

DCIL: Deep Contextual Internal Learning for Image Restoration and Image Retargeting

Indra Deep Mastan and Shanmuganathan Raman
 Indian Institute of Technology Gandhinagar
 Gandhinagar, Gujarat, India
 {indra.mastan, shanmuga}@iitgn.ac.in

Abstract

Recently, there is a vast interest in developing unsupervised methods that are independent of the feature learning from the training data, e.g., deep image prior [24], zero-shot learning [21], and internal learning [20]. These methods are based on the common goal of maximizing the quality of image features learned from a single image despite inherent technical diversity. In this work, we bridge the gap between the various unsupervised approaches above and propose a general framework for image restoration and image retargeting. We use contextual feature learning and internal learning to improvise the structure similarity between the source and the target images. We perform image resizing application in the following setups: classical image resizing using super-resolution, a challenging image resizing where the low-resolution image contains noise, and content-aware image resizing using image retargeting. We also compare our framework with relevant state-of-the-art methods.

1. Introduction

Deep learning based supervised models could implicitly capture the image prior by feature learning on a collection of images [16, 28, 6, 27, 25, 30, 2]. However, deep feature learning using training data could suffer from transformation bias or model collapse [31, 9]. Recently, there is a vast interest in using a convolutional neural network (CNN) to minimize the use of training samples [24, 19, 11, 21, 17, 20]. More specifically, the following unsupervised models are remarkably successful [21, 24, 20]. Unsupervised models could allow image restoration when the degradation process is complex and/or unknown and obtaining realistic data for supervised training is difficult [24].

Unsupervised image feature learning attracts various applications such as image super-resolution, inpainting, and image retargeting. Shocher *et al.* proposed zero-shot super-resolution (ZSSR) which does not use any training dataset

[21]. Another research thread for training-data independent methods is Deep Image Prior (DIP) proposed by Ulyanov *et al.* [24]. DIP bridges the gap between handcrafted image prior based classical methods and CNN based deep prior. It shows that the structure of the encoder-decoder network itself works as the image prior. Later, Raman and Mastan gave a generalization of [24] and showed various aspects of the relationship between network construction and image restoration [17]. For example, skip connections improve super-resolution but adversely affect image inpainting [17].

ZSSR and DIP compute pixel-to-pixel loss (e.g., mean squared error MSE). The pixel-to-pixel loss is limited to the applications which have a spatial correspondence between the pixels of the source and the target images (aligned image data). Mechrez *et al.* have proposed contextual loss for non-aligned image data applications, e.g., style transfer [19, 18]. However, their approach is not completely independent of training samples. We call image feature learning by minimizing the contextual loss as contextual feature learning (CFL).

Recently, Shocher *et al.* proposed Internal-GAN (InGAN) for image retargeting without using any training samples. Image retargeting requires feature transfer when there is no spatial correspondence between pixels of the source image and the target image (non-aligned image data) [20]. InGAN considers image retargeting as a distribution matching problem to take advantage of GAN. Shocher *et al.* observed that the reconstructions suffer from the object partition ambiguity.

There is a technical diversity in the unsupervised methods described above. However, they are all subjected to maximizing the quality of image feature learning from a single image. There are two interesting challenges here. (1) What aspect of the network would help for the task of image generation in the limited contextual understanding due to the lack of feature learning from the training data? (2) What should be the structure of the loss function when the source image and the target image are non-aligned and do not have spatial correspondence?

To better understand the challenges above, let us consider the task of resizing an image. Super-Resolution (SR) scales the entire image, whereas image retargeting resizes the input while preserving the size and the aspect ratio of the local elements [20]. Another image resizing application would be to scale a low-resolution image, which contains noise, termed as Denoising-Super-Resolution (DSR). Image resizing in DSR setting is more general and challenging than image super-resolution as it also removes noise from a low-resolution image.

The training data independent image resizing in various scenarios described above would require a careful design of network and loss function. DIP and ZSSR perform image restoration using pixel-to-pixel loss. Therefore, they are not applicable to image retargeting (non-aligned image data). InGAN performs image retargeting, but is not studied for DSR setting [20].

We propose deep contextual internal learning (DCIL) for image retargeting and image restoration. Our models include structure of the network as the implicit image prior and an image degradation based loss term. The desired image is reconstructed by finding the optimal solution of the model. The key idea is to maximize deep feature learning from a single image by a modular network structure with a generalized loss function which works for aligned and non-aligned image data. We use diverse techniques such as deep prior learning, adversarial learning, and CFL. Deep prior learning fits the generator network by maximizing the likelihood of weights given the corrupted image and restoration model. Adversarial learning and CFL perform distribution matching to generate realistic image patches.

There is an interesting contrast to our objectives. On the one hand, we need an image restoration strategy which enhances the image features present in the input corrupted image. On the other end, image retargeting specific reconstruction, which uses similar image features from input image for synthesizing objects in the output image (distribution matching). Image distribution matching between the corrupted image and the output image could adversely influence output image. Therefore, we provide modularity in the network structure with a task-specific loss function (Sec. 3). The network provides a high impedance to the noise and allows reconstruction of the signal [24]. The loss function influence the quality of learned features by minimizing the dis-similarity in the features of the source image and the output image.

Network construction in DCIL is based on various network components. These components are as follows: network depth, skip connections, a cascade of network input, and network composition. We do not use cascading of network input as it does not provide a significant enhancement of prior learning [17]. We take advantage of residual blocks as it improves the generator output [14, 32].

