
Multitask Learning for Face Forgery Detection: A Joint Embedding Approach

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Multitask learning for face forgery detection has experienced impressive successes
2 in recent years. Nevertheless, the semantic relationships among different forgery
3 detection tasks are generally overlooked in previous methods, which weakens
4 knowledge transfer across tasks. Moreover, previously adopted multitask learning
5 schemes require human intervention on allocating model capacity to each task and
6 computing the loss weighting, which is bound to be suboptimal. In this paper,
7 we aim at automated multitask learning for face forgery detection from a joint
8 embedding perspective. We first define a set of coarse-to-fine face forgery detection
9 tasks based on face attributes at different semantic levels. We describe the ground-
10 truth for each task via a textural template, and train two encoders to jointly embed
11 visual face images and textual descriptions in the shared feature space. In such a
12 manner, the semantic closeness between two tasks is manifested as the distance
13 in the learned feature space. Moreover, the capacity of the image encoder can be
14 automatically allocated to each task through end-to-end optimization. Through joint
15 embedding, face forgery detection can be performed by maximizing the feature
16 similarity between the test face image and candidate textual descriptions. Extensive
17 experiments show that the proposed method improves face forgery detection in
18 terms of generalization to novel face manipulations. In addition, our multitask
19 learning method renders some degree of model interpretation by providing human-
20 understandable explanations.

21 1 Introduction

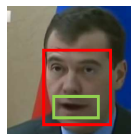
22 The emergence of deep generative models [1, 34, 67, 71] has significantly simplified and automated
23 the process of generating realistic counterfeit face images, popularly known as DeepFake. The
24 prevalence of falsified face images can erode the reliability and credibility of digital visual information.
25 Additionally, the exploitation and manipulation of such technologies pose a threat to individual rights
26 and national security.

27 Traditional DeepFake detectors were largely influenced by classic photo forensics [21] to expose
28 forgery traces by examining statistical anomalies [51, 58], visual artifacts [32, 46, 50, 51, 59],
29 and physical and geometric inconsistencies [15, 33, 35, 56]. With the rapid development of deep
30 learning, there has recently been a growing consensus on exploiting multitask learning for face
31 forgery detection [8, 10, 19, 41, 55, 80, 81]. The underlying assumption is that the primary task
32 (*i.e.*, global face forgery classification) is likely to benefit from other highly relevant auxiliary tasks
33 through knowledge transfer. Representative auxiliary tasks include manipulation type (and degree)
34 classification [10], manipulation parameter estimation [75], blending boundary detection [41], spatial
35 forgery localization [28], face reconstruction [8], and face segmentation [55].

36 The prevailing multitask learning paradigm for face forgery detection follows a discriminative
37 approach, predicting multiple target outputs, one for each task, directly from the input face image.
38 Such a paradigm suffers from two main drawbacks. First, semantic relationships across tasks are
39 overlooked, which weakens knowledge transfer. For example, irrelevant information (*e.g.*, every
40 detail of the face image in face reconstruction [8]) may be transferred across tasks. Second, extensive
41 human expertise should be involved, when determining task-agnostic (and task-specific) model
42 parameters and the loss weightings.

43 In this paper, we explore multitask learning for face forgery detection from a joint embedding
44 perspective [38]. In the joint embedding architecture, both the input and the target output are encoded
45 into latent representations in the shared feature space such that the irrelevant information can be
46 discarded from feature encoding. More importantly, the semantic closeness between two tasks can
47 be naturally modeled as the distance in the learned feature space, which is subsequently end-to-
48 end optimized to facilitate knowledge transfer across multiple tasks. Meanwhile, joint embedding
49 gives us a great opportunity to automate multitask learning in terms of allocating model capacity
50 (*i.e.*, specifying task-agnostic and task-specific model parameters). In the context of face forgery
51 detection, the parameters of the face image encoder are shared across all tasks, whose capacity is
52 dynamically adjusted through end-to-end optimization. In addition, the multitask loss weightings can
53 be automatically computed in either theoretical [45, 65] or empirical [13, 36, 47] ways.

54 More concretely, we first introduce three coarse-to-fine face forgery detection tasks based on face
55 attributes at different semantic levels. Leveraging
56 the recent advances in vision-language correspond-
57 ence as joint embedding [61], we encode the binary
58 labels of the three tasks via textual prompts, and
59 thus the semantic dependencies among tasks can be
60 represented with the textual embeddings in the rep-
61 resentation space. Fig. 1 shows an example, in which
62 we describe a fake face image with a set of coarse-
63 to-fine textual descriptions: 1) “A photo of a fake
64 face,” 2) “A photo of a face with the global attribute
65 of expression altered,” and 3) “A photo of a face with
66 the local attribute of mouth altered.” By jointly embedding the face image and all its associated
67 textual prompts through a popular vision-language model - CLIP [61], face forgery detection can
68 then be performed by maximizing the vision-language correspondence.



- (1) A photo of a fake face
- (2) A photo of a face with the global attribute of expression altered
- (3) A photo of a face with the local attribute of mouth altered

Figure 1: Illustration of a fake face image with its textural descriptions of three coarse-to-fine face forgery detection tasks at different semantic levels.

69 **Our contributions** are threefold. First, we formulate multitask face forgery detection from a joint
70 embedding perspective. Second, we define a set of coarse-to-fine face forgery detection tasks with
71 corresponding textural templates to describe (fake) face images. Compared to previous multitask
72 learning schemes, our instantiation gives rise to a more interpretable face forgery detector. Third,
73 we conduct extensive experiments on five popular face forgery detection datasets, and show that our
74 method performs favorably against state-of-the-art (SOTA) detectors in terms of generalization to
75 novel face manipulations.

76 2 Related Work

77 In this section, we briefly review the literature on face forgery detection, multitask learning, and joint
78 embedding architectures.

79 2.1 Face Forgery Detection

80 Many face forgery detection methods usually explore the specific clues to detect the forgery inspired
81 by the traditional photo forensics [15, 32, 33, 35, 46, 50, 51, 56], in which they detect eye blink-
82 ing [42], head pose [77], pupil shape [24], lipreading [26], statistical anomalies [43, 60, 66, 81],
83 corneal specularities [29], and idiosyncratic behavioral patterns of a well-known person [3]. In
84 recent years, there is a growing consensus of exploiting multitask learning on face forgery detec-
85 tion [8, 10, 41, 55, 81]. Besides the main face forgery classification task, these methods include
86 auxiliary tasks to get performance improvement by knowledge transfer across tasks, such as manipula-
87 tion type (and degree) classification [10], manipulation parameter estimation [75], blending boundary
88 detection [41], spatial forgery localization [28], face reconstruction [8], and face segmentation [55].

