POLITECNICO DI TORINO Repository ISTITUZIONALE

Synthetic Data Pretraining for Hyperspectral Image Super-Resolution

Original

Synthetic Data Pretraining for Hyperspectral Image Super-Resolution / Aiello, E.; Agarla, M.; Valsesia, D.; Napoletano, P.; Bianchi, T.; Magli, E.; Schettini, R. - In: IEEE ACCESS. - ISSN 2169-3536. - ELETTRONICO. - 12:(2024), pp. 65024-65031. [10.1109/ACCESS.2024.3396990]

Availability: This version is available at: 11583/2990365 since: 2024-07-04T13:06:17Z

Publisher: IEEE

Published DOI:10.1109/ACCESS.2024.3396990

Terms of use:

This article is made available under terms and conditions as specified in the corresponding bibliographic description in the repository

Publisher copyright

(Article begins on next page)



Received 21 March 2024, accepted 1 May 2024, date of publication 6 May 2024, date of current version 14 May 2024. Digital Object Identifier 10.1109/ACCESS.2024.3396990

RESEARCH ARTICLE

Synthetic Data Pretraining for Hyperspectral Image Super-Resolution

EMANUELE AIELLO¹, MIRKO AGARLA^{®2}, DIEGO VALSESIA^{®1}, (Member, IEEE), PAOLO NAPOLETANO^{®2}, (Member, IEEE), TIZIANO BIANCHI^{®1}, (Member, IEEE), ENRICO MAGLI^{®1}, (Fellow, IEEE), AND RAIMONDO SCHETTINI^{®2}

¹Department of Electronics and Telecommunications, Politecnico di Torino, 10129 Turin, Italy
 ²Department of Computer Science, Systems and Communication (DISCO), University of Milano-Bicocca, 20126 Milan, Italy

Corresponding author: Emanuele Aiello (emanuele.aiello@polito.it)

This work was supported in part by the Future Artificial Intelligence Research (FAIR) through the European Union Next-GenerationEU (PIANO NAZIONALE DI RIPRESA E RESILIENZA (PNRR)—MISSIONE 4 COMPONENTE 2, INVESTIMENTO 1.3—D.D. 1555, 11/10/2022) under Grant PE00000013. This manuscript reflects only the authors' views and opinions, neither the European Union nor the European Commission can be considered responsible for them.

ABSTRACT Large-scale self-supervised pretraining of deep learning models is known to be critical in several fields, such as language processing, where its has led to significant breakthroughs. Indeed, it is often more impactful than architectural designs. However, the use of self-supervised pretraining lags behind in several domains, such as hyperspectral images, due to data scarcity. This paper addresses the challenge of data scarcity in the development of methods for spatial super-resolution of hyperspectral images (HSI-SR). We show that state-of-the-art HSI-SR methods are severely bottlenecked by the small paired datasets that are publicly available, also leading to unreliable assessment of the architectural merits of the models. We propose to capitalize on the abundance of high resolution (HR) RGB images to develop a self-supervised pretraining approach that significantly improves the quality of HSI-SR models. In particular, we leverage advances in spectral reconstruction methods to create a vast dataset with high spatial resolution and plausible spectra from RGB images, to be used for pretraining HSI-SR methods. Experimental results, conducted across multiple datasets, report large gains for state-of-the-art HSI-SR methods when pretrained according to the proposed procedure, and also highlight the unreliability of ranking methods when training on small datasets.

INDEX TERMS Hyperspectral images, super resolution, synthetic data, self-supervised pretraining, spectral reconstruction.

16 I. INTRODUCTION

Hyperspectral imaging is a powerful technology that captures 17 images across a wide range of the electromagnetic spectrum, 18 revealing insights unattainable in the visible. This advanced 19 imaging technique has diverse applications, ranging from 20 medical diagnostics [1] and agricultural monitoring to 21 ensure food quality, to remote sensing for environmental 22 analysis [2], [3], as well as military applications. The 23 rich spectral information contained in hyperspectral images 24 (HSIs) enables precise material identification and analysis, 25 making it an invaluable tool in these fields. 26

The associate editor coordinating the review of this manuscript and approving it for publication was Bing Li^D.

However, the design of hyperspectral imagers faces 27 significant trade-offs. To achieve a fine spectral resolution 28 and capture a broad range of wavelengths, compromises in 29 the optical and sensor designs must be made that sacrifice 30 spatial resolution in favor of spectral resolution. Moreover, 31 the sheer amount of data produced for a hyperspectral 32 cube can pose challenges in handling, particularly when a 33 rapid frame rate is desired or in certain applications, such 34 as satellite imaging, where computational and transmission 35 resources are limited. 36

This limitation in spatial resolution has thus raised interest in hyperspectral image super-resolution (HSI-SR). Super-resolution techniques are well-established in the RGB imaging domain [4], [5], but their adaptation to the HSI domain is not straightforward. Indeed, one would like to

98

99

100

101

102

103

104

105

106

107

108

109

110

extend techniques developed for RGB images to more 42 carefully account for spatio-spectral correlation and the 43 characteristics of infrared bands. However, the primary 44 obstacle is the scarcity of high-resolution hyperspectral 45 datasets, largely due to the prohibitive costs and logistical 46 challenges in collecting such data. Even worse, different 47 instruments may capture different subsets of wavelengths, 48 rendering the creation of larger datasets as collections from 49 multiple cameras problematic. This lack of extensive, high-50 quality HSI data has slowed down the development and 51 refinement of HSI-SR methods. Most of the current work 52 focuses on the design of novel neural network architectures, 53 potentially exploiting clever priors or layer structures in their 54 operations. On the other hand, it is well known [6], [7] that 55 training on more data is often more impactful than revising 56 architectural design. Moreover, using small datasets, such as 57 the ones in the current literature, poses the risk of producing 58 unreliable scientific results when assessing the merits of a 59 design over another. 60

