

# Supplementary Materials for “Cross-Scale MAE: A Tale of Multiscale Exploitation in Remote Sensing”

Generally speaking, there are three ways to evaluate the effectiveness of a representation learning scheme: 1) use reconstruction error for a *direct* evaluation; 2) use ablation studies to investigate the impact of hyperparameter selection and different components of the proposed network structure; and 3) use downstream tasks like classification accuracy for an *indirect* evaluation. Since Cross-Scale MAE is essentially a self-supervised representation scheme, we follow similar standards. This supplement provides additional quantitative and qualitative (visual) results complementing results presented in the main paper to provide further validation of the claims made. Specifically,

- In Sec. A., we provide both visual and quantitative comparisons (in MSE and SSIM) of the reconstruction performance between baseline and Cross-Scale MAE.
- In Sec. B., we provide further ablation studies of Cross-Scale MAE on the effect of different masking strategies.
- An interesting question we asked ourselves is if the performance gain of the proposed Cross-Scale MAE also generalizes to natural images. In Sec. C., we investigate the performance of Cross-Scale MAE in natural images, using CoCo as training and test sets. We also study its generalization capacity across the different domains of natural imagery and remote sensing imagery.
- In Sec. D., we elaborate on the details of the deployment of xFormers as an efficient backbone (See Sec. 4.3 in the main paper) with further studies regarding the memory efficiency of different attention types.

## A. Direct Evaluation through Multiscale Reconstruction Performance

In this section, we analyze the multiscale reconstruction performance of the proposed Cross-Scale MAE. We compare it with SatMAE, a baseline model, and demonstrate the improvements achieved by Cross-Scale MAE. Both models were pre-trained on the fMoW-RGB dataset, and we assessed their capabilities in handling images at different scales.

To assess the reconstruction performance, we employ the Mean Squared Error (MSE) and the Structural Similarity Index (SSIM) [8, 5] metrics. These metrics quantify the structural difference between images, with a higher SSIM value and lower MSE value indicating better performance.

Fig. 1 demonstrates the improvements achieved across different scales. The first column represents the original input image before masking, while the second column displays the input with a 75% mask applied. The third column exhibits the reconstruction results of the baseline model, while the last column showcases the improved reconstruction achieved by our model, Cross-Scale MAE. Additionally, the corresponding SSIM metric is presented alongside each reconstruction. The red and green boxes over the raw image in the first row showcase the crop locations of the images in the second and third rows, respectively.

We observe a significant reduction in artifacts within the masked portions of the reconstructed images. The artifacts in these regions indicate the presence of uninformative or distorted latent representations. This observation implies that such regions would be ineffective for representation learning and might even negatively impact downstream models attempting to learn from these representations. Notably, at the full scale (first row), our model demonstrates a 38% improvement in the SSIM metric. At a 40% scale (second row), we observe an 18% improvement. Finally, at a 25% scale (last row), our model showcases a 35% improvement.

To further illustrate the effectiveness of Cross-Scale MAE, we provide additional visualizations in Fig. 2. These samples showcase individual images cropped to random scales, highlighting the superior reconstruction achieved by our model compared to the baseline.

To comprehensively evaluate the performance in a multiscale scenario, we present the average metrics over the fMoW-RGB testing set in Fig. 3. Each input image is evaluated 25 times with different random crop scales and masks. This procedure ensures a robust assessment of our model’s performance in a multiscale context.

