# LP-IOANET: EFFICIENT HIGH RESOLUTION DOCUMENT SHADOW REMOVAL

*Konstantinos Georgiadis[1]\*, M. Kerim Yucel[2]\*, Evangelos Skartados[1]\*, Valia Dimaridou[1],*
*Anastasios Drosou[1], Albert Saà-Garriga[2], Bruno Manganelli[2]*

[1] CERTH, Information Technologies Institute, Thessaloniki, Greece
[2] Samsung Research UK

## ABSTRACT

Document shadow removal is an integral task in document enhancement pipelines, as it improves visibility, readability and thus the overall quality. Assuming that the majority of practical document shadow removal scenarios require real-time, accurate models that can produce high-resolution outputs in-the-wild, we propose **L**aplacian **P**yramid with **I**nput/**O**utput **A**ttention **Net**work (**LP-IOANet**), a novel pipeline with a lightweight architecture and an upsampling module. Furthermore, we propose three new datasets which cover a wide range of lighting conditions, images, shadow shapes and viewpoints. Our results show that we outperform the state-of-the-art by a 35% relative improvement in mean average error (MAE), while running real-time in four times the resolution (of the state-of-the-art method) on a mobile device.

*Index Terms*— Shadow Removal, Document Enhancement, Super-Resolution

## 1. INTRODUCTION

The wide spread use of mobile phone cameras has made document digitization significantly practical. Using mobile phone cameras often leads to issues like distortion, blur, noise and shadows cast on documents. Document shadow removal task aims to remove shadows cast on document images in a visually pleasant manner. Despite the recent advances [1, 2, 3, 4], there are issues with existing methods. First, the majority of the existing methods do not often aim for a lightweight solution. Second, most methods do not operate at high resolutions. Third, document images come in many forms, such as text-only colorless documents and colored/figure-heavy documents, which necessitates good performance in-the-wild. A recent work [4] addresses document shadow removal in an end-to-end manner, however, it does so with a large model that does not operate at high resolutions.

The contributions of our work are as follows: we propose i) IOANet, a document shadow removal network with input/output attention for real-time operation, ii) a lightweight upsampling module that encapsulates IOANet, letting us operate at high resolutions, iii) three new datasets which cover various lighting conditions, document types and viewpoints and iv) a two-stage training pipeline that lets us leverage any

---

\* The first three authors contributed equally.

low-resolution dataset for improved generalization. Our LP-IOANet comfortably outperforms the state-of-the-art, runs in real-time on a mobile device in 4 times the resolution of the state-of-the-art. The state-of-the-art method runs out of memory, thus can not be run, even on a 24GB VRAM desktop GPU. The diagram of LP-IOANet is shown in Figure 1.
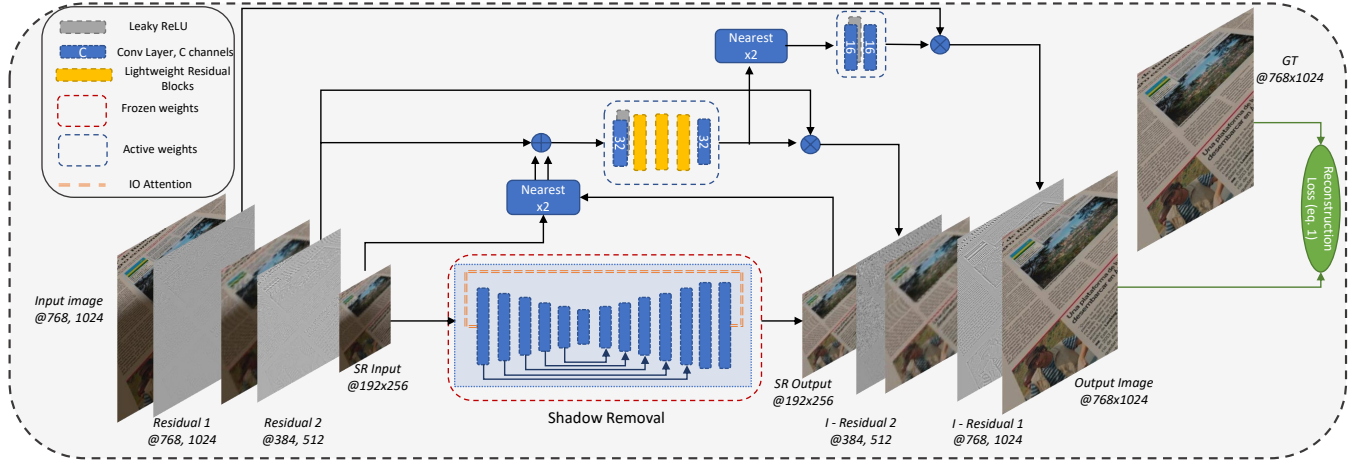
## 2. RELATED WORK

