

# Correlating Head Gestures of Speakers in political debates with Audience Engagement

Divyansh Agarwal (da2629)

January 5, 2015

## 1 Introduction

### 1.1 Problem Statement

The work is based on the hypothesis that certain speaker gestures can convey significant information that are correlated to audience engagement, and that people's engagement can be effectively reported by EEG data. The domain is focused on political debates, particularly the 2012 U.S. Presidential debates. The method works under the hypothesis that neural activity would be most similar across subjects at moments when something piques their interest. Therefore, the method identifies these moments by computing the points of maximum correlation of EEG measurements across subjects and comparing it against a significance derived from a permutation test. In this project we have tried to extend the work done in Correlating Speaker Gestures in Political Debates with Audience Engagement Measured via EEG (John R. Zhang, Jason Sherwin, et al) [1] by testing a new hypothesis, i.e. 'Correlating Head Gestures in Political Debates with Audience Engagement Measured via EEG'.

The data used in the work was collected by John R. Zhang et al. To collect audience engagement data, 6 videos totaling 61 minutes were shown to each of 20 human subjects in an electrostatically shielded room. 28 clips are selected out of the 291 clips in the recorded video of the first (original air date October 3, 2012) and third (original air date October 22, 2012) debates between Obama and Romney: 8 from the first debate and 6 from the third debate featuring Obama, and 8 from the first debate and 6 from the third debate featuring Romney. The total duration of the clips is approximately 30.5 minutes, so each clip is 65 seconds on average. The 28 clips are randomly shuffled and 28 silent versions were made to make the subjects more concentrated on the image features rather than the words and sound. The EEG data were smoothed over a 5-second window and the frequency was finally 1HZ.

### 1.2 Possibilities of other features influencing audience engagement

To identify the relevant features piquing audience attention, an effort was made to see the clips multiple times and empirically identify the moments that captured the attention of the viewer. However, with so many different and concurring movements, it was difficult to ascertain whether these features triggered an

interest collectively or in a stand alone manner. Also, some of the movements were too recurring and might not always trigger interest. Following are some interesting observations about the dataset and the various facial features:

1. Eyebrows: Romney had a prominent eyebrow movement whenever he was emphasizing a point. A rising up of eyebrows was observed.
2. Head movement: changes in direction of head, head shaking and nodding. In first debate obama's static head position lead to boredom. Once again the clips belonging to Romney showed a much greater movement in head movement.
3. Fluency and pace of speaking could be tested as one of the cue capturing attention

### **1.3 Head gestures role in day to day conversation**

- dialogue structuring, nodding and shaking to show agreement and disagreement
- changes in head position right before a current listener takes turn at speaking
- gaze and head direction can be used for referencing to objects or points
- as a reinforcement agent
- quick head movement as a sign of alarm
- mutual gaze - a sign of awareness

## **2 Exploration of affective attributes influencing audience engagement**

During the initial phases of the project, we explored the possibility of correlating speaker's emotion with audience engagement. The heuristic was derived from the simple understanding that in normal human interaction often the emotions conveyed by the speaker do influence our perception and interest, for instance a smiling and confident person can charm the audience easily as opposed to a angry or a disinterested person. The dive into literature led to a number of interesting approaches and challenges.

### **2.1 Description of Affect**

Affect has been described by psychologists in terms of six discrete categories. These basic emotion categories include happiness, sadness, fear, anger, disgust and surprise. The advantage of such a classification is that people use this categorical scheme to describe observed emotional displays in daily life. However, this discrete lists of emotions fails to describe the range of emotions in natural communication settings which involve a lot of subtle and interchanging signals. An alternative to the categorical description is the dimensional description where an affective state is characterized in terms of small number of latent

dimensions rather than in terms of a small number of discrete emotion categories. These dimensions include Activation, Expectation, Power and Valence. Activation (Arousal) is the individual’s global feeling of dynamism or lethargy. Power (Dominance) dimension subsumes two related concepts, power and control. Valence is whether the person feels positive or negative. Expectation is whether the person is expecting something or is caught unaware (surprised). In contrast to categorical representation, dimensional representation enables to capture a range of emotions as a combination of the four dimensions. However, this representation is non-intuitive.[2]

There have been work addressing both categories of emotion description. The efforts on learning classifiers for different discrete emotion categories have been persistent for a while. The Extended Cohn-Kanade Dataset (CK+)[3] provides a complete dataset for action unit and emotion-specified expression. The target expression for each sequence is fully FACS coded and emotion labels have been revised and validated. Baseline result with Active Appearance Model and Support Vector Machine has been provided. The motivation behind the dataset is to remove ambiguity regarding the different methods reporting accuracy and results on different datasets. Though the results are promising for emotions detection but the dataset is meant for exaggerated expressions.

Recently AVEC (Audio/Visual Emotion Challenge)-2011 and 2012[4][5] have made a conscious effort towards advancing emotion recognition systems to be able to deal with naturalistic behavior in large volumes of un-segmented, non-prototypical and non-preselected data as this is exactly the type of data that both multimedia retrieval and human-machine/human-robot communication interfaces have to face in the real world. The Audio/Visual Emotion Challenge 2011/2012 (AVEC2012) data contains 32 continuous, frontal face videos of an interview subject being emotionally engaged. It is more challenging than previous data sets because the data is not acted or sponsored by the interviewer. The subject is speaking with an interviewer. As a result, emotions are subtle and more difficult to detect than in other data sets. AVEC-2011[4] put forth the problem of learning binary classifier for both valence and arousal domain. AVEC - 2012[5] put fourth the regression problem as opposed to the classification problem in AVEC - 2011. But, the results are only able to achieve an accuracy of around 50 percent, thus the idea of correlating speaker’s emotions with audience engagement, though interesting, looks a bit premature.

