

Rethinking Domain Generalization for Face Anti-spoofing: Separability and Alignment

Yiyou Sun^{2*}, Yaojie Liu¹, Xiaoming Liu^{1,3}, Yixuan Li², Wen-Sheng Chu¹

¹Google Research, ²University of Wisconsin-Madison, ³Michigan State University

¹{yaojieliu, xiaomingli, wschu}@google.com, ²{sunnyiou, sharonli}@cs.wisc.edu, ³liuxm@cse.msu.edu

Abstract

This work studies the generalization issue of face anti-spoofing (FAS) models on domain gaps, such as image resolution, blurriness and sensor variations. Most prior works regard domain-specific signals as a negative impact, and apply metric learning or adversarial losses to remove them from feature representation. Though learning a domain-invariant feature space is viable for the training data, we show that the feature shift still exists in an unseen test domain, which backfires on the generalizability of the classifier. In this work, instead of constructing a domain-invariant feature space, we encourage domain separability while aligning the live-to-spoof transition (i.e., the trajectory from live to spoof) to be the same for all domains. We formulate this FAS strategy of separability and alignment (SA-FAS) as a problem of invariant risk minimization (IRM), and learn domain-variant feature representation but domain-invariant classifier. We demonstrate the effectiveness of SA-FAS on challenging cross-domain FAS datasets and establish state-of-the-art performance. Code is available at <https://github.com/sunnyiou/SAFAS>.

1. Introduction

Face recognition (FR) [16] has achieved remarkable success and has been widely employed in mobile access control and electronic payments. Despite the promise, FR systems still suffer from presentation attacks (PAs), including print attacks, digital replay, and 3D masks. As a result, face anti-spoofing (FAS) has been an important topic for almost two decades [3, 35, 45, 47, 66, 74, 76].

In early systems like building access and border control with limited variations (e.g., lighting and poses), simple methods [6, 17, 41] have exhibited promise. These algorithms are designed for the closed-world setting, where

*This work was done during Yiyou Sun’s internship at Google.

¹In statistics, spurious correlation is a mathematical relationship in which multiple events or variables are associated but not causally related.

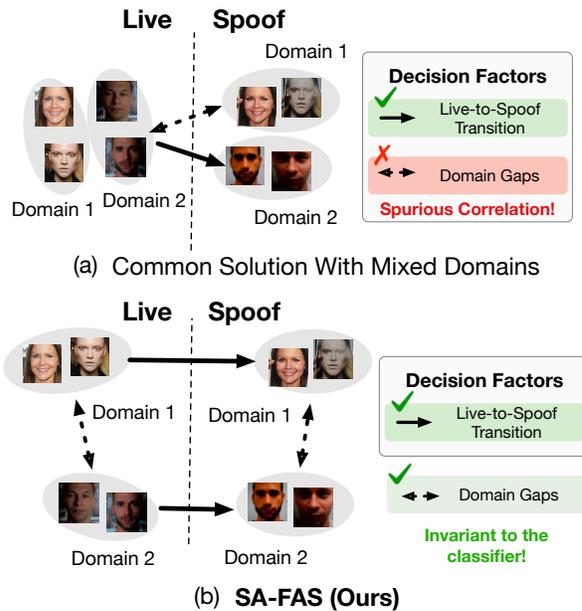


Figure 1. **Cross-domain FAS:** (a) Common FAS solutions aim to remove domain-specific signals and mix domains in one cluster. However, we empirically show domain-specific signals still exists in the feature space, and model might pick domain-specific signals as spurious correlation¹ for classification. (b) Our SA-FAS aims to retain domain signal. Specifically, we train a feature space with two critical properties: (1) **Separability:** Samples from different domains and live/spoof classes are well-separated; (2) **Alignment:** Live-to-spoof transitions are aligned in the same direction for all domains. With these two properties, our method keeps the domain-specific signals invariant to the decision boundary.

the camera and environment are assumed to be the same between train and test. This assumption, however, rarely holds for in-the-wild applications, e.g., mobile face unlock and sensor-invariant ID verification. Face images in those FAS cases may be acquired from wider angles, complex scenes, and different devices, where it is hard for training data to cover all the variations. These differences between training

and test data are termed domain gaps and the FAS solutions to tackle the domain gaps are termed cross-domain FAS.

Learning domain-invariant representation is the main approach in generic domain generalization [70], and has soon been widely applied to cross-domain FAS [30, 43, 44, 60, 67, 72]. Those methods consider domain-specific signals as a confounding factor for model generalization, and hence aim to remove domain discrepancy from the feature representation partially or entirely. Adversarial training is commonly applied so that upon convergence the domain discriminator cannot distinguish which domain the features come from. In addition, some methods apply metric learning to further regularize the feature space, *e.g.*, triplet loss [67], dual-force triplet loss [60], and single-side triplet loss [30].

There are two crucial issues that limit the generalization ability of these methods [30, 43, 44, 60, 67, 72] with domain-invariant feature losses. First, these methods posit a strong assumption that the feature space is perfectly domain-invariant after removing the domain-specific signals from training data. However, this assumption is unrealistic due to the limited size and domain variants of the training data, on which the loss might easily overfit during training. As shown in Fig. 7, the test distribution is more expanded compared to the training one, and the spatial relation between live and spoof has largely deviated from the learned classifier. Second, feature space becomes ambiguous when domains are mixed together. Note that the domain can carry information on certain image resolutions, blurriness and sensor patterns. If features from different domains are collapsed together [54], the live/spoof classifier will undesirably leverage spurious correlations to make the live/spoof predictions as shown in Fig. 1 (a), *e.g.*, comparing live from low-resolution domains to spoof from high-resolution ones. Such a classifier will unlikely generalize to a test domain when the correlation does not exist.

In this work, we rethink feature learning for cross-domain FAS. Instead of constructing a domain-invariant feature space, we aim to find a generalized classifier while explicitly maintaining domain-specific signals in the representation. Our strategy can be summarized by the following two properties:

- **Separability:** We encourage features from different domains and live/spoof classes to be separated which facilitates maintaining the domain signal. According to [4], representations with well-disentangled domain variation and task-relevant features are more general and transferable to different domains.
- **Alignment:** Inspired by [31], we regard spoofing as the process of transition. For similar PA types², the transition process would be similar, regardless of environments and sensor variations. With this assumption,

²This work focuses on print and replay attacks.

we regularize the live-to-spoof transition to be aligned in the same direction for all domains.

We refer to this new learning framework as *FAS with separability and alignment* (dubbed **SA-FAS**), shown in Fig. 1 (b). To tackle the separability, we leverage Supervised Contrastive Learning (SupCon) [33] to learn representations that force samples from the same domain and the same live/spoof labels to form a compact cluster. To achieve the alignment, we devise a novel Projected Gradient optimization strategy based on Invariant Risk Minimization (PG-IRM) to regularize the live-to-spoof transition invariant to the domain variance. With normalization, the feature space is naturally divided into two symmetric half-spaces: one for live and one for spoof (see Fig. 6). Domain variations will manifest inside the half-spaces but have minimal impact to the live/spoof classifier.

We summarize our contributions as three-fold:

- We offer a new perspective for cross-domain FAS. Instead of removing the domain signal, we propose to maintain it and design the feature space based on separability and alignment;
- We first systematically exploit the domain-variant representation learning by combining contrastive learning and effectively optimizing invariant risk minimization (IRM) through the projected gradient algorithm for cross-domain FAS;
- We achieve state-of-the-art performance on widely-used cross-domain FAS benchmark, and provide in-depth analysis and insights on how separability and alignment lead to the performance boost.

2. Related Work

Face Anti-Spoofing Face anti-spoofing attracts growing attention in several thriving directions. Early works exploit spontaneous human behaviors (*e.g.*, eye blinking, head motion) [36, 53] or predefined movements (*e.g.*, head-turning, expression changes) [12]. Later, hand-crafted features are utilized to describe spoof patterns, *e.g.*, LBP [6, 17], HoG [17, 75] and SIFT [55] features. Recently, deep neural networks have been applied to face anti-spoofing. There are classification-based methods [47, 66, 74], regression-based methods [3, 35, 45, 76], and generative models [31, 46, 48, 72]. In addition, the vision transformer also shows promising performance in tackling FAS [23, 29].

Cross-domain FAS Recently, several works explore learning FAS models from multiple domains that generalize to unseen ones. Some methods [15, 27, 39, 69, 80] require data from the target domain to adapt the model (*i.e.*, domain adaptation), while others [30, 34, 58, 60, 63, 72] learn shared features based on adversarial training and triplet loss (*i.e.*, domain generalization). A few methods [11, 61, 71] explore meta-learning to simulate the domain shift at training time. Most previous works regard the domain-specific signals as

a negative impact. Contrastively, our paper first systematically exploits the explicit usage of domain-specific signals by invariant risk minimization in cross-domain FAS.

Domain-invariant Classifier Learning a domain-invariant classifier has always been the focus of machine learning for decades [8, 9, 28, 64] and is also one of the keys to the success of domain generalization. Along this line, kernel-based methods [5, 18, 24, 26, 40, 51] propose to learn a domain-invariant kernel from the training data. Domain adversarial learning [19, 20, 25, 30, 40, 42, 56, 60, 72] adversarially trains the generator and discriminator while the generator is trained to fool the discriminator to learn domain invariant feature representations. Recently, Invariant Risk Minimization (IRM) and its variants [1, 2, 14, 37, 50, 56, 62] seek to directly enforce the optimal classifier on top of the representation space to be the same across all domains. However, IRM is known to be hard to optimize and can fail in non-linear optimization [32, 57]. In this paper, we propose an equivalent objective (PG-IRM) which is easier to optimize and achieve strong performance.

3. Method

We formally introduce the new learning framework, *FAS with separability and alignment* (dubbed **SA-FAS**). The goal is to produce a feature space with two critical properties: (1) *Separability*: We encourage samples from different domains and from different classes to be well-separated; (2) *Alignment*: Live-to-spoof transition³ is aligned in the same direction for all domains. These two properties work jointly: *separability* ensures the awareness of domain variance in the feature space; *alignment* encourages the domain variance to be invariant to its live-vs-spoof hyperplane.

This section is structured as follows: Sec. 3.1 describes the problem setup, followed by the algorithm design of separability (Sec. 3.2) and alignment (Sec. 3.3). Finally, Sec. 3.4 summarizes the training and inference processes.

3.1. Problem Setup

We start by defining the setting of the cross-domain FAS problem. We denote by $\mathcal{X} = \mathbb{R}^d$ the input space and $\mathcal{Y} = \{0 \text{ (live)}, 1 \text{ (spoof)}\}$ the output space. A learner is given access to a set of training data from E domains $\mathcal{E} = \{e^{(1)}, e^{(2)}, \dots, e^{(E)}\}$ and is evaluated on test domain e^* . Let e_i as the domain label for the i -th sample, we denote $\mathcal{D} = \{(\mathbf{x}_i, y_i, e_i)\}_{i=1}^N$ drawn from an unknown joint data distribution \mathcal{P} defined on $\mathcal{X} \times \mathcal{Y} \times \mathcal{E}$. Cross-domain FAS is a special binary classification problem to distinguish live and spoof faces from an unseen domain. The goal is to define a decision function:

$$f : \mathbf{x} \rightarrow \{0 \text{ (live)}, 1 \text{ (spoof)}\},$$

which classifies whether a sample \mathbf{x} from a new domain e^* is live or spoof.

