

Semantic Segmentation and Synthesis for Paintings

Zheyu Liu, Zheying Lu, Qinghui (Luna) Xia

Abstract

What were the views that artists saw 500 years ago? Curious about the sceneries depicted in artworks, we propose a deep learning approach to reconstruct photo-realistic views from paintings. Models that we used in our project were DeepLab (Chen et al., 2016) for semantic segmentation on photos and SPADE (Park et al., 2019) for photo-realistic image synthesis. The novelty of our project focused on improving segmentation performance on landscape paintings. Given that most of the current segmentation techniques and datasets were based on real photos, we conducted data pre-processing by applying style transfers to landscape photos with the AdaIN model (Huang and Belongie, 2017). Our approach achieved reasonably accurate reconstructions for some landscape paintings, as shown in both quantitative and qualitative evaluations.

1 Introduction

With developments in Artificial Intelligence and the rise of interdisciplinary practices in the art industry, it was without doubt that AI-infused artworks had come into view. Machines gained the ability to create music, paintings, photos, and even a novel that won a literature prize (Smithsonian Magazine, 2016). As a new addition to traditional art-making tools, AI-Art also spawned controversy. Danny (Smithsonian Magazine, 2016) stated that such technologies could bring inspirations to authors. But in the long term, relying on machines for inspirations might negatively impact human creativity. Another aspect of the conversation raised doubts over the authenticity of machine-generated articles. How to assess similarities between the generated results and the training sets was a challenge with great moral values.

Sharing strong interests in machine-generated artworks, we researched the applications of Generative Adversarial Networks (GAN) (Creswell et al.,

2018) in such fields, and attempted to find a non-controversial research direction.

There were a lot of existing researches on transforming photos using different styles. Inspired by A Neural Algorithm of Artistic Style (Gatys et al., 2015), which performed style-transfer on a realistic photo to get an artistic photo, we obtained the idea of converting paintings to realistic photos to get a sense of what the painted views would look like in real life. Besides the generation part, evaluation was also an essential task for us. CAN (Elgammal et al., 2017), a model that combined different artistic styles in generating art pieces, proposed quantitative and qualitative metrics to determine if a piece of artwork met certain criteria. Part of their metrics were realized by conducting surveys with groups of people, which was a method we adopted in qualitative assessments.

2 Related Works

CNN Convolutional Neural Networks (Lecun et al., 1998) had been widely applied to computer vision tasks like image classification and image segmentation. Albawi (Albawi et al., 2017) stated that peoples' interests in deeper hidden layers had recently begun to make a great impact on the performance of classical methods in different fields, especially in pattern recognition. With an increasing amount of researches conducted on deep models, GANs (Creswell et al., 2018) gained attention with their excellent performance in computer vision tasks.

Generative Adversarial Models Generative adversarial networks (GANs) (Creswell et al., 2018) are capable of learning deep representations without extensively annotated training data. GAN applications include image synthesis, semantic image editing, style transfer, and image classification. A lot of image generation related topics have conducted researches with GANs, such as the

famous fake face generator (Karras et al., 2018). We also used a GAN for data pre-processing. GAN-related works inspired us to look deeper into image segmentation and image synthesis tasks.

Image Segmentation Image segmentation is a key task in computer vision (Minaee et al., 2021). The broad success of Deep Learning (DL) has prompted the development of new image segmentation approaches leveraging deep learning models. A lot more different architectures arose in this field of study. With the existing architectures, we were able to conduct reasonably accurate image segmentation for real photos. However, we did not find any existing work on segmentation for paintings. In the following sections, we will illustrate more on our approach to solve this problem.

Image Synthesis Developments in Image Synthesis tasks have brought research understandings of image contexts and styles to another level. SPADE (Park et al., 2019) proposed an approach for generating realistic and high-quality photos given simple hand-drawn lines as inputs. We included the model in our project to explore its broader potential in applications.

3 Datasets

For the training dataset, we chose COCO-Stuff (Caesar et al., 2016) as our base dataset and generated Landscape-Dataset and AdaIN-Dataset from COCO-Stuff.

3.1 Training dataset

Our primary goal of this project was to generate synthesized photos from paintings. We planned to train a segmentation model for paintings, and feed its results to SPADE (to be discussed in the following section). Given that SPADE generated synthesis outputs based on pixel-wise semantic segmentation masks from the segmentation model, the segmentation labels (for example, 1 for sky, 2 for trees, etc.) between the two models should match. In the following sections, we will discuss the details of our data-prepossessing part.

