# Estimating speech from lip movement

Jithin D. George, Ronan Keane, Conor Zellmer

March 16, 2017

**Abstract**

The goal of this project is to develop a limited lip reading algorithm for a subset of the english language. We consider a scenario in which no audio information is available. We first develop a method for processing raw video and extracting the position of the lips in each frame. We then prepare the lip data for processing and classify the lips into visemes and phonemes. We then use Hidden Markov Models to predict the words the speaker is saying based on the classification. We use the GRID audiovisual sentence corpus database for our study.

## 1 Introduction and Overview

This report details the use of classification algorithms and Hidden Markov Models to perform visual speech recognition. Videos from the GRID Audiovisual Sentence Corpus database[10] were used in the project. The GRID database contains a set of 1000 videos of a single speaker for visual speech recognition. Each video is approximately 3 seconds in length, and contains a man speaking a collection of words chosen to encompass each sound possible in spoken English. The database includes transcripts of each video, with the locations of each word in the video by frame numbers. This report details the extraction of lip contours from each video, the identification of phonemes being spoken using classification methods, and prediction of speech using Hidden Markov Models (HMM). The lip contours for each video frame were extracted using active contour masks, dynamic mode decomposition, edge detection and color classification, and both Naive Bayes and k-Nearest Neighbors classification methods were used. Classification and prediction was performed based upon both phonemes and visemes. Classification based upon phonemes was accomplished with a maximum accuracy of 11.6103%, while viseme classification had a maximum accuracy of 19.7355%, both using the k-Nearest Neighbors algorithm. For 2 word sequences, prediction based upon phonemes using the HMM was found to have a maximum accuracy of 87.5%, and prediction based upon visemes was found to be at most 92% accurate.

## 2 Theoretical Background

### 2.1 Feature Extraction

Lip reading is a complicated task and there are no "go-to" algorithms for detecting and tracking the position of an individual's lips. We can use Matlab's built in active contour and edge detection to hope to do background/foreground seperation on the video. We can also use DMD to do background/foreground seperation. The assumption is the speaker's face is stationary

enough that it is possible to detect the lips as the foreground. In practice, this is not the case. We can also classify the pixel colors of the video frames and segment the lips based on the idea that the speaker's lip color is different from their skin color. For color classification, we can project the pixel color into the LAB color space to acheive better classification results. In either of these different strategies we also need to isolate the general mouth region of the speaker so that we will not also detect eyes, nose, etc.

## 2.2 Phonemes and Visemes

Phonemes are the smallest identifiable sound unit in the sound system of a language. [6] According to Zhong, et al, phonemes are "basic sound segments that carry linguistic distinction." [7] In theory, visemes are the analogous basic units in the visual domain. However, there is no agreement on a common definition for visemes in practice[8]. In audio speech recognition, phonemes are detected and used to reconstruct speech. In visual speech recognition, only visemes can be detected. Phonemes for the basis for a spoken language, and hence automated lipreading typically employs a mapping from phonemes to visemes. Many such mapping can be found in the literature, but all suffer from the issue of there being more phonemes than visemes, resulting in a many-to-one map.[8] For instance, this project uses 37 phonemes and 11 visemes. See Table (??) for the phoneme to viseme map used in this project.

Table 1: Phoneme to Viseme Map from Lee and York, 2000, via [8].

| Viseme Number | Viseme Label | Associated Phonemes |
|---|---|---|
| 1 | P | b p m |
| 2 | T | d t s z th dh |
| 3 | K | g k n l y hh |
| 4 | CH | jh ch |
| 5 | F | f v |
| 6 | W | r w |
| 7 | IY | iy ih |
| 8 | EH | eh ey ae |
| 9 | AA | aa aw ay ah |
| 10 | A0 | ao oy ow |
| 11 | UH | uh uw |

## 2.3 Hidden Markov Models

A Markov model involves the transition of a particular state to other states based on transition probabilities. A future state is only depends on the current state and not the states before it. Now, consider that at every state, there would be real world observations. These observations are controlled by the emission probabilities at each state.

For example, if we were to represent the ever changing weather, the states would be sunny, rainy or snowing and the observations would be summer clothes, rain boots or snow shoes. We can see that the emission probabilities for each observation is different depending on the state. To be more clear, the emission probabilities depend on the states.

We look at Hidden Markov Models(HMM). We decide that this is a relevant model because the words spoken are those defined by language and thus occur in specific pattern and not randomly. For example, given the first letter of word 'k', the probability that the next letter is a vowel is much higher than it being a consonant. A machine learning algorithm without this would be as inefficient as the initial Enigma machine in the movie "The Imitation Game". HMMs are very popular in the fields of speech [**?**] and gesture recognition [**?**] [**?**].

Although HMMs have fascinating problems related to evaluation and learning, our interests are in decoding. We have a sequence of observations and our aim is to estimate the states that created that. The Viterbi algorithm [**?**] gives us the states that maximize the occurrence of the observations.

So, given the features from the videos, we find the states. The states are the units of words, here chosen to be phonemes.

# 3 Implementation and Development

## 3.1 Feature extraction

### 3.1.1 Obtaining the initial mask

The first step is to determine the general location of the mouth in each frame. Matlab has a built in facial recognition package "Cascade Vision Detector" which uses the Viola-Jones algorithm to detect the facial features of people. We use the Cascade Mouth detector to determine the location of the face in each frame of each movie. This algorithm is not perfect however and we often detect the eyes or chin of the subject in question. We filter out these false detections by looking at the position of the mouth. The speaker is in generally the same location in each movie so we know if a detected mouth position is too low or too high we can throw that region out. We refer to this detected mouth region as the *initial mask*. We need a good initial mask in order to properly filter our results. Without knowing the general mouth region our algorithms will act on the whole face, but we consider only the lips for our project.Therefore, really our objective is to separate the lips from the mouth region.

