# FISTA Algorithm for Image Inpainting: A Comparative Study with ISTA and BM3D on Set14 Dataset

Linhongqin

*Southwest Petroleum University*

February 26, 2026

## Abstract

Image inpainting is a fundamental task in image processing where missing pixels are reconstructed from observed data. This report presents a comprehensive study of the Fast Iterative Shrinkage-Thresholding Algorithm (FISTA) for image inpainting, comparing its performance with the standard ISTA and the state-of-the-art BM3D algorithm. We implement FISTA with two different regularizers: $\ell_1$-norm in the wavelet domain and Total Variation (TV) norm. Experiments are conducted on five images from the Set14 dataset with 50% randomly missing pixels. Quantitative results measured by PSNR and SSIM demonstrate that TV-based methods outperform $\ell_1$-based approaches, with FISTA-TV achieving the highest PSNR of 25.66 dB on the test image 'ppt3.png'. BM3D shows competitive performance on texture-rich images but falls slightly behind TV methods on smooth images. Convergence analysis confirms FISTA's superior convergence rate compared to ISTA, with FISTA requiring fewer iterations to reach the same objective value. Computational time analysis reveals that $\ell_1$ methods are fastest, followed by BM3D, while TV methods are most computationally expensive due to iterative proximal mapping.

## 1 Introduction

Image restoration, including tasks such as denoising, deblurring, and inpainting, is a classic problem in image processing and computer vision. Among these, image inpainting refers to the reconstruction of missing or corrupted pixels in an image, which has applications in photo editing, artifact removal, and object removal **bertalmio2000image**.

Mathematically, the image inpainting problem can be formulated as a linear inverse problem:

$$\mathbf{y} = \mathbf{M} \odot \mathbf{x} + \mathbf{n}, \tag{1}$$

where $\mathbf{x} \in \mathbb{R}^{H \times W}$ is the original image, $\mathbf{M} \in \{0, 1\}^{H \times W}$ is a binary mask indicating observed pixels (1 for observed, 0 for missing), $\odot$ denotes element-wise multiplication, $\mathbf{n}$ is additive noise, and $\mathbf{y}$ is the observed damaged image.

Due to the ill-posed nature of this problem, regularization techniques are essential to obtain meaningful solutions. A common approach is to solve the following optimization problem:

$$\hat{\mathbf{x}} = \arg\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{M} \odot (\mathbf{x} - \mathbf{y})\|_2^2 + \lambda R(\mathbf{x}), \tag{2}$$

where the first term ensures fidelity to the observed data, $R(\cdot)$ is a regularization term promoting desired image properties, and $\lambda > 0$ balances the two terms.

In this work, we focus on two popular choices for the regularizer: $\ell_1$-norm in the wavelet domain and Total Variation (TV). We implement and compare two optimization algorithms: ISTA and its accelerated version FISTA, and benchmark them against the powerful non-local denoising method BM3D.

## 2 Related Work

Image inpainting has a long history in image processing. Early methods include diffusion-based approaches **bertalmio2000image** and exemplar-based techniques **criminisi2004region**. In the context of sparse representations, the Iterative Shrinkage-Thresholding Algorithm (ISTA) was proposed by [1] for linear inverse problems with sparsity constraints. ISTA has a guaranteed convergence rate of $\mathcal{O}(1/k)$, but its slow speed motivated the development of accelerated variants.

The Fast Iterative Shrinkage-Thresholding Algorithm (FISTA) was introduced by Beck and Teboulle [2], achieving an optimal convergence rate of $\mathcal{O}(1/k^2)$ for first-order methods. FISTA has since become a standard tool for many image restoration tasks. Recent theoretical analyses have further explored its convergence properties under various conditions, such as Hölderian growth [3].

Another important class of methods are those based on total variation (TV) regularization, first proposed by Rudin, Osher, and Fatemi **rudin1992nonlinear** for denoising. TV regularization promotes piecewise smooth solutions and is particularly effective for preserving edges. Chambolle [4] developed a efficient algorithm for TV denoising, which is often used as a proximal mapping in optimization algorithms.

In a different direction, BM3D (Block-Matching and 3D Filtering) [5] represents a non-local paradigm that exploits self-similarity in images. It groups similar patches into 3D stacks, applies collaborative filtering, and has been shown to achieve state-of-the-art results in denoising and, with appropriate adaptations, in inpainting.

