

TEXT-TO-IMAGE SYNTHESIS USING KT GAN

Rajaiah M^{1,a)}, Subramanyam N^{2,b)}, Vijaya Kumar C^{3,c)} and
Chandra Kanth P^{4,d)}

¹Professor &HOD, Department of Computer Science Engineering, Audisankara College of Engineering and Technology, Gudur, India

^{2,3}PG Scholar, Department of Computer Science Engineering, Audisankara College of Engineering and Technology, Gudur, India

⁴Associate Professor, Department of Computer Science Engineering, Audisankara College of Engineering and Technology, Gudur, India

^{a)} rajagopal1402@gmail.com

^{b)} Corresponding author: subramanyam.naguru@gmail.com

^{c)} vijay.mba.vsu@gmail.com

^{d)} chandrakanthc4u@gmail.com

ABSTRACT_ This assignment gives a fresh model, Knowledge-Transfer Generative Adversarial Network (KT GAN), for text-to-image generation. We use two new methods an Alternate Attention-Transfer Mechanism (AATM) and a Semantic Distillation Mechanism (SDM), to assist generator higher connect the inter functional hole between textual content and image. The AATM modifies phrase interest weights and interest weights of picture subdivision one after other, to gradually spotlight necessary phrase statistics and enrich small print of generated images. The semantic distillation mechanism makes use of picture encoder skilled in the Image to Image project to information coaching of textual content encoder in the Text to Image process, for producing higher textual content points and greater excellent pictures. By significant investigational testing on two public data sets, KTGAN surpass the existing approach notably, and additionally attains the comparative effects over one-of-a-kind comparison metrics.

1.INTRODUCTION

Text-to-Image synthesis objectives to create a practical picture that is significantly constant with a given textual content, with the aid of studying a bounding between the semantic textual content area and complicated RGB picture area. A key undertaking in generating sensible objects with significant small print is the assorted hole between powerful ideas in textual content descriptions and pixel-level contents of artificial pictures. Building a wonderful generator to fill this area hole is tough. Large no of processes

based totally on Generative Adversarial Networks (GANs) fill the area hole by means of making use of a discriminator to compare the generated text-image pair and the real pair. However, such a discriminator on my own is normally inadequate to mannequin underlying semantic consistency between textual content and photo, and consequently, consequences in semantic or structural blunders in generated pix (view Fig 1, the "Direct T2I" row). Currently, the interest method has been used to solve this problem and this courses the generator to higher in shape positive visible phrases

with corresponding picture sub regions. But the usage of word-level interest by myself does now not make sure world semantic persistency because of the variety between the textual content and photograph methods. Hence, the MirrorGAN fashions Text2Image and Image2Text collectively to beautify world inter functional semantic persistency. However, the Image2Text in MirrorGAN is nonetheless a inter functional creation, which is now not simpler than comparable technology assignment such as I2I task. Hence, the hassle of semantic disparity in between diverse records nonetheless remains same. SEGAN develops a novel contrastive loss and a Semantic Consistency Module (SCM) to higher line up the generated photo and the floor reality in function space. Nevertheless due to the diversified semantic disparity, SEGAN can't pull out fine textual content aspects that can information the synthesis of practical and targeted images.

2.LITERATURE SURVEY

2.1 Automatic adaptation of object detectors to new domains using self-training

AUTHORS: A. RoyChowdhury et al.,

ABSTRACT: The goal of this research is to unsupervised adapt an present object detector to a novel target domain. We think that there are a numerous unlabeled videos in this field. We get labels on the target data automatically through combining high-confidence detections from the current detector with rigid examples obtained by leveraging temporal cues with a tracker. These labels are then used to retrain the real model using the automatically generated labels. We

propose a updated knowledge distillation loss and enquire various methods for assigning soft-labels to target domain training examples. Our method has been empirically tested on difficult face and pedestrian detection works: a face detector trained on wider-Face, which has more quality pictures dragged from the web, is altered to a more amount of surveillance data set; a pedestrian detector trained on transparent, dayhours images from the BDD-100K driving data set is altered to all other scenes, which includes rainy, foggy, and night-time. These findings show the importance of using real-world examples gleaned via tracking, the benefit of using soft-labels via distillation loss vs hard-labels, and also promising performance as a simple method for unsupervised domain adaptation of object detectors with little reliance on hyper-parameters.

