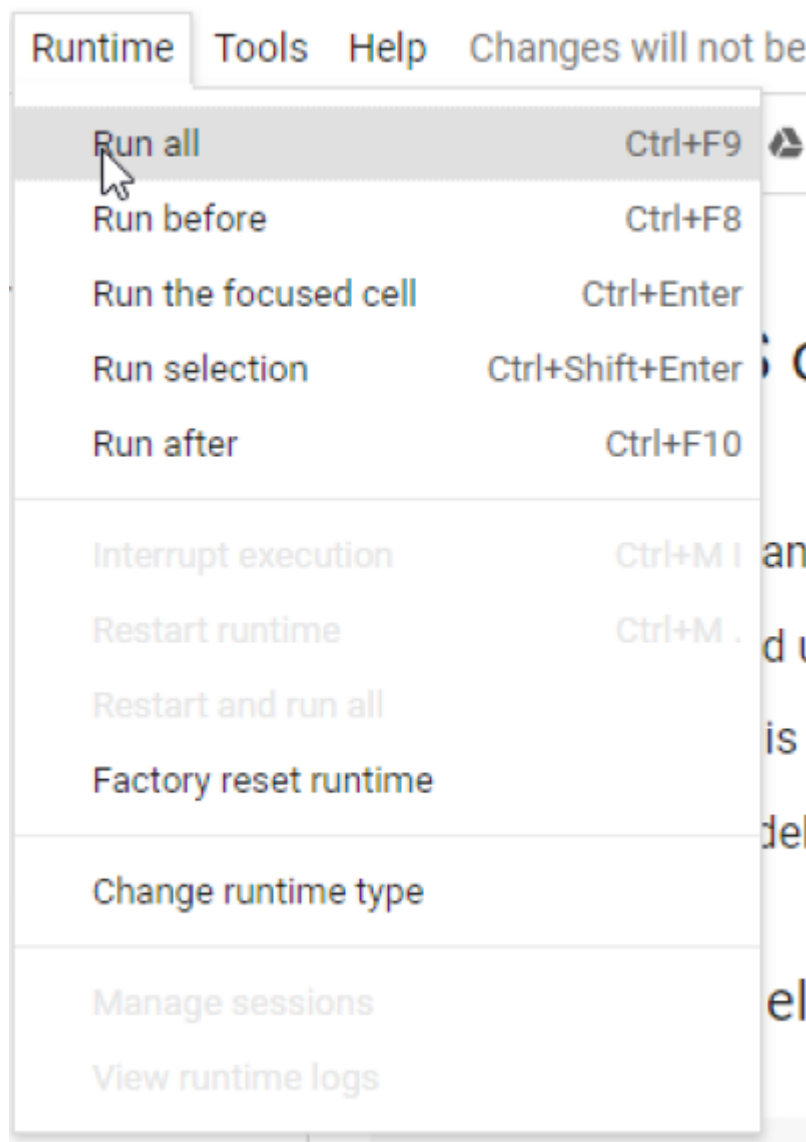
## Text to Speech with Mozilla TTS

Our final example is TTS with Mozilla TTS:

[https://colab.research.google.com/drive/1u\\_16ZzHjKYFn1HNVuA4Qf\\_i2MMFB9oIY?usp=sharing#scrollTo=6LWsNd3\\_M3MP](https://colab.research.google.com/drive/1u_16ZzHjKYFn1HNVuA4Qf_i2MMFB9oIY?usp=sharing#scrollTo=6LWsNd3_M3MP)

You can dive straight into this and use it to generate speech. This example uses Tacotron2 and MultiBand-Melgan models and LJSpeech dataset.

## Run All Cells



## Generate Speech





## Going Further

ML is such a big and fast moving area of research there are countless other ways to explore and learn, here are a few two-minute videos to pique your interest:

- [Video restoration](#)
- [OpenAI Plays Hide and Seek](#)

Make sure you check out the resources in Lynda, which you will have free access to as a State Library of Queensland member

## Links

<https://machinelearningforkids.co.uk/#!/links#top>

<https://experiments.withgoogle.com/collection/ai>

<https://openai.com/blog/>

1)

By Internet Archive Book Images - <https://www.flickr.com/photos/internetarchivebookimages/14776509983/Source> book page:

[https://archive.org/stream/belltelephonemag19amerrich/belltelephonemag19amerrich#page/n78/mode/1upReference\[Fig.4\]](https://archive.org/stream/belltelephonemag19amerrich/belltelephonemag19amerrich#page/n78/mode/1upReference[Fig.4]) The Voder Fascinates the Crowds from: Williams, Thomas W. (January 1940) I. At the New York World's Fair. "Our Exhibits at Two Fairs". Bell Telephone Quarterly XIX (1): 65. "The Voder Fascinates the Crowds - The manipulative skill of the operator's fingers makes the Voder's voice almost too good to be true"; No restrictions, <https://commons.wikimedia.org/w/index.php?curid=43343073>

2)

<https://deepmind.com/blog/article/wavenet-generative-model-raw-audio>

3)

<https://colab.research.google.com/notebooks/intro.ipynb>

4)

[https://simple.wikipedia.org/wiki/Python\\_\(programming\\_language\)](https://simple.wikipedia.org/wiki/Python_(programming_language))

5)

<https://jupyter.org/>

6)

this is not unusual when using a hosted runtime

7)

confusingly we need to call it from bash (with the exclamation

8)

a fancy way of saying option