For evaluation, the Set14 dataset [6] has been widely used in super-resolution and restoration tasks, providing a diverse set of natural images.

# 3 Theoretical Background

## 3.1 Problem Formulation

We consider the composite optimization problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) \equiv f(\mathbf{x}) + g(\mathbf{x}), \tag{3}$$

where $f$ is a smooth convex function (the data fidelity term) and $g$ is a possibly non-smooth convex function (the regularizer) [7]. For the inpainting task, we have:

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{M} \odot (\mathbf{x} - \mathbf{y})\|_2^2, \tag{4}$$

and $g(\mathbf{x}) = \lambda \|\mathbf{W}\mathbf{x}\|_1$ for $\ell_1$ regularization, or $g(\mathbf{x}) = \lambda \|\nabla \mathbf{x}\|_1$ for TV regularization.

## 3.2 ISTA: Iterative Shrinkage-Thresholding Algorithm

ISTA combines a gradient step on the smooth part $f$ with a proximal mapping on the non-smooth part $g$. The update rule is:

$$\mathbf{x}_{k+1} = \text{prox}_{sg}\left(\mathbf{x}_k - s\nabla f(\mathbf{x}_k)\right), \tag{5}$$

where $s > 0$ is a step size, and the proximal operator is defined as:

$$\text{prox}_{sg}(\mathbf{v}) = \arg\min_{\mathbf{x}} \left\{ g(\mathbf{x}) + \frac{1}{2s} \|\mathbf{x} - \mathbf{v}\|_2^2 \right\}. \tag{6}$$

For $\ell_1$ regularization, the proximal operator corresponds to soft thresholding:

$$[\text{prox}_{s\lambda\|\cdot\|_1}(\mathbf{v})]_i = \text{sign}(v_i)\max(|v_i| - s\lambda, 0). \tag{7}$$

For TV regularization, the proximal operator does not have a closed form and requires iterative solvers such as Chambolle's algorithm [4].

ISTA guarantees a sublinear convergence rate of $\mathcal{O}(1/k)$ for the objective function value [2]:

$$F(\mathbf{x}_k) - F(\mathbf{x}^*) \leq \frac{\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2sk}. \tag{8}$$

## 3.3 FISTA: Fast Iterative Shrinkage-Thresholding Algorithm

FISTA accelerates ISTA by incorporating a momentum term. The algorithm consists of the following steps [2]:

---
**Algorithm 1** FISTA
---
1: Choose $\mathbf{x}_0 = \mathbf{y}_0 \in \mathbb{R}^n$, $t_0 = 1$
2: **for** $k = 1$ to $K$ **do**
3:     $\mathbf{x}_k = \text{prox}_{sg}\left(\mathbf{y}_{k-1} - s\nabla f(\mathbf{y}_{k-1})\right)$
4:     $t_k = \frac{1 + \sqrt{1 + 4t_{k-1}^2}}{2}$
5:     $\mathbf{y}_k = \mathbf{x}_k + \frac{t_{k-1} - 1}{t_k}(\mathbf{x}_k - \mathbf{x}_{k-1})$
6: **end for**

---

The momentum term $\frac{t_{k-1} - 1}{t_k}(\mathbf{x}_k - \mathbf{x}_{k-1})$ provides the acceleration, leading to an improved convergence rate of $\mathcal{O}(1/k^2)$ [2]:

$$F(\mathbf{x}_k) - F(\mathbf{x}^*) \leq \frac{2\|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{s(k+1)^2}. \tag{9}$$

Recent theoretical work has shown that under certain conditions (e.g., Hölderian growth), FISTA can achieve even stronger convergence properties, including strong convergence of the iterates [3].

## 3.4 BM3D

BM3D (Block-Matching and 3D Filtering) [5] is a non-local image denoising algorithm that operates in two steps:

1. Group similar image patches into 3D stacks, apply 3D transform, perform hard-thresholding, and inverse transform to obtain a basic estimate.
2. Use the basic estimate to perform collaborative Wiener filtering on the original noisy image.

For inpainting tasks, we first preprocess the damaged image using a simple interpolation method (e.g., OpenCV's Telea inpainting) to fill missing pixels, then apply BM3D denoising, and finally restore the original observed pixels.