### 3.1.2 Active contour

Next we convert the image to grayscale. On the grayscale image we can apply Matlab's built in active contour and edge detection, as well as the DMD algorithm we develop in homework 4. We can experiment with the different segmentation types for Matlab's active contour and edge detection. For active contour, Chan-Vese gives poor results with irregular edges. Edge method for active contour gives smoother regions but irregular shapes. Active contour takes many iterations to properly converge and gives poor results. From the figures of active contour, we see the inconsistent lip region detected. Because we are trying to determine the changes in lip shape over time, active contour is not suitable because the results it produces vary too much. Additionally, because it is extremely slow, active contour is not ideal for our purposes. We need to process a large number of videos so a fast algorithm is preferable. If active contour was more computationally efficient it might be possible to use a large number of iterations for each frame in order to obtain a consistent result.
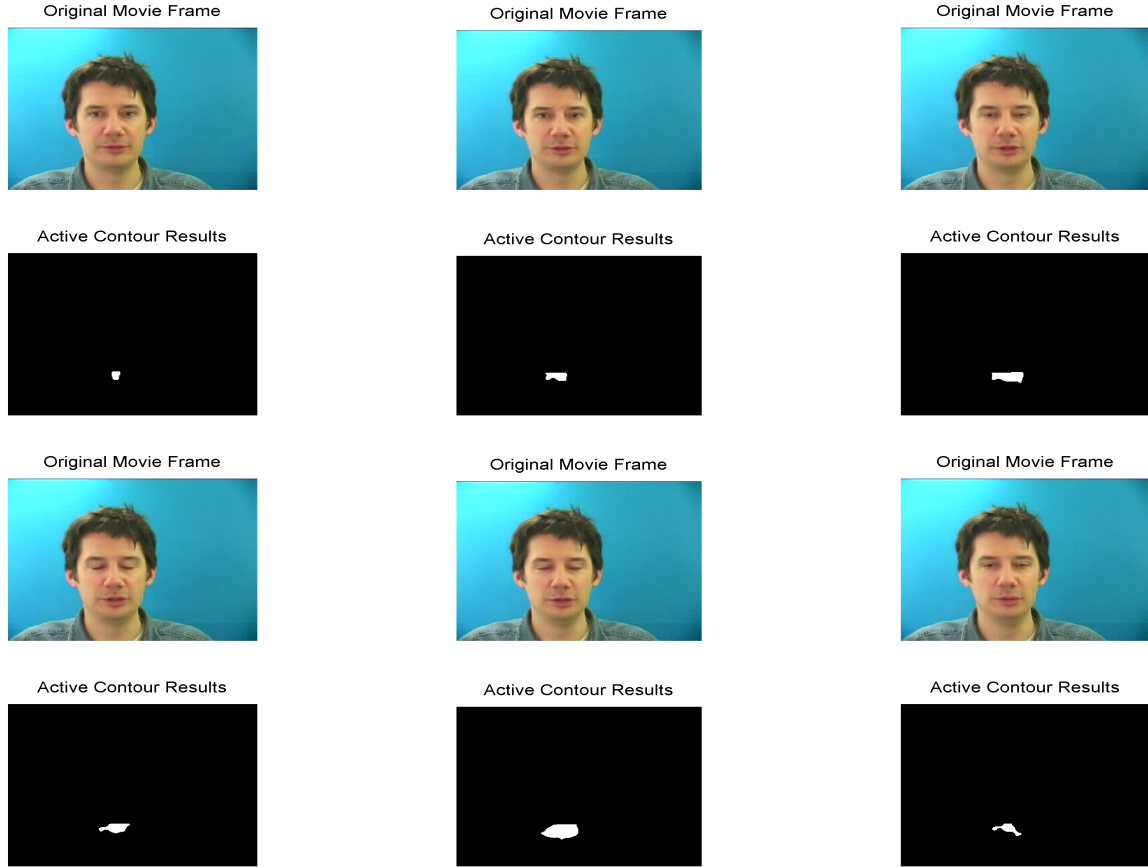
Figure 1: Active Contour Results for selected frames. Active contour provides an inconsistent lip profile because it never properly converges.

### 3.1.3 Canny edge detection

Matlab's built in edge detection has many different methods which produce different results. We consider using Canny, Prewitt and Sobel edge detection.Canny typically finds the "strong" edges whereas Sobel and Prewitt tend to pick out more detail in the image. We are really only interested in the lips, and ideally want only 1 contour, so we decide to use Canny edge detection. We want to detect the outside and inside edge of the lips, but all methods typically detect the inside edge of the lip. When using Matlab's "edge", it is important to consider the threshold. Too low a threshold value will produce other features of the face and possibly noise. Too high a threshold and we will not pick out the entire lip. Therefore we decide to start at a high threshold value, and if we do not find a large enough edge, we lower the threshold and do the detection again. Typically we only need the lower the threshold once or twice. We set a minimum number of detected pixels in the edge to decide whether or not to automatically lower the threshold. Although this may seem expensive, in practice we only need to lower the threshold for certain frames, and so most of the time we only need to do the algorithm once. Additionally, edge is much faster than active contour, so this method is still much faster than active contour, and gives more consistent results in well. Overall the Canny edge detection finds a good edge for the lips but it oftentimes only picks out part of the mouth or sometimes

picks out other features around the mouth as well, like the chin or space between the nose and mouth.
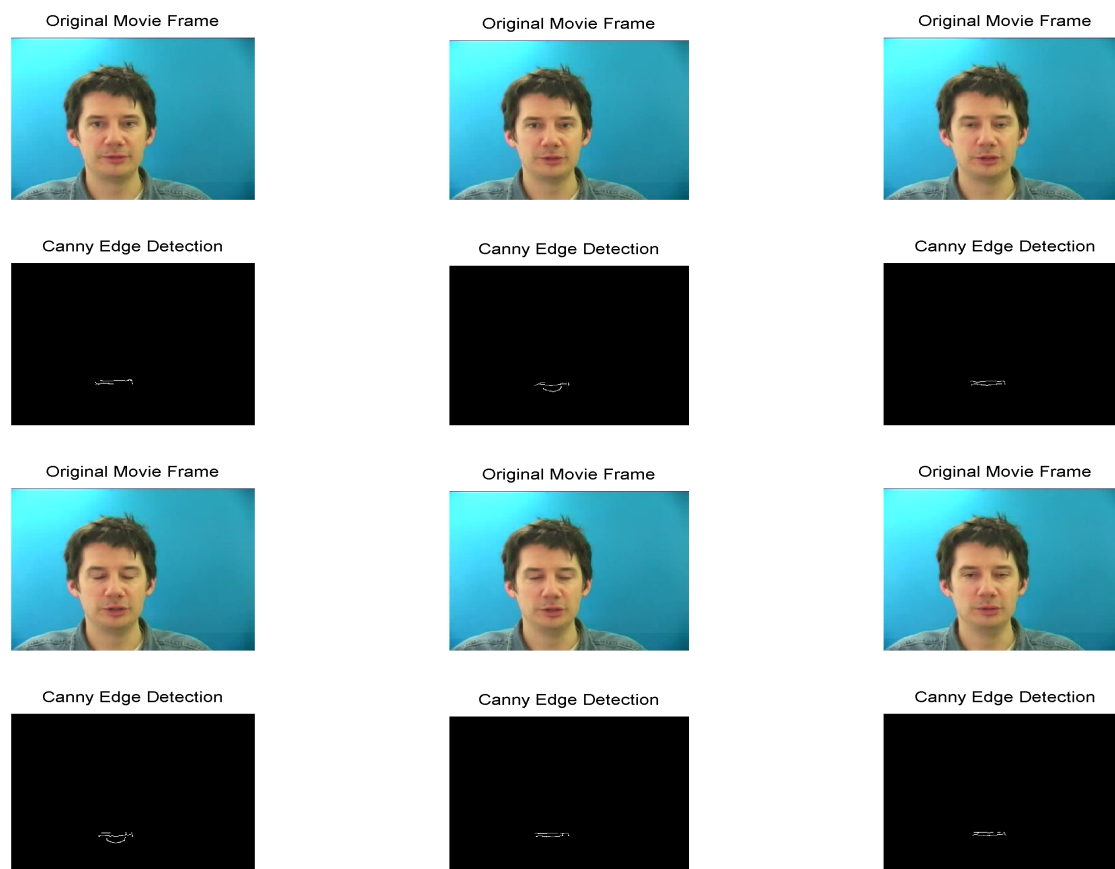