**Document Shadow Removal and Upsampling** Earlier works based on intrinsic images [5, 6] and hand-crafted methods [2, 1, 3] are built on simplifying assumptions and fail to perform effectively. Most recently, a deep-learning based solution BEDSR is proposed [4]. BEDSR assumes document images have a dominant background color, and explicitly estimates this via a background estimation module. This is used to create an attention map of background/foreground pixels, which acts as a shadow mask. Input image and the attention map are fed to an encoder/decoder to remove the shadows. On the other hand, several document specific super-resolution methods, based on CNNs [7] or GANs [8], are shown to be effective and also improve OCR performance.

**Document Shadow Removal Datasets.** Existing real-life document shadow removal datasets are quite small [1, 2, 3, 4] and thus not suitable for training purposes. Synthetically casting shadows on shadow free images is arguably more practical than creating large-scale real-life datasets, as synthetic datasets can eliminate intrinsic data errors (on non-shadow regions), simulate various lighting/occluder settings and can automate labelling process entirely. SDSRD is the largest synthetic dataset formed of 8.3K triplets, created from 970 unique images [4], however, it is not publicly available and even larger datasets are still desirable.

In comparison to the current state-of-the-art BEDSR [4], our work has several key advantages: we i) do not need the background color labels during training, ii) address high resolution output requirement explicitly, iii) enlarge data distribution via our new datasets and iv) do not use large, complex architectures. Our pipeline is also a prime candidate for optimizations via pruning/quantization as it uses simple building blocks, whereas BEDSR requires gradient information [9] in inference, which can be hard to obtain on resource-constrained environments like mobile phones. Such optimizations are interesting for future work, but not within our scope.

**Fig. 1**. Our LP-IOANet pipeline. Following the training of our shadow removal network (red dashed lines) in low-resolution (see Figure 2), we freeze it and train our lightweight upsampling module (dashed blue lines) on our proposed A-BSDD dataset.

## 3. METHODOLOGY

### 3.1. Shadow Removal

Our aim is to mimic the high-fidelity results of two-network setups (removal and refinement) [10], but by using a single, efficient network that implicitly localizes and removes shadows. We base our architecture on [11] due to its strong accuracy/runtime performance and propose to use lightweight attention modules [12, 13] over the input and output (IOA) of the network, and sum their results via a long residual connection from input to output. IOA has the advantage of being parallelizable (i.e. input attention is executed concurrently with the network), makes the network focus on the shadow areas only (i.e. non-shadow areas are copied by the long residual connection) and introduces additional capacity for blending/color-correction with minimal computational overhead. We call the resulting architecture IOANet; its architecture and training details are shown in Figure 2.

### 3.2. Efficient Upsampling

**Preliminaries.** The visual quality of the output is important in documents, since high frequency components (i.e. text) must be preserved after generation and upsampling. We set the target resolution for the shadow removal network to (192×256), which is reminiscent of the aspect ratio of documents. Instead of naively running the network at high-resolution, we aim to efficiently upsample the shadow-removed image four times to high-resolution (768×1024).

