

# Beyond Natural Images: A Benchmark for Cross-Domain Image Reconstruction

Anonymous CVPR submission

Paper ID \*\*\*\*\*

## Abstract

This project will investigate how well state-of-the-art deep learning models for image reconstruction generalize across different visual domains. By evaluating CNN-, Transformer-, and Diffusion-based architectures on super-resolution, denoising, and inpainting tasks across diverse image types—such as natural scenes, text, astronomical, and stylized images—we aim to identify consistent strengths, weaknesses, and failure patterns. The findings will provide insights into cross-domain robustness and guide the development of more adaptable, content-aware reconstruction approaches.

## 1. Introduction

Every day, cameras capture billions of images—landscapes, selfies, night skies, documents, street scenes—and many of them are blurry, noisy, or low-resolution. Image reconstruction models promise to fix that: they take a damaged picture and make it clear again. But here's the catch—the same model that works beautifully on a photo of a mountain often fails on a picture of a starry sky, a page of text, or a sketch. You might get letters that melt together, stars that vanish, or cartoon lines that blur.

Right now, researchers and developers typically train separate models for each type of image. There are models just for faces, others for anime drawings, and others for natural scenes. This works in narrow settings but creates a messy ecosystem of one-off tools. If you want to enhance mixed content—say a photo of a person under the night sky with visible text—there's no reliable, all-purpose solution. Each model specializes but lacks general understanding of what it's looking at.

That's the gap we want to fill. Our goal is to first quantify that gap and test how well today's best reconstruction models actually generalize across different kinds of images—from skies and constellations to text, animals, and human portraits—and figure out why some architectures succeed where others fail. Instead of building yet another model from scratch, we'll take the most advanced existing ones

(CNNs, Transformers, Diffusion models) and evaluate them on multiple visual domains using public datasets. By comparing their strengths and weaknesses, we'll learn what design choices make them robust and where they break.

If we succeed, we'll provide something useful for everyone—from photographers to scientists.

## 2. Planned Method

The central idea is to take State-of-The-Art models and evaluate their out-of-the-box generalization performance on unseen image domains without any fine-tuning. This approach directly tests their inherent robustness. By analyzing model designs that generalize best, we aim to contribute insights that could lead to a SoTA AI for improving image quality across diverse content. The process is broken down into three main stages.

### 2.1. Selection of Tasks and Architectures

We will focus on three core image reconstruction tasks: **Super-Resolution**, **Denoising**, and **In-Painting**. For each task, we will select representative, high-performing architectures that embody different design philosophies. Our selection will include CNN-based models like EDSR or ESRGAN, which excel at learning local features; Transformer-based models like SwinIR, which use attention to capture long-range dependencies and global structure; and Diffusion-based models like SR3, which iteratively reverse a noise process to synthesize photorealistic details.

### 2.2. Assembling a Cross-Domain Evaluation Suite

We will curate a benchmark suite composed of several challenging image domains, as detailed in Table 1. Models will be trained only on the Natural/Scenic domain, using standard benchmarks like DIV2K and Flickr2K, which will serve as our baseline. We will then test their generalization on several out-of-domain datasets. These include the Human/Animal domain (CelebA-HQ, AFHQ) to test reconstruction of fine facial features; the Text/Document domain (TextZoom) to evaluate the preservation of sharp edges and legibility; the Astronomy/Night Sky domain (Hubble, SDSS images) to assess performance on sparse,

038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075

076 high-contrast data; and the Stylized/Line Art domain  
 077 (anime/manga datasets) to test the handling of abstract con-  
 078 tent with sharp lines and flat colors.

Domain	Example Sources / Datasets	Purpose
Natural / Scenic	DIV2K, Flickr2K, COCO	Baseline domain
Human / Animal	CelebA-HQ, AFHQ	Faces, animals
Text / Document	TextZoom, scanned images	Structured symbols
Astronomy	Hubble, SDSS, astrophotos	Sparse, high-contrast
Stylized / Art	Anime/manga, Waifu2x	Strong edges, abstraction

Table 1. Cross-Domain Evaluation Datasets.