Figure 2: Canny Edge Detection results for selected frames. Although it captures the profile of the lips, it does not always capture the entire contour.It also does not show the inside of the lips.

### 3.1.4 Dynamic mode decomposition

Next we use the DMD algorithm we developed in homework 4 to try foreground/background separation. We find that the face moves too much, and there is too much variation in the speaker's skin around the mouth to properly use DMD. We detect many points that are not on the lip. While DMD is good at separating the man's face from the background, it cannot accurately separate the lips from mouth region. It might be possible to alter the brightness and contrast of the separated image using DMD to obtain a better result, but it would still be too inconsistent to obtain good results for the rest of the project.

### 3.1.5 Color classification

Thus far, all of the methods discussed have acted on grayscale frames of the color movie. We now consider a technique for acting on the color frames. If we assume that the lips of the
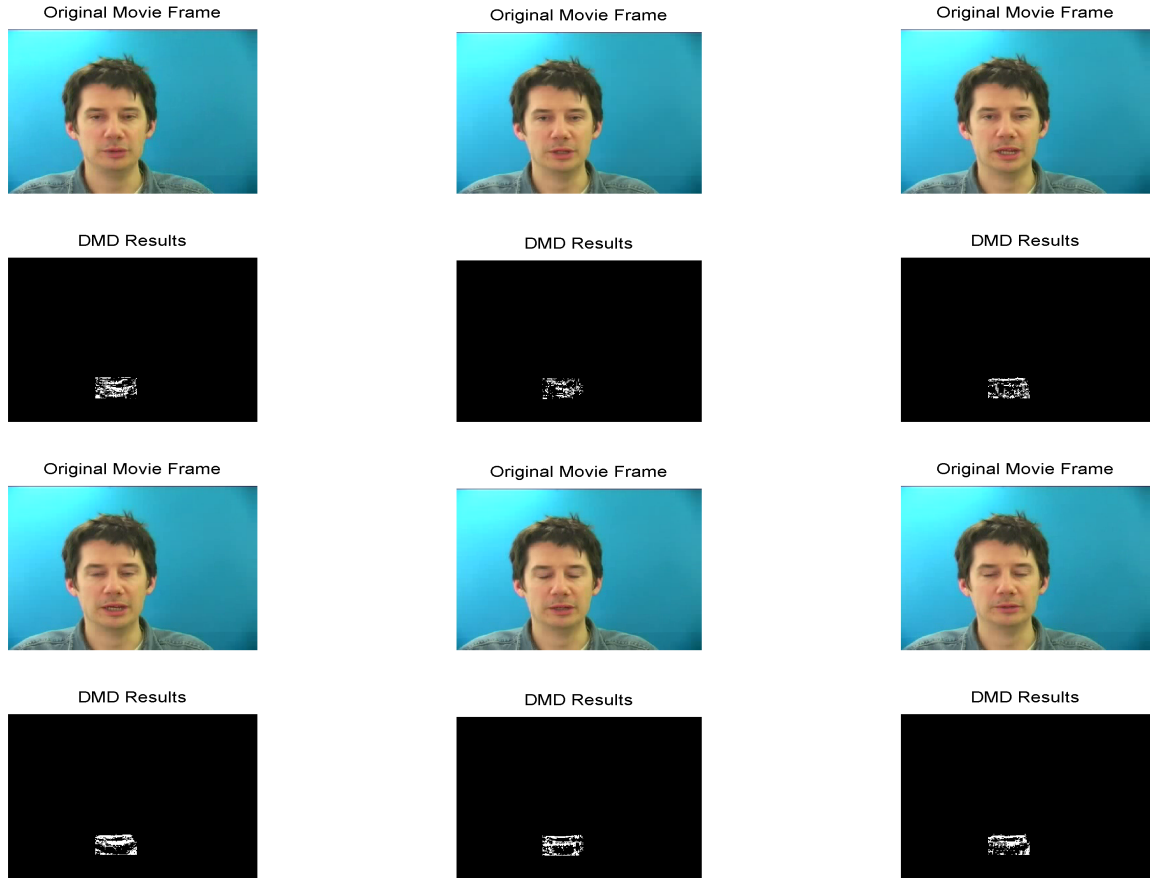
Figure 3: DMD results for selected frames. DMD is not ideal in this situation because the background (the man's face) is not completely stationary.

