

# Super-Resolution for Overhead Imagery Using DenseNets and Adversarial Learning

Marc Bosch  
The Johns Hopkins University  
Applied Physics Laboratory  
marc.bosch.ruiz@jhuapl.edu

Christopher M. Gifford  
The Johns Hopkins University  
Applied Physics Laboratory  
Christopher.Gifford@jhuapl.edu

Pedro A. Rodriguez  
The Johns Hopkins University  
Applied Physics Laboratory  
Pedro.Rodriguez@jhuapl.edu

## Abstract

Recent advances in Generative Adversarial Learning allow for new modalities of image super-resolution by learning low to high resolution mappings. In this paper we present our work using Generative Adversarial Networks (GANs) with applications to overhead and satellite imagery. We have experimented with several state-of-the-art architectures. We propose a GAN-based architecture using densely connected convolutional neural networks (DenseNets) to be able to super-resolve overhead imagery with a factor of up to  $8\times$ . We have also investigated resolution limits of these networks. We report results on several publicly available datasets, including SpaceNet data and IARPA Multi-View Stereo Challenge, and compare performance with other state-of-the-art architectures.

## 1. Introduction

Super-resolution is the task of estimating plausible pixel information given an image and creating a corresponding higher resolution version. Typically, super-resolution methods aim to recover high frequency components of the scene lost in the image acquisition process for an improved perceived quality. The majority of approaches attempt to either produce new pixel values by estimating them from a support region (neighborhood) or through a learned model given many examples of low resolution to high resolution mappings. The latter has been an active area of research for many years within the computer vision and image processing communities [8, 14, 23, 19, 6]. Some of these works have proven successful for recovering detail from low resolution images typically acquired with consumer electronic

cameras. However, the low resolution input of many available solutions still offer enough detail to infer most of the semantics of the scene.

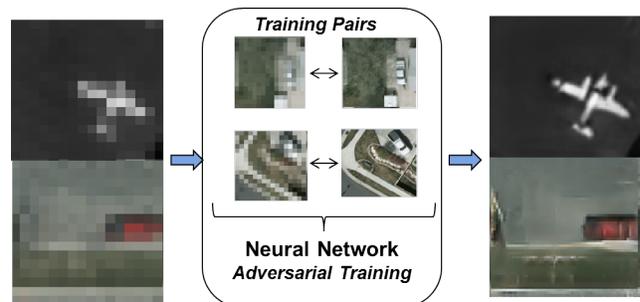


Figure 1. Super-resolution applied to overhead imagery using our system.

In overhead imaging (*i.e.* airborne and satellite), super-resolution has the potential to offer advanced automatic target recognition (ATR) capabilities and improved human exploitation value. Due to the distance from source to target, even a small resolution gain can dramatically improve the end use. Super-resolving an image by an upscale factor of 4 means that from one pixel we need to estimate fifteen new pixels in a  $4 \times 4$  neighborhood. In other words, if the original image resolution represents a ground sampling distance (GSD) of four meters per pixel, the new super-resolved image would be resolved up to one meter per pixel. In this context, many common objects would not have enough detail for analysts to “understand” the scene. Hence, it becomes more a problem of enhancing semantics, objects, and fine-grained features rather than improve image quality.

Super-resolution of overhead imagery can have a significant impact for the space industry. There are many potential

remote sensing applications that would directly benefit from this technology, such as crops and deforestation monitoring, economic activity tracking, space imaging, and various reconnaissance activities. Also, different imagery modalities such as panchromatic electro-optical (EO), hyperspectral (HSI), infrared (IR), and synthetic aperture radar (SAR) can benefit from these advances.

With the increasing attention being given to neural network models, recent works have proposed the use of several network architectures to tackle the super-resolution problem, including [6, 17, 4]. Generative models have been proposed to infer and recover plausible details from low resolution images. Generative Adversarial Networks (GANs) are one of the most popular generative Deep Learning framework for super-resolution. GANs are trained through adversarial training where two engines, namely *generator* and *discriminator*, participate in a two-player game with the goal of improving as the opponent improves. In [17], authors proposed a GAN-based algorithm using neural networks with deep convolution layers to model the low-to-high resolution mappings.

In this work, we have investigated multiple convolutional neural network (CNN) configurations. Further, we propose the integration of a densely connected convolutional network architecture (DenseNet) within the *generator* of a GAN. DenseNets [12] have the particularity of connecting layers in a dense manner to other layers for better feature representation and computational efficiency. Figure 1 shows an example of inputs/outputs of our framework applied to satellite images.

