# Cross-Domain Identification for Thermal-to-Visible Face Recognition

Cedric Nimpa Fondje[1,*]     Shuowen Hu[2]     Nathaniel J. Short[2,3]     Benjamin S. Riggan[1,*]

[1] University of Nebraska-Lincoln, 1400 R St, Lincoln, NE 68588
[2] CCDC Army Research Laboratory, 2800 Powder Mill Rd., Adelphi, MD 20783
[3] Booz Allen Hamilton, 8283 Grennsboro Dr., McLean, VA 22102

*Corresponding authors: cedricnimpa@huskers.unl.edu, briggan2@unl.edu*

## Abstract

*Recent advances in domain adaptation, especially those applied to heterogeneous facial recognition, typically rely upon restrictive Euclidean loss functions (e.g., $L_2$ norm) which perform best when images from two different domains (e.g., visible and thermal) are co-registered and temporally synchronized. This paper proposes a novel domain adaptation framework that combines a new feature mapping sub-network with existing deep feature models, which are based on modified network architectures (e.g., VGG16 or Resnet50). This framework is optimized by introducing new cross-domain identity and domain invariance loss functions for thermal-to-visible face recognition, which alleviates the requirement for precisely co-registered and synchronized imagery. We provide extensive analysis of both features and loss functions used, and compare the proposed domain adaptation framework with state-of-the-art feature based domain adaptation models on a difficult dataset containing facial imagery collected at varying ranges, poses, and expressions. Moreover, we analyze the viability of the proposed framework for more challenging tasks, such as non-frontal thermal-to-visible face recognition.*

## 1. Introduction

Facial recognition (FR) systems have become more ubiquitous, being deployed across diverse services, including popular social media platforms, consumer devices (e.g., smart phones), and local/federal government or law enforcement databases (e.g., DoD Automated Biometric Identification System [4]). The algorithms/models in these FR systems are almost exclusively developed for matching visible spectrum facial imagery, due to the ubiquity of low-cost visible cameras with increasingly high resolution. Advances in visible spectrum FR algorithms, such as enhanced robustness to a varying pose, illumination, expression, resolution, and partial occlusion conditions, can be partially at-

tributed to the proliferation of large-scale FR datasets (e.g. LFW [11], MegaFace Challenge 1 [15], and MegaFace Challenge 2 [18]), computational resources (e.g., GPUs [15]), and deep learning models (e.g., convolutional neural networks or generative adversarial networks).

However, there has been growing interest in heterogeneous facial recognition (HFR), such as matching facial signatures in near infrared (NIR) [28, 16, 17] and thermal infrared [9, 7, 26] images to visible facial signatures, to better facilitate FR under low-light and variable illumination settings (e.g., nighttime FR). Unlike within spectrum matching (e.g., visible-to-visible FR), HFR—especially thermal-to-visible FR—has primarily been studied under significantly limited conditions (e.g., frontal images) due to the increased cost and complexity necessary to collect datasets containing visible and thermal face imagery that is on par with the size and number of conditions consistent with visible FR benchmarks.

Therefore, in this paper, we facilitate the advancement of thermal-to-visible FR research by introducing a new state-of-the-art domain adaptation framework and demonstrate its robustness on a difficult multi-modal face dataset, which incorporates more subjects and conditions than previous datasets. Since existing databases/watch-lists across academia, industry, government, and law enforcement currently enroll only visible face imagery, matching thermal face images with visible face imagery is necessary to perform FR in settings with low or variable light, which is significantly more difficult than conventional visible-to-visible FR in practice. This increased difficulty arises from significant differences between the thermal and visible face signatures [3].

Recently, several domain adaptive machine learning techniques have been used to reduce the modality gap between visible and thermal imagery. Hu et al. [9] demonstrated that using thermal cross-examples in a one-vs-all framework using partial least squares (PLS) classifiers enhanced discriminability when matching histogram of ori-

ented gradient (HOG) features between thermal and visible images. Sarfraz et al. [24] addressed the problem of thermal-to-visible FR using a small (shallow) neural network to perform a direct regression between local features extracted from thermal images and corresponding visible features that showed promising results, especially given the limited data used to train the network. Riggan et al. [20] used a coupled auto-associative model to learn a discriminative common latent subspace between visible and thermal image patches, and later [21] introduced an optimal feature regression and discriminative framework that exploited a combination of elements from [9], [24], and [20]. More recently, there have been domain adaptive techniques, such as [12, 13] that utilize a pre-trained deep learning model and domain adaptive convolutional neural networks (CNNs) to learn more global (contextual) domain-invariant features for thermal-to-visible FR. However, these techniques rely upon polarimetric thermal imagery and a relatively larger face crop size and are more sensitive to pose variations and temporal changes to the face (e.g., hair style, facial hair, aging).

The primary objective of this work is to enhance thermal-to-visible FR performance. Our contributions include:

- modified VGG16 [27] and Resnet50 [8] architectures for feature extraction,
- a new feature mapping sub-network to help bridge the domain gap,
- a new cross-domain identification loss function to relax requirement for precise co-registration and synchronization,
- a new domain invariance loss function provide a type of cross-domain regularization.

