Frame-Level Speech Recognition

FXPFRIMENT

- Build LA-DNN
 - A powerful model for speech recognition and good candidate to replace FF-DNN
 - Helps to remove the vanshing gradient problems in deep network
 - https://ieeexplore.ieee.org/document/7472646
- No Augmentation
- No hyper-parameter tuning however used our past knowledge to set hyper-paramters
- Advanced weight initialization, Label smoothing in loss function, gradient clipping and scheduler are used.
- 1 !pip install torch-summary --quiet
- 1 import torch
- import numpy as np
- import sklearn
- import qc
- 5 import zipfile
- import pandas as pd
- 7 from tqdm.auto import tqdm
- import os 8
- 9 import datetime
- 10 import torchsummary
- device = 'cuda' if torch.cuda.is available() else 'cpu' 11
- 12 print("Device: ", device)

```
→ Device: cuda
```

- 1 # ### If you are using colab, you can import google drive
- 2 from google.colab import drive
- drive.mount('/content/drive')
- 1!pip install --upgrade --force-reinstall --no-deps kaggle==
- 2 !mkdir -p /root/.kaggle
- → Collecting kaggle==1.5.8 Downloading kaggle-1.5.8.tar.gz (59 kB)

- 59.2/59.2 kB 1.8 MB/s eta 0:00:

```
Preparing metadata (setup.py) ... done
   Building wheels for collected packages: kaggle
     Building wheel for kaggle (setup.py) ... done
     Created wheel for kaggle: filename=kaggle-1.5.8-py3-none-any.whl size=73249
     Stored in directory: /root/.cache/pip/wheels/b5/23/bd/d33cbf399584fa44fa049
   Successfully built kaggle
   Installing collected packages: kaggle
     Attempting uninstall: kaggle
       Found existing installation: kaggle 1.7.4.2
       Uninstalling kaggle-1.7.4.2:
        Successfully uninstalled kaggle-1.7.4.2
   Successfully installed kaggle-1.5.8
  1 with open("/root/.kaggle/kaggle.json", "w+") as f:
  2
        f.write('{"username":"dineshbuswala","key":"a67fefaecac
  4 !chmod 600 /root/.kaggle/kaggle.json
  1 # commands to download data from kaggle
  2
  3 !kaggle competitions download -c 11785-hw1p2-f23 --force
  4 !mkdir -p '/content/data'
  5
  6 # !unzip -qo /content/11785-hw1p2-f23.zip -d '/content/data
Downloading 11785-hw1p2-f23.zip to /kaggle/working
                                            3.99G/3.99G [00:24<00:00, 219MB
   100%
   100%
                                            3.99G/3.99G [00:24<00:00, 174MB
  1 !unzip -qo /kaggle/working/11785-hw1p2-f23.zip -d '/content
  1 # ! rm -r /content/data/11-785-f23-hw1p2/test-clean
  2 # ! rm -r /content/11785-hw1p2-f23.zip
      ! ls /content/data/11-785-f23-hw1p2/train-clean-100/mfcc/
→ 28539
  1 ### configuration variables
  2 EPOCHS
                                        = 4
  3 BATCH SIZE
                                        = 2048 * 2
  4 LEFT CONTEXT
                                        = 7
  5 RIGHT CONTEXT
                                        = 7
  6 INITIAL LEARNING RATE
                                        = 1e-3
  7 L2 PENALTY
                                        = 1e-5
  8 STEP SIZE
```

```
9 GAMMA
                                   = 0.1
10 BASE DIRECTORY
                                   = '/content/data/11-785-f23
                                   = BASE_DIRECTORY + 'train-c
11 TRAINING DATA
                                   = BASE DIRECTORY + 'dev-cle
12 EVALUATION DATA
                                   = ['[SIL]', 'AA',
13 PHONEMES
                                                          'D',
                                        'B',
                                                 'CH',
14
                                        'F',
                                                'G',
                                                         'HH',
15
                                                'M',
                                        'L',
                                                         'N',
16
                                                'S',
                                        'R',
                                                         'SH',
17
                                                'W',
                                        ١٧١,
18
                                                         'Y',
                                    = {phoneme: idx for idx, p
19 PHONEMES TO INDEX
20 NUMBER OF NEURONS
                                    = [2048, 2048, 1024, 1024,
21 MODEL DIR
                                    = "/content"
22 CLIP VALUE
                                    = 1.0
23 LABEL SMOOTHING
                                    = 0.01
  1 class AudioDataset(torch.utils.data.Dataset):
       def init (self, root, phonemes = PHONEMES TO INDEX,
  2
                     left context = LEFT CONTEXT,
  3
  4
                     right context = RIGHT CONTEXT):