We formalize network construction and explain the abstract description of the network to simplify the network design in the presence of diverse components (Sec. 3.1).

After network construction, DCIL iteratively minimizes the loss between the source image and the target image. DCIL loss compares image features between the source image and the target image in three ways: pixel-to-pixel, patch-to-patch, and contextual features comparison. The motive behind this loss is to capture better image statistics by comparing the diverse set of image features. (1) Pixel-to-pixel comparison is similar to the MSE based reconstruction loss [24, 20]. (2) Patch-to-patch comparison is made using the adversarial loss [21, 20]. (3) Contextual features comparison between the source and the target image is done using the contextual loss [19].

Adversarial loss generate realistic samples by preserving the distribution of image patches [20]. Contextual loss is motivated to enhance the structural similarity of the objects in the output image [18]. The reconstruction loss ensures the preservation of global image features in the target image.

Image resizing using DSR and SR are both naturally occurring. We corrupt the input image to a high degree to observe the quality of image features captured in DIP [24] and CFL [19]. We show that DCIL generates reconstruction that is comparable to that of the other relevant unsupervised frameworks (Sec. 4.1). Mechrez *et al.* have shown that CFL exhibits natural internal statistics for SR [18]. We illustrate the performance of CFL in the training data-independent setup for SR task (Sec. 4.2).

DCIL performs image retargeting by changing the size and the shape of the generator output. The target image is subjected to preserve the distribution of image patches. Adversarial loss and contextual loss are well suited for the above task. More specifically, the contextual feature learning with the internal patch distribution learning (InGAN) is observed to preserve good object statistics for image retargeting task (Sec. 4.3).

The key contributions of the paper are as follows.

1. We propose a generalized framework (DCIL) for image resize in various scenarios by coupling an internal learning scheme in a *novel* unsupervised contextual feature learning framework (Sec. 3 and Table 2).
2. We verify effectiveness of DCIL by extensive experimentation for DSR, SR, and image retargeting tasks (Sec. 4).
3. To the best of our knowledge, we are the first to study CFL in an unsupervised framework for DSR task (Table 1).
4. DCIL preserves the object structure and alignment in the image retargeting output (Fig. 6 and Fig. 5). We also give ablation studies for understanding various aspects of the loss functions (Sec. 5).

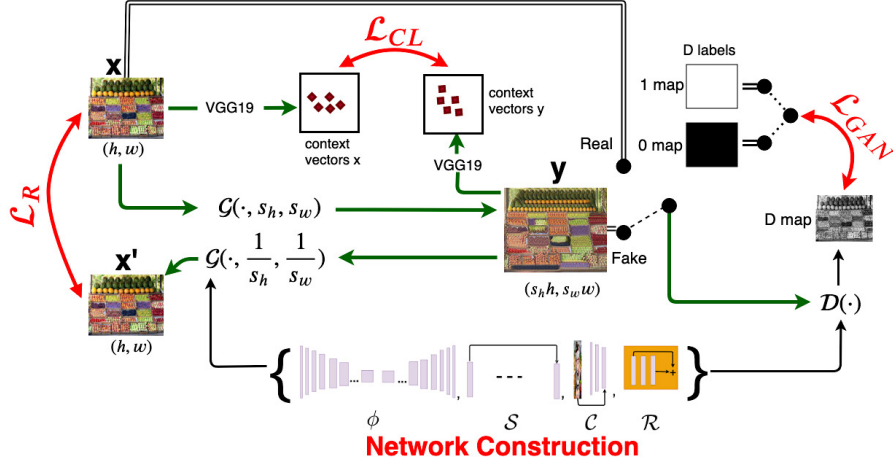


Figure 1: **Deep Contextual Internal Learning (DCIL)**. The figure shows the pictorial representation of the DCIL framework for image retargeting. It contains a generator \mathcal{G} that takes an input image x of size (h, w) and outputs an image y of a different size using the scaling factors (s_h, s_w) . The output of the generator $y = \mathcal{G}(x, s_h, s_w)$ is fed into the discriminator \mathcal{D} and feature extractor pre-trained *VGG19* network [23]. The same framework is used for image restoration where the definitions of the loss functions is different. The idea here is to create the generator and discriminator using network construction module and then iteratively minimize the loss functions (we describe various entities of the pictorial representation above in Sec. 3).

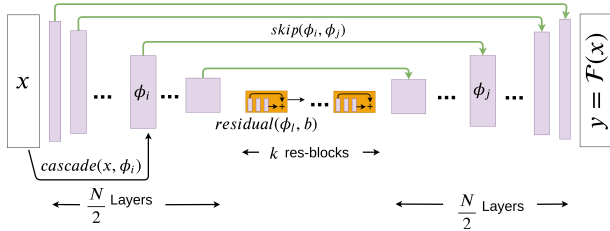


Figure 2: **Network Construction**. The figure illustrate various network components needed for the construction of the network \mathcal{F} , defined in Eq. 1 and Sec. 3.1. ϕ_i denotes the network layers. $\text{skip}(\phi_i, \phi_j)$ denote the link between i^{th} and j^{th} layer. $\text{cascade}(x, \phi_i)$ denotes the cascading of input x at i^{th} layer. Similarly, k residual blocks are shown. The network description is $\mathcal{F} = (\Phi, \mathcal{S}, \mathcal{C}, \mathcal{R})$ (Eq. 1, Eq. 2, and Eq. 3).