### 3 Head Gesture Detection

#### 3.1 Model for identifying key facial landmarks:

An open-source adaptation of CVPR-2014 paper is used to identify 68 key facial points [One Millisecond Face Alignment with an Ensemble of Regression Trees by Vahid Kazemi and Josephine Sullivan][6]. It uses an ensemble of regression trees to estimate the key facial points directly from a sparse subset of pixel intensities. The landmark detection takes approximately milliseconds per frame and achieves accuracy greater than or comparable to state of the art methods on standard dataset.

The open source implementation (in dlib C++ library) uses a small dataset to train a face land-marking model using the program **trainShapePredic-**

**tor.cpp.** The learnt model is represented as a .dat file. For the purpose of the project the pre learnt sp.dat file was used. The Library can be installed using the following link: <http://sourceforge.net/projects/dclib/files/latest/download> and is critical for working of the later parts of the project. The Figure 1 shows the location of various landmarks. Individual frames of the video can be processed to simultaneously track different facial landmarks thereby capturing the various movement of head, eye, nose, mouth, etc as shown in Figure 2.

Listing 1: C++ code to get the facial landmark information for each frame

```

1 #include <dlib/image_processing/frontal_face_detector.h>
2 #include <dlib/image_processing/render_face_detections.h>
3 #include <dlib/image_processing.h>
4 #include <dlib/gui_widgets.h>
5 #include <dlib/image_io.h>
6 #include <iostream>
7 #include <cv.h>
8 #include <highgui.h>
9 #include <dlib/opencv.h>
10 #include <fstream>
11 using namespace dlib;
12 using namespace std;
13 #define SKIP_RATE 1
14
15
16 int main(int argc, char** argv)
17 {
18     try
19     {
20         // The code takes in a shape model file (.dat) and name of the video to
21         if (argc == 1)
22         {
23             cout << "Call this program like this:" << endl;
24             cout << "./face_landmark_detection-ex shape_predictor_68-face_landm
25             return 0;
26         }
27
28         // We need a face detector.
29         frontal_face_detector detector = get_frontal_face_detector();
30         // And we also need a shape_predictor.
31         // loading the model from the shape_predictor_68-face_landmarks.dat
32         shape_predictor sp;
33         deserialize(argv[1]) >> sp;
34
35         ofstream myfile;
36         string filename;
37         cout << "Enter file to process: \n";
38         cin >> filename;
39         string landmarkInfo = filename + ".txt";
40         myfile.open(&landmarkInfo[0]);

```

```

41     int frameCount = 0;
42     string videoName = filename + ".avi";
43     cv::VideoCapture cap( videoName); // open the video file for reading
44     if ( !cap.isOpened() )
45     {
46         cout << "Cannot open the video file" << endl;
47         return -1;
48     }
49     double fps = cap.get(CV_CAP_PROP_FPS); //get the frames per seconds of
50     double totalFramesInVideo = cap.get(CV_CAP_PROP_FRAME_COUNT);
51     double capturedFrames = (int)(totalFramesInVideo/SKIP_RATE);
52     cout << "Frame per seconds : " << fps << ", Total Frames: " << totalFr
53     cv::Mat frame;
54
55     myfile << capturedFrames << endl;
56     while(1)
57     {
58         bool bSuccess = cap.read(frame); // read a new frame from video
59         if (!bSuccess) //if not success, break loop
60         {
61             cout << "Cannot read the frame from video file" << endl;
62             break;
63         }
64         frameCount++;
65
66         myfile << frameCount << endl;
67         cout << "Processing frame " << frameCount << endl;
68         array2d<rgb_pixel> img;
69         assign_image(img, cv_image<bgr_pixel>(frame));
70
71         // Now tell the face detector to give us a list of bounding boxes
72         // around all the faces in the image.
73         std::vector<rectangle> dets = detector(img);
74
75         // Now we will go ask the shape_predictor to tell us the pose of
76         // each face we detected.
77         std::vector<full_object_detection> shapes;
78         full_object_detection shape;
79         // Save the position of each facial landmark for each frame in a fi
80         if(dets.size() > 0)
81         {
82             for (unsigned long j = 0; j < 1; ++j) // Chk for just the first
83             {
84                 shape = sp(img, dets[j]);
85                 for( int parts = 0; parts < shape.num_parts(); parts++)
86                     myfile << shape.part(parts) << endl;
87                 shapes.push_back(shape);
88             }
89         }
90         else