## 1.1. Background

Work in super-resolution has been ongoing for decades now. There have been two main approaches: multi-frame and single frame super-resolution. Among single frame/image super-resolution, we can find baseline algorithms like bicubic interpolation, cubic splines [14], and other local-based approaches that consider local regions of support to estimate new pixels. Some techniques rely on statistical priors from the image [8]. More recently, learning-based approaches have successfully been used to produce improved results. These techniques are based on the self-similarity principle, where it is assumed that many images share similar visual properties at scale. These approaches back-project high-frequency components lost in the low resolution image from similar patches found in other images given a large collection of reference images [19, 7, 23, 24].

As in many other applications in computer vision, Deep Learning techniques have prevailed in the last few years. There have been several proposed systems where the non-linear mapping between low and high resolution images is learned end-to-end using neural networks [6, 17, 4, 16, 22].

Dong et al. proposed the use of Deep CNNs (SRCNN) to learn the high-frequency representation given low resolutions similar to sparse-coding-based techniques with the added advantage of the joint optimization that occurs in a Deep Learning system [6].

In [15], authors proposed a very deep network that exploits the advantages of a recursive architecture like adding convolutions without adding new parameters for improved performance, and also adding skip connections to reduce the effect of the vanishing gradient found in the recursive schemes.

Dahl and colleagues proposed a method that produced very impressive results for the task of super-resolving human faces [4]. Their method is based on a network that operates in a recursive fashion. Each pixel is super-resolved using previously visited pixels. They combine two CNNs, one based on PixelCNN that makes predictions using previous stochastic estimates, and another that behaves as the conditioning network by receiving the low resolution image and generating logits that encode the log-probability of each pixel in the high resolution image.

In [2], authors proposed a method that uses CNNs to characterize high-frequency content of images not present in the source image. They model the conditional mapping of a high-resolution image given its low-resolution version as a Gibbs distribution. Following this work, Ledig et al., in [17], proposed the use of adversarial training and GANs to tackle the recovery of high-level details of an image. They designed a framework that uses residual blocks [11] to generate features that encode plausible missing information. The weights of the network are updated so that the adversarial loss combined with feature matching loss is minimized. They add a feature matching loss to quantify the fidelity of the reconstructed image in terms of perceptual loss by a VGG-19 universal feature extractor.

## 1.2. Contributions

In this paper, we describe the integration of densely connected convolutional networks to a GAN framework for the task of super-resolution for overhead imagery. Our main contributions can be summarized as follows:

- Investigated application of existing state-of-the-art super-resolution models to overhead imagery to determine what is possible with today's models for both panchromatic electro-optical (EO) and multi-band images.
- Proposed a network architecture based on dense blocks repurposed for super-resolution applications under an adversarial training framework.
- Evaluated several loss functions, including feature matching criteria, with the purpose of understand-

ing how transferable pre-trained models are on non-overhead natural images with much different scale and geometry viewpoints with respect to satellite imagery.

- Evaluated several super-resolution gain factors to understand the limits of these techniques, as well as constrain the scope of the problem by conducting experiments on specific semantic categories. We have used imagery with a large diversity of geometric and semantic features to gain an understanding on the limits and capabilities of this technology.

The remainder of the paper is organized as follows: In Section 2 we review the GAN model, in Section 3 we present our proposed method using a DenseNet, in Section 4 we show our experimental results on overhead imagery, and finally, we conclude the paper by offering some remarks from our experimental observations.

## 2. Generative Adversarial Network (GAN) Models

Generative Adversarial Networks (GANs) are a particular case of generative models. Their goal is to learn the probability distribution of the source data using adversarial training. This means applying alternating stochastic gradient updates in a two-player, zero-loss game. In the task of super-resolution, a generative model is required to input new details into the low resolution input. Two engines, *generator* and *discriminator*, are established in a min-max optimization framework. The *generator* aims at generating new samples, *fake data*, from the learned data distribution that look as real as possible. The *discriminator's* goal is to detect such *fake data* among a collection of real training data. At training time, both *generator* and *discriminator* take turns to fool the opponent and to distinguish the fake from real data respectively. As a result, training is successful if the *generator* becomes increasingly better at creating realistic data, and the *discriminator* improves spoofing detection capabilities, which forces the *generator* to become even better in its attempt to blur the line between *fake* and *real* data.