2. Related work

Our approach is related to the training data-independent CNN based methods. We bridge the gap between various unsupervised methods proposed to minimize the use of training samples such as deep image prior [24], contextual learning [18], and internal learning [20]. Unlike the classical image prior [1, 10, 12], the deep prior learning [24] shows that hand-crafted structure of the network work as a prior to capture good image statistics for various image restoration tasks. Zero-shot learning uses the internal recurrence of information inside a single image to collect various image specific statistics for super-resolution [21]. InGAN uses multi-scale patch discriminator for learning the patch

distribution from the source image [20]. DCIL gets network structure insight from [17]. The construction of DCIL loss is related to [24, 18, 20].

3. Deep Contextual Internal Learning

DCIL uses image-conditional GAN to map the input image (as opposed to noise) to a different size target image. It performs the image restoration and retargeting without using training data-set. Fig. 1 shows how the major parts of the DCIL framework (highlighted in red) are connected. We describe how these components are related as follows.

Overview. Given a source image x , the objective is to output a target image y from the target distribution Y . The learning procedure is unsupervised and only uses information from the source x . The image in the target domain is of different size than that of y (image resizing). For example, in Denoising-Super-Resolution (DSR), the source image x is a low resolution noisy image and the target domain Y is a set of high-resolution clean image (image restoration).

DCIL constructs generator \mathcal{G} and discriminator \mathcal{D} using the network components described in Sec. 3.1. The network parameters are randomly initialized. The source image x is fed to the generator $\mathcal{G}(\cdot, s_h, s_w)$, where s_h and s_w are the scaling factors for height and width. Next, we iteratively minimize the total loss (Eq. 6) computed between the generator output $y = \mathcal{G}(x, s_h, s_w)$ (*i.e.*, target image) and input source image x . The total loss consists of contextual loss \mathcal{L}_{CL} for contextual feature learning, adversarial loss \mathcal{L}_{GAN} for internal patch distribution learning, and reconstruction loss \mathcal{L}_R for global features learning. We describe these ma-

for parts of the DCIL below.

3.1. Network Construction

We simplify the network construction based on the major components and abstract out network layer specific details. The structure of the network influences the quality of the image features captured. More specifically, the network structure itself works as a prior in the training data-independent methods [24]. Therefore, the purpose of this component is to provide more modularity and a degree of freedom for the DCIL framework.

Consider a network \mathcal{F} which could be a generator or a discriminator. The formal description of network \mathcal{F} is given in Eq. 1.

$$\mathcal{F} = (\Phi, \mathcal{S}, \mathcal{C}, \mathcal{R}) \quad (1)$$

Here, ϕ denotes the set of network layers, \mathcal{S} denotes the configurations of the skip connections, \mathcal{C} denotes the set of layers for which the cascading of the network inputs is performed, and \mathcal{R} denotes the residual blocks. We discuss the major network components given in Eq. 1 as follows.

- **Network Layers (Φ).** The network layers store image representations. Given a network \mathcal{F} with depth N , let $\Phi = \{\phi_l\}_{l=1}^N$ be the set of layers present in the network. Here, ϕ_i could be a convolution layer, an activation layer, or a batch normalization layer.
- **Skip connections (\mathcal{S}).** The skip link between the layers ϕ_i and ϕ_j , where $i < j$, is made by concatenating the output of the layer ϕ_{j-1} with the output of the layer ϕ_i and then feeding into the layer ϕ_j . Let $\text{skip}(\phi_i, \phi_j) = \text{conv}(\text{conv}(\phi_i) \parallel \phi_{j-1})$ denote the skip link between the layers ϕ_i and ϕ_j . The set of skip connections of \mathcal{F} is given in Eq. 2.

$$\mathcal{S} = \{\text{skip}(\phi_i, \phi_j) : \phi_i, \phi_j \in \Phi; i < j\} \quad (2)$$

To simplify the network description, let us denote \mathcal{S} by the set of tuples where each tuple contains the network layer identifier for the skip connection. For example, $\mathcal{S} = \{(1, N), (2, N - 1)\}$ denotes the two skip connections, $\text{skip}(\phi_1, \phi_N)$ and $\text{skip}(\phi_2, \phi_{N-1})$.

- **Cascading of network input (\mathcal{C}).** It is a procedure to successively resize the network input x and then feed it into the intermediate layer ϕ_i of the network. We have not used cascading of network input in DCIL as it was shown to not significantly improve the performance [17]. We denote this by $\{\}$ in the network description for completeness. We have described it more in the supplementary material.
- **Residual Block (\mathcal{R}).** The residual learning framework helps in training higher depth networks while preventing the vanishing gradients problem [13]. It adds the output of two convolution layers b blocks apart. Let

$\text{residual}(\phi_l, b) = \text{add}(\text{conv}(\phi_{l+b}), \phi_l)$ denotes the output residual block $\{\phi_i\}_{i=l+1}^{l+b}$ of length b . The set of residual blocks \mathcal{R} of \mathcal{F} is defined in Eq. 3.

$$\mathcal{R} = \{\text{residual}(\phi_l, b) : \phi_l \in \Phi, b \in [N]\} \quad (3)$$

To simplify the description, let us denote $\mathcal{R} = [k]$, where k is the number of residual blocks.