89 With the development of deep learning, some advanced networks are employed to facilitate the face
90 forgery detection based on multiple tasks, such as two-stream CNN [82], self-attention model [80],
91 and vision transformers [19]. Additionally, more advanced training strategies are also utilized to
92 enhance the forgery detectors, including adversarial learning [10], reconstruction learning [8], and
93 meta learning [11]. However, the previous learning paradigm and human intervention are sub-optimal
94 for multitask learning on face forgery detection. In this paper, we explore an automated multitask
95 learning method for face forgery detection from the joint embedding perspective, where multiple
96 tasks are encoded into the language prompts, and vision-language correspondence is transferred
97 across tasks as the primary knowledge.

98 2.2 Multitask Learning

99 Multitask learning aims to jointly learn multiple related tasks to improve the generalization perfor-
100 mance of all tasks by leveraging the knowledge contained in each [79]. Two main groups are model
101 parameter sharing and loss weighting. The former involves both manual specifications of shared pa-
102 rameters [4, 22, 37, 54] and learning to determine parameters for specific tasks [52, 64, 68, 74]. Loss
103 weighting is typically divided as follows: Pareto Optimization (PO) methods and weight adaption
104 methods. PO methods formulate multitask learning as a multi-objective optimization [45, 65], and
105 find a Pareto stationary solution for the optimal loss weighting. Weight adoption methods adaptively
106 adjust the loss weights during training based on pre-defined heuristics, such as uncertainty [36],
107 gradient normalization [13], and loss descending rate [47]. In this paper, we consider multitask
108 learning from the joint embedding perspective, in which the semantic closeness between tasks can be
109 manifested as the distance in the learned feature space. Moreover, we assume all parameters in the
110 image encoder are shared, whose capacity is dynamically allocated to each task during end-to-end
111 optimization. We also adopt the method in [47] for dynamic loss weighting.

112 2.3 Joint Embedding Architectures

113 Joint embedding architectures (JEA) [38] aim at learning to output similar embeddings for compatible
114 inputs, x and y , and dissimilar embeddings for incompatible inputs, which is different from the
115 discriminative approaches that predict y directly from x . Becker *et al.* [6] propose the first JEA for
116 maximizing mutual information between representations from two views of the same scene. Later on,
117 Bromley *et al.* [7] propose a contrastive method of JEA for signatures verification. After a long hiatus,
118 JEA has been re-explored in face verification [14] and recognition [69], dimensionality reduction [25],
119 and video feature learning [70]. With the emergence of self-supervised learning, the use of JEA has
120 explored in recent years with methods training on contrastively (*e.g.*, PIRL [53], MoCo [27], and
121 SimCLR [12]) or non-contrastively (*e.g.*, BYOL [23], Barlow Twins [78], and I-JEPA [5]). More
122 recently, the emerging vision-language foundation models [30, 61] can also be grouped into JEA,
123 in which two separate encoders encode the compatible visual (*i.e.*, x) and textual (*i.e.*, y) inputs
124 into similar embeddings and contrast incompatible visual and textual embeddings. In this paper, we
125 use CLIP [61], a joint vision-language model pretrained on massive image-text pairs, to implement
126 the JEA to aid DeepFake detection by vision-language correspondence in the embedding space.
127 Moreover, we end-to-end fine-tune the CLIP in the context of automated multitask learning.

128 3 Method

129 In this section, we present multitask learning for face forgery detection using a joint embedding
130 approach, including preliminaries of the problem formulation, language prompts over multiple tasks,
131 and specifications of loss functions. The main joint embedding framework for face forgery detection
132 is shown in Fig. 2.

133 3.1 Preliminaries

134 Given a face image $x \in \mathbb{R}^N$, a face forgery detector $f_\theta : \mathbb{R}^N \mapsto \mathbb{R}$ aims to predict a binary label
135 y for the authenticity of x , *i.e.*, 0 as the real or 1 as the fake. Considering that existing forged face
136 images are mainly generated by modifying face components/attributes, we include two other related
137 tasks - global face manipulation detection and local face manipulation detection. We consider three

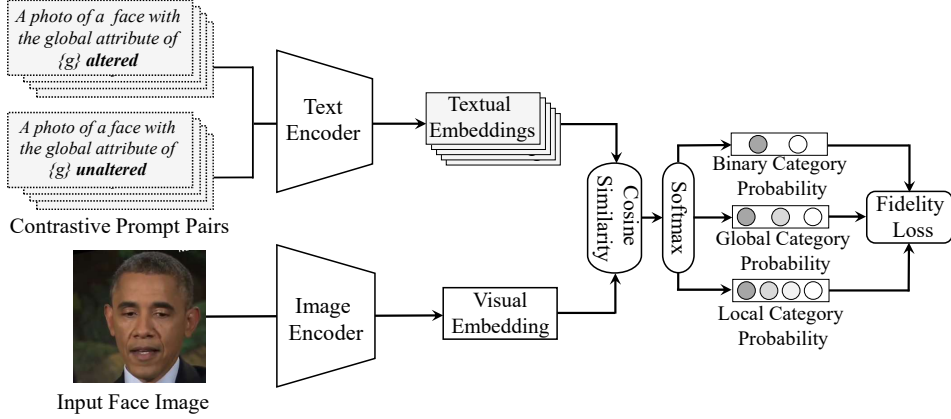


Figure 2: Proposed joint embedding paradigm for multitask face forgery detection.

138 face attributes (*i.e.*, expression, identity, and physical consistency¹) for global face manipulations, and
 139 four face attributes (*i.e.*, eye, illumination, mouth, and nose) for local face manipulations. Notably, a
 140 face image may contain multiple attribute labels.

141 3.2 Multitask Language Prompts

142 For each face attribute label from multiple tasks, we encode the ground-truth labels via language
 143 prompts. In specific, we design textual templates as follows. 1) **binary level**: *a photo of a {c} face*,
 144 where $c \in \mathcal{C} = \{\text{real, fake}\}$; 2) **global-attribute level**: *A photo of a face with the global attribute of*
 145 *{g} altered*, where $g \in \mathcal{G} = \{\text{expression, identity, physical consistency}\}$; and 3) **local-attribute level**:
 146 *A photo of a face with the local attribute of {l} altered*, where $l \in \mathcal{L} = \{\text{eye, illumination, mouth,
 147 *nose}\}*$. Inspired by contrastive methods [27, 53] in the joint embedding architecture, we also introduce
 148 contrastive language prompts, which are opposite in meaning to the original textual templates. Thus,
 149 we can have a contrastive prompts pair for each attribute label, as follows: **global-attribute level**:
 150 *{(1) A photo of a face with the global attribute of {g} altered, (2) A photo of a face with the global*
 151 *attribute of {g} unaltered}*; **local-attribute level**: *{(1) A photo of a face with the local attribute of*
 152 *{l} altered, (2) A photo of a face with the local attribute of {l} unaltered}*. Notably, the binary level
 153 prompts naturally have the property of contrastive prompt pairing. In this way, multiple tasks are
 154 encoded into a text corpus \mathcal{T} , where each language prompt represents a ground-truth label y of the
 155 corresponding task, and their semantic closeness can be learned through joint embedding.

