

000  
 001  
 002  
 003  
 004  
 005  
 006  
 007  
 008  
 009  
 010  
 011  
 012  
 013  
 014  
 015  
 016  
 017  
 018  
 019  
 020  
 021  
 022  
 023  
 024  
 025  
 026  
 027  
 028  
 029  
 030  
 031  
 032  
 033  
 034  
 035  
 036  
 037  
 038  
 039  
 040  
 041  
 042  
 043  
 044  
 045  
 046  
 047  
 048  
 049  
 050  
 051  
 052  
 053

## Supplementary Materials for “Watch Those Words: Video Falsification Detection Using Word-Conditioned Facial Motion”

054  
 055  
 056  
 057  
 058  
 059  
 060  
 061  
 062  
 063  
 064  
 065  
 066  
 067  
 068  
 069  
 070  
 071  
 072  
 073  
 074  
 075  
 076  
 077  
 078  
 079  
 080  
 081  
 082  
 083  
 084  
 085  
 086  
 087  
 088  
 089  
 090  
 091  
 092  
 093  
 094  
 095  
 096  
 097  
 098  
 099  
 100  
 101  
 102  
 103  
 104  
 105  
 106  
 107

Anonymous CVPR submission

Paper ID 11596

## 1. Overview

Here we provide some of the quantitative and qualitative results to support the analysis made in the main paper. In the main paper we compared with the related works using the average AUCs across all individuals. In order to give a better insight into the comparison, we first present per-individual results for each of the related works. We then present a more detailed version of qualitative results using both training and testing datasets.

### 1.1. Comparison with State-Of-The-Art

Shown in Table 1 are the per-individual results for all the related methods that were presented in the main paper.

### 1.2. Videos for Qualitative Analysis

Here we provide the videos used for qualitative analysis of the words presented in the Figure 1 and Figure 6 of the main paper. For Obama, Trump, and Oliver we provide occurrences of the word “hi”, “tremendous”, and “billion” in the real and fake videos. Therefore, there are a total of six videos for this section: Obama\_hi\_{real,fake}.mp4, Trump\_tremendous\_{real,fake}.mp4, and Oliver\_billion\_{real,fake}.mp4. In each video, the output probability of the word-specific classifier is shown in red on the top left corner (a value of 1 is for real and 0 is fake). The occurrences of the words are selected from the training dataset. This is done to demonstrate the facial gestures associated with specific words during training.

In each case, it can be observed that a specific facial gesture is present in real videos which is missing in the fake videos. For example, the occurrences of the word “hi” is associated with an upward head movement which is missing in the fake examples. Similarly, in case of the word “tremendous”, notice the presence of lip rounding and chin raise action in multiple occurrences of the word in real videos, whereas these actions are missing in the fake videos.

XceptionNet					
	Audio Dubbing	Wav2Lip	Impersonator	FaceSwap	in-the-wild
Obama	0.50	0.94	0.74	0.96	0.47
Trump	0.50	0.84	0.70	0.82	0.54
Biden	0.50	0.49	0.69	0.67	0.45
Harris	0.50	0.80	0.48	0.24	-
O’ Brien	0.50	0.69	0.44	0.11	-
Oliver	0.50	0.93	0.26	0.15	-
PWL					
	Audio Dubbing	Wav2Lip	Impersonator	FaceSwap	in-the-wild
Obama	0.5	0.56	0.96	0.96	0.83
Trump	0.5	0.51	0.95	0.94	0.41
Biden	0.5	0.53	0.65	0.66	0.55
Harris	0.5	0.45	0.94	0.94	-
O’ Brien	0.5	0.84	0.69	0.67	-
Oliver	0.5	0.88	0.99	0.93	-
LipForensics					
	Audio Dubbing	Wav2Lip	Impersonator	FaceSwap	in-the-wild
Obama	0.50	1.00	0.83	1.00	0.98
Trump	0.50	1.00	0.68	0.98	0.97
Biden	0.50	0.93	0.15	0.30	0.91
Harris	0.50	1.00	0.08	0.71	-
O’ Brien	0.50	0.96	0.48	0.90	-
Oliver	0.50	0.97	0.39	0.98	-
ID-Reveal					
	Audio Dubbing	Wav2Lip	Impersonator	FaceSwap	in-the-wild
Obama	0.50	0.77	0.81	0.71	0.59
Trump	0.50	0.66	0.92	0.88	0.77
Biden	0.50	0.47	0.75	0.59	0.47
Harris	0.50	0.73	0.98	0.98	-
O’ Brien	0.50	0.66	0.63	0.56	-
Oliver	0.50	0.69	0.98	0.93	-