³The transition can be considered as a path in the high-D manifold.

In our network architecture, function f consists of two components: (1) a deep neural network encoder $\phi : \mathcal{X} \rightarrow \mathbb{R}^m$ that maps the input \mathbf{x} to a l_2 -normalized feature embedding $\mathbf{z} = \phi(\mathbf{x})$; (2) a classifier (via a weight vector) $\beta : \mathbb{R}^m \rightarrow \mathbb{R}$ that maps the m -dimensional embedding \mathbf{z} to a scalar value, where a binary cross-entropy loss can be applied after using a sigmoid function. Because the true distribution of live/spoof data is unknown, the optimization commonly relies on an Empirical Risk Minimization (ERM).

Remark on the terminology: β can be considered as a norm vector of the hyperplane separating live and spoof samples. In the remaining part of the paper, when we use “**live-vs-spoof hyperplane**” or “**hyperplane**”, it has the same meaning as β . Note, “live-to-spoof transition” is an abstract procedure in the image space, while “live-vs-spoof hyperplane” refers to a concrete classifier in the feature space.

Preliminary on Empirical Risk Minimization (ERM):

ERM principle [65] is a ubiquitous strategy that merges data from all training domains and learns a predictor that minimizes an averaged training error. Specifically,

$$\mathcal{L}_{ERM} = \min_{\phi, \beta} \frac{1}{|\mathcal{E}|} \sum_{e \in \mathcal{E}} \mathcal{R}^e(\phi, \beta), \quad (1)$$

where the empirical risk function $\mathcal{R}^e(\phi, \beta)$ for a given environment e is defined by:

$$\mathcal{R}^e(\phi, \beta) \triangleq \mathbb{E}_{(\mathbf{x}_i, y_i, e_i=e) \sim \mathcal{D}} \ell(f(\mathbf{x}_i; \phi, \beta), y_i).$$

Common choices of the loss function $\ell(\cdot, \cdot)$ include cross-entropy loss [30] and L_1 regression loss [22, 45].

However, if samples from different domains are mixed together, ERM can utilize the easiest difference (image resolution, blurriness, camera setting) to differentiate live vs. spoof. Such a classifier will undesirably leverage spurious correlations to make live/spoof predictions [2]. Therefore, the naive strategy can hurt the generalization of the unseen domain. As shown in Fig. 2(a), ERM tends to fit all training data together and fails to learn a domain-invariant classifier with the mixed feature space.

3.2. Separability

We characterize the domain separability as supervised contrastive learning (dubbed SupCon) [33], one of the latest developments for visual representation learning. Unlike other contrastive learning methods [9, 10] that treat the augmented samples as a single class, SupCon aims to learn a representation space that gathers samples with the same labels while repelling samples from different ones. It naturally suits the need for the cross-domain FAS setting, since we treat samples with the same domain and with the same live/spoof label to form a cluster.

Given a training mini-batch $\{\mathbf{x}_i, y_i, e_i\}_{i=1}^b$, we augment [33] the mini-batch as $\{\tilde{\mathbf{x}}_i, \tilde{y}_i, \tilde{e}_i\}_{i=1}^{2b}$, using two

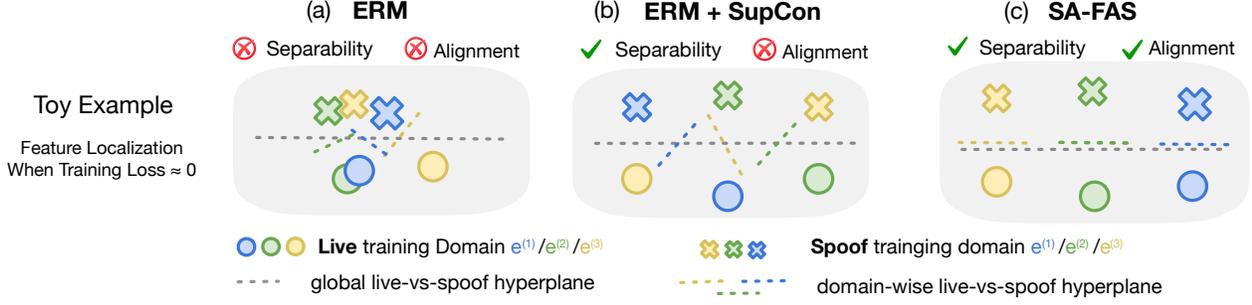


Figure 2. **Optimization objectives:** Illustration of feature space optimized by different objectives: (a) ERM, (b) ERM+SupCon, (c) SA-FAS (ours). Circle/cross denotes live/spoof label; different colors indicate different domains. A UMAP visualization for real data is provided in Appendix (Fig. 8) to support the feature distribution shown in the toy example.

random augmentations $\tilde{\mathbf{x}}_{2i}$ and $\tilde{\mathbf{x}}_{2i-1}$ of inputs \mathbf{x}_i , with $\tilde{y}_{2i-1} = \tilde{y}_{2i} = y_i$, $\tilde{e}_{2i-1} = \tilde{e}_{2i} = e_i$. These images are fed into the network, yielding L_2 -normalized embeddings $\{\mathbf{z}_i\}_{i=1}^{2b}$. The per-batch SupCon loss (separability loss) is defined as:

$$\mathcal{L}_{sep} = \sum_{i=1}^{2b} \frac{-1}{|S(i)|} \sum_{j \in S(i)} \log \frac{\exp(\mathbf{z}_i \cdot \mathbf{z}_j / \tau)}{\sum_{t=1, t \neq i}^{2b} \exp(\mathbf{z}_i \cdot \mathbf{z}_t / \tau)}, \quad (2)$$

where τ is a temperature parameter, i is the index of a sample typically called the *anchor*, $S(i) = \{j \in \{1, \dots, 2b\} : j \neq i, \tilde{y}_j = \tilde{y}_i, \tilde{e}_j = \tilde{e}_i\}$ is the index set of *positive samples* that have the same live/spoof labels and belong to the same domain as the anchor i , and $|S(i)|$ is its cardinality. All the other samples in the mini-batch are referred to as *negative samples*. Since positive samples are pulled together and negative samples are pushed apart, SupCon in Fig. 2(b) is capable of providing more distinguishable feature clusters for different domains and liveness classes, compared to a typical feature space learned by a vanilla ERM in Fig. 2(a).

3.3. Alignment

Fig. 2(b) also shows that separability alone is not sufficient for improving domain generalization. The separated feature clusters can be located in any place in the feature space, and hence the domain-wise optimal hyperplane remains **variant**. In this case, the global classifier can still undesirably incorporate the spurious correlation as the deciding factor as we show in Fig. 1. To tackle this, we naturally investigate the following problem:

How do we regularize a global live-vs-spoof hyperplane to align with domain-wise live-vs-spoof hyperplanes?

We propose to formulate this problem as Invariant Risk Minimization (IRM) [2], which aims to jointly optimize the feature space ϕ and the global live-vs-spoof hyperplane β , where β is also optimal for each domain, shown in Fig. 2(c).

Preliminary on Invariant Risk Minimization (IRM): Specifically, the IRM objective can be formulated as the fol-

lowing constrained optimization problem:

$$\min_{\phi, \beta^*} \frac{1}{|\mathcal{E}|} \sum_{e \in \mathcal{E}} \mathcal{R}^e(\phi, \beta^*) \rightarrow \mathcal{L}_{IRM} \quad (3)$$

$$s.t. \quad \beta^* \in \arg \min_{\beta} \mathcal{R}^e(\phi, \beta), \forall e \in \mathcal{E}. \quad (4)$$

Compared to the ERM (1), IRM enforces an additional constraint (4) to learn the domain-invariant hyperplanes. Specifically, if we define the domain-wise optimal hyperplane as $\beta_e \in \arg \min_{\beta} \mathcal{R}^e(\phi, \beta)$. A sufficient condition for constraint (4) to hold is $\beta_{e(1)} = \dots = \beta_{e(E)} = \beta^*$, which requires consistency between the globally optimal hyperplane and the domain-wise optimal hyperplanes. However, IRM is known to be hard to solve [32, 57] due to the bi-level optimization nature of objective (3) and constraint (4).

Projected Gradient Optimization for IRM (PG-IRM):

We leverage Projected Gradient (PG) algorithm [52] to solve the non-trivial optimization objective (3), termed as PG-IRM. In PG-IRM, we propose to optimize multiple hyperplanes and converge them into a globally one via projected gradient. In Appendix A.1, we provide detailed proof of PG-IRM objective being **equivalent** to IRM. Formally, the objective is rewritten as:

Theorem 1. (PG-IRM objective) For all $\alpha \in (0, 1)$, the IRM objective is equivalent to the following objective:

$$\min_{\phi, \beta_{e(1)}, \dots, \beta_{e(E)}} \frac{1}{|\mathcal{E}|} \sum_{e \in \mathcal{E}} \mathcal{R}^e(\phi, \beta_e) \rightarrow \mathcal{L}_{align} \quad (5)$$

$$s.t. \quad \forall e \in \mathcal{E}, \exists \beta_e \in \Omega_e(\phi), \beta_e \in \Upsilon_{\alpha}(\beta_e),$$

where the parametric constrained set for each environment is simplified as $\Omega_e(\phi) = \arg \min_{\beta} \mathcal{R}^e(\phi, \beta)$, and we define the α -adjacency set:

$$\Upsilon_{\alpha}(\beta_e) = \{v \mid \max_{e' \in \mathcal{E} \setminus e} \min_{\beta_{e'} \in \Omega_{e'}(\phi)} \|v - \beta_{e'}\|_2\} \quad (6)$$

$$\leq \alpha \max_{e' \in \mathcal{E} \setminus e} \min_{\beta_{e'} \in \Omega_{e'}(\phi)} \|\beta_e - \beta_{e'}\|_2 \quad (7)$$

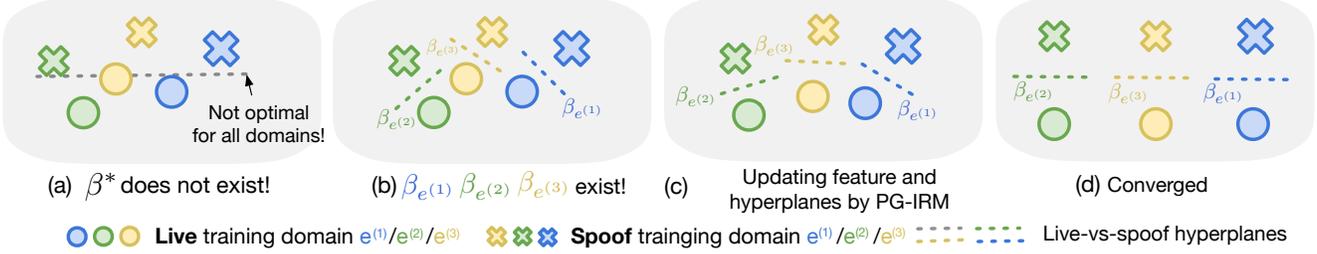


Figure 3. **PG-IRM optimization process:** An illustration of why a vanilla IRM can suffer from an infeasible solution (a), and how the proposed PG-IRM algorithm jointly updates the feature space and multiple hyperplanes towards convergence (b)-(d).

