# Issues with AudioLottery code

May 5, 2024

We tried to replicate the Conformer part of AudioLottery on Google colab and we encountered the following issues. (Link to Colab notebook)

## 1    Problem installing 'ctcdecode' module

We ran the following lines of code to install the module.

```
# get the code
! git clone --recursive https://github.com/parlance/ctcdecode.git
%cd ctcdecode
! pip install .
```

This gave us the following error:



Figure 1: Error

## 1.1 Fix

'ctcdecode' needs torch 1.11.0 while we were using torch 2.3.0. Thus, the following line fixes the error:

```
! pip install torch==1.11.0
```



Figure 2: Fix

# 2 Problem with 'warp-rnnt' module

When we tried to run the main_lth.py file, we got the the following error:



Figure 3: Fix