

# A Generative Approach for Image Registration of Visible-Thermal (VT) Cancer Faces

Catherine Ordun<sup>1,2</sup>, Alexandra Cha<sup>1</sup>, Edward Raff<sup>1,2</sup>, Sanjay Purushotham<sup>1,2</sup>,  
Karen Kwok<sup>3</sup>, Mason Rule<sup>3</sup>, and James Gulley<sup>3</sup>

<sup>1</sup> Booz Allen Hamilton, Washington, DC USA  
{cha\_alexandra, raff\_edward}@bah.com

<sup>2</sup> University of Maryland Baltimore County, Baltimore, MD USA  
{cordun1, psanjay}@umbc.edu

<sup>3</sup> National Cancer Institute, National Institutes of Health, Rockville, MD USA  
{karen.kwok, mason.rule, gulleyj}@nih.gov

**Abstract.** Since thermal imagery offers a unique modality to investigate pain, the U.S. National Institutes of Health (NIH) has collected a large and diverse set of cancer patient facial thermograms for AI-based pain research. However, differing angles from camera capture between thermal and visible sensors has led to misalignment between Visible-Thermal (VT) images. We modernize the classic computer vision task of image registration by applying and modifying a generative alignment algorithm to register VT cancer faces, without the need for a reference or alignment parameters. By registering VT faces, we demonstrate that the quality of thermal images produced in the generative AI downstream task of Visible-to-Thermal (V2T) image translation significantly improves up to 52.5%, than without registration. *Images in this paper have been approved by the NIH NCI for public dissemination.*

**Keywords:** Chronic Pain · Generative AI · Generative Adversarial Network · Diffusion Models · Visible-Thermal Image Translation

## 1 Introduction

Thermal temperatures have been shown to directly correlate with gold standard vital metrics that measure physiological excitement such as electrocardiography (ECG) and galvanic skin response (GSR), making it a promising non-invasive measure for pain assessment [7,18,20,22]. Further, thermograms reveal signs of pain by detecting temperature variations in inflamed tissue such as heat patterns in the supraorbital and periorbital facial areas at the onset of a migraine attack [19], and mean temperature elevation in swollen and tender body regions when analyzing arthritis, tennis elbow, and fibromyalgia [23]. As a result, the Intelligent Sight and Sound (ISS) [15] study under the U.S. National Institutes of Health (NIH), National Cancer Institute (NCI), has collected thermal images of cancer patient faces since October 2020. The study’s objective is to assess the use of AI as a tool to detect chronic cancer pain - an understudied problem in



Fig. 1: ISS VT Facial Pairs Before and After Alignment using a Generative Approach for Image Registration for Six Patients.

machine learning and thermal physiology. For many AI tasks, such as generative modeling, the thermal image must be paired with its corresponding visible image in order to extract multimodal signals. Further, the pairing must be exact since shifts in pixels can lead to significant changes in model prediction and generative AI image quality [11,28]. Multi-spectral image alignment between thermal and visible images is a non-trivial task, that is further exacerbated when a reference for scale, such as deformation parameters, is not available. To address this problem, we train and modify Vista Morph [16], a generative multi-spectral image registration framework, on the ISS VT Facial Dataset of cancer patients. Vista Morph predicts the affine matrix parameters used to align thousands of thermal images, respective to the geometry of their visible pairs, even under extremely warped conditions, leading to registered VT facial pairs shown in Figure 1. In addition, automatic registration saves significant human labor otherwise spent manually aligning images by hand, or relying on proprietary, thermal landmark libraries. This paper offers multiple contributions: 1) We introduce the ISS VT Facial Dataset, which is to the best of our knowledge, the largest cancer VT facial dataset available, useful for thermal physiology, thermal facial recognition, and generative AI studies consisting of 29,461 VT facial pairs. 2) We train and modify the generative Vista Morph image registration algorithm to correct for challenging, severe misalignment between VT faces. 3) We show that registered images generate improved quality of thermal faces than without registration using GANs, as a method of synthetic data augmentation.

