Segmentation Style Discovery: Application to Skin Lesion Images

Kumar Abhishek^{1[0000-0002-7341-9617]}, Jeremy Kawahara^{2[0000-0002-6406-5300]}, and Ghassan Hamarneh^{1[0000-0001-5040-7448]}

¹ School of Computing Science, Simon Fraser University, Canada ² AIP Labs, Budapest, Hungary {kabhishe,hamarneh}@sfu.ca, jeremy@aip.ai

Abstract. Variability in medical image segmentation, arising from annotator preferences, expertise, and their choice of tools, has been well documented. While the majority of multi-annotator segmentation approaches focus on modeling annotator-specific preferences, they require annotator-segmentation correspondence. In this work, we introduce the problem of segmentation style discovery, and propose StyleSeg, a segmentation method that learns plausible, diverse, and semantically consistent segmentation styles from a corpus of image-mask pairs without any knowledge of annotator correspondence. StyleSeg consistently outperforms competing methods on four publicly available skin lesion segmentation (SLS) datasets. We also curate ISIC-MultiAnnot, the largest multi-annotator SLS dataset with annotator correspondence, and our results show a strong alignment, using our newly proposed measure AS², between the predicted styles and annotator preferences. The code and the dataset are available at https://github.com/sfu-mial/StyleSeg.

Keywords: inter-rater variability \cdot image segmentation \cdot dermatology.

1 Introduction

Medical image segmentation is a critical component in medical image analysis pipelines, either as a preprocessing step for subsequent analyses or for treatment planning and image-guided human or robotic intervention. Following the seminal works of Long et al. [15] and Ronneberger et al. [21], there has been tremendous progress in deep learning (DL)-based medical image segmentation [4]. The majority of these works focus on learning to predict a single segmentation for an image. However, variability among experts when segmenting images has been well-documented, and these resulting segmentation masks are the product of latent factors such as ambiguous object boundaries and differences in tools, annotators' skill levels, criteria, and approaches to segmentation, and they capture different annotator segmentation preferences or "styles". Without accommodating these variations, a segmentation model optimized to minimize training error over a variety of human annotations may produce an "average" segmentation. This has motivated research that can be broadly categorized into two classes: methods that model and learn to predict a single "gold standard" segmentation through label aggregation [24,13,18] (SSeg) and methods that predict multiple segmentations to capture the variability of annotations [22,26,12] (MSeg).

MSeg methods rely on modeling annotator-specific preferences, and training them typically requires annotations with annotator-segmentation correspondence. Therefore, given a set of images and a set of annotators, annotatorsegmentation correspondences can be represented as a bipartite graph when every image has been segmented by at least 1 annotator, e.g., LIDC-IDRI [3], or a complete bipartite graph when every image has been segmented by every annotator, e.g., RIGA [2]. However, in the absence of such a correspondence, i.e., a scenario where we have a corpus of images and corresponding annotations without any knowledge of annotator IDs, defining a segmentation style is non-trivial since the latent factors associated with each segmentation are unknown, thus making it challenging to explicitly train a segmentation model to reproduce a particular style. Since we are unable to confirm even the number of unique annotators, we hypothesize that a possible solution for modeling multiannotator segmentations would be discovering unique annotation styles from the dataset alone. Such a discovery-based approach needs to ensure (1) diversity in the discovered styles, (2) segmentation plausibility across all the styles. and (3) semantic consistency of the segmentations across all the images. However, to the best of our knowledge, there is minimal prior work on the discovery and modeling of annotation styles in the absence of annotator correspondence.

We argue that since even experts can (considerably) differ in how they segment, it is only natural that automated models trained thereupon also exhibit this variety. We envision that a segmentation system should produce results that align with the expectations of its (clinical) users, and that these users can vary in their personalization preferences (e.g., a study [10] found that expert dermatologists prefer "tighter" segmentations than their inexperienced counterparts). Moreover, such a system should, with minimal supervision, continue to produce the style that a user expects, thus avoiding constant user involvement with either manual corrections or image-by-image selection of preferred segmentation style.

In this work, we tackle the problem of style discovery and personalization modeling in medical image segmentation without requiring annotator correspondence, and focus our analysis on skin lesion segmentation (SLS). Advancements in DL over the past decade as well as the availability of large publicly available annotated datasets have enabled large strides in SLS [17,23]. Therefore, in this work, we work on style discovery in the context of multiple annotators for SLS, which has not been explored extensively. The majority of previous works focus on SSeg methods: either to select training samples that have high inter-annotator agreement [20] or training ensemble models to handle annotators' variability [18]. More recently, Zepf et al. [25] presented a small-scale (n = 300) analysis of annotation styles in images from the ISIC 2019 dataset based on the granularity of the annotation boundaries. In this work, we make the following contributions: (1) we introduce the problem of segmentation style discovery in the absence of any annotator correspondence and propose a method (StyleSeg) that predicts

multiple plausible, diverse, and semantically consistent segmentation styles, (2) we curate, to the best of our knowledge, the largest multi-annotator SLS dataset (ISIC-MultiAnnot) with annotator-segmentation mapping, and (3) we introduce a new measure (AS^2) for measuring the strength of alignment of the predicted styles with annotator preferences.