speaker have a distinct color from the rest of the speaker's mouth region (i.e. the skin, teeth, inside of the mouth) then we can determine the lip region by classifying the colors of the mouth region. It is possible to do this using the RGB values of each pixel. We instead represent each pixel with it's LAB color space values. The point of LAB colors is that they match how human vision perceives colors. The L stands for lighting and AB are values corresponding to color. By classifying each pixel into one of up to 4 clusters based on it's AB values with a k-means algorithm, we can separate the lips from the rest of the mouth region. The number of clusters created varies per frame because in some frames, the teeth and inner mouth can be grouped in their own clusters. Other frames, the speaker has their mouth closed so only 2 clusters are needed. We can differentiate between the clusters by looking at the calculated centroid positions of each cluster. It is determined experimentally that the lip region has the highest A color value, and it usually has the lowest B color value as well. Color classification is fairly fast and accurate. It gives the thickness as well as the shape of the lips. It's only downside is that it often classifies the left and right sides of the lips into a different cluster. Color classification typically picks out an ellipse in the middle of the lips as opposed to the whole lip region. Overall, color classification is determined to give the most accurate results, and edge detection is a close second. In hindsight, it would have made much more sense to use a classification

6

algorithm, like knn-search, instead of k-means. Although we choose the wrong algorithm to use our results are still acceptable.
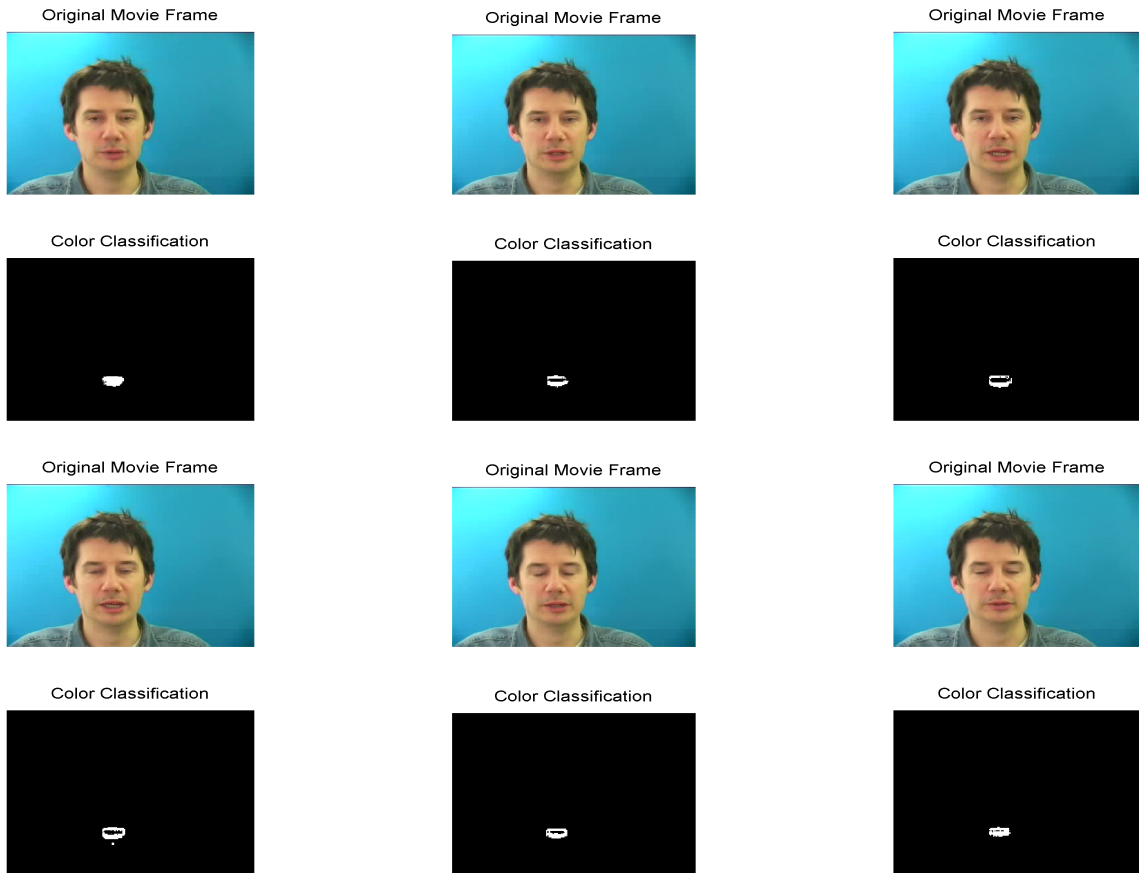


Figure 4: Color Classification for selected frames. Arguably the best overall method for capturing the movement of the lips.

In general the best results come from either more computationally expensive or mathematically sophisticated techniques. For example, training a neural net to classify the lip region of each frame would probably be the most effective technique overall, but would be very computationally costly and require manually determining the lip region for training. Many papers in the literature consider fitting a vector to fit the lip region and then track the lips by moving points of the vector and measuring the change of shape. They then can compute the most likely shape by comparing it to known positions of the lip. This method relies more heavily on analysis but also requires manual training as well. In general, lip tracking is a classification problem where the goal is to classify the pixels that belong to the lip of the speaker.

## 3.2 Extracting Phonemes

Using the nltk library in Python, we convert every word to its constitutive arpabet phonetics. It gives the following output for the words - f', 'see', 'sea', 'compute', 'comput', 'cat'. Only 'comput' fails because it isn't a real word

```
['EH1', 'F']
['S', 'IY1']
['S', 'IY1']
['K', 'AH0', 'M', 'P', 'Y', 'UW1', 'T']
'comput'
['K', 'AE1', 'T']
```

From the words spoken in our videos, we get a set of 36 unique phonemes. The code for this is shown in using **??**. From the transcripts, we extract all the data into a csv file using **??**

## 3.3   Extracting Subtitles and Assigning Phonemes

The Corpus Grid II database contains a *.txt* file for each video containing the words spoken and corresponding frames for each word. These were downloaded and extracted into Matlab using the *textscan* function in Matlab. The words were deconstructed into their phonemes, and phonemes were assigned to each frame to create the training set. The assignment of a phoneme label to each frame was done by assigning each phoneme from a word to an equivalent number of the video frames corresponding to the phonemes in each word.

## 3.4   Phoneme and Viseme Classification

The data matrix was created by reshaping each frame in each video in to a single column. When each column was reshaped, its saved phoneme was checked, and the corresponding viseme index was saved to create the labels for a classification algorithm. The columns for each video frame were then concatenated in order to form the data matrix. The singular value decomposition was computed of the matrix of video frames. Classification was performed on first 30 columns of the $V$ matrix using both a Naive Bayes and a k-nearest neighbors algorithm. Classification was done with a random 75% of the data used for training and the remaining 25% of the data for cross validation.