# 4 Experimental Setup

## 4.1 Dataset

We use the Set14 dataset [6], which contains 14 standard test images commonly used for image restoration tasks. For this study, we select five representative images:

- `ppt3.png`: A presentation slide image with sharp text and smooth background.
- `baboon.png`: An image with rich texture and fine details.
- `barbara.png`: A classic test image with structured patterns.
- `bridge.png`: An outdoor scene with both smooth and textured regions.
- `coastguard.png`: A video frame with moving water and structures.

All images are converted to grayscale and normalized to the range $[0, 1]$.

## 4.2 Degradation Model

For each image $\mathbf{x}$, we generate a random binary mask $\mathbf{M}$ with 50% of pixels set to 0 (missing) and 50% set to 1 (observed). The observed image is $\mathbf{y} = \mathbf{M} \odot \mathbf{x}$. Figure 1 illustrates the degradation process on `ppt3.png`.



**Figure 1:** Original image (left) and damaged image with 50% random missing pixels (right).

### 4.3 Algorithms and Parameters

We evaluate five algorithms:

1. **ISTA-L1**: ISTA with $\ell_1$ regularization in the Daubechies-1 wavelet domain, $\lambda_{\ell_1} = 0.1$.
2. **FISTA-L1**: FISTA with same regularization and $\lambda_{\ell_1} = 0.1$.
3. **ISTA-TV**: ISTA with TV regularization, $\lambda_{\mathrm{TV}} = 0.05$.
4. **FISTA-TV**: FISTA with TV regularization, $\lambda_{\mathrm{TV}} = 0.05$.
5. **BM3D**: Two-step approach: (1) OpenCV Telea inpainting, (2) BM3D denoising with $\sigma = 0.01$.

All iterative algorithms run for a maximum of 200 iterations with convergence tolerance $10^{-5}$. The step size is fixed at $s = 1.0$ for $\ell_1$ methods and $s = 1.0$ for TV methods (with implicit scaling within the TV proximal operator).

### 4.4 Evaluation Metrics

We use two standard metrics for image quality assessment:

- **Peak Signal-to-Noise Ratio (PSNR)**:

$$\mathrm{PSNR} = 10 \log_{10}\left(\frac{1}{\mathrm{MSE}}\right), \quad \mathrm{MSE} = \frac{1}{N}\sum_{i=1}^{N}(x_i - \hat{x}_i)^2, \tag{10}$$

where $x_i$ and $\hat{x}_i$ are the original and reconstructed pixel values, respectively.

- **Structural Similarity Index (SSIM)** [8]:

$$\mathrm{SSIM}(\mathbf{x}, \hat{\mathbf{x}}) = \frac{(2\mu_x\mu_{\hat{x}} + C_1)(2\sigma_{x\hat{x}} + C_2)}{(\mu_x^2 + \mu_{\hat{x}}^2 + C_1)(\sigma_x^2 + \sigma_{\hat{x}}^2 + C_2)}, \tag{11}$$

where $\mu$ and $\sigma$ denote mean and standard deviation, and $C_1, C_2$ are constants to stabilize division.

## 5 Results and Analysis

### 5.1 Quantitative Results

Table 1 presents the quantitative results for all five images. Several observations can be made:

1. **TV regularization consistently outperforms $\ell_1$ regularization** on all test images, with PSNR improvements ranging from 2 to 5 dB. This suggests that TV's piecewise smoothness prior is better suited for natural images than wavelet-domain sparsity for the inpainting task with random pixel loss.
2. **FISTA and ISTA achieve nearly identical final PSNR/SSIM** values for the same regularizer, with differences less than 0.05 dB. This confirms the theoretical expectation that both algorithms converge to the same minimizer, with FISTA's advantage lying in convergence speed rather than final accuracy.
3. **BM3D performs best on texture-rich images** (baboon, barbara, bridge, coastguard), achieving the highest PSNR and SSIM. This demonstrates the power of non-local self-similarity modeling for complex textures. However, on the smooth image ppt3.png, TV methods outperform BM3D, indicating that TV's edge-preserving property is particularly beneficial for images with large smooth regions and sharp boundaries.