```

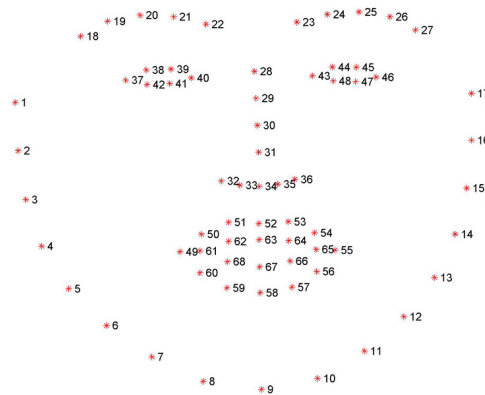
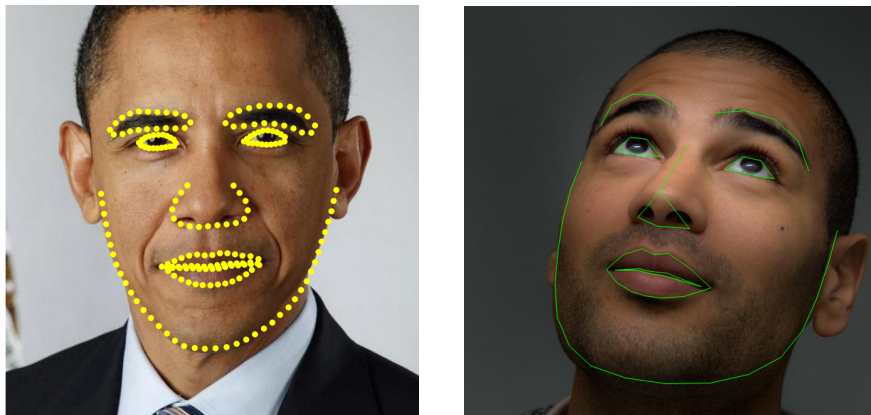


Figure 1: Anatomy of Key facial landmarks (68)



(a) Subfigure

(b) Subfigure

Figure 2: 68-Key facial landmarks annotated

```

91             shapes.push_back(shape);
92
93         }
94         myfile.close();
95
96     }
97     catch (exception& e)
98     {
99         cout << "\nexception thrown!" << endl;
100        cout << e.what() << endl;
101    }
102 }
```

---

## 3.2 Head pose/head normal estimation using facial landmarks

### 3.2.1 Background for human head pose estimation

Human head pose estimation is a hard problem because it suffers from classical problems in computer vision such as viewpoint variation, illumination variation and occlusion. It has been a topic of active research and number of approaches has been tried. For the purpose of the project, we felt that the geometric approaches are best suited. Nonetheless, a succinct analysis for different head pose estimation methods is provided in the following paragraphs. Refer to [7] for detailed discussion.

Appearance based modeling methods generally consider the head pose as a collection of face images from multiple viewpoints. In a training phase, a model is presented with a set of face images of known pose. In an online phase, a query face image is compared to a learned model to find the most faithful match. However without interpolation, these appearance based approaches are only capable of estimating discrete pose locations. Moreover, they require large amounts of training data, and are highly sensitive to small changes in face position and image scaling.

Flexible models such as active appearance models (AAM) show good invariance to head localization error since they adapt to the image and find the location of the facial features. Although precise, using 3D head models and image registration techniques are computationally intensive and require a high image resolution that could slow down realtime performance.

Geometric based approaches can potentially achieve accurate pose measurements by tracking the location of only a few facial features such as the eyes, nose, and mouth. The computational efficiency of tracking a few features and determining the head orientation from their configuration allows real-time performance and is well suited for our application. The downside is that feature detection and tracking are critical to the accuracy of the pose estimation results, and can lead to failure when features are lost. However, the implementation for landmark detection, as discussed in the previous section, is fairly robust and provide a near real time estimates of facial landmarks.

### 3.2.2 Geometric methods for Head pose estimation

There are various ways in which the facial configuration can be exploited for head pose estimation. For the purpose of this project, we adopted the approach of [8].

The facial model is based on the ratios of four world lengths  $L_f, L_n, L_e$  and  $L_m$  as shown in figure 4. The far corners of the eyes and mouth demarcate the facial plane.  $L_f$  and  $L_m$  are measured on the symmetry axis of the facial plane while  $L_e$  is measured perpendicular to the symmetry axis of the facial plane.  $L_n$  is the distance between the nose tip and the nose base. It is measured along the normal to the facial plane, the facial normal. The corresponding distances in the image are denoted by small case letters. The method also required two model ratios:  $R_m = L_m/L_f$  and  $R_n = L_n/L_f$

The symmetry axis is found by joining the midpoint between the eyes to the mouth, and the nose base is located using the model ratio  $R_m = L_m/L_f$ . The angle  $\tau$  is the angle between the projection of face normal in the image and

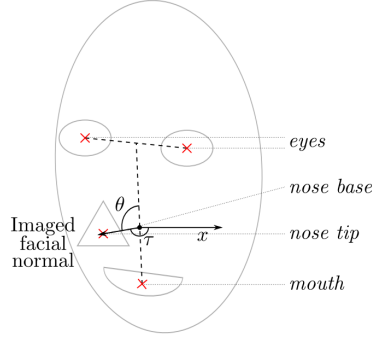


Figure 3: Facial Model [8]

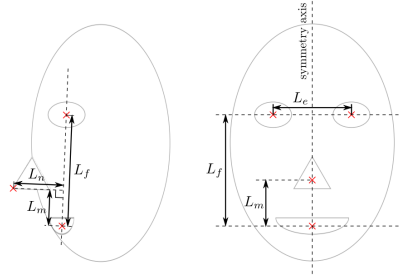


Figure 4: Profile (left) and frontal (right) views of a face. The face model is composed of distance ratios between the eye, nose and mouth positions of a typical face. The facial plane is defined by the locations of the eyes and the mouth. The pitch and yaw angles are estimated by assuming that the face is planar, and the nose is a vector normal to this plane.[8]

the x-axis as can be seen 3. The angle  $\sigma$  which the facial normal makes with the vector  $d$  normal to the image plane is called the slant, as shown in figure 5. This angle may be found by using the image measurements  $\theta, l_n, l_m$  and the model ratio  $R_n = L_n/L_f$ . Please refer to [8] for supporting equations. It can be shown that in camera centered coordinates, the facial normal  $\hat{n}$  is given by:

$$\hat{n} = [\sin \sigma \cos \tau, \sin \sigma \sin \tau, -\cos \sigma]$$

Listing 2: C++ code to get the normal estimates (normalEstimation.cpp)

```

1
2 #include <iostream>
3 #include <cv.h>
4 #include <highgui.h>
5 #include <fstream>
6 #include <vector>
7 #include <utility>
8 #include <string>

