

# TATL: Task Agnostic Transfer Learning for Skin Attributes Detection

Duy M. H. Nguyen<sup>a,b</sup>, Thu T. Nguyen<sup>c</sup>, Huong Vu<sup>d</sup>, Quang Pham<sup>e</sup>, Manh-Duy Nguyen<sup>f</sup>, Binh T. Nguyen<sup>g,h,i,\*</sup>, Daniel Sonntag<sup>a,j</sup>

<sup>a</sup>German Research Center for Artificial Intelligence, Saarbrücken, Germany

<sup>b</sup>Max Planck Institute for Informatics, Germany

<sup>c</sup>University of Louisiana at Lafayette, USA

<sup>d</sup>University of California, Berkeley, USA

<sup>e</sup>School of Computing and Information Systems, Singapore Management University

<sup>f</sup>School of Computing, Dublin City University, Ireland

<sup>g</sup>AISIA Research Lab, Ho Chi Minh City, Vietnam

<sup>h</sup>University of Science, Ho Chi Minh City, Vietnam

<sup>i</sup>Vietnam National University, Ho Chi Minh City, Vietnam

<sup>j</sup>Oldenburg University, Germany

---

## Abstract

Existing skin attributes detection methods usually initialize with a pre-trained Imagenet network and then fine-tune on a medical target task. However, we argue that such approaches are suboptimal because medical datasets are largely different from ImageNet and often contain limited training samples. In this work, we propose *Task Agnostic Transfer Learning (TATL)*, a novel framework motivated by dermatologists' behaviors in the skincare context. TATL learns an attribute-agnostic segmenter that detects lesion skin regions and then transfers this knowledge to a set of attribute-specific classifiers to detect each particular attribute. Since TATL's attribute-agnostic segmenter only detects skin attribute regions, it enjoys ample data from all attributes, allows transferring knowledge among features, and compensates for the lack of training data from rare attributes. We conduct extensive experiments to evaluate the proposed TATL transfer learning mechanism with various neural network architectures on two popular skin attributes detection benchmarks. The empirical results show that TATL not only works well with multiple architectures but also can achieve state-of-the-art performances, while enjoying minimal model and computational complexities. We also provide theoretical insights and explanations for why our transfer learning framework performs well in practice.

*Keywords:* Transfer Learning, Skin Attribute Detection, Encoder-Decoder Architecture.

---

## 1. Introduction

Melanoma is one of the most dangerous types of skin cancer. Even though it only accounts for 1% of all skin cancer cases, it is responsible for the majority of skin cancer deaths (Ward and Farma, 2017). In 2021, it is estimated that there will be 207,390 new cases of melanoma will be diagnosed and 7,180 recent deaths from the disease in the United States alone (Society, 2021).

Moreover, the 5-year relative survival rate for melanoma reduced from 99% for cases diagnosed at a localized stage to 27% for a distant stage (Society, 2021). Therefore, there have been tremendous efforts in detecting the disease in its early stages (Masood and Ali Al-Jumaily, 2013; Curiel-Lewandrowski et al., 2019). One of the most promising technology is *Dermoscopy*, which can generate high-resolution images of skin lesions and allows dermatologists to examine the lesion regions more carefully (Celebi et al., 2019). However, dermoscopy still requires extensive training, which is expensive, time-consuming, error-prone, and might not be widely available (Zalaudek et al., 2008). Therefore, it is important and highly beneficial to develop automatic systems to detect abnormal skin lesions and aid dermatologists during diagnosis (Nunnari et al., 2021b).

For this purpose, the International Skin Imaging Collaboration (ISIC) hosted challenges for automatic melanoma detection based on dermoscopic images (ISIC-2018, ISIC-2017) (Codella et al., 2019, 2018). In this work, we focus on Task 2 of predicting the locations of dermoscopic attributes in an image. In particular, there are *five* dermoscopic attributes that the challenge focused on: Streaks, Globules, Pigment Network, Negative Network, and Milia-like Cysts. Locating these clinically meaningful skin lesion patterns helps detect anomalous regions and provides an explanation for dermatologists to verify and make further diagnoses. For instance, the Negative Network, which consists of relatively light regions and some darker regions, is usually considered a melanoma-specific structure (Pizzichetta et al., 2013). We provide an example of such attributes in Figure 1.

---

\*Corresponding author: ngtbinh@hcmus.edu.vn (Binh T. Nguyen)

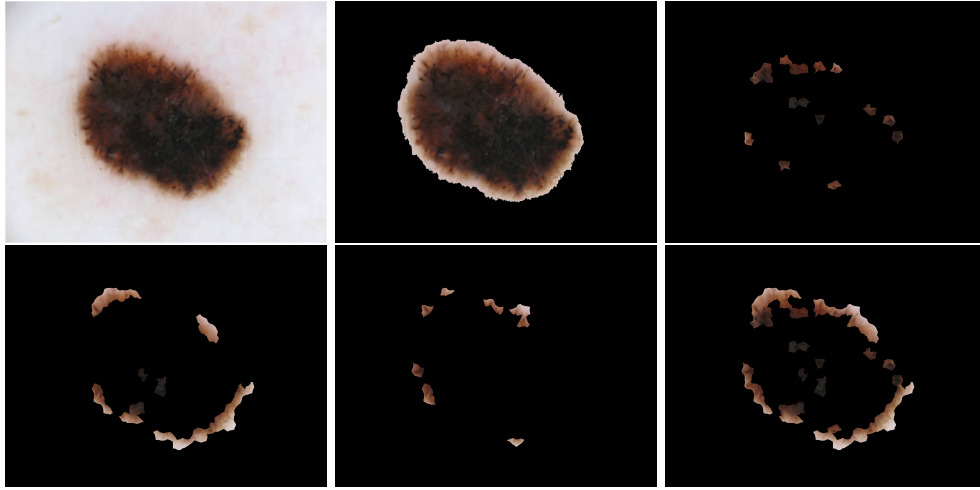


Figure 1: From left to right: the picture of melanoma, its segmentation, and the annotation for globules (upper row), the annotation for Pigment Network, Streaks, and the three’s union (below row). The images are taken from ISIC 2018 (Codella et al., 2019).

Due to its highly competitive nature, many methods in the ISIC challenges are based on the ensemble strategies with different types of ImageNet (ImgNet) pre-trained backbones. For instance, the ISIC 2018’s winner method (Koohbanani et al., 2018) employed a mixture of four pre-trained networks ResNet152 (He et al., 2016), DenseNet169 (Huang et al., 2017), Xception (Chollet, 2017), and ResNetV2 (Szegedy et al., 2017) in the encoding part of the U-Net (Ronneberger et al., 2015a) segmentation models and then performed transfer learning on the ISIC 2018 training set. Although Koohbanani et al. (2018) achieved state-of-the-art performance, it is not an attractive method in practice because of *two* reasons. First, Koohbanani et al. (2018) simultaneously stores five different models in a single GPU in the pre-training phase, which costs a total of over 300 million parameters (Table 7). As a result, this approach cannot be deployed on most **non** high-end GPU cards, which limits its accessibility to practical deployment and its extensibility for future research. Second, recent works in Raghu et al. (2019); Cheplygina (2019); Nguyen et al. (2020) showed that transfer learning from ImageNet might be sub-optimal in many scenarios because medical images are primarily different from ImageNet data. Moreover, most medical datasets, including ISIC-2018 and ISIC-2017, suffer from the scarcity of training data, and the number of instances per attribute is imbalanced, as depicted in Table 2.

To address the problems of massive memory requirement and domain gap between ImageNet and medical images mentioned above, we propose *Task Agnostic Transfer Learning* (TATL), an efficient framework to detect skin attributes in dermoscopic images. TATL’s design is inspired by how dermatologists diagnose in practice: identify abnormal skin regions in the first step and then inspect them more closely in the second step <sup>1,2</sup>. Unlike previous works that try to segment skin attributes on the image directly, TATL introduces an *Attribute-Agnostic Segmenter* that first detects anomalous regions in an image, regardless of their attributes. Then, TATL transfers this segmenter’s knowledge to a set of attributes-specific segmenters (*Target Segmenters*) to detect each specific attribute. Notably, the Attribute-Agnostic Segmenter is task-agnostic because it only identifies abnormal skin regions, possibly including data from all attributes. Therefore, TATL alleviates the lack of training samples by training the segmenter in the first stage. Furthermore, through the experiments, we find out that TATL also facilitates knowledge sharing among the attribute-segmenters, thus enhancing the generalization and stability of the whole system. Furthermore, we also provide theoretical insights reasoning that TATL works by bridging the gap between the target task’s data and the source dataset. The attribute-specific classifiers are particularly initialized from the TATL’s Union-Segmenter, which enjoys a tighter generalization error bound than other methods initializing from ImageNet. This analysis sheds light on the promising performances of TATL.

In summary, our contributions are threefold. Firstly, we propose TATL, a novel transfer learning approach for skin attribute detection. Extensive experiments on the ISIC 2018 and ISIC 2017 benchmarks validate the effectiveness of TATL against state-of-the-art methods while requiring only 1/30 number of parameters compared with the ISIC 2018’s winner. Furthermore, TATL significantly improves several skin diagnosis methods pre-trained on ImageNet, especially for attributes with limited training instances. Secondly, we provide theoretical insights explaining the success of TATL; thereby, TATL has a high ability to reduce domain gaps by shifting from color images (ImageNet) to the medical domain (Skin images). Finally, TATL can provide informative outputs to aid dermatologists in further diagnosing as it can simultaneously predict both abnormal regions

<sup>1</sup><https://www.polyclinic.com/health-wellness-library/find-skin-cancer.html>

<sup>2</sup><https://www.cancer.org/content/dam/CRC/PDF/Public/8825.00.pdf>

and different kinds of skin features, which makes a promising first step towards a more practical and beneficial medical-aid deep learning system.

## 2. Related Work

### 2.1. Transfer Learning for Medical Image Analysis

Medical image analysis is a vital research venue and has a significant impact on practice. However, most medical image datasets have limited training samples and often suffer from the data-imbalanced problem. Therefore, a popular strategy is *transfer learning*, which uses a pre-trained ImageNet model as an initialization to build additional components. Transfer learning is a base of many existing methods (Abràmoff et al., 2016; De Fauw et al., 2018; Gulshan et al., 2016; Rajpurkar et al., 2017), and is a norm for practitioners (Tan et al., 2018). However, recent studies in Cheplygina (2019); He et al. (2020) conducted a large-scale analysis on the benefit of this strategy. They concluded that for medical images, transfer learning based on pre-trained ImageNet is not consistently better than random initialization. One reason is that medical images are vastly different from the ImageNet dataset, resulting in the pre-trained weights that are not helpful for the current task. Another reason is that medical data are often imbalanced and rare due to data privacy. For example, Table 2 illustrates the distributions of each skin characteristic in the ISIC 2017 and ICSIC 2018 datasets, with the rarest class (Streaks) accounting for only 7.98% (113 images) and 3.86% (100 images) of the training data, respectively. In comparison, the most common class (Pigment Network) accounts for 79.03% (1119 images) and 58.67% (1522 images) of the total samples. Fortunately, TATL can address the scarcity of training data in such situations by transferring the knowledge from the *Attribute-Agnostic Segmenter*. Moreover, we since apply the strategy *one class one model* (Buda et al., 2018), each *Target-Segmenter* classifier in TATL only detects one attribute, and can alleviate the data’s imbalance problem.

### 2.2. Self-Supervised Learning In Computer Vision

Self-supervised learning, first mentioned in Schmidhuber (1990), refers to a technique of creating additional tasks for training where the label is also a part of the data (images) rather than a set of separate labels (annotations). This strategy has been a successful pre-training technique in various vision applications, including image colorization (Vondrick et al., 2018; Larsson et al., 2016; Zhang et al., 2016), image inpainting (Pathak et al., 2016; Chen et al., 2020), and video representation (Misra et al., 2016; Lee et al., 2017). In self-supervised learning, a newly created task for pre-training is called the “*pretext task*”, and the main tasks used for fine-tuning are called the “*downstream tasks*”. Various strategies have been proposed to construct the pretext task based on the image rotation (Gidaris et al., 2018), temporal correspondence (Li et al., 2019; Wang et al., 2019b), cross-modal consistency (Wang et al., 2019a) and instance discrimination with contrastive learning (Wu et al., 2018). Recently, in the medical domain, He et al. (2020) successfully applied self-supervised learning in diagnosing COVID19 from CT scans based on contrastive self-supervised learning for reducing the risk of overfitting. Chen et al. (2019) also presented a context restoration framework in which the image is disordered by randomly changing the order of their sub-patches, and then a neural network is trained to recover the original input.

Our TATL approach shares some similarities with self-supervised learning in the sense of building better pretrained models without extra training instances or learning through auxiliary tasks. Here, the Attribute-Agnostic Segmenter plays a role as of learning a pretext task, while detecting the attributes in the Target Segmenters are learning the downstream tasks. However, our work differs from the SSL approaches because we define an auxiliary task through solving  $(x, g(y))$  while SSL methods follow the scheme  $(x, f(x))$  where  $(x, y)$  indicates for the image and corresponding label,  $g$  is an operator on the training label, and  $f$  is another transformation on the image such as rotation (Gidaris et al., 2018), dividing images into sub-patches and suffering their positions (Chen et al., 2019). The construction of  $g$  is specifically designed for the medical domain, which makes TATL’s pretext task closely complements the subsequent downstream tasks. Therefore, if the pretext task of recognizing skin attribute regions can perform well, it will likely facilitate detecting such areas’ attributes. Finally, by providing skin attribute regions or abnormal regions from the pretext task, TATL is helpful to end-users by allowing dermatologists to validate the employed system’s diagnostics.

## 3. Preliminaries

This section aims to formalize our problem setting and outline the dermatologists’ practice to diagnose skin attributes, which later motivates our method.

