

## ✓ Frame-Level Speech Recognition

### GOAL:

Given data which consists of melspectrograms, and phoneme labels for each 28-dimensional vector in the melspectrogram. The task is to predict the label of a particular 28-dimensional vector in an utterance (plus optional context) using a Feed Forward Deep Neural Network (FF-DNN or MLP).

### OBJECTIVES:

- The main goal here is to explore neural networks for speech recognition, mainly focusing on phoneme state labelling.
- Speech recognition is an important field in deep learning with multiple applications. Speech is a fundamental form of human communication.
- Speech recognition involves converting spoken language into written text or data that could be understood by machines.
- Speech data refers to audio recordings of human speech, while phonemes represent the smallest units of sound that can convey a different meaning(bake, take: b,t).
- Spectrograms are visual representations of the acoustic properties of speech signals, i.e. captures the changes in the frequency over time.

## ✓ DATASET DETAILS

- Dataset contains audio recordings (utterances) and their phoneme state (subphoneme) labels.
- The data comes from articles published in the Wall Street Journal (WSJ) that are read aloud and labelled using the original text.
- The dataset provided has speech data in the form of Mel spectrograms.
- The data comprises of:
  - Speech recordings (raw mel spectrogram frames)
  - Frame-level phoneme state labels
- Training data have around 28539 samples

### Phonemes and Phoneme States

- As letters are the atomic elements of written language, phonemes are the atomic elements of speech.
- It is crucial for us to have a means to distinguish different sounds in speech that may or may not represent the same letter or combinations of letters in the written alphabet.
- In the dataset, we will consider a total of 40 phonemes in this language.
- A powerful technique in speech recognition is to model speech as a markov process with unobserved states.
- This model considers observed speech to be dependent on unobserved state transitions. We refer to these unobserved states as phoneme states or subphonemes. For each phoneme, there are 3 respective phoneme states. The transition graph of the phoneme states for a given phoneme is as follows:
- Example: ["+BREATH+", "+COUGH+", "+NOISE+", "+SMACK+", "+UH+", "+UM+", "AA", "AE", "AH", "AO", "AW", "AY", "B", "CH", "D", "DH", "EH", "ER", "EY", "F", "G", "HH", "IH", "IY", "JH", "K", "L", "M", "N", "NG", "OW", "OY", "P", "R", "S", "SH", "SIL", "T", "TH", "UH", "UW", "V", "W", "Y", "Z", "ZH"]
- Hidden Markov Models (HMMs) estimate the parameters of this unobserved markov process (transition and emission probabilities) that maximize the likelihood of the observed speech data.
- We will take model-free approach and classify mel spectrogram frames using a neural network that takes a frame (plus optional context) and outputs class probabilities for all 40 phoneme states. The training data has the corresponding phonemes for this data and we need to train an MLP for predicting. We will be building a multilayer perceptron(MLP) that can effectively recognize and label the phoneme states in the training data. An MLP is a type of neural network that comprises multiple layers of perceptrons, that help it capture the features and patterns of the data.

## Speech Representation

- Raw speech signal (also known as the speech waveform) is stored simply as a sequence of numbers that represent the amplitude of the sound wave at each time step. This signal is typically composed of sound waves of several different frequencies overlaid on top of one another. For human speech, these frequencies represent the frequencies at which the vocal tract vibrates when we speak and produce sound. Since this signal is not very useful for speech recognition if used directly as a waveform, we convert it into a more useful representation called a "melspectrogram" in the feature extraction stage.
- The variation with time of the frequencies present in a particular speech sample are very useful in determining the phoneme being spoken. In order to separate out all the individual

frequencies present in the signal, we perform a variant of the Fourier Transform, called the Short-Time Fourier Transform (STFT) on small, overlapping segments (called frames, each of 25ms) of the waveform. A single vector is produced as the result of this transform. Since we use a stride of 10ms between each frame, we end up with 100 vectors per second of speech. Finally, we convert each vector into a 28-dimensional vector (for further readings <https://haythamfayek.com/2016/04/21/speech-processing-for-machine-learning.html>). For an utterance T seconds long, this leaves us with a matrix of shape (100\*T, 28) known as the melspectrogram. Note that in the dataset provided to you, we have already done all of this pre-processing and provided the final (\*, 28) shaped melspectrograms to you. The data provided consists of these melspectrograms, and phoneme labels for each 28-dimensional vector in the melspectrogram. The task is to predict the label of a particular 28-dimensional vector in an utterance.