Table 1. Accuracy in terms of AUC on 10-second video clips for the six individuals and five different video falsification scenarios. The average AUC across all individuals is given in the last row. From top-bottom are the AUCs for XceptionNet, PWL, LipForensics and ID-Reveal.

### 1.3. Word Analysis for in-the-wild videos

Here we show how the results of our method can be interpreted during the evaluation of a test video. For this we provide four example videos, a real and a fake video of

108	Obama and Trump. The real videos are from test-split of	162
109	real dataset and fake videos are from in-the-wild dataset.	163
110	The videos are named as Obama_{itw, real}_test.mp4 and	164
111	Trump_{itw, real}_test.mp4.	165
112	Given a test video of 10-second length, we show the out-	166
113	put of word-specific classifier for each word. Shown on the	167
114	x-axis of the plot is time and on the y-axis is the proba-	168
115	bility that the word occurrence is real. Shown in orange	169
116	is the probability of the word in the test video and shown	170
117	in the blue is the average real probability of the word in	171
118	real dataset during training. The region in blue indicates the	172
119	standard deviation of training probability. The gaps in the	173
120	plot indicate that the word-specific classifier was missing.	174
121	The current time is indicated by the red dot on the plot and	175
122	the current word is displayed on the top of the video.	176
123	These word-level probabilities, can be used to isolate	177
124	the words which obtain low probability of being real. For	178
125	example, in Obama.itw.test.mp4 many words have a low	179
126	probability of being real with a minimum probability of zero	180
127	for the word “coverage”. Similarly in Trump_itw_test.mp4	181
128	video, the word “protected” has the zero probability of be-	182
129	ing real. Whereas in the videos Obama_real_test.mp4 and	183
130	Trump_real_test.mp4, the real probability for each of the	184
131	words is close to training real dataset (average of 0.8).	185
132	Shown in Figure 1 are the distributions of the 25 facial-	186
133	gesture features for the word “coverage” for Obama. In	187
134	each panel, shown in blue is the distribution of one facial-	188
135	gesture feature in real training videos of Obama. Shown	189
136	with red line is the value of facial-gesture feature in the	190
137	current test video of Obama which in this case is the fake	191
138	video shown in Obama_itw_test.mp4. The word “coverage”	192
139	in this example fake video have an out-of-distribution value	193
140	for AU26 i.e. jaw drop. The out-of-distribution value can	194
141	also be observed for lip-ver motion where the value in the	195
142	fake is lower than any of the value seen during training.	196
143	Similarly, shown in Figure 2 are the distributions of	197
144	the 25 facial-gesture features for the word “protect” for	198
145	Trump. The red line in each panel is the value of facial-	199
146	gesture feature in the fake test video of Trump shown in	200
147	Trump_itw_test.mp4. For the word “protect” the value for	201
148	AU17 (chin raise) and AU23 (lip tightner) in the fake is	202
149	lower than any of the value seen during training.	203
150		204
151		205
152		206
153		207
154		208
155		209
156		210
157		211
158		212
159		213
160		214
161		215

