# Accepted Manuscript

Single Satellite Imagery Simultaneous Super-resolution and Colorization using Multi-task Deep Neural Networks

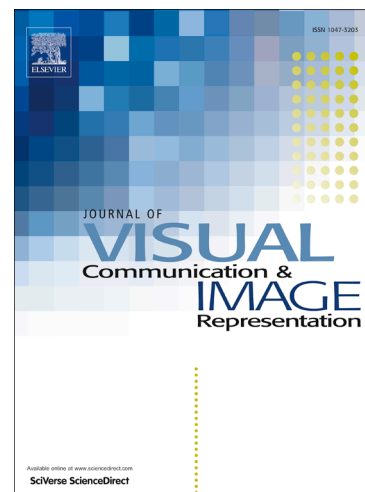Heng Liu, Zilin Fu, Jungong Han, Ling Shao, Hongshen Liu

Please cite this article as: H. Liu, Z. Fu, J. Han, L. Shao, H. Liu, Single Satellite Imagery Simultaneous Super-resolution and Colorization using Multi-task Deep Neural Networks, *J. Vis. Commun. Image R.* (2018), doi: https://doi.org/10.1016/j.jvcir.2018.02.016

# Single Satellite Imagery Simultaneous Super-resolution and Colorization using Multi-task Deep Neural Networks

Heng Liu[a], Zilin Fu[a], Jungong Han[b], Ling Shao[c], Hongshen Liu[a]

[a]*Anhui University of Technology, China, 243032*
[b]*Lancaster University, U.K., LA1 4YW*
[c]*University of East Anglia, U.K., NR4 7TJ*

## Abstract

Satellite imagery is a kind of typical remote sensing data, which holds preponderance in large area imaging and strong macro integrity. However, for most commercial space usages, such as virtual display of urban traffic flow, virtual interaction of environmental resources, one drawback of satellite imagery is its low spatial resolution, failing to provide the clear image details. Moreover, in recent years, synthesizing the color for grayscale satellite imagery or recovering the original color of camouflage sensitive regions becomes an urgent requirement for large spatial objects virtual reality interaction. In this work, unlike existing works which solve these two problems separately, we focus on achieving image super-resolution (SR) and image colorization synchronously. Based on multi-task learning, we provide a novel deep neural network model to fulfill single satellite imagery SR and colorization simultaneously. By feeding back the color feature representations into the SR network and jointly optimizing such two tasks, our deep model successfully achieves the mutual cooperation between imagery reconstruction and image colorization. To avoid color bias, we not only adopt the non-satellite imagery

to enrich the color diversity of satellite image, but also recalculate the prior color distribution and the valid color range based on the mixed data. We evaluate the proposed model on satellite images from different data sets, such as RSSCN7 and AID. Both the evaluations and comparisons reveal that the proposed multi-task deep learning approach is superior to the state-of-the-art methods, where image SR and colorization can be accomplished simultaneously and efficiently.

## 1. Introduction

### 1.1. Motivation

Remote sensing satellite imagery holds the characteristics of extensive coverage, strong macro integrity and consistent imaging scales, which can be widely used for the spatial information related virtual reality applications, such as resource survey virtual interactions, urban traffic virtual analysis, climate change virtual display, and military action virtual deduction. However, due to optical device and imaging sensor limitations coupled with the extreme distance between sensor and sensed object on earth, one natural drawback of satellite imagery is that the spatial resolution is always low, which leads to the inaccurate sensing data due to lack of image details.

On the contrary, high-resolution (HR) satellite imagery, which is very helpful for the realization of large scale spatial information virtual reality (VR), allows extracting the rich details and accurate information at multi-level scales. In order to improve the spatial resolution of the satellite images,

2

the traditional hardware handling method reduces the physical sizes of the charge-coupled device (CCD) or complementary metal oxide semiconductor (CMOS) sensors among sensor fabrication procedure, which will easily generate shot noise that severely degrades the image quality (Yang and Huang, 2010). In addition, manufacturing imaging chips and optical instruments to capture very high-resolution images will cost huge. Thus, it is necessary to exploit signal processing techniques to reconstruct the high-resolution (HR) images from the degraded low-resolution (LR) remote sensing images, which is specifically referred to as satellite imagery SR.

In addition to low resolution, the color of satellite imagery can be easily faded due to inappropriate illumination, exposure and storage. Moreover, satellite imagery sometimes even cannot reflect the actual original color of the observed targets. For example, intentional camouflage is a common means of visually hiding the military facilities or important infrastructure, where the color as well as the appearance of these special targets are always altered and disguised. In order to get clear and accurate knowledge of these objects, recoloring the disguised imagery and enhancing their spatial resolution at the same time become a pressing demand. Therefore, it is necessary to solve these two problems - imagery SR and colorization - simultaneously in one framework. It should be noted that for some military applications, the multi-task simultaneous imagery SR and colorization perhaps is especially significant, for example solider combat VR glasses. Here, to keep consistency with human visual perception, the word 'colorization' throughout the whole text only refers to the visible light 3-bands color operations.

3

## 1.2. Related work

Over the past five years, a considerable number of image SR works have addressed to reconstruct HR satellite imagery from LR inputs. Usually, these methods are divided into two categories: multiple images reconstruction (Pickup, 2007; Zhang et al., 2014; Hung et al., 2016; Zhu et al., 2016; Brodu, 2016; Alvarez-Ramos et al., 2016) or single image SR (Liebel and Körner, 2016; Patrick, 2016).

Bayesian machine learning method was firstly applied for multi-frame super-resolution (Pickup, 2007), which fully utilizes a prior distribution over the super-resolution image. This Bayesian inference method was improved with variation approximation (Hung et al., 2016) to estimate the distribution of HR satellite imagery, the registration parameter, and some other hyper-parameters. Due to possible resolution differences in multi-angle remote sensing images over the same scene, adaptive weighting schemes (Zhang et al., 2014) are utilized to reconstruct HR satellite imagery. In addition, adaptive multi-scale detail enhancement measures (Zhu et al., 2016) were attempted for multiple LR satellite images SR. Moreover, sparse representation (Alvarez-Ramos et al., 2016) has been employed to deal with overlapping blocks for satellite image SR. Recently, the band-specific information is also applied in resolution enhancement for multi-spectral and multi-resolution satellite images (Brodu, 2016), where the independent reflectance of LR bands is preserved in details reconstruction.

Actually, for satellite imaging, even if it is easily to orbit to acquire multi-frame images of the same scene on a regular basis, the imaging scenes will always keep changes due to many uncontrollable factors, such as clouds or

4

<sub>75</sub> snow coverage, objects moving or seasonal alternation. Thus, if there were
<sub>76</sub> no available or reliable multi-frame data, single satellite imagery SR would
<sub>77</sub> become a more challenging problem. Fortunately, recent developments in the
<sub>78</sub> field of deep learning cast a bright way for single remote sensing image SR.
<sub>79</sub> An end-to-end CNN model (Dong et al., 2014, 2016), referred as SRCNN, has
<sub>80</sub> been proposed recently and successfully applied in single image SR. Then,
<sub>81</sub> Liebel and Körner (2016) retrain the SRCNN model for multi-spectral re-
<sub>82</sub> mote sensing imagery SR with a domain-specific data set to introduce the
<sub>83</sub> characteristics of multiple spectral bands. Furthermore, motivated by resid-
<sub>84</sub> ual learning (He et al., 2016), Patrick (2016) proposes to construct a deep
<sub>85</sub> residual network for single satellite imagery SR.

<sub>86</sub> On the other hand, very recently a few works (Larsson et al., 2016; Zhang
<sub>87</sub> et al., 2016) have exploited deep models to address the problem of image
<sub>88</sub> colorization which augment color from gray-scale images. These methods
<sub>89</sub> manage to learn the corresponding color representation or color distribution
<sub>90</sub> by constructing deep networks and training it with ImageNet data set (Rus-
<sub>91</sub> sakovsky et al., 2015). It should be noted that such colorization methods
<sub>92</sub> actually carry out color remapping and do not consider keeping the recon-
<sub>93</sub> struction accuracy of pixels' intensity value between input and output image.
<sub>94</sub> For the accurate understanding and better utilizing of low quality satel-
<sub>95</sub> lite imagery in large spatial related virtual reality applications, in this work
<sub>96</sub> we provide an efficient approach that not only can reconstruct HR satellite
<sub>97</sub> imagery from single LR input, but also is able to simultaneously colorize the
<sub>98</sub> grayscale satellite imagery with appropriate color information. Our contri-
<sub>99</sub> butions can be summarized as:

100  • We propose a multi-task deep neural model to achieve satellite imagery
101  SR and colorization simultaneously. Our multi-task deep model contains two
102  concurrent but not separated task networks - image features of colorization
103  network are fed back to the beginning of the feature representation parts
104  of the SR network and these two kinds of loss are combined for a joint
105  optimization. To the best of our knowledge, this is the first work which
106  explores to achieve satellite imagery SR and colorization cooperatively.