More formally, the adversarial training process in a GAN is described by the following adversarial cost/loss function:

$$Loss_{G,D} = \min_{\psi_G} \max_{\psi_D} (E_{x \sim p_{target}(x)} (\log D_{\psi_D}(x)) + E_{z \sim p_{model}(z)} (\log(1 - D_{\psi_D}(G_{\psi_G}(z)))) \quad (1)$$

with  $x$  being samples from the target distribution (i.e., high-resolution training images), and  $z$  representing the input variables needed in the model estimate (i.e., low-resolution training images) to generate new high-resolution images that approximate to the target high res images  $x$ . During

training, the *discriminator* maximizes the following expression given a batch of generated data from the *generator*,  $G_0(z)$ :

$$\psi_D^{(*)} = \max_{\psi_D} (E_{x \sim p_{target}(x)} (\log D_{\psi_D}(x)) + E_{z \sim p_{model}(z)} (\log(1 - D_{\psi_D}(G_{\psi_{G_0}}(z)))) \quad (2)$$

Next, the *generator* objective is to minimize the following:

$$\psi_G^{(*)} = \arg \min_{\psi_G} E_{z \sim p_{model}(z)} (\log(1 - D_{\psi_D^{(*)}}(G_{\psi_G}(z)))) \quad (3)$$

**Deep Convolutional GAN** Deep Convolutional Generative Adversarial Networks (DCGANs) are a class of GANs implemented, in general, with several convolutional layers. Both the generator and discriminator can be implemented using CNNs and applied to several tasks, such as image generation (generator) or image classification (discriminator) [10, 5]. Two CNN architecture types are found in most of the work using DC-GANs, these are encoder-decoder type architecture and ResNet [11]. Note that both do not necessarily need to be complementary.

In the encoder-decoder approach, the input is passed through a series of layers that downsample the feature maps to be able to represent larger receptive fields. A decoder brings back the information to the full input resolution.

ResNet is made of residual blocks that, at each layer, learn a “residual” mapping,  $R(x)$ , obtained from subtracting the underlying input ( $x$ ) to output ( $y$ ) mapping,  $H(x)$ , from the layer input,  $y = H(x) = R(x) + x$ . These residual mappings are learned and added to the identity mapping (shortcut). ResNet has been proven to outperform other state-of-the-art approaches in many tasks, including object recognition. Shortcut connections allow the network to skip certain layers to avoid vanishing gradients and curse of dimensionality problems.

## 3. Dense Network GANs for Super-Resolution

In [12], the authors introduced a new type of CNN architecture, DenseNet, that allows for a richer description of the visual elements in the scene compared to state-of-the-art networks like ResNet. In short, a dense network is a collection of blocks where each of the block's layers is connected to one another; hence, they are referred to as *dense* blocks [12].

Dense blocks have the advantage over residual blocks in that there is a stronger gradient flow due to direct connection at any layer from all other layers in the network. It also maintains and exposes low complexity features at deeper layers in the network better than ResNet does. These much richer representations of visual attributes enable better parameter and computational efficiency, which better po-

sitions DenseNet for onboard processing in airborne and spaceborne platforms.

### 3.1. DenseNet Generator

Our implementation of a DenseNet is limited to the generator of a GAN. It consists of a dense block with BN-ReLU-Conv ( $3 \times 3$ ) basic unit sequence, preceded by a bottleneck block to reduce the number of input feature maps composed by the BN-ReLU-Conv ( $1 \times 1$ ) layer. In addition, each dense block is separated from its neighboring (dense) block by a transition layer. In the original dense network, the transition layer is made of a convolution layer and a pooling layer. In our super-resolution framework, we modify this layer to account for the resolution increase goal. Rather than performing all the convolutions in the low resolution domain as other approaches do, the high efficiency dense blocks allow us to learn the mappings at several scales. Our transition layer doubles the resolution of the feature maps each time. This in turn introduces smaller receptive fields that learn local information, and, thus, recovers fine-grain details. One obvious benefit of this cascaded process is that the network is trained at several resolution gains, making it a more versatile design.

As shown in [20], a way to connect coarse outputs to dense pixels is using backward convolutions, also known as deconvolutions or transpose convolutions. We use deconvolution with a *stride* = 2 to generate an upscale version ( $2 \times$ ) of the input feature map at each transition layer.

Following [21, 9] implementations, the last stage of the generator is a fully-convolutional network. Figure 2 shows our proposed generator architecture for super-resolution.

