
COGS 181 Final Project

Yue Yin

University of California, San Diego
yuyin@ucsd.edu

Keyi Chen

University of California, San Diego
kec020@ucsd.edu

Abstract

1 The project focuses on image reconstruction, specifically targeting the restoration
2 of images with a centrally located square mask that obscures part of the image.
3 We propose the development of a context encoder, utilizing L2 loss as the primary
4 mechanism, to effectively reconstruct the missing central portion of the image.
5 This initiative is particularly aimed at assisting individuals in repairing damaged or
6 incomplete images, thereby preserving valuable visual information.

7 1 Introduction

8 The project aims to bridge the gap between deep learning theories and their practical applications
9 through the development of a context-based pixel prediction algorithm for image restoration. Cen-
10 tral to our approach are context encoders, a specialized variant of convolutional neural networks,
11 tasked with reconstructing missing image segments by leveraging contextual information from the
12 surrounding pixels. These encoders are designed to assimilate and interpret the content of the input
13 image, enabling them to generate a coherent fill for the absent sections. Upon establishing the
14 encoder framework, our exploration will extend to identifying and assessing various factors that could
15 influence the fidelity and precision of the reconstructed areas. Specifically, we intend to investigate
16 the effects of the training dataset's composition and the spatial characteristics of the missing segments.
17 Through these analyses, we aim to refine our methodology and enhance the predictive performance
18 of our system. Our preliminary results demonstrate the model's capability to generate convincing and
19 satisfactory reconstructions, aligning with our initial expectations.

20 2 Method

21 2.1 Dataset

22 We evaluate our methodology on two distinct datasets: one comprising images of resolution 32×32
23 pixels and the other containing images of resolution 128×128 pixels.

24 **CIFAR-10:** The CIFAR-10 dataset consists of 60,000 color images, each of 32×32 pixels, divided
25 into 10 classes, with 6,000 images per class. The dataset is split into 50,000 training images and
26 10,000 testing images. The ten unique classes include: airplane, automobile, bird, cat, deer, dog, frog,
27 horse, ship, and truck, each being mutually exclusive to ensure no overlap among them. Specifically,
28 the "automobile" category includes sedans and SUVs, excluding pickup trucks, while the "truck"
29 category is confined to large trucks only.

30 **Universal Image Embeddings:** Our second dataset, tailored for the Google Universal Image
31 Embedding competition, comprises over 130,000 images resized to 128×128 pixels and spans 11
32 categories. This dataset is an amalgamation of various sources, encompassing both scraped and
33 meticulously curated images to form a rich collection. The categories feature Apparel from the Deep

34 Fashion Dataset, Artwork, Cars from the Stanford Cars Dataset, Dishes, and others predominantly
 35 sourced from Google. Given the dataset's extensive size and complexity, our evaluation concentrates
 36 on the Landmark category, anticipating it to offer higher inpainting accuracy. Thus, we have refined
 37 the dataset to 30,000 images within this category for a more focused and feasible analysis.

38 **2.2 Network Architecture for 32×32 pixels images:**

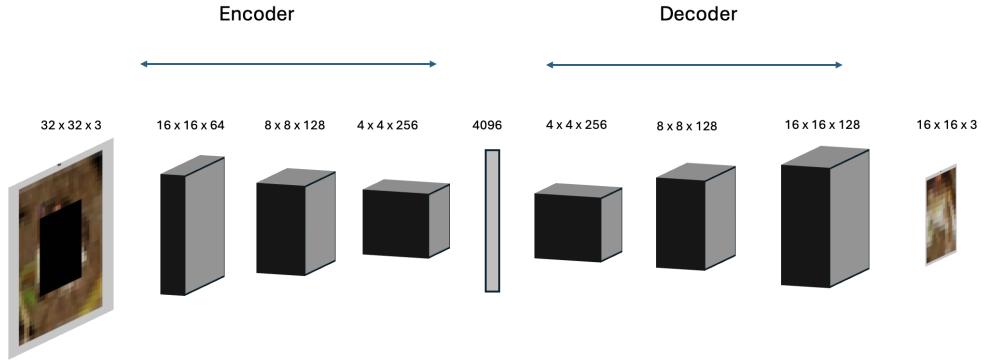


Figure 1: Architecture Visualization for 32×32

39 The proposed encoder-decoder architecture is designed to process images of size 32×32 pixels with
 40 3 color channels (RGB). The model consists of an encoder part, which progressively down-samples
 41 the input image to extract and encode salient features, and a decoder part, which up-samples the
 42 encoded features to reconstruct the image at a target resolution of 16×16 pixels.