Generator. Fig. 2 shows an example construction of the generator $\mathcal{G} : X \rightarrow Y$ to show the four-tuple network description. It is an encoder-decoder network which maps the given source image x to the target image $y = \mathcal{G}(x, s_h, s_w)$, $x \sim X$ and $y \sim Y$. The network description of the generator $\mathcal{G} = (\Phi, \mathcal{S}, \mathcal{C}, \mathcal{R})$ based on Eq. 1 is defined in Eq. 4.

$$\mathcal{G} = \left(\{\phi_l\}_{l=1}^N, \{(i, N - i)\}_{i=2}^{\frac{N}{2}-1}, \{\}, [k] \right) \quad (4)$$

Here, $\Phi = \{\phi_l\}_{l=1}^N$ is the set of network layers. We have defined network components \mathcal{S} and \mathcal{R} in Eq. 2 and Eq. 3. $\{\}$ denotes that cascading of the network input is not performed. We use network configurations defined in Eq. 11 and Eq. 13 for our experiments.

Discriminator. It maps the generated image $\mathcal{G}(x) \in Y$ to a patch discriminator $m \in M$, where each entry in m denotes the probability of a patch coming from the patch distribution of the natural image, i.e., $\mathcal{D} : Y \rightarrow M$. We define discriminator as $\mathcal{D}(z) = \sum_{i=1}^4 w_i D^i(z)$. Here, each D^i is a convolution patch discriminator which outputs a map containing the scores of the image patches to be real. And there are four discriminators. The description of discriminator D^i is given in Eq. 5.

$$D^i = \left(\{d_l^i\}_{l=1}^4, \{\}, \{\}, \{\} \right) \quad (5)$$

Here, D^i is a CNN with 4-layers. The empty set $\{\}$ denote the absence of the network component. Therefore, D^i does not have skip links, no residual blocks, and no cascading of network input. The multiscale discriminator \mathcal{D} matches the patch distribution over a range of patch sizes capturing both the fine-grained details as well as the coarse structures in the image [20]¹.

3.2. Loss Function

Given the source image x , the objective is to generate the target image $\mathcal{G}(x) = y$ from the target domain Y . Total loss function \mathcal{L} minimizes the difference in features of the source image and the target image at different feature representations: pixel-to-pixel comparison (reconstruction loss \mathcal{L}_R), context vectors comparison (contextual loss \mathcal{L}_{CL}),

¹In the supplementary material, we provide more details on the generator and the discriminator network.

and patch-based comparison (adversarial loss \mathcal{L}_{GAN}). The total loss function \mathcal{L} is described in Eq. 6.

$$\mathcal{L} = \lambda_C \mathcal{L}_{CL}(\mathcal{G}(x), x) + \lambda_G \mathcal{L}_{GAN}(\mathcal{G}, \mathcal{D}, x, y) + \lambda_R \mathcal{L}_R(\mathcal{G}, \mathcal{D}, x, y) \quad (6)$$

Here, \mathcal{G} and \mathcal{D} are both CNN described in Sec. 3.1. The terms λ_C , λ_G , and λ_R are the coefficients of the loss functions.

The intuition behind loss in Eq. 6 is that minimizing feature differences at different image representations could help in maximizing the image feature learning from the source image. \mathcal{L}_{CL} compares context vectors to make the distribution of the generator output to be contextually similar to the distribution of the natural images [19]. \mathcal{L}_{GAN} is aimed to output distribution of the image patches, which is indistinguishable from the patch distribution of the natural images. \mathcal{L}_R performs the pixel-to-pixel comparisons between the source image and the target image or an inverse mapping of the target image. It ensures that we do not miss any of the object details in the generator output image. We now describe the loss terms used in Eq. 6 for completeness.

Contextual loss (\mathcal{L}_{CL}). It is used to enhance the contextual features in the reconstruction. The set of context vectors are obtained by feeding image x and y into pre-trained VGG19 ϕ [23]. In Fig. 1, we have pictorially shown context vectors as the output of VGG19. Let $\phi^l(x)$ and $\phi^l(y)$ denote the feature extracted from layer l of the network ϕ . The contextual loss is defined in Eq. 7.

$$\mathcal{L}_{CL}(x, y, l) = -\log CX(\phi^l(x), \phi^l(y)) \quad (7)$$

Here, CX denotes the contextual similarity measure. It is computed by considering cosine distance between the context vectors extracted from the network ϕ [19]. Eq. 7 minimizes dissimilarities between the contextual feature computed from the source image x and the target image y . CX is normalized and lies in the range $[0, 1]$.

Adversarial loss (\mathcal{L}_{GAN}). The purpose of adversarial learning is to synthesize new image features in the output image from the patch distribution of the natural images. It is a sum of the generator loss \mathcal{L}_G and the discriminator loss \mathcal{L}_D . The generator \mathcal{G} and the discriminator \mathcal{D} are both CNN. The generator loss \mathcal{L}_G and the discriminator loss \mathcal{L}_D are used for distribution matching. \mathcal{G} generates the desired image. \mathcal{D} tries to distinguish the output of \mathcal{G} and the source image. Therefore, the generator learns the patch distribution through the interaction with the discriminator. We show the adversarial loss in Eq. 8.

$$\mathcal{L}_{GAN}(\mathcal{G}, \mathcal{D}, x, y) = \mathcal{L}_G(x) + \mathcal{L}_D(x, y) \quad (8)$$

Here, \mathcal{G} outputs the target image $\mathcal{G}(x) = y$. The feature learning in the adversarial framework could suffer

	CL [18]	DIP [24]	DCIL (ours)
BSD100	0.60	0.62	0.63
SET14	0.62	0.66	0.67
SET5	0.64	0.66	0.66

Table 1: **2×Denoising-Super-Resolution.** Performance comparison (SSIM) for 2×SR where low resolution image contains noise with strength $\sigma = 100$.

from mode collapse [3, 4]. The use of multi-scale discriminator prevents it by maximizing feature learning by comparing the reconstruction at multiple scales [20]. We have described the multi-scale discriminator in Sec. 3.1.