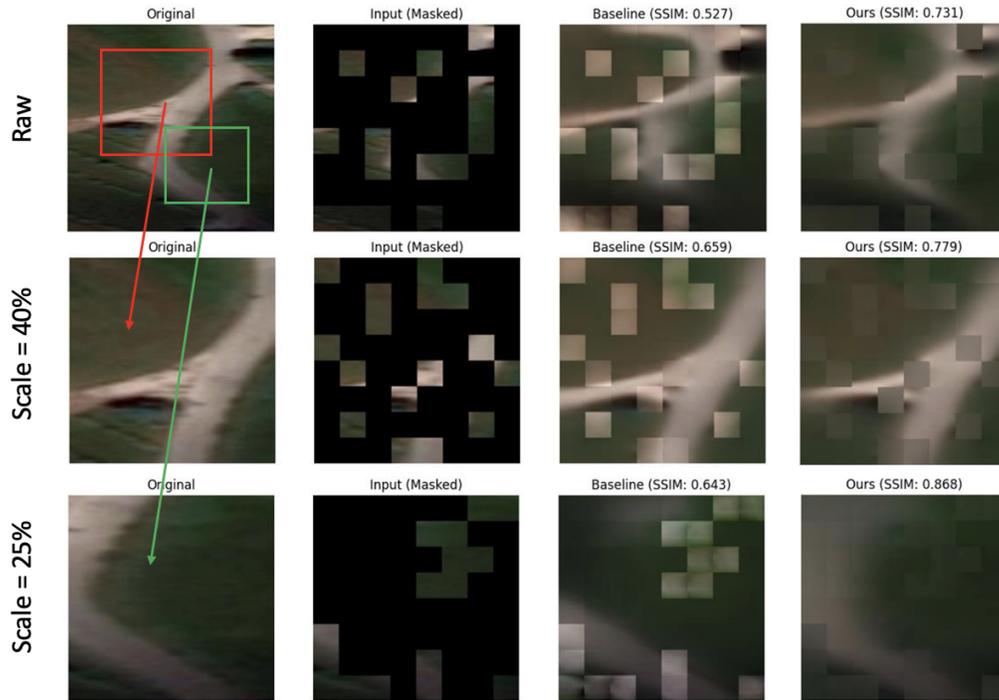


Figure 1: Comparison of Cross-Scale MAE and SatMAE reconstructions at different scales

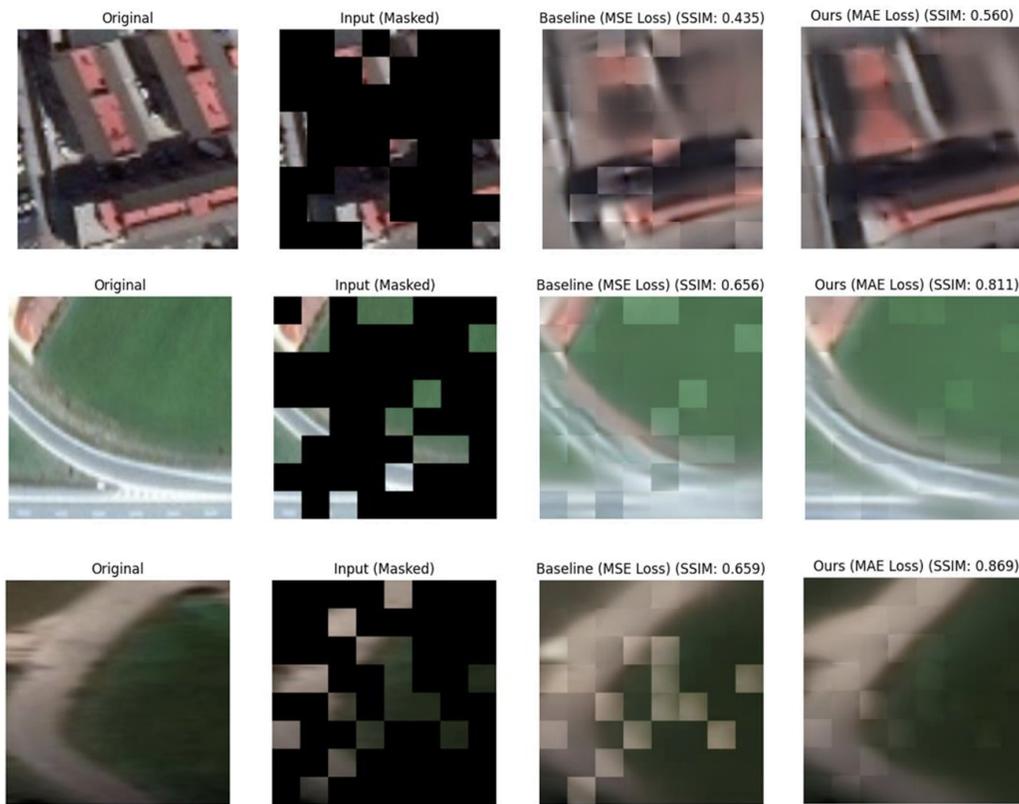


Figure 2: Reconstruction samples with random scales (fMoW test set)