```



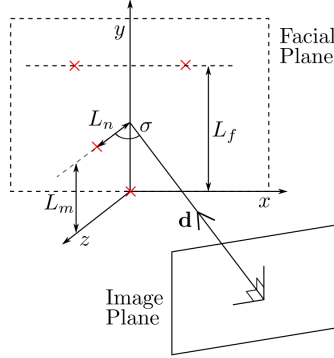


Figure 5: The face centred coordinate system used to calculate the slant angle  $\sigma$ . The facial slant  $\sigma$  is the angle between the vector  $d$  normal to the image plane, and the facial normal located along the  $z$  axis of the facial plane. [8]

```

9  #include <math.h>
10 using namespace std;
11 using namespace cv;
12
13 void facialLandmarks(string landmarkInfo, vector< vector< pair<int,int> > >& la
14     numFramesCaptured)
15 {
16     ifstream infile(&landmarkInfo[0]);
17     //vector< vector< pair<int,int> > > landmarks; // max no of frames 100 and
18     vector<int> frameIndex;
19     //int numFramesCaptured = 0;
20     infile >> numFramesCaptured;
21     cout << numFramesCaptured << endl;
22     for( int i=0; i<numFramesCaptured; i++)
23     {
24         int f;
25         infile >> f;
26         frameIndex.push_back(f);
27         vector< pair<int,int> > myvec;
28         for( int parts=0; parts < NUMPARTS; parts++)
29         {
30             int r,c;
31             string s2,s3;
32             char c1;
33             infile >> c1 >> r >> s2 >> c >> s3;
34             myvec.push_back(make_pair(r,c));
35         }
36         landmarks.push_back(myvec);
37     }
38 }
39 void artifacts(Point eyeMid, Point mouthMid, Point noseTip, Point noseBase,
40     Mat& frame)

```

```

41 {
42     circle(frame, eyeMid, 1 , CV_RGB(255,255,0), 1);
43     circle(frame, mouthMid, 1 , CV_RGB(255,255,0), 1);
44     circle(frame, noseBase, 1 , CV_RGB(255,255,0), 1);
45     circle(frame, noseTip, 1 , CV_RGB(255,255,0), 1);
46     line(frame, noseBase, noseBase + Point(10,0), CV_RGB(0,0,255),1);
47     line(frame, noseBase, noseBase - Point(0,-10), CV_RGB(0,0,255),1);
48 }
49 }
50 void normalEstimates(double nVec[], Point eyeMid, Point mouthMid,
51                     Point noseTip, Point noseBase )
52 {
53     double tilt = acos( (noseTip - noseBase).dot(Point(1,0))
54                        /norm(noseTip - noseBase) );
55     double len_nose = norm(noseTip - noseBase);
56     double len_face = norm(eyeMid - mouthMid);
57     double m1 = pow(len_nose/len_face ,2);
58     double m2 = pow( (noseTip - noseBase).dot(Point(0,-1))/
59                    norm(noseTip - noseBase), 2);
60     double Rn = RATIO_NOSE;
61     double dz = ( Rn*Rn - m1 - 2*m2*Rn*Rn + sqrt( pow( m1 - Rn*Rn, 2)+4*m1*m2*Rn*Rn )
62                / ( 2*Rn*Rn*(1-m2) );
63     dz = sqrt(dz);
64     double sigma = acos( dz );
65     nVec[0] = sin(sigma)*cos(tilt);
66     nVec[1] = sin(sigma)*sin(tilt);
67     nVec[2] = -cos(sigma);
68 }
69 int main()
70 {
71     string filename, landmarkFilename, normalFilename;
72     int numFramesCaptured = 0;
73     vector< vector< pair<int,int> > > landmarks; // max no of frames 100 and on
74     cout << "Video to be processed: \n";
75     cin >> filename;
76     landmarkFilename = filename + ".txt";
77
78     //Get Facial Landmark Information for all the frames
79     facialLandmarks(landmarkFilename, landmarks, numFramesCaptured);
80
81     ofstream normals;
82     normalFilename = filename + "_normal.txt";
83     normals.open(&normalFilename[0]);
84     normals << numFramesCaptured << endl;
85     int frameCount = 0;
86     while( frameCount < numFramesCaptured - 1 )
87     {
88         frameCount++;
89         Point eyeMid( (landmarks.at(frameCount).at(37).first +
90                     landmarks.at(frameCount).at(46).first)/2,