In the case of hyperspectral images, collecting large 61 labeled datasets (such as paired HR-LR HS images) for 62 supervised training can be prohibitive or entirely impos-63 sible, due to the lack of higher resolution cameras at the 64 desired wavelengths. This calls for the development of 65 self-supervised pretraining techniques that can leverage a 66 wealth of unlabeled data so that the small amount of labeled 67 data can be used much more effectively. While techniques 68 following this idea [8], [9] have led to robust and transferable 69 models in natural language processing as well as other fields, 70 a further complication arises with hyperspectral images, 71 i.e., the overall relative scarcity of publicly available HSI 72 products, even without demanding additional pairing with 73 higher resolution data. 74

In response to this challenge, this paper introduces an inno-75 vative approach that pivots on the creation of a large-scale 76 synthetic hyperspectral dataset. Abundant high-resolution 77 RGB data can be found on the Internet and large datasets [10], 78 [11] have already been developed for applications like RGB 79 image generation, restoration, detection, etc. At the same 80 time, spectral reconstruction techniques [12], [13], [14] 81 have recently enjoyed great success in estimating plausible 82 material spectra that extend to the infrared from visible RGB 83 images only. We thus first propose to use spectral recon-84 struction techniques to transform a large-scale RGB dataset 85 into an HSI dataset with, obviously not perfect, but plausible 86 spectral content and high spatial resolution. Then, a spatial 87 super-resolution pretext task requiring to invert an arbitrary 88 degradation model is set up as a pretraining step. Critically, 89 this does not require further data or annotations, as the LR 90 images are spatially degraded from the available ones by 91 simulating the degradation process. Finally, finetuning with 92 paired real HSI data can be performed. 93

Experimental results are conducted on multiple datasets and with three state-of-the-art methods for HSI-SR. We report large gains (up to 2dB in MPSNR) in the quality of the super-resolved images when the proposed pretraining approach is followed. Moreover we conduct an ablation experiment (Sec. IV-C) that proves that pretraining with our synthetic dataset leads to better performance than using RGB images as an auxiliary task [7]. We also would like to remark that our results raise questions about the significance of results assessing merits of neural network design obtained on small datasets. In fact, we see that pretraining on the large dataset affects the relative ranking of the state-of-the-art methods. Moreover, we argue that the large-scale pretraining technique we propose could pave the way for development of bigger and more powerful neural network models.

II. BACKGROUND AND RELATED WORK

A. HYPERSPECTRAL IMAGE SUPER RESOLUTION

Hyperspectral Image Super resolution seeks to increase 111 the spatial resolution of hyperspectral images starting from 112 low-resolution observations. Several methods have been 113 developed to solve this task under various settings. This 114 work is focused on the single hyperspectral image super-115 resolution (SHSR) setting where the LR HS image is the only 116 information available to reconstruct the HR image. This is 117 contrast with other settings in which a co-registered auxiliary 118 image with one or few bands at higher resolution is available 119 as a guide [3], [15], [16]. The SHSR task is generally more 120 interesting due to the wider applicability as it does not require 121 an auxiliary input, as well as more challenging due to its 122 highly ill-posed nature. Several approaches for SHSR have 123 been proposed over the years [7], [17], [18], [19], [20], 124 [21], starting from a pioneering work leveraging a Bayesian 125 prior [17] and, more recently, deep learning methods focused 126 on applying deep neural networks to learn a direct mapping 127 between LR inputs and HR ground truth images. Among 128 them, [22] makes use of 3D convolutions to explore both 129 spatial and spectral correlation. MCNet [23] adopts a mixed 130 convolutional module, that contains a combination of 2D and 131 3D convolutions to mine spatial features of the hyperspectral 132 image and spectral information in contrast to a more 133 computationally expensive fully 3D-convolutional model. 134 SSPSR [24] introduces a spatial-spectral prior network to 135 fully exploit the spatial information and the correlation 136 between spectra. Moreover, given that hyperspectral data 137 are very scarce and have high dimensionality the authors 138 propose to use grouped convolution to increase the training 139 stability. More recently, HSISR [7] proposes the use of RGB 140 super resolution as an auxiliary task in a multi-task training 141 framework, showing how this can be beneficial to the HSI SR 142 task. 143

B. HYPERSPECTRAL DATA SCARCITY

The single image super-resolution task is an ill-posed inverse problem that necessitates a strong prior to be effectively regularized. Traditional handcrafted priors like Bayesian approaches [25] and sparse coding [26] are increasingly being replaced by learning-based approaches and neural networks which require large amounts of data for training. This is one 150

of the main challenges in the hyperspectral domain due to the 151 inherent difficulties and cost of data acquisitions. Commonly 152 used datasets [21], [27], [28] usually have only a small 153 number of images, e.g. CAVE [27] contains 20 images for 154 training while NTIRE2020 [21] has 480 images. This limits 155 the applicability and performance of most SHSR methods 156 in real world cases, where better generalization abilities 157 could be achieved if more data were available. Recent 158 work [29] develops a novel data augmentation procedure 159 to enlarge the number of data during the training phase of 160 hyperspectral super resolution methods. On the other hand, 161 some approaches have attempted to exploit the abundance 162 of RGB images, although in a way that is different from 163 the technique proposed in this paper. Yuan et al. [30] train 164 a single-band SR network on natural images and apply it to 165 HSIs in a band wise manner to exploit the spatial correlations 166 learned on RGB data. This is clearly suboptimal as it does 167 not exploit spectral correlation, and might also be challenged 168 in learning features that are specific to each wavelength. 169 Li et al. [31] develop an RGB-induced feature modulation 170 171 network that exploits features learned from RGB datasets transferring them to the SHSR task. Subsequently, Li et al. [7] 172 proposed a multi-task approach where RGB super-resolution 173 is treated as an auxiliary task to boost the performance of the 174 SHSR task. Their method exploits the correlation between 175 176 RGB and HS image features for the super-resolution task. Our method is orthogonal and possibily complementary to 177 all the previously proposed methods and models in the SHSR 178 landscape. 179