$\mathbf{2}$ Method

Let $\mathcal{X} = \{X_i\}_{i=1}^N$ be a set of images and corresponding segmentation masks $\mathcal{Y} = \{\{Y_{ik}\}_{k=1}^{K_i}\}_{i=1}^N$, where $K_i > 0$ denotes the number of different ways X_i was segmented, without any knowledge of annotator correspondence. The goal is to discover unique annotation "styles" in this data $(\mathcal{X}, \mathcal{Y})$ such that, when given an image X_i , we predict $\{Y_{ij}\}_{j=1}^M$: M unique segmentations of X_i , that are diverse, plausible for X_i , and are of semantically consistent styles across all images.

To this end, we propose StyleSeg (Fig. 1 (a)): a segmentation approach that learns to predict M plausible segmentations that capture a variety of styles from a corpus of images and corresponding masks without any annotator correspondence. StyleSeg consists of two deep learning models that are trained together: (i) a segmentation model f_s , parameterized by Θ_s , that predicts M segmentation masks from image X_i , where $M \in \mathbb{N}$ is a user-specified value,

$$\{\hat{Y}_{ij}\}_{j=1}^{M} = f_s(X_i;\Theta_s),$$
 (1)

and (ii) a style classifier model f_c , parameterized by Θ_c , that predicts a vector $p_i \in \mathbb{R}^M$ of M probabilities,

$$p_i = f_c(X_i, Y_{ik}; \Theta_c), \tag{2}$$

where p_{ij} is the probability that (X_i, Y_{ik}) is of style j. Note that knowing X_i is necessary to define the styles, since the observed segmentations are a product of image content and annotation style.

Of the M predicted segmentations from f_s , we first need to identify the style that is the closest to the ground truth Y_{ik} , and then optimize it to make it even closer. Mathematically, we minimize the loss \mathcal{L}_1 ,

$$\mathcal{L}_1 = L_D(Y_{ik}, Y_{im^*}), \tag{3}$$

$$m^* = \arg\max_{i} \operatorname{Dice}(Y_{ik}, \dot{Y}_{ij}), \tag{4}$$

where $L_D = 1 - \text{Dice}$ and Dice denotes the Dice similarity coefficient [7]. We also require the other predicted styles to still be plausible, i.e., similar to ground truth Y_{ik} . However, requiring all styles to be equally plausible compromises the styles? diversity. Therefore, we make the strength of a style's plausibility requirement proportional to its likelihood of being the predicted style, according to the style classifier f_c (Eqn. 2). To this end, we encourage the weighted sum of predicted segmentations \hat{Y}_{ij} to be similar to ground truth Y_{ik} , where the scalar weights are p_{ij} . Mathematically, we minimize the loss \mathcal{L}_2 ,

$$\mathcal{L}_2 = L_D\left(Y_{ik}, \sum_j^M p_{ij}\hat{Y}_{ij}\right).$$
(5)



(a) StyleSeg overview with M = 3: the segmentation model f_s predicts 3 plausible segmentations of different styles, while the style classifier f_c predicts the style that is the most similar to the ground truth.







(c) Distribution of ISIC-MultiAnnot by number of annotators and segmentations.

Fig. 1: (a) An overview of the proposed method StyleSeg. (b) Inter-annotator variability in the training images. (c) An annotator-wise breakdown of the newly curated ISIC-MultiAnnot dataset.

Weighting the M segmentations by p_i ensures that when p_i has a high entropy (e.g., in the initial training epochs), all styles are encouraged to be similar to Y_{ik} , whereas when p_i has a low entropy, only a subset of the M styles are encouraged to be similar to Y_{ik} , thus enabling a coarse-to-fine style refinement.

Additionally, we employ a cross-entropy loss \mathcal{L}_3 to train the style classifier f_c by learning to predict the style that is the most similar to the ground truth,

$$\mathcal{L}_3 = L_{CE}(p_i, m^*). \tag{6}$$

Finally, we optimize the parameters Θ_s of f_s and Θ_c of f_c using

$$\Theta_s^*, \Theta_c^* = \arg\min_{\Theta_s, \Theta_c} \sum_{i}^{N} \mathcal{L}_{\text{total}}, \tag{7}$$

where

$$\mathcal{L}_{\text{total}} = \mathcal{L}_1 + \mathcal{L}_2 + \mathcal{L}_3. \tag{8}$$

Note that we do not include an explicit style distinctiveness constraint since, in the absence of annotator correspondence, the styles are entangled with the segmentations. Nevertheless, these loss terms used together (Eqn. 8) implicitly encourage the styles to be different as the training progresses, as seen in our results.



(a) Sample outputs of StyleSeg for lesions with and without distinct borders.

(b) StyleSeg produces better segmentations than even the test "ground truth".