156 3.3 Multitask Learning via Joint Embedding

157 **Joint Embedding Formulation.** Given the input face image \mathbf{x} and the set of possible outputs \mathcal{Y} ,
 158 we predict the output by minimizing an energy-based model [39], *i.e.*, $\hat{y} = \arg \min_{y \in \mathcal{Y}} E(\mathbf{x}, y)$, in
 159 the joint embedding architecture. In this paper, we construct E by two encoders: one image encoder
 160 $\mathbf{f}_\phi : \mathbb{R}^N \mapsto \mathbb{R}^K$ for encoding the face image and one text encoder $\mathbf{g}_\varphi : \mathcal{T} \mapsto \mathbb{R}^K$ for encoding the
 161 language prompts, parameterized by ϕ and φ , respectively.

162 The ideal energy landscape of joint embedding satisfies that the energy is low for similar embeddings
 163 of compatible inputs, while energy is high for dissimilar embeddings [39]. Thus, we calculate
 164 the probability of similarity $\hat{p}(\cdot|\mathbf{x})$ between the visual embedding and textual embeddings for the
 165 following optimization. Let $\mathbf{u} \in \mathbb{R}^K$ be the visual embedding, and let $\mathbf{v} \in \mathbb{R}^K$ and $\bar{\mathbf{v}} \in \mathbb{R}^K$ be the
 166 textual embeddings from the two prompts opposing in meaning, we then estimate $\hat{p}(\cdot|\mathbf{x})$ as

$$\hat{p}(\cdot|\mathbf{x}) = \frac{1}{1 + e^{-(s-\bar{s})}}, \quad (1)$$

167 where

$$s = \frac{\langle \mathbf{u}, \mathbf{v} \rangle}{\|\mathbf{u}\| \|\mathbf{v}\|} \quad \text{and} \quad \bar{s} = \frac{\langle \mathbf{u}, \bar{\mathbf{v}} \rangle}{\|\mathbf{u}\| \|\bar{\mathbf{v}}\|}. \quad (2)$$

¹We refer the interested readers to the Appendix for the detailed explanations.

168 $\langle \cdot, \cdot \rangle$ denotes the inner product and $\| \cdot \|$ represents the ℓ_2 -norm. The probability $\hat{p}(\cdot|\mathbf{x})$ is the
 169 abbreviation of $\hat{p}(c|\mathbf{x})$, $\hat{p}(g|\mathbf{x})$, and $\hat{p}(l|\mathbf{x})$ according to a specific task, and a larger probability
 170 indicates a closer match to the corresponding semantic meaning of v .

171 **Losses for Multitask Learning.** We use the statistical distance measure in the form of fidelity
 172 loss [73] to calculate the losses for multitask learning. Given the predicted category probability
 173 $\hat{p}(c|\mathbf{x})$, we design the loss at the **binary level** as

$$\ell_1(\mathbf{x}; \boldsymbol{\theta}) = 1 - \sqrt{p(c|\mathbf{x})\hat{p}(c|\mathbf{x})} - \sqrt{(1-p(c|\mathbf{x}))(1-\hat{p}(c|\mathbf{x}))}, \quad (3)$$

174 where $\boldsymbol{\theta} = \{\phi, \varphi\}$ indicates the learnable parameters in image and language encoders, and $p(c|\mathbf{x}) = 1$
 175 if \mathbf{x} belongs to the c category or otherwise we have $p(c|\mathbf{x}) = 0$. In our setting, a face image can be
 176 assigned with labels regarding one or more global face attribute manipulations, which forms a typical
 177 multi-label classification problem. Therefore, the averaged loss at the **global-attribute level** can be
 178 defined as follows,

$$\ell_2(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \left(1 - \sqrt{p(g|\mathbf{x})\hat{p}(g|\mathbf{x})} - \sqrt{(1-p(g|\mathbf{x}))(1-\hat{p}(g|\mathbf{x}))} \right), \quad (4)$$

179 where $p(g|\mathbf{x}) = 1$ if \mathbf{x} belongs to the g category, otherwise we have $p(g|\mathbf{x}) = 0$. Since the
 180 manipulations over different local face attributes may appear in one face image, we also consider it
 181 as a multi-label classification task, and the loss at the **local-attribute level** is:

$$\ell_3(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{|\mathcal{L}|} \sum_{l \in \mathcal{L}} \left(1 - \sqrt{p(l|\mathbf{x})\hat{p}(l|\mathbf{x})} - \sqrt{(1-p(l|\mathbf{x}))(1-\hat{p}(l|\mathbf{x}))} \right), \quad (5)$$

182 where $p(l|\mathbf{x}) = 1$ if \mathbf{x} belongs to the l category.

183 Given a minibatch of training data \mathcal{B} at the t -th iteration, we evaluate the overall loss function via the
 184 weighted sum of the individual losses in different levels as follows,

$$\ell(\mathcal{B}, t; \boldsymbol{\theta}) = \frac{1}{|\mathcal{B}|} \sum_{\mathbf{x} \in \mathcal{B}} (\lambda_1(t)\ell_1(\mathbf{x}; \boldsymbol{\theta}) + \lambda_2(t)\ell_2(\mathbf{x}; \boldsymbol{\theta}) + \lambda_3(t)\ell_3(\mathbf{x}; \boldsymbol{\theta})). \quad (6)$$

185 Here, the weighting vector $\boldsymbol{\lambda}(t) = [\lambda_1(t), \lambda_2(t), \lambda_3(t)]^\top$ at the t -th iteration is automatically com-
 186 puted according to the relative descending rate [47]:

$$\lambda_i(t) = \frac{3 \exp(w_i(t-1)/\tau)}{\sum_{j=1}^3 \exp(w_j(t-1)/\tau)}, \text{ where } w_i(t-1) = \frac{\ell_i(t-1)}{\ell_i(t-2)}, \quad (7)$$

187 and τ is a fixed temperature parameter.

188 4 Experiments

189 4.1 Experimental Setup

190 **Datasets.** We adopt the widely used FF++ [63] dataset for training. It contains 1,000 real videos,
 191 among which 720 and 140 are used for training and validation, respectively, and the remaining
 192 140 are reserved for testing. All videos are manipulated by four face forgery methods, including
 193 Deepfakes [1], Face2Face [72], FaceSwap [2], and NeuralTextures [71], with three compression levels,
 194 *i.e.*, no compression (denoted as Raw), slight compression with quantization parameter QP = 23
 195 (denoted as C23), and severe compression with QP = 40 (denoted as C40). Following [10, 11, 26],
 196 C23 version is adopted by default in our experiments. We evaluate the generalizability of the
 197 proposed method on four popular DeepFake benchmarks, including FaceShifter (FSh) [40], Celeb-DF
 198 (CDF) [44], DeeperForensics-1.0 (DF-1.0) [31], and DeepFake Detection Challenge (DFDC) [18].