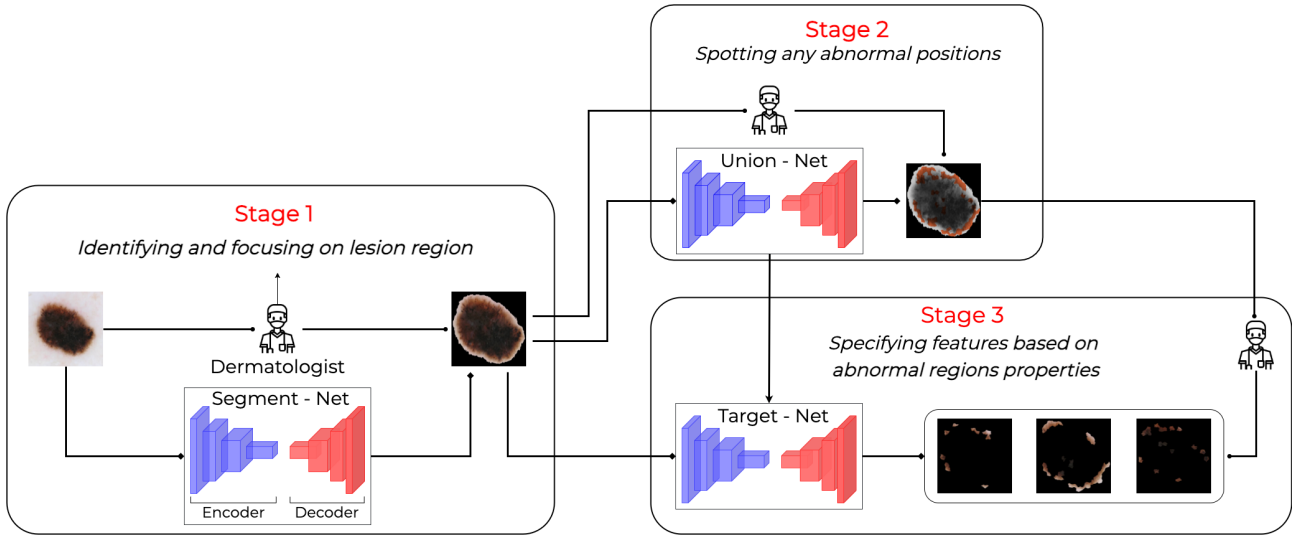


Figure 2: The prediction pipeline for each skin feature is motivated by a dermatologist’s behavior when examining a patient.

### 3.1. Problem Statement and Background

We consider the skin attributes detection problem on a target dataset consisting of training images and their corresponding masks  $\mathcal{D} = \{(x_1, \{y_1^{(i)}\}_{i=1}^{|Y|}) \dots, (x_n, \{y_n^{(i)}\}_{i=1}^{|Y|})\}$ . The detector, parameterized by  $W$ , can be initialized from a pre-trained model on another dataset, which we call the source dataset  $\mathcal{D}^{src}$ . Moreover, each training sample in the target domain consists of an image  $x \in \mathbb{R}^{c \times w \times h}$  and a set of labels  $\{\{y^{(i)}\}_{i=1}^{|Y|}, y^{(i)} \in \mathbb{R}^{w \times h}\}$ , where  $y^{(i)}$  is a binary mask indicating the skin region associated with the  $i$ -th attribute. In this work, we consider **five** different attributes: Globules, Milia, Negative Network, Pigment Network, and Streaks, shorthanded as  $Y = \{G, M, N, P, S\}$ , i.e.,  $|Y| = 5$ . It is worth noting that each sample may not have all the attributes and the label for those missing attributes is the empty mask. The training process can be performed by minimizing the empirical risk:

$$\hat{f}(\{x_j, \{y_j^{(i)}\}_{j=1}^n | W\}) = \arg \min_f \frac{1}{n} \sum_{j=1}^n \sum_{i=1}^{|Y|} \lambda_1 L_{Tversky}(\hat{y}_j^{(i)}, y_j^{(i)}) + \lambda_2 L_{Jaccard}(\hat{y}_j^{(i)}, y_j^{(i)}), \quad (1)$$

where  $\hat{y}_j^{(i)}$  denotes the binary mask prediction of the network on a sample about the  $i$ -th attribute. For each attribute, we use the Tversky loss  $L_{Tversky}$ , which is a generalization of Dice loss (Eelbode et al., 2020; Jadon, 2020; Salehi et al., 2017) and the soft Jaccard loss functions  $L_{Jaccard}$  (Eelbode et al., 2020; Kawahara and Hamarneh, 2018), to penalise the deviation between network’s predictions and the ground-truths. Formally, these loss functions can be calculated as:

$$L_{Tversky}(\hat{y}, y) = 1 - \frac{\alpha + \langle y, \hat{y} \rangle}{\alpha + \langle y, \hat{y} \rangle + \beta \langle 1 - \hat{y}, y \rangle + (1 - \beta) \langle \hat{y}, 1 - y \rangle}, \quad (2)$$

$$L_{Jaccard}(\hat{y}, y) = 1 - \frac{\alpha + \langle y, \hat{y} \rangle}{\alpha + \|y\|_1 + \|\hat{y}\|_1 - \langle y, \hat{y} \rangle}. \quad (3)$$

Here, the prediction  $\hat{y}$  and the ground-truth  $y$  are first re-shaped into a vector form;  $\langle \cdot, \cdot \rangle$  and  $\|\cdot\|_1$  are the inner product and  $L^1$  norm respectively. The parameter  $\alpha$  is used to ensure that the loss functions are not undefined when division by zero or in case  $y = \hat{y} = 0$ . The parameter  $\beta$  in other way controls the magnitude of false positives and false negatives. In our experiment, we choose  $\lambda_1 = \lambda_2 = 0.5$  to balance the importance of two loss functions and parameters  $\alpha = 1, \beta = 0.6$  through validation experiments.

### 3.2. Inspirations from Dermatologists’ Behaviours

Figure 2 depicts a prediction pipeline inspired by the conventional diagnosis process of dermatologists, as discussed in Section 1. In the first step, dermatologists will identify lesion regions by eliminating irrelevant background and rescaling these regions to a higher resolution for better visualization (Stage 1). Following that, they continue to spot any abnormal and clinically relevant sub-areas on the lesion (Stage 2). Finally, by accounting for these factors, doctors diagnose specific skin attributes by comparing various features based on their textures and colors compared to nearby spaces (Stage 3). We argue that identifying lesion and skin attribute regions is crucial since it serves focal points for later steps, and we develop a skin attributes detection framework that closely follows the three-step procedure represented in Figure 2.



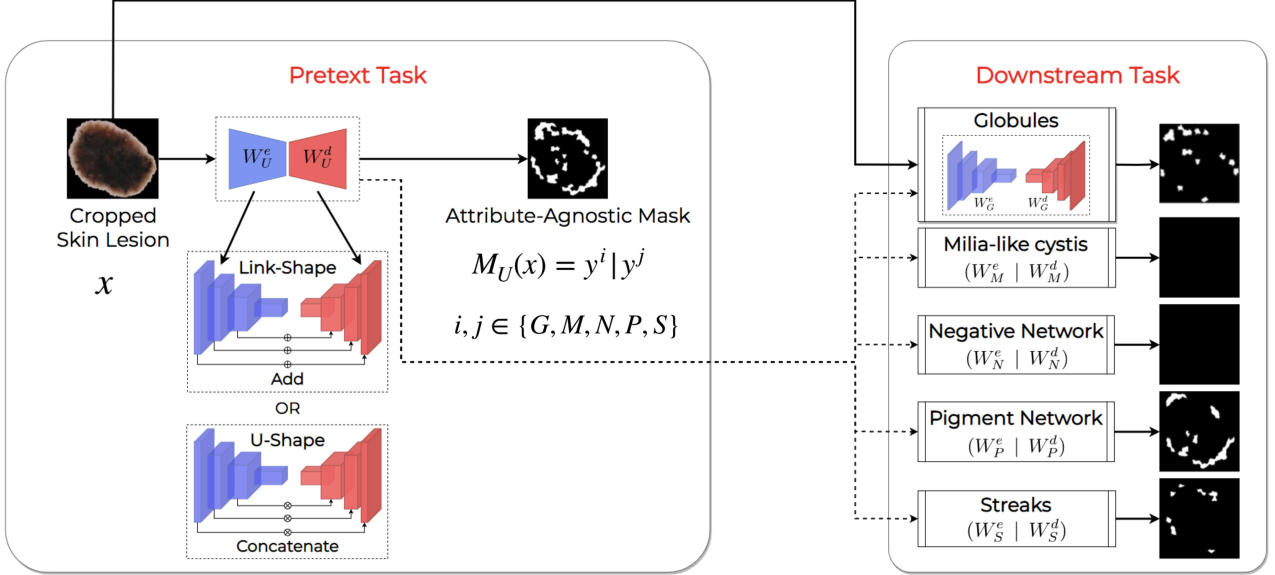


Figure 3: The training procedure of proposed TATL with two steps. Left: learning to recognize regions containing any attributes through attribute-agnostic mask (Pretext Task). Right: training parameters are transferred from Pretext Task to Downstream Task to segment each feature independently.

We realize the diagnosis procedure into a single framework named Task Agnostic Transfer Learning (TATL). First, TATL employs a Segment-Net to segment the lesion regions from an input image. Then, TATL trains an Attribute-Agnostic Segmenter to detect all skin attribute regions in the image, regardless of their attributes, which is inspired by the second step in the procedure. Finally, the parameters of the Attribute-Agnostic Segmenter are utilized as an initialization for the Target Segmenters (Tar-S), which are trained to identify just one specific attribute and are employed as the final step in the method.

TATL not only closely resembles how dermatologists diagnose but also enjoys two additional benefits than conventional approaches. First, TATL provides additional information about the skin attribute regions regardless of their attributes, which can be helpful for dermatologists. Remarkably, such areas reveal variations and commonalities of relevant lesions, thereby reducing subjective errors in the evaluation process. This can be demonstrated by two examples in Figure 9. Here, attributes such as “Negative Network” or “Globules” can be challenging to identify in isolation. In contrast, the union of all attributes provided by TATL can correctly cover those areas. Second, adapting weights trained on skin attribute regions to a specific attribute can guide the network to pay attention to shared features across diverse attributes, thus strengthening trained systems to be more robust and stable. We will empirically demonstrate in detail these properties from Subsection 5.3 to Subsection 5.8.

## 4. Methodology

We now detail our TATL framework and discuss the theoretical properties, which sheds light on its competitiveness. We first provide an overview of our TATL framework in Section 4.1. Then, we discuss the pretext task with the Segment-Net and Attribute-Agnostic Segmenter in Section 4.2, and the downstream task with the Target Segmenter in Section 4.3. Then, we summarize the TATL framework and outline its algorithm in Section 4.4. Lastly, we conclude this section with a theoretical insight in Section 4.5.

### 4.1. Task Agnostic Transfer Learning for Skin lesion Attribute Detection

*The Encoder - Decoder Architecture.* The core component in our TATL framework is the encoder-decoder architecture, which takes an image as input and produces a binary mask as output. While we employ two kinds of encoder-decoder networks in our method, they share the same design as follows. The encoder part *could be any feature extraction layers* in arbitrary architectures such as ResNet152 (He et al., 2016) or EfficientNet (Tan and Le, 2019). For feature extraction purposes, we thus discard non-linear rectification layers in these architectures. For the decoder path, we used layers to up-sample the encoder’s features back to the original input’s dimension. Particularly, to match the encoder’s stages, the decoder consists of an up-sampling layer and a sequence of convolutional blocks where each block has two  $3 \times 3$  convolution filters with activation functions in between. Each stage in the decoder receives a feature map from its immediate preceding layer and a corresponding feature from the encoder’s stage. The two inputs are combined by either the *adding* or *concatenating* operations,

corresponding to the settings of LinkNet (Chaurasia and Culurciello, 2017) and U-Net (Ronneberger et al., 2015b).

*The TATL Framework.* Our TATL framework involves of *three* encoder-decoder networks. The first network, **Segment-Net**, segments and upscales the lesion regions in the original image. Then, the second network, **Attribute-Agnostic Segmenter**, takes the lesion regions as input and learns to segment the skin attribute regions, possibly including any of the five attributes of interests. Finally, for each attribute, a corresponding network, the **Target Segmenter**, is trained to segment that attribute’s regions. Moreover, the Target Segmenter’s decoder also rescales the final mask to match the original image’s dimensions. Therefore, our TATL framework consists of *seven* networks in total, one Segment-Net, one Attribute-Agnostic Segmenter, and five Target Segmenters corresponding to five attributes. Each network uses either the *Link-shape* or *U-shape* architecture with the b0-EfficientNet (Tan and Le, 2019) as the main backbone due to its lightweight property compared to other architectures (Table 7).

We now introduce some notations to detail our framework. We denote  $\{f_{seg}(\cdot, W_{seg}), f_U(\cdot, W_U), f_i(\cdot, W_i)\}$  as the corresponding networks of the Segment-Net, the Attribute-Agnostic Segmenter, and the Target-Segmenter for each attribute  $i \in Y = \{G, M, N, P, S\}$  in Stage 3 respectively (Figure 2). We also use  $W_m = \{W_m^e, W_m^d\}$ ,  $m \in \{Seg, U, G, M, N, P, S\}$  as the total parameters of networks  $f_m$ , where  $\{W_m^e\}$ , and  $\{W_m^d\}$  represent weights of encoder and decoder layers. To train the Segment-Net, we denote  $D_{seg} = \{x, M_{seg}\}$  as the set of images and their corresponding lesion region masks.

#### 4.2. Pretext Task with Segment-Net and Attribute-Agnostic Segmenter

Here, we refer to the TATL *Pretext Task* as the problem of recognizing regions containing any attributes. The pretext training scheme consists of two stages: (i) cropping skin lesion with the Segment-Net; and (ii) segmenting the attribute-agnostic mask with Attribute-Agnostic Segmenter. In the following, we describe the training procedure with the corresponding components in detail.

##### Segment-Net

In Stage 1, we use the Segment-Net on the lesion dataset  $D_{seg}$  to eliminate extraneous skin-based regions and only keep the lesion regions. Especially, given an original image, we re-scale it to a size of  $386 \times 512$  and pass it into the Segment-Net. A bounding box with an offset value of 40 pixels in four directions is used to crop the Segment-Net’s output so that errors in the segment step are not propagated in the later stages. Once this bounding box has been created, it is scaled to match the resolution of the input image and utilized as a new input for the next stage.

##### Attribute-Agnostic Segmenter

The second stage focuses on training the Attribute-Agnostic Segmenter. We first define the *Attribute-Agnostic* region as a region that contains at least one of the attributes in  $Y$ . From this, we define an intermediate dataset of the Attribute-Agnostic as  $D_U = \{x, M_U\}$ , where  $M_U$  is the binary mask corresponding to an image whose value is 1 whenever a pixel is an attribute from  $Y$  (Pretext Task). Note that given an image  $x$  and a set of attributes masks,  $M_U$  is the *union* of all the masks and can be easily constructed by performing bitwise OR operator as:

$$M_U(x) \triangleq y^{(1)} | y^{(2)} | \dots | y^{(|Y|)}, \quad (4)$$

where  $|$  denotes the bitwise OR operator. The dataset  $D_U$  is used to train the Attribute-Agnostic Segmenter such that it can detect the skin attribute regions belonging to any of the attributes. It is important to note that  $D_U$  contains masks covering all attributes, thus ameliorating the negative effect of training data scarcity, especially for minor classes.