Reconstruction loss (\mathcal{L}_R). It is used to maximize the likelihood of randomly initialized network weights. One could define a spatial correspondence in the case of image restoration [24]. However for image retargeting, reconstruction loss in a cycle consistent approach performs well as the generator output does not have spatial correspondence with the source image [20]. Therefore, the two different ways of computation of reconstruction loss are as follows.

\mathcal{L}_R for image restoration is computed between the generator output $\mathcal{G}(x)$ and the source image x , as in Eq. 9.

$$\mathcal{L}_R(\mathcal{G}, s_h, s_w, x, y) = \|\mathcal{G}(x, s_h, s_w) - x\| \quad (9)$$

\mathcal{L}_R for image retargeting is computed between the source image x and the inverse mapping of the generator output $\mathcal{G}(y)$, where $y = \mathcal{G}(x, s_h, s_w)$, as in Eq. 10.

$$\mathcal{L}_R(\mathcal{G}, s_h, s_w, x, y) = \|\mathcal{G}(y, \frac{1}{s_h}, \frac{1}{s_w}) - x\| \quad (10)$$

4. Applications

In this section, we describe image resizing in two different setups. The first setup is image restoration problems. There are two ways for it. DSR where the low-resolution input contains noise. SR where the low-resolution input does not contain noise. The second setup for image resizing is content-aware image retargeting. We describe these applications below.

4.1. Denoising-Super-Resolution.

DSR makes the image resize operation challenging as one has to perform two tasks - image denoising and image super-resolution. The description of the generator network for DSR is given in Eq. 11.

$$\mathcal{G}_1 = \left(\{G_l\}_{l=1}^{10}, \{(i, 10-i)\}_{i=2}^4, \{\}, \{\} \right) \quad (11)$$

Here, $\{G_l\}_{l=1}^{10}$ is the depth-5 encoder-decoder network where $\{G_l\}_{l=1}^5$ are the layers of encoder and $\{G_l\}_{l=6}^{10}$ are

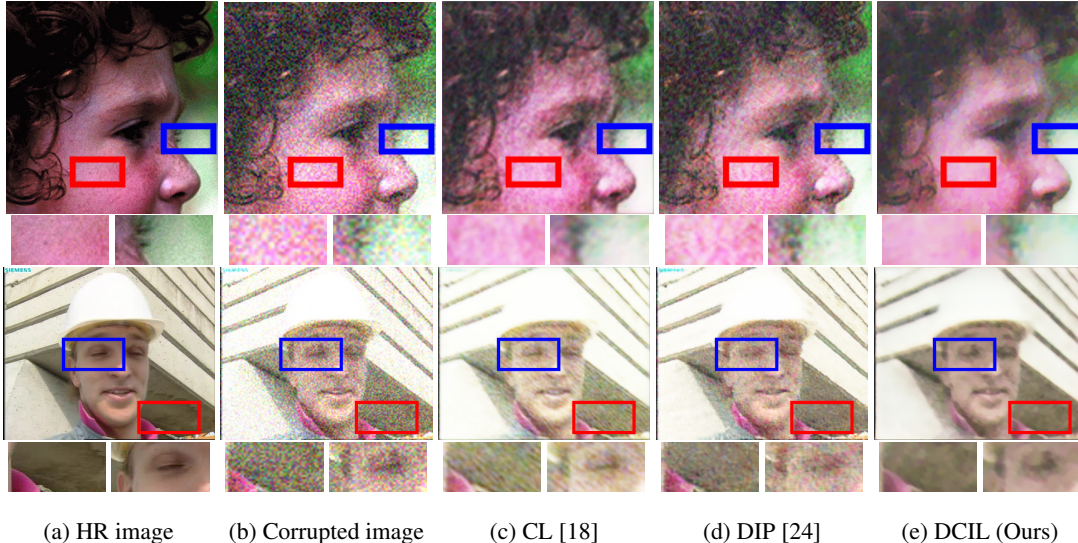


Figure 3: **2×Denoising-Super-Resolution.** The corrupted low resolution images contain noise with strength $\sigma = 100$. CL [18] and DIP [24] create noisy spots in the image restoration output. DCIL (ours) output clean images compared to CL [18], DIP [24] (see the cropped images below the figures).

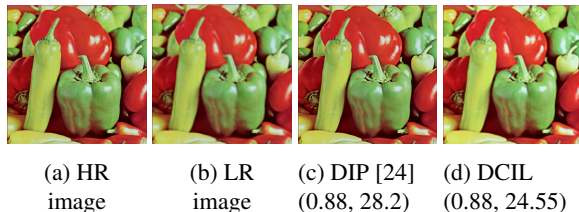


Figure 4: **4× SR comparison** for (SSIM, PSNR) values. DCIL output image is comparable to DIP [24] (the images are best viewed after zooming). It could be observed that a higher PSNR value does not imply a higher perceptual quality [17].

the layers of decoder. There are skip connections from encoder layers to the decoder layer in \mathcal{G}_1 . The cascading of the network input is not performed and there are no residual blocks. Encoder-decoder architecture countermeasures the mode collapse and improves stability [20].