3.1.1 COCO-Stuff

In this project, we selected COCO dataset (Lin et al., 2014) as the dataset for both models to share consistent labels. After observing inference results from DeepLab using different types of paint-



Figure 1: Landscape-Dataset

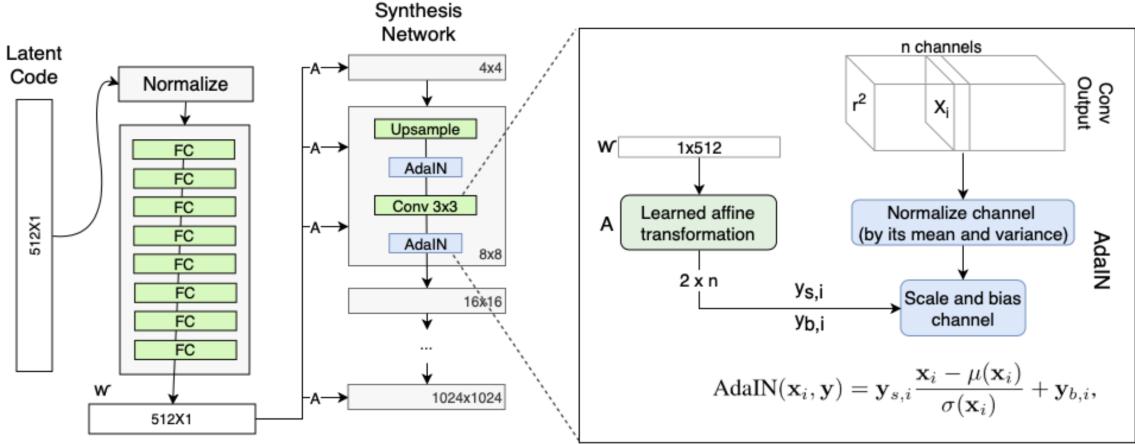
ings, we found that this segmentation architecture performed better on landscape objects. Also, a lot of paintings were about views or landscapes. The COCO-Stuff dataset (Caesar et al., 2016) augmented all 164K images from the popular COCO dataset with pixel-level stuff annotations. These annotations can be used for scene understanding tasks like semantic segmentation, object detection and image captioning.

3.1.2 Landscape-Dataset

Considering of the observations mentioned before, we decided to select landscape photos from the original dataset to construct our Landscape-Dataset. Across all 171 categories in the COCO-Stuff dataset, we manually selected 33 categories as common features of landscape photos. If images had more than five landscape features, we included them as landscape images. By doing this step, We were hoping to better familiarize or bias the segmentation model with landscapes. This dataset consisted of 8191 training samples and 900 testing samples.

3.1.3 AdaIN-Dataset

After compiling Landscape-Dataset, which was a subset extracted from COCO-Stuff dataset with landscape photos, we focused more on the differences between paintings and photos. The existing architectures for segmentation tasks all experimented on real photos. However, we wanted to feed paintings into the model as our inputs. Since there was no existing dataset with pairings of paintings and annotation maps, we believed that conducting style-transfer on the COCO-Stuff dataset with its original annotations could be a solution to this problem. AdaIN (Huang and Belongie, 2017) is a tool for style transfer. It is similar to



The generator's Adaptive Instance Normalization (AdaIN)

Figure 2: AdaIN architecture

the architecture mentioned above - A Neural Algorithm of Artistic Style (Gatys et al., 2015). The AdaIN (Adaptive Instance Normalization) module transfers the encoded information, created by the Mapping Network, into the generated image. The module is added to each resolution level of the Synthesis Network and defines the visual expression of the features in that level:

1. Each channel of the convolution layer output is first normalized to make sure the scaling and shifting of step 3 have the expected effect.
2. The intermediate vector w is transformed using another fully-connected layer (marked as A) into a scale and bias for each channel.
3. The scale and bias vectors shift each channel of the convolution output, thereby defining the importance of each filter in the convolution. This tuning translates the information from w to a visual representation.

Here we chose AdaIN over other approaches because it is faster and has more flexibility for users. The key to the efficiency in such feed-forward method is that AdaIN applies instance Normalization in each layer which normalizes the mean and standard deviation which are computed across spatial dimensions independently for each channel and each sample.

We applied AdaIN style-transfer on roughly half of our Landscape-Dataset to generate AdaIN-Dataset, which consisted of painting-style images. The resulting dataset contained 3994 and 899 images in training and testing sets, respectively. By



Figure 3: AdaIN-dataset

198
199
200
201
202
203
204
205
206
207
208
209
210
211
212

performing this step of data pre-processing, we were aiming to improve segmentation performance on paintings.