**The proposed method.** Laplacian Pyramid Networks [14] decompose an image into a Laplacian pyramid [15], where low-frequency components are fed to an image-to-image translation network in low resolution. High frequency components are adaptively refined via a mask learning network based on all frequency components. This mask is upsampled and finetuned for each resolution level, and all components are used to reconstruct the high resolution output.
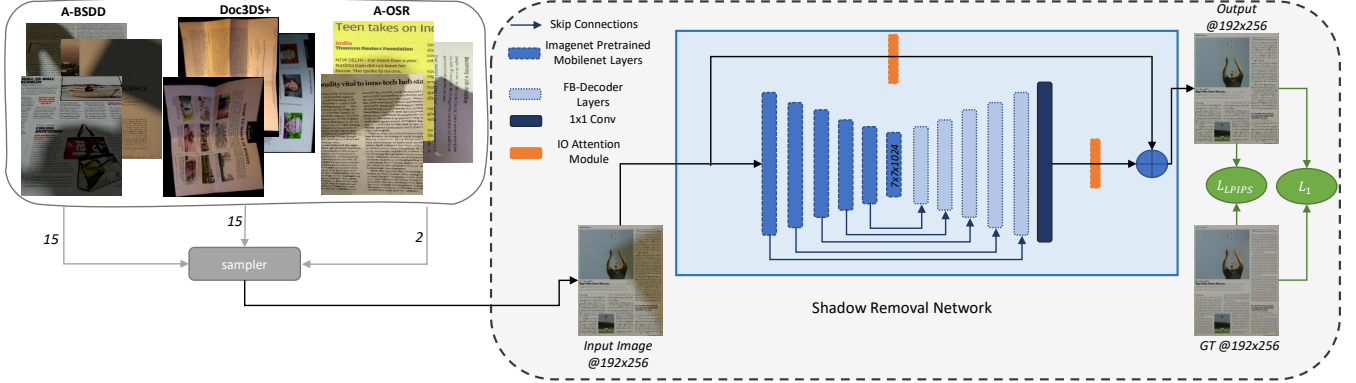
We use a 2-level pyramid where IOANet operates on low resolution (192×256) images. Unlike the original work [14], we train IOANet with low-resolution images and then the rest of the network with high-resolution images (see Section 4.1). The residual refinement network operates on the intermediate resolution of (384×512). This network, in its original version, leads our pipeline to have 22.8 GFLOPs complexity. Unlike [14], we implement the residual refinement network with cheap, depthwise separable convolutions, resulting in 3.82 GFLOPs (called LPTN-lite). We further decrease the width of the network and achieve 1.47 GFLOPs; this is our LP-IOANet pipeline. LP-IOANet is shown in Figure 1.

### 3.3. Datasets

The largest dataset in the literature (SDSRD [4]) is not publicly available, therefore we opt to create our own synthetic datasets. In addition to being able to train on a more diverse dataset, our new datasets allow us to evaluate on a larger distribution, giving us a better insight on models' performance.

**BSDD.** We follow the principles of [4] and create Blender Synthetic Dataset (BSDD) using 1328 unique images. BSDD has 3863 high resolution triplets, split into 3477 and 386 images for training and testing, respectively.

**Doc3DS+.** We leverage the Doc3DShade dataset [16], which is formed of 90K image triplets of shadow images with colored backgrounds, white balanced shadow images and albedo document images. Training a model using colored shadowed images as input and albedo as output may result in neglecting the paper color, whereas using white balanced shadowed as input may restrict the model's view to white paper only. We extract the background color from the input image using color clustering and reapply it on the albedo document images to alleviate such issues. We rotate the resulting images to be closer to our desired A4 document resolution.

**Augmenting the datasets.** A shadow area is not just a darkened version of the original image; natural shadows tend to have different colors and illuminations. Following [10, 17],

**Fig. 2**. Our IOANet shadow removal network. Using all three datasets we propose, we train our network on low-resolution images using a combination of L1 and LPIPS losses. This training stage is followed by stage 2 (see Figure 1).

| Dataset | Size | Unique Images | Resolution |
|---------|------|---------------|------------|
| A-BSDD  | 24082 | 1328 | High |
| Doc3DS+ | 71595 | 9393 | Low |
| A-OSR   | 1410 | 23 | Low |

**Table 1**. Overview of the datasets.

we apply illumination augmentation. Furthermore, we also modify the colour values of the shadows, which gives us a more diverse distribution. With this procedure, we augment the train split of BSDD and the entire OSR [18] dataset and use the augmented versions *A-BSDD* and *A-OSR* in our experiments. Details of our datasets are shown in Table 1.

## 4. EXPERIMENTAL RESULTS

### 4.1. Implementation and Training Details

We adopt a two-stage training regime, where we first train the removal network IOANet (see Figure 2) in low-resolution and then train our upsampling framework LP-IOANet with the IOANet fixed (see Figure 1). We note that a one-stage, end-to-end training is possible, but since upsampling module requires high-resolution data, one-stage training can only be done on A-BSDD dataset. Two-stage training lets us train IOANet on low-resolution data and improves the final result.

In the first stage, we train IOANet for 1000 epochs using Adam [19] with two losses; L1 and LPIPS [20]. Our loss weights are empirically chosen as 10 and 5 for L1 and LPIPS, and we use a mixed training strategy where we sample 15, 15 and 2 images in a batch from A-BSDD, Doc3DS+ and A-OSR, respectively. In the second stage, we freeze IOANet and train the upsampler using L1 loss. We train on A-BSDD for 200 extra epochs. Models are trained using PyTorch [21].