### 079 2.3. Evaluation Framework

080 For each task, we will use the official pre-trained weights  
 081 for the selected architectures, which are typically trained  
 082 on natural image datasets. We will then run inference with  
 083 these models on the test sets from all domains listed above.  
 084 The core analysis will focus on quantifying the performance  
 085 degradation as models move from the familiar "Natural" do-  
 086 main to the unfamiliar target domains.

## 087 3. Planned Experiments

088 Our experiments are designed to systematically evaluate  
 089 how well existing deep learning models for image recon-  
 090 struction generalize across different image domains. The  
 091 evaluation will span three tasks (super-resolution, denois-  
 092 ing, and inpainting) across six distinct image domains.

### 093 3.1. Experimental Setup

094 For each task, we will compare three representative archi-  
 095 tecture families: CNN-based models (ESRGAN, EDSR), a  
 096 Transformer-based model (SwinIR), and a Diffusion-based  
 097 model (SR3 or Stable Diffusion Upscaler). To ground our  
 098 findings, we will include classical methods like bicubic in-  
 099 terpolation and BM3D/DnCNN as baseline references. All  
 100 deep learning models will be used in their pretrained form  
 101 to ensure a fair comparison. Each model will be tested  
 102 on a fixed set of 100 images per domain, with controlled  
 103 degradations applied to simulate realistic conditions, such  
 104 as downsampling, Gaussian noise, or random masking.

### 105 3.2. Metrics and Evaluation

106 To quantify performance, we will use a combination of  
 107 metrics. For pixel-level fidelity, we will use **PSNR** and  
 108 **SSIM**. For perceptual similarity, we will rely on **LPIPS**.  
 109 Domain-specific metrics will also be employed, such as  
 110 **OCR accuracy** for text domains and **star-count consis-**  
 111 **tency** for astronomy images. We will also record runtime  
 112 and GPU memory usage to assess computational efficiency.  
 113 To specifically measure cross-domain robustness, we will  
 114 compute a **Cross-Domain Drop (CDD)** metric, defined as

the normalized performance difference between a model's  
 115 in-domain and out-of-domain results.

### 116 3.3. Success Criteria

117 The project will be considered successful if our experiments  
 118 demonstrate clear, measurable differences in how the model  
 119 families behave across domains. A key indicator of suc-  
 120 cess will be quantifying a significant generalization gap (a  
 121 high CDD value) when models are applied to unfamiliar  
 122 image types. We also aim to find consistent evidence that  
 123 diffusion- and transformer-based models retain higher per-  
 124 ceptual quality under domain shift.

## 125 4. Timeline and Feasibility

126 This project is designed to be feasible within a typical  
 127 semester timeframe. We have broken down the work into  
 128 several phases to ensure steady progress and timely com-  
 129 pletion.

131 **Weeks 1-3: Setup and Baseline.** This phase will involve  
 132 a final literature review, setting up the computational envi-  
 133 ronment, and downloading all datasets. We will also run all  
 134 models on the baseline "Natural/Scenic" domain to estab-  
 135 lish our in-domain performance metrics.

136 **Weeks 4-8: Cross-Domain Experiments.** This is the  
 137 core experimental phase. We will systematically run the  
 138 pre-trained models for all three tasks (super-resolution,  
 139 denoising, inpainting) across the five out-of-domain test  
 140 suites. Results and performance metrics will be carefully  
 141 logged.

142 **Weeks 9-11: Analysis and Interpretation.** During this  
 143 period, we will analyze the collected data, calculate the  
 144 Cross-Domain Drop (CDD) for each model, and generate  
 145 visualizations to compare performance. We will focus on  
 146 identifying failure modes and patterns that explain why cer-  
 147 tain architectures generalize better than others.

148 **Weeks 12-14: Final Report and Presentation.** The fi-  
 149 nal weeks will be dedicated to writing the final project re-  
 150 port, structuring our findings into a coherent narrative, and  
 151 preparing the final presentation.