#### 4.3. Downstream Task with Target Segmenters

Given  $W_U$  learned from the Attribute-Agnostic Segmenter, we can proceed to Stage 3 and train the segmenter for each of the attributes (*Downstream Task*). Different from previous approaches, we initialize the Target-Segmenter’s parameters from the Attribute-Agnostic Segmenter’s parameters as:  $W_i^e \leftarrow W_U^e$  and  $W_i^d \leftarrow W_U^d$  for each type  $i$ - attribute. Lastly, a set of *Target-Segmenters* is trained to segment the attributes.

Having a dedicated network for each attribute is advantageous in alleviating the imbalance training data problem. Moreover, we explore two strategies in training the *Target-Segmenters*, which corresponds to allowing knowledge sharing across attributes or not. First, we *freeze* all the encoders (**TATL-Freeze**) to allow feature sharing across attributes because the encoder is initialized from the Attribute-Agnostic Segmenter. Second, we allow both the encoder and decoder to be updated (**TATL-Non Freeze**), which allows each Target-Segmenter to adapt to their dedicated attribute. We summarize these Pretext (with Attribute-Agnostic Segmenter) and Downstream Tasks in Algorithm 1 and Figure 3. Ablation studies for freeze encoder layers are also discussed in Table 8 of Subsection 5.4.

---

**Algorithm 1:** The TATL Algorithm

---

**Input:** Pre-trained ImageNet from employed backbone  $W_{\text{ImgNet}} = \{W_{\text{ImgNet}}^e, W_{\text{ImgNet}}^d\}$   
The attribute dataset:  $\mathcal{D} = \{(x_1, \{y_1^{(i)}\}_{i=1}^{|Y|}) \dots, (x_n, \{y_n^{(i)}\}_{i=1}^{|Y|})\}$   
**Output:** Trained weights  $W_i = \{W_i^e, W_i^d\}$ ,  $i \in Y = \{G, M, N, P, S\}$   
// Create the attribute-agnostic masks  
1 **for** each image  $x \in \mathcal{D}$  **do**  
2 |  $M_U[x] = y^{(1)} | y^{(2)} | \dots | y^{(|Y|)}$   
3 **end**  
// Learning *Attribute-Agnostic Segmenter*  
4 **Initialise:**  $W_U^e \leftarrow W_{\text{ImgNet}}^e$  and  $W_U^d \leftarrow W_{\text{ImgNet}}^d$   
5 **for** minibatch  $\{x_k, y_k\}$  where  $x_k \in \mathcal{D}$ ,  $y_k \in M_U$  **do**  
6 | Minimize  $\hat{f}_U(\{x_k, y_k\} | W_U)$  using Eq. (1)  
7 | Update  $W_U^e, W_U^d$   
8 **end**  
// Learning *Target-Segmenter* for each attribute  
9 **for** attribute  $i \in Y$  **do**  
10 | **Initialise:**  $W_i^e \leftarrow W_U^e$  and  $W_i^d \leftarrow W_U^d$   
11 | **for** minibatch  $\{x_k, y_k^i\}$  where  $x_k, y_k^i \in \mathcal{D}$  **do**  
12 | | Minimize  $\hat{f}_i(\{x_k, y_k^i\} | W_i)$  using Eq. (1)  
13 | | Update  $W_i^e, W_i^d$   
14 | **end**  
15 **end**  
16 **return**  $W_i^e$  and  $W_i^d$  where  $i \in Y$

---

#### 4.4. Summary of TATL

We summarize the training and inference pipeline of TATL as follows,

*Training Step.* Given a training dataset, TATL training performs the following steps,

1. Training the Segment-Net using dataset  $D_{\text{seg}} = \{x, M_{\text{seg}}\}$ .
2. Training Attribute-Agnostic Segmenter and Target-Segmenters using Algorithm 1 and dataset  $\mathcal{D} = \{(x_1, \{y_1^{(i)}\}_{i=1}^{|Y|}) \dots, (x_n, \{y_n^{(i)}\}_{i=1}^{|Y|})\}$

Note that we only apply a U-Shape encoder-decoder architecture in the first step, while we use both U-shape and Link-shape connections in the second step.

*Inference Step.* Given an input image, it first will be segmented by the Segment-Net and then fed into five different Target-Segmenters to segment five kinds of skin features. To compare with other competitors, we selected the b0-EfficientNet (Tan and Le, 2019) as the network backbone and used either the *adding* (Link-shape) and *concatenating* (U-shape) operations to correlate feature maps obtained from encoder and decoder layers. The final binary map for each skin attribute is produced by averaging the probability estimates from these two network architectures. Note that our pipeline prediction requires a focused lesion image generated by the Segment-Net, which might be influenced by the employed segmentation method. Fortunately, TATL only uses a bounding box with a specific offset around predicted lesion regions, thus making the following steps less susceptible to the resulting segmentation step. We present an experiment in Table 5 to validate this property.

#### 4.5. Theoretical Insights

This section provides theoretical insights to justify our approach using recent results from data-dependent stability of optimization schemes. Particularly, we investigate the model’s generalization to a target domain based on its initialization weights and show that initializing in a TATL fashion gives a tighter generalization error bound than ImageNet initialization via Proposition 1.

First, we introduce used notations, which are briefly summarized in Table 1. We consider the supervised training problem with  $X \subset \mathbb{R}^n$  as the input space and  $Y \subset \mathbb{R}$  as the output space. We also assume that training and testing instances are sampled i.i.d. from a probability distribution  $D$  over  $Z = X \times Y$ . Also, we denote the training set  $S = \{z_i\}_{i=1}^m \stackrel{iid}{\sim} D$ , where a training sample  $z_i$  consists of an input image  $x \in X$  and its corresponding label  $y \in Y$ . We express  $H = \{w_j\}$  as the hypothesis space where  $w_j \in \mathbb{R}^d$  denotes a hypothesis (model) with dimension  $d$  that maps an input instance  $x$  to its corresponding label  $y$ . Lastly, we define a map  $A_S : S \rightarrow H$  as a learning algorithm that returns a hypothesis given a training data set  $S$ .

Table 1: Summary of used notations.

Notations	Definitions
$X, Y$	The input and output spaces
$S$	the training set
$Z$	the joint input-output space ( $Z = X \times Y$ )
$\{\alpha_t\}_{t=1}^T$	step sizes
$I = \{j_t\}_{t=1}^T$	random indices
$D$	the data-generating distribution
$H$	the hypothesis space
$w$	a hypothesis
$A_S$	learning algorithm given the training data $S$
$\epsilon(D, w_1)$	stability function of $D$ and $w_1$

Kuzborskij and Lampert (2018) established a data-dependent aspect of algorithm stability for Stochastic Gradient Descent (SGD) given a training set  $S = \{z_i\}_{i=1}^m \stackrel{iid}{\sim} D$ , step sizes  $\{\alpha_t\}_{t=1}^T$ , random indices  $I = \{j_t\}_{t=1}^T$ , and an initialization weight  $w_1$ , which is sequentially updated as:

$$w_{t+1} = w_t - \alpha_t \nabla \ell(w_t, z_{j_t}). \quad (5)$$

Here,  $\ell(w, z)$  is a loss function, which measures the difference between predicted values and true values with parameters  $w \in H$  on an sample  $z$ . We indicate  $(D, w_1)$  as the data-generating distribution and the initialization point  $w_1$  of SGD,  $\epsilon(D, w_1)$  as a stability function of  $D$  and  $w_1$ .

To characterize a randomized learning algorithm  $A$ , we define its ‘‘On-Average stability’’.

**Definition 1.** (*On-Average stability*). A randomized algorithm  $A$  is  $\epsilon(D, w_1)$ -on-average stable if

$$\sup_{i \in [m]} \left\{ \mathbb{E}_S \mathbb{E}_{S, z} [f(A_S, z) - \ell(A_{S^{(i)}}), z] \right\} \leq \epsilon(D, w_1),$$

where  $S \stackrel{iid}{\sim} D$  and  $S^{(i)}$  is  $S$  copy with  $i$ -th example replaced by  $z \stackrel{iid}{\sim} D$ .

We now have the following theorem (Kuzborskij and Lampert, 2018):

**Theorem 1.** Let  $A$  be  $\epsilon(\theta)$  on average stable, then

$$\mathbb{E}_S \mathbb{E}_A [R(A_S) - \hat{R}_S(A_S)] \leq \epsilon(D, w_1),$$

where  $R(A_S), \hat{R}_S(A_S)$  are risk and empirical risk of  $A_S$  respectively, defined by:

$$R(A_S) = \mathbb{E}_{z \sim D} [\ell(A_S, z)]; \quad \hat{R}_S(A_S) = \frac{1}{m} \sum_{i=1}^m \ell(A_S, z_i).$$

In other words, the generalization of a learning algorithm on unseen data drawn from the same distribution is controlled by its  $\epsilon(D, w_1)$ -on average stable, which depends on initialized weights  $w_1$ . In the following, we will examine the model’s generalization performance through the lens of its training algorithm’s stability.

In the transfer learning setting, Theorem 1 provides a tool to understand the model’s generalization on the target domain, given that it is initialized from one of the pre-trained models on a set of source domains (source tasks). Specifically, we suppose that the target task is characterized by a joint probability distribution  $D^{\text{tgt}}$  and assume that a set of source hypotheses  $\{w_k^{\text{src}}\} \subset H, k \in K$  trained on  $K$  different source tasks. In this paper, we consider two distinct source cases with  $K = \{\text{ImgNet}, \text{TATL}\}$  where ‘‘ImageNet’’ refers to weights trained on ImageNet and ‘‘TATL’’ is our approach of learning the attribute-agnostic mask. Now we are ready to analyze the generalization bound of TATL versus ImageNet initialization strategies by utilizing the results in Kuzborskij and Lampert (2018).

**Proposition 1.** Given a non-convex loss function  $\ell$  and assume that  $\ell(\cdot, z) \in [0, 1]$  has a  $p$ -Lipschitz Hessian,  $\beta$ -smooth and step sizes of a form  $\alpha_t = \frac{c}{t}$  satisfy  $c \leq \min(\frac{1}{\beta}, \frac{1}{4(2\beta \ln(T))^2})$ , then with high probability, the  $\epsilon(D^{\text{tgt}}, w_k^{\text{src}})$  of SGD scheme satisfies:

$$\min_{k \in K} \epsilon(D^{\text{tgt}}, w_k^{\text{src}}) \leq \min_{k \in K} \mathcal{O} \left( \left( 1 + \frac{1}{c \hat{\gamma}_k^-} \right) \hat{R}_S(w_k^{\text{src}})^{\frac{c \hat{\gamma}_k^+}{1 + c \hat{\gamma}_k^+}} \cdot \frac{\sqrt{\log(|K|)}}{m^{\frac{1}{1 + c \hat{\gamma}_k^+}}} \right), \quad (6)$$

where

$$\hat{\gamma}_k^\pm = \hat{\gamma}_k \pm \frac{1}{\sqrt{m}}, \quad (7)$$

$$\hat{\gamma}_k = \frac{1}{m} \sum_{i=1}^m \|\nabla^2 \ell(w_k^{\text{src}}, z_i)\|_2 + \sqrt{\hat{R}_S(w_k^{\text{src}})}, \quad (8)$$

with  $\|\cdot\|_2$  is a spectral norm. Intuitively, Theorem 1 and Proposition 1 suggest that an initialization’s generalization error depends on *two* factors: (i) how well it performs on the target domain without any training, which is characterized by  $\hat{R}_S$ ; and (ii) the loss function’s curvature around this initialization, which is characterized by empirical  $\|\nabla^2 \ell(w_k^{\text{src}}, z)\|_2$  over  $m$  training samples, denoted as  $\hat{\gamma}$ . This result provides an intuitive explanation of why TATL provides a more favorable initialization than the traditional ImageNet pre-trained models. Notably, we will explain why TATL, which initializes the Target-Segmenter from the Attribute-Agnostic Segmenter, can perform better than initializing the segmenter from ImageNet pre-trained models.

Pre-trained ImageNet models are unlikely to perform well on medical images due to the vast diversity between the two domains. Therefore, such models often have a higher empirical error on the target domain  $\hat{R}_S$  and usually lie in *high curvature* regions  $\hat{\gamma}$ . On the other hand, TATL uses an initialization from the Attribute-Agnostic Segmenter, which is pre-trained on self-generated data of the target task. Note that the Attribute-Agnostic Segmenter can detect *any* of the attributes, and therefore enjoy *lower* empirical risk  $\hat{R}_S$  compared to ImageNet models. Moreover, due to its construction, the Attribute-Agnostic Segmenter’s parameter likely lies in a region close to the local minimum of each attribute detector, which enjoys *lower curvature*  $\hat{\gamma}$ . Consequently, TATL exploits the target task’s knowledge to form an initialization with a high probability of attaining lower empirical error and curvature, which translates to a tighter generalization error bound than initializing from pre-trained ImageNet models. We empirically verify this hypothesis by comparing the bound’s values in Eq. (6) of different initialization strategies in Figure 7 of Section 5.8.

## 5. Experiments and Results

### 5.1. Dataset

We conduct experiments on two well-known datasets for skin attributes detection: the ISIC 2017<sup>3</sup> and 2018<sup>4</sup> Task 2 datasets. Table 2 provides a summary of the two datasets. It is worth noting that the ISIC 2017 dataset only contains *four* classes: Streaks, Negative Network, Milia, and Pigment Network, while the ISIC 2018 introduces a new class of Globules. Moreover, both datasets exhibit high data imbalance among the attributes. For example, in the ISIC 2018 dataset, the class “Streaks” only appears in 3.86% of the training data while “Pigment Network” is observed in 58.67% of the training data.

Table 2: Distribution of training images in the ISIC 2017 and 2018 Task 2 datasets.

	Class	Streaks	Negative Network	Globules	Milia	Pigment Network	Total
<b>ISIC - 2017</b>	Number	113	122	NA	475	1119	1416
	Rate (%)	7.98	8.62	NA	33.55	79.03	100%
<b>ISIC - 2018</b>	Number	100	190	602	681	1522	2594
	Rate (%)	3.86	7.32	23.21	26.25	58.67	100%

### 5.2. Experimental Settings

We implemented all experiments using the Pytorch framework (Paszke et al., 2019) on 4 NVIDIA TITAN RTX GPUs. All images were pre-processed by centering and normalizing the pixel density per channel. Besides, we also re-scale all images to the resolution of  $386 \times 512$  in training steps and transform final predictions to a size of  $256 \times 256$  in the evaluation step following the standard of the ISIC challenge. We used the SGD optimizer (Goodfellow et al., 2016) with an initial learning rate of 0.01 and momentum of 0.9 to be consistent with the theory presented in Section 4.5. For TATL, we obtained the Segment-Net by training a b0-EfficientNet backbone with *U-shape* on both the ISIC 2018 and ISIC 2017 Task 1 using the loss function in Eq. (1). Given the segmentation results, we defined a bounding box around the masks with an offset of 40 pixels in four directions to mitigate the segmentation errors, before feeding them to the Attribute-Agnostic Segmenter and

<sup>3</sup><https://challenge.isic-archive.com/landing/2017>

<sup>4</sup><https://challenge2018.isic-archive.com/>

the Target-Segmenters. The Attribute-Agnostic Segmenter and the Target-Segmenters were then trained for 40 epochs with early-stopping after 10 epochs. Finally, we measure our performance and compare it with other baselines using the five-fold cross-validation method and report the average values on Dice and Jaccard coefficients as Koohbanani et al. (2018).