  5
            self.left context = left context
  6
            self.right context = right context
  7
            self.phonemes mapping = phonemes
  8
  9
            self.mfcc dir = os.path.join(root, 'mfcc')
 10
            self.transcript dir = os.path.join(root, 'transcri
 11
 12
 13
           # List and sort mfcc and transcript files
 14
           mfcc names = sorted(os.listdir(self.mfcc dir))
 15
            transcript names = sorted(os.listdir(self.transcri
 16
 17
           # Sanity check
            assert len(mfcc names) == len(transcript names), "
 18
 19
 20
           total frames = 0
 21
 22
           for i in range(len(mfcc names)):
                mfcc path = os.path.join(self.mfcc dir, mfcc n
 23
 24
 25
                mfcc = np.load(mfcc path,
                               allow pickle = False,
 26
                               mmap mode='r')
 27
                total frames += mfcc.shape[0]
 28
 29
                del mfcc
```

```
30
31
           sample mfcc = np.load(os.path.join(self.mfcc dir,
32
          self.mfcc dim = sample mfcc.shape[1]
33
          ### Right Padding is added here automatically
          self.mfccs = np.zeros((total frames + right contex
34
          self.transcripts = [None] * (total_frames + right_
35
36
          ### Release memory
          del sample mfcc, total frames
37
38
          gc.collect()
39
          index = 0
40
41
          for i in range(len(mfcc names)):
              mfcc = np.load(os.path.join(self.mfcc dir, mfc
42
                              allow pickle = False,
43
44
                              mmap mode = 'r')
45
46
               mfcc = (mfcc - mfcc.mean(axis = 1, keepdims =
47
              transcript = np.load(os.path.join(self.transcr
48
49
                                    allow pickle = False,
                                    mmap_mode='r')[1:-1]
50
51
52
               self.mfccs[index: index + mfcc.shape[0]] = mfc
53
               self.transcripts[index: index + len(transcript
54
               index += mfcc.shape[0]
              del mfcc, transcript
55
56
               if i % 1000 == 0:
57
                   qc.collect() ### Release memory
58
59
60
          # Save original dataset length (before adding righ
          self.length = len(self.mfccs) - right context
61
62
63
          # Map transcript phonemes to indices
64
           self.transcripts = [self.phonemes mapping.get(p, -
65
66
          ### Release memory
67
          del mfcc names, transcript names
68
          gc.collect()
69
      def len (self):
70
71
          return self.length
72
      def getitem (self, ind):
73
          if ind < self.left context: # index is less than</pre>
74
```

```
EXPERIMENT II - Colab
14/06/2025, 13:02
                    # zeros need to prepend with it
     75
     76
                    padding = []
                    for _ in range(self.left_context - ind):
     77
                        padding.append(np.zeros((1, self.mfcc dim)
     78
     79
     80
                    for i in range(ind):
                        padding.append(self.mfccs[i][None, :])
     81
     82
     83
                    # Include current and right context frames
                    current = self.mfccs[ind][None, :] # Ensure s
     84
                    right context = [self.mfccs[i][None, :] for i
     85
                    frames = np.concatenate(padding + [current] +
     86
     87
               else: # when index is greater than or equal to le
     88
     89
                    left context frames = self.mfccs[ind - self.le
     90
                    current = self.mfccs[ind][None, :] # Ensure s
                    right context frames = self.mfccs[ind + 1: ind
     91
     92
     93
                    frames = np.concatenate([left context frames,
     94
     95
                frames = frames.flatten() # Flatten to get 1D dat
     96
                frames = torch.FloatTensor(frames) # Convert to t
               phonemes = torch.tensor(self.transcripts[ind], dty
     97
     98
     99
                return frames, phonemes
    100
     1# Create a dataset object using the AudioDataset class for
     2 train data = AudioDataset(TRAINING DATA)
     3
     4 # Create a dataset object using the AudioDataset class for
     5 val data = AudioDataset(EVALUATION DATA)
     6
     1 train loader = torch.utils.data.DataLoader(
     2
          dataset = train data,
     3
          num workers = 2,
          batch_size = BATCH_SIZE,
     4
          pin_memory = True,
     5
          shuffle = True
     6
     7)
     8
     9 val loader = torch.utils.data.DataLoader(
          dataset = val data,
    10
    11
          num workers = 1,
```