107  • In order to avoid color bias in imagery colorization, we incorporate
108  natural images with satellite data to enrich the color diversity and we manage
109  to realize the expectation color distribution learning based on these mixed
110  data.

111  • We introduce a novel multi-scale deep encoder-decoder symmetrical
112  network for satellite imagery SR, where a residual structure is adopted to
113  improve the imagery reconstruction performance.

## 114  2. Methods

### 115  2.1. Overall scheme

116  In order to overcome the processing irrelevancy of existing image SR and
117  colorization methods, our comprehensive consideration is adopting a multi-
118  task optimization strategy which not only can reconstruct the HR satellite
119  imagery but also can colorize the gray scale imagery for proper color infor-
120  mation. In the sense of this, we manage to achieve the cooperative learning
121  tasks for satellite imagery by constructing and training a multi-task deep
122  neural network based on satellite imagery data set.

123  There are several components in our multi-task satellite imagery SR and

6

124 colorization deep model, including multi-scale SR reconstruction, color distri-
125 bution prediction based grayscale colorization, features interaction between
126 SR and colorization parts, and multiple tasks synchronous optimization.

127 Benefiting from the powerful non-linear mapping, SRCNN (Dong et al.,
128 2014, 2016) improves the performance dramatically compared with the tra-
129 ditional SR methods. Since training SRCNN model usually takes a very long
130 time before convergence, Liang et al. (2016) introduce Sobel edge detection
131 so as to capture gradient information to accelerate the training convergence.
132 In fact, the method does reduce the training time but the reconstruction im-
133 provement is rather limited. In addition to image gradient priors, in view of
134 the network depth with residual structure (He et al., 2016) is of crucial im-
135 portance to a remarkable performance improvement, Kim et al. (2016) take
136 twenty convolution layers with residual connection to construct deep network
137 for image SR reconstruction.

138 The negligence of the above deep SR approaches is that the multiple
139 scales image context in SR reconstruction is not fully utilized at all. Consid-
140 ering the fact that image multi-scale contextual information is essential for
141 the image details reconstruction, in this work, we propose to take a multi-
142 scale symmetrical CNN for image SR. In addition, we also introduce residual
143 structure from the LR input to the end of the network so as to improve the
144 reconstruction accuracy.

145 For satellite imagery colorization based on grayscale component input, we
146 employ a structure similar to Zhang's network (Zhang et al., 2016) to produce
147 the corresponding color distribution under the fused (satellite imagery and
148 natural images) data set. Since the color diversity of satellite imagery is very
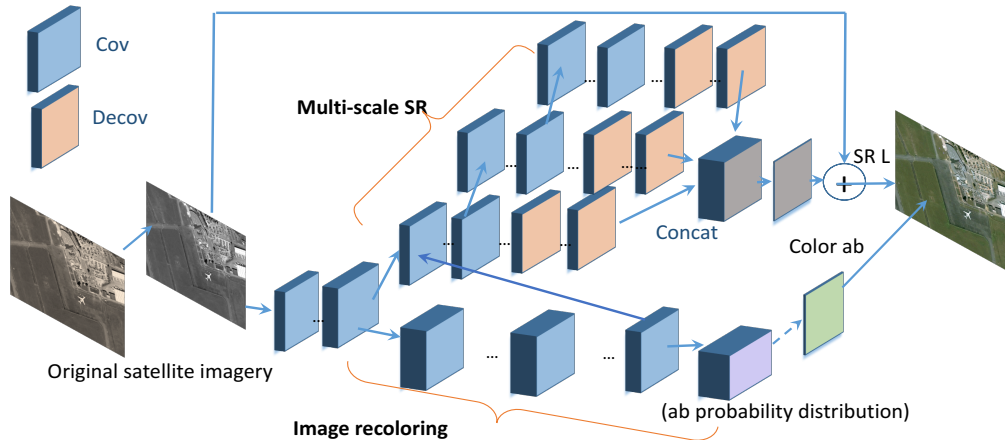
7

Figure 1: The overall multi-task satellite imagery deep SR and colorization model (in the figure, 'ab' refers to the color components of Lab color space).

<sup>149</sup> different from natural image, we recalculate the color statistics prior, instead

<sup>150</sup> of the one available for natural images, and on top of it, we adjust the color

<sup>151</sup> class re-balancing coefficient based on fused data.

<sup>152</sup> In addition, for SR and colorization multiple tasks cooperation, we choose

<sup>153</sup> a feedback strategy. Specifically, the final convolution features of colorization

<sup>154</sup> network are back propagated to the SR network and blended with the input

<sup>155</sup> LR imagery representation together for the HR reconstruction cooperatively.

<sup>156</sup> Our overall multi-task deep SR and colorization model is shown in Fig. 1.

<sup>157</sup> *2.2. Multi-scale learning for satellite imagery SR*

<sup>158</sup> With the capability of the hierarchical feature learning, multi-scale deep

<sup>159</sup> convolution networks appear in literature, including edge detection (Xie and

<sup>160</sup> Tu, 2015), skeleton extraction (Shen et al., 2016) and image dehazing (Ren

<sup>161</sup> et al., 2016). In a recent work (Szegedy et al., 2016), convolution filters

<sup>162</sup> with variable sizes are carefully designed and applied in multiple residual

8

163 connections, which will lead to a very wide inception networks with better
164 learning performance. In general, the common characteristics of these multi-
165 scale works are taking different length convolution branches or different sizes
166 filters to achieve different sizes receptive fields so as to extract the image
167 features at different scales.

168 In addition, without fully connected layers, the fully convolutional net-
169 works (FCNs) containing only convolution and deconvolution layers have
170 been successfully applied to semantic segmentation (Hong et al., 2015) and
171 object detection (Yang et al., 2016). Here, a convolution layer can be inter-
172 preted as an encoder which serves for features extraction and representation
173 while a deconvolution one, named by the decoder, acts as reconstruction.

174 For satellite imagery SR, we adopt a multi-scale deep symmetrical encoder-
175 decoder structure. Obviously, the imagery $f(x)$ will be encoded with multiple
176 scales features by different lengths convolution layers (short for coarse scale
177 and long for fine scale). Through symmetrical decoding, the different lengths
178 deconvolution layers will reconstruct the original imagery based on the multi-
179 scale feature representations in a variety of scales. Actually, for an imagery
180 $f(x)$ in $L^2$ space $R$, the principle of multi-scale encoding and decoding can
181 be formalized by wavelet multi-resolution analysis (MRA) (Mallat, 1999) as:

$$f(x) = \sum_{k \in Z}^{N} a_k^{j_0} \phi_k^{j_0}(x) + \sum_{j=j_0}^{J} \sum_k b_k^j \psi_k^j(x), \tag{1}$$

182 where $j$ is the scale varying from $j_0$ to $J$, $k$ is the index of basis function, and
183 $\{a_k^{j_0}\}$, $\{b_k^j\}$ are coefficients attached to the approximation (scale) function
184 $\phi(x)$ and the detail (wavelet) function $\psi(x)$, respectively. In short, the image
185 $f(x)$ can be viewed as consisting of two components (see Eq. (1)): the low-

9

frequency approximation and the high-frequency detail. When varying the scale j from zero to certain scale, $f(x)$ can be represented as the weighted summation of a series of components at different scales, which contains a low-frequency approximation and several or numerable high-frequency details. From deep learning point of view, Eq. (1) may be treated as a combination of deconvolution (reconstruction) operations at multiple scales. Assuming at each scale, $\tilde{f}_j$, represents a reconstruction of $f(x)$. Thus, according to Eq. (1), if we take a summation function $s$ adding up all encoder-decoder streams, the multi-scale encoder-decoder reconstruction $\tilde{f}(x)$ can be easily represented as:

$$\tilde{f}(x) = s(\tilde{f}_1, \tilde{f}_2, \cdots, \tilde{f}_j, \cdots) \tag{2}$$

Then the optimization target of our multi-scale encoder-decoder learning can be regarded as:

$$\tilde{f} = \arg\min_{\Theta}(\left\| f - \sum_j (F_j^a(y, \Theta_j^a) + F_j^b(y, \Theta_j^b)) \right\|_2^2), \tag{3}$$

where $f$ and $y$ represent the HR image and the corresponding LR image, and $F(\cdot)$ denotes the network reconstruction function. $\Theta$ is the learned parameter of the network and symbols $j$, $a$, $b$ indicate a specific scale, a low-frequency approximation component and a high-frequency component, separately. By taking into account the components of different scales simultaneously, multi-scale learning will partially overcome the deficiency of only considering the energy amplitude recovering (concentrated in low-frequency components) while ignoring the structural details (in high-frequency components).