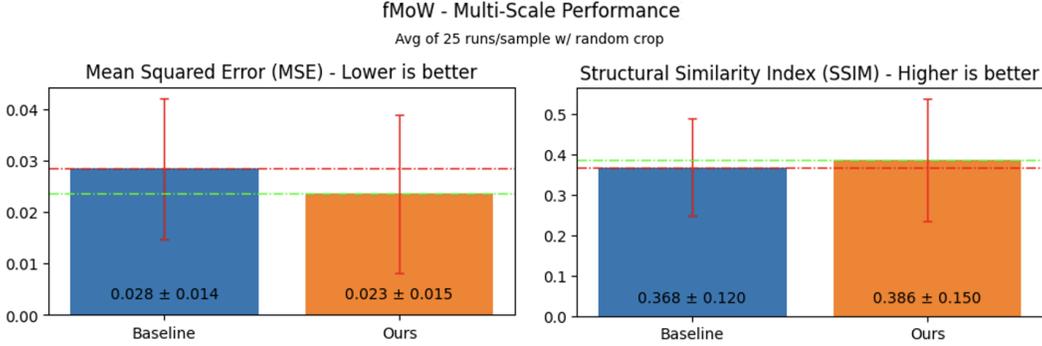


Figure 3: Average metrics comparison of Cross-Scale MAE and SatMAE (fMoW test set)

51 Our extensive evaluation demonstrates that Cross-Scale MAE excels in reconstructing multiscale  
 52 images, surpassing the performance of the baseline model, SatMAE. The observed improvements  
 53 and the mitigation of artifacts in masked portions indicate the superiority of Cross-Scale MAE in  
 54 capturing meaningful representations and enhancing remote sensing image understanding in diverse  
 55 scale conditions.

## 56 B. More Ablation Study

57 In this section, we investigate the impact of mask consistency and masking ratios on the model’s  
 58 representation learning capacity.

59 To assess the quality of the learned representations, we employ a non-parametric k-Nearest Neighbor  
 60 (kNN) classification approach with zero-shot learning. This evaluation method measures the ability  
 61 of the pre-trained model to produce semantically coherent representations, a desired characteristic for  
 62 practical zero-shot classification tasks. Similar evaluation strategies have been employed in other  
 63 notable works [9, 3, 2].

64 In the following subsections, we conduct detailed analyses to examine the effects of mask consistency  
 65 and masking ratios on the performance of Cross-Scale MAE. We employ the RESISC dataset and  
 66 utilize the ViT-Base as the backbone architecture for the evaluations, same as used in the main paper.

### 67 B.1. Effect of Mask Consistency

68 In Cross-Scale MAE, we generate two scale augments from the raw image and have the option to  
 69 apply either a consistent mask, where the patch location of the mask remains fixed across both scale  
 70 images, or different masks with varying patch locations for each scale. In this section, we compare  
 71 the performance of these two cases using kNN with  $k=20$  on the representations with different scale  
 72 ratios.

Table 1: Effect of mask consistency in Cross-Scale MAE on RESISC

| Masking Strategy | kNN 25%      | kNN 50%      | kNN 100%     |
|------------------|--------------|--------------|--------------|
| Consistent       | 0.762        | 0.812        | 0.824        |
| Different        | <b>0.787</b> | <b>0.831</b> | <b>0.853</b> |

73 Table 1 presents the evaluation results, showcasing the effect of mask consistency on Cross-Scale  
 74 MAE performance on the RESISC dataset. The table highlights the kNN accuracy for different scale  
 75 ratios. Notably, we observe that inconsistent masks yield nearly a 2% improvement in performance.  
 76 This improvement may be attributed to the introduction of additional variance during training,  
 77 resulting in a more robust and invariant representation being learned by the model.