Fig. 2: Evaluating StyleSeg on ISIC Archive-Test: diverse and plausible segmentations that are semantically consistent across styles.

3 **Results and Discussion**

Datasets: Similar to previous works [20,18], we train StyleSeg on images obtained from the ISIC Archive [1], specifically images with more than one "ground truth" segmentation. We select 2,261 images that meet this criterion (2,122 with two, 100 with three, 35 with four, and 4 with five segmentations), resulting in 4,704 image-mask pairs. Note that these images exhibit a vast range of interannotator agreement, as evidenced qualitatively (sample images with their masks in Fig. 1 (c)) and quantitatively (pairwise Dice coefficients and Fleiss' kappa in Supp. Mat. Fig. SM1 (a)). We choose Fleiss' kappa [9] over Cohen's kappa [6] used by Ribeiro et al. [20] because the former can be used with multi-rater settings while the latter cannot [19]. We reserve 1,525 image-mask pairs from the ISIC Archive for our validation set. See Supp. Mat. for model architectures and training details. We evaluate on four publicly available datasets: ISIC Archive-Test containing 10,000 dermoscopic images with just one segmentation ground truth per image from the ISIC Archive, DermoFit (1,300 clinical images) [5], SCD (206 dermoscopic images) [11], and PH^2 (200 dermoscopic images) [16].

Competing methods: We train StyleSeg with $M = \{2, 3, 4, 6, 8, 10\}$ and compare it to the following SSeg methods: NaiveTraining: a segmentation model without any annotator-specific knowledge; RandAnnotID [18]: 4 segmentation models, one optimized for each annotator randomly assigned to a mask, LessIs-More [20]: a segmentation model trained on a subset of the masks whose average pairwise Cohen's kappa is above 0.5; and D-LEMA [18]: an ensemble of Bayesian segmentation models. We also compare against an MSeg method MHP (multiple hypothesis prediction) [22] also with $M = \{2, 3, 4, 6, 8, 10\}$.

Qualitative results of StyleSeg on ISIC Archive-Test (Fig. 2 (a)) show plausibility (all segmentations cover the lesion with varying degrees of over- or undersegmentation) as well as semantic consistency across segmentations (e.g., when M = 3, yellow always has a tight and jagged boundary while blue always has a loose boundary). We also provide a quantitative assessment in Supp. Mat. Fig. SM1 (f). Also, observe that in lesions with well-defined borders (top two

Table 1: Dice mean_{std.dev.} comparing StyleSeg to SSeg [20,18] (first 4 rows) and MSeg (MHP [22]) methods. For the latter, we report the mean, median, min., max. of Dice between the ground truth and all the predicted segmentation styles. Note how StyleSeg consistently outperforms all competing methods while also producing more plausible segmentations than MHP. \oslash denotes that the result cannot be reported since D-LEMA's [18] code is not available.