Given a set of LR and HR image pairs $\{f_i, y_i\}_{i=1}^N$, if directly treating the

10

208 input LR image $y_i$ as the low-frequency approximation component of HR

209 image $f_i$ and omitting the high-frequency indicator $b$, the loss function of the

210 proposed multi-scale encoder and decoder learning can be finally denoted as:

$$Loss(\Theta) = \frac{1}{N} \sum_{i=1}^{N} \left\| f_i - \sum_{j}(y_i + F_j(y_i, \Theta_j)) \right\|^2 \tag{4}$$

### 211 2.3. Color distribution prediction based imagery colorization

212 Although the original meaning of colorization refers to adding color to

213 the gray-scale image, colorization for satellite imageries is more of recolor-

214 ing, which indicates enhancing or changing the original color of the input

215 satellite images desiired by specific applications, such as camouflage. The

216 reason behind is that satellite imageries in most cases are already 3-bands

217 color data. In practice, recoloring can also be performed in pure colorization

218 way - extracting the intensity channel and colorizing it. In general, There are

219 two different strategies for gray-scale imagery colorization: direct color pre-

220 diction based on Euclidean color regression loss (Cheng et al., 2015; Iizuka

221 et al., 2016) and multi-modal color distribution prediction based Softmax

222 color classification loss (Zhang et al., 2016; Larsson et al., 2016).

223 Let $x$ denote a gray-scale channel imagery to be colored, assuming in CIE

224 *Lab* color space its associated two channels (i.e., 'ab' color components; all

225 the following 'ab' items keep the same meaning) color is $y \in \mathbb{R}^{h \times w \times 2}$ (where

226 $h, w$ are image dimensions), the objective of color prediction model is to learn

227 a mapping $\tilde{y} = f(x)$ such that the Euclidean loss $L_2(.,.)$ between predicted

228 and ground truth colors is minimized after training:

$$L_2(\tilde{y}, y) = \frac{1}{2h \times w} \sum_{h,w} \|y_{h,w} - \tilde{y}_{h,w}\|_2^2 \tag{5}$$

11

229 Obviously, the Euclidean regression loss will lead the optimal solution $\tilde{y}$
230 to be the mean of all pixels' color of the ground truth image, which favors
231 unsaturated color prediction results. Moreover, the solution does not consider
232 the problem of color plausibility will in fact give inveracious and implausible
233 color results. Thus, Euclidean loss based color prediction way does not handle
234 the ambiguity and multi-modal color distribution well.

235 In this work, for satellite imagery colorization, we can use a deep neural
236 network to learn a mapping $m(x)$ to a color distribution $\tilde{z}$ over possible $ab$
237 color bins ($\tilde{z} \in [0,1]^{h \times w \times q}$, q is the number of color bins ) for a given input $x$.
238 Then we compare the predicted color distribution with the encoded ground
239 truth one and calculate the Softmax cross entropy loss for optimization. We
240 also take color class rebalancing technique to enhance the impact of rare
241 color in the distribution. Finally, we take the $annealed - mean$ technique
242 (Kirkpatrick et al., 1983) to estimate the color of every pixel based on its
243 corresponding color distribution. Our imagery colorization network is sim-
244 ilar to the approach of Zhang et al. (2016) but with two main differences:
245 different means of acquiring the color probability density of satellite imagery
246 data and adopting features interacting feeding back structure for multi-task
247 cooperation. More specially, we calculate the color probability distribution
248 under AID satellite data set (Xia et al., 2017) and fuse it with the prior
249 color probability of ImageNet data set (Russakovsky et al., 2015); we expand
250 the actual number of supported $ab$ color bins under satellite data set to be
251 313 to overcome the problem that the color scope of imagery is not wide
252 enough; finally, two deconvolution layers help to feed back the convolution
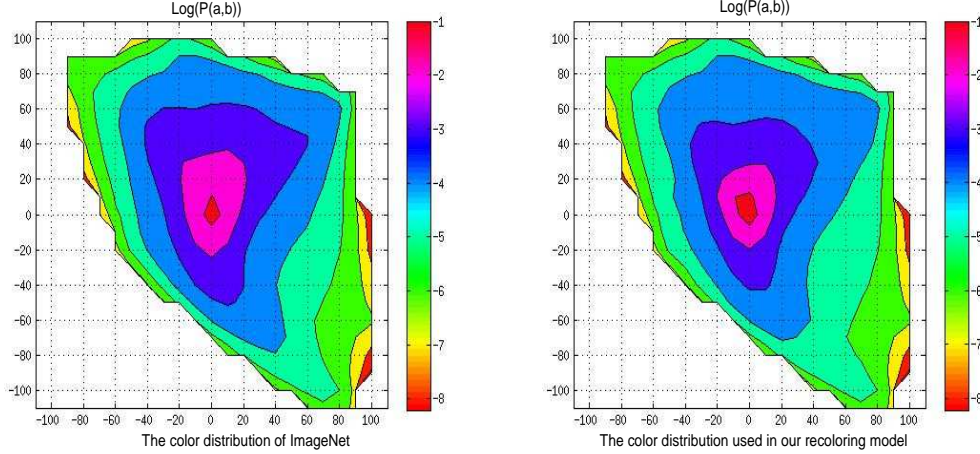253 features of coloriation to the SR network. The original color probability den-

12

Figure 2: The contrast between ImageNet's color distribution (left) and the color distribution used in our colorization/recoloring model (right).

<sup>254</sup> sity distribution of ImageNet data and the corresponding color distribution

<sup>255</sup> used in our colorization/recoloring model are illustrated in Fig. 2 (shown

<sup>256</sup> in log scale). From the figure, we can see the high probability region of our

<sup>257</sup> color probability distribution slightly shrinks, compared to the corresponding

<sup>258</sup> color distribution one of ImageNet data set, which perhaps is due to the color

<sup>259</sup> scope of satellite imagery is much narrower than that of natural images.

<sup>260</sup> With the fused color probability and the expanded $ab$ color bins, each

<sup>261</sup> ground truth color $y$ can be easily encoded to a color vector presentation

<sup>262</sup> $z(\in [0,1]^{h \times w \times q})$ with its nearest neighbor color bins. For whole imagery color

<sup>263</sup> prediction, we define the cross-entropy loss of such color encoding prediction

<sup>264</sup> $L_{ce}(.,.)$ as following:

$$L_{ce}(\tilde{z}, z) = \sum_{h,w} c(z_{h,w}) \sum_{q} z_{h,w,q} log(\tilde{z}_{h,w,q}) \tag{6}$$

<sup>265</sup> Here, $c$ is a loss weighting factor used to consider the effect of the color-

13

266 class rarity. At last, we estimate the final color values $\tilde{y}$ by mapping the

267 probability distribution $\tilde{z}$ through simulated annealing way. The detailed

268 techniques on color rebalance and color estimation can be referred to Zhang

269 et al. (2016).