78 **B.2. Effect of Masking Ratio and Training Time**

79 In addition to exploring the impact of mask consistency, we also investigate the effect of different  
 80 masking ratios on the performance of Cross-Scale MAE. Table 2 reports the evaluation results using  
 81 three different mask ratios: 60%, 75%, and 90%. The kNN accuracy for each mask ratio at various  
 82 scale ratios is measured and compared.

Table 2: Effect of Mask Ratio in Cross-scale MAE on RESISC

| Masking Ratio | KNN 25%      | KNN 50%       | KNN 100%     |
|---------------|--------------|---------------|--------------|
| 60%           | 0.7803       | <b>0.8322</b> | 0.8407       |
| 75%           | <b>0.787</b> | 0.831         | <b>0.853</b> |
| 90%           | 0.7524       | 0.7971        | 0.7977       |

83 We observe interesting patterns from the results presented in Table 2. A relatively low mask ratio  
 84 of 60% still yields excellent performance for remote sensing images, demonstrating competent  
 85 representation learning capabilities. However, employing a high mask ratio of 90% leads to decreased  
 86 performance. This reduction may be attributed to the significant information loss caused by a high  
 87 degree of masking, which affects the model’s ability to capture essential details and features.

88 These findings highlight the importance of carefully selecting an optimal mask ratio to balance  
 89 preserving relevant information and encouraging robust representation learning.

90 Finally, we show the performance of Cross-Scale MAE with different backbones at different training  
 91 epochs in Table 3.

92 **C. Evaluation on Natural Images (CoCo Dataset)**

93 To assess the generalization capabilities of Cross-Scale MAE, we evaluate by pre-training on the  
 94 CoCo2017 dataset, focusing on natural images. This evaluation allows us to validate the effectiveness  
 95 of our model beyond remote sensing images and examine its performance in a different domain. We  
 96 pre-train the Cross-Scale MAE on the CoCo2017 dataset using the following parameter settings:  
 97 ViT-Base as the backbone architecture, a base learning rate of  $5 \times 10^{-3}$ , a weight decay of  $5 \times 10^{-3}$ ,  
 98 and an input size of  $128 \times 128$ . The model is trained for 400 epochs.

99 **C.1. Pre-Training Performance on CoCo Images**

100 We visualize the reconstruction results on the CoCo dataset in Fig.4 and present the corresponding  
 101 evaluation metrics in Fig.5. We compare the performance of Cross-Scale MAE with the baseline  
 102 model, MAE [4].

103 From Fig.5, we observe that Cross-Scale MAE outperforms the baseline model in terms of both  
 104 MSE and SSIM. Additionally, in Fig.5, we notice that the baseline model exhibits more artifacts  
 105 in the reconstruction results at the locations of the masked patches. In contrast, Cross-Scale MAE  
 106 demonstrates a closer representation of the actual distribution of pixel values that should be present  
 107 in those locations.

108 **C.2. Zero-Shot Performance on fMoW-RGB Images**

109 Furthermore, we evaluate the zero-shot reconstruction performance of Cross-Scale MAE on the  
 110 fMoW-RGB dataset using the model pre-trained on CoCo. In this evaluation, we freeze the model  
 111 trained on CoCo and reconstruct images from the fMoW-RGB dataset. It is important to note that  
 112 the model has not seen any images from the fMoW-RGB dataset during its training. The zero-shot

Table 3: Performance with different backbone and training epoch of Cross-Scale MAE (%)

| Epochs    | 50    | 100   | 150   | 200   | 250   | 300   |
|-----------|-------|-------|-------|-------|-------|-------|
| ViT-Base  | 63.72 | 74.13 | 75.42 | 77.73 | 78.55 | 79.25 |
| ViT-Large | 60.01 | 73.84 | 79.42 | 83.09 | 83.39 | 85.34 |