199 **Implementation Details.** To facilitate the multitask learning via joint embedding paradigm, we need
 200 face images associated with the proposed textual templates. In this paper, we adopt FF++ [63] to
 201 enrich the training data. Following the general generation procedures (*i.e.*, detecting face and then
 202 blending two faces according to the region-of-interest mask) in [10, 41], we focus on supplementing
 203 the tampering of “expression” on “eye” and individual face attribute that is linked to “physical
 204 consistency”, *i.e.*, “eye”, “illumination”, “mouth”, and “nose”. Face attribute manipulations associated
 205 with other textual prompts are already included in FF++.

206 As for face pre-processing, we use RetinaFace [17] to detect faces and save the aligned face images
 207 as input with a size of 317×317 . As in [63], we only extract the largest face and use an enlarged
 208 crop, $1.3 \times$ the tight crop produced by the face detector.

209 As for the training, we use CLIP [61] to implement the joint embedding architecture, where we
 210 adopt ViT-B/32 [20] as the visual encoder and GPT-2 [62] with a base size of 63M-parameter as the
 211 text encoder. We then train the model by minimizing the loss using AdamW [49] with a decoupled
 212 weight decay of 1×10^{-3} . The initial learning rate is set to 1×10^{-7} , which changes following a
 213 cosine annealing schedule [48]. The model is optimized for 36 epochs with mini-batches of 32. Data
 214 augmentation strategy is also applied during training, which is a common trick in the face forgery
 215 detection [41, 76, 80], and details can be found in Sec. 4.3. A single NVIDIA RTX 3090 GPU is
 216 used during training.

217 4.2 Comparison with SOTA Methods

218 We compare our method with the several SOTA methods, including Face X-ray [41], PCL [81],
 219 MADD [80], LipForensics [26], RECCE [8], SBI [66], ICT [19], SLADD [10], and OST [11],
 220 to demonstrate its superiority. The test performance on five datasets are listed in Table 1.
 221 Table 1 shows that

222 many methods
 223 do not perform
 224 satisfactorily on
 225 face forgery de-
 226 tection, while the
 227 proposed method
 228 outperforms all
 229 the recent SOTA,
 230 achieving 92.33%
 231 of AUC aver-
 232 aged from five
 233 test datasets and
 234 surpassing the
 235 second best, *i.e.*,
 236 LipForensics, by
 237 2.79% in the term
 238 of Mean AUC over
 239 datasets including
 240 FF++ [63]. For

Table 1: **Comparison results with the SOTA.** All models are developed using the training set of FF++ (or its augmented versions) and tested on the test set of FF++ and other four independent datasets. The evaluation metric we adopt is AUC (%). In the last column are the mean AUC numbers over datasets including / excluding the FF++ test set to emphasize cross-dataset generalization performance. The best results are highlighted in bold.

Method	FF++	CDF	FSh	DF-1.0	DFDC	Mean AUC
Face X-ray [41]	98.37	80.43	92.80	86.80	65.50	84.78 / 81.38
PCL [81]	99.11	81.80	–	99.40	67.50	86.95 / 82.90
MADD [80]	98.97	77.44	97.17	66.58	67.94	81.62 / 77.28
LipForensics [26]	99.90	82.40	97.10	97.60	73.50	89.54 / 87.65
RECCE [8]	99.32	68.71	70.58	74.10	69.06	76.35 / 70.61
SBI [66]	99.64	93.18	97.40	77.70	72.42	88.07 / 85.18
ICT [19]	90.22	85.71	95.97	93.57	76.74	88.44 / 88.00
SLADD [10]	98.40	79.70	–	77.80	76.05	82.99 / 77.85
OST [11]	98.20	74.80	–	93.08	77.73	84.95 / 81.87
Ours	98.49	89.02	98.68	93.38	82.06	92.33 / 90.79

241 cross-dataset generalizability comparison, the proposed method also surpasses the second best (*i.e.*,
 242 ICT) and third best (*i.e.*, LipForensics) by 2.79% and 3.14%, respectively. In addition, we also
 243 have several interesting observations. **First**, all the methods can achieve saturated performance in
 244 FF++ [63], while underperform in the rest datasets, such as CDF [44] and DFDC [18]. This suggests
 245 that the forgery cues in FF++ are easier to spot and overfit by these forgery detectors. **Second**, SBI
 246 reports a very high AUC of 93.18% on CDF, while performing unsatisfactorily on DF-1.0 [31] and
 247 DFDC. Similar results are also demonstrated by PCL, which exhibits an exceedingly high AUC
 248 of 99.40% on DF-1.0 but underperforms in DFDC. This may arise due to the overfitting on the
 249 low-level features, such as statistical inconsistency (*e.g.*, landmark and color mismatch). **Third**, all
 250 methods obtain relatively low scores on DFDC, which we attribute to the domain shift caused by
 251 significantly different filming conditions. However, our method achieves a relative satisfactory result
 252 with a score of 82.06%, surpassing the second best by 4.33%. In summary, the remarkable results
 253 validate the effectiveness and superiority of the proposed joint-embedding-based multitask learning
 254 for DeepFake detection.

255 4.3 Robustness Analysis

256 In this subsection, we study the robustness performance of the proposed method. Following [31], we
 257 consider four popular perturbations (*i.e.*, Patch Substitution (Patch-Sub), additive white Gaussian
 258 Noise contamination (Noise), Gaussian Blurring (Blur), and pixelation), and only four severity levels

Table 2: **Robustness results to low-level image perturbations**, including patch substitution (Patch-Sub), Gaussian noise contamination (Noise), Gaussian blurring (Blur), and pixelation. We constrain the robustness evaluation on the perturbation levels that do not noticeably distort the main face semantics.

Method	Clean AUC	Patch-Sub	Noise	Blur	Pixelation	Mean AUC	Drop Rate
Face X-ray [41]	98.37	97.72	51.13	88.98	92.33	82.54	-16.09%
CNND [76]	99.56	96.25	57.25	92.61	90.10	84.05	-15.58%
LipForensics [26]	99.90	88.63	80.00	96.62	96.63	90.47	-9.44%
Ours (w/o Aug)	98.66	92.47	73.12	55.20	57.17	69.49	-29.57%
Ours	98.49	97.65	82.85	87.31	90.70	89.63	-8.99%

259 (*i.e.*, from level 1 to level 4) are considered in the experiments². Two different models are evaluated
 260 in this section, *i.e.*, our model training without data augmentation (denoted as Ours (w/o Aug)) and
 261 our model training with data augmentation strategy (denoted as Ours). In specific, when training with
 262 data augmentation strategy, each training data
 263 is augmented with a probability of 0.3 by one
 264 randomly chosen perturbation during training,
 265 in which severity level is randomly applied
 266 at level 1 or 2.