### 270 2.4. Joint multi-task learning for satellite imagery SR and colorization

271 Our satellite imagery multi-task deep model actually combines the pro-

272 posed SR network and colorization network for concurrent execution by the

273 convolutional features sharing (see the two front convolution layers in Fig. 1)

274 and the features interaction (see the feedback from colorization network to

275 SR network in Fig. 1). Based on Eq. (4) and Eq. (6), for any low resolution

276 and gray-scale input image $x_i$ ($h, w$ are its height and width), if assuming its

277 HR label image in SR model is $f_i$ and its corresponding ground truth color

278 distribution is $z(x_i)_{h,w}$, then the loss of multi-task joint learning for satellite

279 imagery SR and colorization can be formalized as :

$$
\begin{aligned}
Loss_{total}(\Theta, \tilde{z}) = \frac{1}{N} \sum_{i=1}^{N} (\left\| f_i - \sum_j (x_i + F_j(x_i, \Theta_j)) \right\|^2 \\
+ \eta \sum_{h,w} c(z(x_i)_{h,w}) \sum_q z(x_i)_{h,w,q} log(\tilde{z}(x_i)_{h,w,q})),
\end{aligned}
\tag{7}
$$

280 where $\eta$ is a regularization factor which controls the effects of SR recon-

281 struction loss and color distribution loss in the whole optimization. Obvi-

282 ously, through such joint learning, the procedures of SR reconstruction and

283 multi-modal color prediction will constantly regularize each other and be op-

284 timized simultaneously. When the multi-task model is trained to converge,

285 the acquired solution (the parameters of the deep model) will be an optimal

286 trade-off which not only can reconstruct the low resolution image well but

14

<sub>287</sub> also can map it to a color image with strong sense of reality. Through such <sub>288</sub> joint learning, a satellite imagery with high resolution and visual realistic <sub>289</sub> color can be obtained directly.

## 3. Experiments and discussions

### 3.1. Data sets and evaluation measures

<sub>292</sub> The imageries from AID data set (30 different scene classes with about <sub>293</sub> 200 to 400 samples of size $600 \times 600$ in each class) may be used for SR training <sub>294</sub> while other images from RSSCN7 (Zou et al., 2015) (7 scene categories with <sub>295</sub> 400 samples of size $400 \times 400$ in each class) may be utilized for testing. <sub>296</sub> For satellite imagery colorization, actually the combination of 10000 random <sub>297</sub> selected ImageNet images with 10000 AID satellite imageries is preferred to <sub>298</sub> be applied for colorization training. Also some imageries from RSSCN7 can <sub>299</sub> be regarded as the test data.

<sub>300</sub> As for the quality measurements, for HR reconstruction, well-known PSNR <sub>301</sub> metrics are adopted. For colorization evaluation, visual results are shown in <sub>302</sub> contrast. We notice that in a down-sampled satellite imagery different regions <sub>303</sub> may hold different PSNR values - smooth areas (such as plain or grassland) <sub>304</sub> will get higher PSNR scores than the uneven locations ( such as cross con- <sub>305</sub> nection or zebra line). In general, we take the whole imagery's PSNR for <sub>306</sub> quality evaluation which is a mean of all local regions' PSNR values. Fig. 3 <sub>307</sub> gives an example.

### 3.2. Model training

<sub>309</sub> There are two ways to train the proposed multi-task model: training <sub>310</sub> it from scratch or finetuning it from the colorization model ( Zhang et al.
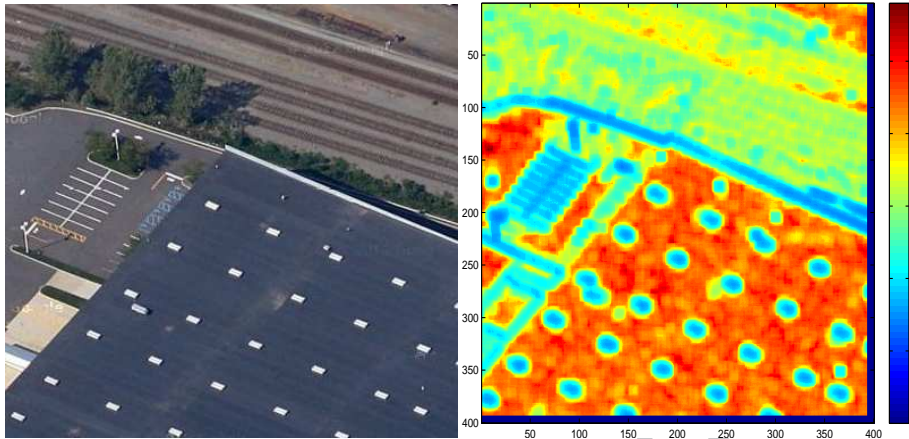
15

Figure 3: A satellite imagery (left) and the PSNR visualization for its down-sampled version (right); different image regions hold different PSNR values.

(2016)) for our multiple tasks. Actually, we tend to train the proposed multi-task model from scratch with a unified data strategy. In this case, about 20,000 images coming from AID satellite data and ImageNet are selected and confused for multi-task model training. Each image is augmented to 8 images by rotation which yields about 160,000 images as training set. Images are cropped into small overlapped patches with a size $96 \times 96$ and a stride of 27. For SR part, the cropped ground truth patches are used as the HR labels and the corresponding LR pairs are acquired by imposing the bi-cubic interpolation twice on the ground truth. For colorization part, the LR color images will be converted to Lab color space and keep the intensity component. The labels of this part are the encoded color distributions in $ab$ color bins of the ground truth.

In the training procedure, we follow the proposal from He et al. (2016) to initialize the weights of all layers. We initially set the learning rate to 0.001 and reduce it by multiplying 0.316 every 100 thousand iterations. Mo-

16

ACCEPTED MANUSCRIPT
326 mentum and weight decay parameters are set to 0.9 and 0.0001, respectively.
327 The regularization coefficient *eta* is set to be 1 at the beginning and can be
328 manually adjusted it to 1.5, emphasizing the impact of image color recover
329 once the gradients of the model become relatively small. The whole deep net-
330 work training is implemented using Adam solver from the Caffe package (Jia
331 et al., 2014) with a batch size of 32. For $4\times$ down-scaling and the confused
332 grayscale satellite data, the model training takes about 1,300,000 iterations
333 before convergence.

334 Our multi-task model can also be trained by finetuning way: training
335 it using ImageNet then finetuning with satellite data. In this second case,
336 twenty thousands images are randomly selected from ImageNet and used
337 to train the multi-task model, then some AID satellite images are taken
338 for model finetuning. All settings and parameters are the same as the first
339 training strategy. However, we found that the second finetuning way is eas-
340 ily inclined to lead to the color deviation (see Fig. 6; more examples can
341 be refereed to Fig .10). The detail configuration of our multi-task SR and
recoloring model is given in Table. 1.

Table 1: Multi-task satellite imagery SR and colorization deep network configuration.

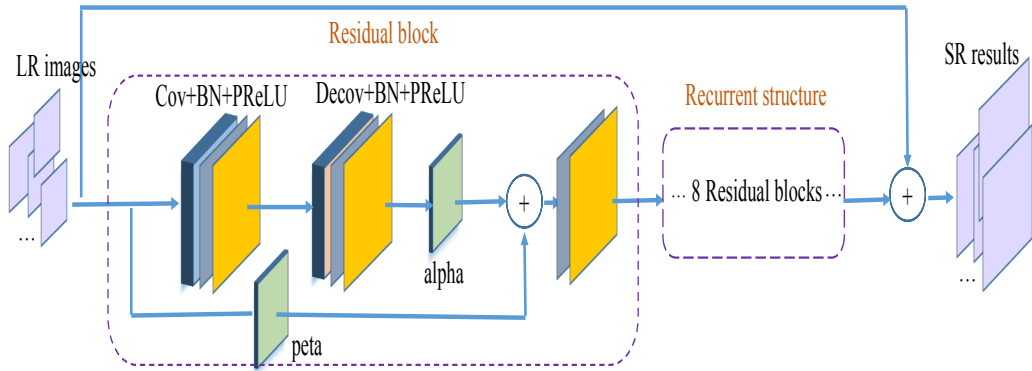| | | | | |
|---|---|---|---|---|
| **Imagery SR** | (Conv3-32)$\times$3 | (Conv3-64)$\times$3 | (Conv3-128)$\times$3 | (BatchNorm)$\times$36 |
| | (Deconv3-32)$\times$9 | (Deconv3-64)$\times$6 | (Deconv3-128)$\times$3 | (Prelu)$\times$36 |
| **Imagery Colorization** | (Conv3-64)$\times$2 | (Conv3-128)$\times$2 | (Conv3-256)$\times$5 | (Conv3-512)$\times$12 |
| | (Conv1-313)$\times$1 | (Deconv3-256)$\times$2 | (Deconv4-64)$\times$2 | (BatchNorm)$\times$7 |

342

17

Figure 4: Another satellite image SR network - deep recurrent residual network.