### 4.2. Evaluation Details and Results

We use PSNR, SSIM and MAE metrics on BSDD dataset for evaluation. We choose BSDD since we can train BEDSR on it (i.e. background colors are available) and because it is the only high-resolution dataset. We report metrics for all, non-shadow and shadow regions separately. We compare with the state-of-the-art BEDSR [4]; we reproduce the implementation

and train it on A-BSDD. We perform evaluation both in low-resolution ($192 \times 256$) and high-resolution ($768 \times 1024$).

Our results are shown in Table 2. We first compare IOANet, its components and BEDSR in low-resolution and train them only on BSDD. IOANet comfortably outperforms BEDSR on all metrics, despite running x10 faster, consuming x30 fewer memory and having x2K times fewer flops. It is also apparent that the input-output attention proposed in our architecture improves the results, with minimal to no overhead in runtime/memory performance. When trained on all three datasets (i.e. full stage 1 training), IOANet shows dramatic improvements, showing the value of our datasets.

In high-resolution, we follow the full two-stage training and the results are even better; BEDSR fails to run in high-resolution and goes out-of-memory on a 24GB VRAM desktop GPU. Our LP-IOANet, on the other hand, achieves comparable results to our IOANet, despite evaluating at four times the resolution. Furthermore, it is still operating comfortably in real-time (around 84 FPS) while being still faster, smaller and less complex than low-resolution BEDSR. Figure 3 shows that our method handles artefacts (1st, 2nd row) and preserves high-frequency content (3rd row), even in high resolutions. Table 3 shows that LP-IOANet reaches up to 20 FPS on mobile devices; it is faster than running IOANet directly on high-resolution, or using LPTN with IOANet.

**Further comparison.** Not many methods focus on runtime performance [3, 4, 5, 6]. Various non-ML methods [1, 2, 3, 18] have hardware-optimized implementations, but none of them utilize GPUs or can scale their results with more data. Our method leverages GPUs and scales its results significantly with more data (see Table 5).

### 4.3. Ablation Studies

**Upsampling.** Table 4 shows three upsampling solutions with different complexities. There is an expected trend here; as the upsampler becomes more complex, we get better results. The timings on desktop GPUs are not that different, but the difference becomes more visible on mobile devices, as LP-IOANet is nearly three times faster than using the original

| Method | BSDD | | | Runtime | Memory | GFLOPs |
|---|---|---|---|---|---|---|
| | MAE ↓ | PSNR ↑ | SSIM ↑ | ms | GB | |
| Input Images - No Removal | 7.2256 / 2.0526 / 24.156 | 23.98 / 13.63 / 13.56 | 0.95 | - | - | - |
| BEDSR [4] | 2.8321 / 2.1534 / 5.0535 | 34.25 / 13.58 / 13.94 | 0.98 | 119.0 | 2.30 | 550 |
| IOANet w/o attention † | 2.7845 / 2.1741 / 4.7823 | 35.66 / 13.85 / 14.05 | 0.98 | **9.6** | **0.076** | **0.25** |
| IOANet w/o attention | 2.7731 / 2.3008 / **4.3190** | 35.35 / 13.82 / **14.12** | 0.98 | **9.6** | **0.076** | **0.25** |
| IOANet | **2.3344 / 1.6414** / 4.6026 | **36.84 / 13.86** / 14.11 | **0.99** | 11.1 | **0.076** | **0.25** |
| IOANet ‡ | **1.7893 / 1.0937 / 4.0659** | 38.76 / 14.08 / 14.26 | **0.99** | 11.1 | **0.076** | **0.25** |
| BEDSR [4] | — / — / — | — / — / — | — | OOM | OOM | $\gg 1K$ |
| LP-IOANet | **1.8003 / 1.1259 / 4.0074** | **38.69 / 14.08 / 14.33** | **0.99** | 12.1 | 0.35 | 1.47 |

**Table 2**. Rows 2-5 (trained on BSDD) show the results of BEDSR, IOANet without attention and IOANet, all in *low-resolution* ($192 \times 256$). Row 6 (‡) shows IOANet trained on all three datasets (results in *low-resolution*). The last two rows show BEDSR and LP-IOANet (trained on all three datasets), all in *high-resolution* ($768 \times 1024$). † indicates training with only L1 loss. Runtime/memory are measured using an RTX 3090 GPU. Note that BEDSR runs out of memory (OOM) in high-resolution.