180 C. SPECTRAL RECONSTRUCTION FROM RGB

Spectral reconstruction is the task of estimating the 181 intensity of light at wavelengths beyond those captured, 182 typically extrapolating information in infrared bands from 183 an RGB input. This task requires to model or learn 184 physically-plausible spectral signatures and to use the limited 185 information in the visible, as well as spatial clues, to guess the 186 spectrum of each pixel at the unseen wavelengths. Traditional 187 methods for this task rely on handcrafted hyperspectral 188 priors [32], [33]. More recently learning based approaches 189 ([12], [13], [14] have been used to learn a direct mapping 190 between RGB images and HS images. Among them, one of 191 the most recent and efficient methods is MST++ [14], that 192 exploits a Transformer-based architecture to process inputs 193 in a multi-scale, spectral-wise manner. The method is based 194 on a spectral-wise multi-head self attention as a basic unit, 195 building on the intuition that HSIs are spatially sparse but 196 spectrally self-similar. The model is built with a U-shaped 197 structure to exploit learned features at different granularities. 198

199 III. METHOD

In this section, we propose a method to enhance the performance of any state-of-the-art neural network for spatial super-resolution of hyperspectral images. The core idea is to pretrain the neural network with a self-supervised superresolution task on a very large dataset of synthetically generated high-resolution hyperspectral images. Since very large datasets of hyperspectral images with consistent band characteristics and high spatial resolutions do not exist, we employ spectral reconstruction techniques to convert an RGB dataset into an HSI one. Finally, finetuning with the few real HSI pairs available yields the best model. 210

A. SYNTHETIC DATA GENERATION

In this phase, we generate synthetic HSI data starting from 212 an RGB dataset by employing a spectral reconstruction 213 technique. Suppose a spectral reconstruction technique is 214 available as a function $\phi : \mathbb{R}^{H \times W \times 3} \to \mathbb{R}^{H \times W \times B}$, where 215 *B* is the desired number of bands at the target wavelengths. 216 Then, we use the spectral reconstructor ϕ on all the images 217 of a large-scale RGB dataset \mathcal{D}_{RGB} to create a synthetic HSI 218 dataset $\mathcal{D}_{HS-synth}$: 219

$$\mathcal{D}_{\text{HS-synth}} = \phi(\mathcal{D}_{\text{RGB}}) \tag{1} 22$$

211

239

The quality of the generated synthetic dataset depends on 221 the ability of the spectral reconstruction method to generate 222 physically-plausible as well as spatially-consistent spectra, 223 where each of the generated bands presents features similar 224 to those of real HSI data at the corresponding wavelength, 225 and is positively correlated with the performance of our 226 pretraining procedure. As a note, most spectral reconstruction 227 methods prioritize distortion over perception in the well-228 known tradeoff [34], leading to spectra that are on average 229 more accurate but do not lie in the distribution of real 230 spectral. It would be interesting to understand if generating 231 data prioritizing being on the real spectral distribution 232 (perception) leads to further improvements in the pretraining 233 framework of this paper, but this is currently outside the scope 234 of this paper and left as future work. 235

In the experiments presented in this paper we employ 236 the state-of-the-art MST++ [14] neural network as spectral 237 reconstructor ϕ . 238

B. PRETRAINING PROCEDURE

The procedure explained in the previous section allowed 240 the creation of a large-scale dataset of hyperspectral images 241 $\mathcal{D}_{HS-synth}$. However, $\mathcal{D}_{HS-synth}$ is just a collection of unlabeled 242 images, so its use for pretraining purposes requires a 243 definition of a suitable self-supervised pretext task from 244 which features can be learned which are useful for the 245 downstream problem our neural network model seeks to 246 solve. Since this paper addresses the downstream problem of 247 HSI-SR, we propose to use a self-supervised formulation of 248 super-resolution as a pretext task for the pretraining phase. 249 In this task, we degrade the HR synthetic HSIs with an 250 arbitrary degradation model that is similar to the degradation 251 model that generates real LR hyperspectral images from the 252 HR originals. A better match between the degradation model 253 used in the pretraining task and the degradation model of real 254 images would result in a more effective pretraining. However, 255 in general, one resorts to supervised training with paired real 256 LR-HR images because the degradation model is unknown 257 and possibly complex, so it might be difficult to approximate
it for the pretraining phase. In this work, we use a simple, but
widely used model, consisting of spatial convolution with a
lowpass kernel, and decimation by a factor *s*. In formulas:

$$I_{\lambda}^{\mathrm{LR}} = \left(K_{\lambda} * I_{\lambda}^{\mathrm{HR}} \right)_{\downarrow s} \tag{2}$$

where I_{λ}^{HR} represents a band at wavelength λ of a high-resolution image in the dataset $\mathcal{D}_{\text{HS-synth}}, K_{\lambda}$ is the filter 263 264 kernel, and I^{LR} is the low-resolution image. For simplicity, 265 one can use the bicubic interpolation kernel for K_{λ} , for all 266 bands. However, if the point spread function of the real 267 optical system is known at each wavelength, then using it 268 for K_{λ} in this pretext task would provide a better pretext task 269 and, possibly, better downstream performance. The pretext 270 task trains the neural network model with a conventional 271 regression loss, such as L1 or Charbonnier [35], between the 272 super-resolved image obtained from I^{LR} and I^{HR} . 273

We remark that using a large-scale RGB dataset with high resolution images to obtain $\mathcal{D}_{\text{HS-synth}}$ is desirable because it allows the model to learn how to restore high-frequency patterns during the pretraining phase.