43 **Encoder:** The encoding section begins with an input layer accepting images of the specified input
 44 shape. It employs a series of convolutional layers, each followed by batch normalization to stabilize
 45 learning and accelerate convergence. The convolutional layers use 64, 128, and 256 filters, respec-
 46 tively, with a kernel size of 3×3 and 'ReLU' activation for non-linearity. Each convolutional layer is
 47 followed by a 2×2 max-pooling layer with 'same' padding to reduce the spatial dimensions by half,
 48 enhancing the model's ability to capture more abstract features at each subsequent layer. The final
 49 layer in the encoder outputs a compact, encoded representation of the input image.

50 **Fully Connected Simulation:** Post-encoding, the model simulates a fully connected layer by first
 51 flattening the encoded feature maps, then passing them through a dense layer with 4096 neurons
 52 and 'ReLU' activation. This dense layer acts as a bottleneck, allowing the network to learn a
 53 compressed representation of the input data. The output of the dense layer is then reshaped back into
 54 a 4-dimensional tensor to be processed by the decoder.

55 **Decoder:** The decoder aims to reconstruct the image from the encoded representation. It mirrors the
 56 encoder's structure but in reverse, utilizing up-sampling and convolutional layers to progressively
 57 increase the spatial resolution of the feature maps. The decoder uses two sets of 3×3 convolutional
 58 layers with 256 and 128 filters, respectively, each followed by batch normalization and a 2×2
 59 up-sampling layer to enlarge the feature maps. The final layer of the decoder is a convolutional layer
 60 with a kernel size of 3×3 , using the 'sigmoid' activation function to output the reconstructed image
 61 at the desired target resolution.

62 The entire model is compiled with the Adam optimizer and mean squared error (MSE) loss function,
 63 suitable for regression tasks like image reconstruction. This architecture is expected to perform image
 64 inpainting or similar tasks by learning to encode the essential features of the input image and then
 65 decode these features to reconstruct or generate the image content.

66 **2.3 Network Architecture for 128×128 pixels images:**

67 This architecture is similar to the previously discussed model, with notable adjustments in the input
 68 and output dimensions, as well as in the configuration of the convolutional layers. Designed to
 69 accommodate images with an input size of $128 \times 128 \times 3$, the model aims to reconstruct these images