| Components | Resolution | Runtime (ms) |
|---|---|---|
| IOANet w/o attention | $768 \times 1024$ | 80.3 |
| IOANet | $768 \times 1024$ | 91.5 |
| LPTN [14] + IOANet | $768 \times 1024$ | 142.6 |
| LPTN-lite + IOANet | $768 \times 1024$ | 84.7 |
| LP-IOANet | $768 \times 1024$ | **57.5** |

**Table 3**. Runtime on a Samsung Galaxy S22 Ultra GPU.

| Upsampler | BSDD | Complexity | Runtime (ms) | |
|---|---|---|---|---|
| | MAE | GFLOPs | GPU | Mobile |
| LPTN [14] | **1.7148 / 1.0594 / 3.8597** | 22.8 | 15.5 | 142.6 |
| LPTN-lite | 1.7537 / 1.0784 / 3.964 | 3.82 | 14.2 | 84.70 |
| Ours (LP) | 1.8003 / 1.1259 / 4.0074 | **1.47** | **12.1** | **57.50** |

**Table 4**. Different upsampling solutions with IOANet. LPTN [14], our LPTN-lite and LP-IOANet. Runtime values are taken on the same hardware as Tables 2 and 3.

LPTN, while being as accurate. We note that all variants perform in high-resolution and comfortably outperform BEDSR.
**Datasets.** Table 5 shows the contribution of each dataset on IOANet performance. Each dataset in the mix introduces visible improvements, despite naive mixing during training.
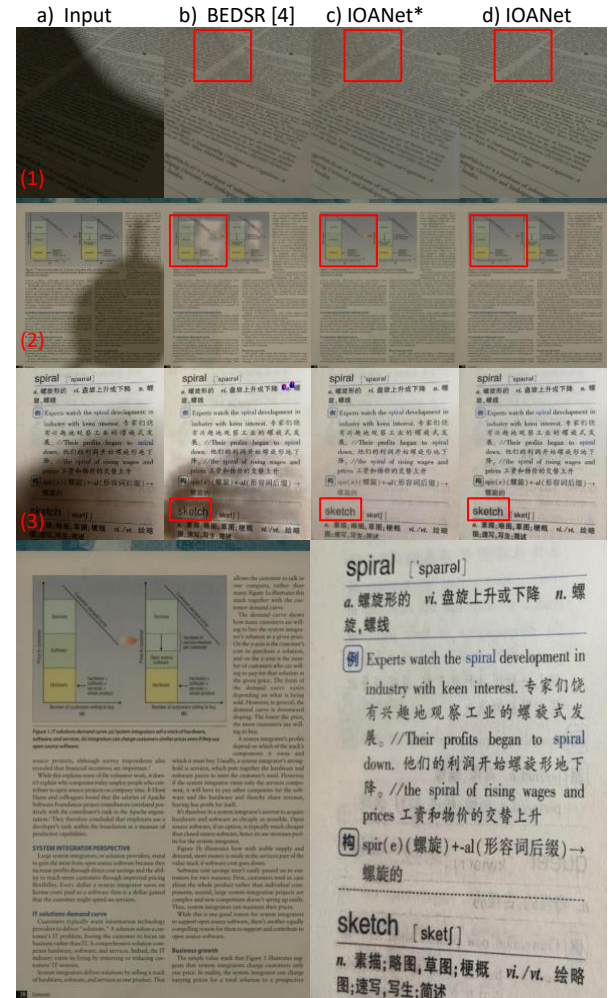**Loss terms.** The second (L1-only) and third rows (L1 + LPIPS) of Table 2 shows that using LPIPS as a loss term improves the results and justifies its addition.

| Training Dataset | MAE | PSNR |
|---|---|---|
| BSDD | 2.3344 / 1.6414 / 4.6026 | 36.84 / 13.76 / 14.11 |
| A-BSDD | 2.1163 / 1.4418 / 4.4220 | 37.43 / 13.91 / 14.14 |
| BSDD & Doc3DS+ | 2.0560 / 1.3761 / 4.2813 | 37.73 / 13.99 / 14.18 |
| All | **1.7893 / 1.0937 / 4.0659** | **38.76 / 14.08 / 14.26** |