## 2 Methodology

### 2.1 ISS VT Facial Dataset

**Capture Protocol.** VT facial images are extracted from videos of cancer patients from the Intelligent Sight and Sound (ISS) Dataset [15] - an ongoing, observational, non-interventional clinical study initiated by the U.S. National Institutes of Health (NIH) National Cancer Institute (NCI). The study is entitled “*A Feasibility Study Investigating the Use of Machine Learning to Analyze*

*Facial Imaging, Voice and Spoken Language for the Capture and Classification of Cancer Pain*" [14].

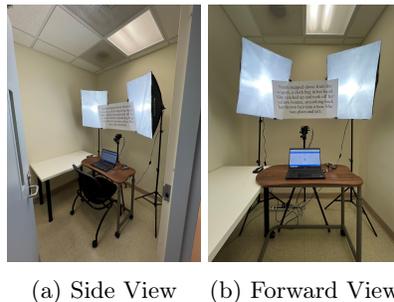


Fig. 2: Experimental Setup to collect VT facial pairs at the NIH clinic.

Videos of participant faces are captured both remotely, using a smartphone, and in the NIH clinic. Smartphone video provides visible spectrum data while the VT pairs are only captured at the NIH clinic. In-clinic, two separate cameras capture video in the visible ( $0.4\text{-}0.7\ \mu\text{m}$ , NexiGo PC Webcam) and thermal ( $7.5\text{-}13.5\ \mu\text{m}$ , Flir Duo Pro R) spectra as shown in Figure 2. Both cameras are positioned in a vertical stack, to minimize parallax effects, and erected on a tripod less than two feet away from the patient in a frontal position. Only the patient and the cameras are present in the room, no other personnel are present during recording. The protocol begins when the patient reads a 10-15 second-long passage selected from a grade 3 reading level text, displayed on the white poster in Figure 2. This is a neutralizing prompt, and is common practice in mood conditioning trials in order to control for an emotionally charged response [1,5,12]. Next, the patient records a video response to the prompt *"Please describe how you feel right now,"* designed to capture narratives about the patient's mood, beliefs and attitudes about their pain, medical conditions and daily activities.

**Data Processing.** From 96 VT video recordings, we extract 29,461 VT pairs across 44 subjects using minimal preprocessing. First, we extract frames from VT videos at a rate of 10 fps using the `ffmpeg` library. Second, we detect and crop faces using a Multi-task Cascaded Convolutional Network (MTCNN) [27]. Lastly, we apply a series of binary thresholding operations to crop the thermal face from the background using `OpenCV`. This crude processing leads to severely misaligned VT facial pairs, as shown in samples per Figure 1 before registration. The incidence of older aged individuals who tend to look down while speaking, as opposed to maintaining frontal head alignment, in addition to low resolution from both sensors, and dynamic-recorded video sessions, makes the ISS VT Facial Dataset more challenging for alignment than other VT facial datasets that are captured in standard settings (i.e. background, prompted poses) such as ARL Devcom [21], Eurecom [13], and Carl [4] that use similar VoX microbolometer thermal sensors.



Fig. 3: Distribution of Pain Classes in the NIH ISS VT Facial Dataset.

**Data Overview.** Samples are shown in Figure 1. We split subjects into non-overlapping train and test sets, so that subjects are seen only once in either group. There are 22,570 VT image pairs in the training set and 6,891 VT image pairs in the test set. The training set consists of 34 subjects (67% male, 33% female) extracted from 75 videos representing 37.6 minutes of video. Shown in Figure 3, approximately 38% of train subjects have No Pain (1), 14% have Low Pain (2), 24% have Moderate Pain (3) and 24% with Severe Pain (4). The test set consists of 10 test subjects (40% male, 60% female) extracted from 21 videos, representing 11.5 minutes of video. In the test set, approximately 50% of the subjects have Low Pain (2), 40% have Severe Pain (4), and 10% have No Pain (1). No test subjects represent Moderate Pain (3) as shown in Figure 3.

## 2.2 Vista Morph Framework

We train a deep generative image registration algorithm called Vista Morph [16] on the ISS VT Facial Dataset, shown in Figure 4a. This framework is specifically designed for VT facial pair alignment and is more accurate than comparable algorithms such as Nemat [2]. It integrates a Spatial Transformer Network (STN) [9] consisting of a Vision Transformer (ViT) [3], and two conditional GANs [8], and is trained using four flows in an end-to-end manner. The objective is an unsupervised image registration approach. The goal is to “cast” the misaligned thermal image into a visible spectrum using GANs, in order to learn deformation parameters in a common spectrum (visible) when passed to the STN algorithm.