Method	ISIC Archive-Test $(n = 10000)$				PH^2 (n = 200)				DermoFit $(n = 1300)$				SCD $(n = 206)$			
	Mean	Median	Minimum	Maximum	Mean	Median	Minimum	Maximum	Mean	Median	Minimum	Maximum	Mean	Median	Minimum	Maximum
NaiveTraining				0.8000.188				0.8800.071				0.8420.199		-		0.7660.108
RandAnnotID [18]		-		0	· ·			0.8970.005	· ·	-		0.8260.004		-		0
LessIsMore [20]		-		0.8150.178	· ·			0.8950.070	· ·	-		0.8540.127		-		0.8040.169
D-LEMA [18]	-	-	-	0	· ·	-	-	0.9200.004	-	-	-	0.8530.003	-	-		0
2-MHP	0.7960.168	0.7270.185	0.7270.195	0.8640.158	0.8500.114	0.7860.162	0.7860.162	0.9140.075	0.7950.149	0.7070.229	0.7070.229	0.8820.089	0.7960.140	0.7130.180	0.7130.180	0.8790.119
2-StyleSeg	$0.814_{o.res}$	0.7600.186	0.7600.186	0.8690.161	0.8780.075	0.8270.102	0.8270.102	0.9290.050	0.8240.128	0.7590.180	0.7590.180	0.8880.000	0.8240.184	0.7540.157	0.7540.157	0.8950.104
3-MHP	0.7720.181	0.7890.184	0.6520.232	0.8760.154	0.7800.185	0.7960.217	0.6250.202	0.9190.075	0.7390.176	0.7670.182	0.5620.289	0.8880.093	0.7150.194	0.7520.217	0.5230.297	0.8690.126
3-StyleSeg	0.8040.169	0.8190.174	0.7130.189	0.8810.154	0.8850.082	0.9000.080	0.8110.137	0.9430.045	0.8170.194	0.8350.136	$0.720_{o.202}$	$0.897 a_{aaa}$	0.8180.151	0.8370.149	0.7160.213	0.9010.120
4-MHP	0.7730.170	0.7520.182	0.6230.225	0.8860.142	0.8300.121	0.8170.164	0.6740.264	0.9330.049	0.7960.194	0.7830.157	0.6360.243	0.9040.072	0.7510.153	0.7260.189	0.5470.251	0.8960.104
4-StyleSeg	0.8040.163	0.7860.177	0.693a.187	0.8890.147	0.8750.084	0.8630.102	0.7760.142	0.9450.042	0.8120.195	0.7940.165	0.6810.221	0.9070.075	0.7860.152	0.7660.179	$0.632_{o.228}$	0.8960.111
6-MHP	0.6470.152	0.7030.202	0.1210.175	0.8860.131	0.7900.072	0.8400.089	0.4000.186	0.9390.028	0.7490.116	0.7770.151	0.4280.169	0.9000.083	0.6490.192	0.7030.181	0.1560.126	0.8810.103
6-StyleSeg	0.7950.178	0.7990.180	0.6480.234	0.8890.154	0.8690.090	0.873a.oss	0.7450.167	0.9480.000	0.8140.136	0.8180.149	$0.651_{o.ess}$	$0.911_{o.oro}$	0.7980.143	0.8060.149	0.6080.244	0.9060.100
8-MHP	0.6250.152	0.6580.215	0.0990.140	0.8960.121	0.7520.072	0.801a.ose	0.2600.148	0.9440.025	0.6980.117	0.7080.158	0.3090.181	0.9080.000	0.6160.195	0.6270.200	0.1340.111	0.8970.092
8-StyleSeg	0.7900.170	0.7980.185	0.5950.232	0.8990.198	0.8750.096	0.8780.113	0.7450.158	0.9500.027	0.8100.141	0.8150.162	$0.632_{o.exp}$	0.9100.075	0.7980.161	0.8120.175	0.5860.241	0.9010.114
10-MHP	0.7060.174	0.7450.204	0.2810.231	0.8940.126	0.7240.183	0.7610.222	0.3390.285	0.9380.029	0.6290.177	0.6670.219	0.1810.210	0.9060.008	0.6900.168	0.7330.218	0.2230.196	0.8980.094
10-StyleSeg	0.7930.173	0.8050.185	0.6030.814	0.8990.144	0.8660.101	0.880	0.6920.179	0.9510.025	0.8010.147	0.8130.166	0.5790.255	0.9180.005	0.7680.181	0.7910.186	$0.513_{o.246}$	0.8850.140

rows), the predicted styles are similar, whereas in lesions with ambiguous borders (bottom two rows), the predictions exhibit considerable diversity. It is also worth noting that several images in ISIC Archive-Test have either incorrect or imprecise "ground truth" masks (Fig. 2 (b)), which leads to incorrect penalization of StyleSeg's accurate predictions during evaluation.

Quantitative results: For SSeg methods, we report the Dice coefficient. For MSeg methods, we report $\max_j(d)$, where $d = \text{Dice}(Y_{ik}, \hat{Y}_{ij})$, to assess the highest agreement, and $\{\max_j(d), \max_j(d), \min_j(d)\}$ to assess the plausibility of all M segmentations. For example, an MSeg model that produces even one poor segmentation will have low scores for $\min_j(d)$, indicating low plausibility.

Table 1 shows that predicting more than one style (StyleSeg, MHP) improves performance $(\max_j(d))$ compared to SSeg methods, and even MSeg methods that predict just two styles (2-StyleSeg, 2-MHP) consistently outperform SSeg methods. Moreover, as M increases, a larger number of diverse segmentations are produced, and the $\max_j(d)$ keeps improving. However, we observe that for three out of the four datasets, the $\max_j(d)$ performance either plateaus or starts to decline as M increases. We posit that after an optimal number of styles, generating more segmentations leads to diversity at the cost of performance, and leave this investigation for future work. Interestingly, datasets that do not have a documented presence of inter-segmentation variability (DermoFit, PH², SCD) also benefit from learning to predict multiple segmentations, indicating style variability in the ground truth masks. A post hoc investigation of DermoFit, for example, confirms the presence of different annotation styles (difference in boundary granularity; Supp. Mat. Fig. SM1 (b)).

Finally, StyleSeg consistently outperforms MHP for all datasets and M, except M = 10 with SCD, and as M increases, the plausibility of MHP models across all the predictions decreases, as evident through the declining $\{\text{mean}_j(d), \text{median}_j(d), \min_j(d)\}$ scores. For example, when modeling 10 styles, the $[\min_j(d), min_j(d)]$ scores.

Table 2: StyleSeg's segmentation agreement (mean_{std.dev}. of Dice_{IASS} and Dice_{ASSS}) and style alignment (AS²) on the 27 annotator preferences in ISIC-MultiAnnot. \mathcal{J} denotes the single style that, for each row, maximizes agreement with the ground truth. As more styles are modeled, Dice_{IASS}, Dice_{ASSS}, and AS² all improve, and all annotator preferences consistently align with a discovered style.