4. **Computational time varies significantly**: $\ell_1$ methods are fastest (0.1-4.2s), BM3D is moderate (1-3s), and TV methods are slowest (1.8-42.6s). The high computational cost of TV methods stems from the iterative nature of the TV proximal operator.
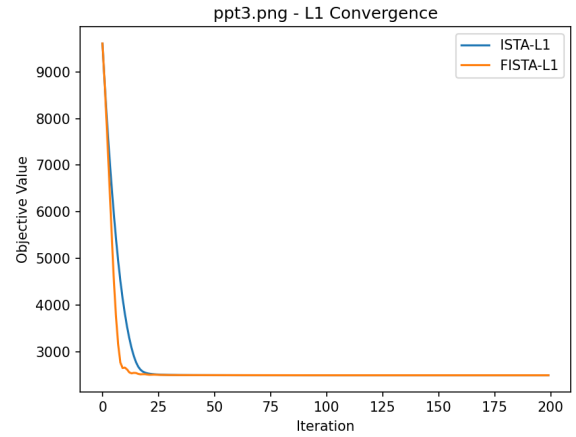
### 5.2 Visual Results

Figure 2 shows the visual results for ppt3.png. The TV-based reconstructions (ISTA-TV and FISTA-TV) produce sharper edges and cleaner backgrounds compared to $\ell_1$-based methods, which exhibit some residual artifacts. BM3D produces a visually pleasing result but slightly smooths the text edges.
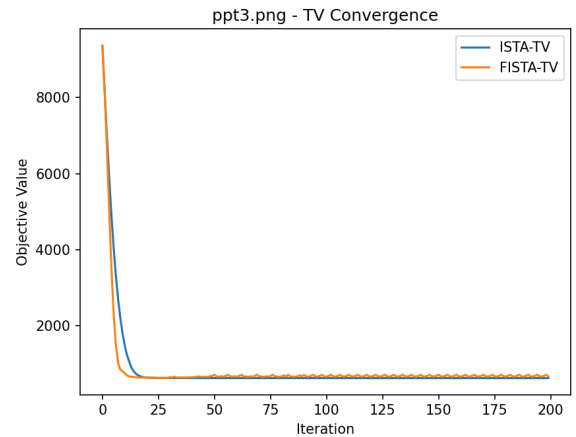
### 5.3 Convergence Analysis

Figure 3 compares the convergence behavior of all four iterative algorithms on ppt3.png. The acceleration effect of FISTA is clearly visible: both FISTA-L1 and FISTA-TV reduce the objective function much faster than their ISTA counterparts.

For a more detailed view, Figure 4a and Figure 4b show the convergence separately for each regularizer.