2.2 Knowledge distillation for end-to-end person search

AUTHORS: M. Bharti, G. Fabio, and A. Sikandar

ABSTRACT: We show how to execute an end-to-end person search using knowledge distillation. In person search, end-to-end approaches are the present state of the art. because they solve both detection and re-identification problems at the same time. Because of a poor detector, these joint optimization approaches exhibit the greatest drop in performance. In a teacher-student scenario, we offer two separate approaches for further supervision of end-to-end person search strategies. The first one is based on cutting-edge object detection knowledge distillation. We utilise this to use a specialised detector to oversee the detector of our person search model at various stages. The next approach

is fresh, simple, and significantly better efficient. Using a pre-computed look-up table of ID attributes, this collects information from a teacher re-identification process. It allows the student to relax while learning identification traits, allowing him or her to focus on the detection assignment. This technique not just helps in the correction of improper detector training in the combined optimization while also enhancing the search for people. However, in this scenario, model compression reduces the performance gap between the teacher and the student. On two benchmark data sets, We show that two current state-of-the-art approaches can be significantly improved using our developed knowledge distillation approach. Furthermore, our application compares the accomplishment of minor and major models in the model compression challenge.

2.3 Controllable text to image generation

AUTHORS: L. Bowen, Q. Xiaojuan, L. Thomas, and H. S. T. Philip

ABSTRACT: We introduce a new controlled text-to-image generative adversarial network in this paper and that can efficiently generate excellent standard images while also controlling components of the image creation based on natural language descriptions. We offer a spatial and channel-wise attention-driven generator that can untangle several visual properties at the word level and allowing the model to concentrate on creating and modifying subregions for the most important terms. By connecting words with picture regions, a word-level discriminator is also presented to provide

fine-grained supervisory feedback, allowing for the guiding of an efficient generator capable of manipulating particular features while not disturbing the creation of additional content. In addition to that, perceptual loss is used to decrease the uncertainty in image creation and to inspire the generator to alter certain properties needed in the updated text. Large scale testing on standard data sets show that our model exceed the current technology and can manage fabricated images effectively making use of natural language descriptions.

3. PROPOSED SYSTEM

Within this paper author is using GAN model (generative Adversarial Network) to convert text to images. In propose paper author modifying GAN with transfer learning to accommodate text with images so generator model get trained on TEXT and discriminator model get trained on images with embed text and when we give any text then generator GAN model will predict equivalent image for given text.

In propose paper an Alternate Attention Transfer Mechanism (AATM) and a Semantic Distillation Mechanism (SDM), to assist the generator in reducing the a inter functional chasm separating text and image. The AATM alternately modifies the attention weights of phrases and Attention weights for image sub-regions to constantly focus important word information and upgrade the quality of generated images. The SDM takes help of image encoder which is trained in the image to image task to assist the training of text encoder which is used in the text to image task for developing more exceptional images.