3.2 Test dataset for the pipeline

To better test our model, we divided the dataset into simple and complex categories depending on different painting styles. For the simple category, we included some of the realistic painting styles. And for the complex models, we included more abstract painting styles such as Impressionism. By dividing our test set into simple and complex categories, we were hoping to improve our model performance on a more specific branch, focusing on the simple baseline while reaching for achievements in the complex cases.

213 4 Methods

214 4.1 Overview

215 The base dataset we used was the COCO-Stuff
216 dataset (Caesar et al., 2016). As mentioned above
217 in the Datasets section, we extracted landscape pho-
218 tos from COCO-Stuff to form Landscape-Dataset,
219 and later generated AdaIN-Dataset by applying
220 style transfers to the landscape images.

221 The training process was conducted in 4 steps.
222 First, we loaded a DeepLab model with pre-trained
223 weights, and trained it on Landscape-Dataset. This
224 step aimed to improve DeepLab segmentation per-
225 formance on landscape photos specifically. We now
226 had our DeepLab-Landscape model. Another vari-
227 ant of DeepLab, named DeepLab-AdaIN was de-
228 veloped similarly by training DeepLab on AdaIN-
229 Dataset. This variant could help us assess whether
230 a modified DeepLab could generate better segmen-
231 tation results on paintings.

232 After training, we ran inferences on the DeepLab
233 variants using the test set from AdaIN-Dataset as in-
234 puts. This would generate segmentation maps from
235 synthesized paintings. We then fed the outputs into
236 SPADE to get synthesized photos.

237 In the following sections, we are going to further
238 discuss the models we used in this project.

239 4.2 Architecture

240 4.2.1 DeepLab

241 DeepLab (Chen et al., 2016) re-purposes a DCNN
242 trained for image classification tasks to semantic
243 segmentation tasks. The model overcomes three
244 challenges.

245 First, DCNNs designed for image classification
246 perform repeated pooling and downsampling at
247 consecutive layers. This causes significantly re-
248duced spatial resolution in feature maps. DeepLab
249 computes the feature maps in a higher sampling
250 rate by upsampling filters in convolutional layers.

251 The second challenge is to segment objects at
252 multiple scales, for which DeepLab resamples a
253 given feature layer at multiple rates prior to convo-
254 lution.

255 The third challenge is caused by the built-in in-
256 variance of DCNNs, which is essential to local im-
257 age transformations and object-centric classifiers.
258 However, it reduces the spatial accuracy of dense
259 prediction tasks. DeepLab manages to capture fine
260 details by using a fully-connected Conditional Ran-
261 dom Field. Figure 4 shows the pipeline of DeepLab
262 for a semantic segmentation task.

263 4.2.2 SPADE

264 SPADE (Park et al., 2019) focuses on a specific
265 form of conditional image synthesis which converts
266 a semantic segmentation mask to a photo-realistic
267 image; they refer to this as semantic image synthe-
268 sis.

269 In conventional network architecture built for
270 semantic synthesis tasks, a stack of convolutional,
271 normalization and nonlinearity layers lead to sub-
272 optimal results because the semantic information
273 from input masks tend to get “washed away” by
274 normalization layers. The solution is the SPatially-
275 Adaptive (DE)normalization, or SPADE, which
276 uses the input layout for modulating activations in
277 normalization layers through a spatially-adaptive,
278 learned transformation.

279 Specifically, the semantic segmentation mask is
280 projected onto an embedding space, and convolved
281 to produce modulation parameters. These are ten-
282 sors with spatial dimensions rather than vectors.
283 Then, the parameters are multiplied and added to
284 the normalized activation element-wise. Each nor-
285 malization layer uses the segmentation mask to
286 modulate layer activations. Figure 5 shows the ar-
287 chitecture of the SPADE generator, which uses a
288 series of SPADE residual blocks.

289 5 Experiments

290 5.1 Naïve-DeepLab

291 Our first experiment started with the pre-trained
292 DeepLab segmentation model and pre-trained
293 SPADE model, which served as the baseline of
294 our experiments for generating synthesized photos
295 from the paintings. This naive pipeline showed
296 that synthesis results were affected by both seg-
297 mentation and synthesis models. The segmentation
298 model defined the boundaries of each element in
299 the input images, outputting labels and maps that
300 would affect the faithfulness of SPADE genera-
301 tion results. More specifically, if the segmentation
302 model performed poorly, such as classifying the sky
303 as the mountain, even though the synthesis model
304 generated mountains with high coherency and qual-
305 ity, the synthesized photo still would not match
306 with the input paintings. To improve the SPADE
307 model, as mentioned in (Park et al., 2019), a better
308 model was trained on the Flickr dataset, which was
309 composed of landscape photos, but both the dataset
310 and the pre-trained model were not published.