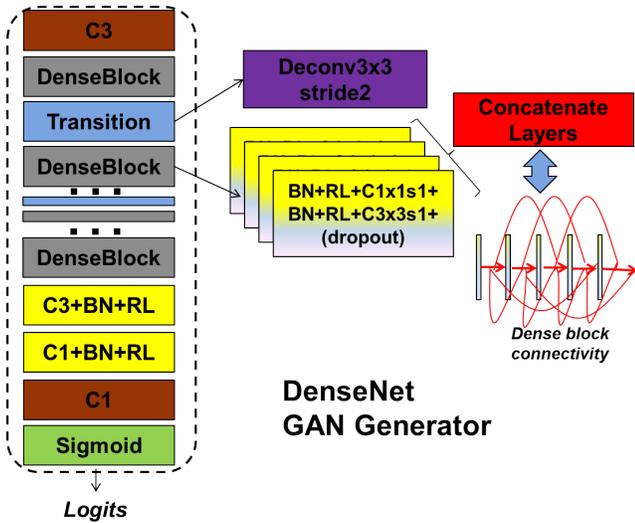


Figure 2. Generator architecture for overhead imagery.

### 3.2. Discriminator

Following other discriminatory models, our discriminator has a relatively shallow configuration with a series of convolutional layers followed by ReLUs and batch normalization blocks. At each layer, the convolutional layer is doubled in a similar fashion as the VGG network [17]. Again, a fully-convolutional subnet is placed in the final layers [21, 9], followed by sigmoid functions that output the decision of real vs. fake input data. Figure 3 shows the architecture of the discriminator used to train the generator.

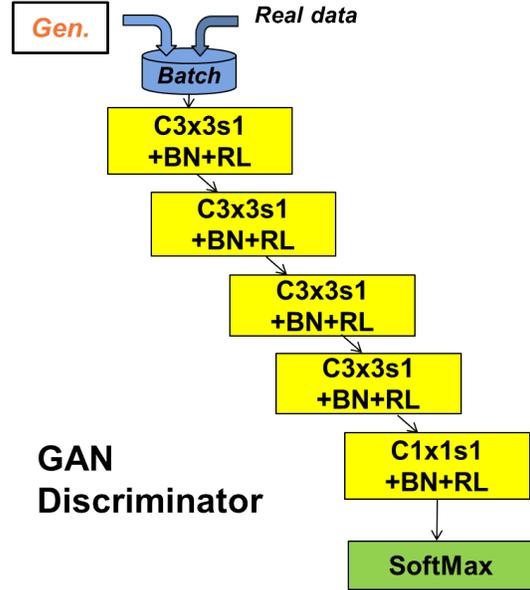


Figure 3. Discriminator architecture for training the generator.

### 3.3. Loss Function

In this work, we have investigated several generator loss types to guide the network to convergence. In general form, the loss can be expressed as a function of the adversarial loss, L1-norm image content loss, and feature matching loss.

$$Loss_{Gen} = \alpha(n)Loss_{adv} + (1 - \alpha(n))((1 - \beta_1) \cdot Loss_{content} + \beta_1 \cdot Loss_{fm}) \quad (4)$$

The  $Loss_{adv}$  is the “vanilla” GAN adversarial loss:

$$Loss_{adv} = E_{z \sim p_{model}(z)}(\log(1 - D_{\psi_D}(G_{\psi_G}(z)))) \quad (5)$$

$Loss_{content}$  represents the L1-norm between the target high-resolution and the generated high-resolution images:

$$Loss_{content} = \|targetHR - gen.HR\|_1 \quad (6)$$

Similarly, the feature matching loss ( $Loss_{fm}$ ) has the goal of describing the visual attribute loss. We use a version of the popular VGG-16 using pre-trained weights on ImageNet as a universal feature extractor. It is computed as:

$$Loss_{fm} = \|T_{\theta_{vgg16}}(targetHR) - T_{\theta_{vgg16}}(gen.HR)\|_1 \quad (7)$$

The parameter  $\alpha$  dynamically modulates the influence of the adversarial loss into the overall loss. Similarly,  $\beta_1$  controls the importance of the feature matching loss. The discriminator only uses the adversarial loss, eq. 1. It is simply a binary cross entropy loss.