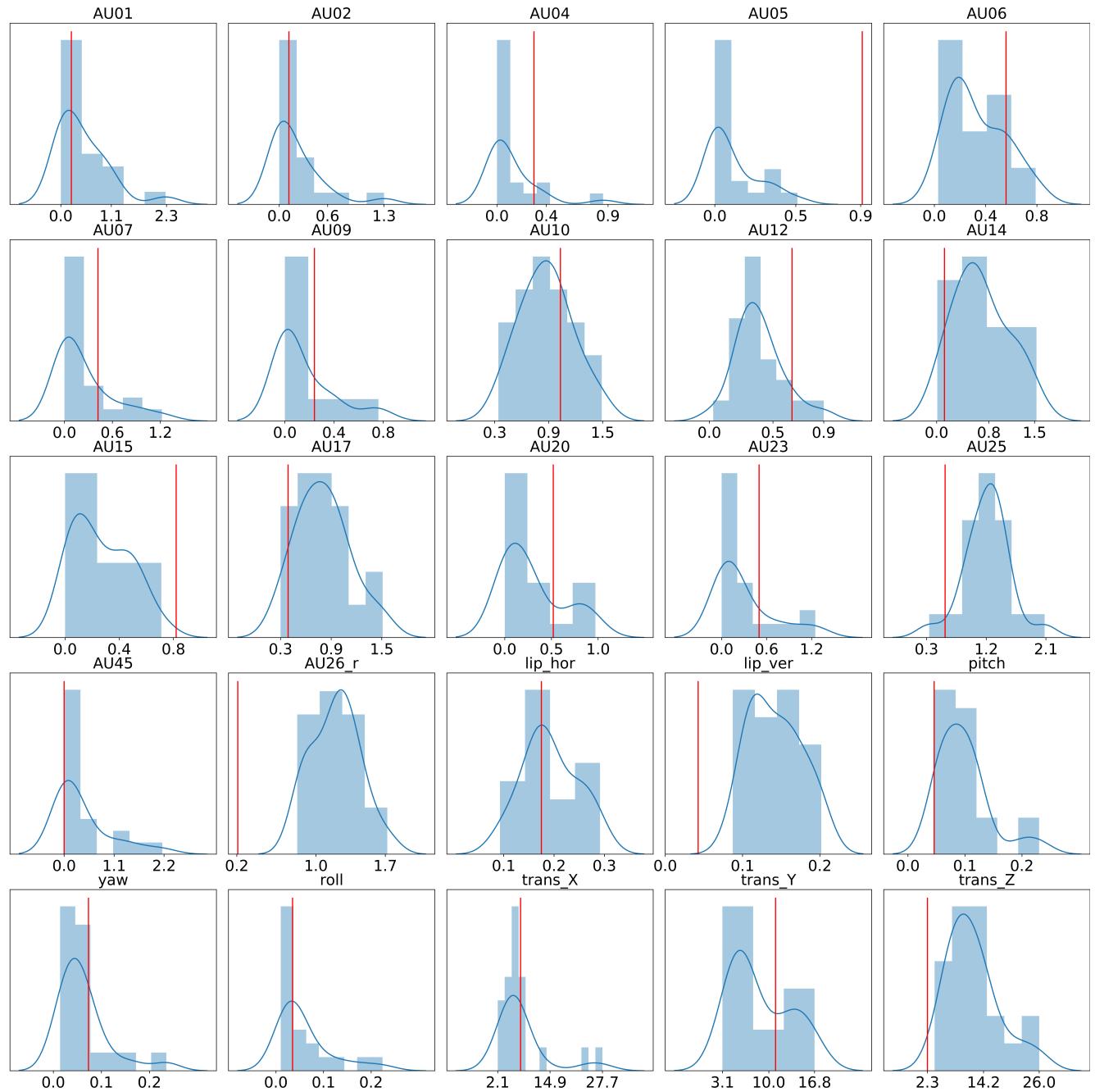


Figure 1. In each panel, shown in blue is the distribution of one facial-gesture feature in real training videos of Obama for the word “coverage”. The name of the facial-gesture feature is given on top of the panel. Shown with red line is the value of the facial feature in the current test video of Obama which in this case is the fake video shown in Obama\_itw\_test.mp4.