311 By observing results from our baseline approach,
312 we also found that synthesis quality for our specific

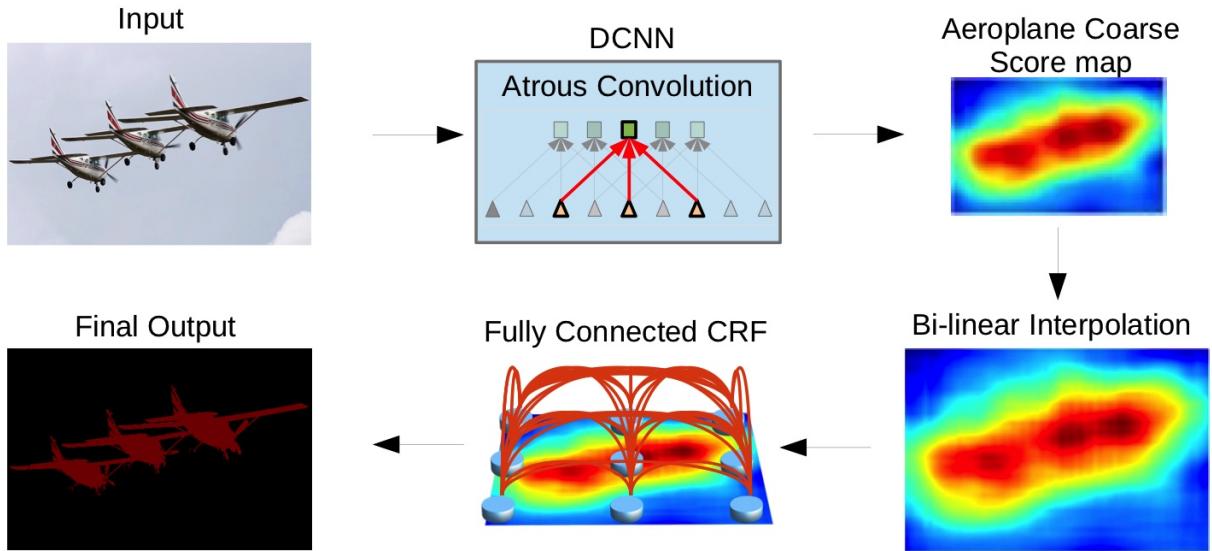


Figure 4: DeepLab architecture

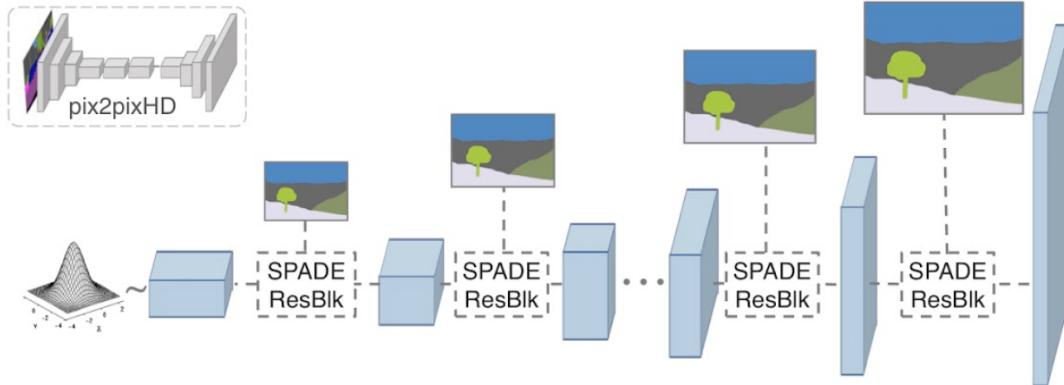


Figure 5: SPADE generator architecture

problem might be correlated with how detailed semantic segmentation was labeled in the datasets. For example, all types of trees, such as pines, willows, and maples, were labeled as trees, which resulted in a partial loss of information after the segmentation step. The synthesized photos hence contained inaccurate tree textures. Consequently, a dataset with more labels should help further improve the quality of the synthesized photos.

We first trained the segmentation model from scratch using Google Cloud Platform, and as mentioned in section 6.1, the experiments crashed after 800 epochs due to being out of CUDA memory. However, the IOU only increased to around 0.02, which is significantly lower than the released pre-trained model on the COCO dataset with an IOU of around 0.32; we thus decided to use the released model from DeepLab, considering the hardware

limitations and the time constraints.