### 3.4. Implementation and Training Details

The generator of the proposed GAN using dense blocks has 5 dense blocks per layer. Each block consists of a bottleneck with a batch normalization block, a ReLU and a convolutional layer with  $1 \times 1$  map size (filter spatial support) and  $stride = 1$ , followed by a combination of batch normalization, ReLU and convolutional layer, with a  $3 \times 3$  map size and  $stride = 1$ . At each of these units, we add 16 feature maps. This is also referred to as the growth rate of the network. The growth rate is kept small to avoid the network from growing too wide and to improve parameter efficiency. The transition layer consists of a deconvolution unit with a  $3 \times 3$  and  $stride = 2$  to account for the upscale.

As mentioned earlier in Section 3.1, the final layers of the generator are an implementation of the fully-convolutional network [21], in particular a layer with convolutional filters with a mapsize of  $3 \times 3$  with  $stride = 1$  followed by batch normalization, and ReLUs. The block is repeated once more but with a bottleneck convolution stage of  $1 \times 1$  size. In this final stage, the size of the tensor is not changed as no new feature maps are added.

The discriminator is a four layer network with [64, 128, 256, 512] feature maps for the respective layers. Each filter in each of the layers has  $3 \times 3$  filter size and  $stride = 2$ . The fully-convolutional network from the generator is mimicked after the four layer stages.

In our reported results, our settings are as follows: Batch size is 16. We initialize the parameter controlling the contribution of each loss type,  $\alpha$ , to 0.95 to guide the network towards perceptual convergence at the initial stages, and decrease at each epoch by a factor of 1.05. Finally, we have evaluated the network for several settings of  $\beta$  to investigate the contribution of each term of the perceptual loss, namely feature matching loss and the pixel content loss. Our goal was to evaluate the correlation of networks trained to extract features trained on non-overhead non-synthetic images like VGG-16 with overhead imagery with such a different viewing geometry.

## 4. Evaluation

In this section, we describe our experiments and present results for image super-resolution for overhead imagery. We compare several methods using an objective quality metric, PSNR, as well as show visual results of the proposed method. For all intents and purposes, we have simulated sensor resolution limitations from overhead imagery as a nearest neighbor down-sampling model.

### 4.1. Satellite Imagery Datasets

We have conducted experiments on several public overhead imagery datasets. These include *SpaceNet Challenge* [3], the *IARPA Multi-View Stereo Satellite Challenge* [1], and the Vehicle Detection in Aerial Imagery dataset (*VEDAI*) [18]. See figure 4 for examples of each.

The SpaceNet dataset was released as part of the *SpaceNet Challenge*. In our experiments, we have used the multi-band images corresponding to the AOI-2 site, which captures several views of the city of Las Vegas. These images have 30cm GSD and were collected with the WorldView-3 sensor from Digital Globe [3]. Another dataset used was gathered from the recent *IARPA Multi-View Stereo (MVS) Challenge*. This dataset consists of 50 WorldView-3 panchromatic images with 30cm GSD over a 50 sqkm area near Buenos Aires in Argentina [1]. We have also extracted chips with airplanes from this imagery to evaluate the algorithms with targeted features. The *VEDAI* dataset was released in 2015 as a benchmark for vehicle detection tasks in aerial imagery. It has more than a thousand images with various objects of interest, including vehicles, boats, tractors, and aircraft. Table 1 summarizes the data used in our experiments.

Table 1. Summary of Datasets.

Dataset	Chip Size	Number of Image Chips
SpaceNet - Las Vegas	$256 \times 256$	45266
IARPA MVS - All	$256 \times 256$	382795
IARPA MVS - Aircraft	$344 \times 344$	1056
VEDAI	$256 \times 256$	3734

### 4.2. Super-Resolution Results

We have conducted our experiments on several overhead images, as described in the previous section. Results are reported using PSNR on a validation subset of images for each dataset. Table 2 summarizes a comparison of our approach with several state-of-the-art algorithms based on other CNN architectures. We have focused our analysis and comparisons on two GAN-based models. First, a system proposed