Given low-resolution noisy image \hat{I} , the loss function for DSR is defined in Eq. 12.

$$\mathcal{L}_1 = \lambda_C \mathcal{L}_{CL}(\mathcal{G}_1(x), y) + \lambda_{G_1} \mathcal{L}_{GAN}(\mathcal{G}_1, \mathcal{D}, \hat{I}, x) + \lambda_{\mathcal{R}} \|\mathcal{G}_1(x) - U_t(\hat{I})\| + \lambda_{TV} \|TV(\mathcal{G}_1(x))\| \quad (12)$$

Here, \mathcal{D} is similar to the one defined in Eq. 5. $U_t(\cdot)$ is the up-sampling operator with the scaling factor as t . λ_{TV} is the coefficient of the Total variation (TV) regularization. TV norm in Eq. 12 reduces the noise from the corrupted image². We have discussed the loss terms of Eq.12 in Sec. 3.2.

²Total variation is a sum of the absolute differences of neighboring pixel values in the input image. It measures the noise in the image.

The adversarial loss \mathcal{L}_{GAN} uses the multi-scale patch discriminator to learn the image features at different resolutions. Intuitively, it utilizes the patch replication across multiple scales to augment feature learning. The contextual loss \mathcal{L}_{CL} improves feature learning at the scale of the target image using context vectors. The reconstruction loss $\|\mathcal{G}_1(x) - U_t(\hat{I})\|$ provides the global features in the resulting output.

In Table 1, we give quantitative comparisons for DSR. The aim is to perform 2×SR with denoising, where noise strength is $\sigma = 100$. The visual comparison for the generated images is provided in Fig. 3. One could observe that we outperform the state-of-the-art methods which we compare with³.

4.2. Super-Resolution.

CNN based SR has been studied in two ways. First, we can use pixel-to-pixel loss, which leads to high PSNR at the price of low perceptual quality [29, 7]. Second, we can use feature space loss or an adversarial loss to achieve higher perceptual quality [16, 15]. CL combines the two training data based approaches above to generate natural-looking images, with good structural similarity [18].

The generator and the discriminator for SR are similar to the ones used in DSR. The loss function for SR is similar to DSR given in Eq. 12 but without TV norm as there is no noise in the input images.

³We use original implementation of contextual loss (github.com/roimehrez/contextualLoss) and DIP (github.com/DmitryUlyanov/deep-image-prior). We generated DIP output images using the default hyper-parameters [24].

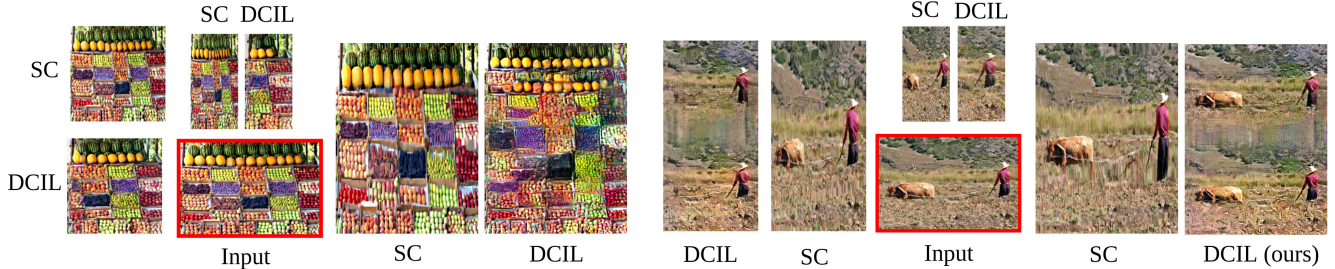


Figure 5: **Image Retargeting (Object Structure)**. The size of the local objects (e.g., fruits and man) confirms the preservation of object structure in the image retargeting output. SC [5] does not preserve the structure of the objects (e.g., the man in the 8th column is deformed). DCIL (ours) preserve the structure of the objects by adding new objects or removing objects.

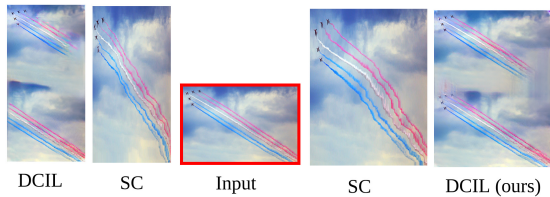


Figure 6: **Image Retargeting (Object Alignment)**. The line-shaped clouds (contrails) produced by aircraft confirms the preservation of the object alignment in the image retargeting output. SC [5] does not preserve the alignment of the objects. DCIL (ours) preserves the alignment of the contrails when increasing height and when increasing width.

We perform $4\times$ SR on BSD100 data set. Fig. 4 shows the perceptual quality comparison for $4\times$ SR. The average SSIM scores on BSD100 dataset are as follows. Mechrez et al. [18]: 0.64, ZSSR [21]: 0.72, DIP [24]: 0.79, and our DCIL: 0.76. We found that MSE based methods capture strong prior for SR compared to the contextual loss and adversarial loss-based frameworks, which is counter-intuitive [24]. We confirm this by the following results (our run)⁴. For Set-5, the score are, Mechrez et al. [18]: 0.86, ZSSR [21]: 0.87, DIP [24]: 0.90, and our DCIL: 0.87. For Set-14, Mechrez et al. [18]: 0.78, ZSSR [21]: 0.76, DIP [24]: 0.81, and our DCIL: 0.79. Our interpretation of this phenomenon is as follows. Pixel-to-pixel comparison is converging to better optima in SR. However, in the case of DSR, a pixel-to-pixel comparison could be over-learning noise with features (Table 1). We believe that the performance of DIP and DCIL could probably be further improvised using hyper-parameter search.