### 3.3. Another SR network and color distribution effect

In addition to the proposed multi-scale structure for satellite image SR, other deep structures can also be utilized to achieve the same target, such as residual skip or recurrent connection. Here, we combine residual skip with recurrent connection to form a deep recurrent residual model for satellite image SR, which is similar to Patrick's one (Patrick, 2016) but with two differences : 1) we add Batch Normalization layer and replace ReLU with PReLU after each convolution layer; 2) we add a direct skip from the input to the end of the network and fix the scale parameters in every residual block (both in bypass connection part and convolution route) while not learning it from training. The architecture of our deep recurrent residual network is illustrated in Fig. 4. According to our observations (a comparison example is shown in Fig. 5), this kind of SR network sometimes can get smoother edges, but many other imagery details will be lost after reconstruction. Thus, we finally opt for the multi-scale deep structure for imagery SR of multiple tasks.

Moreover, different color distributions coming from different training data

18

will lead to diverse colorization effect. For example, using the partial data of ImageNet for color probability distribution calculation will get very different colorization results. Specifically, for satellite image, there exists a trade-off: whether using the satellite image data or using the ImageNet data to derive the color distribution used for colorization. According to our observations, only using the satellite imagery for color distribution acquisition will inevitably cause color bias effect. Fig. 6 gives an typical illustration. Thus, in practice we get the final color distribution by data fusing strategy which is stated in Section 2.3.

### 3.4. Results of Imagery super-resolution and colorization

The proposed multi-task SR and colorization network accepts the LR grayscale satellite imagery as the input and reconstructs it to be a HR version and at the same time maps it to a colorized one. When given a low-resolution color imagery, it should be converted from RGB color space to Lab firstly, then the luminance component is pipelined into the multi-task network and a reconstructed HR and recoloring imagery will be output. Some simultane-



Figure 5: Imagery SR comparison: the recurrent residual network(middle) vs. the multi-scale network(right) with LR satellite image input(left).
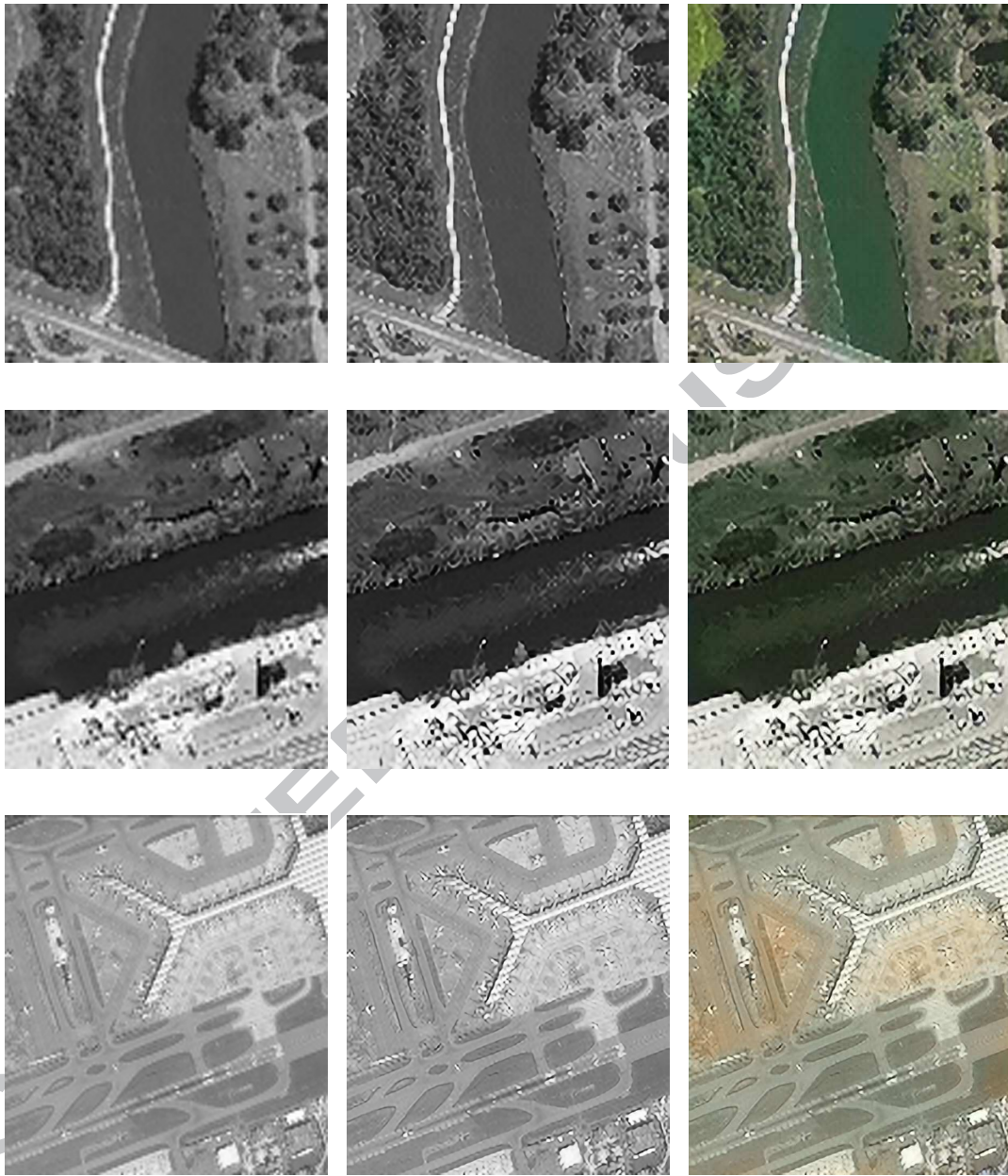
19

Figure 6: Color bias example in imagery colorization when only using satellite data to acquire color probability distribution: (left)LR satellite lake imagery; (middle)the grayscale input; (right)Color biased colorization result.

376 ous super-resolution and colorization results of the proposed multi-task deep
377 learning approach are shown in Fig. 7.

378 *3.5. Comparisons and discussions*

379     Since there is no other related work which pursues single satellite imagery
380 simultaneous SR and colorization, we choose to compare our approach with
381 the state-of-the art methods of two aspects: SRCNN (Dong et al., 2016)
382 and Patrick's method (Patrick, 2016) (the model is realized by ourselves and
383 trained with some images of SpaceNet AOI1 (SpaceNet, 2016)) for single im-
384 agery super-resolution; Zhang's method (Zhang et al., 2016), Iizuka's method
385 (Iizuka et al., 2016) and Larsson's one (Larsson et al., 2016) for single im-
386 agery colorization. We compare and evaluate the effect of super-resolved and
387 colorized imagery not only by subjective visual effect but also with objec-
388 tive PSNR(db) value. For satellite imagery super-resolution, visual results
389 involve the subjective clarity inception of imagery details. As for imagery col-
390 orization, visual results mainly refer to the color consistency and the realism
391 of the objects. These comparisons and experimental results are illustrated

20

(a) LR grayscale      (b) HR reconstruction      (c) SR and colorization

Figure 7: Some results of satellite imagery ('Riverlake' from RSSCN7 and 'Airport' from AID) simultaneous SR and colorization: (a)LR grayscale imagery; (b)Reconstructed HR imagery; (c)Super-resolved and colorized Imagery.
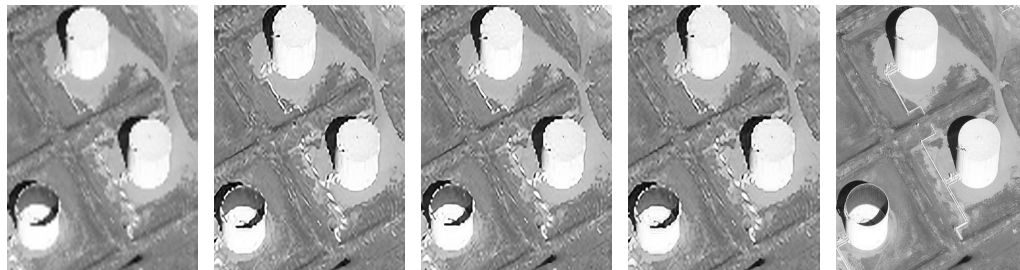
21

in Table 2, Fig. 8 and Fig. 9.

Table 2: The average PSNR (db) comparisions of imagery SR on RSSCN7.