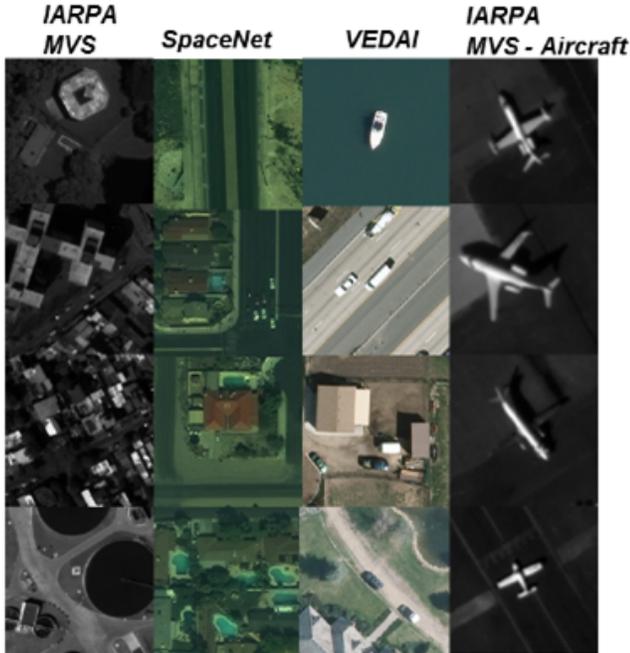


Figure 4. Example samples used in our overhead imagery experiments.

by Ledig et al., in [17], one of the pioneering works on using GANs for super-resolution. We refer to this work as *SR-DCGAN*. Second, an approach introduced in [13] to tackle the image translation problem has been repurposed for super-resolution tasks in our work. Note that this work was not explicitly defined for super-resolution, but rather for a more generic set of image domain mappings. Finally, the last algorithm we have used for comparisons was presented in [4]. It is a non-GAN scheme based on PixelCNN that uses pixel recurrency to predict current samples.

In figures 5, 6, and 7 we show results of the proposed GAN scheme using dense blocks applied to *VEDAI*, *SpaceNet - Las Vegas*, and *IARPA MVS* datasets for a factor of  $4\times$  super-resolution. These results present a more tangible description of the true performance of our system for overhead imagery.

**Limits of Extreme Super-Resolution** We also looked into super-resolution transformations at larger scale factors, *e.g.*,  $> 4\times$ , as well as identifying limits beyond which these systems become unsuited for processing. Figure 8 shows the degradation as we double the resolution factor, or in other words as we downsample the input while aiming for same output resolution.

Finally, for reference we provide results on non-overhead imagery. We used the popular *CelebA* dataset containing many examples of human faces. Figure 9 shows some results of our proposed architecture using DenseNet in the generator. Images have been super-resolved by a fac-

Table 2. Comparison between state-of-the-art network architectures for super-resolution tasks for three overhead imagery datasets. (L1: L1-loss, FM: Feature matching loss, A: Adversarial loss)

Algorithm	VEDAI	SpaceNet Vegas	IARPA MVS	Average
	Quality dBs	Quality dBs	Quality dBs	Quality dBs
SR-DCGAN [17]	29.4	29.6	28.3	29.1
PixelCNN [4]	27.7	29.1	28.3	28.4
pix2pix [13]	29.0	30.8	<b>29.9</b>	29.9
DenseNet GAN				
loss: L1, FM, A (Ours)	29.1	30.7	27.5	29.1
DenseNet GAN				
loss: L1, A (Ours)	<b>29.9</b>	<b>31.3</b>	29.6	<b>30.3</b>

tor of 8, *i.e.*, from  $8 \times 8$  inputs to  $64 \times 64$  outputs.

## 5. Results Discussion

Mapping from lower resolution images to higher resolutions is a specific case of the broader image translation problem. Fine-grained scene details can be lost during the image acquisition process. In overhead imagery, and long-range imaging applications, this translates into information, thus intelligence, loss. Our goal is to produce new plausible information to limit the impact of the imaging system resolution loss for lower cost imagers and long-range imaging applications. We have designed a system that is shown many examples of low-resolution and high-resolution image pairs and asked to learn the non-linear mapping that occurs.

From the results we have shown, we can see that overhead imaging has its own set of challenges (*e.g.*, there is a large variation of features present in the scene). It is difficult to obtain a well-balanced training set that allows to fully model the different viewing geometries of the sensor and variety of objects and geospatially diverse backgrounds present in various scenes. Also, having a nadir view constrains the problem somewhat, as we can see by comparing the results obtained for the *SpaceNet* dataset, figure 6, in contrast to results for the *IARPA MVS* collection with much more angular diversity (figure 7). The proposed super-resolution algorithm can recover information better in the former case than the latter case.

Comparing the results of the proposed method with the other algorithms (see table 2), we can see that adding lay-



Figure 5. Results on the VEDAI dataset. Super-resolution factor of 4.