Flow 1 consists of a Visible-to-Thermal (V2T) conditional GAN ( $GAN_1$ ) that given, the ground-truth visible input ( $A$ ), translates (generates) its fake thermal ( $\hat{B}$ ) pair. Flow 2 uses a Thermal-to-Visible (T2V) conditional GAN ( $GAN_2$ ) to translate the fake visible image ( $\hat{A}_1$ ) from the thermal ground truth ( $B$ ). Both GANs share the same architecture consisting of a small U-NET [24] framework. In Flow 3, the STN shown in Figure 4b, now takes the combination of ( $A, \hat{A}_1$ ) and passes the vector through a ViT as a series of 64 patch embeddings. In the final Flow 4, the same T2V GAN from Flow 2 is used to translate the registered thermal image ( $B_R$ ) to the last intermediate fake thermal image ( $A_2$ ), thereby enforcing a cyclic consistency. Shown in Figure 4c, the encoded output is passed to a MLP which acts as a regressor and is used predict the affine matrix ( $\theta$ ) which

estimates parameters such as scale, rotation, translation, and shear. We modify the original MLP architecture in Figure 4c from two Linear-ReLU blocks to five. Finer control to learn deformation parameters is afforded when the network is deeper to localize geometry of finer features. After predicting the parameters of the affine matrix, the STN samples each pixel in a deformation grid and applies it onto the real thermal image ( $B$ ). The end result is the registered thermal face,  $B_R$ . For further details regarding the generator and discriminator architectures, training regimen, losses, and ablation studies for robustness, we direct the reader to our previous works in [16,17].

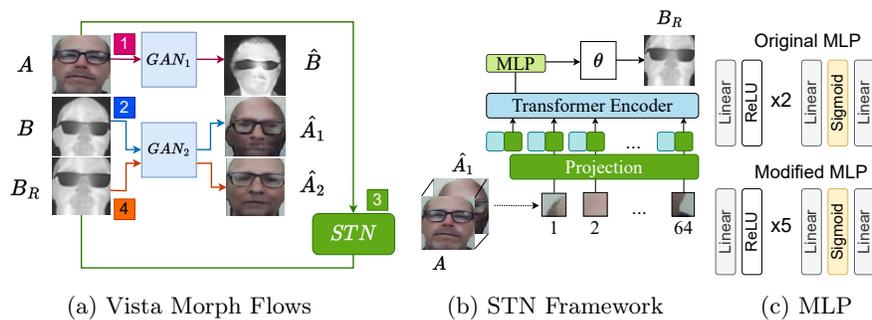


Fig. 4: Vista Morph Framework - Training Flows, STN, and MLP Regressor.

### 2.3 Experiments

**Image Registration.** The objective of this experiment is to align the thermal facial image with respect to the scale and geometry of the visible face, since no thermal reference is provided. We train the Vista Morph algorithm from scratch, using the entire 22,570 VT pairs, and test on 6,891 pairs. We use the PyTorch library and train on 8 Tesla V100 GPUs in parallel using automatic mixed precision for 10 epochs using a `batch_size=64`. To score registration accuracy, we use conventional metrics for unsupervised image registration - Normalized Cross Correlation (NCC) and Structural Similarity Index Measure (SSIM) [26] of the edge maps (e.g. morphological gradients of the visible and thermal images), in addition to Mutual Information (MI) [10,25].

**V2T Image Translation.** The objective of this experiment is to evaluate the quality of generated thermal faces before and after image registration to prove that the quality of generated thermal faces declines without registration leading to artifacts and incorrect identity. We align the entire ISS VT Facial Dataset for all 29,641 pairs using the trained Vista Morph algorithm. We then train a V2T conditional GAN called VTF-GAN [17] on the registered (Vista Morph) and the unregistered (original) training data. We use the PyTorch library and train on 8 Tesla V100 GPUs in parallel using automatic mixed precision for 300 epochs using a `batch_size=64`. We use two common generative metrics that

measure image quality to score results of the generated thermal face - the Frechet Inception Distance (FID) [6] and LPIPS [29].