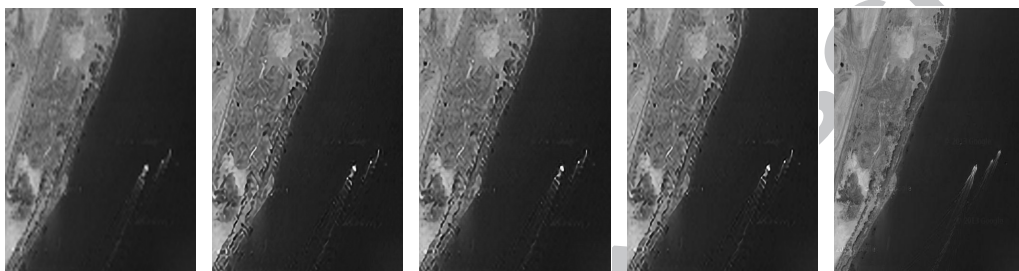| Bicubic | SRCNN | Patric | Our multi-scale SR |
|---------|-------|--------|--------------------|
| 27.85   | 28.63 | 28.86  | 29.07              |

From Tabel 2 and Fig. 8, we can easily see that, our multi-scale SR approach can get superior super-resolved results even at different imagery scenario compared to SRCNN and Patrick's one not only in visual effect but also in PSNR value. Meanwhile, from the Fig. 9, it shows that the colorization effect of Zhang's method is too saturated and unnatural whereas Iizuka's and Larsson's are too light and almost equivalent to without colorization. Obviously, compared to these colorization methods, our colorization approach can get more natural and appropriate colorization effects on the whole, though which may be different from the groundtruth ones.

In addition, for fair play we also finetune Zhang's method (its visual performance ranks second in Fig. 9) with satellite data and compare the corresponding colorization results with ours. Some comparisons are shown in Fig. 10. From the figure, it is clear that even the fine-tuned Zhang's model still fails to provide acceptable colorization effect (color is monotonous or biased), whereas the proposed multi-task approach is always able to get satisfactory results.
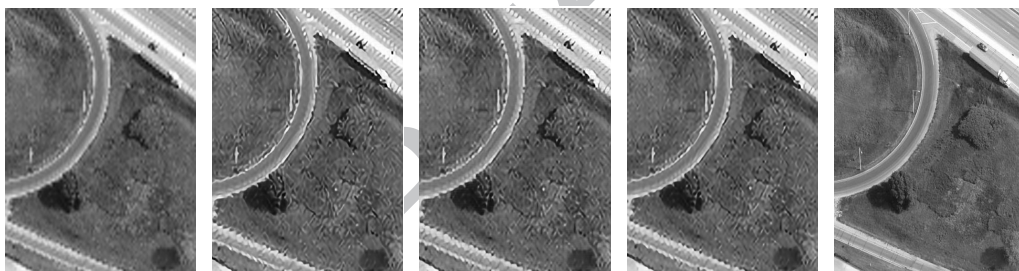
To sum up, our multi-task imagery SR and colorization approach can not only provide subtle imagery details but also make the overall color style be coordinated and natural to visual sensation. For many applications, such as those in image synthesis, the ultimate test of colorization and super-resolution is how compelling the colors and the resolution look to a human observer. Thus, from the perspective of human perception, we also introduce

22

(a) Bicubic    (b) SRCNN:24.20 (c) Patric's:24.33    (d) Ours:24.42    (e) Industry

(a) Bicubic    (b) SRCNN:28.35 (c) Patric's:29.06    (d) Ours:29.30    (e) Riverlake
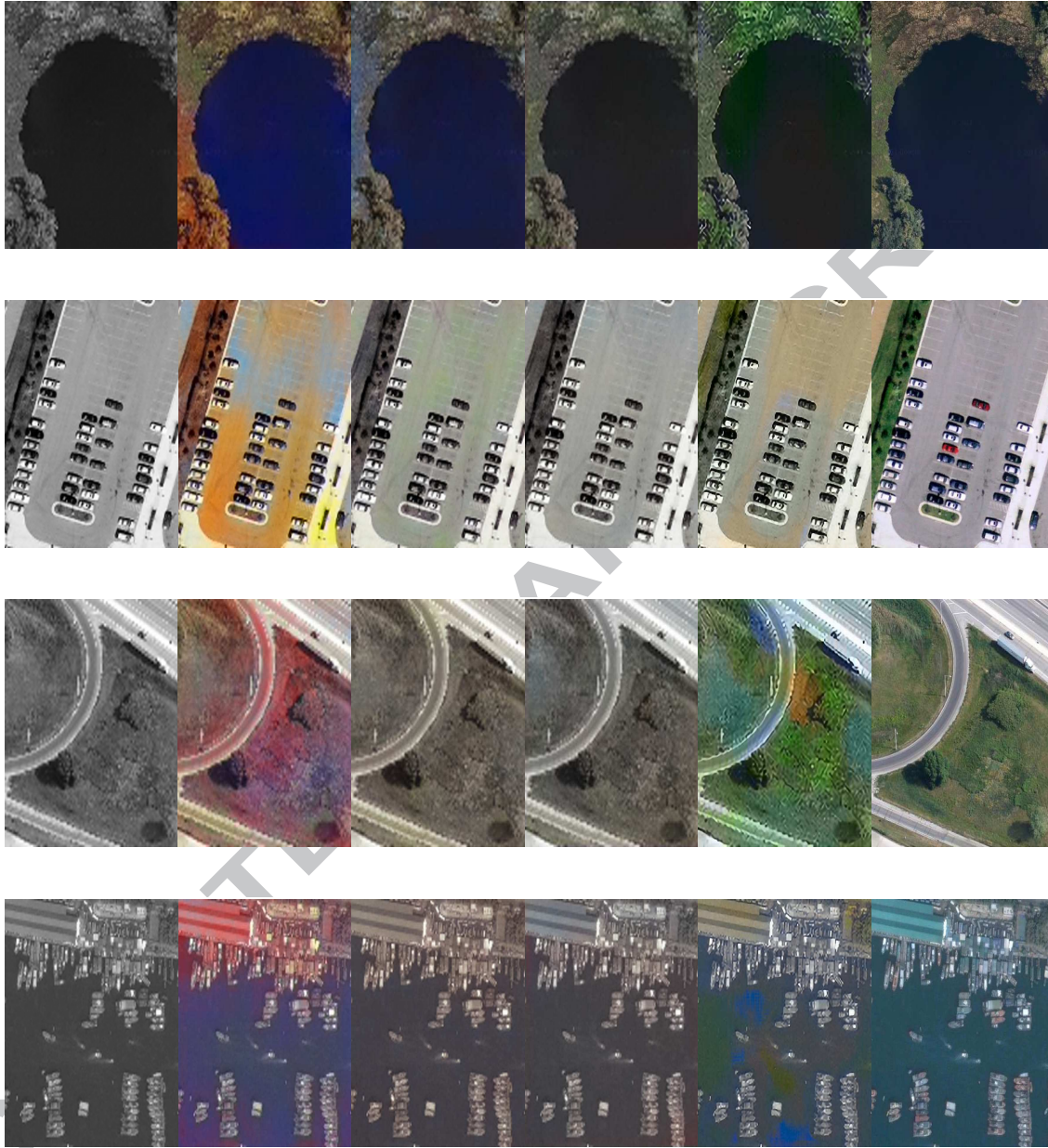
(a) Bicubic    (b) SRCNN:24.19 (c) Patric's:24.75    (d) Ours:24.95    (e) Grass
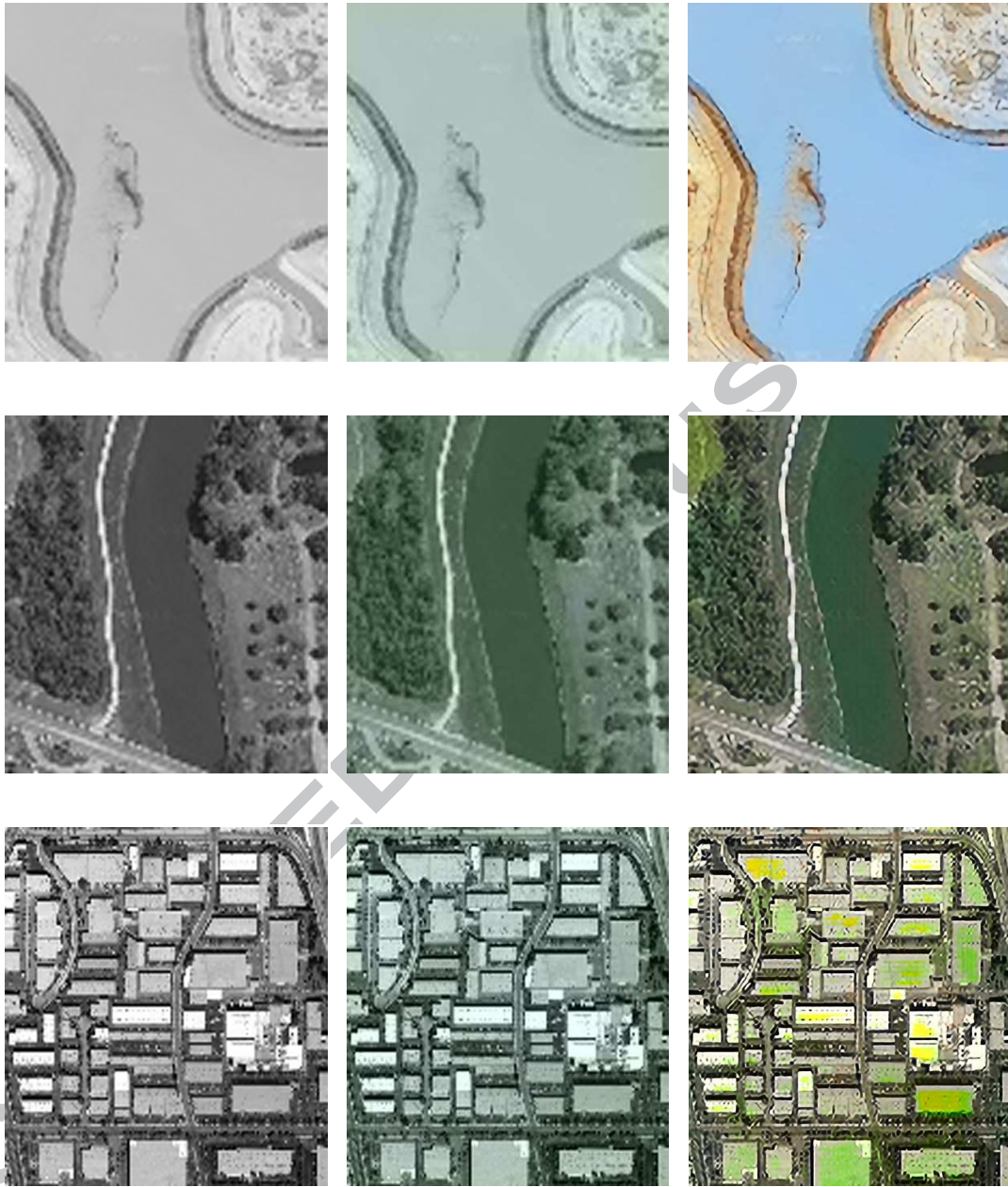
(a) Bicubic    (b) SRCNN:21.17 (c) Patric's:21.55    (d) Ours:22.07    (e) Airplane