Figure 4: Reconstruction performance on the CoCo Dataset

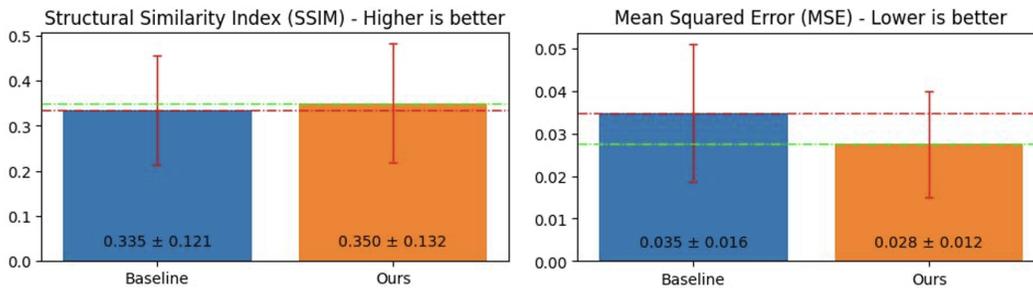


Figure 5: Multiscale reconstruction performance on the CoCo dataset

113 reconstruction results are displayed in Fig. 6. These results indicate that the learned representations by  
 114 Cross-Scale MAE generalize well, as the zero-shot reconstruction still produces meaningful outputs  
 115 on the fMoW-RGB dataset.

116 The evaluation of natural images demonstrates the versatility of Cross-Scale MAE, showcasing its  
 117 ability to capture meaningful representations and generalize effectively across different domains.  
 118 These findings highlight the potential of our model to enhance image understanding and reconstruction  
 119 tasks in various applications beyond remote sensing imagery.

## 120 D. Timm vs. xFormers MAE Backbones

121 This section presents the findings of an ablation study conducted as the initial step of our research.  
 122 The study aimed to establish an efficient and flexible backbone for the final implementation by  
 123 optimizing the original Masked Auto-Encoder for improved training time and a smaller memory  
 124 footprint. The ultimate goal was to enable feasible end-to-end training and inference on a single  
 125 GPU.