278 C. FINETUNING PROCEDURE

262

Subsequent to the pretraining phase, we proceed to the 279 finetuning stage, which follows exactly the same procedure 280 that supervised training would. In this stage, the network, 281 initialized with the pretrained parameters is further trained 282 on real hyperspectral data. In general, a domain gap will exist 283 between the synthetic data and the real data in terms of image 284 features. The finetuning process adapts the network to the 285 characteristics of the real-world data. However, this operation 286 is significantly more data-efficient, as the network already 287 knows how to extract low-level features that are relevant to the 288 super-resolution task. The finetuning stage will also correct 289 discrepancies in the degradation model between the pretext 290 task and the real world. 291

292 IV. EXPERIMENTS

293 A. SETTING

294 a: MODELS AND SAMPLING

We evaluate the proposed pretraining solution on three 295 state-of-the-art methods for hyperspectral super-resolution, 296 MCNet [23], SSPSR [24] and HSISR [7]. Our study investi-297 gates super-resolution factors of $\times 4$ and $\times 8$. For $\times 4$, we train 298 on non-overlapping 64×64 pixel patches cropped from the 299 original images, while for $\times 8$ we use larger 128×128 pixel 300 patches. Both sets of patches are degraded via bicubic 301 interpolation to create their corresponding low-resolution 302 HSI counterparts. 303

304 b: DATASETS

We evaluate the state-of-the-art algorithms on three main datasets commonly used for benchmarking hyperspectral super-resolution, namely, the CAVE dataset [27], the Harvard dataset [28], and the NTIRE 2020 dataset [21] The images TABLE 1. Quantitative results (×4 super-resolution).

Dataset	Method	Pretext	Finetune	MPSNR ↑	RMSE↓	ERGAS↓
NTIRE2020	HSISR [7]	-	✓	38.9642	0.0150	2.0650
		\checkmark	-	35.1876	0.0224	3.0351
		\checkmark	\checkmark	39.8843	0.0137	1.8886
	SSPSR [24]			38.0740	0.0164	2.2539
		\checkmark	-	34.9169	0.0226	3.1501
		\checkmark	\checkmark	39.5264	0.0142	1.9592
				38.0248	0.0168	2.2834
	MCNet [23]	\checkmark	-	40.0617	0.0132	1.8379
		\checkmark	\checkmark	40.0631	0.0132	1.8368
	Bicubic			34.7401	0.0235	3.1901
	HSISR [7]	-	\checkmark	42.7645	0.0114	3.3346
		\checkmark	-	38.5010	0.0176	5.1675
		\checkmark	\checkmark	42.7746	0.0112	3.3374
	SSPSR [24]			40.9131	0.0144	4.0406
CAVE		\checkmark	-	34.9800	0.0251	7.9823
		\checkmark	\checkmark	42.2938	0.0118	3.5755
	MCNet [23]			40.7385	0.0146	4.1659
		\checkmark	-	41.0221	0.0136	4.0295
		\checkmark	\checkmark	43.5819	0.0105	3.0634
	Bicubic			38.7380	0.0185	5.2719
Harvard	HSISR [7]	-	\checkmark	40.9317	0.0132	3.0128
		\checkmark		34.7720	0.0236	5.5637
		\checkmark	\checkmark	40.1527	0.0130	2.9041
	SSPSR [24]			40.3209	0.0142	3.2274
		\checkmark		33.4518	0.0327	7.7063
		\checkmark	\checkmark	39.9613	0.0132	2.9660
	MCNet [23]			40.1873	0.0147	3.2606
		\checkmark		38.8096	0.0151	3.3738
		\checkmark	\checkmark	40.3471	0.0127	2.8224
	Bicubic			38.8975	0.0167	3.8069

in the CAVE and NTIRE 2020 datasets consist of 31 bands 309 spanning from 400 nm to 700 nm, with intervals of 10 nm. 310 The images in the Harvard dataset consist of 31 bands but 311 range from 420 nm to 720 nm. The CAVE dataset comprises 312 32 images, each with dimensions of 512×512 pixels. For 313 the evaluation, we allocate 20 images for training and 10 for 314 testing. Regarding the Harvard dataset, it comprises a total of 315 50 images, with 40 allocated for training and 10 for testing. 316 The NTIRE 2020 dataset consists of 480 images, we assign 317 400 for training and 80 for testing. 318

For the super-resolution pretraining task, we employ a subset of the Large Scale Dataset for Image Restoration (LSDIR) [6]. The dataset is composed of 87,141 RGB images, where we randomly select 20,000 and 5,000 images for the train and test set, respectively. The images are resized to match the resolution of 512×512 pixels. The synthetic HSI dataset $\mathcal{D}_{\text{HS-synth}}$ is obtained following the procedure presented in III-A.

c: EVALUATION METRICS

To assess the performance of all methods, we employ three commonly used metrics: Root Mean Squared Error (RMSE), which measures the average squared difference between 330

319

320

321

322

323

324

325

326



FIGURE 1. Mean Absolute Error visualization for different methods with and without the proposed pretraining strategy on an NTIRE2020 test image (RGB bands shown on the left). The first row shows baseline methods (left-to-right MPSNR: 40.06 dB, 39.71 dB, 40.12 dB), while the second row shows synthetic pretraining followed by finetuning (left-to-right MPSNR: 40.66 dB, 40.76 dB, 40.64 dB).