## 3.5   Word Classification using HMMs

Once a phoneme classifier has been created, we apply it onto the whole dataset of 1000 videos to obtain 1000 phoneme sequences of length 74. If we use phonemes, we will get sequences of numbers which are between 1 and 37. If we use visemes, we get sequences with numbers between 1 and 11.

Our objective is to map sequence of phonemes/visemes to words. This is where we utilise the capability of Hidden Markov Models to incorporate the temporal dynamics of the phoneme changes and thus "remember". So, even if our phoneme classification is erroneous, the patterns generated by the classifier can be mapped onto words resulting in higher accuracy.

We have the 1000 sequences.From the video dataset, **??** extracts sequences corresponding to each word (using regular expressions). This leaves us with a list of sequences for each word. For each word, the list is divided up into training and test sets. We train an HMM for each word. This means that each word has a unique transition matrix and observation matrix associated with it. All the training is done in **??**. The functions used for training the HMM are from Kevin Murphy's HMM toolbox [**?**].

Then, when we have an unknown sequence, we run it through the HMMs of all the words and find the "loglikelihood" that it might be that word. When we find the HMM for which

it has the highest "loglikelihood", it means that the sequence is most likely to be the word associated with that HMM.

# 4 Computational Results

## 4.1 Phoneme and Viseme Classification

For phoneme identification, the classification using a k-nearest neighbors algorithm only 11.6103% accuracy on average on the cross validation over 176 trials. For viseme identification, the k-nearest neighbors method had 19.7355% accuracy over 30 trials. More trials were not performed due to the Matlab *knnsearch* function being computationally intensive.

The naive bayes classificaiton algorithm applied to identificaiton of phonemes had an average accuracy of 3.49% over 1000 trials. For viseme classification, the naive bayes algorithm had an average accuracy of 9.1357% over 1000 trials.

## 4.2 Word Classification using HMMs

The identification of the right word definitely depends on the words it is compared against.

| Word Classification using 37 phonemes | | |
|---|---|---|
| Word | Set | Accuracy |
| bin | bin , blue | 87.5 % |
| blue | bin, blue | 36 % |
| blue | red, blue | 76 % |
| four | four, white | 60 % |
| bin | bin , white | 62.5 % |
| five | blue , five | 60 % |
| red | red , eight | 72 % |
| bin | bin , blue, green | 75 % |
| green | green, white, five | 44 % |
| five | five, blue, four, white | 50 % |
| green | green, white, five, red | 28 % |
| bin | bin , blue, green , red | 75 % |
| bin | bin , blue, red, white | 56.2 % |
| blue | bin , blue, red, white | 28 % |
| five | four , five, red, white | 45 % |
| bin | bin , blue, green , red, eight | 75 % |
| bin | bin , blue, green , red, white | 50 % |
| four | four , five, green , red, white | 30 % |
| five | four , five, green , red, white | 40 % |

| Word Classification using 11 Visemes | | |
|---|---|---|
| Word | Set | Accuracy |
| bin | bin , blue | 50 % |
| blue | bin, blue | 68 % |
| blue | red, blue | 92 % |
| four | four, white | 75 % |
| bin | bin , white | 68.8 % |
| five | blue , five | 65 % |
| red | red , eight | 65 % |
| bin | bin , blue, green | 25 % |
| four | bin , four, green | 35 % |
| red | red , white, green | 35 % |
| green | green, white, five | 32 % |
| green | green, white, blue | 36 % |
| five | five, blue, four, white | 50 % |
| green | green, white, five, red | 32 % |
| bin | bin , blue, green , red | 12.5 % |
| bin | bin , blue, red, white | 12.5% |
| blue | bin , blue, red, white | 68 % |
| five | four , five, red, white | 50 % |
| bin | bin , blue, green , red, eight | 12.5 % |
| bin | bin , blue, green , red, white | 12.5 % |
| four | four , five, green , red, white | 30 % |
| five | four , five, green , red, white | 25 % |

# 5 Summary and Conclusions

This report detailed the classification of phonemes and visemes based upon visual data of a person speaking, and speech prediction based upon identified phonemes and visemes using Hidden Markov Models. The results presented in the report were significantly limited by the quality of the detected lip contours. The videos used in this project were of low quality, making the accurate detection of lip contours difficult. The GRID database contains higher quality versions of the videos, however the database of higher quality videos are significantly larger file sizes. The hardware used for this project did not have the computing power to handle the larger database. The results are further limited by the assignment of an equal number of video frames to each phoneme from the given frame locations of each word. A more sophisticated method to determine which video frames corresponded to which phoneme would improve results. The problem of visual speech recognition is complex, and is usually solved using Long Short-Term Memory (LSTM) recurrent neural network and large video datasets. Considering the complexity and small relative size of the database used, the results presented here are reasonable.

# References

[1] J. Proctor, S. Brunton and J. N. Kutz, Dynamic mode decomposition with control, arXiv:1409.6358.

[2] Rabiner, Lawrence R. "A tutorial on hidden Markov models and selected applications in speech recognition." Proceedings of the IEEE 77.2 (1989): 257-286.

[3] Forney, G. David. "The viterbi algorithm." Proceedings of the IEEE 61.3 (1973): 268-278.

[4] Yang, Jie, and Yangsheng Xu. Hidden markov model for gesture recognition. No. CMU-RI-TR-94-10. CARNEGIE-MELLON UNIV PITTSBURGH PA ROBOTICS INST, 1994.

[5] Starner, Thad E. Visual Recognition of American Sign Language Using Hidden Markov Models. MASSACHUSETTS INST OF TECH CAMBRIDGE DEPT OF BRAIN AND COGNITIVE SCIENCES, 1995.

[6] Hassanat, Ahmad B., 'Visual Words for Automatic Lip- Reading.' PhD diss., University of Buckingham, 2009.

[7] J. Zhong, W. Chou, and E. Petajan, 'Acoustic Driven Viseme Identification for Face Animation.' Bell Laboratories. Murray Hill, NJ. IEEE 0-7803-378. Aug. 1997.

[8] L. Cappelletta and N. Harte. 'Phoneme-to-Viseme Mapping for Visual Speech.' Department of Electronic and Electrical Engineering, Trinity College Dublin, Ireland. May 2012.