21

22 23 for m in [self.U] + [self.V] + [self.T]:

if isinstance(m, torch.nn.Linear):

Kaiming initialization for weights

```
torch.nn.init.kaiming_uniform_(m.weight, no
24
25
26
      def forward(self, x):
27
          t x = self.T(x)
          u x = self.U(x)
28
          v x = self.V(self.activation(u x))
29
          return v_x + t_x
30
31
32
33
34 class Network(torch.nn.Module):
      def __init__(self, input_size):
35
          super(Network, self). init ()
36
37
38
          # Neurons in each layer: input -> hidden(s) -> outp
          self.neurons = [input size] + NUMBER OF NEURONS + [
39
40
41
          layers = []
          for in features, out features in zip(self.neurons[:
42
               layers.append(LADNNLayer(in features, out featu
43
44
45
          # Final laver
          layers.append(LADNNLayer(self.neurons[-2], self.neu
46
47
          # Combine all into a sequential model
48
          self.model = torch.nn.Sequential(*layers)
49
50
51
          # Apply weight initialization
          self. initialize weights()
52
53
      def initialize weights(self):
54
          print('Initialization of weights using Kaiming')
55
          for m in self.model:
56
               if isinstance(m, torch.nn.Linear):
57
                   # Kaiming initialization for weights
58
                   torch.nn.init.kaiming uniform (m.weight, no
59
60
61
62
      def forward(self, x):
          out = self.model(x)
63
64
          return out
65
 1 INPUT SIZE = (LEFT CONTEXT + RIGHT CONTEXT + 1) * 28
 2 model = Network(INPUT SIZE).to(device)
```

3 # Pass the input size as a tuple, without the batch dimensi 1 torchcummary cummary/model / TNDHT CT7E \\

→ Initialization of weights using Kaiming Initialization of weights using Kaiming

Layer (type)	Output Shape	Param #
Linear-1	[-1, 2048]	862,208
Linear-2	[-1, 210]	88,410
ReLU-3	[-1, 210]	0
Linear-4	[-1, 2048]	432,128
LADNNLayer-5	[-1, 2048]	0
Linear-6	[-1, 2048]	4,196,352
Linear-7	[-1, 1024]	2,098,176
ReLU-8	[-1, 1024]	0
Linear-9	[-1, 2048]	2,099,200
LADNNLayer-10	[-1, 2048]	0
Linear-11	[-1, 1024]	2,098,176
Linear-12	[-1, 1024]	2,098,176
ReLU-13	[-1, 1024]	0
Linear-14	[-1, 1024]	1,049,600
LADNNLayer-15	[-1, 1024]	0
Linear-16	[-1, 1024]	1,049,600
Linear-17	[-1, 512]	524,800
ReLU-18	[-1, 512]	0
Linear-19	[-1, 1024]	525,312
LADNNLayer-20	[-1, 1024]	0
Linear-21	[-1, 512]	524,800
Linear-22	[-1, 512]	524,800
ReLU-23	[-1, 512]	0
Linear-24	[-1, 512]	262,656
LADNNLayer-25	[-1, 512]	0
Linear-26	[-1, 256]	131,328
Linear-27	[-1, 256]	131,328
ReLU-28	[-1, 256]	65 703
Linear-29	[-1, 256]	65,792 0
LADNNLayer-30 Linear-31	[-1, 256] [-1, 256]	65,792
Linear-31 Linear-32	[-1, 230]	32,896
ReLU-33	[-1, 128]	0
Linear-34	[-1, 126]	33,024
LADNNLayer-35	[-1, 256]	0
Linear-36	[-1, 40]	10,280
Linear-37	[-1, 128]	32,896
ReLU-38	[-1, 128]	0
Linear-39	[-1, 40]	5,160
LADNNLayer-40	[-1, 40]	0
/	r = / · · · · · ·	