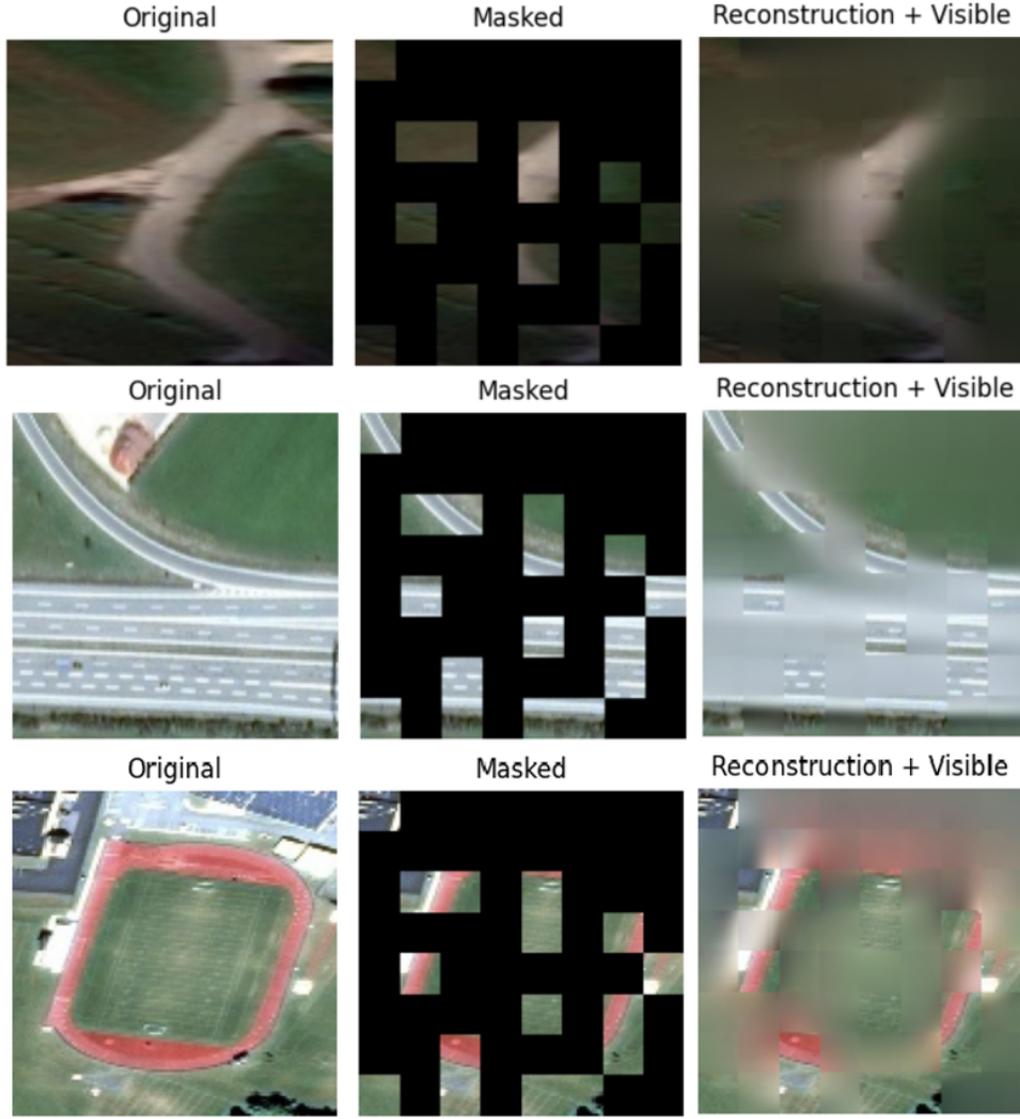


Figure 6: Zero-shot reconstruction in fMoW with Cross-Scale MAE pre-trained on CoCo

126 Initially, the MAE implementation relied on the Timm library for its various components. However,  
 127 during our research, we came across the xFormers[7] library, which offers customizable building  
 128 blocks and cutting-edge components yet to be available in mainstream libraries like PyTorch. xForm-  
 129 ers claims to be built with efficiency in mind, delivering fast and memory-efficient performance[7].

### 130 D.1. Baseline Efficiency Benchmark

131 We re-implemented the original MAE using this library to evaluate the potential benefits of using  
 132 xFormers. This allowed us to compare the performance of the Timm implementation against the  
 133 xFormers implementation based on PyTorch 1.13.1. Additionally, we considered the recently released  
 134 PyTorch 2.0.0, which promised improved efficiency and optimizations for transformer models. It is  
 135 important to note that xFormers did not yet support the newly released PyTorch version during our  
 136 experiments.

137 We conducted several benchmark experiments to evaluate the performance of the original Timm  
 138 implementation using PyTorch 1.13.1 and PyTorch 2.0.0 against our xFormers implementation on  
 139 PyTorch 1.13.1. These evaluations were performed on an NVIDIA RTX A6000 GPU.

Table 4: Time/Step -  $224 \times 224$ , patch16, batch256

| Model     | Timm (1.13.1) | xFormers (1.13.1) | Timm (2.0.0) |
|-----------|---------------|-------------------|--------------|
| ViT-Base  | 0.4849        | <b>0.4144</b>     | 0.4685       |
| ViT-Large | 0.7939        | <b>0.7143</b>     | 0.7584       |

Table 5: Memory Usage -  $224 \times 224$ , patch16, batch256

| Model     | Timm (1.13.1) | xFormers (1.13.1) | Timm (2.0.0) |
|-----------|---------------|-------------------|--------------|
| ViT-Base  | 22639         | <b>19020</b>      | 22223        |
| ViT-Large | 34225         | <b>30739</b>      | 32974        |

140 The results revealed the superior performance of xFormers when working with an input resolution  
 141 of  $224 \times 224$ . Comparing memory usage and time per step, xFormers outperformed the Timm  
 142 implementation on PyTorch 1.13.1. Specifically, the Vision Transformer (ViT) Base achieved  
 143 a 17% increase in speed, while ViT Large demonstrated an 11% improvement with xFormers.  
 144 Regarding memory efficiency, xFormers showcased a 19% enhancement for ViT Base and an 11%  
 145 improvement for ViT Large. It is worth noting that PyTorch 2.0.0 also provided some speed and  
 146 memory improvements at this input resolution, often falling between the performance of Timm and  
 147 xFormers on PyTorch 1.13.1.