267 To begin, we first evaluate the robustness for
 268 the model without data augmentation. We
 269 find that the CLIP-based model is sensitive to
 270 the perturbations to images, which we argue
 271 that the vision-language correspondence is
 272 corrupted by perturbations. We then evaluate
 273 the model training with data augmentation.
 274 In Table 2, we find that training with a slight
 275 data augmentation can alleviate the model
 276 sensitivity to the perturbations, and achieve
 277 a satisfactory performance on average. More-
 278 over, the model of Ours also maintains a sat-
 279 isfactory performance on pixelation and Blur.
 280 It is noteworthy that CNND [76] and Face
 281 X-ray [41] also augment their training data
 282 by compression and blurring during training,
 283 thus leading to good robustness to perturba-
 284 tions of pixelation and Blur. Fig. 3 demon-
 285 strates the effect of increasing the severity for
 286 each perturbation, where we compare with
 287 Xception [63], CNND, PatchForensics [9],
 288 Face X-ray, and LipForensics [26]. It can be
 289 observed that the proposed method maintains a
 290 good performance against the perturbations by
 Patch-Sub and Noise, while other methods suffer
 from the Noise, and LipForensics also suffers
 from the Patch-Sub.

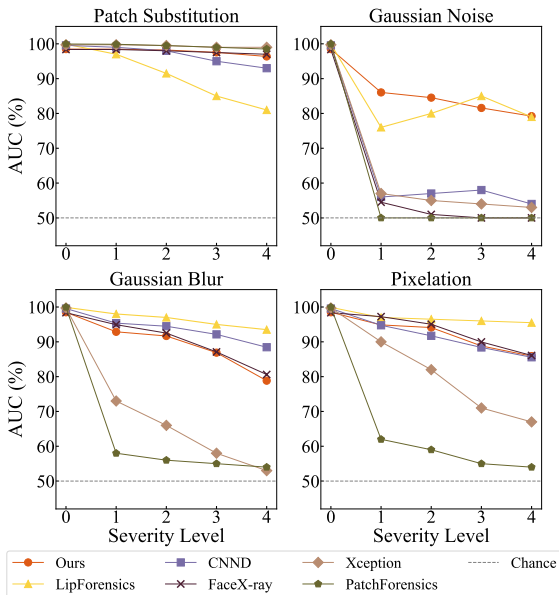


Figure 3: Robustness results in terms of AUC. Models are trained on the train set of FF++ and tested on perturbed test sets. Zoom in for clearer comparison.

291 4.4 Ablation Studies

292 **Joint Embedding Framework.** We conducted a series of ablations to verify the instantiated joint
 293 embedding framework by CLIP [61]. We first (1) evaluate the pretrained CLIP, and then (2) fine-tune
 294 it with the frozen text encoder on FF++ [63]. The following ablations adopt the same training
 295 procedure, while differing in two alternatives: (3) using equal task weights for multiple tasks instead
 296 of dynamic loss weighting; (4) training without the contrastive prompt pairs, *i.e.*, no contrastive
 297 textual descriptions are used during training. From Table 3, we can observe that freezing language
 298 encoder negatively affects the generalization performance, which we believe is because forgery-
 299 related concepts have not been sufficiently captured during the pretraining stage of CLIP. We also
 300 find that utilizing contrastive prompts can improve generalization, further indicating the contrasting

²The perturbations on severity level 5 often make the face semantically unrecognized, leading meaningless to detect its authenticity.

operation can benefit the joint embedding methods [12, 27]. Moreover, including the dynamic loss weighting scheme is advantageous as it not only yields a slight improvement compared to using equal task weights but also frees us from the burdensome task of hyper-parameter tuning.

Textual Templates. In this subsection, we investigate how the textual template design affects the model performance. We try three different alternatives from single task to three tasks: (5) binary-level text templates, *i.e.*, single task formulation only considering the label of real or fake; (6) two-level separate text templates, *i.e.*, two-level-task formulation, where we consider the separate templates describing the overall authenticity and global face attributes; and (7) the joint text templates putting together labels from three tasks, *e.g.*, “A photo of a {fake} face with the global attribute of {expression} and the local attribute of {mouth} are altered”. The joint probability over multiple tasks can be computed from the similarities between the image embedding and all candidate tex-

Table 3: **Ablation Studies.** Baseline denotes the single-task formulation w/o contrastive textual pairing nor data augmentation, optimized for the BCE loss.

Model Variant	CDF	FSh	DF-1.0	DFDC	Mean AUC
(1) Pretrained CLIP	65.38	51.04	53.38	55.56	56.34
(2) Frozen g_φ	90.56	98.92	91.22	80.19	90.22
(3) Equal Weights	88.32	98.77	92.93	82.27	90.57
(4) w/o Contrastive Pair	87.89	98.34	93.30	81.27	90.20
(5) Binary Templates	85.03	98.42	93.33	81.58	89.59
(6) Two-Levels	87.57	98.47	93.74	80.81	90.15
(7) Joint Templates	88.05	98.42	94.21	81.31	90.50
(8) ViT-B/16	88.13	99.62	93.30	82.30	90.84
(9) ViT-L/14	90.78	99.95	98.60	86.22	93.89
(10) BCE Loss	86.45	98.35	93.40	80.81	89.75
(11) Probabilistic Loss	87.81	98.41	91.55	81.18	89.74
Ours (Baseline)	71.63	98.19	89.94	74.02	83.44
Ours (w/o Aug)	85.53	98.82	93.95	80.41	89.68
Ours (Default)	89.02	98.68	93.38	82.06	90.79

tual embeddings. Then, we marginalize the joint distribution to obtain the marginal probability for each task. From Table 3, we can observe that the performance of the model using joint templates is inferior to that of the model using separate templates (*i.e.*, Ours (Default)), indicating that separate templates for each task are more conducive for learning the semantic closeness between two face forgery detection tasks in joint embedding. On the other hand, less tasks (*i.e.*, single task and two tasks) result in the inferior performance. Notably, benefiting from the joint embedding, the model using binary templates also achieves comparable results on generalization, though it only classifies the overall authenticity of the face.

Encoder Architecture. In this subsection, we investigate other visual encoders with different settings and model sizes. In specific, we choose (8) ViT-B/16 [20] and (9) ViT-L/14 [20]. As shown in Table 3, two alternative ViT-based architectures achieve better results on generalization. However, the larger model will result in both computationally more expensive and time-consuming.

Multitask Objective. In this subsection, we study how different optimization objectives affect the performance. As a reference, we first replace the fidelity loss functions with (10) binary cross entropy loss (BCE Loss). We also adopt the (11) hierarchical probabilistic loss [16] to jointly formulate multi-level classification tasks under a hierarchical label semantic graph. The relative similarity score (*i.e.*, $s - \bar{s}$), as a raw score, for each node in the label hierarchy, will be converted into marginal probabilities for loss computation. From Table 3, we observed that the proposed method outperforms the variant trained with BCE loss, thus providing evidence for the effectiveness of the designed fidelity losses. Furthermore, Table 3 shows that fidelity loss yields better performance than the hierarchical probabilistic loss, suggesting that implicitly learning the semantic dependencies may be better than explicitly encoding the prior knowledge in the label hierarchy graph in advance.