### 5.3. Comparison Against Other Approaches

We compare our method against the winner of ISIC 2017 (Kawahara and Hamarneh, 2018) and ISIC 2018 (Koohbanani et al., 2018) and report the Dice and Jaccard index in Table 4 and Table 3 respectively. Here the results of the 1-st method in ISIC 2018 are extracted from their original paper (Koohbanani et al., 2018) while we use the source code published by (Kawahara and Hamarneh, 2018) and run-again experiments using our setting with the five-fold cross-validation method.

Table 3: Jaccard and Dice metrics on the ISIC2018 challenge. Blue and Red colors are the best results in Jaccard and Dice for each attribute, respectively. Stage 1: segmenting the lesion region, Stage 2: training the Attribute-Agnostic Segmenter, Stage 3: training the Target-Segmenters.

Method	ISIC2018 1st		TATL Stage 3		TATL Stage 2, 3		TATL Stage 1, 3		TATL Stage 1, 2, 3		TATL w/o ImgNet Stage 1, 2, 3	
	Jaccard	Dice	Jaccard	Dice	Jaccard	Dice	Jaccard	Dice	Jaccard	Dice	Jaccard	Dice
<b>Pigment Net.</b>	0.563	0.720	0.532	0.691	0.580	0.721	0.542	0.701	0.584	0.730	0.565	0.721
<b>Globules</b>	0.341	0.508	0.308	0.471	0.368	0.546	0.332	0.467	0.379	0.552	0.359	0.528
<b>Milia-like cysts</b>	0.171	0.289	0.141	0.252	0.161	0.268	0.142	0.264	0.172	0.288	0.161	0.277
<b>Negative Net.</b>	0.228	0.371	0.149	0.260	0.269	0.403	0.194	0.348	0.283	0.438	0.280	0.437
<b>Streaks</b>	0.156	0.270	0.139	0.241	0.254	0.394	0.135	0.224	0.254	0.401	0.263	0.416
<b>Average</b>	0.292	0.432	0.254	0.383	0.326	0.466	0.269	0.401	0.334	0.482	0.326	0.476

Table 4: Jaccard and Dice metrics on the ISIC2017 challenge. Blue and Red colors are best results in Jaccard and Dice for each attribute, respectively. Stage 1: segmenting the lesion region, Stage 2: training the Attribute-Agnostic Segmenter, Stage 3: training the Target-Segmenters.

Method	ISIC2017 1st		TATL Stage 3		TATL Stage 2, 3		TATL Stage 1, 3		TATL Stage 1, 2, 3		TATL w/o ImgNet Stage 1, 2, 3	
	Jaccard	Dice	Jaccard	Dice	Jaccard	Dice	Jaccard	Dice	Jaccard	Dice	Jaccard	Dice
<b>Pigment Net.</b>	0.389	0.556	0.426	0.597	0.499	0.667	0.497	0.665	0.516	0.681	0.473	0.639
<b>Milia-like cysts</b>	0.119	0.215	0.072	0.127	0.091	0.157	0.101	0.172	0.108	0.188	0.092	0.168
<b>Negative Net.</b>	0.201	0.333	0.126	0.218	0.191	0.310	0.147	0.251	0.213	0.334	0.234	0.380
<b>Streaks</b>	0.192	0.321	0.139	0.233	0.203	0.336	0.139	0.232	0.215	0.346	0.209	0.345
<b>Average</b>	0.225	0.356	0.191	0.294	0.246	0.367	0.221	0.330	0.263	0.387	0.252	0.383

### TATL Variations

Due to the high competitiveness of skincare challenges, we utilize both U-shape and Link-shape architectures with b0-EfficientNet (Tan and Le, 2019) as the backbone for the Attribute-Agnostic Segmenter and Target-Segmenters, then taking the average probability predictions. For a comprehensive comparison, we include *five* variants of four TATL corresponding to removing components in the TATL framework:

- the vanilla encoder-decoder architecture but without the Segment-Net and the Agnostic-Attribute Segmenter (*TATL Stage 3*);
- a variant that performs the second and last stage of our TATL: train first the Attribute-Agnostic segmenter on the original images and then a set of Target-Segmenters (*TATL Stage 2, 3*);
- a variant that performs the first and last stage of our TATL: segment the lesion regions and then the attributes (*TATL Stage 1, 3*);
- our full TATL framework that performs all three stages (*TATL Stage 1, 2, 3*);
- our full TATL framework initialized from scratch (*TATL w/o ImgNet Stage 1, 2, 3*).

### Overall performance

Our TATL shows competitive performances against other approaches on both benchmarks and metrics. Notably, our method presents substantial improvements over other baselines on attributes with the least amount of training data. For example, in Table 3, **TATL Stage 1, 2, 3** for Negative Network achieved the Jaccard index of 0.283, which is 5.5% higher compared to ISIC 2018 winner (0.228) (Koohbanani et al., 2018). Also, this TATL setting for the Streak feature has the Dice index of 0.401, which is 13.1% higher compared to ISIC 2018 winner (0.270) (Koohbanani et al., 2018).

### Ablation Study of TATL’s Sub-stages

This experiment aims to investigate the contribution of the first and second stages to the final performance, given that Stage 3 is always enabled for supervised training. Table 3 and Table 4 demonstrate that enabling Stage 2, i.e., **TATL Stage 2, 3** results in a significant improvement in most skin attributes compared to using Segment-Net (**TATL Stage 1, 3**). For example, the Average Dice score in Table 3 of **TATL Stage 3** increased from 0.383 to 0.466 with **TATL Stage 2, 3** while only attaining at 0.401 with **TATL Stage 1, 3**. This result thus emphasizes the critical role of learning the Attribute-Agnostic Segmenter in our framework.

In addition, progressively adding Stage 1 to the Stage 2 and 3 model further improves the results. For example, in Table 4, when using a pre-trained ImageNet, adding Stage 1 increased the Averaged Dice score from 0.367 to 0.387. Overall, the findings support our hypothesis that all three phases contribute to TATL’s competitive performance.

### The Influence of using Pre-trained ImageNet in TATL

We investigate the advantages of employing pre-trained models in TATL (Algorithm 1) by examining the *TATL w/o ImgNet Stage 1, 2, 3*. The outcomes suggest that employing pre-trained models in TATL favours classes with more training data. Remarkably, for ISIC 2017’s Pigment Network, the attribute with the most training samples, pre-trained models improves the Jaccard index from 0.473 to 0.516 (9.09% relative improvement). On the other hand, the contributions of pre-trained models on attributes with limited training samples such as Streaks and Negative Network are much less significant, e.g., ISIC 2017’s Negative Network Jaccard index increased from 0.209 to 0.215 (2.87% relative improvement). In general, we conclude that the improvements observed in minor classes come from our TATL framework rather than the pre-trained ImageNet; however, the TATL version with pre-trained ImageNet, on average, outperforms the TATL version using random weights.

Table 5: The average performance of the skin feature segmentation step with various lesion inputs generated by various networks and offset settings. Dice score is used to measure accuracy.

Method	Performance of Segmented Skin Lesion	Performance of Segmented Skin Features			
		Offset = 0	Offset = 20	Offset = 40	Offset = 60
U-Net	0.783	0.358	0.362	0.387	0.376
SegNet	0.824	0.365	0.371	0.387	0.376
Mask-RCNN	0.876	0.371	0.372	0.387	0.376

### The Influence of Cropping Lesion Segmentation on the Final Result

Our inference step requires the segmented lesion region to eliminate less relevant parts for the later phase. This task is handled by the Segment-Net (Stage 1). The lesion regions then are cropped by a bounding box with an offset value in four directions to generate input for the next step. In Table 3 and Table 4, we presented the ablation study for Stage 1 derived from *U-Shape* using a b0-EfficientNet backbone with an offset of 40 pixels. We now investigate how much the errors in Stage 1 can propagate to the final predictions by varying *segmentation methods* and *offset values*.

We conducted tests on ISIC 2017 in which four different models were trained to segment four skin features using the same configuration in Stage 2 and 3 but will take distinct inputs in Stage 1 created by various networks as U-Net (Ronneberger et al., 2015b), SegNet (Badrinarayanan et al., 2017), and Mask-RCNN (He et al., 2017). In addition, we changed offset values ranging from 0 to 60 pixels with a 20-pixel step. Table 5 depicts the results of various approaches when changing these factors, in which we used the Dice score to compute accuracy for all experiments.

We observe that applying an offset 0 pixel reduces the accuracy of the subsequent step because segmentation errors lead to the loss of some essential information, particularly at the image’s border locations. It also explains why improved performance in the feature segmentation stage results from higher accuracy in the lesion segmentation step. When increasing the magnitude of offset to 20 pixels, all methods are improved; for instance, the U-Net case rises from 0.358 to 0.362. Moreover, with the offset 40 and 60 pixels, margins between baselines are no longer available, and the performance of two later cases is better than offset 20 pixels. Table 5 also



presents a trade-off in selecting large offset values. In particular, a considerable value of 60 pixels can reduce efficiency compared to 40 pixels because the image may involve more unrelated data. In summary, we conclude that adding an optimized offset value of 40 pixels around segmented lesion areas allows our inference step to be stable despite lesion segmentation perturbations.

#### Performance of Each Network Backbone

This experiment examines the contribution of each component network to the overall performance of TATL. In particular, we evaluate our method’s performance using just U-shape or Link-shape based on the b0-EfficientNet backbone and compare them to networks used in the ISIC-2018-1st: ResNet-151, Resnet-v2, and DenseNet-169. Furthermore, we also include the b0-EfficientNet performance in the ISIC-2018 challenge for overall comparison. Table 6 highlights the acquired data, with blue and red representing the best Jaccard and Dice scores, respectively. In this table, our two variants, labeled as *U-Eff(TATL)* and *L-Eff(TATL)*, are the results of employing the U-shape and the Link-shape, respectively.

Table 6: The performance of ISIC2018-1st using a single backbone compared to our framework on the ISIC 2018. Blue and Red colors are the best values in Jaccard and Dice metrics.

Method	Pigment Net.		Globules		Milia-like cysts		Negative Net.		Streaks		Average	
	Jaccard	Dice	Jaccard	Dice	Jaccard	Dice	Jaccard	Dice	Jaccard	Dice	Jaccard	Dice
ResNet-151	0.527	0.690	0.304	0.466	0.144	0.257	0.149	0.260	0.125	0.222	0.245	0.379
ResNet-v2	0.539	0.706	0.310	0.473	0.159	0.274	0.189	0.318	0.121	0.216	0.264	0.397
DenseNet-169	0.538	0.699	0.324	0.490	0.158	0.273	0.213	0.351	0.134	0.236	0.273	0.410
b0-EfficientNet	0.554	0.713	0.324	0.497	0.157	0.272	0.213	0.351	0.139	0.242	0.277	0.415
U-Eff (TATL)	0.565	0.722	0.373	0.549	0.157	0.271	0.268	0.421	0.243	0.390	0.321	0.471
L-Eff (TATL)	0.562	0.719	0.356	0.522	0.168	0.287	0.292	0.452	0.252	0.393	0.326	0.475

We can see that all of the top results came from one of our techniques, with the U-shape architecture gaining first place in the Pigment Network and Globules features and the Link-shape architecture taking second. The baseline with the b0-EfficientNet backbone, on the other hand, appeared to outperform the other approaches. Considering skin characteristics with a large amount of training data, such as Pigment Network (79.03%) or Milia-like cysts (33.55%), our model enhanced the Jaccard of the b0-EfficientNet from 55.4% to 56.5% and from 15.7% to 16.8%, respectively. On Pigment Network, the Dice coefficient was improved by 0.9%, and on Milia-like cysts, it was improved by 1.6%. The lower the number of images, the greater the margin of improvement made by our model through the transfer learning step in TATL. For example, with the Streaks, our L-Eff(TATL) was 25.2% and 39.3% in Jaccard and Dice, correspondingly, which were 11.3% and 15.1% higher than the best baseline results with EfficientNet backbone.

Overall, TATL with the Link-shape structure performed the best across all network backbones, followed by TATL with the U-shape with a minor margin. Furthermore, these configurations surpass all remaining baselines with a large margin, thereby proving the benefit of using TATL.

#### Comparison of Network Parameters

While achieving state-of-the-art performances, our method enjoys a significant reduction in the number of parameters. We provide the number of trainable parameters on different architectures in Table 7. Notably, compared to the winner of ISIC 2017 (Kawahara and Hamarneh, 2018) and ISIC 2018 challenge (Koohbanani et al., 2018), our method has **1.4 to 2.33 times** and **30 to 50 times** fewer parameters during training. Consequently, our TATL consumes less GPU memory and thus can be trained with higher image resolution or employed in mobile devices with low memory costs.

#### 5.4. TATL Works Well with Various Network Architectures

In Section 5.3, we demonstrated that TATL achieved promising results using the b0-EfficientNet backbone. In this experiment, we explore the robustness of TATL to different network architectures beyond EfficientNets. Particularly, we compare ImageNet initialization against TATL on *five* different backbone networks (VGG16, ResNet151, ResNet-v2, DenseNet-169, and EfficientNet-b0) and evaluate the performances on the Negative and Streaks attributes because they are the most challenging ones with the least training samples. In addition, we consider the following settings:

- The first one, denoted as *TATL (FE)*, was to apply the TATL technique but froze the encoder part and only update weights of the decoder module while training for a specific skin attribute.
- The second configuration, denoted as *TATL (NF)*, was similar to the former but allowed updating the parameters in the encoder.

Table 7: Number of parameters in each architecture in the ISIC challenge.

Architecture	Number of Parameters
VGG16	138,357,544
ResNet-151	60,419,944
ResNet-v2	55,873,736
DenseNet-169	14,307,880
EfficientNetb0	5,330,571
ISIC2018-1st	308,747,840
ISIC2017-1st	14,780,929
<b>Our (EfficientNet, U-shape)</b>	10,115,501
<b>Our (EfficientNet, Link-shape)</b>	6,096,333

Table 8: The average Dice coefficients of different backbones using our method and the transfer learning based on ImageNet using the datasets ISIC 2017 and 2018. The bold and red-marked values are the highest and second-highest of each architecture. FE indicates freezing the encoder, NF indicates non-freezing and training both the encoder and decoders, L-Shape: Link-Shape.