Compared to state-of-the-art feature-based domain adaptation methods [24, 21, 25], the proposed framework (Section 4) achieves enhanced conventional thermal-to-visible FR performance using an expanded multi-modal face dataset from the U.S. Army CCDC Army Research Laboratory [29], which contains frontal imagery (visible and polarimetric thermal) with neutral (baseline) and variable expressions. Moreover, we present new state-of-art results on a dataset that contains significantly more visible and thermal image pairs from a larger population of subjects, which is evaluated using imagery under varying pose and expression conditions. We intentionally exclude polarimetric thermal face signatures from our evaluations, since polarimetric thermal imaging is still an emerging area, whereas conventional thermal imagers have been widely deployed for military and homeland security applications, and is even becoming more widely available for commercial applications.

## 2. Heterogeneous Face Recognition

In this section, HFR models that have demonstrated some success are briefly reviewed so that we may compare and contrast these approaches with our proposed methodology (Section 4). In particular, we discuss two general approaches to HFR: transfer learning and domain adaptation.

### 2.1. Transfer Learning

Goodfellow et al. [5] describe transfer learning as a situation where knowledge learned in one setting is exploited to improve generalization in another setting. Often, the knowledge (features, weights, etc.) is transferred from previously trained models for one task to modified models for a different task in order to mitigate issues with over-fitting on the new task.

Intuitively, it seems that some shared information exists between corresponding visible and thermal face images, but the optimal form of knowledge transfer is not obvious. Instance, feature, parameter, and relational knowledge are four different types of knowledge that may be transferred between domains [19], where feature and parameter transfer learning are the most common.

The simplest form of parameter transfer learning includes weight sharing for deep neural networks, where usually the low-level parameters are assumed to be more applicable to both domains and are shared between two models. However, the high-level, more application specific, features are less generalizable and are fine-tuned to the desired task. More complex forms of parameter transfer learning, such as enforcing parameters to be related via a linear transform [23], may also be used. However, when domain gaps are very large (as with thermal and visible images), parameter constraints may be too restrictive to effectively learn a common representation for cross-domain FR. Therefore, we use domain adaptation.

### 2.2. Domain Adaptation

Unlike transfer learning, domain adaptation assumes a large shift between the source (image) distributions. The large perceptual gap between thermal and visible images is due to the fact that visible images are acquired using reflected light off of objects/faces and thermal images are acquired using "emitted" thermal radiation emanating from objects/faces. Sharing low-level weights may not be optimal since visible and thermal facial imagery exhibit a highly non-linear relationship and contain disparate information—thermal imagery has less high frequency and geometric details compared to visible imagery.

#### 2.2.1 Deep Perceptual Mapping

Sarfraz et al. [24] introduced the Deep Perceptual Mapping (DPM) specifically for thermal-to-visible FR. With DPM,

thermal-to-visible FR is assimilated as a regression problem, where a multilayer neural network directly regresses features (e.g., DSIFT or HOG) by minimizing

$$J_{dpm}(\boldsymbol{\Theta}) = \sum \|\mathbf{y} - f_{dpm}(\mathbf{x}; \boldsymbol{\Theta})\|^2, \qquad (1)$$

where $\mathbf{y}$ denotes a thermal feature vector, $f_{dpm}(\mathbf{x}; \boldsymbol{\Theta})$ denotes the DPM estimate for the thermal feature vector from a given visible feature vector ($x$), and $\boldsymbol{\Theta}$ is the set of trainable model parameters.

In [24], a three-layer DPM is optimized to predict thermal DSIFT features given the corresponding DSIFT features from the visible domain. Specifically, the authors show that extracting DSIFT features followed by principal components analysis (PCA) dimensionality reduction was effective for thermal-to-visible FR on the UND X1 database [2]. Image representations are constructed by concatenating local feature vectors extracted from overlapping image patches. A gallery is constructed from the estimated thermal representations of the visible gallery images, and matching is performed by computing the cosine similarity between the actual image representation from a thermal probe image and the predicted thermal feature representations from the gallery.

Similarly, our framework incorporates some compression in order to extract the most representative information from visible and thermal feature embedding representations. However, instead of DSIFT features, we integrate features extracted from deep neural networks that exhibit an effective number of features (i.e., channels), receptive fields size, and non-linearity, while also alleviating potential for over-fitting.

### 2.2.2 Coupled Neural Networks (CpNN)

An alternative approach for thermal-to-visible FR is a coupled neural network (CpNN), which learns how to extract common latent features between visible and thermal imagery by optimizing

$$J_{cpnn}(\boldsymbol{\Theta}_x, \boldsymbol{\Theta}_y) = \|f(\mathbf{x}, \boldsymbol{\Theta}_x) - g(\mathbf{y}; \boldsymbol{\Theta}_x)\|^2, \qquad (2)$$

where $f(\mathbf{x}, \boldsymbol{\Theta}_x)$ and $g(\mathbf{y}, \boldsymbol{\Theta}_y)$ are mapped visible and thermal features, respectively. Unlike DPM, this approach attempts to find two mappings such that features are sufficiently close in the mean square sense.

CpNNs train two encoders—one for the visible domain and one for the thermal domain—to produce similar features for corresponding inputs. The primary difference between CpNNs and DPMs is that CpNNs learn how to extract similar features, rather than extracting features from each domain and learning a mapping explicitly between these features [13].