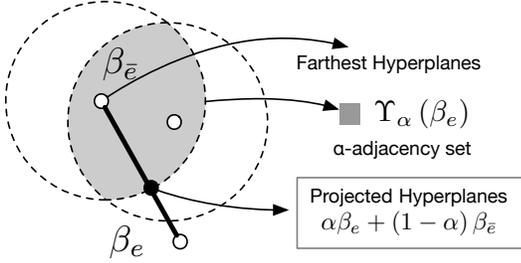


Figure 4. **Euclidean projection:** Illustration of Euclidean projection (solid black dot) to the α -adjacency set $\Upsilon_\alpha(\beta_e)$. Detailed proof and steps are provided in Alg. 2 in Appendix A.2.

Fig. 3 shows the intuition of the optimization process. For a 3-domain case, PG-IRM starts with a shared feature space ϕ and 3 separate hyperplanes $\beta_{e(1)}$, $\beta_{e(2)}$, $\beta_{e(3)}$ for each domain (Fig. 3(b)). After each projected gradient descent, the hyperplanes move closer with feature space jointly updated (Fig. 3(c)). Upon convergence, $\beta_{e(1)}$, $\beta_{e(2)}$, $\beta_{e(3)}$ become nearly identical (Fig. 3(d)), satisfying the IRM constraint $\beta^* = \beta_{e(1)} = \beta_{e(2)} = \beta_{e(3)}$ for all domains. We provide two main insights of our PG-IRM algorithm (see more details in Appendix A):

1. **Optimizing multiple hyperplanes:** Compared to the conventional IRM that optimizes a single hyperplane, it is easier to converge for PG-IRM that optimizes multiple hyperplanes (*i.e.*, one for each domain). Shown in Fig. 3(a-b), for the same feature space from the intermediate optimization stage, the solution β^* to conventional IRM may not exist and the optimization has to be terminated. In contrast, $\beta_{e(1)}, \dots, \beta_{e(E)}$ **always** exists (Fig. 3(b)) which makes solving for multiple hyperplanes more viable.
2. **Pushing hyperplanes to be closer:** To align $\beta_{e(1)}$, $\beta_{e(2)}$ and $\beta_{e(3)}$, PG-IRM updates domain-wise hyperplanes by interpolating with other hyperplanes. It can be mathematically considered as projecting the parameters of a hyperplane into the α -adjacency set $\Upsilon_\alpha(\beta_e)$ as we illustrated in Fig. 4.

Remark (why PG is not applicable to IRM): The PG algorithm can be infeasible for the conventional IRM,

Algorithm 1 Training pipeline for SA-FAS

- 1: **Input:** Training data $\mathcal{D} = \{(\mathbf{x}_i, y_i, e_i)\}_{i=1}^N$, network encoder ϕ , classifiers $\beta_{e(1)}, \dots, \beta_{e(E)}$, learning rate γ , alignment parameter α , alignment starting epoch T_α .
- 2: **for** t in $0, 1, \dots, T$ **do**
- 3: **Data Prep.:** Sample and augment a mini-batch.
- 4: **Forward/Backward:** Calculate gradient by \mathcal{L}_{all} .
- 5: **for** $e \in \mathcal{E}$ **do**
- 6: $\tilde{\beta}_e^{t+1} = \beta_e^t - \gamma \nabla_{\beta_e^t} \mathcal{L}_{all}$ \triangleright SGD
- 7: select $\beta_{\bar{e}}^t$ with $\bar{e} = \operatorname{argmax}_{e' \in \mathcal{E} \setminus e} \|\tilde{\beta}_e^{t+1} - \beta_{e'}^t\|_2$
- 8: $\alpha' = 1 - \mathbf{1}_{t > T_\alpha} (1 - \alpha)$ \triangleright α' is 1 when $t \leq T_\alpha$
- 9: $\beta_e^{t+1} = \alpha' \tilde{\beta}_e^{t+1} + (1 - \alpha') \beta_{\bar{e}}^t$ \triangleright Interpolation
- 10: **end for**
- 11: Update $\phi^{t+1} = \phi^t - \gamma \nabla_{\phi^t} \mathcal{L}_{all}$. \triangleright Update encoder
- 12: **end for**

as the solution set for (4) can be **empty** and is thus non-projectable. Our PG-IRM objective in Eq. (7) contains a non-empty α -adjacency set $\Upsilon_\alpha(\beta_e)$, and guarantees being projectable by simple linear interpolation.

3.4. Training and inference

Overall losses Considering the contrastive loss Eqn. (2), the overall objective (dubbed as SA-FAS) can be written as:

$$\begin{aligned} \min_{\phi, \beta_{e(1)}, \dots, \beta_{e(E)}} \mathcal{L}_{align} + \lambda \mathcal{L}_{sep} &\rightarrow \mathcal{L}_{all} \quad (8) \\ \text{s.t. } \forall e \in \mathcal{E}, \exists \beta_e \in \Omega_e(\phi), \beta_e \in \Upsilon_\alpha(\beta_e), & \end{aligned}$$

where λ is the coefficient for the loss term. The overall training pipeline is provided in Alg. 1.

Inference At the inference stage, we use the mean hyperplane from $\beta_{e(1)}, \dots, \beta_{e(E)}$ to get the final score. Specifically, the output is given by

$$f(\mathbf{x}) = \mathbb{E}_{e \in \mathcal{E}} [\beta_e^T \phi(\mathbf{x})].$$

Note that upon convergence, the cosine distance between any two of $\beta_{e(1)}, \dots, \beta_{e(E)}$ is very close to 1, *i.e.*, $\beta_{e(1)} \approx \dots \approx \beta_{e(E)}$. This observation is verified in Appendix C, with an ablation (converged angles *vs.* different α) in Appendix D.

Method (%)	OCI→M		OMI→C		OCM→I		ICM→O	
	HTER ↓	AUC ↑	HTER ↓	AUC ↑	HTER ↓	AUC ↑	HTER ↓	AUC ↑
MMD-AAE [40]	27.08	83.19	44.59	58.29	31.58	75.18	40.98	63.08
MADDG [60]	17.69	88.06	24.50	84.51	22.19	84.99	27.98	80.02
SSDG-M [30]	16.67	90.47	23.11	85.45	18.21	94.61	25.17	81.83
DR-MD-Net [68]	17.02	90.10	19.68	87.43	20.87	86.72	25.02	81.47
RFMeta [61]	13.89	93.98	20.27	88.16	17.30	90.48	16.45	91.16
NAS-FAS [77]	19.53	88.63	16.54	90.18	14.51	93.84	13.80	93.43
D2AM [11]	12.70	95.66	20.98	85.58	15.43	91.22	15.27	90.87
SDA [71]	15.40	91.80	24.50	84.40	15.60	90.10	23.10	84.30
DRDG [44]	12.43	95.81	19.05	88.79	15.56	91.79	15.63	91.75
ANRL [43]	10.83	96.75	17.83	89.26	16.03	91.04	15.67	91.90
SSAN-M [72]	10.42	94.76	16.47	90.81	14.00	94.58	19.51	88.17
SSDG-R [30]	7.38	97.17	10.44	95.94	11.71	96.59	15.61	91.54
SSAN-R [72]	6.67	98.75	10.00	96.67	8.88	96.79	13.72	93.63
PatchNet [66]	7.10	98.46	11.33	94.58	13.40	95.67	11.82	95.07
SA-FAS (Ours)	5.95	96.55	8.78	95.37	6.58	97.54	10.00	96.23

Table 1. **Comparisons with SoTA methods:** Cross-domain face anti-spoofing is evaluated among four popular benchmark datasets: CASIA (C), Idiap Replay (I), MSU-MFSD (M), and Oulu-NPU (O). Methods are compared at their best performance following the commonly used evaluation process [30]. ↑ indicates larger values are better, and ↓ indicates smaller values are better.

Method (%)	OCI→M			OMI→C			OCM→I			ICM→O			
	HTER ↓	AUC ↑	TPR95 ↑	HTER ↓	AUC ↑	TPR95 ↑	HTER ↓	AUC ↑	TPR95 ↑	HTER ↓	AUC ↑	TPR95 ↑	
SSDG-R [30]	14.65	^{1.21} / 91.93	^{1.35} / 53.68	^{2.56} / 28.76	^{0.89} / 80.91	^{1.10} / 41.47	^{2.68} / 22.84	^{1.14} / 78.67	^{1.31} / 50.80	^{5.95} / 15.83	^{1.29} / 92.13	^{0.96} / 66.54	^{4.00} / 25.72
SSAN-R [72]	21.79	^{3.68} / 84.06	^{3.78} / 51.91	^{4.28} / 26.44	^{2.91} / 78.84	^{2.83} / 45.36	^{4.29} / 35.39	^{8.04} / 70.13	^{9.03} / 64.00	^{2.70} / 25.72	^{3.74} / 79.37	^{4.69} / 36.75	^{5.19} / 23.49
PatchNet [66]	25.92	^{1.13} / 83.43	^{0.87} / 38.75	^{8.31} / 36.26	^{1.98} / 71.38	^{1.89} / 19.22	^{3.85} / 29.75	^{2.76} / 80.53	^{1.35} / 54.25	^{2.18} / 23.49	^{1.80} / 84.62	^{1.92} / 39.39	^{6.83} / 11.29
SA-FAS (Ours)	14.36	^{1.10} / 92.06	^{0.53} / 55.71	^{4.82} / 19.40	^{0.66} / 88.69	^{0.67} / 50.53	^{3.60} / 11.48	^{1.10} / 95.74	^{0.55} / 77.05	^{3.26} / 11.29	^{0.32} / 95.23	^{0.24} / 73.38	^{1.64} / 25.72

Table 2. **Evaluation upon convergence:** Evaluation of cross-domain face anti-spoofing among CASIA (C), Idiap Replay (I), MSU-MFSD (M), and Oulu-NPU (O) databases. Methods are compared at their mean/std performance based on the last 10 epochs.

4. Experiments

4.1. Experimental setups

Datasets and protocols We evaluate on four widely used datasets: Oulu-NPU (O) [7], CASIA (C) [79], Idiap Replay attack (I) [13], and MSU-MFSD (M) [73]. Following prior works, we treat each dataset as one domain and apply the leave-one-out test protocol to evaluate their cross-domain generalization. Specifically, we refer **OCI→M** to be the protocol that trains on Oulu-NPU, CASIA, Idiap Replay attack and tests on MSU-MFSD. **OMI→C**, **OCM→I** and **ICM→O** are defined in a similar fashion.

Implementation details The input images are cropped using MTCNN [78] and resized to 256×256. For fair comparisons with SoTA methods [30, 66, 72], we use the same ResNet-18 backbone. We train the network with SGD optimizer and an initial learning rate of 5e-3, which is decayed by 2 at epoch 40 and 80 and the total training epoch is 100 in most set-ups⁴. We set the weight decay as 5e-4 and the batch size as 96 for each training domain. For SA-FAS hyperparameters, we set $\alpha=0.995$, $\lambda=0.1$ and $T_a=20$.