Architecture	Methods	ISIC 2017				ISIC 2018			
		Negative		Streaks		Negative		Streaks	
		U-shape	L-Shape	U-shape	L-Shape	U-shape	L-Shape	U-shape	L-Shape
<b>VGG-16</b>	TATL (FE)	<b>0.273</b>	0.178	<b>0.281</b>	0.279	<b>0.303</b>	0.225	0.266	<b>0.272</b>
	TATL (NF)	<b>0.253</b>	0.191	<b>0.282</b>	0.278	<b>0.309</b>	0.252	0.262	<b>0.279</b>
	ImageNet	0.231	0.232	0.254	0.225	0.235	0.244	0.207	0.254
<b>ResNet151</b>	TATL (FE)	0.231	<b>0.279</b>	0.283	<b>0.285</b>	<b>0.357</b>	0.33	<b>0.317</b>	0.288
	TATL (NF)	0.248	<b>0.289</b>	0.281	<b>0.289</b>	<b>0.344</b>	0.315	<b>0.324</b>	0.286
	ImageNet	0.239	0.244	0.175	0.194	0.275	0.235	0.201	0.174
<b>ResNet-v2</b>	TATL (FE)	0.294	0.256	<b>0.293</b>	0.248	0.445	0.406	<b>0.324</b>	0.299
	TATL (NF)	<b>0.30</b>	0.279	<b>0.299</b>	0.252	<b>0.458</b>	0.413	<b>0.329</b>	0.300
	ImageNet	<b>0.313</b>	0.28	0.226	0.198	0.432	<b>0.460</b>	0.235	0.240
<b>DenseNet169</b>	TATL (FE)	<b>0.339</b>	0.284	0.299	<b>0.303</b>	0.292	<b>0.367</b>	0.338	<b>0.368</b>
	TATL (NF)	<b>0.307</b>	0.288	0.296	<b>0.306</b>	0.353	<b>0.389</b>	0.342	<b>0.346</b>
	ImageNet	0.241	0.227	0.205	0.194	0.285	0.377	0.210	0.216
<b>Eff. Net-b0</b>	TATL (FE)	0.289	0.295	<b>0.346</b>	0.321	0.401	<b>0.445</b>	0.384	<b>0.395</b>
	TATL (NF)	<b>0.297</b>	<b>0.355</b>	<b>0.334</b>	0.332	0.421	<b>0.440</b>	0.359	<b>0.410</b>
	ImageNet	0.286	0.218	0.259	0.219	0.355	0.392	0.199	0.220

- The last setting (*ImageNet*) was not to apply the transfer learning process and train from scratch using weights pre-trained on the ImageNet dataset.

For each model, we use five different backbone architectures and two different convolution network shapes (U-shape and Link-shape). We also rerun three times for each configuration and measure the corresponding performance with the five-fold cross-validation technique to estimate average results. This configuration results in a total of **600 experiments** to be examined, which provides a comprehensive analysis of our TATL.

Table 8 reports the results of the experiment and shows that applying TATL could help improve all backbone performance except the ResNet-v2 with the Negative attribute. However, the difference between the Dice values, in this case, was not noticeable (less than 1%). In contrast, the TATL could boost the Dice by nearly 13% when using DenseNet-169 with U-shape to segment Streaks regions in the ISIC2018 dataset, and 11.4% when using ResNet151 with Link-shape in a similar task. On the backbones such as DenseNet-169 + U-shape, ResNet151 + Link-shape, TATL consistently provided significant improvements. In summary, we find that our proposed TATL transfer learning could operate with various network architectures, thereby demonstrating its efficacy and generalizability.

### 5.5. The Stability of TATL under The Influences of Data Size

Using pre-trained ImageNet models have been a common practice for many computer vision applications because of the diversity in the ImageNet dataset, making the pre-trained models stable and can detect local relationships. In contrast, TATL does not use additional data sources and only relies on the dataset at hand to create and learn the skin attribute regions. As a result, TATL’s performance is subject to the amount of training data for the current task. Nevertheless, the results in Section 5.3 showed that TATL had achieved promising results on standard benchmarks using all the labeled data provided. In this section, we explore the stability of TATL under the effect of different data size.

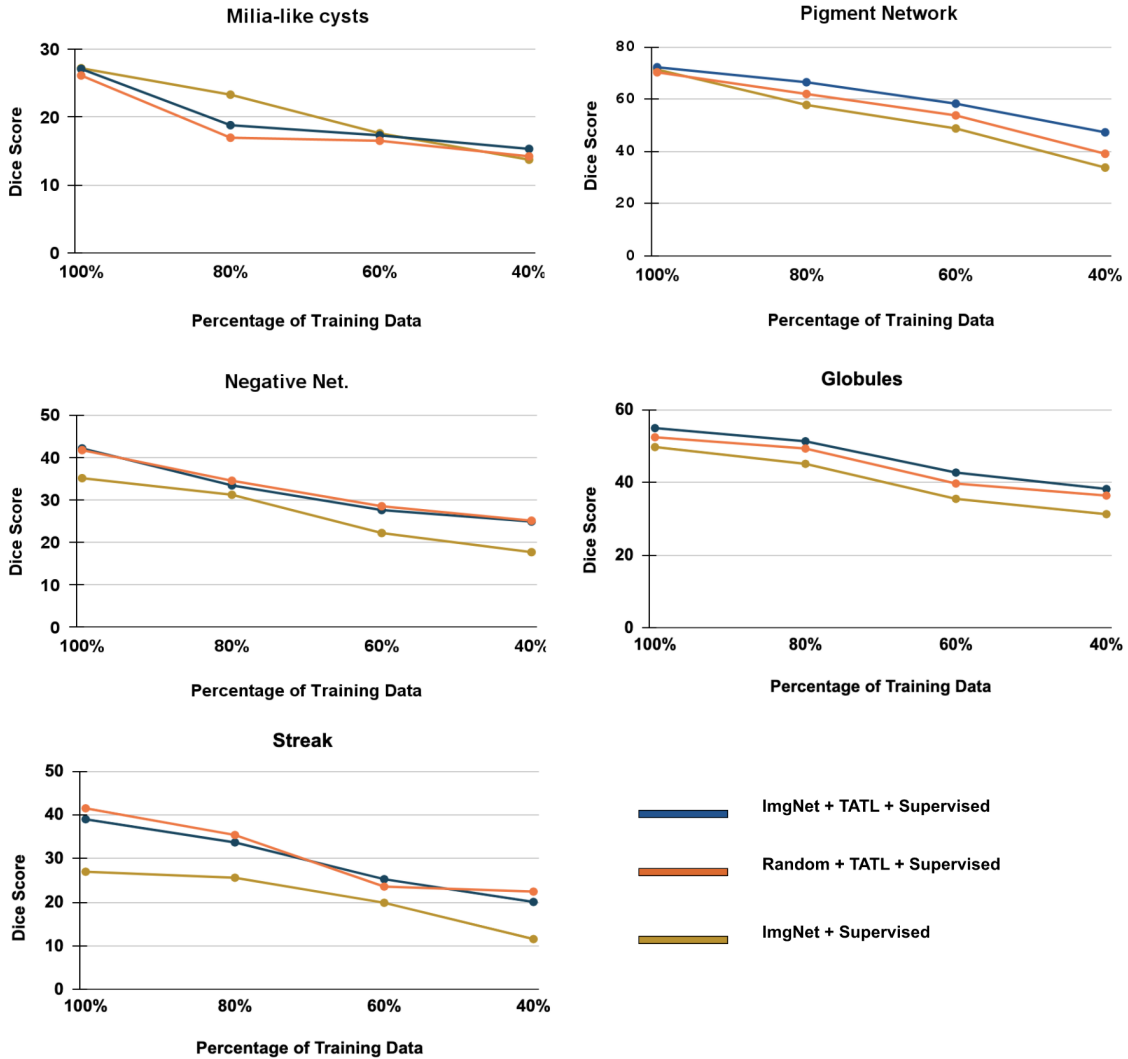


Figure 4: The stability of TATL on five skin lesion attributes in ISIC 2018 under various sample sizes.

To test the stability of TATL, for each skin feature, its testing data are reserved for evaluation, and we subsample a portion of the remaining data for training. Particularly, we vary the amount of training data for each skin feature from 100% (the original experiments in Section 5.3) to only 40% and consider three competitors:

- TATL using a pre-trained ImageNet model;
- TATL with a random initialization;
- Only a pre-trained ImageNet model.

We report the results in Figure 4. In most cases, reducing the amount of training data results in worse performances for all methods. Moreover, both TATL versions achieved better performances than only using a pre-trained ImageNet model, except for the Millia feature at 80% training data. Interestingly, the two TATL versions achieved similar performances on skin attributes with limited training data (Negative Net and Streak), while the gap was larger on the attributes with more training data (Pigment Network). However, overall, both TATL variants achieved better results than the standard strategy of using a pre-trained ImageNet. Furthermore, TATL with a pre-trained ImageNet model achieved the best-averaged results across attributes.

### 5.6. Exploring TATL's Features Generality and Convergence Rate

As discussed in Section 3.2, our TATL is inspired by the dermatologists' behaviors, which motivates the learning of skin attribute regions. In this experiment, we examine the relationship between skin attribute regions' features learned by the Attribute-Agnostic Segmenter and the specific features learned for each particular attribute of Target-Segmenters. We hypothesize that there are additional benefits of performing supervised learning on attributes whose features are similar to those obtained from the Attribute-Agnostic learning step. To validate this hypothesis, we consider two TATL variants:

- *Downstream TATL*: The standard TATL framework (Figure 3), where the model for each skin attribute is initialized from the Attribute-Agnostic step;
- *Pretext TATL*: A TATL variant which directly uses the model trained in the Attribute-Agnostic learning step to inference for each skin attribute, *without additional supervised learning in downstream task* (Figure 3 with only the left block).

We run experiments on the ISIC 2017 benchmark with the two aforementioned TATL variants and report the Dice score for each attribute after the learning procedures in Figure 5 as well as the loss curves in Figure 6 with respect to the number of training steps. Note that the Pigment Network enjoys the most training samples amongst the four attributes, while the remaining three, Streaks, Negative Network, and Milia, are considered the minority classes.

Figure 5 shows that for the Pigment Network feature, the improvement of *Downstream TATL* compared to *Pretext TATL* is almost neglectable, which shows that the benefit of additional supervised training is minor for the attributes with ample training data. In contrast, the minor attributes’ performance gaps are much more significant, suggesting TATL brings more benefits to the classes with limited training samples. We further verify this result in Figure 6 where the validation loss for the Pigment Network attribute plateaus very early on (at around 300 iterations) while it further decreases for other minor attributes. Overall, these experiments show that the downstream fine-tuning step is particularly beneficial for the small classes. Therefore, one can infer that TATL, with all three stages, can significantly boost the performance on the classes with limited training data, which are more challenging to improve performance.

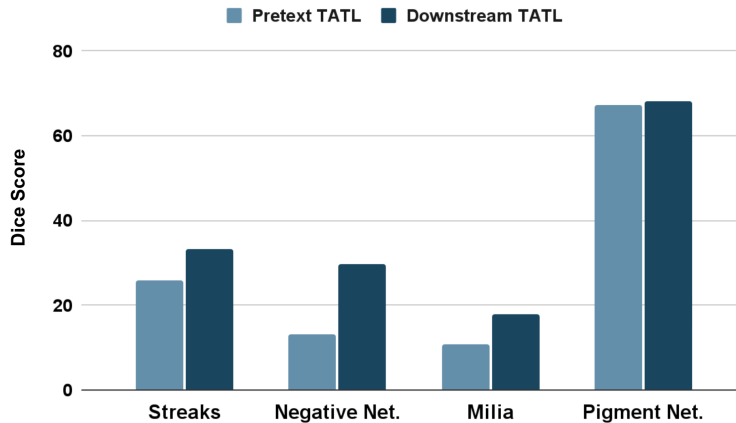


Figure 5: The performance comparison of four skin lesion attributes in ISIC 2017 using (i) TATL trained directly on Pretext Task, and (ii) TATL after trained on the Downstream Task.

### 5.7. Comparison Between TATL and Other Training Paradigms

Table 9: Average performance over different skin attributes on ISIC-2018 with TATL compared to other training strategies. In each encoder-decoder type, Blue and Red colors are the best results in Jaccard and Dice coefficients, respectively.

Methods	U-shape		Link-shape	
	Dice	Jaccard	Dice	Jaccard
TATL	<b>0.471</b>	<b>0.321</b>	<b>0.475</b>	<b>0.326</b>
Pretrained ImageNet	0.422	0.286	0.419	0.285
Rotation-based SSL	0.429	0.283	0.417	0.284
Image Context-based SSL	0.446	0.298	0.426	0.301
Data Augmentation	0.435	0.294	0.403	0.276
Attention-based Method	0.441	0.296	0.416	0.285

In this experiment, we explore the benefits of TATL compared to other training paradigms, namely self-supervised learning (SSL), training with data augmentation, and attention-based methods. Note that our TATL is related to the SSL paradigm in terms of deriving better pre-trained models through solving auxiliary tasks. For the SSL baselines, we implement an additional pre-training phase to replace TATL’s first and second stages. Particularly, we consider the task of predicting the rotation angle applied on the input (Gidaris et al., 2018) or reconstructing the image after scrambling the pixels (Gidaris et al., 2018). We also consider the strategy of

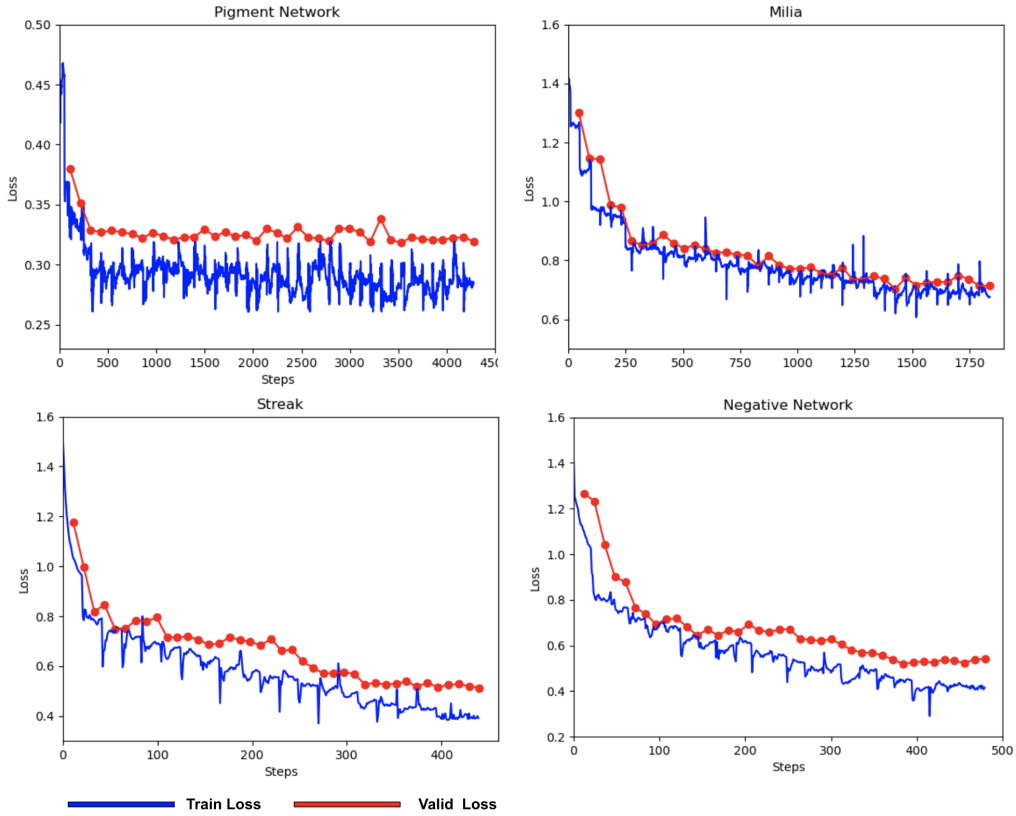


Figure 6: The training and validation curves of four skin characteristics in ISIC 2017 with respect to the number of training steps.

supervised training with data augmentation to increase the total training samples, and using an attention-based architecture which can suppress irrelevant regions in an input image while focusing salient features useful for a specific task (Oktay et al., 2018).