## PERFORMANCE METRIC

- Accuracy
- Confusion Matrix
- Precision and Recall

## ✓ EXPERIMENT

- Build FF-DNN only without any batchnorm or dropout
- Different right and left context are added so that capture more variability in the speaker utterances smoothing
- Sparse MLP is created using the loss function
  - Reference: <https://ieeexplore.ieee.org/document/5734801>
- 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 ! ls /content/
```

```
1 !pip install torch-summary --quiet
```

```
1 import torch
2 import numpy as np
3 import sklearn
4 import gc
```

```
5 import zipfile
6 import pandas as pd
7 from tqdm.auto import tqdm
8 import os
9 import datetime
10 import torchsummary
11 device = 'cuda' if torch.cuda.is_available() else 'cpu'
12 print("Device: ", device)
```

Device: cuda

```
1# # using colab for including google drive to save model ch
2# from google.colab import drive
3# 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 2.7 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')
3
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-hwlp2-f23.zip to /kaggle/working
100%|██████████| 3.97G/3.99G [00:24<00:00, 195MB/s]
100%|██████████| 3.99G/3.99G [00:24<00:00, 173MB/s]

```

```
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
```

```
1 ! 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 = 2
9 GAMMA = 0.1
10 BASE_DIRECTORY = '/content/data/11-785-f23'
11 TRAINING_DATA = BASE_DIRECTORY + 'train-c
12 EVALUATION_DATA = BASE_DIRECTORY + 'dev-cle
13 PHONEMES = ['[SIL]', 'AA', 'AE'
14             'B', 'CH', 'D',
15             'F', 'G', 'HH',
16             'L', 'M', 'N',
17             'R', 'S', 'SH',
18             'V', 'W', 'Y',
19 PHONEMES_TO_INDEX = {phoneme: idx for idx, p
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):
2     def __init__(self, root, phonemes = PHONEMES_TO_INDEX,
3                 left_context = LEFT_CONTEXT,
4                 right_context = RIGHT_CONTEXT):
5
6         self.left_context = left_context
7         self.right_context = right_context
8         self.phonemes_mapping = phonemes
9
```

```
10 self.mfcc_dir = os.path.join(root, 'mfcc')
11 self.transcript_dir = os.path.join(root, 'transcri
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
18 assert len(mfcc_names) == len(transcript_names), "
19
20 total_frames = 0
21
22 for i in range(len(mfcc_names)):
23     mfcc_path = os.path.join(self.mfcc_dir, mfcc_n
24
25     mfcc = np.load(mfcc_path,
26                   allow_pickle = False,
27                   mmap_mode='r')
28     total_frames += mfcc.shape[0]
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
34 self.mfccs = np.zeros((total_frames + right_context
35 self.transcripts = [None] * (total_frames + right_
36 ### Release memory
37 del sample_mfcc, total_frames
38 gc.collect()
39
40 index = 0
41 for i in range(len(mfcc_names)):
42     mfcc = np.load(os.path.join(self.mfcc_dir, mfc
43                   allow_pickle = False,
44                   mmap_mode = 'r')
45
46     mfcc = (mfcc - mfcc.mean(axis = 1, keepdims =
47
48     transcript = np.load(os.path.join(self.transcr
49                   allow_pickle = False,
50                   mmap_mode='r'))[1:-1]
51
52     self.mfccs[index: index + mfcc.shape[0]] = mfc
53     self.transcripts[index: index + len(transcript
54     index += mfcc.shape[0]
```

```
55         del mfcc, transcript
56
57         if i % 1000 == 0:
58             gc.collect() ### Release memory
59
60         # Save original dataset length (before adding right
61         self.length = len(self.mfccs) - right_context
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
70     def __len__(self):
71         return self.length
72
73     def __getitem__(self, ind):
74         if ind < self.left_context: # index is less than
75             # zeros need to prepend with it
76             padding = []
77             for _ in range(self.left_context - ind):
78                 padding.append(np.zeros((1, self.mfcc_dim)
79
80             for i in range(ind):
81                 padding.append(self.mfccs[i][None, :])
82
83             # Include current and right context frames
84             current = self.mfccs[ind][None, :] # Ensure s
85             right_context = [self.mfccs[i][None, :] for i
86             frames = np.concatenate(padding + [current] +
87
88         else: # when index is greater than or equal to le
89             left_context_frames = self.mfccs[ind - self.le
90             current = self.mfccs[ind][None, :] # Ensure s
91             right_context_frames = self.mfccs[ind + 1: ind
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
97             phonemes = torch.tensor(self.transcripts[ind], dtype
98
```