Annotator + Tool	Seg.	1-StyleSeg 2-StyleSeg				3	-StyleSeg	4-StyleSeg			
+ Experience	Count	Dice _{ISSS}	Dice _{ISSS}	Dice _{ASSS}	\mathcal{J}	Dice _{ISSS}	Dice _{ASSS}	I	Dice _{ISSS}	Dice _{ASSS}	J
A00+T2+E	1573	0.8920.089	0.9230.061	0.9130.087	2	0.9440.049	0.9130.106	3	0.9440.044	0.9140.111	1
A00+T2+N	1305	0.716 <i>o.soz</i>	0.7610.295	0.728 <i>o.sos</i>	2	0.7930.287	0.7270.313	3	0.7900.290	0.7260.304	3
A01+T1+N	6	0.5590.362	0.7660.152	0.7660.152	1	0.7540.132	0.7410.125	2	0.8190.106	0.7670.118	2
A01+T3+E	297	0.9000.104	0.9150.095	0.8970.107	2	0.9270.075	0.9000.097	1	0.9310.067	0.9040.090	3
A01+T3+N	2148	0.8290.185	0.8570.167	0.8170.170	1	0.8690.159	0.8360.178	1	0.8760.148	0.8360.175	3
A02+T1+E	1742	0.8440.177	0.8800.140	0.8560.159	1	0.8860.132	0.8540.159	1	0.8950.112	0.8590.148	4
A02+T3+E	468	0.8560.172	0.8890.167	0.8830.175	2	0.8990.161	0.8740.188	3	0.9030.146	0.8900.160	1
A03+T1+E	1622	0.7780.168	0.8450.117	0.8270.137	1	0.8540.111	0.8240.145	2	0.8810.095	0.8230.132	4
A03+T3+E	260	0.8910.116	0.9120.086	0.8760.173	2	0.9230.089	0.8680.150	1	0.9320.074	0.8740.163	3
A04+T1+E	992	0.8500.158	0.8800.131	0.8600.149	1	0.8880.132	0.8660.153	2	0.9060.108	0.8560.157	4
A04+T1+N	61	0.7600.242	0.8400.152	0.8230.164	1	0.8370.162	0.7860.201	1	0.8270.206	0.7890.226	4
A04+T3+E	913	0.9120.088	0.9390.054	0.9340.065	2	0.9480.047	0.9260.069	1	0.9510.045	0.9320.065	3
A04+T3+N	90	0.8770.096	0.9100.068	0.9050.070	2	0.9280.031	0.9080.044	3	0.9260.052	0.9130.055	1
A05+T1+E	752	0.8150.205	0.8620.165	0.8370.179	1	0.8730.162	0.8270.184	1	0.8820.147	0.8410.177	4
A05+T3+E	742	0.8750.129	0.9030.109	0.8910.118	2	0.9160.098	0.8780.120	1	0.9190.091	0.8910.108	1
A06+T1+E	10	0.8240.187	0.902 <i>0.03</i> 7	0.8850.070	1	0.9090.034	0.8890.049	2	0.9090.039	0.8800.065	4
A06+T3+E	24	0.8620.079	0.9160.055	0.9160.055	2	0.9340.031	0.9230.031	3	0.9330.041	0.9290.040	1
A07+T1+E	67	0.8200.157	0.8770.124	0.8670.150	1	0.8900.108	0.8620.157	2	0.8970.104	0.8620.149	4
A07+T1+N	251	0.8370.141	0.8920.085	0.8790.104	1	0.9030.067	0.8750.114	2	0.9050.070	0.8730.101	4
A07+T3+E	12	0.9250.055	0.9380.019	0.9370.019	2	0.9390.020	0.9160.055	1	0.9470.016	0.9320.017	1
A07+T3+N	39	0.8630.177	0.9180.061	0.9130.071	2	0.9330.037	0.8990.148	3	0.934 <i>0.039</i>	0.9140.079	1
A08+T1+E	26	0.6660.225	0.7500.161	0.6800.242	2	0.7470.197	0.6530.260	1	0.7930.134	0.6660.261	1
A08+T3+E	111	0.6050.230	0.6680.197	0.6260.210	1	0.6770.206	0.6280.218	2	0.7350.166	0.6690.203	2
A09+T1+E	30	0.8150.121	0.8410.098	0.7840.156	1	0.8730.089	0.8330.118	2	0.8840.076	0.8120.119	4
A09+T1+N	1	0.9530.000	0.9270.000	0.9270.000	2	0.9550.000	0.9550.000	1	0.9470.000	0.9470.000	3
A09+T3+E	10	0.9000.074	0.9180.054	0.9180.054	2	0.9330.038	0.9090.044	1	0.9370.045	0.9190.040	3
A09+T3+N	3	0.8940.070	0.9110.058	0.9110.058	2	0.9570.015	0.9570.015	3	0.9440.050	0.9440.050	1
AS ² (Eqn.	9)	-		0.2990.208		(0.3470.237	0.4660.296			

 $\max_j(d)$] range across 10 segmentations for 10,000 test images in ISIC Archive-Test is [0.281, 0.894] for 10-MHP and [0.603, 0.899] for 10-StyleSeg, meaning all the predicted segmentations are more plausible for the latter. We attribute this improvement to the plausibility constraint (\mathcal{L}_2 in Eqn. 5), which penalizes predicted segmentations that considerably deviate from the ground truth.