As mentioned above, synthesis quality depended on both the segmentation and synthesis models. While the synthesis model could be improved through training on a Flickr dataset, the details were not released. Hence, our experiments focused on improving the segmentation model, and repurposing it to produce more accurate segmentation maps on AdaIN-Dataset.

5.2 DeepLab-Landscape

Our project aimed to synthesize photos from a set of paintings we selected, which mainly consisted of landscape contents. Since a large proportion of the COCO-Stuff dataset contained indoor objects, which was out of scope for our problem, we decided to experiment with a smaller and more specific dataset of landscapes. We envisioned a

313
314
315
316
317
318
319
320
321

322
323
324
325
326
327
328
329
330

331
332
333
334
335
336
337
338
339

340
341
342
343
344
345
346
347

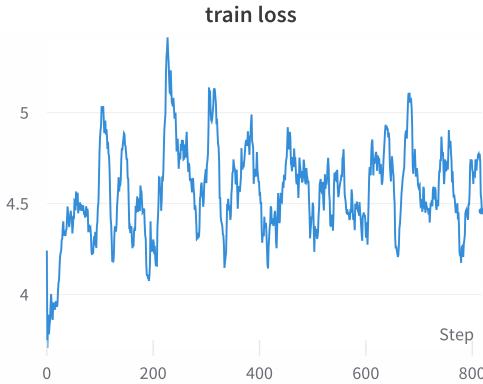


Figure 6: Deeplab-landscape training loss



Figure 7: Deeplab-AdaIn training loss

348 segmentation model more suitable to the problem,
 349 if we trained it on a dataset with themes that closely
 350 aligned with the common themes in our paintings.
 351 Therefore, we extracted a dataset with around 5000
 352 images from the COCO-Stuff dataset following
 353 the approach mentioned in 3.1.2. As shown in the
 354 training loss plotted in figure 6, the training loss
 355 did not converge. We then further tested the saved
 356 model on the test set (composed of the landscape
 357 photos). The result in 6.1.1 indicates that training
 358 in these selected landscape photos did not improve
 359 the semantic segmentation on landscape photos.

360 5.3 DeepLab-AdaIN

361 The experiments described in 5.2 attempted to
 362 improve the model through training on a dataset
 363 whose styles and contents conformed with the
 364 paintings better. As mentioned in 3.1.3, we
 365 did style transfers using the AdaIN model on
 366 Landscape-Dataset extracted from COCO-Stuff.
 367 Our baseline pipeline, as mentioned in 6.1, worked
 368 well with paintings in the simple category, but
 369 yielded relatively poor results on the test dataset.
 370 We thus decided to investigate if training on a
 371 dataset of Impressionism-like paintings would help
 372 the model better identify the objects in paintings.
 373 Given that the annotations of the painting segmen-
 374 tation were rather expensive, we used AdaIN to
 375 automate such a style transfer process on the land-
 376 scape dataset because the annotations of the seg-
 377 mentation were already available. The training loss
 378 of this experiment shown in figure 7 indicates that
 379 the training loss converged from 10 to 6 after 1000
 380 epochs.

381 6 Evaluation

382 6.1 Single Images

383 6.1.1 Quantitative Evaluation

384 Following the pipeline described above, we trained
 385 several models and compared their testing perfor-
 386 mances with metrics provided by DeepLab (Chen
 387 et al., 2016). The percentage measures of the mean
 388 pixel intersection-over-union (mIOU), frequency
 389 weighted IOU, mean accuracy, and pixel accuracy
 390 across all dataset classes on Landscape-Dataset are
 391 shown in Table 1. The same metrics from testing
 392 on AdaIN-Dataset are shown in Table 2.

393 We evaluated the naive DeepLab model with
 394 and without its pre-trained weights as baselines
 395 for model comparisons. A direct evaluation of the
 396 pre-trained DeepLab yielded 35.99% and 14.30%
 397 on our landscape and AdaIN-processed datasets
 398 respectively. This set the highest performance on
 399 vanilla landscape segmentation, and second highest
 400 on non-photo images (paintings) following our
 401 guess that DeepLab performance would degrade in
 402 such cases. Training DeepLab from scratch using
 403 the landscape dataset seemed to quickly saturated
 404 to a high loss, resulting in a significantly reduced
 405 testing mIOU value of 2.505%. After noticing
 406 the performance drop, we decided only to used
 407 DeepLab with pretrained weights in developing
 408 variants.