Figure 6. Results on the SpaceNet Las Vegas dataset. Super-resolution factor of 4.

ers of DenseNet blocks into the GAN generator improves the performance of other configurations for this task. These results are aligned with findings and claims found in [12], where dense blocks are superior schemes compared to more traditional convolutional and residual layers due to being able to expose lower-level features at deeper layers of the network. This is certainly an advantage for low-level tasks such as super-resolution. We did not observe any improvement by adding more dense blocks or increasing the growth rate. We tried to keep the complexity of the network as low as possible with no noticeable loss in performance.

As perhaps expected, the network performed worse when using VGG-16 pre-trained on Imagenet images as a universal feature extractor to compute the feature matching loss component of the overall loss function, as shown in table 2 (last two rows). There is little to no correlation between features found in a natural image to the features found in an overhead image. A standard universal feature extractor does not efficiently capture visual attributes at such different geometric viewpoints.

Figure 8 shows the performance of the network when resolution of the input image is reduced further. At  $8\times$  factors (fourth row in the figure) it is still capable to recover and properly estimate the content in the image. However, at lower input resolutions (second from the bottom row) the network is completely “hallucinating” the wrong content.

These types of exercises are useful, though, to show limitations, feasibility, and generalization of the solution.

We have also observed that, when training data is constrained to particular semantic categories (*e.g.*, aircraft dataset extracted from the *IARPA MVS* images), the capability of the network to train the probability distribution is quite good even with a small number of samples. One way of achieving meaningful improvements on generic datasets is to pre-process the low resolution image with a functional model and add the (coarse) semantic results to the GAN, using a Conditional GAN (cGAN) framework.

We feel that these results reveal the current potential of these methods. Overall, the results are encouraging, as the network is capable of recovering details to improve the interpretability of the image, not just for improved quality purposes but for functionality of the imagery and potential impact on automatic target recognition (ATR) applications. As we enter an era where data quality becomes as important, if not more important, as the algorithm design, there is a challenge ahead of optimizing the methodology of how to train these networks for super-resolution tasks, as well as understanding limitations, expectations, and feasibility to recover information loss.



Figure 7. Results on the IARPA MVS dataset. Super-resolution factor of 4.

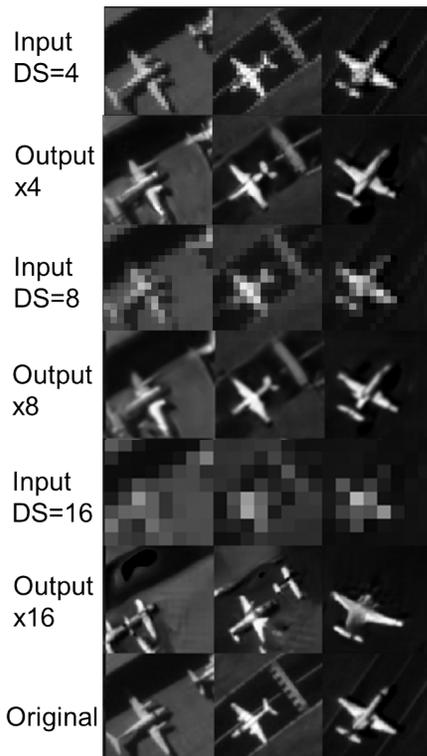


Figure 8. Results of the proposed system as the resolution factor is doubled.

## 6. Conclusion

Given recent progress in super-resolution using Deep Learning, overhead imagery is one field that can leverage these advances and use it for improved information exploitation capabilities. In this exploratory work, we have explored state-of-the-art super-resolution neural network models and proposed an architecture based on dense blocks to carry out the task of increasing semantic meaning by adding plausible realistic information in the scene. Our model aims at learning the probability distribution mapping

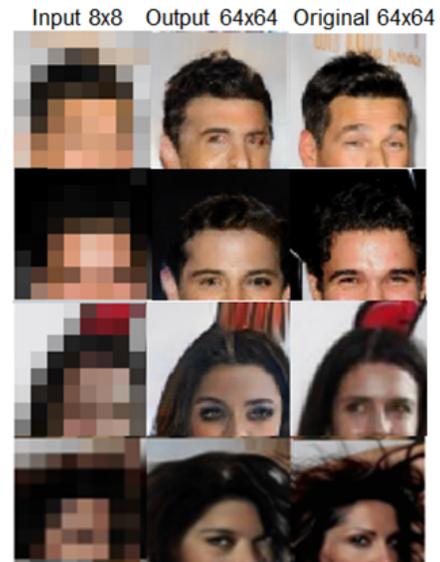


Figure 9. Results of the proposed system on human faces.

from the low-resolution to high-resolution using adversarial training from many exemplar data sets and then applying it to never before seen data. We have compared several generative models and architectures with several publicly available satellite and airborne image datasets (panchromatic electro-optical (EO) and multi-band images), and have shown what is realistically possible with today's tools. We have also shown that a GAN framework with a collection of modified dense blocks in the generator can outperform state-of-the-art models that have been proposed for natural images.