3.1 METHODOLOGY

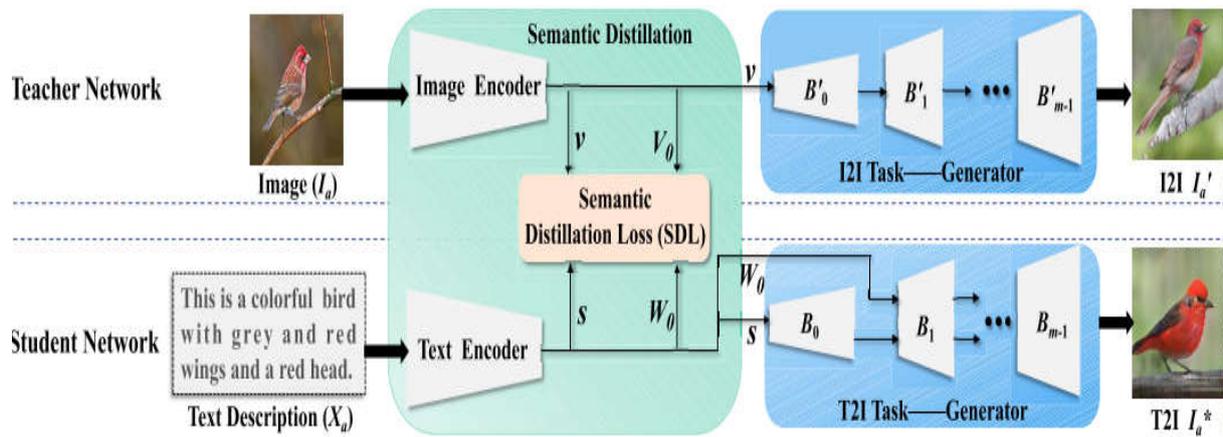


Fig 1: Work Flow

We propose a new Knowledge-Transfer Generative Adversarial Network (KTGAN) with two new mechanisms for Text to Image synthesis: (1) a Semantic Distillation Mechanism (SDM) which includes the usage of image encoder to guide text encoder for better image quality. (2) an Alternate Attention Transfer Mechanism (AATM) to identify higher important words from the text.

The following are the paper's main contributions:

- (i) We developed a Semantic Distillation Mechanism which has a new Semantic distillation loss function (SDL) and using this we are able to guide the text to image task using image to image task for better results.
- (ii) We introduced an Alternate Attention Transfer Mechanism to frequently update the attention weights and increase the image quality.
- (iii) We tested our KT-GAN on two different data sets CUB-Bird and large

In 'birds/birdname/.txt' file contains text for each bird and this you can see in below screen

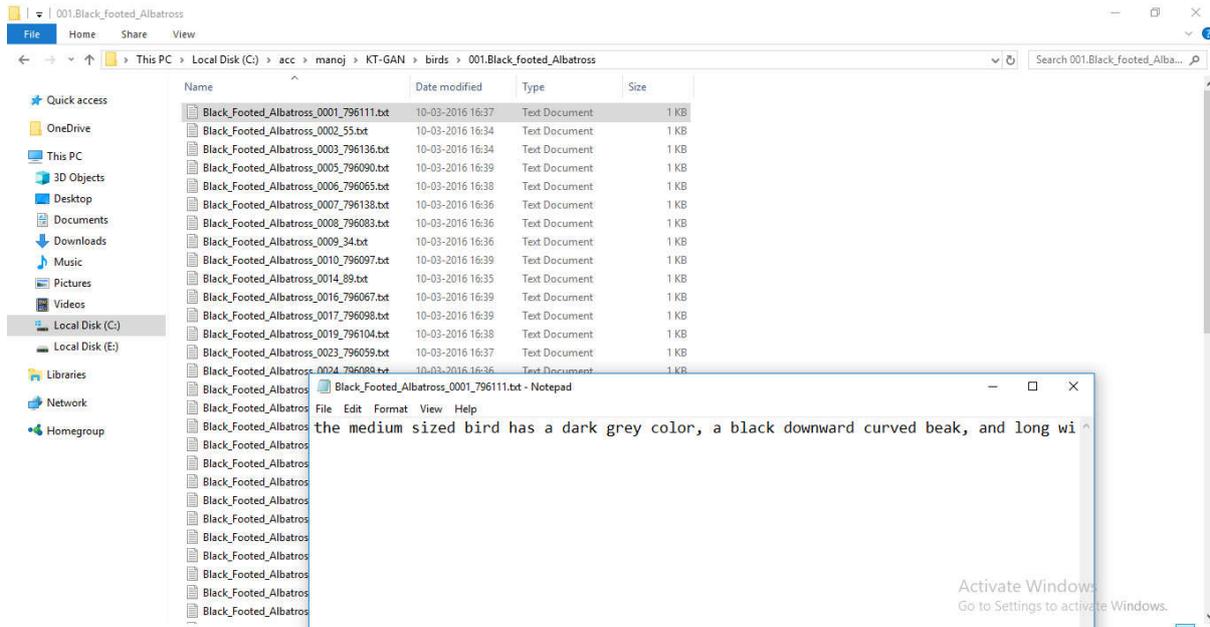
scale MS-COCO. Experimental effects and study shows the importance of KT-GAN and extensively expanded overall effect in contrast towards preceding modern strategies on all 4 comparison metrics.

We created the following modules to help us carry out this project.