Similar to [20, 21], [12] uses coupled networks to perform polarimetric thermal-to-visible FR. However, they use global average pooling with a VGG-like CNN architecture instead of local-features to form the common embedding representations. The motivation for global average pooling is to remove parameter intensive fully connected layers, which can potentially lead to over-fitting. However, one disadvantage is that global pooling assumes that the effective receptive field for the network is sufficiently large to provide enough contextual information to perform facial verification or identification, which requires additional layers and parameters. Moreover, in practice, it is generally easier to minimize the apparent modality gaps over relatively small image regions, but at the risk of needing more applications of local mappings to provide holistic image representation. Therefore, our proposed approach aims to learn local features with the largest context possible, but without the use of global pooling.

### 2.2.3 Generative Adversarial Network (GANs)

Generative Adversarial Networks [6] are composed of two parts: a generator ($G(\mathbf{z})$) and discriminator ($D(\mathbf{x})$), which are trained to minimize and maximize Eq. 3 with respect to $G$ and $D$, respectively.

$$\mathbb{E}_{\mathbf{x}}[\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z}}[1 - \log D(G(\mathbf{z}))]. \qquad (3)$$

The generator aims to confuse the discriminator by randomly synthesizing realistic samples from some underlying distribution. In the context of HFR, instead of generating random, realistic faces from a random vector, $\mathbf{z}$, conditional GANs (CGANs) are used to synthesize specific images from domain corresponding to a condition input, i.e., an image from another domain.

Similar to Eq. 3, CGANs optimize

$$\mathbb{E}_{\mathbf{x}}[\log D(\mathbf{x}|\mathbf{y})] + \mathbb{E}_{\mathbf{z}}[1 - \log D(G(\mathbf{z}|\mathbf{y}))], \qquad (4)$$

using the same minimax optimization used for GANs. Specifically for HFR, much like image-to-image translation techniques [14, 1] that perform image style transfer, [29] generates high-quality visible face images from given polarimetric thermal face imagery. The disadvantage to this type of approach is computational complexity of the model. In [22], a similar network reported an average run time of more than 270 seconds per image. Therefore, such a computationally complex approach is not suitable for real-time FR, but may be a useful tool for analyst and less time critical activities.

Even GAN-based architectures, which use encoder-decoder networks for the generators, rely on robust, discriminative embedding representations. For example, in addition to the adversarial loss (e.g., Eq. 4), [29] uses Euclidean loss functions between visible and thermal features to enforce similarity between feature embeddings. This not only ensures photo-realism at the generators output, but also

similar generated visible latent features from a given thermal image. Therefore, we expect that our proposed approach could also be used as a complementary loss to provide additional guidance for GANs. However, this is beyond the scope of this paper.

## 3. Problem Definition

Let $\mathbf{V}$ and $\mathbf{T}$ each be the set of $n$-dimensional descriptive feature vectors from the visible domain and the thermal domain, respectively. These vectors may be considered as either hand-crafted features (e.g., DSIFT) or features from trainable neural networks. Let $\mathbf{v_i} \in \mathbf{V}$ and $\mathbf{t_i} \in \mathbf{T}$, where the index $i$ denotes a pair of images $\mathbf{p}_i = (\mathbf{v_i}, \mathbf{t_i})$ corresponding to the same subject. Without loss of generality, the pairs need not be precisely co-registered or synchronized.

Given a set of training pairs with known identities (i.e., labels), $y$, our goal is to find a mapping $\mathbf{f}_t : \mathbf{V} \to \mathbf{T}$ (or $\mathbf{f}_v : \mathbf{T} \to \mathbf{V}$) such that: $\mathbf{f}_t(\mathbf{v}_i) \approx \mathbf{t}_i$ (or $\mathbf{f}_v(\mathbf{t}_i) \approx \mathbf{v}_i$). In either scenario, the visible representations are used as the reference in the gallery and the thermal representations are the test samples, which means both forms may be considered thermal-to-visible FR. The key difference is whether the mapping is applied during enrollment or matching.

During deployment (i.e., testing), only the identities corresponding to the visible representations of enrolled subjects (i.e., gallery) are known. Thus, we aim to learn a generalizable mapping $\mathbf{f}_t$ (or $\mathbf{f}_v$) that optimally discriminates genuine pairs and imposters pairs for the purposes of thermal-to-visible face identification.

The problem with thermal-to-thermal (or visible-to-visible) FR is that the features, $\mathbf{t}$ (or $\mathbf{v}$), and classifiers, $\hat{\mathbf{y}}(\mathbf{t}; \boldsymbol{\Theta}_t)$ (or $\hat{\mathbf{y}}(\mathbf{v}; \boldsymbol{\Theta}_v)$), rely on distinct phenomenology. The within-spectrum classifiers are learned by minimizing the cross-entropy,

$$\mathcal{L}(\boldsymbol{\Theta}_x) = -\sum \mathbf{y} \log \hat{\mathbf{y}}(\mathbf{x}; \boldsymbol{\Theta}_x), \qquad (5)$$

between labels and predictions. In Eq. 5, $\mathbf{x}$ denotes either $\mathbf{t}$ or $\mathbf{v}$, depending on the scenario. Therefore, in principle, we also want to learn visible features that are optimally discriminative when fed into a thermal classifier (or vice versa). So, in addition to Eq. 5, we also want to minimize

$$\mathcal{L}(\boldsymbol{\Phi}_{x'}) = -\sum \mathbf{y} \log \hat{\mathbf{y}}(\mathbf{f}_x(\mathbf{x}'; \boldsymbol{\Phi}_{x'}); \boldsymbol{\Theta}_x), \qquad (6)$$

where $\mathbf{x}'$ denotes features from the opposite domain of $\mathbf{x}$ (e.g., if $\mathbf{x} = \mathbf{t}$ then $\mathbf{x}' = \mathbf{v}$). $\boldsymbol{\Phi}_{x'}$ denotes the trainable parameters for the mapping $\mathbf{f}_x(\cdot)$.