## References

- [1] M. Bosch, Z. Kurtz, S. Hagstrom, and M. Brown. A multiple view stereo benchmark for satellite imagery. In *Proceedings of the Applied Imagery Pattern Recognition Workshop (AIPR)*, Washington, DC, USA, 2016.

- [2] J. Bruna, P. Sprechmann, and Y. LeCun. Super-resolution with deep convolutional sufficient statistics. In *Proceedings of the International on Learning Representations (ICLR)*, 2016.
- [3] CosmiQWorks, DigitalGlobe, and NVIDIA. Spacenet. <http://explore.digitalglobe.com/spacenet>, 2016.
- [4] R. Dahl, M. Norouzi, and J. Shlens. Pixel recursive super resolution. *CoRR*, abs/1702.00783, 2017.
- [5] E. Denton, S. Chintala, A. Szlam, and R. Fergus. Deep generative image models using a laplacian pyramid of adversarial networks. In *Proceedings of the Conference on Neural Information Processing Systems (NIPS)*, Montreal, Canada, 2015.
- [6] C. Dong, C. C. Loy, K. He, and X. Tang. Image super-resolution using deep convolutional networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(2):295–307, 2016.
- [7] W. Freeman, T. Jones, and E. Pasztor. Example based super-resolution. *Proceedings of the IEEE Computer Graphics and Applications*, 22(2):56–65, 2002.
- [8] W. Freeman and E. Pasztor. Markov networks for superresolution. In *Proceedings of the Annual Conference on Information Sciences and Systems (CISS)*, pages 1841–1848, 2000.
- [9] D. Garcia. Image super-resolution through deep learning. <https://github.com/david-gpu/srez>, 2016.
- [10] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2672–2680, 2014.
- [11] 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 (CVPR)*, Las Vegas, NV, USA, 2016.
- [12] G. Huang, Z. Liu, and K. Q. Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Honolulu, HI, USA, 2017.
- [13] P. Isola, J. Zhu, T. Zhou, and A. A. Efros. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Honolulu, HI, USA, 2017.
- [14] R. Keys. Cubic convolution interpolation for digital image processing. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 29(6):1153–1160, 1981.
- [15] J. Kim, J. K. Lee, and K. M. Lee. Deeply-recursive convolutional network for image super-resolution. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Las Vegas, NV, USA, 2016.
- [16] W. Lai, J. Huang, N. Ahuja, and M. Yang. Deep laplacian pyramid networks for fast and accurate super-resolution. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Honolulu, HI, USA, 2017.
- [17] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. P. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi. Photo-realistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Honolulu, HI, USA, 2017.
- [18] S. Razakarivony and F. Jurie. Vehicle detection in aerial imagery. *Journal on Visual Communication and Image Representation*, 34(C):187–203, Jan. 2016.
- [19] S. Schuler, C. Leistner, and H. Bischof. Fast and accurate image upscaling with super-resolution forests. In *Proceedings of the IEEE Conference Computer Vision and Pattern Recognition (CVPR)*, Boston, USA, 2015.
- [20] E. Shelhamer, J. Long, and T. Darrell. Fully convolutional networks for semantic segmentation. *CoRR*, abs/1605.06211, 2016.
- [21] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. A. Riedmiller. Striving for simplicity: The all convolutional net. *CoRR*, abs/1412.6806, 2014.
- [22] Y. Tai, J. Yang, and X. Liu. Image super-resolution via deep recursive residual network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Honolulu, HI, USA, 2017.
- [23] J. Yang, J. Wright, T. S. Huang, and Y. Ma. Image super-resolution via sparse representation. *IEEE Transactions on Image Processing*, 19(11):2861–2873, 2010.
- [24] R. Zeyde, M. Elad, and M. Protter. On single image scale-up using sparse-representations. In *Proceedings of the International Conference on Curves and Surfaces*, pages 711–730, 2012.