Figure 8: Visual and PSNR (db) comparisons of super-resolved ($4\times$) images for 'Industry', 'Riverlake', 'Grass', and 'Airplane' grayscale satellite imagery from RSSCN7 by (a)Bicubic, (b)SRCNN, (c)Patric's method, and (d)Our multi-scale SR method, respectively.

23

(a) Input      (b) Zhang's      (c) Iizuka's      (d) Larsson's      (e) Ours      (f) Groundtruth

Figure 9: Visual comparisons of satellite imagery colorization for 'Riverlake' from RSSCN7, 'Parking' from AID, 'Grass' from RSSCN7, 'Port' from AID: (a)Input LR grayscale; (b)Colorization by Zhang's; (c)Colorization by Iizuka's; (d)Colorization by Larrsson's; (e)Colorization by the proposed multi-task approach; (f)Groundtruth imagery.

(a) LR grayscale      (b) Finetued Zhang's model      (c) Our method

Figure 10: Colorization comparisons between ours and the finetuned Zhang's model for 'Riverlak' from RSSCN7, 'Industrial' from AID: (a) Input LR grayscale; (b) Imagery colorization with finetuned Zhang's model; (c) Imagery colorization using the proposed multi-task approach.

25

subjective evaluation measure to show the performance of our multi-task approach. We ran a real vs. fake two-alternative forced choice experiment on campus. Totally 30 people participated in such survey and they were shown eight pairs of satellite imageries from RSSCN7 and AID, which contain natural scene -river or lake and military sensitive images - airport or parking lots. Each pair consisted of a satellite imagery next to a re-colorized and super-resolved version, produced by either our algorithm or others. Participants were asked to discriminate the imageries and choose the one they believed contained fake colors or resolution generated by a computer program and the comparisons. Each experimental session contains eight tests (each test for only one algorithm besides ground truth: four tests for colorization and three for super-resolution ) and the result of each choice is recorded and no feedback was given during all eight test pairs. To ensure that all algorithms were tested in equivalent conditions (i.e. time of day, demographics, etc.), all experiment sessions were posted simultaneously and distributed to campus in an i.i.d. fashion. These satellite imagery subjective results are shown in Table 3.

To check that participants do understand the connotation of the task, additional experimental tests were carried out - the two images of each pair were both derived from random baseline described above. Participants successfully identified these random synthesis as fake 91% of the time, indicating that they understood the task and were paying attention. The ground truth satellite imagerires are 'd162','d164','d023','d294','a007' from RSSCN7 and 'ariport228','airport108','parking176' from AID. We also compare the average PSNR value of such eight super-resolved imageries in Table 3.

26

Table 3: Satellite imagery colorization and SR subjective results.

| Method | | Model | | PSNR(db) | Labeled Real(%) |
|---|---|---|---|---|---|
| SR | Colorization | Params(MB) | Runtime(ms) | | |
| Ground Truth | | – | – | – | 47 |
| Random | | – | – | – | 9.0 |
| SRCNN(Dong et al., 2016) | | 0.3(mat file) | 115 | 24.31 | 18.1 |
| Patric's(Patrick, 2016)* | | 13.6 | 242 | 24.78 | 19.2 |
| Our multi-scale SR | | 1.1 | 141 | **25.05** | **20.8** |
| | Zhang's(Zhang et al., 2016) | 128.9 | 570 | – | 26.6 |
| | Iizuka's(Iizuka et al., 2016) | 694.7 | 360 | – | 24.5 |
| | Larsson'sLarsson et al. (2016) | 516.0 | 440 | – | 25.2 |
| | Our imagery clorization | 129.0 | 570 | – | **28.4** |
| Our multi-task SR and Colorization | | 131.6 | 390 | **25.05** | **29.7** |

*: We realize and train its caffe version.

⁴⁴⁰ From the table, it is clear that our multi-task approach fooled partici-
⁴⁴¹ pants on about 30% of tests, which is significantly higher than all compared
⁴⁴² imagery colorization or SR algorithms. These results validate the effective-
⁴⁴³ ness and applicability of the proposed multi-task model for satellite imagery
⁴⁴⁴ simultaneous colorization and SR. In addition, it is interesting to catch that
⁴⁴⁵ image color perhaps plays more important role than the resolution when we
⁴⁴⁶ try to perceive satellite imageries visually.

## 4. Conclusions

⁴⁴⁸ In this work, for satellite imagery virtual reality applications, by present-
⁴⁴⁹ ing a novel multi-task deep learning model, we have achieved simultaneous
⁴⁵⁰ satellite imagery SR and colorization. The proposed multi-scale SR deep
⁴⁵¹ structure can reconstruct LR imagery with high-frequency details and the
⁴⁵² given imagery colorization engine can efficiently recover realistic color im-
⁴⁵³ agery for a grayscale input. Through features interaction of different task
⁴⁵⁴ networks and simultaneous optimization, the experimental results and com-

27

455 parisons based on the satellite imagery data sets show that the proposed
456 multi-task approach outperforms the state-of-the-art methods and will get
457 better imagery SR and colorization effect.

458 Future work will focus on two aspects: introducing satellite image clas-
459 sification (Gong et al., 2017) structure for multi-task learning; investigating
460 the possibility of applying our multi-task deep neural model to other applica-
461 tions, such as saliency detection (Zhang et al., 2017b,c), image retrieval (Guo
462 et al., 2017; Lin et al., 2017) and activity recognition (Zhang et al., 2017a;
463 Han et al., 2012).

## Acknowledgment

## Conflict of Interest

470 The authors declare that there is no conflict of interest.

## Reference

472 Alvarez-Ramos, V., Ponomaryov, V., Reyes-Reyes, R., Gallegos-Funes, F.,
473 2016. Satellite image super-resolution using overlapping blocks via sparse
474 representation. In: Physics and Engineering of Microwaves, Millimeter and
475 Submillimeter Waves (MSMW), 2016 9th International Kharkiv Sympo-
476 sium on. IEEE, pp. 1–4.