Evaluation metrics We evaluate the model performance using three standard metrics: Half Total Error Rate (HTER), Area Under Curve (AUC), and True Positive Rate (TPR95) at a False Positive Rate (FPR) 5%. While HTER and AUC assess the theoretical performance, TPR at a certain FPR is adept at reflecting how well the model performs in practice.

⁴Due to the smaller training data size of **ICM**, we let the **ICM→O** to train for 300 epochs and decay at epoch 120 and 240.

4.2. Cross-domain performance

Tab. 1 summarizes our comparison with an extensive collection of recent studies, including SoTA methods: PatchNet [66], SSAN [72] and SSDG [30]. SA-FAS outperforms the rivals by a significant margin on cross-domain FAS benchmarks. In particular, we improve upon the best baseline [72] by 2.30% in HTER in the setting **OCM→I**, which is more than **25%** improvement.

Comparison upon convergence Note that the performance in Tab. 1 follows the convention in [30], which is reported on the training snapshot (*e.g.*, epoch 16) with the lowest test error. While this setting may manifest the best performance from the model, the results can significantly fluctuate on the test set and hard to reflect the generalization performance when a test set is unavailable (shown in Appendix B). To provide a more fair setting, we propose to report the average performance from the **last 10 epochs** upon convergence. In our case, the stopping criterion is either (1) the binary classification loss for live/spoof is smaller than 1e-3 for consecutive 10 epochs, or (2) the epoch number reaches max limit, whichever comes first.

In Tab. 2, we compare with SoTA methods in this setting, and provide three key observations: (1) The numbers are way worse than the ones in Tab. 1 across all methods, indicating the best model selected by conventional lowest test errors [30] has large randomness. This also shows that cross-domain FAS is far less-solved than expected. (2) The standard deviation in Tab. 2 denotes how stable each method

Method	HTEr↓	Average AUC↑	TPR95↑
SimCLR [9]	22.53 ^{1.31}	84.42 ^{1.04}	51.14 ^{3.44}
SimSiam [10]	18.89 ^{0.97}	89.93 ^{0.80}	56.62 ^{2.88}
Triplet [59]	18.75 ^{2.31}	88.11 ^{2.30}	50.53 ^{8.76}
SupCon (SSDG) [30]	17.91 ^{1.05}	90.10 ^{0.68}	61.98 ^{2.87}
SupCon [33]	17.03 ^{1.73}	90.68 ^{1.29}	56.72 ^{5.06}
ERM	17.22 ^{1.26}	90.21 ^{1.38}	58.62 ^{3.77}
DANN [21]	17.93 ^{1.02}	90.66 ^{0.56}	58.66 ^{3.14}
IRM-v1 [2]	17.41 ^{0.77}	91.16 ^{0.52}	60.98 ^{2.10}
VREx [38]	25.02 ^{1.92}	80.65 ^{2.20}	45.12 ^{3.78}
IB-IRM [1]	17.57 ^{0.74}	91.71 ^{0.51}	62.16 ^{2.35}
PG-IRM (Ours)	15.58 ^{0.96}	92.03 ^{0.62}	63.31 ^{2.59}
SA-FAS (Ours)	14.25 ^{0.79}	92.93 ^{0.49}	64.16 ^{3.33}

Table 3. **Ablation study:** The averaged performance is computed over all four cross-domain settings.

performs. Most methods can converge to a relatively stable status, while methods with an adversarial loss (e.g., [72]) have a relatively larger standard deviation, indicating adversarial loss might trigger more unstable training. (3) In our setting, SA-FAS still largely outperforms SoTA [30, 66, 72], which further validates the superiority of our method. Our method is also the most stable compared to SoTA, with the smallest standard deviation. We proceed by analyzing why traditional methods are less favorable in cross-domain FAS and why our methods perform better.

5. Ablation and Discussion

5.1. Effectiveness of loss components

Our overall objective function Eq. (8) consists of two parts: (a) Separability loss (\mathcal{L}_{sep}) for feature space; and (b) Alignment loss (\mathcal{L}_{align}) for regularizing the classifier. We ablate the contribution of each component in Tab. 3.

Separability loss We consider the two most common strategies used in the contrastive learning community (i.e., SimCLR [9], SimSiam [9]) and one in face recognition (i.e., Triplet loss [59]). We also provide the comparison of SupCon with the clustering policy used in SSDG [30]⁵. All losses are directly applied on the penultimate layer’s feature $\phi(\mathbf{x})$ with the same hyper-parameters, and all final classifications are supervised by ERM. We observe that the SupCon loss used in our framework outperforms other rivals. This validates the effectiveness of the domain-wise separable feature space for cross-domain FAS.

Alignment loss IRM objective is known to be hard to optimize. Other than the proposed PG-IRM, existing works IRM-v1 [2], IB-IRM [1], and VREx [37] alternatively consider a Lagrangian form:

$$\min_{\phi, \beta^*} \frac{1}{|\mathcal{E}|} \sum_{e \in \mathcal{E}} \left[\mathcal{R}^e(\phi, \beta^*) + \lambda \|\nabla_{\beta^*} \mathcal{R}^e(\phi, \beta^*)\|_2^2 \right]. \quad (9)$$

⁵SSDG assumes live samples in all domains form one cluster, and spoof samples in each domain respectively form the other three clusters.

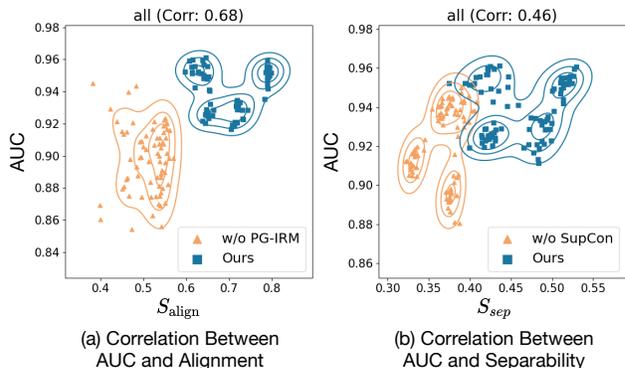


Figure 5. **Correlation of performance and SA-FAS:** Correlation between the test performance AUC and two properties measure (S_{align}/S_{sep}). Each dot represents one snapshot during the training stage in all four cross-domain settings. We provide separate figures for each setting in Appendix (Fig. 13).

In Tab. 3, we compare PG-IRM with the baseline ERM as well as other IRM alternatives, and our method shows a better overall performance. This shows further evidence that the Lagrangian penalty term can be ineffective, especially in the non-linear case [32, 57]. In comparison, PG-IRM optimizes the IRM objective directly with Projected Gradient, which clearly distinguishes it from existing methods.

Overall, the ablation studies suggest all components in our framework are indispensable to enhancing the generalization ability of cross-domain spoof detection.

5.2. Separability and alignment analysis

SA-FAS aims to produce a feature space with two critical properties: Separability and Alignment. In this section, we empirically investigate if these two properties can lead to a better generalization performance. Specifically, we provide two corresponding measures based on the learned classifiers and the extracted feature vector \mathbf{z} of samples from the **test** domain. We define the separability score as:

$$S_{sep} = 1 - \cos(\mathbb{E}_{spoof}[\mathbf{z}], \mathbb{E}_{live}[\mathbf{z}]),$$

where we measure the cosine angle between the center of live/spoof features. A separated feature space naturally leads to a small cosine value and thus a larger S_{sep} score. For the alignment score, we define:

$$S_{align} = \mathbb{E}_{e \in \mathcal{E}} [\cos(\beta_e, \mathbb{E}_{spoof}[\mathbf{z}] - \mathbb{E}_{live}[\mathbf{z}])], \quad (10)$$

where the trajectory from the center of spoof to live is treated as an *oracle vector*, which we measure how close it is with the norm vector of β of the learned hyperplane.

With the measure of two properties, we show their correlation to the generalization performance (i.e., AUC) in Fig. 5. We see that S_{sep} and S_{align} are positively related

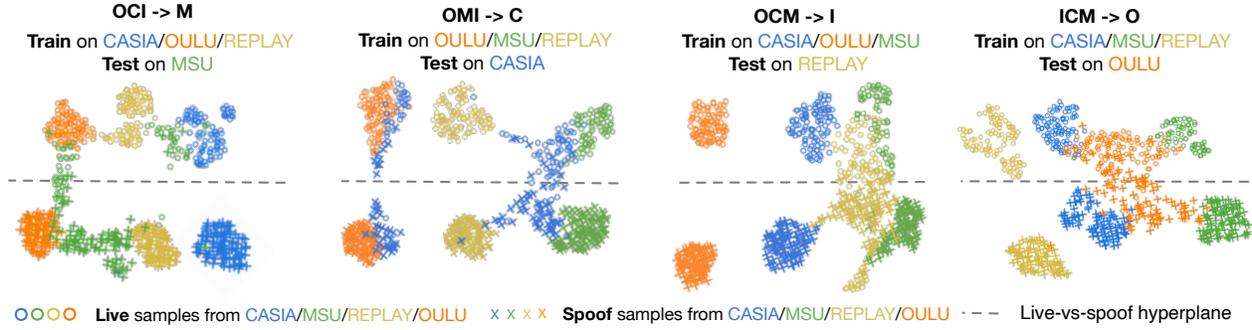


Figure 6. **Feature learned in different domains:** UMAP [49] projection of the penultimate layer of ResNet-18 trained with SA-FAS in the cross-test setting of face anti-spoofing datasets. The dotted line shows the decision boundary derived from training samples in 2D space.

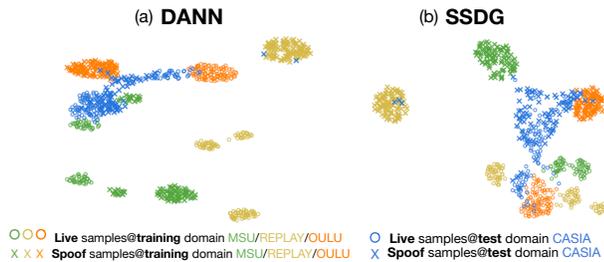


Figure 7. **Features of DANN vs. SSDG:** UMAP [49] visualization of the penultimate layer of ResNet-18 trained with DANN [21] and SSDG [30] in the cross-test setting of OMI→C.

to their test performance. It validates that these two properties are beneficial for a domain-invariant classifier. Specifically, Fig. 5(a) compares the setting with and without PG-IRM. Using PG-IRM leads to a higher alignment score and AUC, which verifies that PG-IRM can better align the live-vs-spoof hyperplanes for the unseen domain and improve the generalization ability. Similarly, Fig. 5(b) compares the setting with and without SupCon. The results validate that SupCon can lead to better separability in the feature space which benefits the classification.

5.3. UMAP visualization

Fig. 6 first provides UMAP [49] visualization of SA-FAS feature space from the penultimate layer. We see that the hyperplane between live samples and spoof samples is consistent across different training domains and also transferable to unseen test domains. For instance, in the setting of OMI→C, the test live samples in blue circles can be separated from the test spoof samples in blue crosses by the hyperplane. Another interesting finding is that some CASIA samples in blue are closer to OULU with **high resolution** and some are closer to MSU or REPLAY with **low resolution**, which reflects the fact that CASIA is a mixed dataset with both low and high resolution images. These findings validate that the domain gap (resolution) manifests in a way that is invariant to the live-vs-spoof hyperplane.