409 Table 1 and 2 also show the testing performance
 410 of our landscape- and painting-tuned DeepLab vari-
 411 ants. Our original intention behind fine-tuning
 412 DeepLab by re-using a subset of COCO-Stuff in
 413 the training set was to increase its exposure to land-
 414 scapes. However, this variant did not outperform
 415 the baseline, possibly because the landscape dataset
 416 was too small comparing to what DeepLab was

	mean-IOU	frequency weighted IOU	accuracy	pixel accuracy
Naïve-DeepLab-Pretrained	35.99	59.64	44.72	71.69
DeepLab-Landscape	32.97	53.33	41.34	69.04

Table 1: Mean IOU, frequency weighted IOU, mean accuracy, pixel accuracy (%) of different models tested on Landscape-Dataset. The original DeepLab model with pretrained weights (Naïve-DeepLab-Pretrained) achieved the highest mIOU.

	mean-IOU	frequency weighted IOU	accuracy	pixel accuracy
Naïve-DeepLab-Pretrained	14.30	27.99	21.02	38.14
DeepLab-Landscape	13.16	25.39	19.41	37.86
DeepLab-AdaIN	23.46	46.75	30.83	63.19

Table 2: Mean IOU, frequency weighted IOU, mean accuracy, pixel accuracy (%) of different models on AdaIN-dataset. Our re-trained DeepLab variant (DeepLab-AdaIN) performed the best.



Figure 8: Naïve-DeepLab pipeline on simple paintings

417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478

trained on. DeepLab-Landscape performed worse on paintings due to its compromised segmentation accuracy. We therefore changed directions and designed the DeepLab-AdaIN model variant independently from the Landscape dataset. This variant directly trained DeepLab on the style-transferred landscape paintings generated using AdaIN, and actually achieved a $1.64 \times$ improvement on mIOU when compared to vanilla DeepLab on paintings.

6.1.2 Human Evaluation

Our Naïve-DeepLab approach can achieve our goal in simple paintings where the objects are close to the real-world objects as shown in figure 8. However, it does not work well for more paintings, such as impressionism paintings, whose object textures or colors differ from the real world. To evaluate if the deepLab-AdaIN model improves the result with the complex category of the paintings, human assisted evaluation is introduced.

The quantitative evaluation where we compared the metrics between the experiments mentioned above claims a significant improvement after we trained the segmentation model on the style transferred landscape photos dataset, which indicates its potential in better synthesizing photos from the paintings. However, to confirm the improvement

of the segmentation model trained based on the style transferred landscape photos dataset, we leveraged human assistance for the evaluation. Figure 9 shows a few randomly selected examples from our test result, wherein each example, photo 1 is the result denoted by the naive deepLab pipeline as mentioned in 5.1, and photo 2 is the result denoted by the deepLab-AdaIN pipeline as mentioned in 5.3. We then designed five questions using Google form to collect responses on the question, to what extent these synthesized photos resemble the original painting. We asked the users to select from the five options; photo 1 resembles better significantly, photo 1 resembles better slightly, cannot tell, photo 2 resembles better significantly, and photo 2 resembles better slightly. To analyze these results quantitatively, we encoded these options as -1, -0.3, 0, 0.3, 1, where the vote for significantly better is considered as one vote for the corresponding pipeline, and the vote for slightly better is considered as a partial vote. In consequence, if the average of the responses on an image is positive, the majority of the responses believe that the photos synthesized by the deepLab-AdaIN pipeline as mentioned in 5.3 are better; otherwise, the photos synthesized by the naive deepLab pipeline are better. The average score for each example is shown in figure 9 is shown in the following Table 3. We can then find that all examples have a positive average vote score, which indicates that most of the responses believe that the synthesized photos from these paintings using the approach defined in 5.3 is better than the approach defined in 5.1 for all of these examples. The second row and the third row in Table 3 indicate the percentage of responses in favor of the synthesized photos from the approach defined in