A new multi-annotator dataset: Next, we propose ISIC-MultiAnnot, a new multi-annotator SLS dataset curated from the ISIC Archive that, to the best of our knowledge, is the largest such dataset to contain annotator correspondence. The annotator-segmentation mapping in ISIC-MultiAnnot forms an incomplete bipartite graph, i.e., not every image has been segmented by every annotator. ISIC-MultiAnnot contains 12,951 images segmented by 10 annotators, resulting in 13,555 image-mask pairs (breakdown in Fig. 1 (c)). Unlike other multi-annotator datasets, the variability in ISIC-MultiAnnot's segmentations stems from three annotation pipeline factors: the annotator (10 annotator IDs: "A00"–"A09"), the tool used ("T1", "T2", "T3"), and the expertise of the manual reviewer ("expert" or "novice") [18], resulting in 27 unique annotator preferences, which we use for our evaluation. We measure StyleSeg performance in two settings: (i) **image-adaptive style selection (IASS)**: for every image,

we find the style that maximizes the agreement with ground truth, measured as $\text{Dice}_{\text{IASS}} = \max_j(\text{Dice}(Y_{ik}, \hat{Y}_{ij}))$, and (ii) a more challenging **annotator-specific** style selection (ASSS): we find a single style, fixed across all images, that maximizes agreement with ground truth, measured as $\text{Dice}_{\text{ASSS}} = \text{Dice}(Y_{ik}, \hat{Y}_{i\mathcal{J}})$, where $\mathcal{J} = \arg\max_j(\sum_i \text{Dice}(Y_{ik}, \hat{Y}_{ij}))$. Note that $\text{Dice}_{\text{ASSS}} \leq \text{Dice}_{\text{IASS}}$.

Quantitative results of StyleSeg on ISIC-MultiAnnot (Table 2; additional results in Supp. Mat. Fig. SM1 (c)) show that each of the discovered styles presents a high agreement with almost all the annotator preferences. A notable outlier is "A08", and upon manual inspection, we found a large number of ground truth segmentations to be incorrect to varying degrees (Supp. Mat. Fig. SM1 (e)), which explains the lower evaluation performance. Similar to Table 1, modeling even two styles yields better performance than one style. Moreover, as M increases, the newly discovered styles continue to show increasing usefulness, since all of them consistently align with one or more annotator preferences, meaning that they are able to capture the diversity in segmentations with an increasing level of granularity.

As an additional experiment to assess whether the learned styles are able to model tool-specific ("T1", "T2", "T3") latent factors, we separate the segmentations into three groups based on the tool, pass each corresponding image through a trained 3-StyleSeg model, and determine the predicted segmentation style with the highest overlap. We observe that the most commonly chosen style within a group is unique for each of the three tool groups, suggesting that differences among the three tools are learned within the three styles.

A new style alignment measure: When choosing only one style to evaluate each annotator preference, it is important to note that a particular style could be assigned as the chosen style for a certain annotator even if it best fits either 100% (perfect alignment) of the images or just slightly above random chance (weak alignment). We propose to measure this Annotator-Style Alignment Strength (AS^2) as 1 - 'normalized Shannon entropy of annotator-style assignment', i.e.,

$$AS^{2} = 1 - \frac{-\sum_{i=1}^{M} q_{i} \log_{2} q_{i}}{-\sum_{i=1}^{M} \frac{1}{M} \log_{2} \frac{1}{M}},$$
(9)

where q_i is the vector of fractions of segmentations assigned to each style (e.g., q = [0.70, 0.15, 0.15] for a 3-StyleSeg model that assigns 100 images from a certain annotator preference as 70:15:15 images for styles 1, 2, and 3, yielding $AS^2 = 0.255$). Note that AS^2 is 0 for uniform assignments, and increases logarithmically approaching 1 as assignments become more consistent. Our results in Table 2 show that AS^2 values do not decrease as M increases, meaning that learning to model more styles is not detrimental to segmentation quality and indeed captures more diversity. Additional results are presented in Supp. Mat. Fig. SM1 (d, g).

4 Conclusion

We formulated the problem of segmentation style discovery in the absence of annotator correspondence. We showed how our proposed method, StyleSeg, discovers segmentation styles that are diverse, semantically consistent, and more plausible than those generated by competing methods, as evaluated on four public skin lesion segmentation (SLS) datasets. We also curated ISIC-MultiAnnot, the largest multi-annotator SLS dataset with 13,555 image-mask pairs from 10 annotators from ISIC Archive, and showed how StyleSeg consistently achieves high agreement with the annotator preferences, as measured through the Dice coefficient as well as a newly proposed measure, Annotator-Style Alignment Strength, for measuring annotator-style alignment. Future work would include an explicit "disentanglement" of annotation styles from image content and approaches to find the optimal number of styles in a segmentation dataset.