Beyond numerical and visual results, the superiority of

domain-variant feature space can also be validated by theoretical support. Specifically, the estimated error bound for binary classification in domain generalization [5] becomes larger if (M, n) is replaced with $(1, Mn)$, where M is the domain number and n is the training set size per domain. It indicates that separately training datasets from different domains is better than pooling them into one mixed dataset.

DANN [21] and SSDG [30] visualization We also compare the feature space of methods that aim to remove domain-specific signals from its feature representation. DANN [21] leverages the adversarial loss to encourage the backbone to provide a domain-invariant feature. Fig. 7(a) shows that the domain gap yet still broadly exists, especially for the test data from an unseen domain, which backfires on the generalizability of the classifier. Similarly, SSDG [30] learns a partial domain-invariant feature space where all live samples are clustered in one group while spoof samples are kept to be domain-dispersed. Although the degradation direction aligns better between train and test, compared to DANN, the domain gap still exists for live training samples as shown in Fig. 7(b). These findings further validate the necessity of regularizing the live-vs-spoof hyperplanes to be consistent across different domains.

6. Conclusion

This paper provides a new learning framework SA-FAS that learns domain-variant features but domain-invariant decision boundaries for cross-domain FAS. Our framework is naturally motivated, which facilitates invariant decision boundaries and learning distinguishable representations. We provide important theoretical insights that IRM objectives can be equivalently optimized by the PG with an alternative objective. Experiments show that SA-FAS can notably improve performance compared to the current best methods, establishing state-of-the-art. We also discuss the limitation of our work in Appendix E. We hope this paper will inspire more future works in incorporating domain-specific signals in FAS feature representation, and also extending this idea to broader domain generalization tasks.

References

- [1] Kartik Ahuja, Ethan Caballero, Dinghuai Zhang, Jean-Christophe Gagnon-Audet, Yoshua Bengio, Ioannis Mitliagkas, and Irina Rish. Invariance principle meets information bottleneck for out-of-distribution generalization. *Advances in Neural Information Processing Systems*, 34:3438–3450, 2021. 3, 7
- [2] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*, 2019. 3, 4, 7
- [3] Yousef Atoum, Yaojie Liu, Amin Jourabloo, and Xiaoming Liu. Face anti-spoofing using patch and depth-based CNNs. In *2017 IEEE International Joint Conference on Biometrics (IJCB)*, pages 319–328. IEEE, 2017. 1, 2
- [4] Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. Analysis of representations for domain adaptation. *Advances in neural information processing systems*, 19, 2006. 2
- [5] Gilles Blanchard, Aniket Anand Deshmukh, Ürun Dogan, Gyemin Lee, and Clayton Scott. Domain generalization by marginal transfer learning. *The Journal of Machine Learning Research*, 22(1):46–100, 2021. 3, 8
- [6] Zinelabidine Boulkenafet, Jukka Komulainen, and Abdenour Hadid. Face anti-spoofing based on color texture analysis. In *2015 IEEE international conference on image processing (ICIP)*, pages 2636–2640. IEEE, 2015. 1, 2
- [7] Zinelabidine Boulkenafet, Jukka Komulainen, Lei Li, Xiaoyi Feng, and Abdenour Hadid. Oulu-npu: A mobile face presentation attack database with real-world variations. In *2017 12th IEEE international conference on automatic face & gesture recognition (FG 2017)*, pages 612–618. IEEE, 2017. 6
- [8] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Proceedings of Advances in Neural Information Processing Systems*, 2020. 3
- [9] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *Proceedings of the international conference on machine learning*, pages 1597–1607. PMLR, 2020. 3, 7
- [10] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15750–15758, 2021. 3, 7
- [11] Zhihong Chen, Taiping Yao, Kekai Sheng, Shouhong Ding, Ying Tai, Jilin Li, Feiyue Huang, and Xinyu Jin. Generalizable representation learning for mixture domain face anti-spoofing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 1132–1139, 2021. 2, 6
- [12] Girija Chetty. Biometric liveness checking using multimodal fuzzy fusion. In *International Conference on Fuzzy Systems*, pages 1–8. IEEE, 2010. 2
- [13] Ivana Chingovska, André Anjos, and Sébastien Marcel. On the effectiveness of local binary patterns in face anti-spoofing. In *2012 BIOSIG-proceedings of the international conference of biometrics special interest group (BIOSIG)*, pages 1–7. IEEE, 2012. 6
- [14] Yo Joong Choe, Jiyeon Ham, and Kyubyong Park. An empirical study of invariant risk minimization. *arXiv preprint arXiv:2004.05007*, 2020. 3
- [15] Debayan Deb, Xiaoming Liu, and Anil Jain. Unified detection of digital and physical face attacks. In *Proceedings of the 17th International Conference on Automatic Face and Gesture Recognition (FG)*, 2023. 2
- [16] Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin loss for deep face recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4690–4699, 2019. 1
- [17] Tiago de Freitas Pereira, André Anjos, José Mario De Martino, and Sébastien Marcel. LBP-TOP based countermeasure against face spoofing attacks. In *Asian Conference on Computer Vision*, pages 121–132. Springer, 2012. 1, 2
- [18] Chuang Gan, Tianbao Yang, and Boqing Gong. Learning attributes equals multi-source domain generalization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 87–97, 2016. 3
- [19] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning*, pages 1180–1189. PMLR, 2015. 3
- [20] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1):2096–2030, 2016. 3
- [21] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1):2096–2030, 2016. 7, 8
- [22] Anjith George and Sébastien Marcel. Deep pixel-wise binary supervision for face presentation attack detection. In *2019 International Conference on Biometrics (ICB)*, pages 1–8. IEEE, 2019. 3
- [23] Anjith George and Sébastien Marcel. On the effectiveness of vision transformers for zero-shot face anti-spoofing. In *2021 IEEE International Joint Conference on Biometrics (IJCB)*, pages 1–8. IEEE, 2021. 2
- [24] Muhammad Ghifary, David Balduzzi, W Bastiaan Kleijn, and Mengjie Zhang. Scatter component analysis: A unified framework for domain adaptation and domain generalization. *IEEE transactions on pattern analysis and machine intelligence*, 39(7):1414–1430, 2016. 3
- [25] Rui Gong, Wen Li, Yuhua Chen, and Luc Van Gool. Dlow: Domain flow for adaptation and generalization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2477–2486, 2019. 3
- [26] Thomas Grubinger, Adriana Birlutiu, Holger Schöner, Thomas Natschläger, and Tom Heskes. Domain generalization based on transfer component analysis. In *International Work-Conference on Artificial Neural Networks*, pages 325–334. Springer, 2015. 3

- [27] Xiao Guo, Yaojie Liu, Anil Jain, and Xiaoming Liu. Multi-domain learning for updating face anti-spoofing models. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XIII*, pages 230–249. Springer, 2022. [2](#)
- [28] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9729–9738, 2020. [3](#)
- [29] Hsin-Ping Huang, Deqing Sun, Yaojie Liu, Wen-Sheng Chu, Taihong Xiao, Jinwei Yuan, Hartwig Adam, and Ming-Hsuan Yang. Adaptive transformers for robust few-shot cross-domain face anti-spoofing. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XIII*, pages 37–54. Springer, 2022. [2](#)
- [30] Yunpei Jia, Jie Zhang, Shiguang Shan, and Xilin Chen. Single-side domain generalization for face anti-spoofing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8484–8493, 2020. [2](#), [3](#), [6](#), [7](#), [8](#), [15](#), [16](#)
- [31] Amin Jourabloo, Yaojie Liu, and Xiaoming Liu. Face despoofing: Anti-spoofing via noise modeling. In *Proceedings of the European conference on computer vision (ECCV)*, pages 290–306, 2018. [2](#)
- [32] Pritish Kamath, Akilesh Tangella, Danica Sutherland, and Nathan Srebro. Does invariant risk minimization capture invariance? In *International Conference on Artificial Intelligence and Statistics*, pages 4069–4077. PMLR, 2021. [3](#), [4](#), [7](#)
- [33] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. In *Advances in Neural Information Processing Systems*, volume 33, pages 18661–18673, 2020. [2](#), [3](#), [7](#)
- [34] Taewook Kim and Yonghyun Kim. Suppressing spoof-irrelevant factors for domain-agnostic face anti-spoofing. *IEEE Access*, 9:86966–86974, 2021. [2](#)
- [35] Taewook Kim, YongHyun Kim, Inhan Kim, and Daijin Kim. Basn: Enriching feature representation using bipartite auxiliary supervisions for face anti-spoofing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019. [1](#), [2](#)
- [36] Klaus Kollreider, Hartwig Fronthaler, Maycel Isaac Faraj, and Josef Bigun. Real-time face detection and motion analysis with application in “liveness” assessment. *IEEE Transactions on Information Forensics and Security*, 2(3):548–558, 2007. [2](#)
- [37] David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghuai Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapolation (rex). In *International Conference on Machine Learning*, pages 5815–5826. PMLR, 2021. [3](#), [7](#)
- [38] David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghuai Zhang, Remi Le Priol, and Aaron Courville. Out-of-distribution generalization via risk extrapolation (rex). In *International Conference on Machine Learning*, pages 5815–5826. PMLR, 2021. [7](#)
- [39] Haoliang Li, Wen Li, Hong Cao, Shiqi Wang, Feiyue Huang, and Alex C Kot. Unsupervised domain adaptation for face anti-spoofing. *IEEE Transactions on Information Forensics and Security*, 13(7):1794–1809, 2018. [2](#)
- [40] Haoliang Li, Sinno Jialin Pan, Shiqi Wang, and Alex C Kot. Domain generalization with adversarial feature learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5400–5409, 2018. [3](#), [6](#)
- [41] Lei Li, Xiaoyi Feng, Zinelabidine Boulkenafet, Zhaoqiang Xia, Mingming Li, and Abdenour Hadid. An original face anti-spoofing approach using partial convolutional neural network. In *2016 Sixth International Conference on Image Processing Theory, Tools and Applications (IPTA)*, pages 1–6. IEEE, 2016. [1](#)
- [42] Ya Li, Xinmei Tian, Mingming Gong, Yajing Liu, Tongliang Liu, Kun Zhang, and Dacheng Tao. Deep domain generalization via conditional invariant adversarial networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 624–639, 2018. [3](#)
- [43] Shubao Liu, Ke-Yue Zhang, Taiping Yao, Mingwei Bi, Shouhong Ding, Jilin Li, Feiyue Huang, and Lizhuang Ma. Adaptive normalized representation learning for generalizable face anti-spoofing. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 1469–1477, 2021. [2](#), [6](#)
- [44] Shubao Liu, Ke-Yue Zhang, Taiping Yao, Kekai Sheng, Shouhong Ding, Ying Tai, Jilin Li, Yuan Xie, and Lizhuang Ma. Dual reweighting domain generalization for face presentation attack detection. *arXiv preprint arXiv:2106.16128*, 2021. [2](#), [6](#)
- [45] Yaojie Liu, Amin Jourabloo, and Xiaoming Liu. Learning deep models for face anti-spoofing: Binary or auxiliary supervision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 389–398, 2018. [1](#), [2](#), [3](#)
- [46] Yaojie Liu and Xiaoming Liu. Spoof trace disentanglement for generic face anti-spoofing. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(3):3813–3830, 2023. [2](#)
- [47] Yaojie Liu, Joel Stehouwer, Amin Jourabloo, and Xiaoming Liu. Deep tree learning for zero-shot face anti-spoofing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4680–4689, 2019. [1](#), [2](#)
- [48] Yaojie Liu, Joel Stehouwer, and Xiaoming Liu. On disentangling spoof trace for generic face anti-spoofing. In *European Conference on Computer Vision*, pages 406–422. Springer, 2020. [2](#)
- [49] Leland McInnes, John Healy, and James Melville. UMAP: Uniform manifold approximation and projection for dimension reduction. *arXiv preprint arXiv:1802.03426*, 2018. [8](#)
- [50] Jovana Mitrovic, Brian McWilliams, Jacob Walker, Lars Buesing, and Charles Blundell. Representation learn-