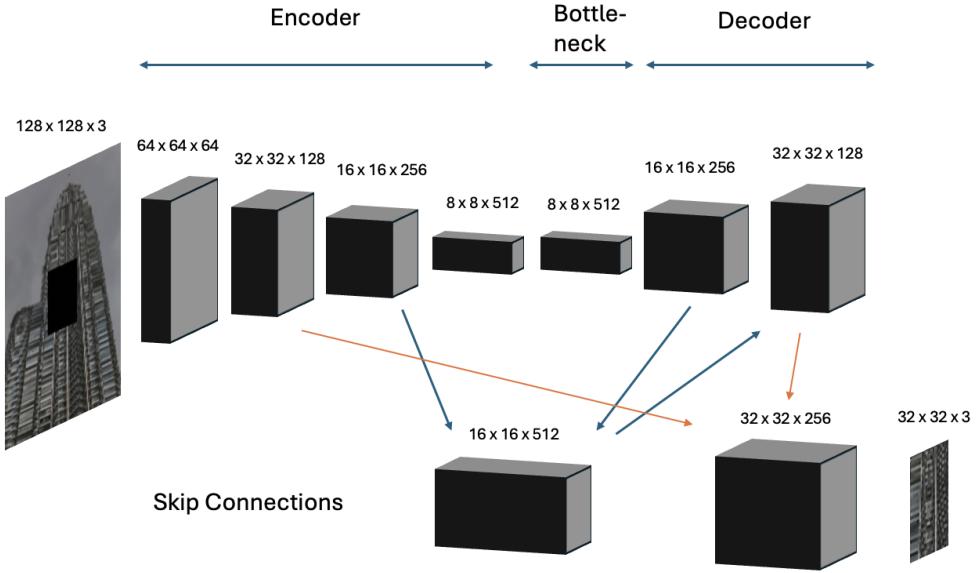


Figure 2: Architecture Visualization for 32×32

70 at a reduced resolution of $32 \times 32 \times 3$. The modified arrangement of convolutional layers further
 71 distinguishes this version, tailoring it to the specific requirements of dimensionality reduction and
 72 feature preservation within a different spatial context.

73 **Encoder:** The encoder segment initiates with an Input layer receiving 128×128 RGB images. It
 74 features a series of Convolutional layers, each followed by Batch Normalization to improve training
 75 stability and convergence speed. The convolutional layers are configured with 64, 128, 256, and
 76 512 filters respectively, all utilizing a kernel size of 3×3 and 'ReLU' activation for introducing
 77 non-linearity. Each Convolutional layer is followed by a 2×2 MaxPooling layer with 'same' padding,
 78 systematically reducing the spatial dimensions by half and thereby encoding the input image into a
 79 more compact and abstract representation. The final layer of the encoder outputs a condensed feature
 80 map of size $8 \times 8 \times 512$, effectively capturing the essential features of the input image.

81 **Bottleneck:** At the bottleneck of the network, another Convolutional layer with 512 filters is applied
 82 to the encoded feature map, maintaining the spatial dimensions while potentially enabling the model
 83 to refine the feature abstraction further.

84 **Decoder:** The decoder component, designed to reconstruct the image, begins with an UpSampling
 85 layer to increase the spatial dimensions of the feature map from the bottleneck. This is followed
 86 by a series of Convolutional layers, each with Batch Normalization, to progressively refine the
 87 upsampled features. The architecture incorporates skip connections from the encoder to the decoder,
 88 using the Concatenate function, which helps in preserving spatial details by combining feature maps
 89 from the encoder with those in the decoder. This approach mitigates the information loss typically
 90 associated with deep networks and aids in more accurate reconstruction of the output image. The
 91 decoder up-samples and processes the feature maps back to a size of 32×32 , followed by a final
 92 Convolutional layer that maps the features to the desired output dimension of $32 \times 32 \times 3$, employing
 93 a 'sigmoid' activation to ensure the output pixel values are within a normalized range.

94 Similar to the previous one, the model employs the Adam optimizer and Mean Squared Error (MSE)
 95 as the loss function, suitable for the regression nature of the image reconstruction task.

96 **3 Experiments**

97 In our experimental analysis, we conducted training and testing across three distinct scenarios to
 98 evaluate the versatility and effectiveness of our model in varying inpainting contexts.

99 **3.1 General Class**

100 In our initial phase of experimentation, the model was tasked with inpainting across a diverse array
 101 of image categories, having been trained on the full CIFAR-10 dataset, which encompasses 60,000
 102 images sized at 32×32 pixels, in which we have 50,000 images for training. This broad-based
 103 approach aimed to evaluate the model's capability in handling a wide variety of general classes.

104 Figure 3 below illustrates the progression of training and validation losses throughout the training
 105 period. Notably, a sharp decrease in the validation loss is observed within the first two epochs, after
 106 which the validation loss stabilizes, indicating a swift adaptation of the model to the general features
 107 of the dataset early in the training process.