Acknowledgments. The authors are grateful for the computational resources provided by NVIDIA Corporation and Digital Research Alliance of Canada (formerly Compute Canada). Partial funding for this project was provided by the Natural Sciences and Engineering Research Council of Canada (NSERC RGPIN/06752-2020).

Disclosure of Interests. The authors have no competing interests to declare.

References

- 1. International Skin Imaging Collaboration: Melanoma Project. https://www. isic-archive.com/, [Online. Accessed February 01, 2024]
- 2. Almazroa, A., Alodhayb, S., Osman, E., Ramadan, E., Hummadi, M., Dlaim, M., Alkatee, M., Raahemifar, K., Lakshminarayanan, V.: Agreement among ophthalmologists in marking the optic disc and optic cup in fundus images. International Ophthalmology 37, 701-717 (2017)
- 3. Armato III, S.G., McLennan, G., Bidaut, L., McNitt-Gray, M.F., Meyer, C.R., Reeves, A.P., Zhao, B., Aberle, D.R., Henschke, C.I., Hoffman, E.A., et al.: The lung image database consortium (LIDC) and image database resource initiative (IDRI): a completed reference database of lung nodules on CT scans. Medical Physics **38**(2), 915–931 (2011)
- 4. Asgari Taghanaki, S., Abhishek, K., Cohen, J.P., Cohen-Adad, J., Hamarneh, G.: Deep semantic segmentation of natural and medical images: A review. Artificial Intelligence Review 54, 137–178 (2021)
- 5. Ballerini, L., Fisher, R.B., Aldridge, B., Rees, J.: A color and texture based hierarchical K-NN approach to the classification of non-melanoma skin lesions. In: Celebi, M.E., Schaefer, G. (eds.) Color Medical Image Analysis. vol. 6, pp. 63–86. Springer Netherlands (2013)
- 6. Cohen, J.: A coefficient of agreement for nominal scales. Educational and Psychological Measurement **20**(1), 37–46 (1960)
- 7. Dice, L.R.: Measures of the amount of ecologic association between species. Ecology 26(3), 297-302 (1945)
- 8. Ercal, F., Chawla, A., Stoecker, W.V., Lee, H.C., Moss, R.H.: Neural network diagnosis of malignant melanoma from color images. IEEE Transactions on Biomedical Engineering **41**(9), 837–845 (1994)
- 9. Fleiss, J.L.: Measuring nominal scale agreement among many raters. Psychological Bulletin **76**(5), 378 (1971)

- 10 Abhishek et al.
- Fortina, A.B., Peserico, E., Silletti, A., Zattra, E.: Where's the naevus? interoperator variability in the localization of melanocytic lesion border. Skin Research and Technology 18(3), 311–315 (2012)
- Glaister, J., Amelard, R., Wong, A., Clausi, D.A.: MSIM: Multistage illumination modeling of dermatological photographs for illumination-corrected skin lesion analysis. IEEE Transactions on Biomedical Engineering 60(7), 1873–1883 (2013)
- Ji, W., Yu, S., Wu, J., Ma, K., Bian, C., Bi, Q., Li, J., Liu, H., Cheng, L., Zheng, Y.: Learning calibrated medical image segmentation via multi-rater agreement modeling. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12341–12351 (2021)
- Kats, E., Goldberger, J., Greenspan, H.: A soft STAPLE algorithm combined with anatomical knowledge. In: Medical Image Computing and Computer Assisted Intervention-MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part III 22. pp. 510–517. Springer (2019)
- Kawahara, J., Hamarneh, G.: Fully convolutional neural networks to detect clinical dermoscopic features. IEEE Journal of Biomedical and Health Informatics 23(2), 578–585 (2018)
- Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3431–3440 (2015)
- Mendonça, T., Ferreira, P.M., Marques, J.S., Marcal, A.R.S., Rozeira, J.: PH² a dermoscopic image database for research and benchmarking. In: IEEE Engineering in Medicine and Biology Society. pp. 5437–5440 (2013)
- Mirikharaji, Z., Abhishek, K., Bissoto, A., Barata, C., Avila, S., Valle, E., Celebi, M.E., Hamarneh, G.: A survey on deep learning for skin lesion segmentation. Medical Image Analysis p. 102863 (2023)
- Mirikharaji, Z., Abhishek, K., Izadi, S., Hamarneh, G.: D-LEMA: Deep learning ensembles from multiple annotations-application to skin lesion segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1837–1846 (2021)
- Powers, D.M.W.: The problem with kappa. In: Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics. pp. 345–355 (2012)
- Ribeiro, V., Avila, S., Valle, E.: Less is more: Sample selection and label conditioning improve skin lesion segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops. pp. 738–739 (2020)
- Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: Medical Image Computing and Computer-Assisted Intervention-MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, part III 18. pp. 234–241. Springer (2015)
- Rupprecht, C., Laina, I., DiPietro, R., Baust, M., Tombari, F., Navab, N., Hager, G.D.: Learning in an uncertain world: Representing ambiguity through multiple hypotheses. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 3591–3600 (2017)
- Tschandl, P., Sinz, C., Kittler, H.: Domain-specific classification-pretrained fully convolutional network encoders for skin lesion segmentation. Computers in Biology and Medicine 104, 111–116 (2019)
- Warfield, S.K., Zou, K.H., Wells, W.M.: Simultaneous truth and performance level estimation (STAPLE): An algorithm for the validation of image segmentation. IEEE Transactions on Medical Imaging 23(7), 903–921 (2004)