```

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
7
8
9 train_loader = torch.utils.data.DataLoader(
10     dataset      = train_data,
11     num_workers  = 2,
12     batch_size   = BATCH_SIZE,
13     pin_memory   = True,
14     shuffle      = True
15 )
16
17 val_loader = torch.utils.data.DataLoader(
18     dataset      = val_data,
19     num_workers  = 1,
20     batch_size   = BATCH_SIZE,
21     pin_memory   = True,
22     shuffle      = False
23 )
24
25 print("Train dataset samples = {}, batches = {}".format(train_loader.get_num_samples(), train_loader.get_num_batches()))
26 print("Validation dataset samples = {}, batches = {}".format(val_loader.get_num_samples(), val_loader.get_num_batches()))

```

⇒ Train dataset samples = 36091157, batches = 8812  
 Validation dataset samples = 1928204, batches = 471

```

1 # # Testing code to check if data loaders are working as ex
2 # total_batches = len(train_loader)
3
4 # for i, (frames, phoneme) in enumerate(train_loader):
5 #     if i < 10 or i >= total_batches - 10:
6 #         print(f"Batch {i + 1}/{total_batches}")
7 #         print("Frames:", frames)
8 #         print("Phoneme:", phoneme)
9 #         print("-" * 50)
10

```

```

1 class Network(torch.nn.Module):
2     def __init__(self, input_size):
3         super(Network, self).__init__()
4

```



```

5         # Neurons in each layer: input -> hidden(s) -> outp
6         self.neurons = [input_size] + NUMBER_OF_NEURONS + [
7
8         layers = []
9         for in_features, out_features in zip(self.neurons[:
10             layers.append(torch.nn.Linear(in_features, out_
11             layers.append(torch.nn.ReLU())
12
13         # Final layer (no activation)
14         layers.append(torch.nn.Linear(self.neurons[-2], sel
15
16         # Combine all into a sequential model
17         self.model = torch.nn.Sequential(*layers)
18
19         # Apply weight initialization
20         self._initialize_weights()
21
22     def _initialize_weights(self):
23         print('Initialization of weights using Kaiming')
24         for m in self.model:
25             if isinstance(m, torch.nn.Linear):
26                 # Kaiming initialization for weights
27                 torch.nn.init.kaiming_uniform_(m.weight, no
28
29
30     def forward(self, x):
31         out = self.model(x)
32         return out
33

```

```

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
4 torchsummary.summary(model, (INPUT_SIZE,))

```

Initialization of weights using Kaiming

Layer (type)	Output Shape	Param #
Linear-1	[-1, 2048]	862,208
ReLU-2	[-1, 2048]	0
Linear-3	[-1, 2048]	4,196,352
ReLU-4	[-1, 2048]	0
Linear-5	[-1, 1024]	2,098,176
ReLU-6	[-1, 1024]	0
Linear-7	[-1, 1024]	1,049,600
ReLU-8	[-1, 1024]	0
Linear-9	[-1, 512]	524,800

ReLU-10	[-1, 512]	0
Linear-11	[-1, 256]	131,328
ReLU-12	[-1, 256]	0
Linear-13	[-1, 256]	65,792
ReLU-14	[-1, 256]	0
Linear-15	[-1, 40]	10,280

```
=====
Total params: 8,938,536
Trainable params: 8,938,536
Non-trainable params: 0
```

```
-----
Input size (MB): 0.00
Forward/backward pass size (MB): 0.11
Params size (MB): 34.10
Estimated Total Size (MB): 34.21
-----
```