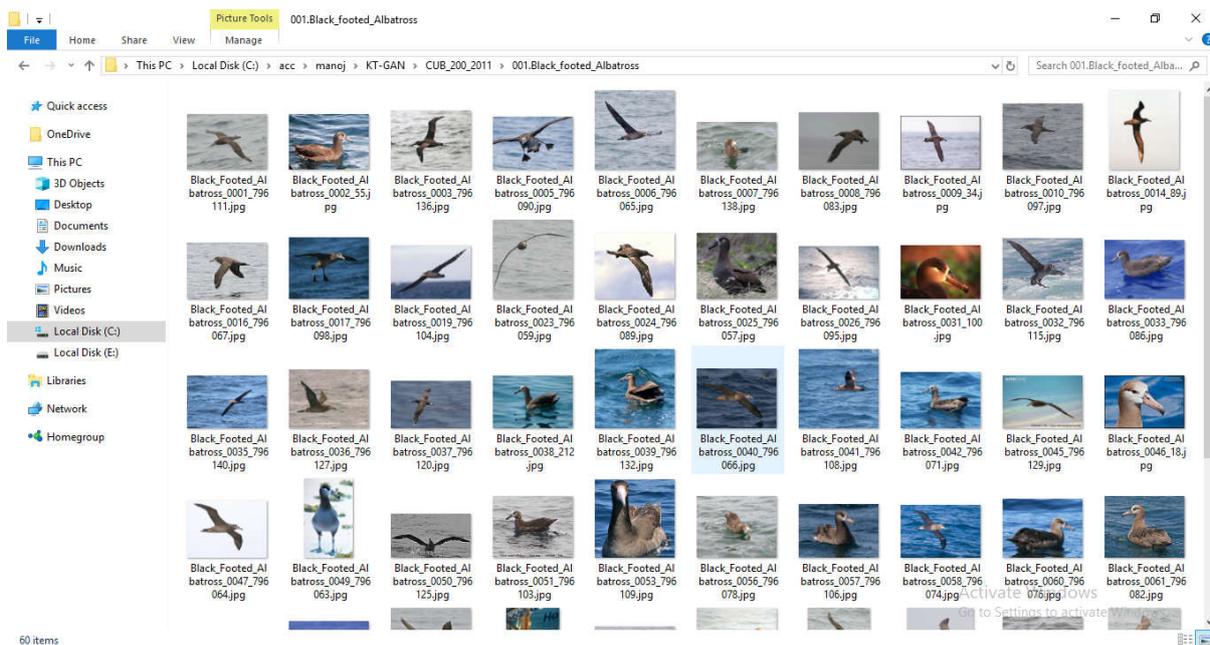
- 1) Upload CUB-Bird Data set: using this module we will upload dataset folders to application
- 2) Generate & Load KT-GAN Model: using this module we will read all images and then generate KT-GAN model
- 3) Generate Image from Text: using this module we will input TEXT and then KT-GAN will generate image from that text

3.2 ABOUT DATA SET

To implement this project author has used CUB-Bird dataset which contains TEXT and images and by using both TEXT and images we will train KT-GAN model and in below screen we are showing dataset details



In above screen we can see bird description text for each bird and in below screen we can see images of all those birds and this images you can find inside 'CUB_200_2011/bird_name/' folder like below screen



So by using above TEXT and images we will train KT-GAN model.

4.RESULTS AND DISCUSSION



Fig 2: In above screen, model is generated and in text field enter some bird description to get image



Fig 3: In above screen in text field I entered bird description as ‘bird with a white breast and a black crown and black webbed feet’ and then press ‘Generate Image from Text’ button to get below output