- 25. Zepf, K., Petersen, E., Frellsen, J., Feragen, A.: That label's got style: Handling label style bias for uncertain image segmentation. In: The Eleventh International Conference on Learning Representations (2023), https://openreview.net/forum? id=wZ2SVhOTzBX
- 26. Zhang, L., Tanno, R., Bronik, K., Jin, C., Nachev, P., Barkhof, F., Ciccarelli, O., Alexander, D.C.: Learning to segment when experts disagree. In: Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part I 23. pp. 179–190. Springer (2020)

1

Supplementary Material

Implementation Details

The style classification model $f_c : (X_i, Y_{ik}) \in \mathbb{R}^{224 \times 224 \times 4} \rightarrow p_i \in \mathbb{R}^M$ is an ImageNet-pretrained ResNet-50 model with two modifications: it takes 4channels (i.e., a concatenation of X_i and Y_{ik}) as input and produces a *M*-class prediction. The segmentation model $f_s : X_i \in \mathbb{R}^{224 \times 224 \times 3} \rightarrow \hat{Y}_i \in \mathbb{R}^{224 \times 224 \times M}$ is an ImageNet-pretrained VGG-16 model with the fully connected layers removed, and multi-scale features resized, concatenated, and passed through a Conv2D layer with a sigmoid activation for binary mask prediction [14]. The images and masks were resized to 224×224 spatial resolution using nearestneighbor interpolation. All models were trained for 10 epochs with a batch size of 4 using the Adam optimizer and a learning rate of 5e-5, and the epoch with the lowest loss $\mathcal{L}_{\text{total}}$ (Eqn. 8) on the validation set was used for evaluation. All models were trained on an Ubuntu 20.04 workstation with AMD Ryzen 9 5950X, 32 GB of RAM, and NVIDIA RTX 3090 GPU, running Python 3.10.13 and PyTorch 2.1.2. The PyTorch implementation of StyleSeg and more details about ISIC-MultiAnnot are available at https://github.com/sfu-mial/StyleSeg.



(a) Distribution of image-level pairwise Dice coefficient and Fleiss' kappa values in the 4,704 training image-mask pairs from ISIC Archive. Note how a considerable number of images have poor inter-annotator agreement (Dice <0.2 and Fleiss' kappa < 0).



(c) StyleSeg segmentation performance (Dice) on ISIC-MultiAnnot, reported at annotator-level ("A00"–"A09"). Note how the Dice values are strongly correlated (Pearson correlation coefficient ρ) with the number of styles (M).

Fig. SM1: Supplementary Figures



(b) Undocumented annotation style variability (varying boundary granularity) in DermoFit's [5] ground truth masks.



(d) AS^2 values remain high as M increases, meaning modeling more styles captures larger diversity.

1.8. Surre Subbromonoury 1.8a



(e) Images and "ground truth" segmentation masks from annotator "A08" in ISIC-MultiAnnot. Note the poor quality of segmentation, which in turn, affects evaluation.



(f) Evaluating semantic consistency of segmentation styles on ISIC Archive-Test. For all the 10,000 images, the shape features (area and perimeter of segmentation contours) of StyleSeg outputs are calculated, normalized per lesion w.r.t. the first style (red dot at (1.0, 1.0)), and the points and their kernel density estimates are colored by their style. Even for large values of M, the styles remain distinct. Note that the styles vary in their contour area indicating under- and over-segmentation of skin lesions. Also, for a certain contour area (A), the contour perimeter (P) varies, implying the styles differ in higher-order features such as border irregularity index and contour compactness, which is expected since both border irregularity = $\frac{P^2}{4\pi A}$ [8] and compactness = $\frac{\text{convex hull area}}{A}$ are also functions of perimeter and area.



(g) Annotator-Style Alignment Strength (AS²) values for a variety of fractions of segmentations q (Eqn. 9) with M = 10. AS² is 1 in case of a perfect assignment and 0 in case of a uniform assignment. AS² increases logarithmically as q becomes more concentrated: for example, note that even for a low-entropy assignment such as $[0.9, 0.1, 0.0, \dots 0.0]$, the Annotator-Style Alignment Strength drops to 0.86.

Fig. SM1: (continued) Supplementary Figures.