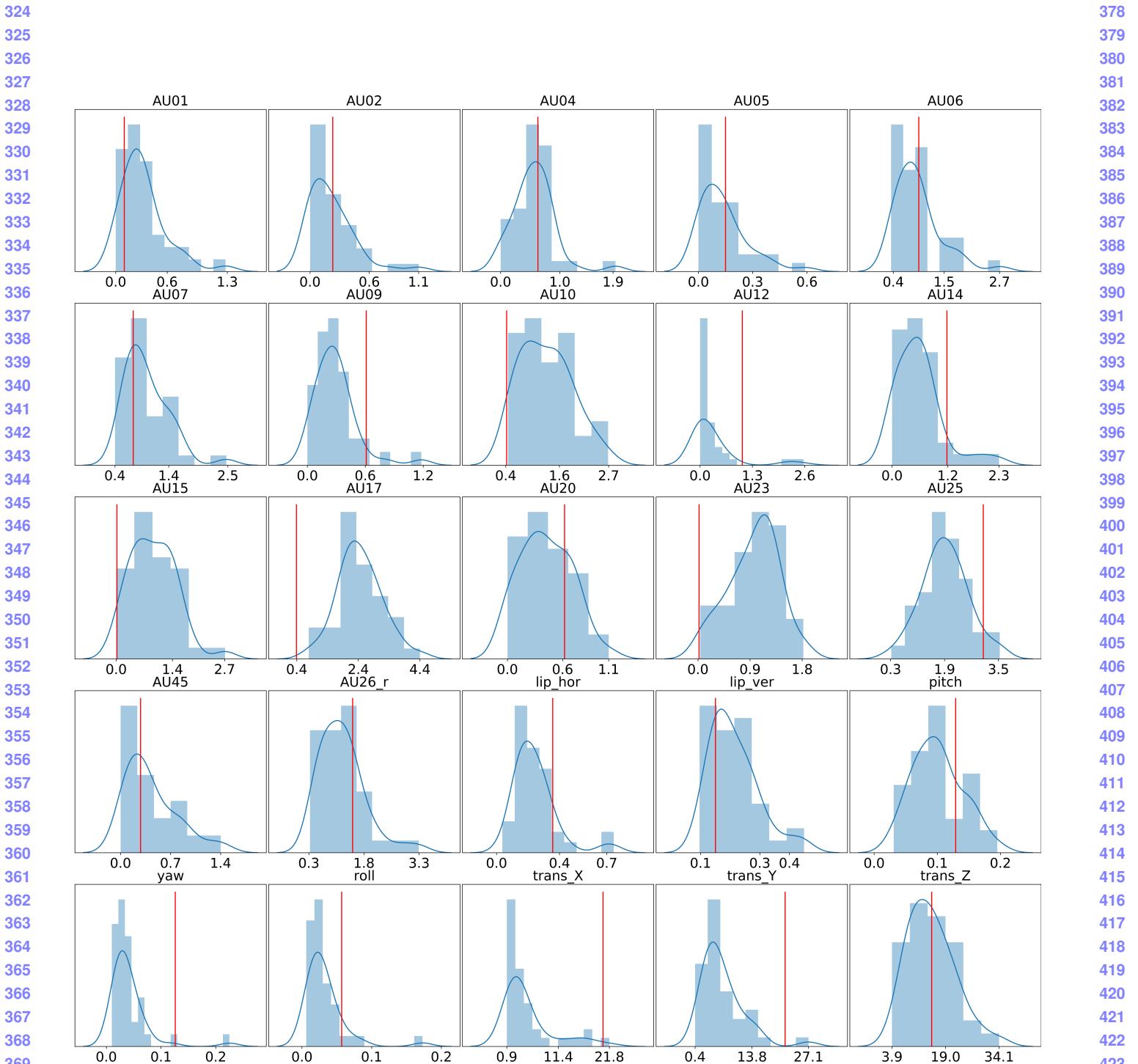


Figure 2. In each panel, shown in blue is the distribution of one facial-gesture feature in real training videos of Trump for the word “protect”. The name of the facial-gesture feature is given on top of the panel. Shown with red line is the value of the facial feature in the current test video of Trump which in this case is the fake video shown in `Trump_itw_test.mp4`.