In summary, given the same U-shape architecture and pre-trained b0-EfficientNet as the network backbone, we compare TATL against the standard Pre-trained ImageNet models and the following training strategies:

- Self-supervised with the image rotation-based method (Gidaris et al., 2018);
- Self-supervised with the image context restoration method (Chen et al., 2019);
- Supervised training with data-augmentations, in which we use random rotation, flip, shift, brightness, or zoom;
- U-Eff network with attention gates as proposed in Oktay et al. (2018).

Table 9 presents the experiment results with the metrics computed by averaging all skin attributes in the ISIC-2018 challenge. In general, both the attention approach and the image reconstruction SSL can provide marginal improvements to the traditional Imagenet initialization on the U-Shape design. However, our TATL still outperforms such strategies on both evaluation metrics and architecture designs. This evidence confirms our finding that transferring knowledge from the Attribute-Agnostic Segmenter is beneficial for the skin-attribute segmentation task.

### 5.8. Generalization Bound of TATL Compared to Other Strategies

In Section 5.7, we demonstrated that TATL could outperform standard training paradigms such as using a pre-trained model and SSL methods. This section explores how the theoretical insights developed in Section 4.5 supports the empirical success of TATL. Recall that from Proposition 1, one can infer that given the same conditions of employed SGD algorithm and other hyper-parameters such as learning rate or the number of epochs, *a model is expected to achieve small generalization errors on a testing set if its errors on the corresponding training data measured at the point of initialization (without any supervised learning) is small.* Given that TATL achieved lower testing errors in experiments (Table 8 and 9), we now conduct a test to explicitly verify if this results correspond to a lowest TATL’s generalization error bound values compared with other strategies.

In particular, we estimate generalization error’s bound in the right side of Eq. (6) on the ISIC 2018 with four different cases: our TATL, pre-trained ImageNet, Rotation-based SSL (Gidaris et al., 2018), and Image



Context Reconstruction-based SSL (Chen et al., 2019). We do not compare with the Attention-based method and Data Augmentation approach since both use the same Pre-trained ImageNet. For each method, we use the U-Eff network and run a full pass over all training samples of each attribute to estimate the spectral norm of the Hessian matrix and the empirical risk  $\hat{R}_S$ . Here, the largest eigenvalue is approximated by the power iteration method (Solomon, 2015). We set  $K = 4$ ,  $c = 0.01$  for all attributes and present the relative relations among categories in Figure 7.

It is noteworthy that our TATL acquired the lowest generalization error values for all skin attributes, especially with Streaks and Negative Net. These observations are compatible with our experiment, in which we outperformed other training strategies (as shown in Table 9) and surpassed the Pre-trained ImageNet by a wide margin for the two characteristics Negative and Streak (Table 8, at *Eff. Net-b0* row and *ISIC 2018* column). In conclusion, we argue that the TATL’s efficacy could be demonstrated in both experimental and theoretical settings.

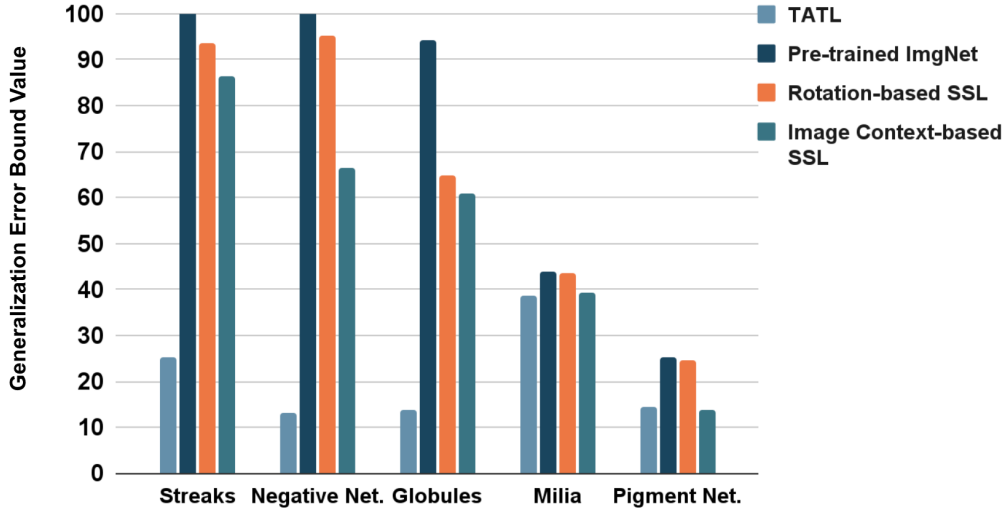


Figure 7: Generalization error bound value (lower is better) in Proposition 1 measured on ISIC-2018 for different initialization methods. Values are scaled with a factor of **100** for better visualization.

### 5.9. Visualization

Figure 8 illustrates some sample results of our proposed TATL model. The ground-truth segmentation was highlighted in green, and our prediction was marked with red. Regarding attributes with many training samples such as Globules or Pigment Network, TATL has a better segmentation covering most ground truth areas. Furthermore, although Streaks and Negative Network’s prediction missed some injured regions, the result still captured the primary matter location.

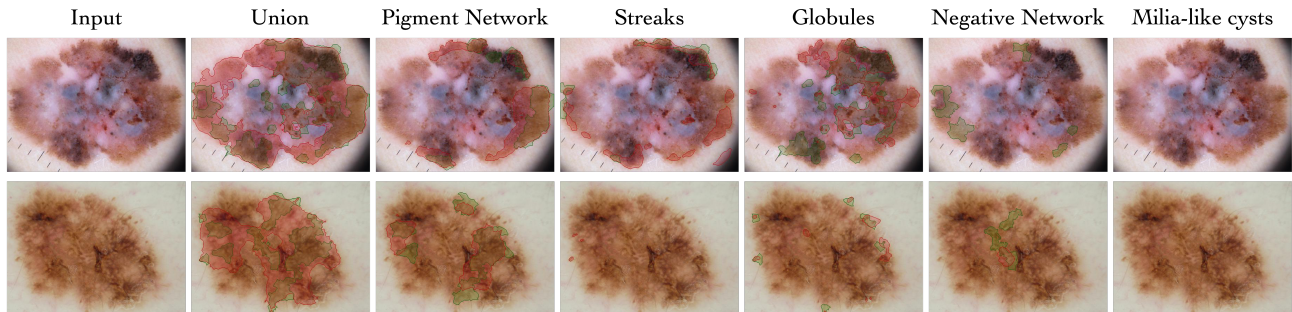


Figure 8: An example of the TATL results in two typical skin images (zoom in for better view) where the green indicates the ground truth, and the red illustrates our prediction.

Our model also provides the benefit of extra information for end-users through the predicted union regions. For instance, we show, in Figure 9, the correlation between the union regions with typical skin attributes for the same example in Figure 8. For each skin attribute, the ground truth is in green color, and we draw them over the binary maps to indicate the union positions’ prediction. It can be seen that predicted union could cover

both large regions as in Pigment Network and small disconnected regions as Negative Network (Figure 9a) and Globules (Figure 9b). This result can especially be useful under the scenarios where the dermatologists could not detect deformed or disconnected regions. In such cases, the union region provided by TATL can become helpful by highlighting the region of interest to assist the dermatologists. As a result, TATL outputs not only can speed up the diagnosis process and but also help the dermatologists diagnose better, which is critical because they, not our model, will make the final diagnosis.

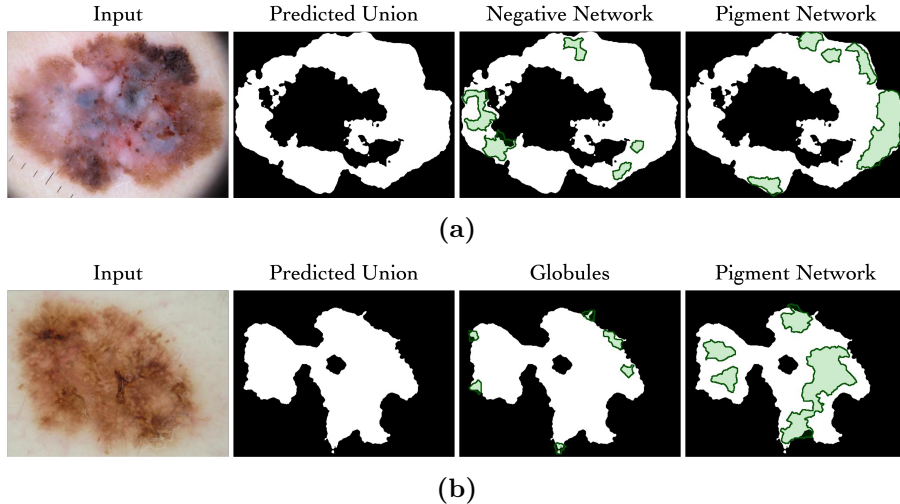


Figure 9: An illustration demonstrating the usefulness of our union segmentation, which can help dermatologists know which area they should focus on. Our predicted union can cover most of the ground truth of attribute regions which are depicted in the green areas.

## 6. Discussion and Future Work

Our work proposes a novel strategy to initialize the attribute segmenters’ parameters using an attribute-agnostic segmenter trained on abnormal skin regions. We empirically demonstrate this benefit over the traditional strategy of using the ImageNet pretrained models. From the promising results, we outline several potential and interesting directions for future research.

*Generalization to Other Medical Image Analysis Tasks.* We develop TATL to address the skin-attribute detection problem specifically. It would be interesting to test the TATL’s generalization capabilities to other medical image analysis tasks, where using pretrained Imagenet models is likely to be suboptimal. For example, similar tasks such as brain lesion segmentation (Hu et al., 2018; Duy et al., 2018; Mallick et al., 2019; Nguyen et al., 2017), abnormal chest detection (Hashir et al., 2020; Ibrahim et al., 2021; Nguyen et al., 2021), or diabetic retinopathy lesion segmentation share similar characteristics to our problem setting: the data are often imbalance and classes share semantic features that can be leveraged to improve the overall performance. Therefore, it is of interest to explore the applications of TATL in such tasks and make possible adjustments.

*Real-world Applications Using TATL.* Our ultimate goal is to develop a model that not only makes predictions but also provides helpful information and assists dermatologists in making the final decisions. Our TATL framework realizes this goal by providing a mask of skin attribute regions regardless of their attributes, compensating for inaccurate predictions of later stages, especially on minor classes. A promising future direction for TATL is integrating it in an online learning setting with human-in-the-loop (Nunnari and Sonntag, 2021; Nunnari et al., 2021a). Particularly, a model is trained to detect some diseases and then deployed to a real-world environment with a stream of data and feedback from dermatologists and patients. In such scenarios, the model can continuously improve its performance by accumulating the attribute-agnostic information via the dermatologists’ feedback and then transferring it to the target segmenters, allowing for a fast adaptation to newer patients and more accurate predictions over time.

*A Holistic Medical Image Analysis Method Beyond TATL.* Intuitively, TATL works by achieving a tighter generalization error bound compared to other initialization strategies. However, the theoretical result in Proposition 1 only bounds using the initialization parameters. In practice, additional aspects can affect the model’s generalization, such as (i) the number of source tasks (training classes in our case); (ii) which properties among those tasks that can be safely transferred; and (iii) beyond an initialization, which mechanisms allow for a successful knowledge transfer. Such properties are not yet rigorously studied, and exploring them can potentially



provide a holistic method for medical image analysis: a method not only starts with a quality initialization but also exploits the complex relationship of medical images to improve its performance over time. Such a method can provide accurate detection and assist dermatologists in diagnosing rare diseases more precisely, which results in effective treatments at a lower cost.

## 7. Conclusion

We have investigated the limitations of the common fine-tuning strategy in state-of-the-art skin attributes detection methods. We show that such strategies are not optimal when the current task is largely different from ImageNet and contains limited training data. This limitation motivated us to develop TATL, a novel transfer learning method that exploits all attribute data to train the agnostic segmenter. By transferring the agnostic segmenter’s knowledge to each attribute classifier, TATL alleviates issues of training data scarcity, especially for small classes, and allows knowledge sharing among attribute models. Through extensive experiments on the ISIC 2017 and ISIC 2018 benchmarks, we demonstrate the efficacy of TATL over existing state-of-the-art methods. Moreover, TATL is proven to work effectively with various backbone networks while enjoying minimal model and computational complexity. Finally, we present theoretical insights that demonstrate that TATL works in practice by bridging the domain gap via the task-agnostic segmenter, thus leading to competitive performance.

## 8. Acknowledgement

This research has been supported by the Ki-Para-Mi project (BMBF, 01IS1903-8B), the pAItient project (BMG, 2520DAT0P2), and the Endowed Chair of Applied Artificial Intelligence, Oldenburg University. Binh T. Nguyen is funded by Vietnam National University Ho Chi Minh City (VNU-HCM) under grant number NCM2019-18-01. We would like to thank Dr. Fabrizio Nunnari (German Research Centre for Artificial Intelligence, Germany) and Dr. Paul Swoboda (Max Planck Institute for Informatics, Germany) for their valuable discussions.