In following sections, we assume that $\mathbf{x} = \mathbf{t}$ then $\mathbf{x}' = \mathbf{v}$ to describe our approach.

## 4. Proposed Domain Adaptation Framework

Our proposed domain adaptation framework (Figure 1) is composed of four main components:

- feature extraction using truncated deep neural network for both visible and thermal imagery (Section 4.1),
- our proposed Residual Spectral Transform (RST) that bridges the remaining gap between visible and thermal features (Section 4.2),
- our proposed cross-domain identification loss that enhances holistic discriminability, especially when matching visible and thermal image representations (Section 4.3),
- our proposed domain invariance loss that discourages domain predictability from visible or thermal features (Section 4.4).

Each of these components are discussed in depth in the following subsections.

### 4.1. Feature Extraction

Feature representations, $\mathbf{v}$ and $\mathbf{t}$, are initially extracted from images using pre-trained neural networks, such as VGG16 and Resnet50 architecture. However, since high-level features tend to be less transferable (i.e., high frequency details in visible spectrum facial imagery are lacking in the inherently smoother thermal facial imagery), we intentionally truncate the networks at the optimal depth in order to simultaneously obtain the largest receptive field and the most similar feature response.

For the VGG16 network, we determined experimentally that the output of the third convolution-pooling block provided the most discriminative visible and thermal information for thermal-to-visible FR. Similarly, we found that the output of third residual block for the Resnet50 architecture provided the most robust thermal-to-visible FR performance. Interestingly, the layers in which we truncate both of these networks result in feature maps with spatial dimensions of $25 \times 25$ ($H \times W$). The feature maps do, however, have different number of channels ($C$); VGG16 has 256 channels and ResNet50 has 512 channels. Similar to [24], we include a compression layer that is shared between visible and thermal representations. This not only reduces the number of parameters in the subsequent layers, but also helps to eliminate factors associated with noise. The complete experiment details and results are provided in section 5.2.

### 4.2. Residual Spectral Transform

After extracting the most similar features possible, there is still a significant domain gap between the visible and thermal representations. Therefore, we augment the visible (or thermal) network with our proposed Residual Spectral Transform (RST), which is the sub-network shown in Figure 1.

The RST is a residual block that aims to preserve as much discriminabilty from the truncated networks as possible while transforming features between thermal and visible
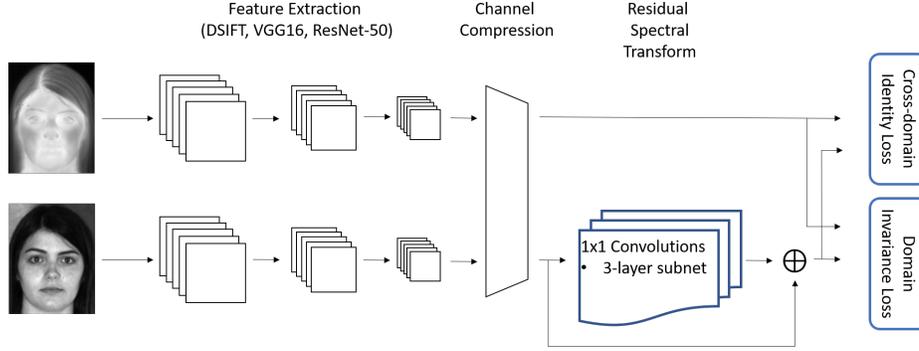
Figure 1. Proposed domain adaptation model.

domains. This sub-network transforms the features using three 1x1 convolutional layers:

$$\mathcal{F}(\mathbf{u}) = \text{Conv}_c \circ \text{Conv}_{200} \circ \text{Conv}_{200}(\mathbf{u}), \qquad (7)$$

where $\text{Conv}_k(\cdot)$ denotes the 1x1 convolutional layers with $k$ units and uses hyperbolic tangent activation function, and $c$ denotes the number of channels which depends on the selected network architecture or features. Then, the RST is given by

$$RST(\mathbf{u}) = \mathcal{F}(\mathbf{u}) + \mathbf{u}. \qquad (8)$$

### 4.3. Cross-domain Identification Loss

To preserve the identity of the images during FR, we use a cross-domain identification loss which is based on combining (5) and (6). Therefore, our proposed loss function aims to simultaneously learn a discriminative thermal network (features and classifier) and optimal visible-to-thermal RST by minimizing

$$\mathcal{L}_{xID} = -\sum \left[ \mathbf{y} \log \left\{ \hat{\mathbf{y}}(\mathbf{t}; \boldsymbol{\Theta}_t) \cdot \hat{\mathbf{y}}(\hat{\mathbf{t}}; \boldsymbol{\Theta}_t) \right\} \right], \quad (9)$$

where $\hat{\mathbf{t}} = \mathbf{f}_t(\mathbf{v}; \Phi_v)$.