Total params: 18,942,890 Trainable params: 18,942,890

Non-trainable params: 0

```
1 criterion = torch.nn.CrossEntropyLoss(label smoothing = LAB
  2 # We use CE because the task is multi-class classification
  3
  4 optimizer = torch.optim.Adam(model.parameters(),
                                 lr = INITIAL LEARNING RATE,
  5
                                 weight decay = L2 PENALTY) #De
  6
 7 scheduler = torch.optim.lr scheduler.StepLR(
       optimizer, step size = STEP SIZE, gamma = GAMMA
  9)
 10 # Refer - https://pytorch.org/docs/stable/notes/amp example
    torch.cuda.empty_cache()
  2 qc.collect()
→ 101
     def train(model, dataloader, optimizer, criterion):
  1
  2
  3
         model.train()
         tloss, tacc = 0, 0 # Monitoring loss and accuracy
  4
         batch bar = tqdm(total=len(train loader), dynamic ncc
  5
  6
  7
         for i, (frames, phonemes) in enumerate(dataloader):
  8
  9
             ### Initialize Gradients
 10
             optimizer.zero grad()
 11
 12
             ### Move Data to Device (Ideally GPU)
             frames = frames.to(device)
 13
             phonemes = phonemes.to(device)
 14
 15
 16
             ### Forward Propagation
 17
             logits = model(frames)
 18
 19
             ### Loss Calculation
 20
             loss = criterion(logits, phonemes)
 21
             ### Backward Propagation
 22
             loss.backward()
 23
 24
 25
             ### Clip gradients
             torch.nn.utils.clip_grad_norm_(model.parameters(),
 26
 27
 28
             ### Gradient Descent
```

```
29
            optimizer.step()
30
31
                    += loss.item()
            tloss
                    += torch.sum(torch.argmax(logits, dim= 1) =
32
            tacc
33
            batch bar.set postfix(loss="{:.04f}".format(float())
34
                                  acc="{:.04f}%".format(float()
35
            batch bar.update()
36
37
38
            ### Release memory
            del frames, phonemes, logits
39
            torch.cuda.empty cache()
40
41
        batch bar.close()
42
        tloss /= len(train loader)
43
        tacc /= len(train loader)
44
45
46
        return tloss, tacc
 1 def eval(model, dataloader):
 2
 3
      model.eval() # set model in evaluation mode
      vloss, vacc = 0, 0 # Monitoring loss and accuracy
 4
      batch bar = tqdm(total=len(val loader), dynamic ncols
 5
 6
 7
      for i, (frames, phonemes) in enumerate(dataloader):
 8
          ### Move data to device (ideally GPU)
 9
                      = frames.to(device)
10
          frames
                     = phonemes.to(device)
11
          phonemes
12
13
          # makes sure that there are no gradients computed a
14
          with torch.inference mode():
15
              ### Forward Propagation
              logits = model(frames)
16
              ### Loss Calculation
17
              loss = criterion(logits, phonemes)
18
19
20
          vloss += loss.item()
                  += torch.sum(torch.argmax(logits, dim= 1) =
21
          vacc
22
23
          batch bar.set postfix(loss="{:.04f}".format(float(v
24
                                 acc="{:.04f}%".format(float(v
25
26
          batch bar.update()
```

```
27
 28
            ### Release memory
 29
            del frames, phonemes, logits
            torch.cuda.empty cache()
 30
 31
 32
       batch bar.close()
       vloss /= len(val_loader)
 33
               /= len(val_loader)
 34
       vacc
 35
 36
       return vloss, vacc
  1 best_model_path = os.path.join(MODEL_DIR, "best_model.pt")
  2 \text{ best acc} = -np.inf
  3 for epoch in range(EPOCHS):
       ### clean up memory before computation
       torch.cuda.empty cache()
  5
       gc.collect()
  6
  7
  8
  9
       print(f"\nEpoch {epoch + 1}/{EPOCHS}")
 10
 11
       curr lr = float(optimizer.param groups[0]['lr'])
       train loss, train acc = train(model, train loader, opti
 12
       val loss, val acc = eval(model, val loader)
 13
 14
 15
        print(f"\tTrain Acc: {train acc*100:.2f}%\tTrain Loss:
       print(f"\tVal Acc: {val acc*100:.2f}%\tVal Loss: {v
 16
 17
 18
       # Save model at every epoch
       epoch model path = os.path.join(MODEL DIR, f"model at e
 19
       torch.save(model.state dict(), epoch model path)
 20
 21
 22
       # Save best model
 23
        if val acc > best acc:
 24
            best acc = val acc
            torch.save(model.state dict(), best model path)
 25
 26
            print(f"Updated Best Model at: {best model path}")
 27
 28
       ### take step in adjusting the learning rate
 29
       scheduler.step()
 30
 31
\overline{\Rightarrow}
```

```
Epoch 1/4
   Train: 0%| | 0/8812 [00:00<?, ?it/s] Val: 0%| | 0/471 [00:00<?, ?it/s]
          Train Acc: 67.04% Train Loss: 13.0031 LR: 0.0010000 Val Acc: 69.21% Val Loss: 1.0357
   Updated Best Model at: /content/best model.pt
   Epoch 2/4
   Train: 0%|
                    | 0/8812 [00:00<?, ?it/s]
   Val: 0% | 0/471 [00:00<?, ?it/s]
           Train Acc: 72.84% Train Loss: 0.9069 LR: 0.0010000 Val Acc: 70.96% Val Loss: 0.9690
   Updated Best Model at: /content/best model.pt
   Epoch 3/4
   Train: 0%
                      | 0/8812 [00:00<?, ?it/s]
   Val: 0% | 0/471 [00:00<?, ?it/s]
          Train Acc: 77.41% Train Loss: 0.7553 LR: 0.0001000 Val Acc: 74.43% Val Loss: 0.8521
   Updated Best Model at: /content/best model.pt
   Epoch 4/4
   Train: 0%
                      | 0/8812 [00:00<?, ?it/s]
   Val: 0%| | 0/471 [00:00<?, ?it/s]
          Train Acc: 78.37% Train Loss: 0.7237 LR: 0.0001000 Val Acc: 74.58% Val Loss: 0.8479
   Undated Rest Model at: /content/hest model nt
  1
  2 # Load best model after training
  3 model.load state dict(torch.load('/content/best model.pt'))
  4
<All keys matched successfully>
  1 from sklearn.metrics import classification report, confusion
      model.eval() # Set model in evaluation mode
  3 predicted = []
      groundtruth = []
  5
  6
      for frames, phonemes in val loader:
  7
  8
           # Move data to device
  9
           frames = frames.to(device)
 10
           phonemes = phonemes.to(device)
 11
 12
           # Disable gradient calculation
 13
           with torch.inference mode():
 14
                logits = model(frames)
 15
 16
           predict = torch.argmax(logits, dim = 1)
```