4.3. Image Retargeting.

It is a content-aware image resizing operation which aims to output image with a different size, smaller or larger,

⁴For BSD100, we have used SSIM values provided in [18]. For Set-5 and Set-14 dataset, we have used unsupervised implementation of [18] for a fair comparison. SSIM values for [24] is computed by our run.

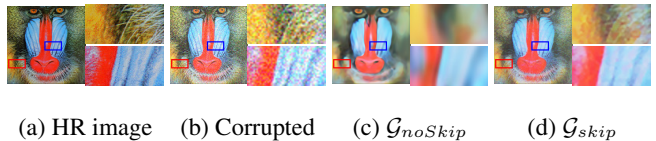


Figure 7: **Ablation Study**. $2\times$ Denoising-Super-Resolution with noise value $\sigma = 100$. The network with skip connections \mathcal{G}_{skip} performed better than the network without skip connections \mathcal{G}_{noSkip} . The above experiment study skip connections for multiple corruptions, unlike [24, 17] (the images are best viewed after zooming).



Figure 8: **Failure Example**. The aim is to preserve the object context when performing image retargeting. SC [5] deforms the object (i.e., man). InGAN does not partition the object well [20]. DCIL partition the object, but the image feature of the elbow is not well-formed (the images are best viewed after zooming).

and with a different aspect ratio. Image retargeting is performed in various ways. There are methods which are aimed to preserve only the salient objects and discarding/extending the object background (e.g., [8, 26]). Other methods (e.g., [20, 22]) including our DCIL preserve the local sizes/aspect-ratios of the local objects while resizing the image. The replication/reduction of the objects is desired to fill the scene with similar image features.

Suppose input image x is of size (h, w) . The scaling factors s_h and s_w are used as the input for retargeting. The retargeting objective is to output image y with the size $(s_h h, s_w w)$. The description of the generator network for image retargeting is given in Eq. 13.

$$\mathcal{G}_3 = \left(\{G_l\}_{l=1}^{10}, \{(i, 10 - i)\}_{i=2}^4, \{\}, \{6\} \right) \quad (13)$$

Here, $\{G_l\}_{l=1}^{10}$ are the layers of the depth-5 encoder-decoder network. The network is equipped with skip connections from the layers of the encoder to the layers of the decoder. There are six residual blocks. The discriminator network is similar to the one defined in the Eq. 5.

The loss function for image retargeting is given in Eq. 14.

$$\mathcal{L}_3 = \lambda_C \mathcal{L}_{CL}(\mathcal{G}_3(x), y) + \lambda_{GAN} \mathcal{L}_{GAN}(\mathcal{G}_3, \mathcal{D}, x, y) + \lambda_{\mathcal{R}} \|\mathcal{G}_3(\mathcal{G}_3(x)) - y\| \quad (14)$$

Here, λ_C , λ_{GAN} , and $\lambda_{\mathcal{R}}$ are the scaling factors. The adversarial loss \mathcal{L}_{GAN} and the contextual loss \mathcal{L}_{CL} both matches the distribution of image patch of the source image and the target images. Distribution matching is the essential requirement for image retargeting [20]. Also, they both work for the non-aligned image data of the source and the target images, unlike DIP [24].

We compute an automorphism as the cycle consistency check to preserve all the object details in the synthesized output [20]. The automorphism retargets the generator output back to its source domain. Then we could perform the pixel comparison using reconstruction loss. It preserves the global image features in the retargeted image.

Our DCIL uses contextual learning to preserve the object features and object alignment in the image retargeting output better than Seam-Carving (SC) [5] as shown in Fig. 5 and Fig. 6. DCIL maximizes the feature learning and performs comparably to InGAN (Fig. 8).

	ZSSR [21]	CL [18]	DIP [24]	InGAN [20]	DCIL (ours)
DSR	✘	✓	✓	✘	✓
SR	✓	✓	✓	✓	✓
Retargeting	✘	✘	✘	✓	✓

Table 2: The table shows the comparison between various frameworks. DCIL (ours) is a generalized framework which performs all the tasks and generates images comparable to the other methods. We provide the extended version of Table 2 and the implementation details of DCIL in the supplementary material.

User study. We conducted a user study to evaluate the image retargeting results. We collected feedback from 58 human experts with a total of 290 votes. Each subject is asked to vote the perceptually better images constrained to the preservation of object properties. The percentage of the votes for SC [5] is 35%. Our DCIL got 65% votes. The user study shows that DCIL performs good image retargeting.

5. Ablation Study and Limitations

The limitations of the DCIL framework are due to the lack of contextual understanding by feature learning from a single image. Fig. 7 shows that the network with skip connection outperforms the network without skip connections.



Figure 9: **Image Inpainting.** This shows the results for region inpainting (the images are best viewed after zooming).

Therefore, one needs to carefully design an application-specific network to maximize feature learning [17]. Fig. 8 shows that contextual feature learning of DCIL leverages adversarial learning of InGAN for object partitioning limitations in image retargeting. However, the perceptual quality in the presence of object replications could be further improvised.