Eq. 9 is minimized by alternately optimizing parameters $\boldsymbol{\Theta}_t$ and $\boldsymbol{\Phi}_v$, which learns discriminative information for the thermal domain while also performing cross-domain identification.

### 4.4. Domain Invariance Loss

While Eq. 9 provides discriminative information, it is still possible to over-fit, resulting in inconsistent performance between training and inference. In order to provide some regularization for better stability, we add a parallel domain classifier that estimates the probabilities that a given image representation comes from a visible image or a thermal image. However, since our ultimate goal is to bridge the representational gap between visible and thermal imagery, optimal domain discriminability is not desirable because this would imply that there is still a detectable difference

between visible and thermal feature representations. Therefore, we add a regularizing loss function that encourages gradient-driven parameter updates to maintain poor domain predictability.

Let a domain detector be denoted as $\mathcal{D}(\cdot)$, which produces two probabilities, $P(vis|\mathbf{h})$ and $P(thm|\mathbf{h})$, which are the probabilities that the feature representation $\mathbf{h}$ is computed from a visible ($vis$) or thermal ($thm$) image, respectively. Intuitively, we prefer to learn visible and thermal image representations such that $\mathcal{D}$ produces $P(vis|\mathbf{h}) \approx P(thm|\mathbf{h})$ for any $\mathbf{h} \in \{\mathbf{T}, \mathbf{V}\}$. Therefore, we add the following loss function to our proposed domain adaptation framework:

$$\mathcal{L}_{\mathcal{D}} = -\sum \alpha \left\{ \log \mathcal{D}(\mathbf{t}) + \log \mathcal{D}(\hat{\mathbf{t}}) \right\}, \qquad (10)$$

where $\alpha$ denotes the target probabilities. In this case, $\alpha \in \mathbb{R}^2$ should be an approximately uniform distribution (i.e., all elements are close to $0.5$), which encourages poor domain prediction.

Therefore, the total loss function combines Eqs. 9 and 10, which yields

$$\mathcal{L}_{total} = (1 - \lambda)\mathcal{L}_{xID} + \lambda\mathcal{L}_{\mathcal{D}}, \qquad (11)$$

### 4.5. Implementation Details

Pre-processing for thermal-to-visible FR applications typically includes image registration of corresponding visible and thermal images, image filtering, and cropping. Using fiducial landmarks, including the center of eyes, base of nose, and mouth corners (i.e., 5-point alignment), the facial images are aligned to canonical coordinates using a similarity transform. In practice, the landmarks can be automatically detected. However, the datasets (section 5) used in this study include manually annotated landmarks for every image.

Next, a Difference of Gaussians (DoG) filter is applied to the registered images (visible and thermal). DoG filtering enhances the common edges between visible and thermal

Table 1. Feature ablation study on protocol 1 showing Rank-1 identification (ID) rate and feature map dimensions (Dims.)

| Method | Rank-1 ID (%) | Feature Map Dims. ($H \times W \times C$) |
|---|---|---|
| DOG+image_patch | 2.83 | $24 \times 24 \times 400$ |
| DSIFT | 8.33 | $24 \times 24 \times 128$ |
| DOG+DSIFT | **10.50** | $\mathbf{24 \times 24 \times 128}$ |
| DOG+VGG16(block5) + Avg. Pool | 12.67 | $1 \times 1 \times 512$ |
| DOG+VGG16(block5) | 15.17 | $6 \times 6 \times 512$ |
| DOG+VGG16(block4) | 46.67 | $12 \times 12 \times 512$ |
| DOG+VGG16(block3) | **58.83** | $\mathbf{25 \times 25 \times 256}$ |
| DOG+VGG16(block2) | 23.17 | $50 \times 50 \times 128$ |
| DOG+Resnet50 (5c) + Avg. Pool | 9.00 | $1 \times 1 \times 2048$ |
| DOG+Resnet50 (5c) | 10.00 | $7 \times 7 \times 2048$ |
| DOG+Resnet50 (4f) | 19.50 | $13 \times 13 \times 1024$ |
| DOG+Resnet50 (3d) | **70.83** | $\mathbf{25 \times 25 \times 512}$ |
| DOG+Resnet50 (2c) | 33.17 | $50 \times 50 \times 256$ |

imagery, and has been shown to help for thermal-to-visible matching [9].

After filtering, all the images are cropped around the eyes, nose, and mouth to a $200 \times 200$ pixel image, in order to be consistent with past studies [21] and to limit computational and memory requirements. Additionally, this "tight" crop ensures that the network is attributing the identification with facial features rather than factors that may change frequently over time, such as hair style. "Better" performance on benchmarks may be achieved when using a larger crop. However, larger crops may under-perform when significant variations, such as changes in hair style or image background, are observed.

When training our framework—the truncated VGG16 and Resnet50 architectures with our proposed RST—using the proposed cross-domain identification and domain invariance loss functions, $\lambda$ is set to 0.25. Also, we used a compression ratio of 50%—meaning the dimensionality is reduced by half. Our compression ratio is consistent with the number of principal components used with DPM [24, 25]. For reproducibility, our code can be found on https://git.unl.edu/ece-unl-images-lab/cross-domain-identification.