4.5 Discussion: Vision-Language Correspondence

Human-Understandable Interpretation. The proposed joint embedding approach enjoys the vision-language correspondence, which naturally provides model interpretations by providing human-understandable explanations. Fig. 4 shows some examples of FF++ [63], in which Deepfakes [1] indicate the identity swap, leading all local parts of the face are fake; and NeuralTextures [71] modify the expression in the mouth part. Take an example of NeuralTextures, the texts with a probability over 50% include “fake”, “expression”, and “mouth”. Hence, we consider this face image to be fake

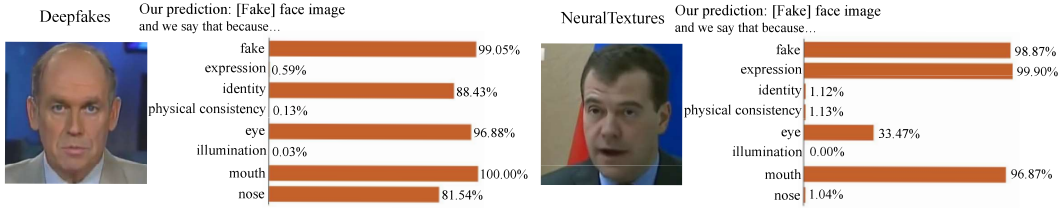
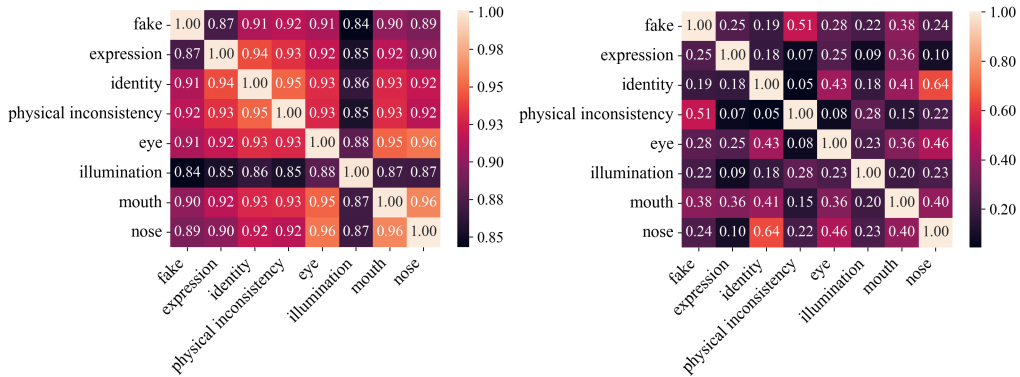


Figure 4: Bar charts of the similarity scores between the visual image and the textual descriptions a form of human-understandable explanations.

356 because the model’s prediction relies on the following three textual prompts: “a photo of a fake face”,
 357 “a photo of a face with the global attribute of expression altered”, and “a photo of a face with the
 358 local attribute of mouth altered”. More examples can be found in Appendix.

359 **Semantic Closeness across Tasks.** We show the semantic closeness across tasks by a correlation
 360 matrix in Fig. 5, in which each entry is represented by the cosine similarity between two textual
 361 embeddings from the language prompts depicting the specific tasks. From Fig. 5, we can observe
 362 that the text encoder of the pretrained CLIP has not sufficiently captured the semantic closeness
 363 across tasks and treats most tasks equally, further verifying the results of the variant with frozen text
 364 encoder in Table 3. After joint embedding learning on the forged faces, the semantic closeness across
 365 tasks can be sufficiently learned, e.g., the concept of “identity” forgery is more related to the “nose”,
 “mouth”, and “eye”, thus improving the performance of multitask learning for face forgery detection.



(a) By Text Encoder of Pretrained CLIP

(b) By Text Encoder of Fine-tuned CLIP

Figure 5: Illustration of semantic closeness across tasks before and after fine-tuning.

367 5 Conclusion and Limitations

368 **Conclusion.** In this paper, we consider multitask learning for face forgery detection from the joint
 369 embedding perspective. We have designed a set of coarse-to-fine language prompts to represent
 370 multiple tasks for face forgery detection. We then take an automated multitask learning scheme to train
 371 two encoders to joint embed visual face images and textual descriptions. Thus, semantic closeness
 372 across tasks is manifested as the distance in the learned feature space, thus improving multitask
 373 learning. From extensive experiments, vision-language correspondence after joint embedding shows
 374 great promise to support better face forgery detection by maximizing the feature similarity between the
 375 face image and candidate textual prompts, verifying the effectiveness and superiority of the proposed
 376 method. Moreover, the joint embedding scheme also renders some degree of model interpretation in
 377 a human-friendly way.

378 **Limitations.** The proposed method relies on the assumption that the forged faces are generated with
 379 the blending operation [41]. Thus, it may perform unsatisfactorily when fake face images are totally
 380 synthesized by GAN- or diffusion-model-based methods. Additionally, our model is image-based,
 381 though it can handle video-based DeepFake by sampling frames for prediction, it may fail when
 382 encountering the fake video manipulated by only lowering the frame rate [57].