FIGURE 2. Visualization of spectra of three pixels from a super-resolved image from the NTIRE2020 test set. Ground truth: continuous line, Baseline: dashed line, Pretraining+Baseline (ours): dotted line. Best viewed zoomed.

predicted and actual values: 331

332

339

$$\text{RMSE} = \sqrt{\frac{1}{NB} \sum_{i=1}^{N} \sum_{\lambda=1}^{B} (I_{i,\lambda}^{\text{true}} - I_{i,\lambda}^{\text{pred}})^2}; \quad (3)$$

N is the total number of pixels in each of the B bands, $I_{i,\lambda}^{\text{true}}$ 333 and $I_{i,\lambda}^{\text{pred}}$ are the values of the *i*-th pixel in the λ -th band for 334 the ground truth and predicted images, respectively; 335

Erreur Relative Globale Adimensionnelle de Synthese 336 (ERGAS), a dimensionless indicator of overall reconstruction 337 error frequently used in HSI fusion: 338

$$\text{ERGAS} = 100s \sqrt{\frac{1}{B} \sum_{\lambda=1}^{B} \left(\frac{\text{RMSE}_{\lambda}}{\mu_{\lambda}}\right)^{2}}; \qquad (4)$$

 $RMSE_{\lambda}$ is the RMSE for each band, s represents the 340 upsampling factor (e.g., 4 for \times 4 upsampling) and μ_{λ} is the 341 mean value for the spectral band. 342

Multi-scale Peak Signal-to-Noise Ratio (MPSNR) pro-343 vides a composite measure of the reconstruction fidelity. 344

PSNR_{$$\lambda$$} = 10 log₁₀ $\left[\frac{MAX_{\lambda}}{MSE_{\lambda}}\right]$; (5)

where MAX $_{\lambda}$ is the maximum possible value in band λ (e.g., 255 for 8-bit images). The MPSNR is the average of the 347 $PSNR_{\lambda}$ over all bands.

d: IMPLEMENTATION DETAILS

In the pretraining stage, we follow the author's implementation of each method using our synthetically generated LSDIR dataset. We train each model for 4 epochs and select the model with the lowest RMSE on the validation set. Then, the pretrained model is used as the starting configuration for the next training phase that involves the three selected datasets. For this phase, we still use the authors' implementations for all the methods.

The original version of the HSISR method exploits auxiliary RGB images and semi-supervised learning, as described by the authors [7]. For our experiments in Sec. IV-B, we keep the procedure for the baseline HSISR assessment, while we remove it when we use the proposed pretraining.

B. EXPERIMENTAL RESULTS

We evaluate state-of-the-art methods in the $\times 4$ and $\times 8$ super-364 resolution setups, presenting the results of each model 365

346

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

389

390

391

392

393

394

395

396

397

398

399

400

401

407

TABLE 2. Quantitative results (×8 super-resolution).

Dataset	Method	Pretext	Finetune	MPSNR ↑	RMSE↓	ERGAS↓
		-	~	33.4557	0.0263	3.8437
	HSISR [7]	\checkmark	-	31.3115	0.0342	4.7410
		\checkmark	\checkmark	33.4447	0.0279	3.8028
				31.7896	0.0326	4.4952
NTIRE2020	SSPSR [24]	\checkmark	-	30.6348	0.0367	5.0694
		\checkmark	\checkmark	33.2168	0.0285	3.8903
	MCNet [23]			31.9629	0.0327	4.4169
		\checkmark	-	31.6321	0.0336	4.6053
		\checkmark	\checkmark	33.4515	0.0279	3.7922
	Bicubic			29.9589	0.0396	5.4594
	HSISR [7]	-	\checkmark	37.3532	0.0206	6.0027
		\checkmark	-	35.0913	0.0248	7.9551
		\checkmark	\checkmark	37.7347	0.0197	5.8021
	SSPSR [24]		~~~	35.8896	0.0248	7.0394
CAVE		\checkmark	-	34.5132	0.0269	8.0961
CAVE		\checkmark	\checkmark	37.6007	0.0202	5.9358
	MCNet [23]	-	~~~~	34.3116	0.0280	10.2985
		\checkmark	-	35.3778	0.0228	5.0354
		\checkmark	\checkmark	37.8668	0.0198	5.7969
	Bicubic			34.2221	0.0304	8.4350
		-	\checkmark	37.3546	0.0201	4.5448
Harvard	HSISR [7]	\checkmark	-	33.8785	0.0263	6.4212
		\checkmark	\checkmark	36.1885	0.0208	4.5457
	SSPSR [24]			36.4563	0.0228	4.9978
		\checkmark	-	33.6097	0.0266	6.5786
		✓	√	35.9873	0.0212	4.6509
	MCNet [23]		_	36.3921	$0.0\overline{2}3\overline{4}$	5.0572
		\checkmark	-	35.3778	0.0228	5.0354
		\checkmark	\checkmark	36.3761	0.0203	4.4320
	Bicubic			35.7409	0.0249	5.4772

both with and without pretraining using our synthetic data, 366 followed by finetuning on the target dataset. Table 1 and 367 Table 2 present the results for the $\times 4$ and $\times 8$ scenarios, 368 respectively. For each model and dataset, three experiments 369 are reported: i) "finetune only" is the baseline, i.e., the 370 model as published in the literature; ii) "pretext only" is 371 when only the pretraining phase on the synthetic dataset 372 is performed without finetuning on the target dataset; iii) 373 "pretext+finetuning" is the full method with pretraining on 374 synthetic data and finetuning on the target dataset. 375

We can first notice that the domain gap between the 376 synthetic and real datasets can, in general, limit the per-377 formance of using only the pretraining approach without 378 finetuning, albeit some cases (e.g., MCNet on NTIRE2020) already report an improvement over the baseline. In general, 380 pretraining on the synthetically generated data followed by 381 finetuning provides the best results, sometimes with large 382 margins, only occasionally not reporting an improvement 383 over all the three metrics. 384

We can also notice that, while HSISR [7] is generally 385 considered the state-of-the-art approach, providing the best 386 results in the baseline setting, this is no longer true after 387 large-scale pretraining. Indeed, MCNet after pretraining and 388

TABLE 3. Impact of the number of pretraining data, on HSISR finetuned with NTIRE2020 dataset (x4 Super-Resolution).