```

```

91         (landmarks.at(frameCount).at(37).second +
92             landmarks.at(frameCount).at(46).second)/2 );
93
94     Point mouthMid((landmarks.at(frameCount).at(49).first +
95                     landmarks.at(frameCount).at(55).first)/2,
96                     (landmarks.at(frameCount).at(49).second +
97                     landmarks.at(frameCount).at(55).second)/2 );
98     Point noseTip(landmarks.at(frameCount).at(30).first ,
99                  landmarks.at(frameCount).at(30).second );
100    Point noseBase( 0.6*mouthMid.x + 0.4*eyeMid.x ,
101                   0.6*mouthMid.y+ 0.4*eyeMid.y);
102
103    //To plot various facial artifacts
104    artifacts( eyeMid, mouthMid, noseTip, noseBase, frame);
105
106    // Normal Vector Estimation
107    double nVec[3];
108    normalEstimates( &nVec[0], eyeMid, mouthMid, noseTip, noseBase );
109
110    // IF nVec[0] = NaN then assign a default 0 value
111    if( nVec[0] != nVec[0])
112        normals << 0 << " " << 0 << " " << 0 << endl;
113    else
114        normals << nVec[0] << " " << nVec[1] << " " << nVec[2] << endl;
115    }
116    normals.close();
117 }

```

---

### 3.2.3 Head movement estimation: Speed and Acceleration

The velocity and acceleration estimation is done in  $y$  any  $x$  direction, i.e., along the face symmetry axis and the axis perpendicular to it on the facial plane by utilizing the  $y$  and  $x$  component of the estimated normal vector. The normal estimation is done at 30 fps. However, we need the sample the data to 1 Hz since the EEG data is sampled at 1 Hz. Also, the difference between two frames is insignificant ( $\frac{1}{30}$  th of a second) to register any significant movement between two consecutive frames. Thus, the speed is sampled at 1 Hz and is given by:

$$speed = abs(avg. disp. in first 15 frames - avg. disp in last 15 frames)$$

By extension, the acceleration in both directions are sampled at 1 Hz by utilizing the sampled speed. A couple of different experiments/approaches were tried for sampling the data. However, as part of the project we had difficulty in ascertaining the usefulness of one sampling method over other due to lack of quantitative measures. Finally, the above method was adopted for sampling.

Listing 3: C++ code to get the head motion estimates and output video (motionDetection.cpp)

---

```

1
2 #include <iostream>

```

```

3  #include<cv.h>
4  #include<highgui.h>
5  #include<fstream>
6  #include<vector>
7  #include<utility>
8  #include<string>
9  #include<math.h>
10 using namespace std;
11 using namespace cv;
12 #define XVELOCITY 0.08
13 #define WINDOW 15
14
15 void displayOnFrame( double xVelocity, double yVelocity, double xAcc, double yA
16 {
17     char xvelStr[100], xaccStr[100];
18     char yvelStr[100], yaccStr[100];
19     sprintf(xvelStr,"X Velocity: %lf ", xVelocity );
20     sprintf(xaccStr,"X Acceleration: %lf ", xAcc );
21     sprintf(yvelStr,"Y Velocity: %lf ", yVelocity );
22     sprintf(yaccStr,"Y Acceleration: %lf ", yAcc );
23     if( xVelocity > 0.05)
24         putText(frame, xvelStr, Point(20,20), CV_FONT_NORMAL, 0.5,
25             Scalar(255,255,0),1,1);
26     else
27         putText(frame, xvelStr, Point(20,20), CV_FONT_NORMAL, 0.5,
28             Scalar(255,255,255),1,1);
29     if( yVelocity > 0.05)
30         putText(frame, yvelStr, Point(300,20), CV_FONT_NORMAL, 0.5,
31             Scalar(255,255,0),1,1);
32     else
33         putText(frame, yvelStr, Point(300,20), CV_FONT_NORMAL, 0.5,
34             Scalar(255,255,255),1,1);
35     if( xAcc > 0.05)
36         putText(frame, xaccStr, Point(20,40), CV_FONT_NORMAL, 0.5,
37             Scalar(255,255,0),1,1);
38     else
39         putText(frame, xaccStr, Point(20,40), CV_FONT_NORMAL, 0.5,
40             Scalar(255,255,255),1,1);
41     if( yAcc > 0.05)
42         putText(frame, yaccStr, Point(300,40), CV_FONT_NORMAL, 0.5,
43             Scalar(255,255,0),1,1);
44     else
45         putText(frame, yaccStr, Point(300,40), CV_FONT_NORMAL, 0.5,
46             Scalar(255,255,255),1,1);
47 }
48 void velocity1(string velocity1Filename, int &numFramesProcessed, vector< vector<
49 {
50     int numSeconds = numFramesProcessed / 30;
51     ofstream velocityInfo(&velocity1Filename[0]);
52     velocityInfo << numSeconds << endl;