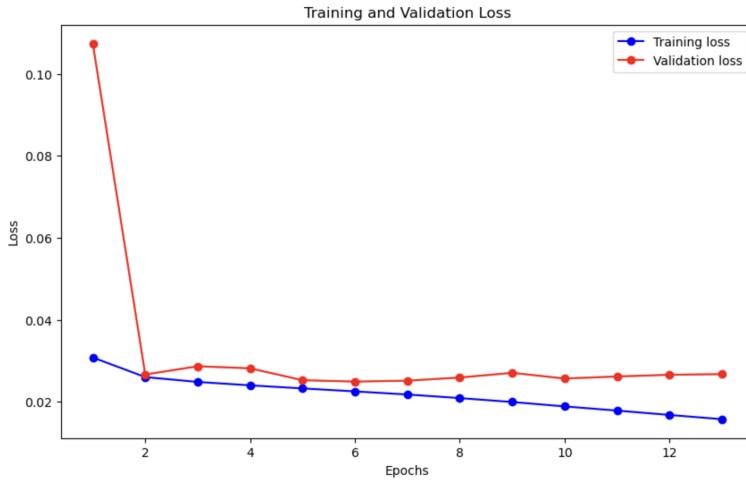


Figure 3: General class training/testing loss history.

108 Figure 4 and Figure 5 provide a visual comparison between the original 32×32 images (ground
 109 truth) and the corresponding inpainted outputs produced from the masked inputs. These figures are
 110 divided into two segments: the left segment displays the original and inpainted images from the
 111 training dataset, whereas the right segment focuses on the testing dataset. This bifurcation allows for
 112 a detailed assessment of the model's inpainting effectiveness, both in scenarios where it has been
 113 directly trained on the data (training) and where it encounters unseen data (testing).



Figure 4: 32×32 general training images.

Figure 5: 32×32 general testing images.

114 **3.2 Specific Class**

115 Subsequently, we narrowed our focus to specialized inpainting tasks, training our model exclusively
116 with images from a specific class with the hope of seeing an enhanced performance in targeted
117 scenarios.

118 Initially, we concentrated on the "truck" category from CIFAR-10, engaging in inpainting experiments
119 with 32×32 images. This subset comprised 5,000 training images and 1,000 testing images, a
120 significantly smaller dataset compared to the general class experiment. This reduction in data volume
121 was instrumental in assessing the model's efficiency and accuracy under constrained conditions.

122 The loss trends, illustrated in Figure 6, depicted for the truck-class-specific model training show
123 a precipitous decline in training loss, suggesting rapid early learning, which stabilizes as epochs
124 progress. Meanwhile, the validation loss initially spikes, indicating potential overfitting or model
125 instability, but then settles into a steady state, closely paralleling the training loss. This pattern
126 indicates effective learning and generalization, despite early fluctuations in validation loss.

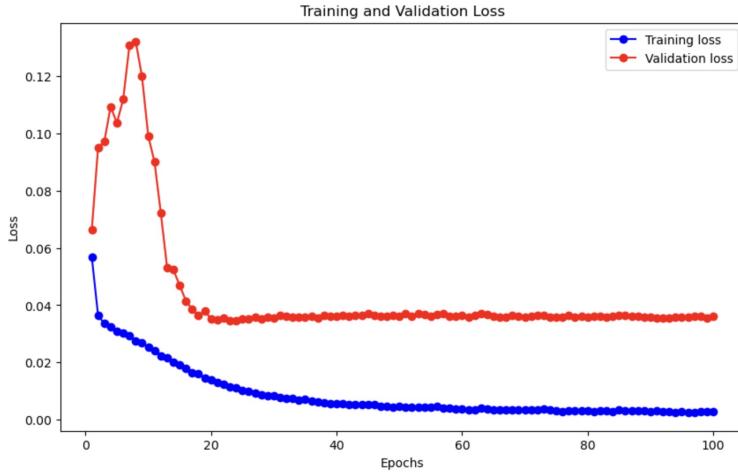


Figure 6: 32×32 specific truck training/testing loss.

127 Figures 7 and 8 offer a comparative analysis of the original 32×32 truck images against the model's
128 inpainted outputs. The closely matched inpainted images, discernible to the naked eye, underscore
129 the model's adeptness at reconstructing specific class features, indicating its robust performance
130 despite the limited dataset size.



Figure 7: 32×32 truck training images.

Figure 8: 32×32 truck testing images.

131 Additionally, we extended our investigation to a more complex dataset by selecting the "landmark"
132 class from Universal Image Embeddings, applying our model to 128×128 images to test its
133 performance on high-resolution inpainting tasks, utilizing a dataset comprised of 24000 training
134 images and 6000 testing images.