# pretext images	MPSNR ↑	RMSE↓	ERGAS↓
2500	39.7836	0.01381	1.9134
5000	39.8645	0.01374	1.9021
10000	39.8777	0.01366	1.8933
20000	39.8843	0.01369	1.8886

TABLE 4. Effectiveness of auxiliary training tasks (HSISR, × 4 super-resolution).

Task	MPSNR ↑	RMSE↓	ERGAS↓
None	38.3149	0.0154	2.2069
RGB-SR+SSL [7]	38.9642	0.0150	2.0650
Proposed	39.8843	0.0137	1.8886

finetuning seem to display the best overall performance. This points out a limitation of the current literature in assessing the merits of model design on small datasets, which may lead to unreliable results, as we demostrate.

Fig. 1 reports a visual comparison for one non-cherrypicked image from the NTIRE2020 test set. We visualize the mean absolute error for the different methods with and without the proposed pretraining strategy. Moreover, in Figure 2 we plot the spectra of randomly selected pixels in the super-resolved image for each method. The proposed pretraining approach yields models that are able to more faithfully reproduce the original spectrum.

C. ABLATION STUDIES

In this section, we first study the impact of the amount of 402 synthetic data used during the proposed pretraining stage. 403 For this experiment we pretrain the same model (HSISR [7]) 404 with a variable number of synthetic data and we finetune 405 each pretrained model on NTIRE2020 dataset, results are 406 reported in Table 3. Our experiments show that increasing the number of data improves the performance with a diminished 408 return over 10K synthetic images. We hypothesize that this 409 may be due to the limited representational capacity of current 410 architectures, being designed to work with a smaller amount 411 of data. 412

Then, we evaluate the effectiveness of the proposed 413 pretraining strategy vis-à-vis an alternative approach using 414 RGB images as an auxiliary task, i.e., the procedure used 415 in [7]. Table 4 shows the performance of the HSISR 416 architecture under three different conditions. First, training 417 without any auxiliary task reports the worst performance 418 across all metrics. The semi-supervised procedure with 419 auxiliary RGB images of [7] improves performance (about 420 +0.6 dB in MPSNR), but it can be noticed that the proposed 421 pretraining strategy is the most effective (about +1.5 dB 422 improvement in MPSNR). 423

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

424 V. CONCLUSION AND DISCUSSION

In this study, we have demonstrated the significant impact 425 of large-scale synthetic data pretraining in the realm of 426 hyperspectral image super-resolution. Our approach, lever-427 ages models for spectral reconstruction to create a large HSI 428 dataset from RGB images. When employed for a pretraining 429 phase with a suitable pretext task, large improvements in 430 the quality of super-resolved images have been observed 431 on a number of datasets and state-of-the-art models. This 432 work not only presents a viable solution to the data 433 limitation in HSI SR but also sets a precedent for future 434 research in synthetic hyperspectral data. We hope that our 435 methodology will inspire further exploration and innovative 436 applications in the field of hyperspectral imaging, extending 437 beyond super-resolution tasks to a broader spectrum of 438 problems. 439

440 ACKNOWLEDGMENT

(Emanuele Aiello and Mirko Agarla contributed equally to
 this work.)

443 **REFERENCES**

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

- G. Lu and B. Fei, "Medical hyperspectral imaging: A review," *J. Biomed. Opt.*, vol. 19, no. 1, Jan. 2014, Art. no. 010901.
- [2] M. J. Khan, H. S. Khan, A. Yousaf, K. Khurshid, and A. Abbas, "Modern trends in hyperspectral image analysis: A review," *IEEE Access*, vol. 6, pp. 14118–14129, 2018.
- L. Wu and D.-W. Sun, "Advanced applications of hyperspectral imaging technology for food quality and safety analysis and assessment: A review—Part I: Fundamentals," *Innov. Food Sci. Emerg. Technol.*, vol. 19, pp. 1–14, Jul. 2013.
- [4] J. Liang, J. Cao, G. Sun, K. Zhang, L. Van Gool, and R. Timofte, "SwinIR: Image restoration using Swin Transformer," in *Proc. IEEE/CVF Conf. Comput. Vis. Workshops (ICCVW)*, Montreal, BC, Canada, Oct. 2021, pp. 1833–1844.
 - [5] S. W. Zamir, A. Arora, S. Khan, M. Hayat, F. S. Khan, and M. Yang, "Restormer: Efficient transformer for high-resolution image restoration," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2022, pp. 5718–5729.
 - [6] Y. Li, K. Zhang, J. Liang, J. Cao, C. Liu, R. Gong, Y. Zhang, H. Tang, Y. Liu, D. Demandolx, R. Ranjan, R. Timofte, and L. Van Gool, "LSDIR: A large scale dataset for image restoration," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW)*, Jun. 2023, pp. 1775–1787.
 - [7] K. Li, D. Dai, and L. Van Gool, "Hyperspectral image super-resolution with RGB image super-resolution as an auxiliary task," in *Proc. IEEE/CVF Winter Conf. Appl. Comput. Vis. (WACV)*, Jan. 2022, pp. 4039–4048.
 - [8] R. Balestriero, M. Ibrahim, V. Sobal, A. Morcos, S. Shekhar, T. Goldstein, F. Bordes, A. Bardes, G. Mialon, Y. Tian, A. Schwarzschild, A. G. Wilson, J. Geiping, Q. Garrido, P. Fernandez, A. Bar, H. Pirsiavash, Y. LeCun, and M. Goldblum, "A cookbook of self-supervised learning," 2023, *arXiv:2304.12210.*
 - [9] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, "Masked autoencoders are scalable vision learners," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2022, pp. 15979–15988.
- I0] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet:
 A large-scale hierarchical image database," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2009, pp. 248–255.
- [11] C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta, C. Mullis, M. Wortsman, P. Schramowski, S. Kundurthy, K. Crowson, L. Schmidt, R. Kaczmarczyk, and J. Jitsev, "LAION-5B: An open large-scale dataset for training next generation image-text models," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 35, 2022, pp. 25278–25294.
- [12] S. Galliani, C. Lanaras, D. Marmanis, E. Baltsavias, and K. Schindler,
 "Learned spectral super-resolution," 2017, *arXiv:1703.09470*.