- ing via invariant causal mechanisms. *arXiv preprint arXiv:2010.07922*, 2020. 3
- [51] Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant feature representation. In *International Conference on Machine Learning*, pages 10–18. PMLR, 2013. 3
- [52] Jorge Nocedal and Stephen J Wright. *Numerical optimization*. Springer, 1999. 4
- [53] Gang Pan, Lin Sun, Zhaohui Wu, and Shihong Lao. Eyeblick-based anti-spoofing in face recognition from a generic webcam. In *2007 IEEE 11th international conference on computer vision*, pages 1–8. IEEE, 2007. 2
- [54] Vardan Papayan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40):24652–24663, 2020. 2
- [55] Keyurkumar Patel, Hu Han, and Anil K Jain. Secure face unlock: Spoof detection on smartphones. *IEEE transactions on information forensics and security*, 11(10):2268–2283, 2016. 2
- [56] Mohammad Mahfujur Rahman, Clinton Fookes, Mahsa Bakhtashmotlagh, and Sridha Sridharan. Correlation-aware adversarial domain adaptation and generalization. *Pattern Recognition*, 100:107124, 2020. 3
- [57] Elan Rosenfeld, Pradeep Kumar Ravikumar, and Andrej Risteski. The risks of invariant risk minimization. In *International Conference on Learning Representations*, 2020. 3, 4, 7
- [58] Suman Saha, Wenhao Xu, Menelaos Kanakis, Stamatios Georgoulis, Yuhua Chen, Danda Pani Paudel, and Luc Van Gool. Domain agnostic feature learning for image and video based face anti-spoofing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 802–803, 2020. 2
- [59] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015. 7
- [60] Rui Shao, Xiangyuan Lan, Jiawei Li, and Pong C Yuen. Multi-adversarial discriminative deep domain generalization for face presentation attack detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10023–10031, 2019. 2, 3, 6
- [61] Rui Shao, Xiangyuan Lan, and Pong C Yuen. Regularized fine-grained meta face anti-spoofing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 11974–11981, 2020. 2, 6
- [62] Anoopkumar Sonar, Vincent Pacelli, and Anirudha Majumdar. Invariant policy optimization: Towards stronger generalization in reinforcement learning. In *Learning for Dynamics and Control*, pages 21–33. PMLR, 2021. 3
- [63] Joel Stehouwer, Amin Jourabloo, Yaojie Liu, and Xiaoming Liu. Noise modeling, synthesis and classification for generic object anti-spoofing. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7294–7303, 2020. 2
- [64] Aaron Van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv e-prints*, pages arXiv–1807, 2018. 3
- [65] Vladimir Vapnik. Principles of risk minimization for learning theory. *Advances in neural information processing systems*, 4, 1991. 3
- [66] Chien-Yi Wang, Yu-Ding Lu, Shang-Ta Yang, and Shang-Hong Lai. PatchNet: A simple face anti-spoofing framework via fine-grained patch recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20281–20290, 2022. 1, 2, 6, 7, 15
- [67] Guoqing Wang, Hu Han, Shiguang Shan, and Xilin Chen. Improving cross-database face presentation attack detection via adversarial domain adaptation. In *2019 International Conference on Biometrics (ICB)*, pages 1–8. IEEE, 2019. 2
- [68] Guoqing Wang, Hu Han, Shiguang Shan, and Xilin Chen. Cross-domain face presentation attack detection via multi-domain disentangled representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6678–6687, 2020. 6
- [69] Guoqing Wang, Hu Han, Shiguang Shan, and Xilin Chen. Unsupervised adversarial domain adaptation for cross-domain face presentation attack detection. *IEEE Transactions on Information Forensics and Security*, 16:56–69, 2020. 2
- [70] Jindong Wang, Cuiling Lan, Chang Liu, Yidong Ouyang, Tao Qin, Wang Lu, Yiqiang Chen, Wenjun Zeng, and Philip Yu. Generalizing to unseen domains: A survey on domain generalization. *IEEE Transactions on Knowledge and Data Engineering*, 2022. 2
- [71] Jingjing Wang, Jingyi Zhang, Ying Bian, Youyi Cai, Chunmao Wang, and Shiliang Pu. Self-domain adaptation for face anti-spoofing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 2746–2754, 2021. 2, 6
- [72] Zhuo Wang, Zezheng Wang, Zitong Yu, Weihong Deng, Jiahong Li, Tingting Gao, and Zhongyuan Wang. Domain generalization via shuffled style assembly for face anti-spoofing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4123–4133, 2022. 2, 3, 6, 7, 15
- [73] Di Wen, Hu Han, and Anil K Jain. Face spoof detection with image distortion analysis. *IEEE Transactions on Information Forensics and Security*, 10(4):746–761, 2015. 6
- [74] Jianwei Yang, Zhen Lei, and Stan Z Li. Learn convolutional neural network for face anti-spoofing. *arXiv preprint arXiv:1408.5601*, 2014. 1, 2
- [75] Jianwei Yang, Zhen Lei, Shengcai Liao, and Stan Z Li. Face liveness detection with component dependent descriptor. In *2013 International Conference on Biometrics (ICB)*, pages 1–6. IEEE, 2013. 2
- [76] Zitong Yu, Xiaobai Li, Xuesong Niu, Jingang Shi, and Guoying Zhao. Face anti-spoofing with human material perception. In *European conference on computer vision*, pages 557–575. Springer, 2020. 1, 2
- [77] Zitong Yu, Jun Wan, Yunxiao Qin, Xiaobai Li, Stan Z Li, and Guoying Zhao. NAS-FAS: Static-dynamic central difference

network search for face anti-spoofing. *IEEE transactions on pattern analysis and machine intelligence*, 43(9):3005–3023, 2020. 6

- [78] Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE signal processing letters*, 23(10):1499–1503, 2016. 6
- [79] Zhiwei Zhang, Junjie Yan, Sifei Liu, Zhen Lei, Dong Yi, and Stan Z Li. A face antispoofing database with diverse attacks. In *2012 5th IAPR international conference on Biometrics (ICB)*, pages 26–31. IEEE, 2012. 6
- [80] Qianyu Zhou, Ke-Yue Zhang, Taiping Yao, Ran Yi, Kekai Sheng, Shouhong Ding, and Lizhuang Ma. Generative domain adaptation for face anti-spoofing. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part V*, pages 335–356. Springer, 2022. 2

A. Detailed Proof

In the main paper, we propose to use a project gradient algorithm to efficiently optimize the hard IRM objective. In this section, we provide a formal proof composed of two main steps: (1) We show in Section A.1 that the original IRM objective is equivalent to the PG-IRM objective shown in Theorem 1. (2) In Section A.2, we show that the PG-IRM objective can be efficiently optimized by the project gradient descent algorithm illustrated in Alg. 2.

A.1. PG-IRM objective is equivalent to IRM

As a recap of our learning setting, a learner is given access to a set of training data from E environments $\mathcal{E} = \{e^{(1)}, e^{(2)}, \dots, e^{(E)}\}$ and the IRM objective is the following constrained optimization problem:

$$\min_{\phi, \beta^*} \frac{1}{|\mathcal{E}|} \sum_{e \in \mathcal{E}} \mathcal{R}^e(\phi, \beta^*) \quad (11)$$

$$\text{s.t. } \beta^* \in \arg \min_{\beta} \mathcal{R}^e(\phi, \beta) \quad \forall e \in \mathcal{E}, \quad (12)$$

where the risk function for a given domain/distribution e is:

$$\mathcal{R}^e(\phi, \beta) \doteq \mathbb{E}_{(\mathbf{x}_i, y_i, e_i=e) \sim \mathcal{D}} \ell(f(\mathbf{x}_i; \phi, \beta), y_i).$$

Theorem. (Recap of Theorem 1) For all $\alpha \in (0, 1)$, the IRM objective is equivalent to the following objective:

$$\min_{\phi, \beta_{e^{(1)}}, \dots, \beta_{e^{(E)}}} \frac{1}{|\mathcal{E}|} \sum_{e \in \mathcal{E}} \mathcal{R}^e(\phi, \beta_e) \quad (13)$$

$$\text{s.t. } \forall e \in \mathcal{E}, \exists \beta_e \in \Omega_e(\phi), \beta_e \in \Upsilon_{\alpha}(\beta_e), \quad (14)$$

where the parametric constrained set for each environment is simplified as

$$\Omega_e(\phi) = \arg \min_{\beta} \mathcal{R}^e(\phi, \beta),$$

and we define

$$\begin{aligned} \Upsilon_{\alpha}(\beta_e) &= \{v \mid \min_{\forall e' \in \mathcal{E} \setminus e, \beta_{e'} \in \Omega_{e'}(\phi)} \|v - \beta_{e'}\|_2 \\ &\leq \alpha \min_{\forall e' \in \mathcal{E} \setminus e, \beta_{e'} \in \Omega_{e'}(\phi)} \|\beta_e - \beta_{e'}\|_2\} \end{aligned} \quad (15)$$

Proof. The constraint (12) means that the β^* is the optimal linear classifier at all environments, which is equivalent to saying that β^* lies in the joint of the optimal solution set in each environment. Equivalently, we can formalize the optimization target (11) as a parametric constrained optimization problem with constrain:

$$\beta^* \in \bigcap_{e \in \mathcal{E}} \underbrace{\arg \min_{\beta} \mathcal{R}^e(\phi, \beta)}_{\Omega_e(\phi)}, \quad (16)$$

where the parametric constrained set for each environment is $\Omega_e(\phi) = \arg \min_{\beta} \mathcal{R}^e(\phi, \beta)$ (Note that $\Omega_e(\phi)$ can be a set with cardinality bigger than 1, since the optimal linear classifier may not be unique). The constraint (16) implies that β^* lies in the joint set of $\Omega_e(\phi)$, which also means that there is an element in each $\Omega_e(\phi)$ equal to β^* . We refer to such element to be $\beta_e \in \Omega_e(\phi)$, and we have the alternative form:

$$\forall e \in \mathcal{E}, \exists \beta_e \in \Omega_e(\phi), \beta^* = \beta_e \quad (17)$$

Equivalently,

$$\forall e \in \mathcal{E}, \exists \beta_e \in \Omega_e(\phi), \beta_e \in \bigcap_{e' \in \mathcal{E} \setminus e} \underbrace{\Omega_{e'}(\phi)}_{\text{by (16) and (17)}} \quad (18)$$

The interpretation of constraint (18) is that — for all environments, there is a hyperplane in the optimal set $\Omega_e(\phi)$ that also lies in the intersection of other environments' optimal set ($\bigcap_{e' \in \mathcal{E} \setminus e} \Omega_{e'}(\phi)$). Now we rewrite the optimization target (3) as:

$$\min_{\phi, \beta_{e^{(1)}}, \dots, \beta_{e^{(E)}}} \frac{1}{|\mathcal{E}|} \sum_{e \in \mathcal{E}} \mathcal{R}^e(\phi, \beta_e) \quad (19)$$

$$\text{s.t. } \forall e \in \mathcal{E}, \exists \beta_e \in \Omega_e(\phi), \beta_e \in \bigcap_{e' \in \mathcal{E} \setminus e} \Omega_{e'}(\phi) \quad (20)$$

In this way, we can get rid of finding a unique β^* , but instead optimizing multiple linear classifiers $\beta_{e^{(1)}}, \dots, \beta_{e^{(E)}}$, which is easier to optimize in a relaxed form as we will show next.

