

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”

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

Anonymous CVPR submission

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

Paper ID 11596

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

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.mp4](#), [Obama_hi_fake.mp4](#), [Trump_tremendous_real.mp4](#), [Trump_tremendous_fake.mp4](#), and [Oliver_billion_real.mp4](#), [Oliver_billion_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

108 real and a fake video of Obama and Trump. The
109 real videos are from test-split of real dataset and fake
110 videos are from in-the-wild dataset. The videos are
111 named as: `Obama_itw_test.mp4`, `Obama_real_test.mp4`,
112 `Trump_itw_test.mp4`, and `Trump_real_test.mp4`.
113

114 Given a test video of 10-second length, we show the out-
115 put of word-specific classifier for each word. Shown on the
116 x-axis of the plot is time and on the y-axis is the proba-
117 bility that the word occurrence is real. Shown in orange
118 is the probability of the word in the test video and shown
119 in the blue is the average real probability of the word in
120 real dataset during training. The region in blue indicates the
121 standard deviation of training probability. The gaps in the
122 plot indicate that the word-specific classifier was missing.
123 The current time is indicated by the red dot on the plot and
124 the current word is displayed on the top of the video.

125 These word-level probabilities, can be used to isolate
126 the words which obtain low probability of being real. For
127 example, in `Obama_itw_test.mp4` many words have a low
128 probability of being real with a minimum probability of zero
129 for the word “coverage”. Similarly in `Trump_itw_test.mp4`
130 video, the word “protected” has the zero probability of be-
131 ing real. Whereas in the videos `Obama_real_test.mp4` and
132 `Trump_real_test.mp4`, the real probability for each of the
133 words is close to training real dataset (average of 0.8).

134 Shown in Figure 1 are the distributions of the 25 facial-
135 gesture features for the word “coverage” for Obama. In
136 each panel, shown in blue is the distribution of one facial-
137 gesture feature in real training videos of Obama. Shown
138 with red line is the value of facial-gesture feature in the
139 current test video of Obama which in this case is the fake
140 video shown in `Obama_itw_test.mp4`. The word “coverage”
141 in this example fake video have an out-of-distribution value
142 for AU26 i.e. jaw drop. The out-of-distribution value can
143 also be observed for lip-ver motion where the value in the
144 fake is lower than any of the value seen during training.

145 Similarly, shown in Figure 2 are the distributions of
146 the 25 facial-gesture features for the word “protect” for
147 Trump. The red line in each panel is the value of facial-
148 gesture feature in the fake test video of Trump shown in
149 `Trump_itw_test.mp4`. For the word “protect” the value for
150 AU17 (chin raise) and AU23 (lip tightner) in the fake is
151 lower than any of the value seen during training.