- [13] Z. Xiong, Z. Shi, H. Li, L. Wang, D. Liu, and F. Wu, "HSCNN: CNN-based hyperspectral image recovery from spectrally undersampled projections," in *Proc. IEEE Int. Conf. Comput. Vis. Workshops (ICCVW)*, Oct. 2017, pp. 518–525.
- [14] Y. Cai, J. Lin, Z. Lin, H. Wang, Y. Zhang, H. Pfister, R. Timofte, and L. V. Gool, "MST++: Multi-stage spectral-wise transformer for efficient spectral reconstruction," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW)*, Jun. 2022, pp. 744–754.
- [15] Q. Xie, M. Zhou, Q. Zhao, Z. Xu, and D. Meng, "MHF-Net: An interpretable deep network for multispectral and hyperspectral image fusion," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 44, no. 3, pp. 1457–1473, Mar. 2022.
- [16] J.-F. Hu, T.-Z. Huang, L.-J. Deng, H.-X. Dou, D. Hong, and G. Vivone, "Fusformer: A transformer-based fusion network for hyperspectral image super-resolution," *IEEE Geosci. Remote Sens. Lett.*, vol. 19, pp. 1–5, 2022.
- [17] N. Akhtar, F. Shafait, and A. Mian, "Bayesian sparse representation for hyperspectral image super resolution," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2015, pp. 3631–3640.
- [18] J. Li, R. Cui, B. Li, R. Song, Y. Li, Y. Dai, and Q. Du, "Hyperspectral image super-resolution by band attention through adversarial learning," *IEEE Trans. Geosci. Remote Sens.*, vol. 58, no. 6, pp. 4304–4318, Jun. 2020.
- [19] O. Sidorov and J. Y. Hardeberg, "Deep hyperspectral prior: Single-image denoising, inpainting, super-resolution," in *Proc. IEEE/CVF Int. Conf. Comput. Vis. Workshop (ICCVW)*, Oct. 2019, pp. 3844–3851.
- [20] J. Jiang, C. Wang, X. Liu, K. Jiang, and J. Ma, "From less to more: Spectral splitting and aggregation network for hyperspectral face super-resolution," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW)*, Jun. 2022, pp. 266–275.
- [21] B. Arad et al., "NTIRE 2020 challenge on spectral reconstruction from an RGB image," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW)*, Jun. 2020, pp. 1806–1822.
- [22] S. Mei, X. Yuan, J. Ji, Y. Zhang, S. Wan, and Q. Du, "Hyperspectral image spatial super-resolution via 3D full convolutional neural network," *Remote Sens.*, vol. 9, no. 11, p. 1139, Nov. 2017.
- [23] Q. Li, Q. Wang, and X. Li, "Mixed 2D/3D convolutional network for hyperspectral image super-resolution," *Remote Sens.*, vol. 12, no. 10, p. 1660, May 2020.
- [24] J. Jiang, H. Sun, X. Liu, and J. Ma, "Learning spatial-spectral prior for super-resolution of hyperspectral imagery," *IEEE Trans. Comput. Imag.*, vol. 6, pp. 1082–1096, 2020.
- [25] H. Irmak, G. B. Akar, and S. E. Yuksel, "A MAP-based approach for hyperspectral imagery super-resolution," *IEEE Trans. Image Process.*, vol. 27, no. 6, pp. 2942–2951, Jun. 2018.
- [26] H. Huang, J. Yu, and W. Sun, "Super-resolution mapping via multidictionary based sparse representation," in *Proc. IEEE Int. Conf. Acoust.*, *Speech Signal Process. (ICASSP)*, May 2014, pp. 3523–3527.
- [27] F. Yasuma, T. Mitsunaga, D. Iso, and S. K. Nayar, "Generalized assorted pixel camera: Postcapture control of resolution, dynamic range, and spectrum," *IEEE Trans. Image Process.*, vol. 19, no. 9, pp. 2241–2253, Sep. 2010.
- [28] A. Chakrabarti and T. Zickler, "Statistics of real-world hyperspectral images," in *Proc. CVPR*, Jun. 2011, pp. 193–200.
- [29] N. Aburaed, M. Q. Alkhatib, S. Marshall, J. Zabalza, and H. Al Ahmad, "Hyperspectral data scarcity problem from a super resolution perspective: Data augmentation analysis and scheme," in *Proc. IEEE Int. Geosci. Remote Sens. Symp. (IGARSS)*, Jul. 2023, pp. 5057–5060.
- [30] Y. Yuan, X. Zheng, and X. Lu, "Hyperspectral image superresolution by transfer learning," *IEEE J. Sel. Topics Appl. Earth Observ. Remote Sens.*, vol. 10, no. 5, pp. 1963–1974, May 2017.
- [31] Q. Li, M. Gong, Y. Yuan, and Q. Wang, "RGB-induced feature modulation network for hyperspectral image super-resolution," *IEEE Trans. Geosci. Remote Sens.*, vol. 61, pp. 1–11, 2023, Art. no. 5512611.
- [32] B. Arad and O. Ben-Shahar, "Sparse recovery of hyperspectral signal from natural RGB images," in *Proc. 14th Eur. Conf. Comput. Vis.*, Amsterdam, The Netherlands. Cham, Switzerland: Springer, Oct. 2016, pp. 19–34.
- [33] J. Wu, J. Aeschbacher, and R. Timofte, "In defense of shallow learned spectral reconstruction from RGB images," in *Proc. IEEE Int. Conf. Comput. Vis. Workshops (ICCVW)*, Oct. 2017, pp. 471–479.