## References

- Abràmoff, M.D., Lou, Y., Erginay, A., Clarida, W., Amelon, R., Folk, J.C., Niemeijer, M., 2016. Improved automated detection of diabetic retinopathy on a publicly available dataset through integration of deep learning. *Investigative ophthalmology & visual science* 57, 5200–5206.
- Badrinarayanan, V., Kendall, A., Cipolla, R., 2017. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE transactions on pattern analysis and machine intelligence* 39, 2481–2495.
- Buda, M., Maki, A., Mazurowski, M.A., 2018. A systematic study of the class imbalance problem in convolutional neural networks. *Neural Networks* 106, 249–259.
- Celebi, M.E., Codella, N., Halpern, A., 2019. Dermoscopy Image Analysis: Overview and Future Directions. *IEEE Journal of Biomedical and Health Informatics* 23, 474–478. URL: <https://ieeexplore.ieee.org/document/8627921/>, doi:10.1109/JBHI.2019.2895803.
- Chaurasia, A., Culurciello, E., 2017. Linknet: Exploiting encoder representations for efficient semantic segmentation, in: *2017 IEEE Visual Communications and Image Processing (VCIP)*, IEEE. pp. 1–4.
- Chen, L., Bentley, P., Mori, K., Misawa, K., Fujiwara, M., Rueckert, D., 2019. Self-supervised learning for medical image analysis using image context restoration. *Medical image analysis* 58, 101539.
- Chen, T., Kornblith, S., Norouzi, M., Hinton, G., 2020. A simple framework for contrastive learning of visual representations, in: *International conference on machine learning*, PMLR. pp. 1597–1607.
- Cheplygina, V., 2019. Cats or cat scans: transfer learning from natural or medical image source data sets? *Current Opinion in Biomedical Engineering* 9, 21–27.
- Chollet, F., 2017. Xception: Deep learning with depthwise separable convolutions, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1251–1258.
- Codella, N., Rotemberg, V., Tschandl, P., Celebi, M.E., Dusza, S., Gutman, D., Helba, B., Kalloo, A., Liopyris, K., Marchetti, M., et al., 2019. Skin lesion analysis toward melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (isic). *arXiv preprint arXiv:1902.03368* .

- Codella, N.C., Gutman, D., Celebi, M.E., Helba, B., Marchetti, M.A., Dusza, S.W., Kalloo, A., Liopyris, K., Mishra, N., Kittler, H., et al., 2018. Skin lesion analysis toward melanoma detection: A challenge at the 2017 international symposium on biomedical imaging (isbi), hosted by the international skin imaging collaboration (isic), in: 2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018), IEEE. pp. 168–172.
- Curiel-Lewandrowski, C., Novoa, R.A., Berry, E., Celebi, M.E., Codella, N., Giuste, F., Gutman, D., Halpern, A., Leachman, S., Liu, Y., Liu, Y., Reiter, O., Tschandl, P., 2019. Artificial Intelligence Approach in Melanoma, in: Fisher, D.E., Bastian, B.C. (Eds.), *Melanoma*. Springer New York, New York, NY, pp. 1–31. URL: [https://doi.org/10.1007/978-1-4614-7322-0\\_43-1](https://doi.org/10.1007/978-1-4614-7322-0_43-1), doi:10.1007/978-1-4614-7322-0\_43-1.
- De Fauw, J., Ledsam, J.R., Romera-Paredes, B., Nikolov, S., Tomasev, N., Blackwell, S., Askham, H., Glorot, X., O’Donoghue, B., Visentin, D., et al., 2018. Clinically applicable deep learning for diagnosis and referral in retinal disease. *Nature medicine* 24, 1342–1350.
- Duy, N.H.M., Duy, N.M., Truong, M.T.N., Bao, P.T., Binh, N.T., 2018. Accurate brain extraction using active shape model and convolutional neural networks. *arXiv preprint arXiv:1802.01268* .
- Eelbode, T., Bertels, J., Berman, M., Vandermeulen, D., Maes, F., Bisschops, R., Blaschko, M.B., 2020. Optimization for medical image segmentation: theory and practice when evaluating with dice score or jaccard index. *IEEE Transactions on Medical Imaging* 39, 3679–3690.
- Gidaris, S., Singh, P., Komodakis, N., 2018. Unsupervised representation learning by predicting image rotations. *CoRR* abs/1803.07728. URL: <http://arxiv.org/abs/1803.07728>, [arXiv:1803.07728](https://arxiv.org/abs/1803.07728).
- Goodfellow, I., Bengio, Y., Courville, A., Bengio, Y., 2016. *Deep learning*. volume 1. MIT press Cambridge.
- Gulshan, V., Peng, L., Coram, M., Stumpe, M.C., Wu, D., Narayanaswamy, A., Venugopalan, S., Widner, K., Madams, T., Cuadros, J., et al., 2016. Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. *Jama* 316, 2402–2410.
- Hashir, M., Bertrand, H., Cohen, J.P., 2020. Quantifying the value of lateral views in deep learning for chest x-rays, in: *Medical Imaging with Deep Learning*, PMLR. pp. 288–303.
- He, K., Gkioxari, G., Dollár, P., Girshick, R., 2017. Mask r-cnn, in: *Proceedings of the IEEE international conference on computer vision*, pp. 2961–2969.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778.
- He, X., Yang, X., Zhang, S., Zhao, J., Zhang, Y., Xing, E., Xie, P., 2020. Sample-efficient deep learning for covid-19 diagnosis based on ct scans. *medRxiv* .
- Hu, Z., Tang, J., Wang, Z., Zhang, K., Zhang, L., Sun, Q., 2018. Deep learning for image-based cancer detection and diagnosis- a survey. *Pattern Recognition* 83, 134–149.
- Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q., 2017. Densely connected convolutional networks, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 4700–4708.
- Ibrahim, A.U., Ozsoz, M., Serte, S., Al-Turjman, F., Yakoi, P.S., 2021. Pneumonia classification using deep learning from chest x-ray images during covid-19. *Cognitive Computation* , 1–13.
- Jadon, S., 2020. A survey of loss functions for semantic segmentation, in: *2020 IEEE Conference on Computational Intelligence in Bioinformatics and Computational Biology (CIBCB)*, IEEE. pp. 1–7.
- Kawahara, J., Hamarneh, G., 2018. Fully convolutional neural networks to detect clinical dermoscopic features. *IEEE journal of biomedical and health informatics* 23, 578–585.
- Koohbanani, N.A., Jahanifar, M., Tajeddin, N.Z., Gooya, A., Rajpoot, N., 2018. Leveraging transfer learning for segmenting lesions and their attributes in dermoscopy images. *arXiv preprint arXiv:1809.10243* .
- Kuzborskij, I., Lampert, C., 2018. Data-dependent stability of stochastic gradient descent, in: *International Conference on Machine Learning*, PMLR. pp. 2815–2824.
- Larsson, G., Maire, M., Shakhnarovich, G., 2016. Learning representations for automatic colorization, in: *European conference on computer vision*, Springer. pp. 577–593.
- Lee, H.Y., Huang, J.B., Singh, M., Yang, M.H., 2017. Unsupervised representation learning by sorting sequences, in: *Proceedings of the IEEE International Conference on Computer Vision*, pp. 667–676.

- Li, X., Liu, S., De Mello, S., Wang, X., Kautz, J., Yang, M.H., 2019. Joint-task self-supervised learning for temporal correspondence, in: *Advances in Neural Information Processing Systems*, pp. 318–328.
- Mallick, P.K., Ryu, S.H., Satapathy, S.K., Mishra, S., Nguyen, G.N., Tiwari, P., 2019. Brain mri image classification for cancer detection using deep wavelet autoencoder-based deep neural network. *IEEE Access* 7, 46278–46287.
- Masood, A., Ali Al-Jumaily, A., 2013. Computer Aided Diagnostic Support System for Skin Cancer: A Review of Techniques and Algorithms. *International Journal of Biomedical Imaging* 2013, 1–22. URL: <http://www.hindawi.com/journals/ijbi/2013/323268/>, doi:10.1155/2013/323268.
- Misra, I., Zitnick, C.L., Hebert, M., 2016. Shuffle and learn: unsupervised learning using temporal order verification, in: *European Conference on Computer Vision*, Springer. pp. 527–544.
- Nguyen, D., Nguyen, D., Vu, H., Nguyen, B., Nunnari, F., Sonntag, D., 2021. An attention mechanism with multiple knowledge sources for covid-19 detection from ct images, in: *AAAI 2021 Workshop on Trustworthy AI for Healthcare*.
- Nguyen, D.M., Vu, H.T., Ung, H.Q., Nguyen, B.T., 2017. 3d-brain segmentation using deep neural network and gaussian mixture model, in: *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*, IEEE. pp. 815–824.
- Nguyen, D.M.H., Ezema, A., Nunnari, F., Sonntag, D., 2020. A visually explainable learning system for skin lesion detection using multiscale input with attention u-net, in: *German Conference on Artificial Intelligence (Künstliche Intelligenz)*, Springer. pp. 313–319.
- Nunnari, F., Alam, H.M.T., Sonntag, D., 2021a. Anomaly detection for skin lesion images using replicator neural networks, in: *International Cross-Domain Conference for Machine Learning and Knowledge Extraction*, Springer. pp. 225–240.
- Nunnari, F., Kadir, M.A., Sonntag, D., 2021b. On the overlap between grad-cam saliency maps and explainable visual features in skin cancer images, in: *International Cross-Domain Conference for Machine Learning and Knowledge Extraction*, Springer. pp. 241–253.
- Nunnari, F., Sonntag, D., 2021. A software toolbox for deploying deep learning decision support systems with xai capabilities, in: *Companion of the 2021 ACM SIGCHI Symposium on Engineering Interactive Computing Systems*, pp. 44–49.
- Oktay, O., Schlemper, J., Folgoc, L.L., Lee, M., Heinrich, M., Misawa, K., Mori, K., McDonagh, S., Hammerla, N.Y., Kainz, B., et al., 2018. Attention u-net: Learning where to look for the pancreas. *arXiv preprint arXiv:1804.03999* .
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al., 2019. Pytorch: An imperative style, high-performance deep learning library. *arXiv preprint arXiv:1912.01703* .
- Pathak, D., Krahenbuhl, P., Donahue, J., Darrell, T., Efros, A.A., 2016. Context encoders: Feature learning by inpainting, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2536–2544.
- Pizzichetta, M.A., Talamini, R., Marghoob, A.A., Soyer, H.P., Argenziano, G., Bono, R., Corradin, M.T., De Giorgi, V., Gonzalez, M.A., Kolm, I., et al., 2013. Negative pigment network: an additional dermoscopic feature for the diagnosis of melanoma. *Journal of the American Academy of Dermatology* 68, 552–559.
- Raghu, M., Zhang, C., Kleinberg, J., Bengio, S., 2019. Transfusion: Understanding transfer learning for medical imaging, in: *Advances in neural information processing systems*, pp. 3347–3357.
- Rajpurkar, P., Irvin, J., Zhu, K., Yang, B., Mehta, H., Duan, T., Ding, D., Bagul, A., Langlotz, C., Shpanskaya, K., et al., 2017. Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning. *arXiv preprint arXiv:1711.05225* .
- Ronneberger, O., Fischer, P., Brox, T., 2015a. U-Net: Convolutional Networks for Biomedical Image Segmentation, in: Navab, N., Hornegger, J., Wells, W.M., Frangi, A.F. (Eds.), *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*. Springer International Publishing, Cham. volume 9351, pp. 234–241. URL: [http://link.springer.com/10.1007/978-3-319-24574-4\\_28](http://link.springer.com/10.1007/978-3-319-24574-4_28), doi:10.1007/978-3-319-24574-4\_28.
- Ronneberger, O., Fischer, P., Brox, T., 2015b. U-net: Convolutional networks for biomedical image segmentation, in: *International Conference on Medical image computing and computer-assisted intervention*, Springer. pp. 234–241.

- Salehi, S.S.M., Erdogmus, D., Gholipour, A., 2017. Tversky loss function for image segmentation using 3d fully convolutional deep networks, in: International workshop on machine learning in medical imaging, Springer. pp. 379–387.
- Schmidhuber, J., 1990. Making the world differentiable: On using self-supervised fully recurrent neural networks for dynamic reinforcement learning and planning in non-stationary environments .
- Society, A.C., 2021. Cancer facts & figures 2021. <https://www.cancer.org/content/dam/cancer-org/research/cancer-facts-and-statistics/annual-cancer-facts-and-figures/2021/cancer-facts-and-figures-2021.pdf>. Accessed: 2021-08-02.
- Solomon, J., 2015. Numerical algorithms: methods for computer vision, machine learning, and graphics. CRC press.
- Szegedy, C., Ioffe, S., Vanhoucke, V., Alemi, A.A., 2017. Inception-v4, inception-resnet and the impact of residual connections on learning, in: Thirty-first AAAI conference on artificial intelligence.
- Tan, C., Sun, F., Kong, T., Zhang, W., Yang, C., Liu, C., 2018. A survey on deep transfer learning, in: International conference on artificial neural networks, Springer. pp. 270–279.
- Tan, M., Le, Q.V., 2019. Efficientnet: Rethinking model scaling for convolutional neural networks. arXiv preprint arXiv:1905.11946 .
- Vondrick, C., Shrivastava, A., Fathi, A., Guadarrama, S., Murphy, K., 2018. Tracking emerges by colorizing videos, in: Proceedings of the European conference on computer vision (ECCV), pp. 391–408.
- Wang, X., Huang, Q., Celikyilmaz, A., Gao, J., Shen, D., Wang, Y.F., Wang, W.Y., Zhang, L., 2019a. Reinforced cross-modal matching and self-supervised imitation learning for vision-language navigation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 6629–6638.
- Wang, X., Jabri, A., Efros, A.A., 2019b. Learning correspondence from the cycle-consistency of time, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 2566–2576.
- Ward, W.H., Farma, J.M., 2017. Cutaneous melanoma: etiology and therapy [internet] .
- Wu, Z., Xiong, Y., Yu, S.X., Lin, D., 2018. Unsupervised feature learning via non-parametric instance discrimination, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 3733–3742.
- Zalaudek, I., Kittler, H., Marghoob, A.A., Balato, A., Blum, A., Dalle, S., Ferrara, G., Fink-Puches, R., Giorgio, C.M., Hofmann-Wellenhof, R., et al., 2008. Time required for a complete skin examination with and without dermoscopy: a prospective, randomized multicenter study. Archives of dermatology 144, 509–513.
- Zhang, R., Isola, P., Efros, A.A., 2016. Colorful image colorization, in: European conference on computer vision, Springer. pp. 649–666.