Table 6: Time/Step -  $128 \times 128$ , patch16, batch512

| Model     | Timm (1.13.1) | xFormers (1.13.1) | Timm (2.0.0)  |
|-----------|---------------|-------------------|---------------|
| ViT-Base  | 0.2948        | 0.2820            | <b>0.2796</b> |
| ViT-Large | 0.5245        | 0.5047            | <b>0.4986</b> |

Table 7: Memory Usage -  $128 \times 128$ , patch16, batch512

| Model     | Timm (1.13.1) | xFormers (1.13.1) | Timm (2.0.0) |
|-----------|---------------|-------------------|--------------|
| ViT-Base  | 12213         | 12003             | <b>11805</b> |
| ViT-Large | 20891         | 21060             | <b>19601</b> |

148 When utilizing an input size of  $128 \times 128$ , the performance differences were subtle. The original  
 149 Timm implementation on PyTorch 2.0.0 exhibited a slight advantage over the xFormers imple-  
 150 mentation on PyTorch 1.13.1. The time per step showed a mere 1% improvement with the Timm  
 151 implementation on PyTorch 2.0.0. In terms of memory usage, there was an approximate 7% improve-  
 152 ment with the Timm implementation on PyTorch 2.0.0 compared to the xFormers implementation on  
 153 PyTorch 1.13.1. Although present, these differences were less substantial than those when using a  
 154 higher input resolution.

155 The results of the ablation study underscored the significance of using xFormers as the backbone for  
 156 our final implementation. Not only did xFormers provide enhanced flexibility through its customizable  
 157 building blocks and cutting-edge components, but it also demonstrated superior speed and memory  
 158 efficiency performance.

## 159 D.2. Effect of Attention Type

160 Our experiments primarily used the Scaled Dot Product (SDP) attention mechanism, a common  
 161 choice for transformer architectures. However, attention mechanisms can significantly influence  
 162 a model’s efficiency and reconstruction accuracy. The initial implementation, using the Timm  
 163 library, only supported SDP attention. In contrast, the xFormers library—our final choice for  
 164 implementation—provides an expanded selection of attention mechanisms, allowing for a more  
 165 comprehensive examination of how different attention types affect model performance and efficiency.

166 The Fourier Mix attention [6], which integrates the generalized Fourier integral theorem into the  
 167 dot-product attention step of the standard transformer, showed significant improvements in both  
 168 speed and memory consumption. Compared to SDP attention, Fourier Mix attention was 34% faster  
 169 and consumed approximately 44% less memory. Incorporating Fourier Mix attention addresses the



Figure 7: Memory efficiency of different attention types

170 traditional SDP attention’s limitations, capturing complex interactions among the features of the queries and keys more effectively and reducing redundancy between attention heads [6].  
 171

172 On the other hand, the Local attention [1] mechanism provided only a minor speed improvement but  
 173 also slightly increased memory consumption for our current architecture. Local attention offers a  
 174 novel approach to managing long sequences by dividing attention into global and local components  
 175 to facilitate the efficient processing of long input sequences. Due to this, we suspect this mechanism’s  
 176 benefits would be a lot more noticeable in even larger architectures with much larger embedding  
 177 dimensions.