```

```

53
54     int frameCount = 0;
55     cout << numFramesProcessed << endl;
56     double prevXvelocity = 0;
57     double prevYvelocity = 0;
58     while(frameCount < numFramesProcessed)
59     {
60         double Xdisp1 = 0, Xdisp2 = 0;
61         double Ydisp1 = 0, Ydisp2 = 0;
62         for( int i = 0; i < WINDOW; i++)
63         {
64             Xdisp1 = Xdisp1 + nVector.at(frameCount).at(0);
65             Ydisp1 = Ydisp1 + nVector.at(frameCount).at(1);
66             frameCount++;
67         }
68         for( int i = 0; i < WINDOW; i++)
69         {
70             Xdisp2 = Xdisp2 + nVector.at(frameCount).at(0);
71             Ydisp2 = Ydisp2 + nVector.at(frameCount).at(1);
72             frameCount++;
73         }
74         double Xvelocity = abs(Xdisp1-Xdisp2)/WINDOW;
75         double Yvelocity = abs(Ydisp1-Ydisp2)/WINDOW;
76         double Xacc = Xvelocity - prevXvelocity;
77         double Yacc = Yvelocity - prevYvelocity;
78         velocityInfo << frameCount/30 << " " << Xvelocity << " " << Yvelocity
79         prevXvelocity = Xvelocity;
80         prevYvelocity = Yvelocity;
81     }
82     velocityInfo.close();
83 }
84 void outputVideo(string inputVideoname, string outputVideoname, string velocity)
85 {
86
87     // open the video files for reading and writing
88     cv::VideoCapture cap(inputVideoname);
89     cv::VideoWriter output(outputVideoname, cap.get(CV_CAP_PROP_FOURCC), cap.get(CV_CAP_PROP_FPS),
90         cv::Size(cap.get(CV_CAP_PROP_FRAME_WIDTH), cap.get(CV_CAP_PROP_FRAME_HEIGHT)));
91     if (!output.isOpened())
92     {
93         std::cout << "!!! Output video could not be opened" << std::endl;
94         return;
95     }
96
97     if ( !cap.isOpened() )
98     {
99         cout << "Cannot open the video file" << endl;
100         return;
101     }
102

```

```

103 // Reading from velocityInfo file
104 ifstream velocityInfo(&velocityFilename[0]);
105 int numSeconds;
106 vector< vector<double> > motionVector;
107 if (!velocityInfo.is_open() )
108 {
109     cout << "ERROR: Cant open normal information file \n";
110     exit(1);
111 }
112 velocityInfo >> numSeconds;
113 char c1;
114 for( int i = 0; i < numSeconds; i++)
115 {
116     vector<double> myvec;
117     double val;
118     for( int j = 0; j < 5; j++ )
119     {
120         velocityInfo >> val;
121         myvec.push_back(val);
122     }
123     motionVector.push_back(myvec);
124 }
125 velocityInfo.close();
126 // Writing Video
127 double fps = cap.get(CV_CAP_PROP_FPS); //get the frames per seconds of the
128 double width = cap.get(CV_CAP_PROP_FRAME_WIDTH);
129 double height = cap.get(CV_CAP_PROP_FRAME_HEIGHT);
130 cout << "Height: " << height << " Width: " << width;
131 cout << " Frame per seconds : " << fps << endl;
132 cv::namedWindow("MyVideo"); //create a window called "MyVideo"
133 cv::Mat frame;
134 int frameCount = 0;
135 while(1)
136 {
137
138     int second = frameCount / 30;
139     bool bSuccess = cap.read(frame); // read a new frame from video
140     if (!bSuccess) //if not success, break loop
141     {
142         cout << "Cannot read the frame from video file" << endl;
143         break;
144     }
145     if( frameCount >= numSeconds*30)
146         break;
147     //if( frameCount % 5 != 0 )
148     // continue;
149     displayOnFrame( motionVector.at(second).at(1), motionVector.at(second).
150         motionVector.at(second).at(3), motionVector.at(second).at(4),
151     frame);
152     output.write(frame);

```

```

152         frameCount++;
153     }
154     output.release();
155     cap.release();
156 }
157 }
158 int main()
159 {
160     string videoname;
161     cout << "Video name: \n";
162     cin >> videoname;
163     string inputVideoname = videoname + ".avi";
164     string outputVideoname = videoname + "_out.avi";
165     string normalFilename = videoname + "_normal.txt";
166     string velocity1Filename = videoname + "_velocity1.txt";
167     cout << "Files: " << inputVideoname << " " << outputVideoname << "
" << normalFilename << " " << velocity1Filename << endl;
168
169     //Getting the nVec's
170     int numFramesCaptured;
171     vector< vector<double> > nVector;
172     ifstream normalInfo(&normalFilename[0]);
173     if(!normalInfo.is_open())
174     {
175         cout << "ERROR: Cant open normal information file \n";
176         exit(1);
177     }
178     normalInfo >> numFramesCaptured;
179     for( int i = 0; i < numFramesCaptured-1; i++)
180     {
181         vector<double> myvec;
182         char c1;
183         double val;
184         for( int j = 0; j < 3; j++ )
185         {
186             normalInfo >> val;
187             myvec.push_back(val);
188         }
189         normalInfo >> c1;
190         nVector.push_back(myvec);
191     }
192     normalInfo.close();
193     cout << numFramesCaptured << endl;
194
195     int numFramesProcessed = numFramesCaptured - 1 - (numFramesCaptured-1)%30;
196     velocity1(velocity1Filename, numFramesProcessed, nVector);
197     outputVideo(inputVideoname, outputVideoname, velocity1Filename);
198
199 }

```

---

Listing 4: Matlab code for correlation between the head motion estimates and output video (headEEGcorrelation.m)