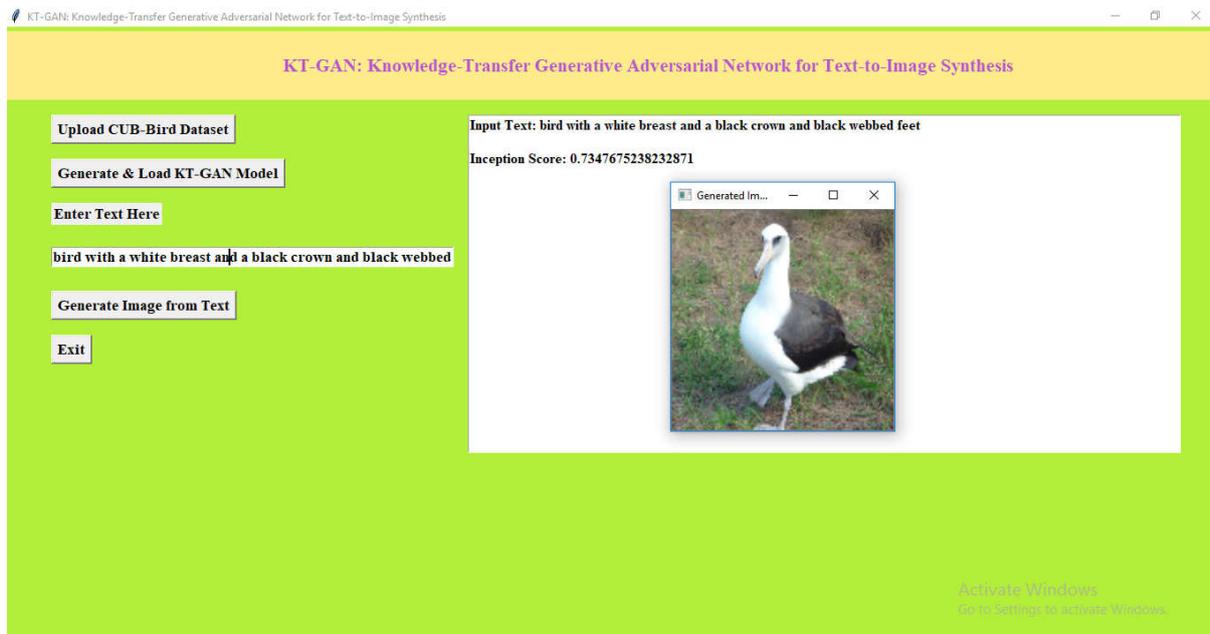


Fig 4: For given bird description we got above image and now try another description for bird