# 5. Experiments & Results

In this section, we first describe the protocols for our experimental analysis. Then, we describe an ablation study used to determine how to best truncate the feature extraction models for thermal-to-visible FR. Next, we provide comprehensive analysis using two different protocols, including new analysis that examines non-frontal to frontal matching for thermal-to-visible FR. Lastly, we perform a second ablation study that examines the impact with and without the proposed domain invariance loss function.
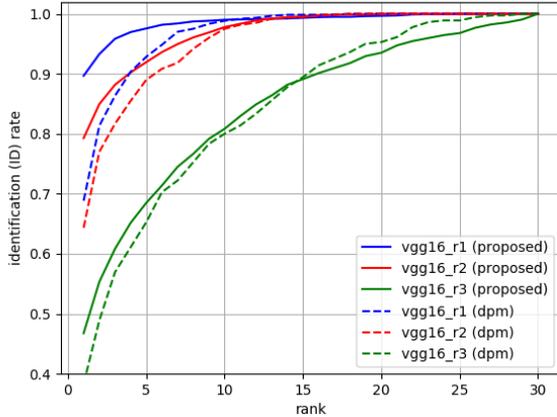
## 5.1. Protocols

For experimental analysis, we use three separate datasets/protocols compiled by the CCDC Army Research Laboratory.
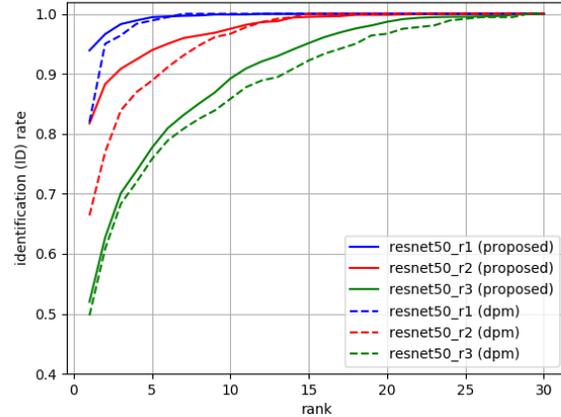
The first dataset [10], referred to as "protocol 1" in [29], contains 2880 corresponding visible and polarimetric thermal image pairs from 60 unique subjects. This collection contains imagery acquired with neutral expressions (baseline) and varying expressions at multiple standoff distances: $2.5m$ (r1), $5.0m$ (r2), and $7.5m$ (r3). Using protocol 1, we train our proposed domain adaptation framework using only imagery at r1 (without any data augmentation) from 30 different subjects, and we evaluate using a 30 subject gallery with 4 baseline visible images per subjects (consistent with [10, 21]).

The second dataset, referred to as "protocol 2" in [29], augments protocol 1 with an additional 1581 visible and polarimetric thermal image pairs from 51 subjects. Therefore, protocol 2 contains 4461 image pairs from 111 subjects. The additional 51 subject collection includes both baseline and varying expressions, but only acquired imagery at a single range equivalent to r1. Complete details regarding this collection, including how to obtain the data, can be found in [10, 21]. Using protocol 2, we train our proposed domain adaptation framework using only imagery at r1 (without any data augmentation) from 81 different subjects, and we evaluate using a 30 subject gallery with 4 baseline visible images per subjects (consistent with [10, 21]).

The third dataset, which we refer to as "protocol 3," contains visible and thermal imagery from a separate collection (or volume) of 126 subjects. There are a total of 5176 images (2198 visible and 2978 thermal) available from a single range consistent with r1. However, unlike protocol 2, this collection contains faces with neutral expressions (baseline), varying expression, and varying pose (yaw in the range of $\pm60°$). Using protocol 3, we train our pro-

(a)



(b)

posed domain adaptation framework (without any data augmentation) using imagery from 96 different subjects, and we evaluate using a 30 subject gallery with 4 frontal, visible images per subject. We intentionally use the same size gallery as protocol 2 for comparability.

## 5.2. The Feature Ablation Study

In this section, using protocol 1, we analyze the impact on thermal-to-visible FR for different features: image patches, DSIFT, VGG16, and Resnet50. For VGG16 and Resnet50, we explore extracting features at various layers to determine optimal feature map size for the respective architectures. All features are extracted from the $200 \times 200$ visible and thermal images and compressed using principal component analysis (PCA) to 64 principal components.

The VGG16 architecture is composed of five convolution-pooling blocks (block1, ..., block5) and several fully connected layers. The Resnet50 architecture is composed of one convolution layer and four residual blocks, and each residual block is composed of multiple residual layers. For example, residual block 2 contains three residual layers (2a, 2b, 2c), residual block 3 contains four residual layers (3a, 3b, 3c, 3d), residual block 4 contains six residual layers (4a, 4b, 4c, 4d, 4e, 4f), and residual block 5 contains three residual layers (5a, 5b, 5c). The final output (prior to the classifier) is average pooling layer that reduces the spatial dimensions of the features maps to $1 \times 1$.

Table 1 shows the rank-1 identification (ID) rate when matching thermal probe image representations with visible gallery image representations (from protocol 1) and the feature map dimensions (prior to compression). These image representations are matched using the cosine similarity measure. This table shows that better thermal-to-visible

FR is achieved when extracting features from intermediate layers of the pre-trained VGG16 and Resnet50 networks, i.e., block3 and 3d, respectively. This appears to indicate that extracting high-level features from the pre-trained networks depends too much on visible texture to perform well for cross-domain matching. Also, lower level features have too small of a receptive field and insufficient context to perform cross-domain matching. Given the aligned and cropped images, the best network performances seem to be attained when the feature maps' spatial dimensions are approximately $25 \times 25$.