---

```

1 % remove the first line from velocity_1.txt file
2 file = input('Filename: ', 's');
3
4 EEG = csvread(strcat(file, '_G.txt'));
5 % fp = fopen( strcat( file , '_velocity1.txt '));
6 % velocityInfo = textscan(fp, '%d %f %f %f %f');
7
8 f1 = fopen(strcat(file, '_normal.txt'));
9 positionInfo = textscan(f1, '%f %f %f');
10 meanX = mean(positionInfo{1})
11 stdX = std(positionInfo{1})
12
13 figure, plot(smooth(positionInfo{1}, 10));
14 title('X position over all frames (Moving Avg)');
15 % disp(['In X dimension, Mean: ', num2str(meanX), ' and standard Deviation: ',
16
17
18 % figure, plot(smooth(positionInfo{2}));
19 % title('Y position over all frames (Moving Avg)');
20 % disp(['In Y dimension, Mean: ', num2str(mean(positionInfo{2})), ' and standar
21
22
23 PC1 = EEG(1, :);
24 % velocityX = velocityMode(positionInfo, 10, 10, 30);
25 velocityX = velocityMean(positionInfo, 10, 10, 30);
26
27 figure;
28 subplot(2, 1, 1);
29 plot(PC1(1: size(velocityX, 2)));
30 title('PC1');
31
32 subplot(2, 1, 2);
33 plot(velocityX);
34 title('X velocity');
35 line([0, size(velocityX, 2)], [mean(velocityX), mean(velocityX)], 'Color', [1, 0, 0]);
36 saveas(gcf, strcat(file, '_velocityX-EEG.fig'));
37 meanVelocityX = mean(velocityX)
38 stdVelocityX = std(velocityX)
39 CorrelationVelX = corr2(PC1(1: size(velocityX, 2)), velocityX)

```

---

### 3.3 Results and Observations

The results have been pretty disappointing as we fail to observe any consistent pattern. The best case correlation with EEG data was 0.32 for clip TPO1.avi (The clips have been renamed but an index is provided mapping each file to its original name) and for some clips ( obama2, obama3 ) even a negative correlation is observed suggesting that the rise in one signal coincides with fall in other



Filename	X-Mean Pos.	X-Std. Dev. Pos	X-Mean Speed	X S.D. Speed	Correlation
TPO1	- 0.086	0.13	0.15	0.13	0.32
TPO2	0.013	0.13	0.13	0.09	0.18
TPO4	0.060	0.11	0.12	0.06	0.20
obama1	-0.030	0.23	0.13	0.09	0.12
obama2	0.200	0.15	0.17	0.11	-0.16
obama3	0.060	0.15	0.18	0.12	-0.08
obama4	0.12	0.14	0.16	0.11	0.15
romney5	-0.300	0.21	0.18	0.19	0.02
romney7	0.34	0.20	0.14	0.16	0.21
romney8	-0.320	0.16	0.15	0.15	0.17

Table 1: Results

signal. The results for the video clips are in Table:1. In Fig: 6 - 12, the red line corresponds to mean speed in x-direction.

We cannot confidently comment that the head gestures does not influence audience engagement because:

- Small set of people, i.e only twenty people, might be too small a sample set
- Small size of clips - 65 seconds is too small an interval to make an inference
- A different sampling algorithm can be applied to capture the head movement from position estimates of face landmarks.
- We are not sure how good the speed (velocity) and acceleration features are in capturing the higher level features of head movement, i.e nodding and shaking. It would be interesting to perform some experiments to preprocess the raw head movement signal.
- We understand in the Zang et al paper, a different correlation algorithm was applied, therefore it might be a good idea to experiment with different correlation algorithms, particularly something that is concerned only with the growth and decay of signals at corresponding points rather than their magnitudes.

## 4 References

1. John R. Zhang, Jason Sherwin et al. Correlating Speaker Gestures in Political Debates with Audience Engagement Measured via EEG. Proceedings of the ACM International Conference on Multimedia, 387-396.
2. Zhihong Zeng, Maja Pantic et al. A Survey of Affect Recognition Methods: Audio, Visual, and Spontaneous Expressions. IEEE Transactions on Pattern Analysis And Machine Intelligence