135 The training and validation loss trends for the landmark class, illustrated in Figure 9, reveal a
136 consistent decrease in training loss accompanied by a fluctuating but overall decreasing trend in

137 validation loss. This indicates the model's learning progression, albeit with varying degrees of
 138 adjustment to the validation dataset over time.

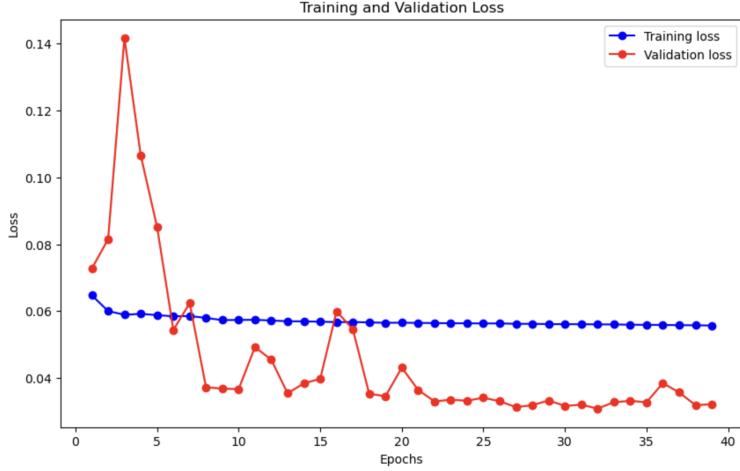


Figure 9: 128 x 128 landmark training/testing loss.

139 Figures 10 and 11 show the original 128 × 128 landmark images with the model's inpainted outputs.
 140 The outputs, while capturing the general structure, exhibit a noticeable blur within the inpainted
 141 regions, reflecting a comparative decline in performance for this dataset size.

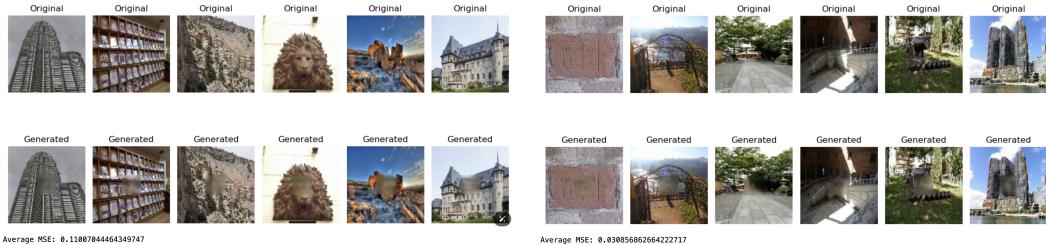


Figure 10: 128 x 128 landmark training images.

Figure 11: 128 x 128 landmark testing images.

142 3.3 General Class & Specific Class Trained Model with the Truck Images

143 Our experimentation extended to comparing the performance of models trained under different
 144 conditions: one with a dataset encompassing all classes within the CIFAR-10 collection and another
 145 with a dataset focused exclusively on the truck category. This comparison aimed to ascertain the
 146 impact of dataset specificity on the quality of inpainted images, particularly those belonging to the
 147 truck class.

148 Figure 12, presented on the left, illustrates the inpainting results for truck images generated by the
 149 model trained on the entire CIFAR-10 dataset, achieving a Mean Squared Error (MSE) of 0.0356.
 150 Conversely, Figure 13, displayed on the right, showcases the results from the model trained solely
 151 on truck class images from the CIFAR-10 dataset, with a slightly higher MSE of 0.0360. Despite
 152 the marginal difference in MSE, a visual assessment of the inpainted images reveals a discernible
 153 improvement in quality for the model trained specifically on truck images, which might be attributed
 154 to the model's refined understanding and familiarity with the specific features and characteristics of
 155 truck images, facilitating a more accurate and detailed inpainting process.