**Table 5**. BSDD evaluation when IOANet is trained with different datasets. *All* refers to A-BSDD, A-OSR and Doc3DS+.

## 5. CONCLUSION

We propose LP-IOANet, an end-to-end, lightweight high-resolution document shadow removal solution. It consists of IOANet, a shadow removal architecture, encapsulated within a lightweight upsampler. We also propose three new datasets, which helps generalize our model to in-the-wild scenarios. Our results show that LP-IOANet comfortably outperforms



**Fig. 3**. Visualization of a) input image, and output of b) BEDSR [4], c) IOANet without attention (marked with *), d) our IOANet and e) our LP-IOANet. Differences are shown with red boxes, better viewed when zoomed in.

the existing state-of-the-art, while running on a mobile-device in real-time in high resolution.

# 6. REFERENCES

[1] Steve Bako, Soheil Darabi, Eli Shechtman, Jue Wang, Kalyan Sunkavalli, and Pradeep Sen, "Removing shadows from images of documents," in *Asian Conference on Computer Vision*. Springer, 2016, pp. 173–183.

[2] Seungjun Jung, Muhammad Abul Hasan, and Changick Kim, "Water-filling: An efficient algorithm for digitized document shadow removal," in *Asian Conference on Computer Vision*. Springer, 2018, pp. 398–414.

[3] Netanel Kligler, Sagi Katz, and Ayellet Tal, "Document enhancement using visibility detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2374–2382.

[4] Yun-Hsuan Lin, Wen-Chin Chen, and Yung-Yu Chuang, "Bedsr-net: A deep shadow removal network from a single document image," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12905–12914.

[5] Qingxiong Yang, Kar-Han Tan, and Narendra Ahuja, "Shadow removal using bilateral filtering," *IEEE Transactions on Image processing*, vol. 21, no. 10, pp. 4361–4368, 2012.

[6] Michael S Brown and Y-C Tsoi, "Geometric and shading correction for images of printed materials using boundary," *IEEE Transactions on Image Processing*, vol. 15, no. 6, pp. 1544–1554, 2006.

[7] Ram Krishna Pandey and AG Ramakrishnan, "Language independent single document image super-resolution using cnn for improved recognition," *arXiv preprint arXiv:1701.08835*, 2017.

[8] Xujun Peng and Chao Wang, "Building super-resolution image generator for ocr accuracy improvement," in *International Workshop on Document Analysis Systems*. Springer, 2020, pp. 145–160.

[9] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 618–626.

[10] Hieu Le and Dimitris Samaras, "Physics-based shadow image decomposition for shadow removal," *IEEE Transactions on Pattern Analysis & Machine Intelligence*, , no. 01, pp. 1–1, 2021.

[11] Mehmet Kerim Yucel, Valia Dimaridou, Anastasios Drosou, and Albert Saa-Garriga, "Real-time monocular depth estimation with sparse supervision on mobile," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2021, pp. 2428–2437.

[12] Qibin Hou, Daquan Zhou, and Jiashi Feng, "Coordinate attention for efficient mobile network design," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 13713–13722.

[13] Mehmet Kerim Yücel, Valia Dimaridou, Bruno Manganelli, Mete Ozay, Anastasios Drosou, and Albert Saà-Garriga, "Lra&ldra: Rethinking residual predictions for efficient shadow detection and removal," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 4925–4935.

[14] Jie Liang, Hui Zeng, and Lei Zhang, "High-resolution photorealistic image translation in real-time: A laplacian pyramid translation network," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 9392–9400.

[15] Peter J Burt and Edward H Adelson, "The laplacian pyramid as a compact image code," in *Readings in computer vision*, pp. 671–679. Elsevier, 1987.

[16] Sagnik Das, Hassan Ahmed Sial, Ke Ma, Ramon Baldrich, Maria Vanrell, and Dimitris Samaras, "Intrinsic decomposition of document images in-the-wild," in *British Machine Vision Conference 2020*, 2020.

[17] Hieu Le and Dimitris Samaras, "Shadow removal via shadow image decomposition," in *The IEEE International Conference on Computer Vision (ICCV)*, 2019.

[18] Bingshu Wang and CL Chen, "Local water-filling algorithm for shadow detection and removal of document images," *Sensors*, vol. 20, no. 23, pp. 6929, 2020.

[19] Diederik P Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," in *ICLR (Poster)*, 2015.

[20] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 586–595.

[21] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al., "Pytorch: An imperative style, high-performance deep learning library," *Advances in neural information processing systems*, vol. 32, 2019.