```
1 class SparseLoss(torch.nn.Module):
2     def __init__(self, model, lambda_l1=1e-4):
3         super(SparseLoss, self).__init__()
4         self.ce_criterion = torch.nn.CrossEntropyLoss(label
5         self.lambda_l1 = lambda_l1
6         self.model = model
7
8         # Collect indices of Linear layers, skipping first
9         self.linear_layer_indices = [
10             i for i, layer in enumerate(self.model.model[1:
11             if isinstance(layer, torch.nn.Linear)
12         ]
13
14     def forward(self, logits, targets):
15         ce_loss = self.ce_criterion(logits, targets)
16
17         sparse_loss = 0.0
18         for i in self.linear_layer_indices:
19             weight = self.model.model[i].weight # [out_fea
20
21             # Neuron-wise aggregation: sum inputs per neuro
22             neuron_weights = weight.sum(dim=1) # [out_feat
23
24             # sum(log(1 + (neuron)^2))
25             sparse_loss += torch.log1p(neuron_weights.pow(2
26
27         total_loss = ce_loss + self.lambda_l1 * sparse_loss
28         return total_loss
29
```

```
1 # criterion = torch.nn.CrossEntropyLoss(label_smoothing = L
```

```
2 # We use CE because the task is multi-class classification
3 criterion = SparseLoss(model, lambda_l1 = 1e-4)
4 optimizer = torch.optim.Adam(model.parameters(),
5                               lr = INITIAL_LEARNING_RATE,
6                               weight_decay = L2_PENALTY) #De
7 scheduler = torch.optim.lr_scheduler.StepLR(
8     optimizer, step_size = STEP_SIZE, gamma = GAMMA
9 )
10 # Refer - https://pytorch.org/docs/stable/notes/amp_example

1 torch.cuda.empty_cache()
2 gc.collect()
```

→ 101

```
1 def train(model, dataloader, optimizer, criterion):
2
3     model.train()
4     tloss, tacc = 0, 0 # Monitoring loss and accuracy
5     batch_bar = tqdm(total=len(train_loader), dynamic_nco
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)
13        frames = frames.to(device)
14        phonemes = phonemes.to(device)
15
16        ### Forward Propagation
17        logits = model(frames)
18
19        ### Loss Calculation
20        loss = criterion(logits, phonemes)
21
22        ### Backward Propagation
23        loss.backward()
24
25        ### Clip gradients
26        torch.nn.utils.clip_grad_norm_(model.parameters(),
27
28        ### Gradient Descent
29        optimizer.step()
30
```

```

31         tloss    += loss.item()
32         tacc     += torch.sum(torch.argmax(logits, dim= 1) =
33
34         batch_bar.set_postfix(loss="{:.04f}".format(float(t
35                               acc="{:.04f}%".format(float(t
36         batch_bar.update()
37
38         ### Release memory
39         del frames, phonemes, logits
40         torch.cuda.empty_cache()
41
42     batch_bar.close()
43     tloss    /= len(train_loader)
44     tacc     /= len(train_loader)
45
46     return tloss, tacc

```

```

1 def eval(model, dataloader):
2
3     model.eval() # set model in evaluation mode
4     vloss, vacc = 0, 0 # Monitoring loss and accuracy
5     batch_bar    = tqdm(total=len(val_loader), dynamic_ncols
6
7     for i, (frames, phonemes) in enumerate(dataloader):
8
9         ### Move data to device (ideally GPU)
10        frames      = frames.to(device)
11        phonemes    = phonemes.to(device)
12
13        # makes sure that there are no gradients computed a
14        with torch.inference_mode():
15            ### Forward Propagation
16            logits   = model(frames)
17            ### Loss Calculation
18            loss     = criterion(logits, phonemes)
19
20        vloss    += loss.item()
21        vacc     += torch.sum(torch.argmax(logits, dim= 1) =
22
23
24        batch_bar.set_postfix(loss="{:.04f}".format(float(v
25                               acc="{:.04f}%".format(float(v
26        batch_bar.update()
27
28        ### Release memory

```

```
29         del frames, phonemes, logits
30         torch.cuda.empty_cache()
31
32     batch_bar.close()
33     vloss /= len(val_loader)
34     vacc /= len(val_loader)
35
36     return vloss, vacc

```

  