[9] https://www.cs.ubc.ca/ murphyk/Software/HMM/hmm.html

[10] J. Barker, et al, 'The GRID Audiovisual Sentence Corpus.' University of Sheffield. 2013. http://spandh.dcs.shef.ac.uk/gridcorpus/

# A  MATLAB Code

## A.1  Contours.m

```matlab
obj=VideoReader('vid1.mpg');
vidFrames = read(obj);
numFrames = get(obj,'numberOfFrames');
[mov]= getmovout(vidFrames,numFrames-1);
X=frame2im(mov(50));
A=rgb2gray(X);
mask = zeros(size(A));
mask(400:450,320:400) = 1;
bw = activecontour(A,mask,300);
figure, imshow(bw), title('Active Contour mask');
[~, threshold] = edge(A, 'sobel');
fudgeFactor = .5;
BWs = edge(A,'sobel', threshold * fudgeFactor);
figure, imshow(BWs), title('binary gradient mask');
```

## A.2  Assign Labels

```matlab
function [ labels ] = assignlabels2(cropvid,frameLocs,words)
%UNTITLED10 Summary of this function goes here
%   Detailed explanation goes here
%Inputs:
%cropped videos (struct)
%frameLocs
% words
%
%
% Outputs:
% labels: cell array of cells containing the phoneme at each frame
%
%Begin Function

numVids = length(cropvid);

labels = cell(numVids,1);

for k = numVids:-1:1
    v = cropvid{k};
    fL = floor(frameLocs{k}/1000);
    lab = cell(size(v,3),1);
    for j = 1:length(fL)-1
        %Break this word into phonemes
        ph = assignphoneme(words{k}(j));
        %Get number of phonemes in this word
        numPh = length(ph);
        %get frame indices for this word to be distributed accross
        x1 = fL(j);
        x2 = fL(j+1);
        %number of frames per phoneme
        nf = round((x2-x1)/numPh);
        xprev = x1;
        if xprev == 0
```

12

```matlab
35          xprev = 1; %make sure 0 index isnt called
36      end
37      for i = 1:numPh−1
38          xnext = xprev + nf; %overwrite next
39          for ii = xprev:xnext−1
40              lab{ii} = ph{i};
41          end
42          xprev = xnext; %overwrite prev
43      end
44      for ii = xprev:size(v,3)
45          lab{ii} = ph{end};
46      end
47      end
48      labels{k} = lab; %store in output
49      clear lab %clear temp variable
50 end
51 end
```

## A.3  SVD and Classify

```matlab
1
2 % Required functions:
3
4 % lipcrop
5 % assignlabels
6 % assignphoneme
7
8 %Required variables to already be in workspace:
9 %
10 % frameLocs
11 % words
12 %Begin Script:
13
14 clearvars −except lipread frameLocs words
15 %% Crop, Assign Labels, create individual data matrices
16
17 %Crop videos based on contour mask
18 cropVid = lipcrop(lipread);
19
20 %Assign labels to each frame
21 labels = assignlabels2(cropVid,frameLocs,words);
22
23 %Sort and create data arrays
24 [X,tags,vinds] = hmmdata(cropVid,labels);
25 numPix = size(X,1);
26 numVids = size(X,2);
27
28 %% SVD
29
30 [U,S,V] = svd(X,'econ');
31
32 %% Create Classification Matrices
33
34 for k = 1:1000
35     q = randperm(numVids);
36     qind = round(numVids*0.75);
```

```matlab
37      q1 = q(1:qind);
38      q2 = q(qind+1:end);
39
40      trainData = V(q1,1:30);
41      testData = V(q2,1:30);
42      trainTags = tags(q1)';
43      testTags = tags(q2)';
44
45      hmmNBdata = V(:,1:30);
46
47      nb = NaiveBayes.fit(trainData,trainTags);
48      pre = nb.predict(testData);
49
50      acc(k) = 100*sum(pre==testTags)/length(pre);
51  end
52
53  %compute accuracy
54
55  %disp(['Accuracy was ' num2str(acc) '%'])
56
57  knnind = knnsearch(trainData,testData);
58  acc2 = 100*sum(knnind==testTags)/length(knnind);
59  % disp(['Accuracy was ' num2str(acc2) '%'])
```

## A.4   Creation of Data Matrix and Classification Labels

```matlab
1  function [X,tags,vinds] = hmmdata(cropVid,labels)
2  %UNTITLED Summary of this function goes here
3  %   Detailed explanation goes here
4
5  kmax =size(cropVid{1}(:,:,:),3);
6  numPix = size(cropVid{1}(:,:,:),1)*size(cropVid{1}(:,:,:),2);
7  numVids = kmax*length(cropVid);
8
9  %Initialize Data Matrix
10 counter = 1;
11 vinds = cell(length(cropVid),1);
12 X = zeros(numPix,36749);
13 for j = length(cropVid):-1:1
14     thisLab = labels{j};
15     thisVid = cropVid{j}(:,:,:);
16     numv = 1;
17     for k = 1:size(thisVid,3)
18         %Get DM index of this phoneme
19         phonemeInd = checkviseme(thisLab(k));
20
21         %If SIL, store in SIL array
22         if phonemeInd <= 0
23             continue
24         end
25         %this phoneme is a regular phoneme:
26
27         %Increment to store in right place
28
29         %Get Frame to store
30         thisFrame = thisVid(:,:,k);
```

14

```matlab
31
32          %Reshape frame
33          xframe = reshape(thisFrame,numPix,1);
34
35          %Store frame in corresponding DM
36          X(:,counter) = xframe;
37          tags(counter) = phonemeInd;
38          vinds{j}(numv) = counter;
39          counter = counter +1;
40          numv = numv + 1;
41      end
42 end
43 end
```