6. Discussion

DCIL is completely unsupervised and does not use training samples. It is different than the supervised methods RCAN [30] and DRLN [2], which use training data to perform image restoration. DCIL exploits the inherent self-similarity present in the source image. Ulyanov *et al.* have shown that self-similarity prior emerged because of the convolutional operations tend to impose self-similarity in the generated images [24]. DCIL incorporates image prior using the network structure implicitly. Similar to the DSR, SR, and image retargeting, it could also perform image inpainting (Fig. 9). It is due to the self-similarity prior captured by DCIL helps to perform inpainting task. The quality of the deep prior for the various tasks depends upon the learning procedure. The network initially learns the image feature, but then it tends to over-learn the noise from the corrupted input [17]. The learning procedure is generally more tricky when we perform distribution matching using GAN and CL. However, an exhaustive hyper-parameter search helped us in the above scenario.

7. Conclusion

DCIL fits a randomly-initialized untrained generator. The structure of the network and the loss function are the main tools for unsupervised approaches described in the paper. We performed image resizing in many challenging scenarios. The performance depends upon the high correlation between the features of the source and the target images. For example, in the presence of high corruption due to noise in the source image, the performance of various methods degrade. We believe that it would be interesting to investigate the image statistics captured by DCIL for the other single image applications, *e.g.*, image inpainting.

References

- [1] M. Aharon, M. Elad, A. Bruckstein, et al. K-svd: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Transactions on signal processing*, 54(11):4311, 2006.
- [2] S. Anwar and N. Barnes. Densely residual laplacian super-resolution. *arXiv preprint arXiv:1906.12021*, 2019.
- [3] M. Arjovsky and L. Bottou. Towards principled methods for training generative adversarial networks. In *5th International Conference on Learning Representations, ICLR 2017*, 2017.
- [4] S. Arora, R. Ge, Y. Liang, T. Ma, and Y. Zhang. Generalization and equilibrium in generative adversarial nets (gans). In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 224–232. JMLR.org, 2017.
- [5] S. Avidan and A. Shamir. Seam carving for content-aware image resizing. In *ACM Transactions on graphics (TOG)*, volume 26, page 10. ACM, 2007.
- [6] S. A. Bigdeli and M. Zwicker. Image restoration using autoencoding priors. *arXiv preprint arXiv:1703.09964*, 2017.
- [7] Y. Blau and T. Michaeli. The perception-distortion tradeoff. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6228–6237, 2018.
- [8] D. Cho, J. Park, T.-H. Oh, Y.-W. Tai, and I. So Kweon. Weakly-and self-supervised learning for content-aware deep image retargeting. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4558–4567, 2017.
- [9] A. Creswell, T. White, V. Dumoulin, K. Arulkumaran, B. Sengupta, and A. A. Bharath. Generative adversarial networks: An overview. *IEEE Signal Processing Magazine*, 35(1):53–65, 2018.
- [10] W. Dong, L. Zhang, G. Shi, and X. Li. Nonlocally centralized sparse representation for image restoration. *IEEE transactions on Image Processing*, 22(4):1620–1630, 2012.
- [11] Y. Gandelsman, A. Shocher, and M. Irani. Double-dip: Un-supervised image decomposition via coupled deep-image-priors. *arXiv preprint arXiv:1812.00467*, 2018.
- [12] S. Gu, Q. Xie, D. Meng, W. Zuo, X. Feng, and L. Zhang. Weighted nuclear norm minimization and its applications to low level vision. *International journal of computer vision*, 121(2):183–208, 2017.
- [13] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [14] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [15] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*, pages 694–711. Springer, 2016.
- [16] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi. Photo-realistic single image super-resolution using a generative adversarial network. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [17] I. D. Mastan and S. Raman. Multi-level encoder-decoder architectures for image restoration. *arXiv preprint arXiv:1905.00322*, 2019.
- [18] R. Mechrez, I. Talmi, F. Shama, and L. Zelnik-Manor. Learning to maintain natural image statistics. *arXiv preprint arXiv:1803.04626*, 2018.
- [19] R. Mechrez, I. Talmi, and L. Zelnik-Manor. The contextual loss for image transformation with non-aligned data. *European Conference on Computer Vision (ECCV)*, 2018.
- [20] A. Shocher, S. Bagon, P. Isola, and M. Irani. Internal distribution matching for natural image retargeting. *IEEE international conference on computer vision*, 2019.
- [21] A. Shocher, N. Cohen, and M. Irani. zero-shot super-resolution using deep internal learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3118–3126, 2018.
- [22] D. Simakov, Y. Caspi, E. Shechtman, and M. Irani. Summarizing visual data using bidirectional similarity. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2008.
- [23] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [24] D. Ulyanov, A. Vedaldi, and V. Lempitsky. Deep image prior. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [25] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, and B. Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [26] L. Wolf, M. Guttman, and D. Cohen-Or. Non-homogeneous content-driven video-retargeting. In *2007 IEEE 11th International Conference on Computer Vision*, pages 1–6. IEEE, 2007.
- [27] C. Yang, X. Lu, Z. Lin, E. Shechtman, O. Wang, and H. Li. High-resolution image inpainting using multi-scale neural patch synthesis. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, 2017.
- [28] K. Zhang, W. Zuo, S. Gu, and L. Zhang. Learning deep cnn denoiser prior for image restoration. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [29] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 586–595, 2018.
- [30] Y. Zhang, K. Li, K. Li, L. Wang, B. Zhong, and Y. Fu. Image super-resolution using very deep residual channel attention networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 286–301, 2018.
- [31] Y. Zhang, Y. Zhang, and W. Cai. Separating style and content for generalized style transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, volume 1, 2018.

- [32] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.