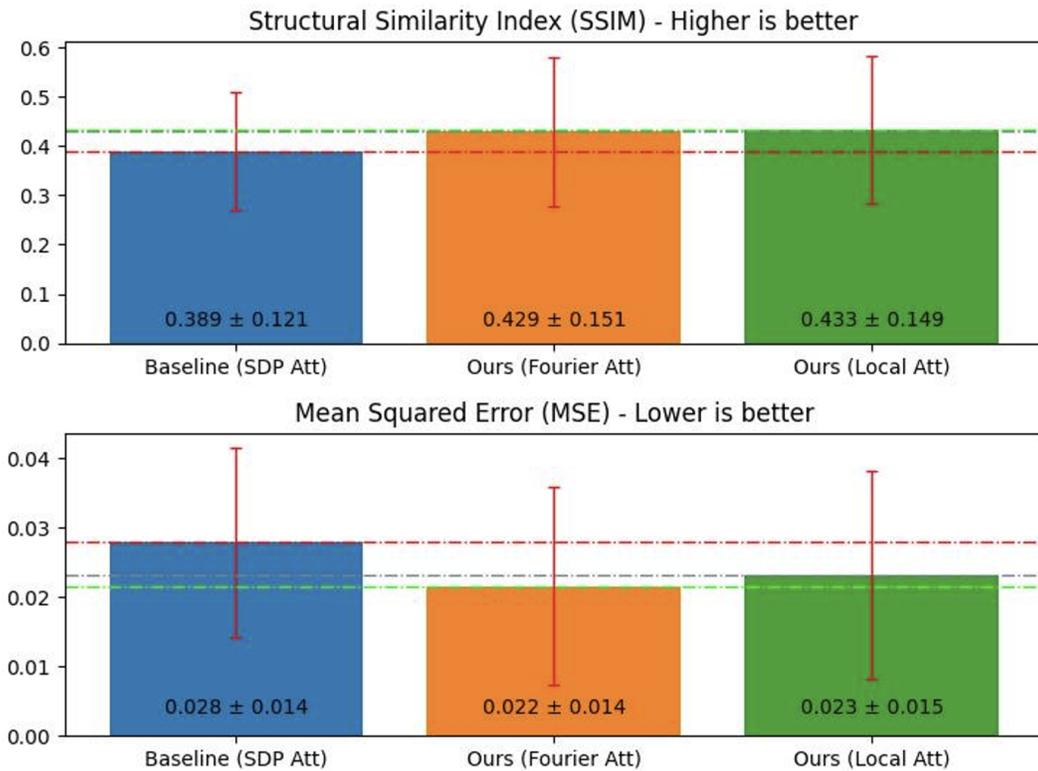


Figure 8: Performance comparison of attention types in Cross-Scale MAE

178 Following these observations, we further tested the Fourier Mix and Local attention mechanisms,  
179 assessing their reconstruction performance against the SDP attention baseline. As shown in Figure 8,  
180 both Fourier Mix and Local attention mechanisms demonstrated a 10% improvement in the Structural  
181 Similarity Index (SSIM) metric, a method for comparing image similarities essential for multiscale  
182 performance. These attention mechanisms also significantly improved the Mean Squared Error (MSE)  
183 metric, which quantifies the average squared differences between estimated and actual values.

184 Our findings underscore the potential of exploring alternative attention mechanisms to enhance  
185 efficiency and performance. The xFormers library, with its diverse attention options, provides an  
186 opportunity to tailor attention mechanism selection to specific applications, leading to substantial  
187 performance gains.

## 188 References

- 189 [1] I. Beltagy, M. E. Peters, and A. Cohan. Longformer: The long-document transformer, 2020.
- 190 [2] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin. Emerging properties in  
191 self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer  
192 Vision*, pages 9650–9660, 2021.
- 193 [3] X. Chen and K. He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF  
194 Conference on Computer Vision and Pattern Recognition*, pages 15750–15758, 2021.
- 195 [4] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick. Masked autoencoders are scalable vision  
196 learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages  
197 16000–16009, 2022.
- 198 [5] A. Hore and D. Ziou. Image quality metrics: Psnr vs. ssim. In *2010 20th international conference on pattern  
199 recognition*, pages 2366–2369. IEEE, 2010.
- 200 [6] J. Lee-Thorp, J. Ainslie, I. Eckstein, and S. Ontanon. Fnet: Mixing tokens with fourier transforms, 2022.
- 201 [7] B. Lefaudeux, F. Massa, D. Liskovich, W. Xiong, V. Caggiano, S. Naren, M. Xu, J. Hu, M. Tintore,  
202 S. Zhang, P. Labatut, and D. Haziza. xformers: A modular and hackable transformer modelling library.  
203 <https://github.com/facebookresearch/xformers>, 2022.
- 204 [8] Z. Wang, E. Simoncelli, and A. Bovik. Multiscale structural similarity for image quality assessment. In *The  
205 Thirty-Seventh Asilomar Conference on Signals, Systems Computers, 2003*, volume 2, pages 1398–1402  
206 Vol.2, 2003. doi: 10.1109/ACSSC.2003.1292216.
- 207 [9] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin. Unsupervised feature learning via non-parametric instance  
208 discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages  
209 3733–3742, 2018.