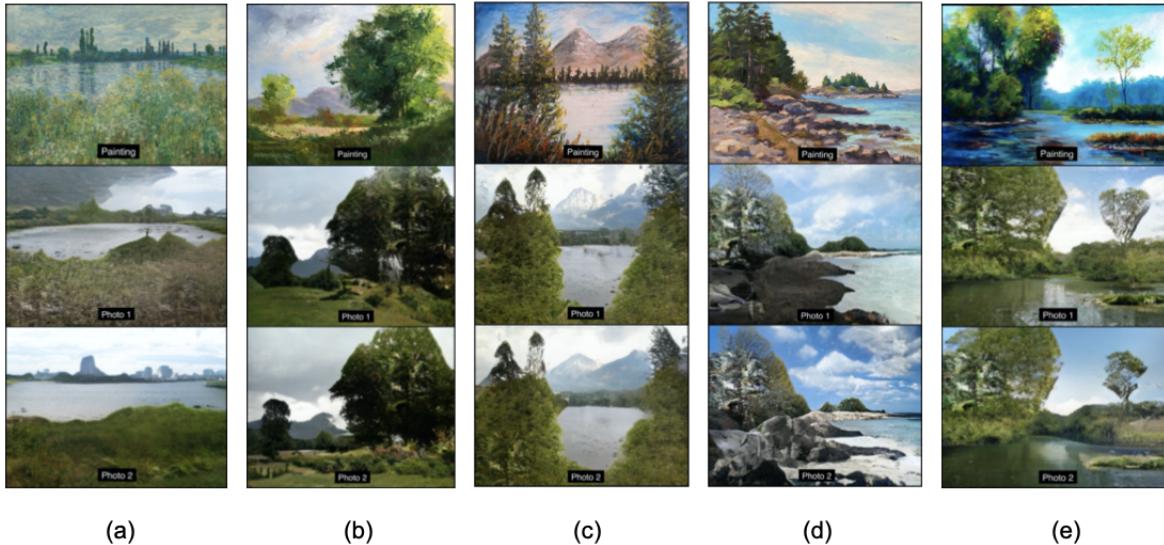


Figure 9: Comparisons of synthesized photos from naïve deepLab pipeline and deepLab-AdaIN pipeline

	example a	example b	example c	example d	example e
avg vote score	0.02	0.505	0.1	0.625	0.34
% responses in favor of photo 2	35	85	50	90	70
% responses not in favor of photo 1	50	90	60	90	85

Table 3: Human assisted evaluation

479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
5.3 and the percentage of responses not in favor of
the synthesized photos from the approach defined
in 5.1. According to these examples, the most con-
troversial one is the example (a), as shown in the
figure 9, and we discussed with some of the inter-
viewees who believe the photos generated using
the naive approach are better; these interviewees
believe that the trees in the distance are not cor-
rectly identified in the photos synthesized using the
deepLab-AdaIN pipeline. However, from us as de-
veloper’s viewpoint, the style transfer has resolved
the problem in the naive pipeline, where the sky is
misclassified as the mountain due to the texture of
the painting. Although the trees in the distance are
classified as buildings, these unusual shapes of the
trees were not seen in the training set. Thus we still
consider the example (a) as an improvement.

496 6.2 Images in Sequence (Video Frames)

497 In addition to the synthesized photos from single
498 paintings, we explored photos synthesis from a se-
499 quence of frames in animated videos as another use
500 case. Each frame in the animated videos could be
501 classified as an simple category instance in our test
502 dataset. This experiment is to explore if the syn-
503 thethesized photos from the video clips are consistent

504 with their shared area.

505 As shown in figure 10, although the segmenta-
506 tion model, even using the naive approach, can
507 segment the animated video frames with relatively
508 high quality, this experiment still indicates two po-
509 tential problems of approach. The first problem
510 is that some of the information is lost after the
511 segmentation; for example, ravines between the
512 mountains are not generated in our result because
513 this information is lost after the segmentation stage.
514 Such a problem might be solved if the labeling on
515 the segmentation is more abundant, given that there
516 is no label like ravines on the COCO dataset. Ad-
517 ditionally, when we generate a new video based
518 on a sequence of synthesized photos, it is not as
519 consistent as the original animated video. To be
520 more specific, only the top part of the mountain
521 changes across different frames. The rest of the
522 mountain is almost not changing, which is reason-
523 able because our training was not trained on any
524 sequences of the video clips, and the consistency
525 between the video frames was not considered in
526 our architecture.

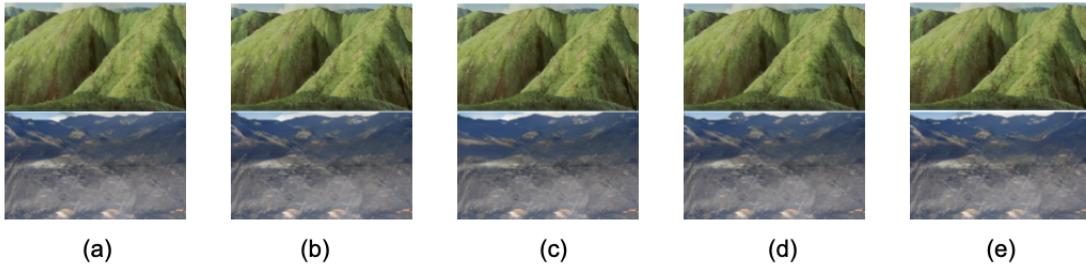


Figure 10: Synthesized photos from a sequence of animated video

527 7 Future Works

528 For future works, we would like to continue to
 529 improve the quality of the photos synthesis from
 530 paintings by resolving the current issues in our
 531 pipeline. For example, we may add another tone
 532 transfer stage after the SPADE outputs so that the
 533 tone of the synthesized photo would be closer to
 534 the original paintings. In addition to that, the
 535 experiments mentioned in 6.2 exposed a problem in
 536 a sequence of video frames, and we would like to
 537 resolve this issue by training with a sequence of
 538 images for the SPADE model.