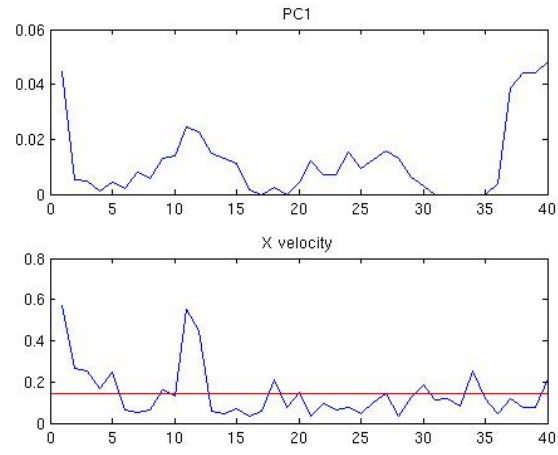


Figure 6: TPO1, Correlation: 0.32

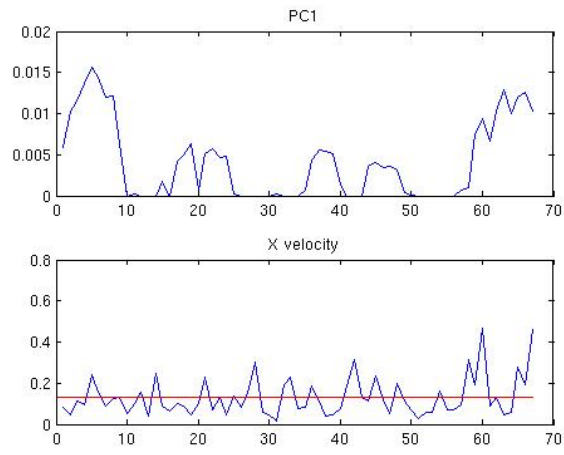


Figure 7: TPO2, Correlation: 0.18

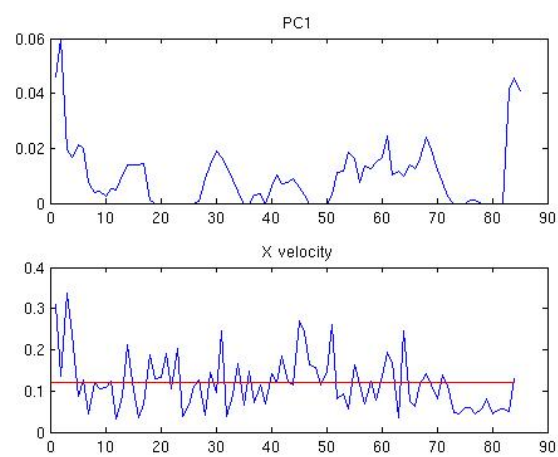


Figure 8: TPO4, Correlation: 0.20

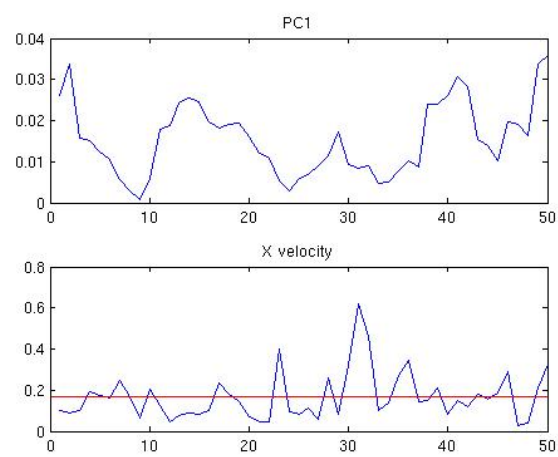


Figure 9: obama2, Correlation: -0.16

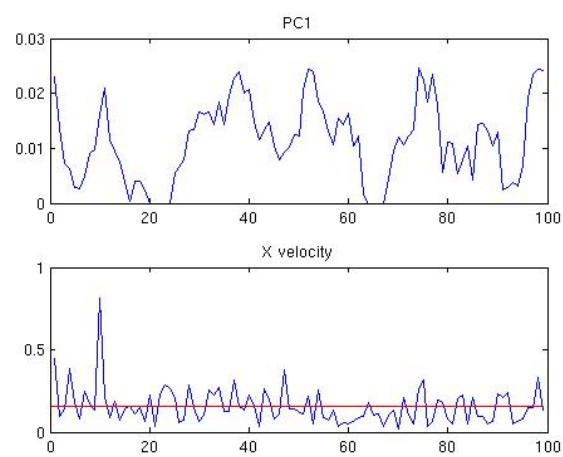


Figure 10: obama4, Correlation: 0.15

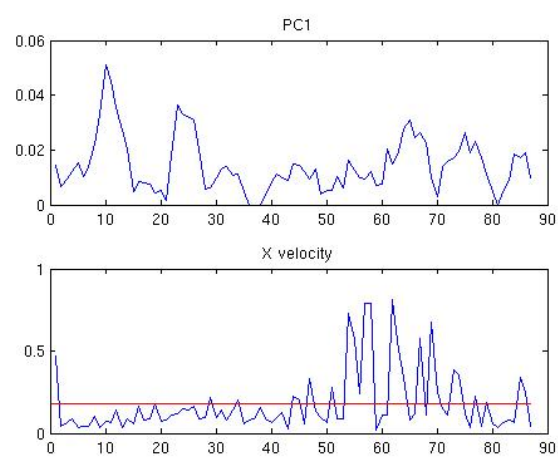


Figure 11: romney5, Correlation: 0.02

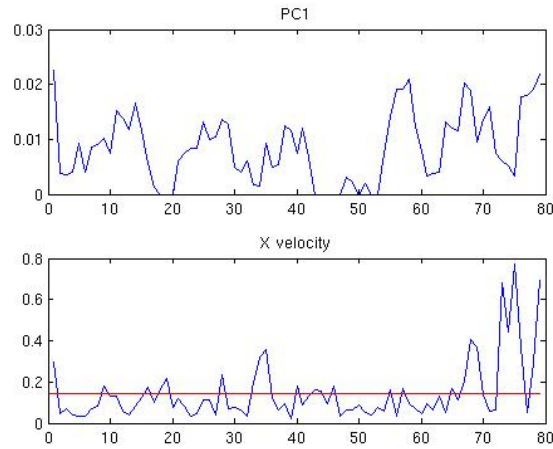


Figure 12: romney8, Correlation: 0.17

3. Patrick Lucey, Jeffrey F. Cohn, Takeo Kanade, Jason Saragih, Zara Ambadar. The Extended Cohn-Kanade Dataset (CK+): A complete dataset for action unit and emotion-specified expression. (CVPR - 2010).
4. Bjorn Schuller, Michel Valstar, Maja Pantic et al. AVEC 2011 – The First International Audio/Visual Emotion Challenge.
5. Bjorn Schuller, Michel Valstar, Maja Pantic et al. AVEC 2012 – The Continuous Audio/Visual Emotion Challenge.
6. Kazemi, Vahid, and Sullivan Josephine. "One Millisecond Face Alignment with an Ensemble of Regression Trees." (CVPR - 2014).
7. Murphy-Chutorian, E and Trivedi, M.M. . Head Pose Estimation in Computer Vision: A Survey. Pattern Analysis and Machine Intelligence, IEEE Transactions, 607-626.
8. A. H. Gee and R. Cipolla, Determining The Gaze Of Faces In Images, Image and Vision Computing - 1994, vol 12, 639 - 647.