152
153
154
155
156
157
158
159
160
161

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

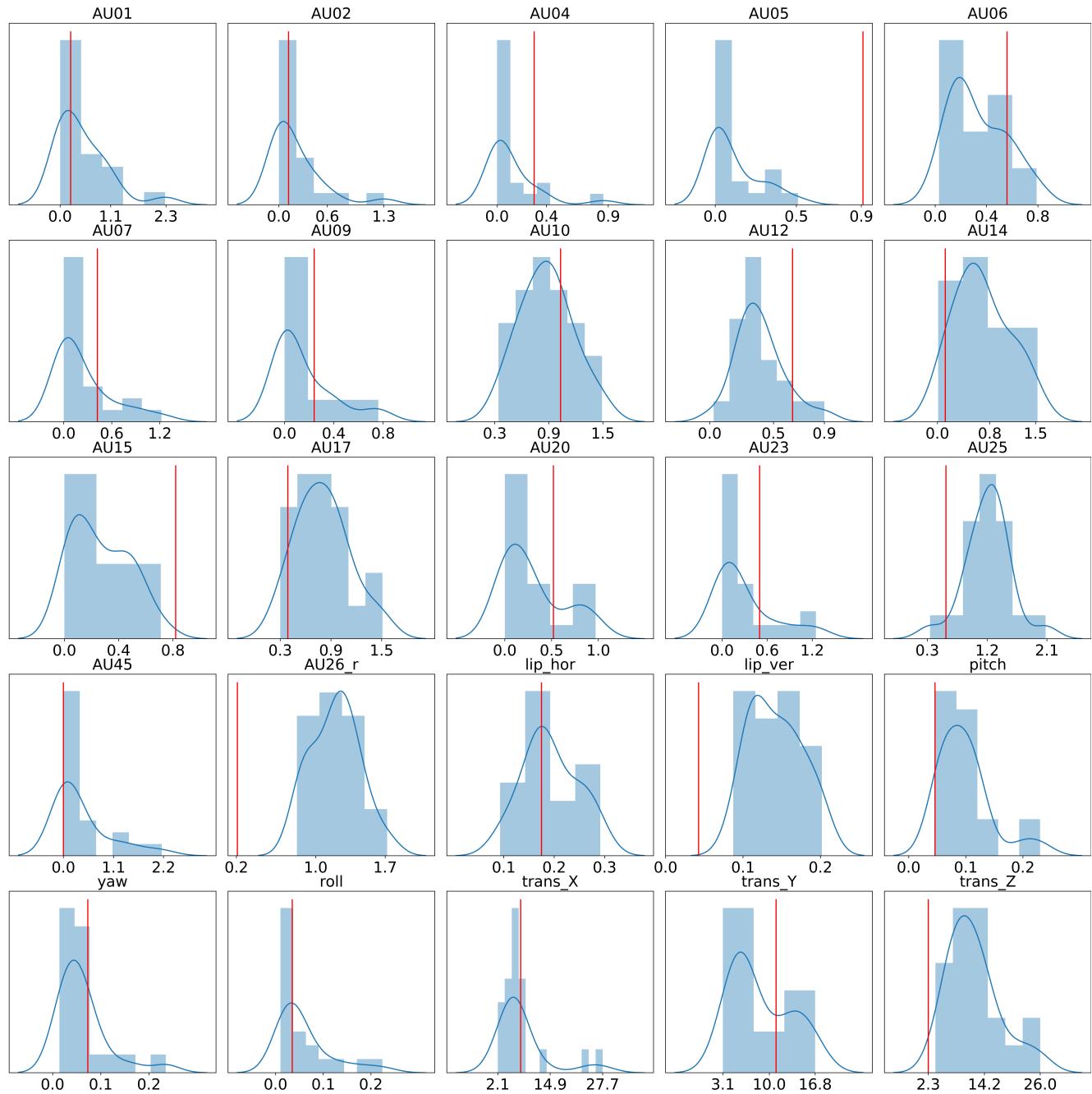
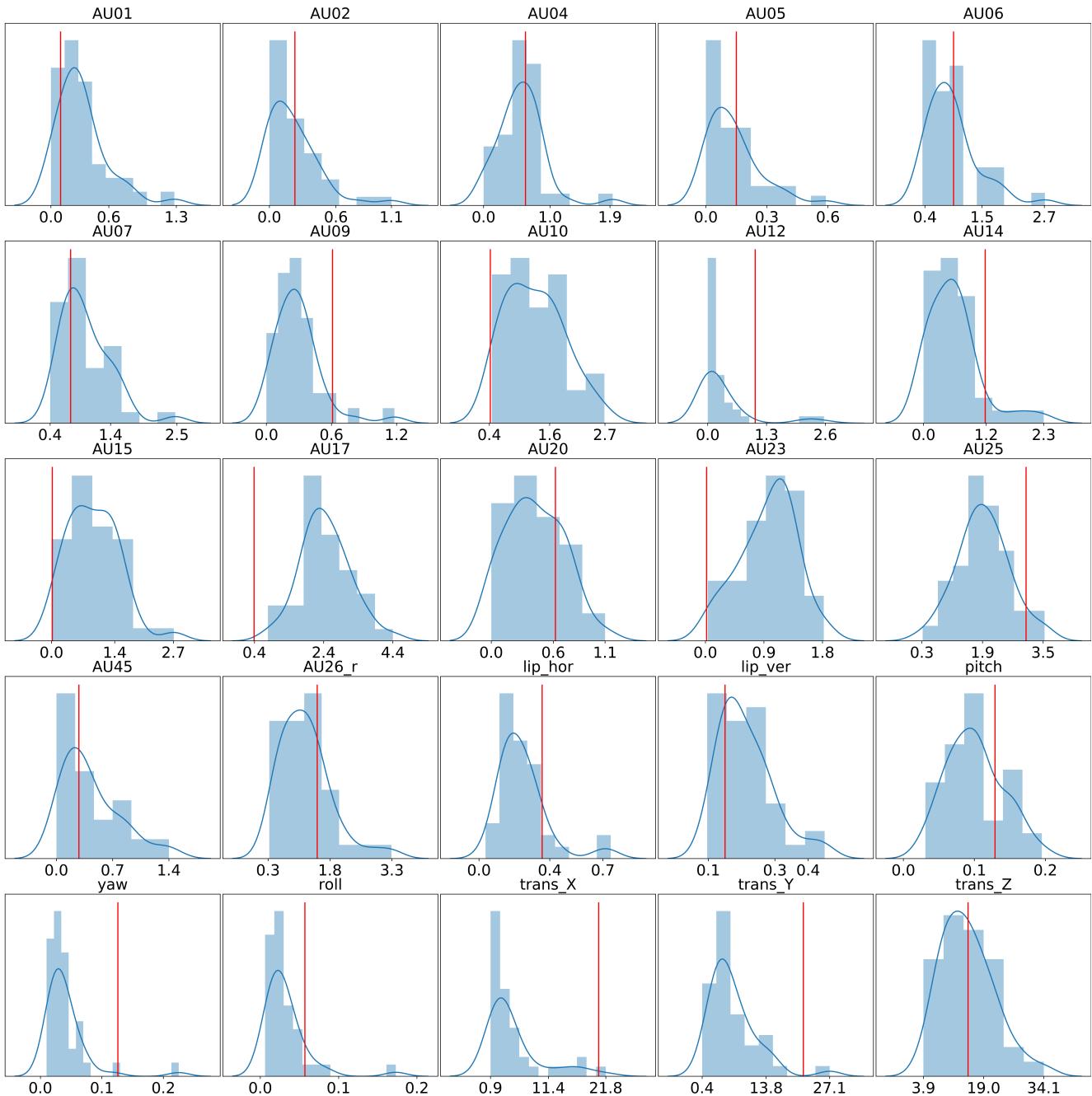


Figure 1. In each panel, shown in blue is the distribution of one facial-gesture feature in real training videos of Obama for the word “coverage”. The name of the facial-gesture feature is given on top of the panel. Shown with red line is the value of the facial feature in the current test video of Obama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



371 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
372 “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
373 current test video of Trump which in this case is the fake video shown in Trump_itw_test.mp4.