## A.5    makethehmm2.m

```matlab
 1 O = 36;
 2 Q = 37;
 3
 4 prior1 = normalise(rand(Q,1));
 5 transmat1 = mk_stochastic(rand(Q,Q));
 6 obsmat1 = mk_stochastic(rand(Q,O));
 7 binmat = binmatrix(1:200,1:7);
 8 binmat(binmat==0)=4;
 9 [LL, priorbin, transmatbin, obsmatbin] = dhmm_em(binmat, prior1, transmat1,
       obsmat1, 'max_iter', 500);
10 bluemat = bluematrix(1:200,1:7);
11 bluemat(bluemat==0)=4;
12 [LL, priorblue, transmatblue, obsmatblue] = dhmm_em(bluemat, prior1, transmat1,
       obsmat1, 'max_iter', 500);
13 whitemat = whitematrix(1:200,1:6);
14 whitemat(whitemat==0)=4;
15 [LL, priorwhite, transmatwhite, obsmatwhite] = dhmm_em(whitemat, prior1,
       transmat1, obsmat1, 'max_iter', 500);
16 greenmat = greenmatrix(1:200,1:6);
17 greenmat(greenmat==0)=4;
18 [LL, priorgreen, transmatgreen, obsmatgreen] = dhmm_em(greenmat, prior1,
       transmat1, obsmat1, 'max_iter', 500);
19 redmat = redmatrix(1:200,1:6);
20 redmat(redmat==0)=4;
21 [LL, priorred, transmatred, obsmatred] = dhmm_em(redmat, prior1, transmat1,
       obsmat1, 'max_iter', 500);
22 eightmat = eightmatrix(1:80,1:6);
23 eightmat(eightmat==0)=4;
24 [LL, prioreight, transmateight, obsmateight] = dhmm_em(eightmat, prior1,
       transmat1, obsmat1, 'max_iter', 500);
25 fourmat = fourmatrix(1:80,1:7);
26 fourmat(fourmat==0)=4;
27 [LL, priorfour, transmatfour, obsmatfour] = dhmm_em(fourmat, prior1, transmat1,
       obsmat1, 'max_iter', 500);
28 fivemat = fivematrix(1:80,1:8);
29 fivemat(fivemat==0)=4;
30 [LL, priorfive, transmatfive, obsmatfive] = dhmm_em(fivemat, prior1, transmat1,
       obsmat1, 'max_iter', 500);
31
32 setmat = setmatrix(1:200,1:8);
```

```matlab
setmat(setmat==0)=4;
[LL, priorset, transmatset, obsmatset] = dhmm_em(setmat, prior1, transmat1,
    obsmat1, 'max_iter', 500);

laymat = laymatrix(1:200,1:11);
laymat(laymat==0)=4;
[LL, priorlay, transmatlay, obsmatlay] = dhmm_em(laymat, prior1, transmat1,
    obsmat1, 'max_iter', 500);

placemat = placematrix(1:200,1:7);
placemat(placemat==0)=4;
[LL, priorplace, transmatplace, obsmatplace] = dhmm_em(placemat, prior1,
    transmat1, obsmat1, 'max_iter', 500);
```