- 384 [1] Deepfakes. <https://github.com/deepfakes/faceswap>.
- 385 [2] FaceSwap. <https://github.com/MarekKowalski/FaceSwap>.
- 386 [3] S. Agarwal, H. Farid, Y. Gu, M. He, K. Nagano, and H. Li. Protecting world leaders against deep fakes. In *CVPRW*, pages 38–45, 2019.
- 387
- 388 [4] A. Argyriou, T. Evgeniou, and M. Pontil. Multi-task feature learning. In *NIPS*, pages 1–13, 2006.
- 389 [5] M. Assran, Q. Duval, I. Misra, P. Bojanowski, P. Vincent, M. Rabbat, Y. LeCun, and N. Ballas. Self-supervised learning from images with a joint-embedding predictive architecture. *arXiv preprint arXiv:2301.08243*, 2023.
- 390
- 391
- 392 [6] S. Becker and G. E. Hinton. Self-organizing neural network that discovers surfaces in random-dot stereograms. *Nature*, 355(6356):161–163, 1992.
- 393
- 394 [7] J. Bromley, I. Guyon, Y. LeCun, E. Säckinger, and R. Shah. Signature verification using a “Siamese” time delay neural network. In *NIPS*, pages 737–744, 1993.
- 395
- 396 [8] J. Cao, C. Ma, T. Yao, S. Chen, S. Ding, and X. Yang. End-to-end reconstruction-classification learning for face forgery detection. In *CVPR*, pages 4113–4122, 2022.
- 397
- 398 [9] L. Chai, D. Bau, S.-N. Lim, and P. Isola. What makes fake images detectable? Understanding properties that generalize. In *ECCV*, pages 103–120, 2020.
- 399
- 400 [10] L. Chen, Y. Zhang, Y. Song, L. Liu, and J. Wang. Self-supervised learning of adversarial example: Towards good generalizations for DeepFake detection. In *CVPR*, pages 18710–18719, 2022.
- 401
- 402 [11] L. Chen, Y. Zhang, Y. Song, J. Wang, and L. Liu. OST: Improving generalization of DeepFake detection via one-shot test-time training. In *NIPS*, pages 1–14, 2022.
- 403
- 404 [12] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, pages 1597–1607, 2020.
- 405
- 406 [13] Z. Chen, V. Badrinarayanan, C.-Y. Lee, and A. Rabinovich. GradNorm: Gradient normalization for adaptive loss balancing in deep multitask networks. In *ICML*, pages 794–803, 2018.
- 407
- 408 [14] S. Chopra, R. Hadsell, and Y. LeCun. Learning a similarity metric discriminatively, with application to face verification. In *CVPR*, pages 539–546, 2005.
- 409
- 410 [15] V. Conotter, J. F. O’Brien, and H. Farid. Exposing digital forgeries in ballistic motion. *IEEE TIFS*, 7(1):283–296, 2011.
- 411
- 412 [16] J. Deng, N. Ding, Y. Jia, A. Frome, K. Murphy, S. Bengio, Y. Li, H. Neven, and H. Adam. Large-scale object classification using label relation graphs. In *ECCV*, pages 48–64, 2014.
- 413
- 414 [17] J. Deng, J. Guo, E. Ververas, I. Kotsia, and S. Zafeiriou. RetinaFace: Single-shot multi-level face localisation in the wild. In *CVPR*, pages 5203–5212, 2020.
- 415
- 416 [18] B. Dolhansky, J. Bitton, B. Pflaum, J. Lu, R. Howes, M. Wang, and C. C. Ferrer. The DeepFake detection challenge (DFDC) dataset. *arXiv preprint arXiv:2006.07397*, 2020.
- 417
- 418 [19] X. Dong, J. Bao, D. Chen, T. Zhang, W. Zhang, N. Yu, D. Chen, F. Wen, and B. Guo. Protecting celebrities from DeepFake with identity consistency transformer. In *CVPR*, pages 9468–9478, 2022.
- 419
- 420 [20] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, pages 1–12, 2020.
- 421
- 422
- 423 [21] H. Farid. Image forgery detection: A survey. *IEEE SPM*, 26(2):16–25, 2009.
- 424
- 425 [22] Y. Gao, J. Ma, M. Zhao, W. Liu, and A. L. Yuille. NDDR-CNN: Layerwise feature fusing in multi-task cnns by neural discriminative dimensionality reduction. In *CVPR*, pages 3205–3214, 2019.
- 426
- 427 [23] J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. Richemond, E. Buchatskaya, C. Doersch, B. Avila Pires, Z. Guo, M. Gheshlaghi Azar, B. Piot, K. Kavukcuoglu, R. Munos, and M. Valko. Bootstrap your own latent: A new approach to self-supervised learning. In *NIPS*, pages 21271–21284, 2020.
- 428
- 429 [24] H. Guo, S. Hu, X. Wang, M.-C. Chang, and S. Lyu. Eyes tell all: Irregular pupil shapes reveal GAN-generated faces. In *ICASSP*, pages 2904–2908, 2022.
- 430
- 431 [25] R. Hadsell, S. Chopra, and Y. LeCun. Dimensionality reduction by learning an invariant mapping. In *CVPR*, pages 1735–1742, 2006.
- 432
- 433 [26] A. Haliassos, K. Vougioukas, S. Petridis, and M. Pantic. Lips don’t lie: A generalisable and robust approach to face forgery detection. In *CVPR*, pages 5039–5049, 2021.
- 434
- 435 [27] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, pages 9729–9738, 2020.
- 436
- 437 [28] Y. He, B. Gan, S. Chen, Y. Zhou, G. Yin, L. Song, L. Sheng, J. Shao, and Z. Liu. ForgeryNet: A versatile benchmark for comprehensive forgery analysis. In *CVPR*, pages 4360–4369, 2021.
- 438
- 439 [29] S. Hu, Y. Li, and S. Lyu. Exposing GAN-generated faces using inconsistent corneal specular highlights. In *ICASSP*, pages 2500–2504, 2021.
- 440
- 441 [30] C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. Le, Y.-H. Sung, Z. Li, and T. Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, pages 4904–4916, 2021.
- 442
- 443
- 444 [31] L. Jiang, R. Li, W. Wu, C. Qian, and C. C. Loy. DeeperForensics-1.0: A large-scale dataset for real-world face forgery detection. In *CVPR*, pages 2889–2898, 2020.
- 445