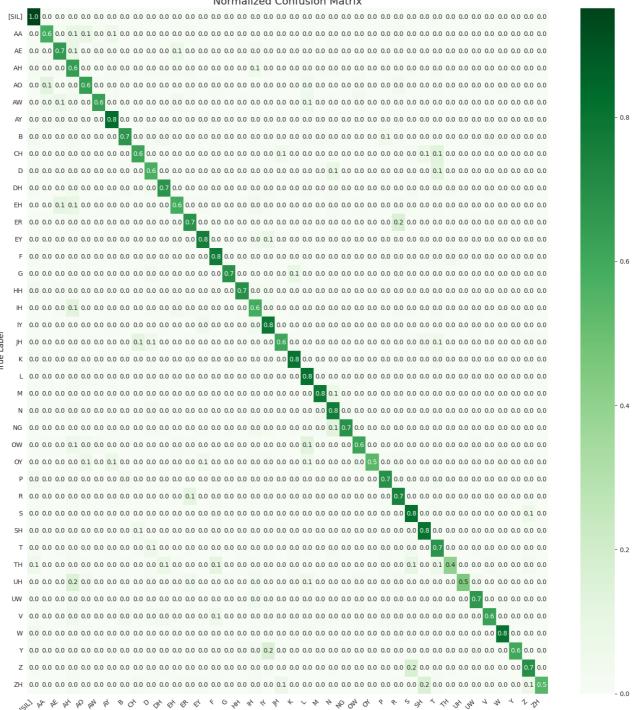
```
# Detach and move to CPU for evaluation
 18
 19
            predicted.extend(predict.detach().cpu().tolist())
 20
            groundtruth.extend(phonemes.detach().cpu().tolist())
 21
 22
           # Release memory
 23
            del frames, phonemes, logits, predict
 24
            torch.cuda.empty cache()
 25
 26
      # Print classification report
 27
      print(classification report(
 28
            groundtruth,
 29
            predicted,
 30
           target names = PHONEMES # Skipping SOS and EOS tokens
 31
      ) )
 32
\overline{\Rightarrow}
                  precision
                              recall f1-score
                                                  support
           [SIL]
                       0.94
                                 0.95
                                           0.94
                                                   319908
                                 0.59
              AA
                       0.59
                                           0.59
                                                    29688
              ΑE
                       0.65
                                 0.66
                                           0.66
                                                    49298
              AH
                       0.61
                                 0.64
                                           0.63
                                                   123734
              Α0
                       0.66
                                 0.64
                                           0.65
                                                    29340
              AW
                       0.72
                                 0.63
                                           0.67
                                                    20274
              AY
                       0.79
                                 0.83
                                           0.81
                                                    49332
                       0.69
                                 0.68
                                           0.68
              В
                                                    23607
              CH
                       0.65
                                 0.62
                                           0.63
                                                    12644
               D
                       0.65
                                 0.58
                                           0.61
                                                    62763
              DH
                       0.70
                                 0.68
                                           0.69
                                                    37100
              EΗ
                       0.62
                                 0.59
                                           0.60
                                                    47112
              ER
                       0.69
                                 0.68
                                           0.69
                                                    54928
              EY
                       0.75
                                 0.76
                                           0.76
                                                    36184
              F
                       0.72
                                 0.78
                                           0.75
                                                    37562
               G
                       0.74
                                 0.69
                                           0.71
                                                    13541
              ΗН
                       0.75
                                 0.70
                                           0.72
                                                    34813
              ΙH
                       0.63
                                 0.59
                                           0.61
                                                    74887
              ΙY
                       0.77
                                 0.79
                                           0.78
                                                    70861
              JH
                       0.69
                                 0.61
                                           0.65
                                                     8730
                                           0.79
               K
                       0.77
                                 0.80
                                                    47016
               L
                       0.73
                                 0.82
                                           0.77