One key challenge for this optimization problem is that there is no guarantee that $\bigcap_{e' \in \mathcal{E} \setminus e} \Omega_{e'}(\phi)$ is non-empty for a feature extractor ϕ and β_e . We therefore relax the optimization target as:

$$\min_{\phi, \epsilon, \beta_{e^{(1)}}, \dots, \beta_{e^{(E)}}} \frac{1}{|\mathcal{E}|} \sum_{e \in \mathcal{E}} \mathcal{R}^e(\phi, \beta_e) \quad (21)$$

$$\text{s.t. } \forall e \in \mathcal{E}, \exists \beta_e \in \Omega_e(\phi), \max_{e' \in \mathcal{E} \setminus e} \|\beta_e - \Omega_{e'}(\phi)\|_2 \leq \epsilon, \quad (22)$$

relax $\beta_e \in \bigcap_{e' \in \mathcal{E} \setminus e} \Omega_{e'}(\phi)$

where we define the l_2 distance between a vector β and a set Ω as: $\|\beta - \Omega\|_2 = \min_{v \in \Omega} \|\beta - v\|_2$.

Practically, ϵ can be set to be any variable converging to 0 during the optimization stage. Without losing the generality, we change the constraint (22) to the following form:

$$\begin{aligned} \forall e \in \mathcal{E}, \exists \beta_e \in \Omega_e(\phi), \\ \max_{e' \in \mathcal{E} \setminus e} \min_{\beta_{e'} \in \Omega_{e'}(\phi)} \|\beta_e - \beta_{e'}\|_2 \leq \\ \alpha \max_{e' \in \mathcal{E} \setminus e} \min_{\beta_{e'} \in \Omega_{e'}(\phi)} \|\beta_e - \beta_{e'}\|_2, \end{aligned} \quad (23)$$

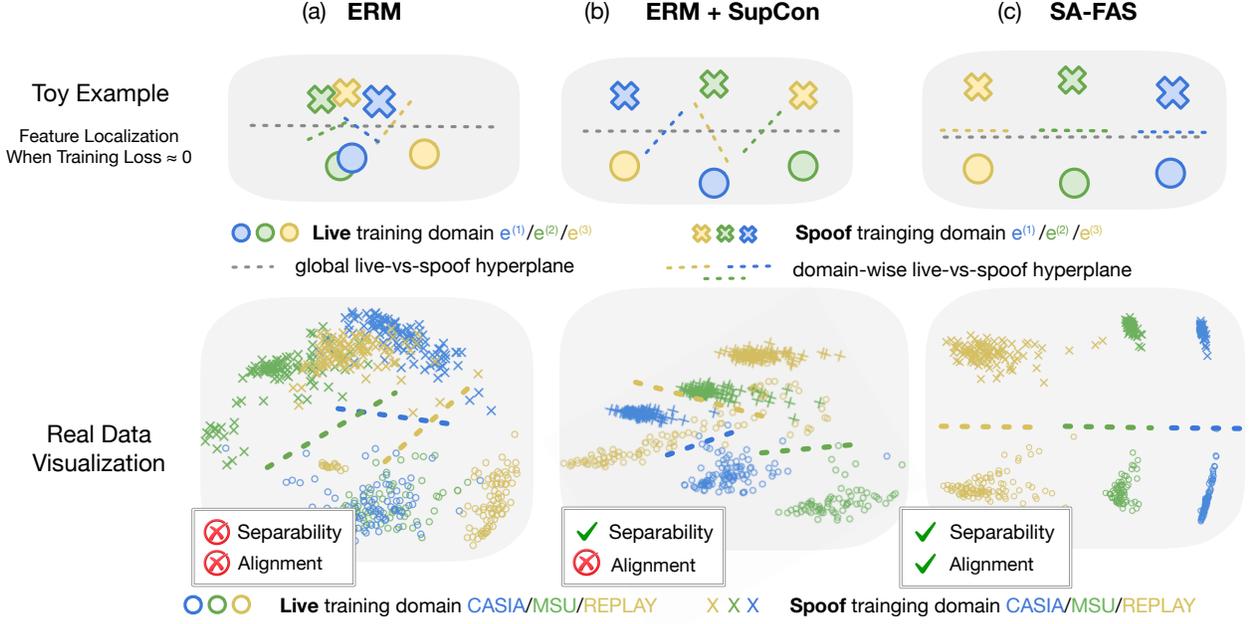


Figure 8. Illustration of feature space with three optimization objectives (ERM/ERM+SupCon/SA-FAS). For each objective, the first row shows the feasible solution in toy examples where each domain with a live/spoof label is represented by one circle/cross. The second row shows the visualization of real data via linear projection. The visualization is conducted by inserting and scattering the features from a 2-dimensional hidden layer between the penultimate layer and the final output layer.

where $\alpha \in (0, 1)$. Note that constraint (23) will be satisfied only when $\max_{e' \in \mathcal{E} \setminus e} \min_{\beta_{e'} \in \Omega_{e'}(\phi)} \|\beta_e - \beta_{e'}\|_2 = 0$. Therefore, constraint (23) is equivalent to constraint (18), and thus equivalent to the original constraint (12).

If we let the set

$$\begin{aligned} \Upsilon_\alpha(\beta_e) &= \{v \mid \max_{e' \in \mathcal{E} \setminus e} \min_{\beta_{e'} \in \Omega_{e'}(\phi)} \|v - \beta_{e'}\|_2 \\ &\leq \alpha \max_{e' \in \mathcal{E} \setminus e} \min_{\beta_{e'} \in \Omega_{e'}(\phi)} \|\beta_e - \beta_{e'}\|_2\} \end{aligned} \quad (24)$$

Then the constraint (23) can be simplified to

$$\forall e \in \mathcal{E}, \exists \beta_e \in \Omega_e(\phi), \beta_e \in \Upsilon_\alpha(\beta_e) \quad (25)$$

□

A.2. Projected Gradient Optimization for PG-IRM objective

We proceed with introducing how the Projected Gradient Descent can effectively optimize the PG-IRM objective. We start by introducing the background of the Projected Gradient Descent algorithm.

Projected Gradient Descent is commonly applied in constrained optimization, which aims to find a point θ achieving the smallest value of some loss function \mathcal{L} subject to the

requirement that θ is contained in the feasible set Ω . Formally, the objective can be written as:

$$\min_{\theta \in \Omega} \mathcal{L}(\theta)$$

If we minimize the objective $\mathcal{L}(\theta)$ by gradient descent, we have

$$(GD) \quad \theta := \theta - \gamma \nabla \mathcal{L}(\theta),$$

where γ is the step size. However, it is not guaranteed that the updated θ still falls into the set Ω . The projected gradient descent (PGD) algorithm is designed to project the solution back in the feasible set. Formally,

$$(PGD) \quad \theta := P_\Omega(\theta - \gamma \nabla \mathcal{L}(\theta)),$$

where the $P_\Omega(\cdot)$ is defined as the Euclidean Projection:

$$P_\Omega(u) = \arg \min_{v \in \Omega} \|u - v\|_2$$

In the PG-IRM objective, we have the constraint set $\Omega = \Upsilon_\alpha(\beta_e)$, we show in the next Lemma 2 that the Euclidean Projection from β_e to $\Upsilon_\alpha(\beta_e)$ is equivalent to the linear interpolation between β_e and the farthest hyperplane $\beta_{\bar{e}}$ for environment \bar{e} .

Lemma 2. *Given that*

Algorithm 2 PG-IRM

Initialize $\phi, \beta_{e(1)}, \dots, \beta_{e(\mathcal{E})}$, learning rate γ , alignment parameter α , alignment starting epoch T_a .

for t in $0, 1, \dots$, **do**

Run forward pass and calculate the gradient.

for $e \in \mathcal{E}$ **do**

$$\tilde{\beta}_e^{t+1} = \beta_e^t - \gamma \nabla_{\beta_e^t} \mathcal{L}_{PG-IRM}$$

$$\alpha' := 1 - \mathbf{1}_{t > T_a} (1 - \alpha)$$

select $\beta_{\bar{e}}^t$ with $\bar{e} = \operatorname{argmax}_{e' \in \mathcal{E} \setminus e} \|\tilde{\beta}_e^{t+1} - \beta_{e'}^t\|_2$

$$\beta_e^{t+1} = \alpha' \tilde{\beta}_e^{t+1} + (1 - \alpha') \beta_{\bar{e}}^t$$

end for

Update $\phi^{t+1} = \phi^t - \gamma \nabla_{\phi^t} \mathcal{L}_{PG-IRM}$.

end for

$$\begin{aligned} \Upsilon_\alpha(\beta_e) &= \{v \mid \max_{e' \in \mathcal{E} \setminus e} \min_{\beta_{e'} \in \Omega_{e'}(\phi)} \|v - \beta_{e'}\|_2 \\ &\leq \alpha \max_{e' \in \mathcal{E} \setminus e} \min_{\beta_{e'} \in \Omega_{e'}(\phi)} \|\beta_e - \beta_{e'}\|_2\} \end{aligned}$$

We have:

$$P_{\Upsilon_\alpha(\beta_e)}(\beta_e) = \alpha \beta_e + (1 - \alpha) \beta_{\bar{e}},$$

where $\beta_{\bar{e}}$ is selected with $\bar{e} = \operatorname{argmax}_{e' \in \mathcal{E} \setminus e} \|\beta_e - \beta_{e'}\|_2$.