156 This observation suggests that the focused training on truck images endows the model with a more
 157 refined and specialized understanding of truck features, enabling it to perform inpainting tasks
 158 with greater precision and detail. The nuanced knowledge acquired from the truck-specific dataset
 159 compensates for the slight increase in MSE, highlighting the importance of dataset relevance and



Figure 12: 32×32 truck images from general class trained model.



Figure 13: 32×32 truck images from specific truck class trained model.

160 specificity in enhancing inpainting performance. This finding underscores the potential benefits of
161 tailoring training datasets to the specific subject matter of interest, particularly in applications where
162 high fidelity and detail in inpainting outcomes are paramount.

163 3.4 Hyper-parameters Tuning

164 The fine-tuning of hyper-parameters plays a pivotal role in optimizing the performance of machine
165 learning models. In this phase of our experiment, we concentrated on two crucial hyper-parameters:
166 the optimizer and the learning rate. Through systematic testing and evaluation, we aimed to identify
167 the optimal combination that would enhance the efficiency and accuracy of our model across the
168 various inpainting tasks.

169 *Optimizer Selection:* We experimented with several popular optimizers, including Stochastic Gradient
170 Descent (SGD), Adam, AdamW, and RMSprop, to determine their impact on the model's performance.
171 Among these, Adam emerged as the slightly superior choice, demonstrating marginally better overall
172 performance compared to the others. This outcome aligns with Adam's well-documented balance of
173 speed and convergence efficiency, making it a robust choice for a wide range of deep learning tasks.

174 *Learning Rate Exploration:* The learning rate is a critical parameter that influences the speed and
175 stability of the training process. We explored learning rates beyond the commonly used default
176 of 0.001, testing alternatives such as 0.01 and 0.0001 to gauge their effects on model training and
177 convergence. Our findings indicated that the default learning rate of 0.001 outperformed the others,
178 striking an optimal balance between rapid convergence and the avoidance of overshooting minimal
179 loss landscapes.

180 Based on these investigations, we adopted Adam as our optimizer with a learning rate of 0.001
181 for all subsequent model implementations. This combination was identified as the best practice,
182 given its consistent performance advantages across our experiments. By integrating these optimized
183 hyper-parameters, we aimed to further refine our model's capability to accurately perform inpainting
184 tasks, thereby enhancing its overall effectiveness and reliability.

185 3.5 Metric Evaluation

186 In our experimental evaluations, we employed the Mean Squared Error (MSE) metric to assess the
187 performance of our model. This evaluation was conducted on two distinct datasets: CIFAR-10
188 and a collection of Universal Image Embeddings featuring 128×128 resolution landmark images.
189 Specifically, our analysis focused on calculating the MSE solely within the regions covered by the
190 generated masks, comparing these areas directly with the corresponding ground truth mask regions to
191 gauge the accuracy of our model's reconstructions.

MSE Evaluation	32×32 all classes	32×32 truck	128×128 landmark
training	0.0216	0.0025	0.11007
testing	0.02493	0.03596	0.03086

Table 1: Mean Squared Error for 32×32 and 128×128 image datasets

192 The data in Table 1 reveals that the model trained on 32×32 images across all classes achieves a lower
193 Mean Squared Error (MSE) than the model trained exclusively on 32×32 images of the truck category.

194 Despite this quantitative outcome, qualitative assessments—based on visual inspection by human
195 evaluators—suggest that the truck-specific model yields more visually convincing reconstructions
196 compared to the general model.

197 **4 Limitation and Future Work**

198 **4.1 Limitation**

199 Our study encountered several limitations that point to areas for future exploration and improvement.
200 These limitations primarily revolve around the resolution of inpainted images and the computational
201 resources required for training on high-resolution datasets.

202 *High-Resolution Prediction Challenges:* One of the notable challenges we faced was achieving
203 high-resolution predictions in both 32×32 and 128×128 image scenarios. Despite our efforts to
204 enhance image resolution through the introduction of adversarial loss, we observed no significant
205 improvement in this regard. This limitation suggests that the current approach may not effectively
206 capture the finer details necessary for high-resolution inpainting, indicating a potential area for further
207 investigation.