539 8 Conclusion

540 To sum up, our initial approach with the naive seg-
 541 mentation model and the SPADE model worked
 542 for paintings depicting objects that were close
 543 to real-world objects. However, since more
 544 paintings lacked such nice traits, our naive ap-
 545 proach performed relatively poorly on them. We
 546 thereby trained our segmentation model using style-
 547 transferred paintings, and improved the quality
 548 of the synthesized photos. Given a shortage of
 549 existing datasets with segmentation annotations
 550 for paintings, as well as the high costs of man-
 551 ual segmentation, our segmentation model did not
 552 use human-created paintings with segmentation
 553 for training. With a dataset consisted of machine-
 554 generated style-transferred photos, our model man-
 555 aged to synthesize photos from paintings with
 556 relatively high quality.

557 9 Appendix

558 [1] Part of our approach is derived from the follow-
 559 ing implementations:

560 <https://github.com/NVlabs/SPADE>;
 561 <https://github.com/kazuto1011/deeplab-pytorch>;
 562 <https://github.com/naoto0804/pytorch-AdaIN>

[2] Please refer to our wandb project for the loss
 curve and the evaluation metrics:

<https://wandb.ai/qx2217/painting-synthesis/table?workspace=user-qx2217>

[3] Please refer to Google form for the human
 evaluation results:

https://docs.google.com/forms/u/1/d/1ehFnRNyAQMTv_j7dvWUA9xT7P2TGRNj0rfiRzoYsIs/edit?usp=sharing

572 References

Saad Albawi, Tareq Abed Mohammed, and Saad Al-Zawi. 2017. [Understanding of a convolutional neural network](#). In *2017 International Conference on Engineering and Technology (ICET)*, pages 1–6.

Holger Caesar, Jasper R. R. Uijlings, and Vittorio Ferrari. 2016. [Coco-stuff: Thing and stuff classes in context](#). *CoRR*, abs/1612.03716.

Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille. 2016. [Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs](#). *CoRR*, abs/1606.00915.

Antonia Creswell, Tom White, Vincent Dumoulin, Kai Arulkumaran, Biswa Sengupta, and Anil A. Bharath. 2018. [Generative adversarial networks: An overview](#). *IEEE Signal Processing Magazine*, 35(1):53–65.

A. Elgammal, Bingchen Liu, Mohamed Elhoseiny, and Marian Mazzone. 2017. Can: Creative adversarial networks, generating "art" by learning about styles and deviating from style norms. In *ICCC*.

Leon A. Gatys, Alexander S. Ecker, and Matthias Bethge. 2015. [A neural algorithm of artistic style](#). *CoRR*, abs/1508.06576.

Xun Huang and Serge J. Belongie. 2017. [Arbitrary style transfer in real-time with adaptive instance normalization](#). *CoRR*, abs/1703.06868.

Tero Karras, Samuli Laine, and Timo Aila. 2018. [A style-based generator architecture for generative adversarial networks](#). *CoRR*, abs/1812.04948.

602 Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. 1998.
603 **Gradient-based learning applied to document recog-**
604 **nition.** *Proceedings of the IEEE*, 86(11):2278–2324.

605 Tsung-Yi Lin, Michael Maire, Serge J. Belongie,
606 Lubomir D. Bourdev, Ross B. Girshick, James Hays,
607 Pietro Perona, Deva Ramanan, Piotr Dollár, and
608 C. Lawrence Zitnick. 2014. **Microsoft COCO: com-**
609 **mon objects in context.** *CoRR*, abs/1405.0312.

610 Shervin Minaee, Yuri Y. Boykov, Fatih Porikli, Anto-
611 nio J Plaza, Nasser Kehtarnavaz, and Demetri Ter-
612 zopoulos. 2021. **Image segmentation using deep**
613 **learning: A survey.** *IEEE Transactions on Pattern*
614 *Analysis and Machine Intelligence*, pages 1–1.

615 Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and
616 Jun-Yan Zhu. 2019. **Semantic image synthe-**
617 **sis with spatially-adaptive normalization.** *CoRR*,
618 abs/1903.07291.

619 S. Smithsonian Magazine. 2016. An ai-written novella
620 almost won a literary prize.