Proof. We give the proof in an intuitive way shown in Figure 9. Specifically, the feasible region $\Upsilon_\alpha(\beta_e)$ can be regarded as an intersection of several hyper-spheres centered with all domain-wise live-vs-spoof hyperplanes $\beta_{e'}$. The radius is given by the α multiplying the distance to the farthest hyperplane $\beta_{\bar{e}}$. Therefore the Euclidean projection of β_e to the feasible set simultaneously lies on the surface of the hypersphere and the line segments between β_e and $\beta_{\bar{e}}$. It can be easily verified that

$$P_{\Upsilon_\alpha(\beta_e)}(\beta_e) = \alpha \beta_e + (1 - \alpha) \beta_{\bar{e}},$$

satisfies the given criteria. \square

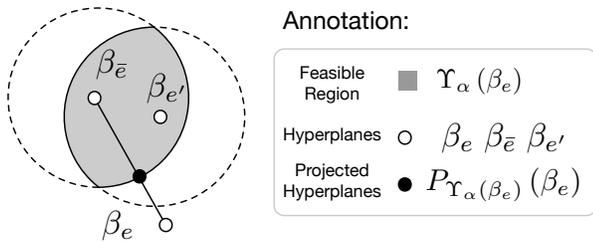


Figure 9. Illustration of the Euclidean projection results (solid black dot) to the feasible set $\Upsilon_\alpha(\beta_e)$.

Main results. When we have the projected form on the constraint set, deriving the optimization strategy is thus straightforward. As shown in Alg. 2, we first calculate the gradient of hyperplanes for all domains

$$\tilde{\beta}_e^{t+1} = \beta_e^t - \gamma \nabla_{\beta_e^t} \mathcal{L}_{PG-IRM}.$$

We then select the farthest domain-wise hyperplanes $\beta_{\bar{e}}$ from other environments. The final projection results are thus given by

$$\beta_e^{t+1} = \alpha' \tilde{\beta}_e^{t+1} + (1 - \alpha') \beta_{\bar{e}},$$

as we demonstrated in Lemma 2.

Remark on the T_a . In the first T_a epochs, we let the feature encoder ϕ and domain-wise hyperplanes β_e trained in a standard way. The goal is to ensure that the hyperplanes β_e will reach or be close to the minimum of the domain-wise empirical risk, and we have:

$$\beta_e \in \Omega_e(\phi).$$

In Alg. 2, we use an additional parameter α' to manifest this procedure:

$$\alpha' := 1 - \mathbf{1}_{t > T_a} (1 - \alpha)$$

Specifically, when $t < T_a$, $\alpha' = 1$, which means the original gradient descent algorithm is applied. When $t > T_a$, $\alpha' = \alpha$, the projected gradient descent takes charge.

B. Why do we need a fair setting?

By visualizing the line plot of the HTER performance over 100 training epochs in Fig. 10, we realize the test performance on the unseen domain is highly testset-dependent and unstable especially in the early epochs. Therefore, the best number reported commonly adopted in existing literature [30, 66, 72] usually happens in an unpredictable earlier epoch. Such “best” snapshot is also hard to be selected by validation strategy because we have zero information regarding the test domain. As an alternative, we noticed that the test performance is more stable in the last 10 epochs upon convergence, which motivates us to propose using a fairer comparison strategy introduced in Section 4.

C. Convergence of PG-IRM

Recall that in PG-IRM, we optimize multiple linear classifiers simultaneously $\beta_{e(1)}, \beta_{e(2)}, \beta_{e(3)}$ and gradually align them during training. In this section, we would like to verify if PG-IRM indeed regularizes domain classifiers to be close to each other and finally converges to the same one $\beta^* = \beta_{e(1)} = \beta_{e(2)} = \beta_{e(3)}$. Empirically, we use the averaged cosine distance between domain classifiers to measure the distance between them:

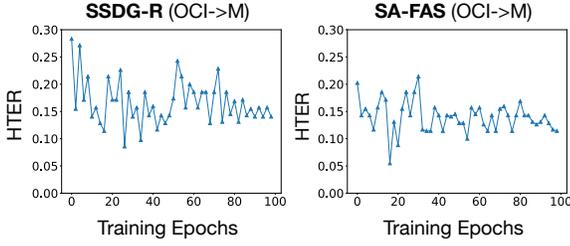


Figure 10. The line plot of the HTER performance tested on MSU dataset when trained on CASIA, Replay and OULU with SSDG-R [30] and SA-FAS over 100 training epochs.

$$S_{\cos} = \mathbb{E}_{e, e' \in \mathcal{E}, e \neq e'} [\cos(\beta_e, \beta_{e'})]$$

As shown in Fig. 11, the averaged cosine value between domain classifiers diminishes gradually and finally converges to 1, which suggests that they converge to a β^* that is aligned for all domains.

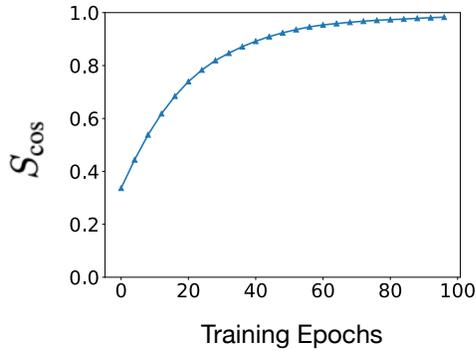


Figure 11. The line plot of the S_{\cos} when trained on CASIA, Replay and OULU with PG-IRM over 100 training epochs.

D. Sensitivity Analysis

In this section, we perform the sensitivity analysis of hyper-parameter settings for SA-FAS in Fig. 12. The performance comparison in the bar plot for each hyper-parameter is reported by fixing other hyper-parameters. In the figure, we observe that the performance of SA-FAS is less sensitive to the learning rate and the alignment starting epoch compared with the maximum gap of 1.2% in the given range. We also notice that choosing the right alignment parameter α is more important, since a proper α ensures the domain-wise decision boundaries are aligned not too fast and not too slow. In the extreme case, if $\alpha = 0$, it degenerates to the ERM after epoch T_a and if $\alpha = 1$, the domain-wise boundaries will never get aligned with each other. In summary, our algorithm does not require heavy

hyper-parameter tuning as long as it falls into a reasonable range.

E. Limitation

Our work has two limitations. Firstly, our framework assumes the dataset collected from each domain contains both live and spoof data. For example, SA-FAS can not handle the training data with live samples only from domain A and spoof samples only from domain B. Secondly, SA-FAS may cause extra computation costs when the domain amount is very large since we set up one hyperplane for each domain.

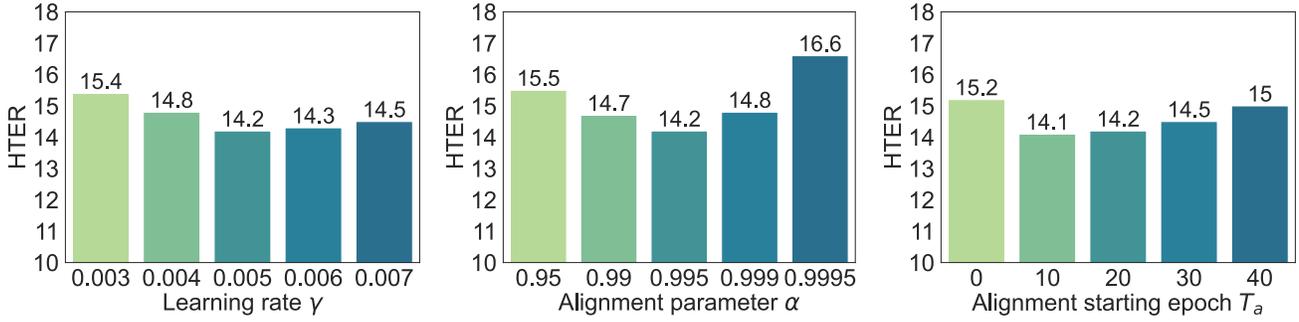
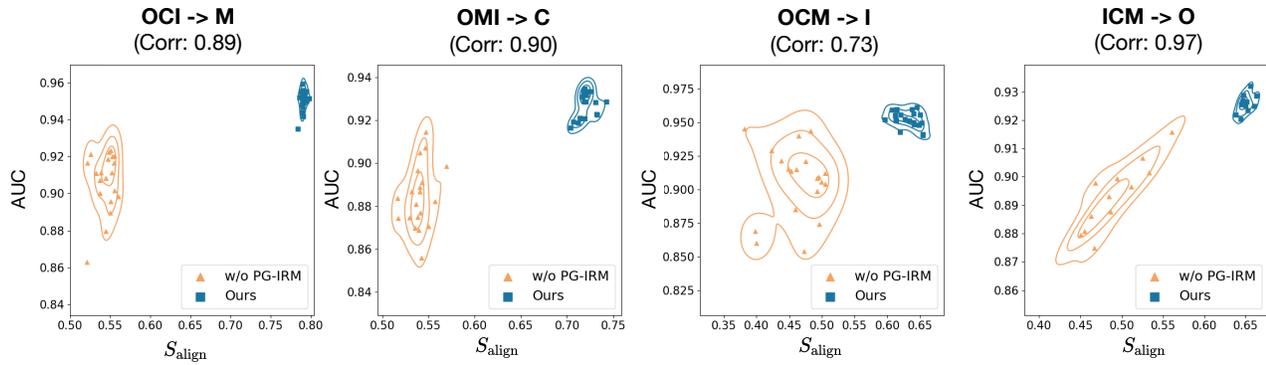
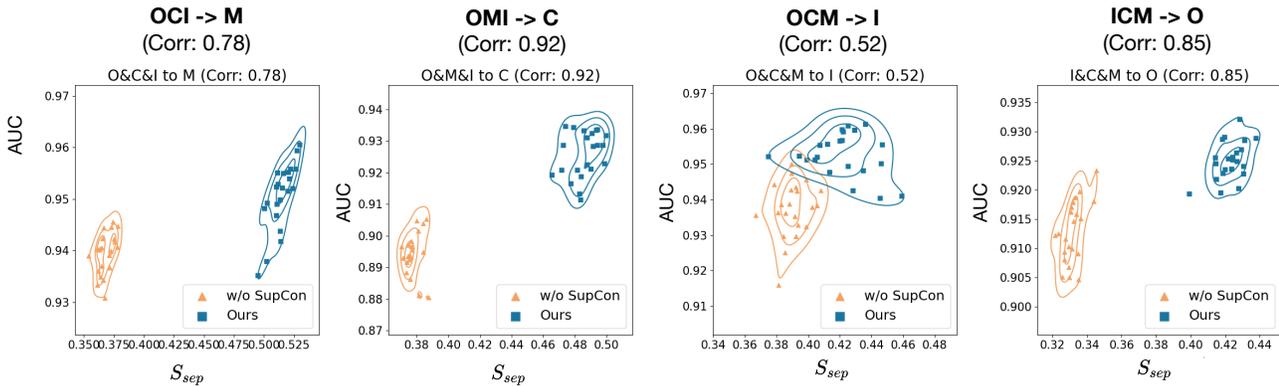


Figure 12. Sensitivity analysis of hyper-parameters: learning rate γ , alignment parameter α , alignment starting epoch T_a . The HTER is reported on the mean performance based on the last 10 epochs. The middle bar in each plot corresponds to the hyperparameter value used in our main experiments.



(a) Correlation Between AUC and Alignment



(b) Correlation Between AUC and Separability

Figure 13. Correlation between the test performance AUC and two properties measure. Each dot represents one snap-shot during the training stage in four cross-domain settings.