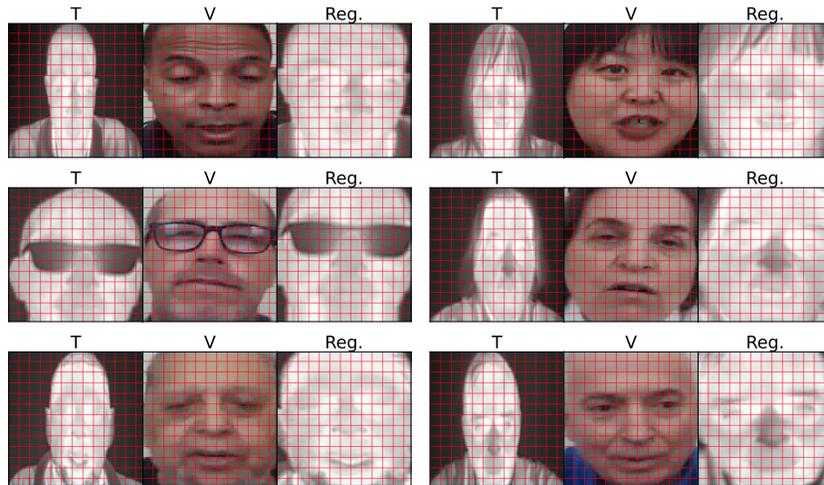


Fig. 5: Image Registration Samples for Six Patients.

	Registration Scores			V2T GAN Scores	
	SSIM Edges ( $\uparrow$ )	NCC Edges ( $\uparrow$ )	Mut. Info. ( $\uparrow$ )	FID ( $\downarrow$ )	LPIPS ( $\downarrow$ )
Before Reg.	0.691	0.003	0.227	121.715	0.399
After Reg.	<b>0.752 (8.9%)</b>	<b>0.136 (53.4x)</b>	<b>0.299 (31.6%)</b>	<b>79.680 (-52.8%)</b>	<b>0.342 (-16.8%)</b>

Table 1: Quantitative Scores for Image Registration and V2T Image Translation.

### 3 Results

**Image Registration.** Quantitative results shown in Table 1 demonstrates that alignment of VT faces improve significantly after registration. SSIM scores increase by 8.9%, NCC increases by 53.4 times, and Mutual Information gains at 31.6%. Registration samples are shown in Figure 5 for six test subjects. Grid lines support visual inspection and indicate that the registered thermal face (“Reg.”) is well aligned to the scale of the (V)isible ground truth, compared to the original, warped, (T)hermal image. Notice that despite non-frontal head poses such as looking down or upwards in addition to right and left head tilts, that the thermal faces are well aligned. In addition, Figure 1 shows that facial obfuscation with eyeglasses does not interfere with registration. Additional evidence is provided in Figure 6 that show difference maps between the visible and thermal image before and after registration. Notice that before registration

("Before Reg."), the red (visible) map is obscured under the blue (thermal) map. However, after registration ("After Reg."), the blue and red maps are completely superimposed.



Fig. 6: Difference Maps Visualizing Alignment Between (V)isible and (T)hermal Pairs Before and After Registration.

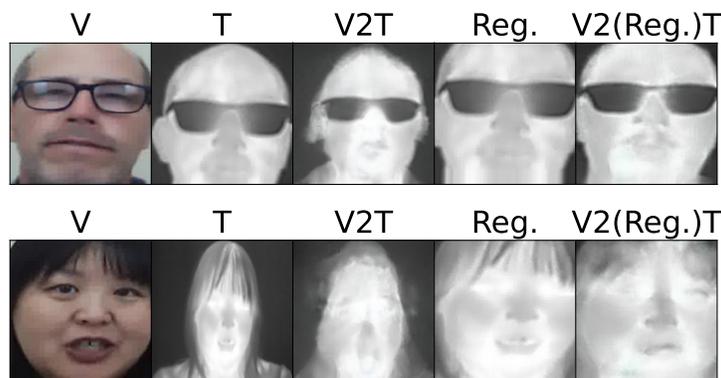


Fig. 7: V2T Image Translation Samples for Two Patients.