## Appendix A. Additional Experiment Results

This Appendix provides additional results of our experiments in Section 5.3. Particularly, Figures A.10 and A.11 provide a visual comparison amongst different methods on the ISIC 2018 and ISIC 2017 challenges, respectively. Lastly, Tables A.10, A.11, and A.12, A.13 provide the standard deviation of the Jaccard and Dice on the ISIC 2018 and 2017 challenges respectively.

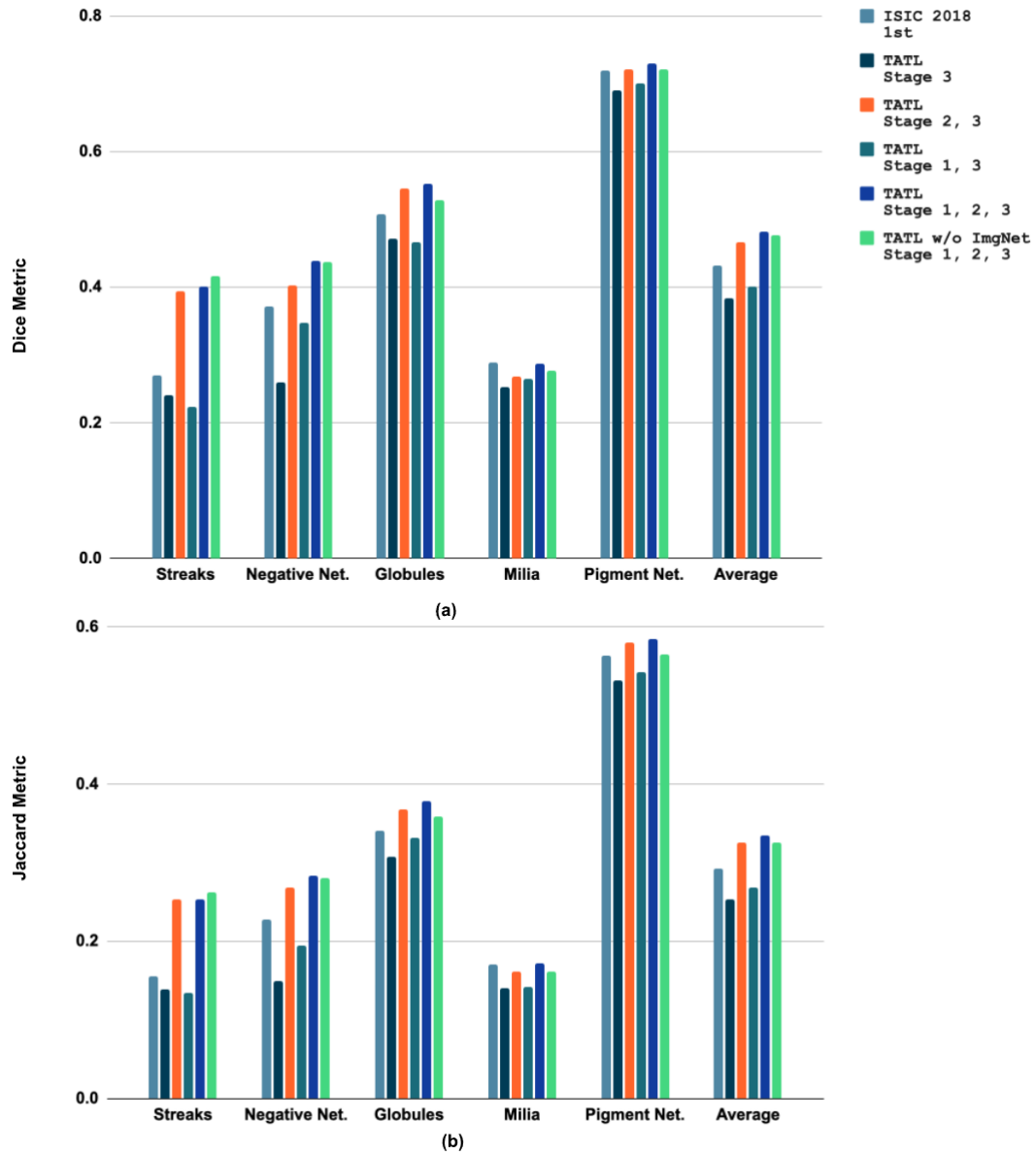


Figure A.10: ISIC-2018 Challenge Performance Visualization with respect to results in Table 3

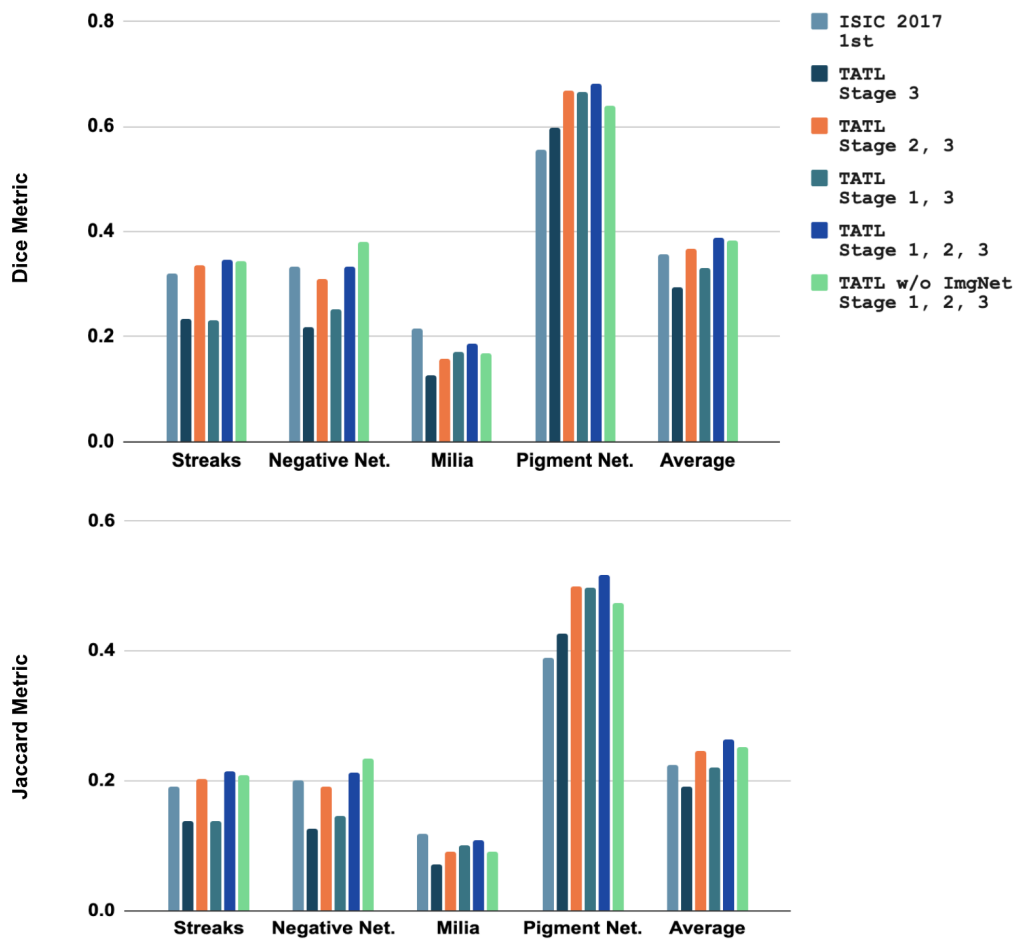


Figure A.11: ISIC-2017 Challenge Performance Visualization with respect to results in Table 4

Table A.10: Standard deviation of the Jaccard metric on the ISIC2018 challenge. The best results are in bold font. Stage 1: segmenting the lesion region, Stage 2: training the Attribute-Agnostic Segmenter, Stage 3: training the Target-Segmenters.

Method	ISIC2018 1st	TATL Stage 3	TATL Stage 2, 3	TATL Stage 1, 3	TATL Stage 1, 2, 3	TATL w/o ImgNet Stage 1, 2, 3
<b>Pigment Net.</b>	0.563	0.532 $\pm$ 0.016	0.580 $\pm$ 0.017	0.542 $\pm$ 0.016	<b>0.584 <math>\pm</math>0.016</b>	0.565 $\pm$ 0.017
<b>Globules</b>	0.341	0.308 $\pm$ 0.023	0.368 $\pm$ 0.025	0.332 $\pm$ 0.024	<b>0.379 <math>\pm</math>0.023</b>	0.359 $\pm$ 0.023
<b>Milia-like cysts</b>	0.171	0.141 $\pm$ 0.016	0.161 $\pm$ 0.015	0.142 $\pm$ 0.014	<b>0.172 <math>\pm</math>0.015</b>	0.161 $\pm$ 0.015
<b>Negative Net.</b>	0.228	0.149 $\pm$ 0.013	0.269 $\pm$ 0.014	0.194 $\pm$ 0.015	<b>0.283 <math>\pm</math>0.013</b>	0.280 $\pm$ 0.014
<b>Streaks</b>	0.156	0.139 $\pm$ 0.044	0.254 $\pm$ 0.042	0.135 $\pm$ 0.043	0.254 $\pm$ 0.045	<b>0.263 <math>\pm</math>0.047</b>
<b>Average</b>	0.292	0.254 $\pm$ 0.022	0.326 $\pm$ 0.023	0.269 $\pm$ 0.022	<b>0.334 <math>\pm</math>0.022</b>	0.326 $\pm$ 0.023

Table A.11: Standard deviation of the Dice metric on the ISIC2018 challenge. The best results are in bold font. Stage 1: segmenting the lesion region, Stage 2: training the Attribute-Agnostic Segmenter, Stage 3: training the Target-Segmenters.

Method	ISIC2018 1st	TATL Stage 3	TATL Stage 2, 3	TATL Stage 1, 3	TATL Stage 1, 2, 3	TATL w/o ImgNet Stage 1, 2, 3
<b>Pigment Net.</b>	0.720	0.691 $\pm$ 0.013	0.721 $\pm$ 0.013	0.701 $\pm$ 0.013	<b>0.730 <math>\pm</math>0.014</b>	0.721 $\pm$ 0.014
<b>Globules</b>	0.508	0.471 $\pm$ 0.024	0.546 $\pm$ 0.025	0.467 $\pm$ 0.027	<b>0.552 <math>\pm</math>0.024</b>	0.528 $\pm$ 0.025
<b>Milia-like cysts</b>	<b>0.289</b>	0.252 $\pm$ 0.024	0.268 $\pm$ 0.025	0.264 $\pm$ 0.025	0.288 $\pm$ 0.025	0.277 $\pm$ 0.027
<b>Negative Net.</b>	0.371	0.260 $\pm$ 0.015	0.403 $\pm$ 0.014	0.348 $\pm$ 0.015	<b>0.438 <math>\pm</math>0.015</b>	0.437 $\pm$ 0.016
<b>Streaks</b>	0.270	0.241 $\pm$ 0.062	0.394 $\pm$ 0.058	0.224 $\pm$ 0.062	0.401 $\pm$ 0.061	<b>0.416 <math>\pm</math>0.064</b>
<b>Average</b>	0.432	0.383 $\pm$ 0.028	0.466 $\pm$ 0.027	0.401 $\pm$ 0.028	<b>0.482 <math>\pm</math>0.028</b>	0.476 $\pm$ 0.029

Table A.12: Standard deviation of the Jaccard metric on the ISIC2017 challenge. The best results are in bold font. Stage 1: segmenting the lesion region, Stage 2: training the Attribute-Agnostic Segmenter, Stage 3: training the Target-Segmenters.

Method	ISIC2017 1st	TATL Stage 3	TATL Stage 2, 3	TATL Stage 1, 3	TATL Stage 1, 2, 3	TATL w/o ImgNet Stage 1, 2, 3
<b>Pigment Net.</b>	0.389	0.426 $\pm$ 0.035	0.499 $\pm$ 0.033	0.497 $\pm$ 0.037	<b>0.516 <math>\pm</math>0.036</b>	0.473 $\pm$ 0.035
<b>Milia-like cysts</b>	<b>0.119</b>	0.072 $\pm$ 0.025	0.091 $\pm$ 0.028	0.101 $\pm$ 0.025	0.108 $\pm$ 0.024	0.092 $\pm$ 0.029
<b>Negative Net.</b>	0.201	0.126 $\pm$ 0.044	0.191 $\pm$ 0.037	0.147 $\pm$ 0.038	0.213 $\pm$ 0.033	<b>0.234 <math>\pm</math>0.035</b>
<b>Streaks</b>	0.192	0.139 $\pm$ 0.033	0.203 $\pm$ 0.045	0.139 $\pm$ 0.037	<b>0.215 <math>\pm</math>0.039</b>	0.209 $\pm$ 0.050
<b>Average</b>	0.225	0.191 $\pm$ 0.034	0.246 $\pm$ 0.036	0.221 $\pm$ 0.034	<b>0.263 <math>\pm</math>0.033</b>	0.252 $\pm$ 0.037

Table A.13: Standard deviation of the Dice metric on the ISIC2017 challenge. The best results are in bold font. Stage 1: segmenting the lesion region, Stage 2: training the Attribute-Agnostic Segmenter, Stage 3: training the Target-Segmenters.

Method	ISIC2017 1st	TATL Stage 3	TATL Stage 2, 3	TATL Stage 1, 3	TATL Stage 1, 2, 3	TATL w/o ImgNet Stage 1, 2, 3
<b>Pigment Net.</b>	0.556	0.597 $\pm$ 0.028	0.667 $\pm$ 0.033	0.665 $\pm$ 0.031	<b>0.681 <math>\pm</math>0.034</b>	0.639 $\pm$ 0.030
<b>Milia-like cysts</b>	<b>0.215</b>	0.127 $\pm$ 0.032	0.157 $\pm$ 0.037	0.172 $\pm$ 0.035	0.188 $\pm$ 0.035	0.168 $\pm$ 0.034
<b>Negative Net.</b>	0.333	0.218 $\pm$ 0.037	0.310 $\pm$ 0.038	0.251 $\pm$ 0.046	0.334 $\pm$ 0.039	<b>0.380 <math>\pm</math>0.036</b>
<b>Streaks</b>	0.321	0.233 $\pm$ 0.047	0.336 $\pm$ 0.062	0.232 $\pm$ 0.049	<b>0.346 <math>\pm</math>0.054</b>	0.345 $\pm$ 0.066
<b>Average</b>	0.356	0.294 $\pm$ 0.036	0.367 $\pm$ 0.043	0.330 $\pm$ 0.040	<b>0.387 <math>\pm</math>0.041</b>	0.383 $\pm$ 0.042