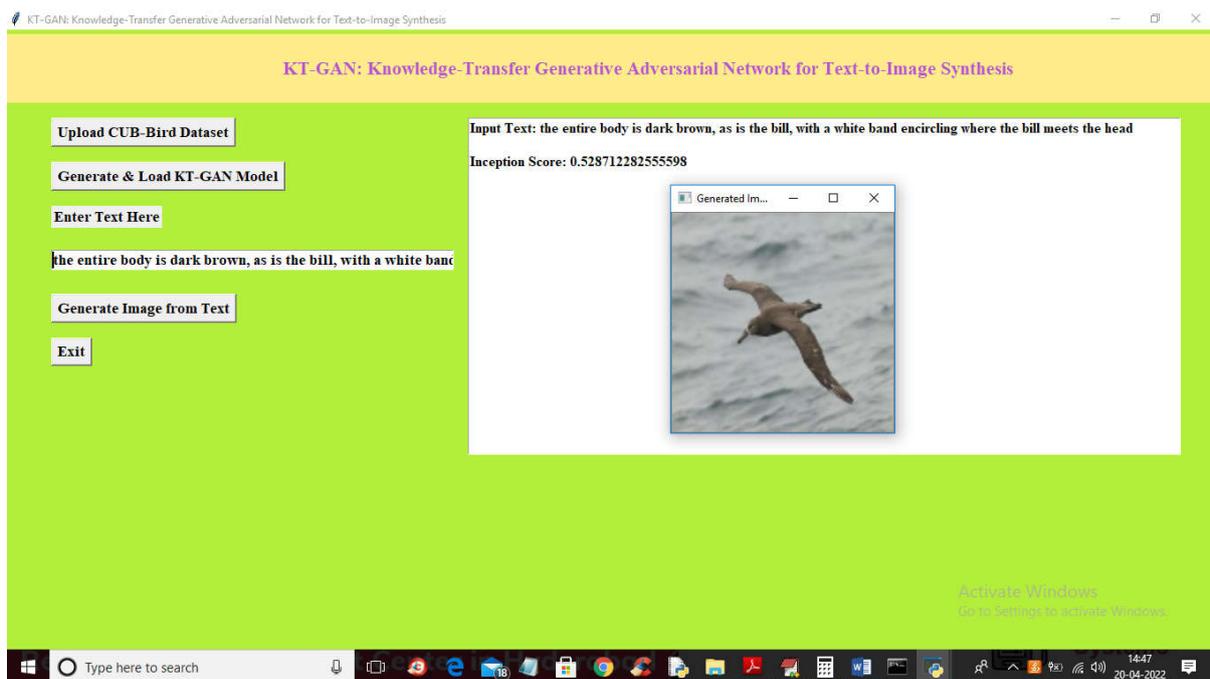


Fig 5: For above description we got the bird image as shown in the above figure

5.CONCLUSION

Here in this, we introduced a modern approach for generating images from text and this is done using the two main mechanisms namely Alternate Attention Transfer Mechanism (AATM) and

Semantic Distillation Mechanism (SDM). Using these two mechanisms we had completed this Knowledge-Transfer Generative Adversarial Network (KT-GAN) for Text-to-Image (T2I) synthesis. The SDM uses image encoder trained in image-to-image task to guide the text

encoder for text to image task. This involves the knowledge flow from image encoder to text encoder which leads to the improvement in the quality of images generated. The AATM is used to assign weights to the words and through this we are able to identify the important words. Using these two mechanisms we are able to decrease the heterogeneous gap and able to generate better quality images. The results of KTGAN shows that the performance of this is better compared to the previously used methods.

REFERENCES

- [1] A. RoyChowdhury et al., “Automatic adaptation of object detectors to new domains using self-training,” in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), Jun. 2019, pp. 780–790.
- [2] M. Bharti, G. Fabio, and A. Sikandar, “Knowledge distillation for end-to-end person search,” in Proc. BMVC, 2019, pp. 1–16.
- [3] L. Bowen, Q. Xiaojuan, L. Thomas, and H. S. T. Philip, “Controllable text-to-image generation,” in Proc. NeurIPS, 2019, pp. 2065–2075.
- [4] C. Ledig et al., “Photo-realistic single image super-resolution using a generative adversarial network,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), Jul. 2017, pp. 4681–4690.
- [5] F. Faghri, J. David Fleet, J. R. Kiros, and S. Fidler, “Vse++: Improved visual-semantic embeddings,” in Proc. BMVC, 2018, pp. 1–9.
- [6] M. Nicolás Guil Francisco Castro and J. Manuel Marín-Jiménez, “End-to-end incremental learning,” in Proc. ECCV, Sep. 2018, pp. 233–248.
- [7] F. Tung and G. Mori, “Knowledge distillation based on similarity,” in Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV), Oct. 2019, pp. 1365–1374.
- [8] H. Geoffrey, V. Oriol, and D. Jeff, “Distilling the knowledge in a neural network,” in Proc. NeurIPS Workshop, 2015, pp. 1–9.
- [9] J. Ian Goodfellow et al., “Generative adversarial nets,” in Proc. NeurIPS, 2014, pp. 2672–2680.
- [10] G. Yin, B. Liu, L. Sheng, N. Yu, X. Wang, and J. Shao, “Semantics disentangling for text-to-image generation,” in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), Jun. 2019, pp. 2327–2336.
- [11] H. Zhang et al., “StackGAN++: Synthesis of realistic images using stacked generative adversarial networks,” IEEE Trans. Pattern Anal. Mach. Intell., vol. 41, no. 8, pp. 1947–1962, Aug. 2019.
- [12] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, “GANs trained by a two time-scale update rule converge to a local nash equilibrium,” in Proc. NeurIPS, 2017, pp. 6626–6637.
- [13] H. Tan, X. Liu, X. Li, Y. Zhang, and B. Yin, “Adversarial nets with semantic enhancements for text-to-image synthesis,” in Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV), Oct. 2019, pp. 10501–10510.
- [14] W. Huang, Y. Xu, and I. Oppermann, “Realistic image generation using region-phrase attention,” 2019, arXiv:1902.05395. [Online]. Available: <http://arxiv.org/abs/1902.05395>

[15] J. Li, K. Fu, S. Zhao, and S. Ge, “Spatiotemporal knowledge distillation for efficient estimation of aerial video saliency,” *IEEE Trans. Image Process.*, vol. 29, pp. 1902–1914, 2020