```
1 best_model_path = os.path.join(MODEL_DIR, "best_model.pt")
2 best_acc = -np.inf
3 for epoch in range(EPOCHS):
4     ### clean up memory before computation
5     torch.cuda.empty_cache()
6     gc.collect()
7
8     print(f"\nEpoch {epoch + 1}/{EPOCHS}")
9
10    curr_lr = float(optimizer.param_groups[0]['lr'])
11    train_loss, train_acc = train(model, train_loader, opti
12    val_loss, val_acc = eval(model, val_loader)
13
14    print(f"\tTrain Acc: {train_acc*100:.2f}%\tTrain Loss:
15    print(f"\tVal Acc: {val_acc*100:.2f}%\tVal Loss: {v
16
17    # Save model at every epoch
18    epoch_model_path = os.path.join(MODEL_DIR, f"model_at_e
19    torch.save(model.state_dict(), epoch_model_path)
20
21    # Save best model
22    if val_acc > best_acc:
23        best_acc = val_acc
24        torch.save(model.state_dict(), best_model_path)
25        print(f"Updated Best Model at: {best_model_path}")
26
27    ### take step in adjusting the learning rate
28    scheduler.step()
29
30
```



```
Epoch 1/4
Train: 0%|          | 0/8812 [00:00<?, ?it/s]
Val: 0%|          | 0/471 [00:00<?, ?it/s]
      Train Acc: 70.53%      Train Loss: 1.0100      LR: 0.0010000
      Val Acc: 70.38%      Val Loss: 1.0031
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: 75.57%      Train Loss: 0.8290      LR: 0.0010000
      Val Acc: 71.77%      Val Loss: 0.9586
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: 79.85%      Train Loss: 0.6779      LR: 0.0001000
      Val Acc: 74.03%      Val Loss: 0.8763
Updated Best Model at: /content/best_model.pt
```

```
1 gc.collect()
```

```
⇒ Val: 0%|          | 0/471 [00:00<?, ?it/s]
      Train Acc: 81.11%      Train Loss: 0.6383      LR: 0.0001000
      Val Acc: 73.02%      Val Loss: 0.8873
```

```
1 # Load best model after training
```

```
2 model.load_state_dict(torch.load('/content/best_model.pt'))
```

```
⇒ <All keys matched successfully>
```

```
1 from sklearn.metrics import classification_report, confusion_matrix
2 model.eval() # Set model in evaluation mode
3 predicted = []
4 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)
17
18    # Detach and move to CPU for evaluation
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

```



	precision	recall	f1-score	support
[SIL]	0.93	0.95	0.94	319908
AA	0.59	0.58	0.58	29688
AE	0.64	0.66	0.65	49298
AH	0.62	0.63	0.62	123734
AO	0.65	0.64	0.65	29340
AW	0.69	0.64	0.66	20274
AY	0.77	0.83	0.80	49332
B	0.68	0.67	0.68	23607
CH	0.65	0.60	0.62	12644
D	0.66	0.56	0.60	62763
DH	0.69	0.68	0.68	37100
EH	0.60	0.59	0.60	47112
ER	0.67	0.71	0.69	54928
EY	0.73	0.76	0.74	36184
F	0.71	0.79	0.75	37562
G	0.73	0.69	0.71	13541
HH	0.73	0.70	0.71	34813
IH	0.62	0.59	0.60	74887
IY	0.77	0.78	0.78	70861
JH	0.67	0.63	0.65	8730
K	0.76	0.80	0.78	47016
L	0.75	0.78	0.76	65902
M	0.76	0.77	0.77	44728
N	0.74	0.76	0.75	94541
NG	0.70	0.69	0.69	19327
OW	0.68	0.63	0.65	30755
OY	0.69	0.57	0.63	3861
P	0.71	0.71	0.71	34131
R	0.69	0.73	0.71	62686
S	0.80	0.82	0.81	101184
SH	0.79	0.82	0.80	17628
T	0.68	0.67	0.68	97390
TH	0.45	0.39	0.42	9247
UH	0.64	0.47	0.54	6286
UW	0.73	0.66	0.69	26691
V	0.70	0.62	0.66	27440
W	0.79	0.80	0.80	37697
Y	0.68	0.63	0.66	9669
Z	0.72	0.68	0.70	54850
ZH	0.74	0.48	0.58	869

accuracy			0.74	1928204
macro avg	0.70	0.68	0.69	1928204
weighted avg	0.74	0.74	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
6
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,
14             fmt=".1f",
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
```







1