## 5.3. Proposed Framework Versus DPM

Using both truncated VGG16 and Resnet50 architectures (Section 5.2), we compared our proposed RST method (trained using Eq. 11) with DPM (trained using Eq. 1). Figure 5 shows the cumulative match characteristic (CMC) curves for protocol 2. First, note that our proposed domain adaptation framework exceeds the performance of DPM across all ranges and conditions when using both the truncated VGG16 and Resnet50 networks. Secondly, it is important to note that DSIFT + DPM [21] achieves 84%, 75%, 58% performance for r1, r2, and r3 and our proposed framework achieves 94.2%, 81.7%, and 52.0%. Thus, our approach improves performance for r1 and r2, but r3 performance is somewhat lacking still. This may be due to either a trade-off between DSIFT and neural networks based features (which can be more sensitive to low-quality images) or the fact that we only train using r1 imagery. Thus, it may be possible to see gains by incorporating multiple resolutions or data augmentation during training in order to boost the r3 performance.

We also compared the DPM and our proposed domain adaptation framework on protocol 3, which includes more

subject and more variations (particularly pose). Table 2 shows the rank-1 ID for the truncated VGG16 and Resnet50 architectures under both variable expression and pose conditions. It is important to note that the gallery is composed of only frontal imagery. Unsurprisingly, the performance drops when matching non-frontal thermal imagery to the frontal visible gallery. However, our proposed framework still achieves better rank-1 performance than DPM.

Table 2. Protocol 3 Rank-1 ID performance for pose and expression variations

| Condition | Method | Rank-1 ID (%) |
|---|---|---|
| Expression | Resnet50 + DPM | 91.56 |
| | VGG16 + DPM | 82.29 |
| | Resnet50 + Proposed | **96.00** |
| | VGG16 + Proposed | **84.00** |
| Pose | Resnet50 + DPM | 24.42 |
| | VGG16 + DPM | 21.33 |
| | Resnet50 + Proposed | **29.91** |
| | VGG16 + Proposed | **21.38** |

### 5.4. Effect of Domain Invariance Loss

This ablation study aims to empirically assess the effect of the domain invariance loss function (Eq. 10). Using protocol 2, we consider our proposed model for scenario 1—finding a mapping $\mathbf{f}_t : \mathbf{V} \rightarrow \mathbf{T}$—and scenario 2—finding a mapping $\mathbf{f}_v : \mathbf{T} \rightarrow \mathbf{V}$. In each scenario, we compare the rank-1 performance at r1 both with ($\lambda = 0.25$) and without ($\lambda = 0$) the domain invariance loss. Additionally, we also consider both the truncated VGG16 and Resnet50 models.

Table 3 shows that using VGG16 with the domain invariance loss improves performance by 17.50% and 31.33% under scenario 1 and scenario 2, respectively. Also, the results show that using Resnet50 with the domain invariance loss improves performance by 5.84% and 1.66% under scenario 1 and scenario 2, respectively. This ablation study demonstrates that the domain invariance loss is an important aspect of the proposed framework.

Table 3. Effect of domain invariance loss for scenario 1 and scenario 2

| Scenario | Domain Invariance | Resnet50 | VGG16 |
|---|---|---|---|
| Scenario 1 | no | 88.33% | 65.83% |
| | yes | 94.17% | 83.33% |
| Scenario 2 | no | 89.17% | 49.17% |
| | yes | 90.83% | 80.50% |

## 6. Conclusion

In this paper, we proposed a new domain adaptation framework, which used truncated deep neural network with a new RST sub-network. This framework was trained using our new cross-domain identification and domain invariance loss function. Compared with DPM, which is trained using a Euclidean loss function between thermal and visible features, our framework shows significant improvements across multiple ranges, poses, and expressions. Also, we demonstrated significant improvement to state-of-the-art features extraction methods, like DSIFT and DPM, and show that our domain invariance loss function plays an important role in achieving robust thermal-to-visible FR.

Most importantly, we have developed a framework that alleviates the restrictive need for precisely registered and synchronized imagery because we introduce two loss functions that operated at the task-level rather than the feature-level. We hypothesize with larger datasets (on the order of 300+ subjects and 500,000+ thermal-visible pairs) will provide further generalization to the challenging non-frontal to frontal matching for thermal-to-visible FR.

## 7. Acknowledgments

## References

[1] D. Berthelot, T. Schumm, and L. Metz. BEGAN: Boundary equilibrium generative adversarial networks. *arXiv preprint arXiv:1703.10717*, 2017. 3

[2] X. Chen, P. J. Flynn, and K. W. Bowyer. Ir and visible light face recognition. *,Comput. Vis. Image Und.*, 99:332–358, 2005. 3

[3] J. Choi, S. Hu, S. S. Young, and L. S. Davis. Thermal to visible face recognition. In S. O. Southern, A. H. J. Kolk, K. N. Montgomery, C. W. Taylor, B. V. K. V. Kumar, S. Prabhakar, and A. A. Ross, editors, *Sensing Technologies for Global Health, Military Medicine, Disaster Response, and Environmental Monitoring II; and Biometric Technology for Human Identification IX*, volume 8371, pages 252 – 261. International Society for Optics and Photonics, SPIE, 2012. 1