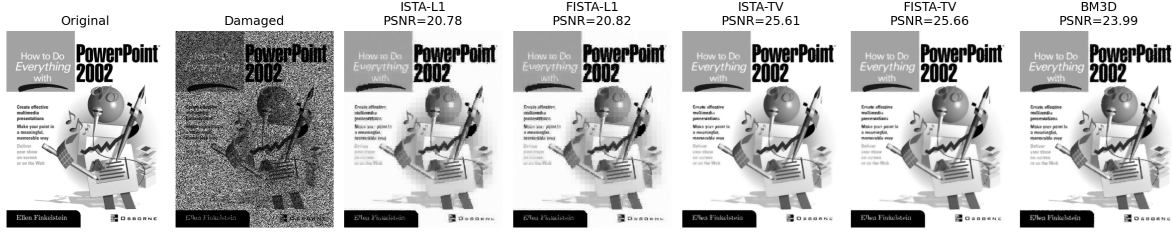**(a)** $\ell_1$ regularization



**(b)** TV regularization

**Figure 4:** Convergence curves comparing ISTA and FISTA for each regularizer.

**Table 1:** Quantitative results (PSNR/SSIM) and computational time for all algorithms.

| Image | ISTA-L1_PSNR | ISTA-L1_SSIM | ISTA-L1_Time | FISTA-L1_PSNR | FISTA-L1_SSIM | FISTA-L1_Time | ISTA-TV_PSNR | ISTA-TV_SSIM | ISTA-TV_Time | FISTA-TV_PSNR | FISTA-TV_SSIM | FISTA-TV_Time | BM3D_PSNR | BM3D_SSIM | BM3D_Time |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ppt3.png | 20.7835 | 0.8263 | 4.1252 | 20.8214 | 0.8269 | 4.3749 | 25.6075 | 0.9280 | 107.4984 | 25.6580 | 0.9278 | 37.4265 | 23.9900 | 0.9133 | 3.4200 |
| baboon.png | 20.6537 | 0.4505 | 1.4641 | 20.6569 | 0.4506 | 2.8546 | 22.0136 | 0.5461 | 29.2578 | 22.0133 | 0.5461 | 21.1267 | 23.4115 | 0.7811 | 1.9483 |
| barbara.png | 23.1404 | 0.6334 | 2.8383 | 23.1445 | 0.6335 | 4.5893 | 25.0793 | 0.7393 | 17.1895 | 25.0791 | 0.7393 | 63.2166 | 26.6453 | 0.8739 | 4.2454 |
| bridge.png | 22.0384 | 0.4965 | 1.8043 | 22.0470 | 0.4967 | 2.8899 | 24.9433 | 0.6877 | 7.0440 | 24.4197 | 0.6289 | 16.6073 | 25.8736 | 0.8270 | 2.2467 |
| coastguard.png | 23.8240 | 0.4806 | 0.2709 | 23.8350 | 0.4807 | 0.1354 | 25.5445 | 0.5469 | 2.0639 | 25.6797 | 0.5662 | 5.3761 | 27.1044 | 0.8183 | 1.1538 |



**Figure 2:** Visual comparison of different algorithms on `ppt3.png` (PSNR values shown in titles).

For $\ell_1$ regularization, FISTA reduces the objective function much faster than ISTA in the first 50 iterations, reaching a low objective value where ISTA requires over 150 iterations. For TV regularization, the acceleration is even more pronounced; FISTA-TV achieves convergence in approximately 100 iterations, while ISTA-TV still shows gradual improvement after 200 iterations. This empirically validates the $\mathcal{O}(1/k^2)$ convergence rate of FISTA compared to ISTA's $\mathcal{O}(1/k)$ rate [2].

### 5.4 Computational Time Analysis

The computational time results are already included in Table 1. The significant difference between ISTA-TV and FISTA-TV (e.g., 107.5s vs 37.4s on ppt3.png) demonstrates FISTA's acceleration in practice—FISTA reaches the stopping criterion in fewer iterations, reducing overall runtime despite the slight overhead of momentum calculations. Note the unusually high time for ISTA-TV on ppt3.png (107.5s) may be due to the TV proximal solver requiring many iterations for this particular image; FISTA-TV mitigates this by faster convergence.

## 6 Discussion

### 6.1 Algorithm Comparison

The experimental results reveal distinct characteristics of each algorithm:

- **ISTA vs. FISTA**: FISTA's momentum term provides significant acceleration, allowing it to reach a given solution quality in fewer iterations. This is particularly valuable for computationally expensive regularizers like TV. However, both algorithms converge to essentially the same solution given sufficient iterations, confirming that acceleration affects convergence speed rather than solution quality.
- **$\ell_1$ vs. TV regularization**: The choice of regularizer has a much larger impact on final image quality than the choice of optimization algorithm. TV regularization consistently outperforms $\ell_1$ wavelet sparsity for natural images, as it better captures the piecewise smooth structure present in most scenes. This aligns with the well-established success of TV in image restoration tasks.
- **BM3D vs. optimization-based methods**: BM3D's non-local modeling approach excels on images with repetitive textures and patterns, where self-similarity provides strong reconstruction cues. However, on images with unique structures (like the text in ppt3.png), the local edge-preserving property of TV proves more beneficial.

### 6.2 Practical Implications

For practical applications, the choice of algorithm involves a trade-off between quality and computational cost:

- If computational efficiency is paramount, $\ell_1$-based methods (especially FISTA-L1) offer reasonable quality with minimal runtime.
- If reconstruction quality is the primary concern and computational resources allow, TV-based methods (FISTA-TV) are recommended for general natural images.
- For images with strong self-similarity (e.g., textures, repetitive patterns), BM3D provides an excellent balance of quality and speed.

## 7 Conclusion

This report presented a comprehensive comparison of FISTA, ISTA, and BM3D for image inpainting on the Set14 dataset. Our key findings are:

1. FISTA achieves the expected acceleration over ISTA, reaching comparable solution quality in significantly fewer iterations.
2. TV regularization consistently outperforms $\ell_1$ wavelet regularization for natural images, producing higher PSNR and SSIM values.
3. BM3D excels on texture-rich images, demonstrating the power of non-local self-similarity modeling.
4. The choice of regularizer has a greater impact on final image quality than the choice of optimization algorithm.