28

477 Brodu, N., 2016. Super-resolving multiresolution images with band-
478 independant geometry of multispectral pixels. arXiv preprint
479 arXiv:1609.07986.

480 Cheng, Z., Yang, Q., Sheng, B., 2015. Deep colorization. In: Proceedings of
481 the IEEE International Conference on Computer Vision. pp. 415–423.

482 Dong, C., Loy, C. C., He, K., Tang, X., 2014. Learning a deep convolutional
483 network for image super-resolution. In: European Conference on Computer
484 Vision. Springer, pp. 184–199.

485 Dong, C., Loy, C. C., He, K., Tang, X., 2016. Image super-resolution using
486 deep convolutional networks. IEEE Transactions on Pattern Analysis and
487 Machine Intelligence 38 (2), 295–307.

488 Gong, C., Han, J., Lu, X., 2017. Remote sensing image scene classification:
489 Benchmark and state of the art. Proceedings of the IEEE 105 (10), 1865–
490 1883.

491 Guo, Y., Ding, G., Liu, L., Han, J., Shao, L., 2017. Learning to hash with
492 optimized anchor embedding for scalable retrieval. IEEE Transactions on
493 Image Processing 26 (3), 1344–1354.

494 Han, J., Pauwels, E., de Zeeuw, P., de With, P., 2012. Employing a rgb-
495 d sensor for real-time tracking of humans across multiple re-entries in a
496 smart environment. IEEE Transactions on Consumer Electronics 58 (2),
497 255–263.

498 He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image

29

499 recognition. In: Proceedings of the IEEE Conference on Computer Vision
500 and Pattern Recognition. pp. 770–778.

501 Hong, S., Noh, H., Han, B., 2015. Decoupled deep neural network for semi-
502 supervised semantic segmentation. In: Advances in Neural Information
503 Processing Systems. pp. 1495–1503.

504 Hung, L. Q., Giang, D. T., Quang, N. N., 2016. Superresolution method
505 approach for vietnam remote sensing imagery. International Journal of
506 Remote Sensing Applications 6 (0), 118–126.

507 Iizuka, S., Simo-Serra, E., Ishikawa, H., 2016. Let there be color!: joint end-
508 to-end learning of global and local image priors for automatic image col-
509 orization with simultaneous classification. ACM Transactions on Graphics
510 (TOG) 35 (4), 110.

511 Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R.,
512 Guadarrama, S., Darrell, T., 2014. Caffe: Convolutional architecture for
513 fast feature embedding. In: Proceedings of the 22nd ACM international
514 conference on Multimedia. ACM, pp. 675–678.

515 Kim, J., Kwon Lee, J., Mu Lee, K., 2016. Accurate image super-resolution
516 using very deep convolutional networks. In: Proceedings of the IEEE Con-
517 ference on Computer Vision and Pattern Recognition. pp. 1646–1654.

518 Kirkpatrick, S., Gelatt, C. D., Vecchi, M. P., et al., 1983. Optimization by
519 simulated annealing. science 220 (4598), 671–680.

520 Larsson, G., Maire, M., Shakhnarovich, G., 2016. Learning representations

30

521 for automatic colorization. In: European Conference on Computer Vision.
522 Springer, pp. 577–593.

523 Liang, Y., Wang, J., Zhou, S., Gong, Y., Zheng, N., 2016. Incorporating
524 image priors with deep convolutional neural networks for image super-
525 resolution. Neurocomputing 194, 340–347.

526 Liebel, L., Körner, M., 2016. Single-image super resolution for multispectral
527 remote sensing data using convolutional neural networks. In: XXIII ISPRS
528 Congress proceedings. pp. 883–890.

529 Lin, Z., Ding, G., Han, J., Wang, J., 2017. Cross-view retrieval via
530 probability-based semantics-preserving hashing. IEEE Transactions on Cy-
531 bernetics 47 (12), 4342–4355.

532 Mallat, S., 1999. A wavelet tour of signal processing. Academic press.

533 Patrick, H., 2016. Super-resolution on satellite imagery using
534 deep learning part 1. https://medium.com/the-downlinq/
535 super-resolution-on-satellite-imagery-using-deep-learning-part-1,
536 accessed April 29, 2017.

537 Pickup, L. C., 2007. Machine learning in multi-frame image super-resolution.
538 Oxford University.

539 Ren, W., Liu, S., Zhang, H., Pan, J., Cao, X., Yang, M.-H., 2016. Single im-
540 age dehazing via multi-scale convolutional neural networks. In: European
541 Conference on Computer Vision. pp. 154–169.

Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al., 2015. Imagenet large scale visual recognition challenge. International Journal of Computer Vision 115 (3), 211–252.

Shen, W., Zhao, K., Jiang, Y., Wang, Y., Zhang, Z., Bai, X., 2016. Object skeleton extraction in natural images by fusing scale-associated deep side outputs. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 222–230.

SpaceNet, 2016. https://amazonaws-china.com/cn/public-datasets/spacenet/, accessed April 30, 2017.

Szegedy, C., Ioffe, S., Vanhoucke, V., Alemi, A., 2016. Inception-v4, inception-resnet and the impact of residual connections on learning. arXiv preprint arXiv:1602.07261.

Xia, G.-S., Hu, J., Hu, F., Shi, B., Bai, X., Zhong, Y., Zhang, L., Lu, X., 2017. Aid: A benchmark data set for performance evaluation of aerial scene classification. IEEE Transactions on Geoscience and Remote Sensing.

Xie, S., Tu, Z., 2015. Holistically-nested edge detection. In: European Conference on Computer Vision. pp. 1395–1403.

Yang, J., Huang, T., 2010. Image super-resolution: Historical overview and future challenges. In: Milanfar, P. (Ed.), Super-resolution imaging. CRC Press, pp. 1–33.

Yang, J., Price, B., Cohen, S., Lee, H., Yang, M.-H., 2016. Object contour

564  detection with a fully convolutional encoder-decoder network. In: Proceed-
565  ings of the IEEE Conference on Computer Vision and Pattern Recognition.
566  pp. 193–202.

567  Zhang, B., Yang, Y., Chen, C., Yang, L., Han, J., Shao, L., 2017a. Ac-
568  tion recognition using 3d histograms of texture and a multi-class boosting.
569  IEEE Transactions on Image Processing 26 (10), 4648–4660.

570  Zhang, D., Han, J., Jiang, L., Ye, S., Chang, X., 2017b. Revealing event
571  saliency in unconstrained video collection. IEEE Transactions on Image
572  Processing 26 (4), 1746–1758.

573  Zhang, D., Meng, D., Han, J., 2017c. Co-saliency detection via a self-paced
574  multiple-instance learning framework. IEEE Trans. on Pattern Analysis
575  and Machine Intelligence 39 (5), 865–878.

576  Zhang, H., Yang, Z., Zhang, L., Shen, H., 2014. Super-resolution reconstruc-
577  tion for multi-angle remote sensing images considering resolution differ-
578  ences. Remote Sensing 6 (1), 637–657.

579  Zhang, R., Isola, P., Efros, A. A., 2016. Colorful image colorization. In:
580  European Conference on Computer Vision. Springer, pp. 649–666.

581  Zhu, H., Song, W., Tan, H., Wang, J., Jia, D., 2016. Super resolution recon-
582  struction based on adaptive detail enhancement for zy-3 satellite images.
583  In: ISPRS Annals of Photogrammetry, Remote Sensing and Spatial Infor-
584  mation Sciences. Vol. III-7. pp. 213–217.

585 Zou, Q., Ni, L., Zhang, T., Wang, Q., 2015. Deep learning based feature selec-
586   tion for remote sensing scene classification. IEEE Geoscience and Remote
587   Sensing Letters 12 (11), 2321–2325.

Highlights

- We propose a multi-task deep neural model to achieve satellite imagery SR and colorization simultaneously. To the best of our knowledge, this is the first work which explores to achieve satellite imagery SR and colorization cooperatively.

- We incorporate natural images with satellite data to enrich the color diversity in imagery colorization and we manage to realize the expectation color distribution learning to avoid color bias in colorization.

- We introduce a novel multi-scale deep encoder-decoder symmetrical network for satellite imagery SR, where a residual structure is adopted to improve the imagery reconstruction performance.