## A.6  binsorter.m

```matlab
binind  =0;
blueind  =0;
whiteind  =0;
greenind  =0;
redind  =0;
layind  =0;
placeind  =0;
setind  =0 ;
zeroind  =0;
oneind  =0;
twoind  =0;
eightind =0;
fourind=0;
fiveind=0;
sevenind=0;
againind=0;
pleaseind=0;
soonind=0;
expr = '[blps][bgrw][abwi][abcdefghijklmnopqrstuvwxyz][0123456789][asnp].align';
expr1= 'b[bgrw][abwi][abcdefghijklmnopqrstuvwxyz][0123456789][asnp].align'; %bin
expr2= '[blps]b[abwi][abcdefghijklmnopqrstuvwxyz][0123456789][asnp].align'; %
    blue
expr3= '[blps]w[abwi][abcdefghijklmnopqrstuvwxyz][0123456789][asnp].align'; %
    white
expr4= '[blps]g[abwi][abcdefghijklmnopqrstuvwxyz][0123456789][asnp].align'; %
    green
expr5= '[blps]r[abwi][abcdefghijklmnopqrstuvwxyz][0123456789][asnp].align'; %red
expr6= 'l[bgrw][abwi][abcdefghijklmnopqrstuvwxyz][0123456789][asnp].align'; %lay
expr7= 'p[bgrw][abwi][abcdefghijklmnopqrstuvwxyz][0123456789][asnp].align'; %
    place
expr8= 's[bgrw][abwi][abcdefghijklmnopqrstuvwxyz][0123456789][asnp].align'; %set
expr9 = '[blps][bgrw][abwi][abcdefghijklmnopqrstuvwxyz]0[asnp].align';  % zero
expr10 = '[blps][bgrw][abwi][abcdefghijklmnopqrstuvwxyz]1[asnp].align';  % one
expr11 = '[blps][bgrw][abwi][abcdefghijklmnopqrstuvwxyz]2[asnp].align';  % two
expr12 = '[blps][bgrw][abwi][abcdefghijklmnopqrstuvwxyz]8[asnp].align';  % eight
expr13 = '[blps][bgrw][abwi][abcdefghijklmnopqrstuvwxyz]4[asnp].align';  % four
expr14 = '[blps][bgrw][abwi][abcdefghijklmnopqrstuvwxyz]5[asnp].align';  % five
expr15 = '[blps][bgrw][abwi][abcdefghijklmnopqrstuvwxyz]7[asnp].align';  % seven
expr16 = '[blps][bgrw][abwi][abcdefghijklmnopqrstuvwxyz][0123456789]a.align'; %
    again
```

```matlab
expr17 = '[blps][bgrw][abwi][abcdefghijklmnopqrstuvwxyz][0123456789]p.align'; %
    please
expr18 = '[blps][bgrw][abwi][abcdefghijklmnopqrstuvwxyz][0123456789]s.align'; %
    soon
for i = 1:1000
j=w(i,1);
video = predvis(i); % check this
vid = video{1};

framvec=  frameLocs(i);
framevec=framvec{1};

%for bin
s=regexp(j,expr1);
if isempty(s{1})
else
binind = binind +1;
a = floor(framevec(2)/1000);
b = floor(framevec(3)/1000);
c = b-a +1 ;
binmatrix(binind, 1:c) = vid(1:c);
end
%for blue
s=regexp(j,expr2);
if isempty(s{1})
else
blueind = blueind +1;
start = floor(framevec(2)/1000);
a = floor(framevec(3)/1000);
b = floor(framevec(4)/1000);
c1 = a-start +1 ;
c2 = b-start +1 ;
bluematrix(blueind, 1:c2-c1+1) = vid(c1:c2);
end;

%for white
s=regexp(j,expr3);
if isempty(s{1})
else
whiteind = whiteind +1;
start = floor(framevec(2)/1000);
a = floor(framevec(3)/1000);
b = floor(framevec(4)/1000);
c1 = a-start +1 ;
c2 = b-start +1 ;
whitematrix(whiteind, 1:c2-c1+1) = vid(c1:c2);
end;

%for green
s=regexp(j,expr4);
if isempty(s{1})
else
greenind = greenind +1;
start = floor(framevec(2)/1000);
a = floor(framevec(3)/1000);
b = floor(framevec(4)/1000);
c1 = a-start +1 ;
```

```matlab
91  c2 = b-start +1 ;
92  greenmatrix(greenind, 1:c2-c1+1) = vid(c1:c2);
93  end;
94
95  %for red
96  s=regexp(j,expr5);
97  if isempty(s{1})
98  else
99  redind = redind +1;
100 start = floor(framevec(2)/1000);
101 a = floor(framevec(3)/1000);
102 b = floor(framevec(4)/1000);
103 c1 = a-start +1 ;
104 c2 = b-start +1 ;
105 redmatrix(redind, 1:c2-c1+1) = vid(c1:c2);
106 end;
107
108 %for lay
109 s=regexp(j,expr6);
110 if isempty(s{1})
111 else
112 layind = layind +1;
113 start = floor(framevec(2)/1000);
114 a = floor(framevec(2)/1000);
115 b = floor(framevec(3)/1000);
116 c1 = a-start +1 ;
117 c2 = b-start +1 ;
118 laymatrix(layind, 1:c2-c1+1) = vid(c1:c2);
119 end;
120
121 %for place
122 s=regexp(j,expr7);
123 if isempty(s{1})
124 else
125 placeind = placeind +1;
126 start = floor(framevec(2)/1000);
127 a = floor(framevec(2)/1000);
128 b = floor(framevec(3)/1000);
129 c1 = a-start +1 ;
130 c2 = b-start +1 ;
131 placematrix(placeind, 1:c2-c1+1) = vid(c1:c2);
132 end;
133
134 %for set
135 s=regexp(j,expr8);
136 if isempty(s{1})
137 else
138 setind = setind +1;
139 start = floor(framevec(2)/1000);
140 a = floor(framevec(2)/1000);
141 b = floor(framevec(3)/1000);
142 c1 = a-start +1 ;
143 c2 = b-start +1 ;
144 setmatrix(setind, 1:c2-c1+1) = vid(c1:c2);
145 end;
146
147 %for zero
```

```matlab
148 s=regexp(j,expr9);
149 if isempty(s{1})
150 else
151 zeroind = zeroind +1;
152 start = floor(framevec(2)/1000);
153 a = floor(framevec(6)/1000);
154 b = floor(framevec(7)/1000);
155 c1 = a-start +1 ;
156 c2 = b-start +1 ;
157 zeromatrix(zeroind, 1:c2-c1+1) = vid(c1:c2);
158 end;
159
160 %for one
161 s=regexp(j,expr10);
162 if isempty(s{1})
163 else
164 oneind = oneind +1;
165 start = floor(framevec(2)/1000);
166 a = floor(framevec(6)/1000);
167 b = floor(framevec(7)/1000);
168 c1 = a-start +1 ;
169 c2 = b-start +1 ;
170 onematrix(oneind, 1:c2-c1+1) = vid(c1:c2);
171 end;
172
173 %for two
174 s=regexp(j,expr11);
175 if isempty(s{1})
176 else
177 twoind = twoind +1;
178 start = floor(framevec(2)/1000);
179 a = floor(framevec(6)/1000);
180 b = floor(framevec(7)/1000);
181 c1 = a-start +1 ;
182 c2 = b-start +1 ;
183 twomatrix(twoind, 1:c2-c1+1) = vid(c1:c2);
184 end;
185
186 %for eight
187 s=regexp(j,expr12);
188 if isempty(s{1})
189 else
190 eightind = eightind +1;
191 start = floor(framevec(2)/1000);
192 a = floor(framevec(6)/1000);
193 b = floor(framevec(7)/1000);
194 c1 = a-start +1 ;
195 c2 = b-start +1 ;
196 eightmatrix(eightind, 1:c2-c1+1) = vid(c1:c2);
197 end;
198
199 %for four
200 s=regexp(j,expr13);
201 if isempty(s{1})
202 else
203 fourind = fourind +1;
204 start = floor(framevec(2)/1000);
```

```matlab
205 a = floor(framevec(6)/1000);
206 b = floor(framevec(7)/1000);
207 c1 = a-start +1 ;
208 c2 = b-start +1 ;
209 fourmatrix(fourind, 1:c2-c1+1) = vid(c1:c2);
210 end;
211
212 %for five
213 s=regexp(j,expr14);
214 if isempty(s{1})
215 else
216 fiveind = fiveind +1;
217 start = floor(framevec(2)/1000);
218 a = floor(framevec(6)/1000);
219 b = floor(framevec(7)/1000);
220 c1 = a-start +1 ;
221 c2 = b-start +1 ;
222 fivematrix(fiveind, 1:c2-c1+1) = vid(c1:c2);
223 end;
224
225 %for seven
226 s=regexp(j,expr15);
227 if isempty(s{1})
228 else
229 sevenind = sevenind +1;
230 start = floor(framevec(2)/1000);
231 a = floor(framevec(6)/1000);
232 b = floor(framevec(7)/1000);
233 c1 = a-start +1 ;
234 c2 = b-start +1 ;
235 sevenmatrix(sevenind, 1:c2-c1+1) = vid(c1:c2);
236 end;
```

# B Python Code

## B.1 Phonemes.py

```python
import nltk

arpabet = nltk.corpus.cmudict.dict()
k =[j for j in 'abcdefghijklmnopqrstuvwxyz']
t= ['again', 'soon', 'now', 'please','bin', 'lay', 'place','set', 'blue', \
  'green','red','white' ,'at','by', 'with', 'in', 'zero','one', 'two',\
  'three','four','five','six','seven','eight','nine'];
g = k+t
ph =[]
for word in g:
        wl =arpabet[word]
        myString = ' '.join(str(r) for v in wl for r in v)
        print( word+' :'+ myString)
        for w in wl:
                ph = ph +w
uniqueph = set(ph)
```

## B.2  Transcripts.py

```python
import csv
import os
os.chdir("align")
beach = os.listdir()
with open("tes.csv", "w") as f:
        for sand in beach[:-1]:
                text_file = open(sand, "r")
                lines = text_file.read().split(',')
                k = lines[0]
                g = k.split()
                writer = csv.writer(f)
                writer.writerow(g)
```