[4] J. M. Gilmore. Department of Defense (DOD) Automated Biometric Identification System (ABIS) version 1.2. Technical Report 16-F-0250, Operational Test & Evaluation (OT&E), May 2015. 1

[5] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016. http://www.deeplearningbook.org. 2

[6] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Gen-

erative adversarial nets. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*, NIPS14, page 26722680, Cambridge, MA, USA, 2014. MIT Press. 3

[7] K. P. Gurton, A. J. Yuffa, and G. W. Videen. Enhanced facial recognition for thermal imagery using polarimetric imaging. *Opt. Lett.*, 39(13):3857–3859, Jul 2014. 1

[8] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. 2

[9] S. Hu, J. Choi, A. L. Chan, and W. R. Schwartz. Thermal-to-visible face recognition using partial least squares. *J. Opt. Soc. Am. A*, 32(3):431–442, Mar 2015. 1, 2, 6

[10] S. Hu, N. J. Short, B. S. Riggan, C. Gordon, K. P. Gurton, M. Thielke, P. Gurram, and A. L. Chan. A polarimetric thermal database for face recognition research. In *2016 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 187–194, 2016. 6

[11] G. B. Huang, M. Ramesh, T. Berg, and E. Learned-Miller. Labeled faces in the wild: A database for studying face recognition in unconstrained environments. Technical Report 07-49, University of Massachusetts, Amherst, October 2007. 1

[12] S. M. Iranmanesh, A. Dabouei, H. Kazemi, and N. M. Nasrabadi. Deep cross polarimetric thermal-to-visible face recognition. In *2018 International Conference on Biometrics (ICB)*, pages 166–173, 2018. 2, 3

[13] S. M. Iranmanesh and N. M. Nasrabadi. Attribute-guided deep polarimetric thermal-to-visible face recognition. In *2019 International Conference on Biometrics (ICB)*, pages 1–8, 2019. 2

[14] P. Isola, J. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5967–5976, 2017. 3

[15] I. Kemelmacher-Shlizerman, S. M. Seitz, D. Miller, and E. Brossard. The megaface benchmark: 1 million faces for recognition at scale. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4873–4882, 2016. 1

[16] B. Klare and A. K. Jain. Heterogeneous face recognition: Matching nir to visible light images. In *2010 20th International Conference on Pattern Recognition*, pages 1513–1516, 2010. 1

[17] S. Z. Li, D. Yi, Z. Lei, and S. Liao. The casia nir-vis 2.0 face database. In *2013 IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 348–353, 2013. 1

[18] A. Nech and I. Kemelmacher-Shlizerman. Level playing field for million scale face recognition. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3406–3415, 2017. 1

[19] S. J. Pan and Q. Yang. A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(10):1345–1359, 2010. 2

[20] B. S. Riggan, C. Reale, and N. M. Nasrabadi. Coupled auto-associative neural networks for heterogeneous face recognition. *IEEE Access*, 3:1620–1632, 2015. 2, 3

[21] B. S. Riggan, N. J. Short, and S. Hu. Optimal feature learning and discriminative framework for polarimetric thermal to visible face recognition. In *2016 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1–7, 2016. 2, 3, 6, 7

[22] B. S. Riggan, N. J. Short, M. S. Sarfraz, S. Hu, H. Zhang, V. M. Patel, S. Rasnayaka, J. Li, T. Sim, S. M. Iranmanesh, and N. M. Nasrabadi. Icme grand challenge results on heterogeneous face recognition: Polarimetric thermal-to-visible matching. In *2018 IEEE International Conference on Multimedia Expo Workshops (ICMEW)*, pages 1–4, 2018. 3

[23] A. Rozantsev, M. Salzmann, and P. Fua. Beyond sharing weights for deep domain adaptation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(4):801–814, 2019. 2

[24] M. S. Sarfraz and R. Stiefelhagen. Deep perceptual mapping for thermal to visible face recognition. In X. Xie, M. W. Jones, and G. K. L. Tam, editors, *Proceedings of the British Machine Vision Conference 2015, BMVC 2015, Swansea, UK, September 7-10, 2015*, pages 9.1–9.11. BMVA Press, 2015. 2, 3, 4, 6

[25] M. S. Sarfraz and R. Stiefelhagen. Deep perceptual mapping for cross-modal face recognition. *International Journal of Computer Vision*, 122(3):426–438, 2017. 2, 6

[26] N. Short, S. Hu, P. Gurram, K. Gurton, and A. Chan. Improving cross-modal face recognition using polarimetric imaging. *Opt. Lett.*, 40(6):882–885, Mar 2015. 1

[27] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015. 2

[28] D. Yi, R. Liu, R. Chu, Z. Lei, and S. Z. Li. Face matching between near infrared and visible light images. In S.-W. Lee and S. Z. Li, editors, *Advances in Biometrics*, pages 523–530, Berlin, Heidelberg, 2007. Springer Berlin Heidelberg. 1

[29] H. Zhang, B. Riggan, S. Hu, N. Short, and V. Patel. Synthesis of high-quality visible faces from polarimetric thermal faces using generative adversarial networks. *International Journal of Computer Vision*, 127:1–18, 03 2019. 2, 3, 6