**V2T Image Translation.** Table 1 demonstrates that the quality of thermal faces improves after registration where FID scores improve by 52.8% and LPIPS scores improve at 16.8%. In Figure 7, we show results of the V2T image translation experiments. When training the GAN on the original, misaligned set, the generated thermal is shown in "V2T". Notice a lack of perceptual clarity, and subjective identity, as well as artifacts. However, when the GAN is trained using the Vista Morph registered VT pairs ("V" and "Reg."), the "V2(Reg.)T" generated thermal faces are higher quality and more similar to the original subject.

**Limitations.** Image registration is less successful when faces are both highly warped and obfuscated by masks or a combination of glasses and hats. Future works can leverage a Fourier Loss per our work in [16] that is integrated into the Vista Morph framework, used to register VT pairs in No- and Low-Light settings by learning signal frequencies (low, high edges) in addition to the spatial (pixel) domain. Validation in-clinic using a ground-truth thermal sensor should also be conducted to verify accurate heat distribution.

In addition to the VTF-GAN, we also trained a conditional diffusion model called VTF-Diff for V2T image translation [17]. The FID (96.420, +20.8%) and



Fig. 8: Unsuccessful VT Image Registration Samples for Two Subjects.

LPIPS scores (0.440, +22.4%) were not as competitive as the VTF-GAN results, since most images were distorted. However, we show in Figure 9, a limited number of successful generated thermal faces (“Diff.”). Although these results preserve spatial geometry, the diffusion results are discolored. This implies inaccurate distribution of heat (light pixels) and cold (dark pixels) compared to the ground truth (“Reg.”), which can lead to misleading medical assessment. As a result, the VTF-GAN algorithm is preferable over VTF-Diff as a method to augment the ISS VT Facial dataset with additional thermal faces.

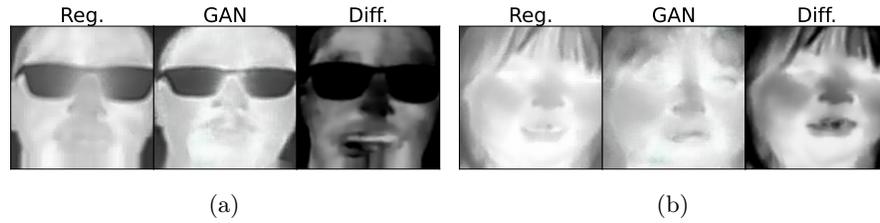


Fig. 9: V2T Image Translation with Diffusion Model (“Diff”) Leads to Inaccurate Heat Distribution and Artifacts.

## 4 Conclusion

We introduce the NIH ISS Visible-Thermal (VT) Facial Dataset, the largest VT cancer dataset to our knowledge, consisting of 29,000 VT pairs for AI pain research. We register all images using a novel, generative approach for multi-spectral facial imagery and modify the regressor network for finer estimation of alignment parameters. We show that generative tasks such as V2T image translation improve markedly after pairs are registered leading to improved image quality and resolution. Further, our approach demonstrates an effective method for applying two image modalities towards the investigation of cancer pain assessment.

**Acknowledgements** The authors would like to thank Elizabeth Lamping, Katherine Lee-Wisdom, NIH NCI clinic partners (Prostate, Thymoma, Phase 1, Neurofibroma, HIV, GI clinics) for the clinical patient protocol, in addition to Alex Hanson at Booz Allen Hamilton for providing computational support. The authors also thank the patients and their families for informed consent.