                                                    65902
               М
                       0.77
                                 0.78
                                           0.78
                                                    44728
                       0.74
                                           0.75
               Ν
                                 0.77
                                                    94541
              NG
                       0.73
                                 0.69
                                           0.71
                                                    19327
              0W
                       0.71
                                 0.62
                                           0.66
                                                    30755
              0Y
                       0.77
                                 0.54
                                           0.63
                                                     3861
               Р
                       0.75
                                 0.70
                                           0.72
                                                    34131
               R
                       0.69
                                 0.75
                                           0.72
                                                    62686
               S
                       0.80
                                 0.82
                                           0.81
                                                   101184
              SH
                       0.81
                                 0.81
                                           0.81
                                                    17628
              Τ
                       0.69
                                 0.68
                                           0.68
                                                    97390
              TH
                       0.44
                                 0.40
                                           0.42
                                                     9247
              UH
                       0.64
                                 0.45
                                           0.53
                                                     6286
              UW
                       0.73
                                 0.67
                                           0.70
                                                    26691
               V
                       0.71
                                 0.63
                                           0.67
                                                    27440
                       0.78
                                 0.81
                                           0.80
                                                    37697
```

```
Υ
                  0.70
                         0.60
                                 0.65
                                          9669
            Ζ
                  0.71
                          0.70
                                  0.70
                                          54850
           ZH
                  0.71
                          0.55
                                   0.62
                                            869
                                   0.75 1928204
      accuracy
                 0.71 0.68
                                  0.69 1928204
     macro avq
   weighted avg
                  0.74
                         0.75
                                  0.74
                                        1928204
  1 import seaborn as sns
  2 import matplotlib.pyplot as plt
  3 sns.set style("darkgrid")
  1 # Compute confusion matrix
  2 cm = confusion matrix(groundtruth, predicted)
  3
  4 # Normalize confusion matrix by row (i.e., by true labels)
  5 cm normalized = cm.astype('float') / cm.sum(axis = 1, keepd
  7 # Replace NaNs (from division by zero, if any row sum is 0)
  8 cm normalized = np.nan to num(cm normalized)
  9
 10 # Plot
 11 plt.figure(figsize=(15, 16))
 12 sns.heatmap(cm normalized,
 13
                annot=True,
                fmt=".1f",
 14
 15
                cmap="Greens",
 16
                xticklabels=PHONEMES,
 17
                yticklabels=PHONEMES,
 18
                cbar kws={'label': 'Proportion'})
 19
 20 plt.title('Normalized Confusion Matrix', fontsize=16)
 21 plt.xlabel('Predicted Label', fontsize=12)
 22 plt.ylabel('True Label', fontsize=12)
 23 plt.xticks(rotation=45, ha='right', fontsize=10)
 24 plt.yticks(rotation=0, fontsize=10)
 25 plt.tight layout()
 26 plt.show()
 27
\rightarrow
```

EXPERIMENT II - Colab

14/06/2025, 13:02

Normalized Confusion Matrix



https://colab.research.google.com/#fileId=https%3A//storage.googleapis.com/kaggle-colab-exported-notebooks/cs16mtech1...

Predicted Label