- 446 [32] M. K. Johnson and H. Farid. Exposing digital forgeries through chromatic aberration. In *ACM MM&Sec*,
447 pages 48–55, 2006.
- 448 [33] M. K. Johnson and H. Farid. Exposing digital forgeries in complex lighting environments. *IEEE TIFS*,
449 2(3):450–461, 2007.
- 450 [34] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila. Analyzing and improving the image
451 quality of StyleGAN. In *CVPR*, pages 8110–8119, 2020.
- 452 [35] E. Kee, J. F. O’Brien, and H. Farid. Exposing photo manipulation from shading and shadows. *ACM TOG*,
453 33(5):165:1–165:21, 2014.
- 454 [36] A. Kendall, Y. Gal, and R. Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry
455 and semantics. In *CVPR*, pages 7482–7491, 2018.
- 456 [37] I. Kokkinos. UberNet: Training a universal convolutional neural network for low-, mid-, and high-level
457 vision using diverse datasets and limited memory. In *CVPR*, pages 6129–6138, 2017.
- 458 [38] Y. LeCun. A path towards autonomous machine intelligence version 0.9. 2, 2022-06-27. 2022.
- 459 [39] Y. LeCun, S. Chopra, R. Hadsell, M. Ranzato, and F. Huang. A tutorial on energy-based learning.
460 *Predicting Structured Data*, 1(0):1–59, 2006.
- 461 [40] L. Li, J. Bao, H. Yang, D. Chen, and F. Wen. Advancing high fidelity identity swapping for forgery
462 detection. In *CVPR*, pages 5074–5083, 2020.
- 463 [41] L. Li, J. Bao, T. Zhang, H. Yang, D. Chen, F. Wen, and B. Guo. Face X-ray for more general face forgery
464 detection. In *CVPR*, pages 5001–5010, 2020.
- 465 [42] Y. Li, M.-C. Chang, and S. Lyu. In Ictu Oculi: Exposing AI created fake videos by detecting eye blinking.
466 In *WIFS*, pages 1–7, 2018.
- 467 [43] Y. Li and S. Lyu. Exposing DeepFake videos by detecting face warping artifacts. In *CVPRW*, pages 46–52,
468 2019.
- 469 [44] Y. Li, X. Yang, P. Sun, H. Qi, and S. Lyu. Celeb-DF: A large-scale challenging dataset for DeepFake
470 forensics. In *CVPR*, pages 3207–3216, 2020.
- 471 [45] X. Lin, H.-L. Zhen, Z. Li, Q.-F. Zhang, and S. Kwong. Pareto multi-task learning. In *NIPS*, pages
472 12037–12047, 2019.
- 473 [46] Z. Lin, R. Wang, X. Tang, and H.-Y. Shum. Detecting doctored images using camera response normality
474 and consistency. In *CVPR*, pages 1087–1092, 2005.
- 475 [47] S. Liu, E. Johns, and A. J. Davison. End-to-end multi-task learning with attention. In *CVPR*, pages
476 1871–1880, 2019.
- 477 [48] I. Loshchilov and F. Hutter. SGDR: Stochastic gradient descent with warm restarts. In *ICLR*, pages 1–13,
478 2017.
- 479 [49] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In *ICLR*, pages 1–10, 2019.
- 480 [50] S. Lyu. Estimating vignetting function from a single image for image authentication. In *ACM MM&Sec*,
481 pages 3–12, 2010.
- 482 [51] S. Lyu, X. Pan, and X. Zhang. Exposing region splicing forgeries with blind local noise estimation. *IJCV*,
483 110(2):202–221, 2014.
- 484 [52] A. Mallya, D. Davis, and S. Lazebnik. Piggyback: Adapting a single network to multiple tasks by learning
485 to mask weights. In *ECCV*, pages 72–88, 2018.
- 486 [53] I. Misra and L. v. d. Maaten. Self-supervised learning of pretext-invariant representations. In *CVPR*, pages
487 6707–6717, 2020.
- 488 [54] I. Misra, A. Shrivastava, A. Gupta, and M. Hebert. Cross-stitch networks for multi-task learning. In *CVPR*,
489 pages 3994–4003, 2016.
- 490 [55] H. H. Nguyen, F. Fang, J. Yamagishi, and I. Echizen. Multi-task learning for detecting and segmenting
491 manipulated facial images and videos. In *BTAS*, pages 1–8, 2019.
- 492 [56] J. F. O’Brien and H. Farid. Exposing photo manipulation with inconsistent reflections. *ACM TOG*,
493 31(1):4:1–4:11, 2012.
- 494 [57] D. O’Sullivan. Doctored videos shared to make Pelosi sound drunk viewed millions of times on social me-
495 dia. <https://edition.cnn.com/2019/05/23/politics/doctored-video-pelosi/index.html>,
496 2019. Date of access: May 12, 2023.
- 497 [58] A. C. Popescu and H. Farid. Exposing digital forgeries by detecting traces of resampling. *IEEE TSP*,
498 53(2):758–767, 2005.
- 499 [59] A. C. Popescu and H. Farid. Exposing digital forgeries in color filter array interpolated images. *IEEE TSP*,
500 53(10):3948–3959, 2005.
- 501 [60] Y. Qian, G. Yin, L. Sheng, Z. Chen, and J. Shao. Thinking in frequency: Face forgery detection by mining
502 frequency-aware clues. In *ECCV*, pages 86–103, 2020.
- 503 [61] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin,
504 J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from natural language
505 supervision. In *ICML*, pages 8748–8763, 2021.
- 506 [62] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. Language models are unsupervised
507 multitask learners. 2019.
- 508 [63] A. Rossler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Nießner. FaceForensics++: Learning to
509 detect manipulated facial images. In *ICCV*, pages 1–11, 2019.

- 510 [64] S. Ruder, J. Bingel, I. Augenstein, and A. Søgaard. Latent multi-task architecture learning. In *AAAI*, pages
511 4822–4829, 2019.
- 512 [65] O. Sener and V. Koltun. Multi-task learning as multi-objective optimization. In *NIPS*, pages 525–536,
513 2018.
- 514 [66] K. Shiohara and T. Yamasaki. Detecting deepfakes with self-blended images. In *CVPR*, pages 18720–18729,
515 2022.
- 516 [67] Y. Song and S. Ermon. Generative modeling by estimating gradients of the data distribution. In *NIPS*, page
517 11918–11930, 2019.
- 518 [68] X. Sun, R. Panda, R. Feris, and K. Saenko. AdaShare: Learning what to share for efficient deep multi-task
519 learning. In *NIPS*, pages 8728–8740, 2020.
- 520 [69] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. DeepFace: Closing the gap to human-level performance
521 in face verification. In *CVPR*, pages 1701–1708, 2014.
- 522 [70] G. W. Taylor, I. Spiro, C. Bregler, and R. Fergus. Learning invariance through imitation. In *CVPR*, pages
523 2729–2736, 2011.
- 524 [71] J. Thies, M. Zollhöfer, and M. Nießner. Deferred neural rendering: Image synthesis using neural textures.
525 *ACM TOG*, 38(4):1–12, 2019.
- 526 [72] J. Thies, M. Zollhöfer, M. Stamminger, C. Theobalt, and M. Nießner. Face2Face: Real-time face capture
527 and reenactment of RGB videos. In *CVPR*, pages 2387–2395, 2016.
- 528 [73] M.-F. Tsai, T.-Y. Liu, T. Qin, H.-H. Chen, and W.-Y. Ma. FRank: A ranking method with fidelity loss. In
529 *ACM SIGIR*, pages 383–390, 2007.
- 530 [74] M. Wallingford, H. Li, A. Achille, A. Ravichandran, C. Fowlkes, R. Bhotika, and S. Soatto. Task adaptive
531 parameter sharing for multi-task learning. In *CVPR*, pages 7561–7570, 2022.
- 532 [75] S.-Y. Wang, O. Wang, A. Owens, R. Zhang, and A. A. Efros. Detecting Photoshopped faces by scripting
533 Photoshop. In *ICCV*, pages 10072–10081, 2019.
- 534 [76] S.-Y. Wang, O. Wang, R. Zhang, A. Owens, and A. A. Efros. CNN-generated images are surprisingly easy
535 to spot...for now. In *CVPR*, pages 8695–8704, 2020.
- 536 [77] X. Yang, Y. Li, and S. Lyu. Exposing Deep Fakes using inconsistent head poses. In *ICASSP*, pages
537 8261–8265, 2019.
- 538 [78] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. Barlow Twins: Self-supervised learning via redundancy
539 reduction. In *ICML*, pages 12310–12320, 2021.
- 540 [79] Y. Zhang and Q. Yang. A survey on multi-task learning. *IEEE TKDE*, 34(12):5586–5609, 2021.
- 541 [80] H. Zhao, W. Zhou, D. Chen, T. Wei, W. Zhang, and N. Yu. Multi-attentional deepfake detection. In *CVPR*,
542 pages 2185–2194, 2021.
- 543 [81] T. Zhao, X. Xu, M. Xu, H. Ding, Y. Xiong, and W. Xia. Learning self-consistency for deepfake detection.
544 In *ICCV*, pages 15023–15033, 2021.
- 545 [82] P. Zhou, X. Han, V. I. Morariu, and L. S. Davis. Two-stream neural networks for tampered face detection.
546 In *CVPRW*, pages 1831–1839, 2017.