- [34] Y. Blau and T. Michaeli, "The perception-distortion tradeoff," in 559 Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., Jun. 2018, 560 pp. 6228-6237. 561
- [35] A. Bruhn, J. Weickert, and C. Schnörr, "Lucas/Kanade meets 562 Horn/Schunck: Combining local and global optic flow methods," 563 564

Int. J. Comput. Vis., vol. 61, no. 3, pp. 1-21, Feb. 2005.



EMANUELE AIELLO received the bachelor's degree in electronics and communications engineering and the master's degree (Hons.) in telecommunications from Politecnico di Torino, where he is currently pursuing the Ph.D. degree in artificial intelligence. He is also a Teaching Assistant with Politecnico di Torino. He has gained practical experience through prestigious internships, as a Research Scientist Intern at Meta. His research interest includes multimodal deep learning.



MIRKO AGARLA received the bachelor's and master's degrees in computer science from the University of Milano-Bicocca. He is currently pursuing the Ph.D. degree in artificial intelligence with Politecnico di Torino. He is also a Teaching Assistant and a Student Mentor with the University of Milano-Bicocca. He was a Research Assistant with the IDIAP Research Institute, Huawei Research, and the University of Milano-Bicocca in the field of AI for cutting-edge research with

practical applications. His primary research interests include quality control in Industry 4.0, with a focus on object detection, defect segmentation, and hyperspectral imaging. In addition, his research covers image and video processing techniques, addressing quality enhancement, super-resolution, and defenses against adversarial attacks.



DIEGO VALSESIA (Member, IEEE) received the Ph.D. degree in electronic and communication engineering from Politecnico di Torino, in 2016. He is currently an Assistant Professor with the Department of Electronics and Telecommunications (DET), Politecnico di Torino. His main research interests include the processing of remote sensing images and deep learning for inverse problems in imaging. He is a member of the EURASIP Technical Area Committee for Signal

and Data Analytics for Machine Learning and a member of the ELLIS Society. He was a recipient of the IEEE ICIP 2019 Best Paper Award and the IEEE Multimedia 2019 Best Paper Award. He is an Associate Editor of IEEE TRANSACTIONS ON IMAGE PROCESSING, for which he received the 2023 Outstanding Editorial Board Member Award.



PAOLO NAPOLETANO (Member, IEEE) has been an Associate Professor of computer science with the Department of Informatics, Systems and Communication, University of Milano-Bicocca, since 2021. His main research interests include artificial intelligence, machine and deep learning, computer vision, pattern recognition, intelligent sensors, biological signal processing, and human-machine systems. He is a member of the European Laboratory for Learning and Intelligent

Systems (ELLIS Society). He is an Associate Editor of IEEE JOURNAL OF BIOMEDICAL AND HEALTH INFORMATICS, Neurocomputing (Elsevier), IET Signal Processing, Sensors (MDPI), and Smart Cities (MDPI). He is the Chair of the IEEE CTSoc Machine Learning, Deep Learning and AI in CE (MDA) Technical Committee (TC).



TIZIANO BIANCHI (Member, IEEE) received the M.Sc. degree (Laurea) in electronic engineering and the Ph.D. degree in information and telecommunication engineering from the University of Florence, Italy, in 2001 and 2005, respectively. From 2005 to 2012, he was a Research Assistant with the Department of Electronics and Telecommunications, University of Florence. In 2012, he joined Politecnico di Torino as an Assistant Professor, where he is currently an Associate

Professor. He has authored more than 100 papers in international journals and conference proceedings. His research interests include multimedia security technologies, multimedia forensics, and the processing of remote-sensing images. He was a recipient of the IEEE Multimedia 2019 Best Paper Award and the 2021 and 2022 Best Associate Editor Award of the Journal of Visual Communication and Image Representation. He is currently an Associate Editor of IEEE TRANSACTIONS ON INFORMATION FORENSICS AND SECURITY and a Senior Area Editor of the Journal of Visual Communication and Image Representation.



ENRICO MAGLI (Fellow, IEEE) received the M.Sc. and Ph.D. degrees from the Politecnico di Torino, Italy, in 1997 and 2001, respectively. He is currently a Full Professor with Politecnico di Torino, where he leads the Image Processing and Learning Group, performing research in the fields of deep learning for image and video processing, image compression, and image forensics for multimedia and remote sensing applications. He is a fellow of the ELLIS Society for the Advancement

of Artificial Intelligence in Europe. He was a recipient of the IEEE Geoscience and Remote Sensing Society 2011 Transactions Prize Paper Award, the IEEE ICIP 2015 Best Student Paper Award (as a Senior Author), the IEEE ICIP 2019 Best Paper Award, the IEEE Multimedia 2019 Best Paper Award, and the 2010 and 2014 Best Associate Editor Award of IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS FOR VIDEO TECHNOLOGY. He is an Associate Editor of IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS FOR VIDEO TECHNOLOGY and EURASIP Journal on Image and Video Processing. He has been an IEEE Distinguished Lecturer, from 2015 to 2016.



RAIMONDO SCHETTINI is currently a Full Professor with the University of Milano-Bicocca, Italy, leading the Imaging and Vision Laboratory. With more than 30 years of experience, he has published extensively in color imaging and image processing, supervised numerous Ph.D. students, and led research projects in collaboration with prominent companies. He holds fellowships with the International Association of Pattern Recognition (IAPR), Asia-Pacific Artificial Intelligence

Association (AAIA), and International Artificial Intelligence Industry Alliance (AIIA). He is listed on Stanford University's World Ranking Scientists List. He serves as the Editor-in-Chief for the Journal of Imaging (MDPI).

Open Access funding provided by 'Politecnico di Torino' within the CRUI CARE Agreement