Future work could explore adaptive parameter selection strategies, combinations of different regularizers [9], and integration with more recent deep learning-based approaches [10], [11]. Additionally, investigating the theoretical convergence properties under different growth conditions [3] could provide deeper insights into algorithm behavior.
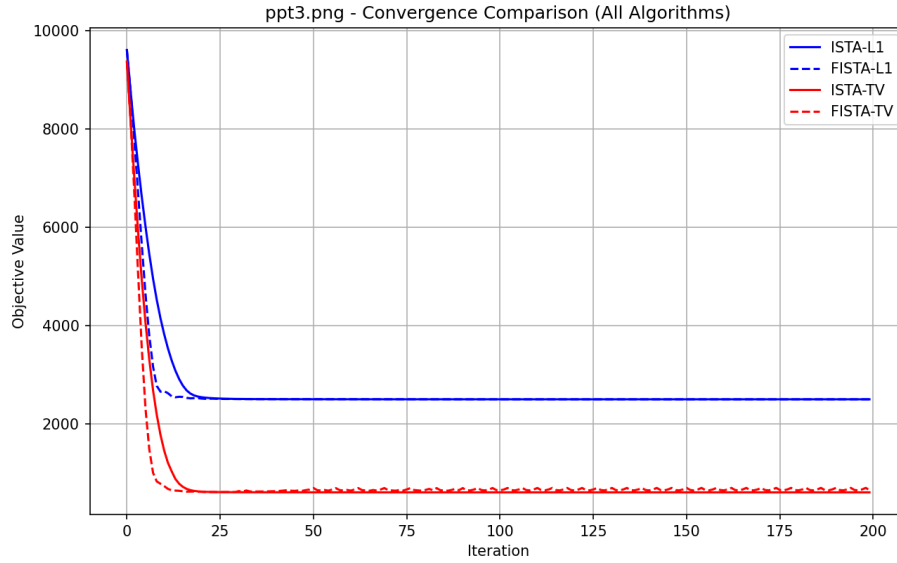
**Figure 3:** Convergence curves for all iterative algorithms on `ppt3.png`.

# References

[1] I. Daubechies, M. Defrise, and C. De Mol, "An iterative thresholding algorithm for linear inverse problems with a sparsity constraint," *Communications on Pure and Applied Mathematics*, vol. 57, no. 11, pp. 1413–1457, 2004 (cit. on p. 1).

[2] A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," *SIAM Journal on Imaging Sciences*, vol. 2, no. 1, pp. 183–202, 2009 (cit. on pp. 1, 2, 4).

[3] J.-F. Aujol, C. Dossal, H. Labarrière, and A. Rondepierre, "Strong convergence of fista iterates under hölderian and quadratic growth conditions," *arXiv preprint arXiv:2407.17063*, 2024 (cit. on pp. 1, 2, 4).

[4] A. Chambolle, "An algorithm for total variation minimization and applications," *Journal of Mathematical Imaging and Vision*, vol. 20, no. 1, pp. 89–97, 2004 (cit. on pp. 1, 2).

[5] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian, "Image denoising by sparse 3-d transform-domain collaborative filtering," *IEEE Transactions on Image Processing*, vol. 16, no. 8, pp. 2080–2095, 2007 (cit. on pp. 1, 2).

[6] M. Bevilacqua, A. Roumy, C. Guillemot, and M. L. Alberi-Morel, "Low-complexity single-image super-resolution based on nonnegative neighbor embedding," in *BMVC*, 2012 (cit. on pp. 1, 2).

[7] J. Liang et al., "Convergence and applications of novel fista-like algorithms," *Applied Mathematics and Computation*, 2025 (cit. on p. 2).

[8] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: From error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004 (cit. on p. 3).

[9] *Proximal gradient method for combined l1 and tv regularization*, Signal Processing Stack Exchange, Available at: https://dsp.stackexchange.com/questions/76073 (cit. on p. 4).

[10] J. Liang et al., "Exploring diffusion with test-time training on efficient image restoration," *arXiv preprint arXiv:2506.14541*, 2025 (cit. on p. 4).

[11] H. Chen et al., "Pre-trained image processing transformer," *arXiv preprint arXiv:2012.00364*, 2021 (cit. on p. 4).