## References

1. Apolinário-Hagen, J., Fritsche, L., Bierhals, C., Salewski, C.: Improving attitudes toward e-mental health services in the general population via psychoeducational information material: A randomized controlled trial. *Internet interventions* **12**, 141–149 (2018)
2. Arar, M., Ginger, Y., Danon, D., Bermano, A.H., Cohen-Or, D.: Unsupervised multi-modal image registration via geometry preserving image-to-image translation. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 13410–13419 (2020)
3. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020)
4. Espinosa-Duró, V., Faundez-Zanuy, M., Mekyska, J.: A new face database simultaneously acquired in visible, near-infrared and thermal spectrums. *Cognitive Computation* **5**(1), 119–135 (2013)
5. Fink-Lamotte, J., Widmann, A., Fader, J., Exner, C.: Interpretation bias and contamination-based obsessive-compulsive symptoms influence emotional intensity related to disgust and fear. *PloS one* **15**(4), e0232362 (2020)
6. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: GANs trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems* **30** (2017)
7. Ioannou, S., et al.: Thermal infrared imaging in psychophysiology: potentialities and limits. *Psychophysiology* **51**(10), 951–963 (2014)
8. Isola, P., et al.: Image-to-image translation with conditional adversarial networks. In: *CVPR* (2017)
9. Jaderberg, M., , et al.: Spatial transformer networks. *NeurIPS* **28** (2015)
10. Kern, J.P., Pattichis, M.S.: Robust multispectral image registration using mutual-information models. *IEEE Transactions on Geoscience and Remote Sensing* **45**(5), 1494–1505 (2007)
11. Kong, L., Lian, C., Huang, D., Hu, Y., Zhou, Q., et al.: Breaking the dilemma of medical image-to-image translation. *Advances in Neural Information Processing Systems* **34**, 1964–1978 (2021)
12. Livesay, J.R., Porter, T.: Emg and cardiovascular responses to emotionally provocative photographs and text. *Perceptual and motor skills* **79**(1), 579–594 (1994)
13. Mallat, K., et al.: A benchmark database of visible and thermal paired face images across multiple variations. In: *BIOSIG*. pp. 1–5. *IEEE* (2018)
14. (NCI), N.C.I.: Machine learning to analyze facial imaging, voice and spoken language for the capture and classification of cancer/tumor pain - full text view (Jun 2020), <https://clinicaltrials.gov/ct2/show/NCT04442425>
15. Ordun, C., Cha, A.N., Raff, E., Gaskin, B., Hanson, A., Rule, M., Purushotham, S., Gulley, J.L.: Intelligent sight and sound: A chronic cancer pain dataset. *arXiv preprint arXiv:2204.04214* (2022)
16. Ordun, C., Raff, E., Purushotham, S.: Vista-morph: Unsupervised image registration of visible-thermal facial pairs. *arXiv preprint arXiv:2306.06505* (2023)
17. Ordun, C., Raff, E., Purushotham, S.: When visible-to-thermal facial gan beats conditional diffusion. *arXiv preprint arXiv:2302.09395* (2023)
18. Ordun, C., et al.: The use of AI for thermal emotion recognition: A review of problems and limitations in standard design and data. *AAAI* (2020)

19. Pavlidis, I., Garza, I., Tsiamyrtzis, P., Dcosta, M., Swanson, J.W., Krouskop, T., Levine, J.A.: Dynamic quantification of migrainous thermal facial patterns—a pilot study. *IEEE Journal of Biomedical and Health Informatics* **23**(3), 1225–1233 (2018)
20. Pavlidis, I., et al.: The imaging issue in an automatic face/disguise detection system. In: *IEEE Computer Vision Beyond the Visible Spectrum*. pp. 15–24 (2000)
21. Poster, D., et al.: A large-scale, time-synchronized visible and thermal face dataset. In: *WACV*. pp. 1559–1568 (2021)
22. Puri, C., et al.: Stresscam: non-contact measurement of users’ emotional states through thermal imaging. In: *CHI’05* (2005)
23. Ring, E., Ammer, K.: Infrared thermal imaging in medicine. *Physiological measurement* **33**(3), R33 (2012)
24. Ronneberger, O., et al.: U-net: Convolutional networks for biomedical image segmentation. In: *MICCAI* (2015)
25. Russakoff, D.B., Tomasi, C., Rohlfing, T., Maurer, C.R.: Image similarity using mutual information of regions. In: *Computer Vision-ECCV 2004: 8th European Conference on Computer Vision, Prague, Czech Republic, May 11-14, 2004. Proceedings, Part III* 8. pp. 596–607. Springer (2004)
26. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* **13**(4), 600–612 (2004)
27. Zhang, K., Zhang, Z., Li, Z., Qiao, Y.: Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE Signal Processing Letters* **23**(10), 1499–1503 (2016)
28. Zhang, R.: Making convolutional networks shift-invariant again. In: *ICML*. pp. 7324–7334. PMLR (2019)
29. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: *CVPR*. pp. 586–595 (2018)