208 *Training Time and Computational Resources:* The training process for 128×128 images was
209 particularly time-consuming, highlighting the need for high-performance computational resources.
210 Our limited access to such resources restricted our ability to experiment with larger datasets and more
211 complex model architectures.

212 **4.2 Future Work**

213 *Adding Adversarial Loss:* Given the inconclusive results with adversarial loss in our current setup,
214 future work could explore different strategies for integrating adversarial components. Adjusting
215 the balance between reconstruction and adversarial loss, experimenting with various adversarial
216 architectures, or implementing more advanced adversarial training techniques may yield better
217 outcomes in enhancing image resolution.

218 *Optimizing Model Architecture:* Future efforts could focus on redesigning the model architecture to
219 reduce the number of parameters without compromising the quality of inpainting. Simplifying the
220 network or employing techniques like pruning or knowledge distillation may enable efficient training
221 on less powerful machines while maintaining or improving performance.

222 *Mask Size in High-Resolution Images:* For 128×128 images, we experimented with a mask size of
223 64×64 , which proved too large, contributing to more blurred outputs, therefore we ended up to use a
224 32×32 mask size. Future research could involve optimizing mask size relative to image dimensions
225 to balance the difficulty of the inpainting task with the model’s ability to generate coherent and
226 detailed reconstructions.

227 *Hardware, Epochs and Model Depth:* Our experiments were constrained by GPU limitations, resulting
228 in fewer training epochs and a reluctance to design deeper architectures. Future studies could secure
229 access to more powerful GPUs, thus being able to increase the number of training epochs, which
230 allows for more thorough learning and exploring deeper, more complex architectures to enhance the
231 model’s capacity for high-resolution image inpainting.

232 In summary, while we have navigated through a series of challenges and identified certain limitations
233 in our current approach, these hurdles pave the way for a broad spectrum of exciting research
234 opportunities. By addressing these limitations through advanced adversarial techniques, architectural
235 optimizations, and leveraging superior computational resources, future work has the potential to
236 significantly advance our current results in image inpainting.

237 **5 Conclusion**

238 In summary, our study embarked on a comprehensive exploration of image inpainting, using the
239 CIFAR-10 and Universal Image Embeddings datasets to examine the effectiveness of our model
240 across a range of scenarios. Our investigation spanned general and specific class inpainting, alongside

241 meticulous hyper-parameter optimization, revealing insights into the model's ability to perform
242 feature learning and reconstruction for images of varying resolutions.

243 Despite encountering challenges such as achieving high-resolution predictions and limitations due
244 to computational resources, our research highlighted several avenues for further exploration. The
245 attempt to incorporate adversarial loss, despite not delivering expected results, suggests the need for
246 further refinement to improve image resolution. Additionally, exploring computational efficiency
247 through model architectural adjustments presents a promising direction for facilitating high-resolution
248 inpainting.

249 Future research holds significant potential to enhance the field of image inpainting through advanced
250 adversarial methods, model optimization, and leveraging improved computational strategies. Address-
251 ing the limitations identified in our work, future efforts can lead to substantial advancements in image
252 inpainting, enabling more detailed and realistic reconstructions across diverse imaging scenarios.

253 Our commitment to this project is demonstrated through our efforts in detailed architectural visualiza-
254 tion and tackling the challenges of image inpainting, affirming our belief in the contribution of our
255 work to the field. Looking forward, the insights and foundation established by our study are poised to
256 support further advancements in image inpainting technology.

257 **6 Bonus Point**

258 We believe we deserve the bonus point because we devoted a significant amount of time to creating
259 the architectural visualization figures.

260 **References**

- 261 [1] Singh, R. (2023) Google Universal Image Embeddings (128x128). Available at: <https://www.kaggle.com/datasets/rhtsingh/google-universal-image-embeddings-128x128>.
- 263 [2] Krizhevsky, A. (2009) CIFAR-10 (Canadian Institute for Advanced Research) dataset. Available at: <https://www.cs.toronto.edu/~kriz/cifar.html>.
- 265 [3] Dumoulin, V., Shlens, J., and Kudlur, M. (2016) A Learned Representation for Artistic Style. Available at:
266 <https://arxiv.org/pdf/1604.07379.pdf>.