

767 **Appendix**

768 We provide all additional details for our paper in the follow-  
769 ing sections.

- 770 • **Border Impact.** We discuss the limitations and potential  
771 future follow-up work.
- 772 • **Details of the Implementation.** We provide additional  
773 details of model setup, training schedules.
- 774 • **Ablation Studies.** We provide additional ablation study re-  
775 sults, including masking strategies, model size, and object-  
776 mask ratio.
- 777 • **Discussions.** We address additional questions about the  
778 usage of additional data, the generalization capability of  
779 our proposed tokenization objective, as well as impact of  
780 auxiliary Gan loss.

781 **A. Broader Impact**

782 **Limitations and future work.** While our method improves  
783 semantic reasoning, there are still some failure cases (Fig-  
784 ure 8). For example, when using fine-grained object masking  
785 during pre-training—where the mask follows the exact shape  
786 of objects—the model may “cheat” by overfitting to the mask  
787 shape. In such cases, it quickly learns to fill in the masked  
788 area without acquiring meaningful representations. To re-  
789 solve this issue, we expand the mask to the bounding box.  
790 In future work, we aim to develop a more structured and ro-  
791 bust tokenizer to enhance the model’s reasoning capabilities.  
792 Our object masks are coarse and can be produced by multi-  
793 ple mechanisms; nevertheless, object discovery quality and  
794 compute cost remain practical considerations. In addition,  
795 we acknowledge the cost of segmentation overhead, but in  
796 our respectful opinion, our pipeline should be viewed as a  
797 proof-of-concept, and the performance gain is strong enough  
798 to justify studying it.

799 **Ethics Statement.** We ensure that our approach adheres to  
800 all legal and ethical guidelines throughout its development,  
801 with no violations. Fair compensation was provided to all  
802 annotators and graduate students involved in this work. The  
803 problems used in our study were collected from publicly ac-  
804 cessible exams<sup>1</sup> and resources licensed under CC Licenses<sup>2,3</sup>.  
805 This research is conducted solely for academic purposes,  
806 and we strictly prohibit any commercial use of the results.  
807 Additionally, the spurious captions generated in Section 4  
808 are limited to problem-solving contexts and pose no harm to  
809 individuals.

810 **Reproducibility statement.** We are committed to efficient  
811 and reproducible research. All code, datasets, and models  
812 will be publicly released.

<sup>1</sup><https://gate2025.iitr.ac.in/>

<sup>2</sup><https://www.allaboutcircuits.com/worksheets/>

<sup>3</sup><https://ocw.mit.edu/>

**B. Additional Implementation Details** 813

**Mask generation and preprocessing.** To efficiently gener- 814  
ate object masks, we leverage off-the-shelf [26], a popular 815  
unsupervised segmentation model, to infer scene-centric im- 816  
ages (where many objects are present). This step yields a 817  
set of binary object masks, which we then convert into the 818  
COCO RLE (Run-Length Encoding) format. Note that this 819  
step can be done either online (during the forward pass of 820  
each batch) or beforehand. Here we test both and empirically 821  
find the pre-processing step crucial as it saves roughly  $3\times$  822  
GPU hours as shown in Table 9. This solution is scalable 823  
as more data can be generated directly using the pre-trained 824  
SAM model. 825

Model	Pre-Processing	Training Cost
MIM (w. Obj Rep)	✓	3.6 ( $-2.7\times$ )
MIM (w. Obj Rep)	×	9.8
MIM+VQGAN(w. Obj Rep)	✓	5.1 ( $-2.5\times$ )
MIM+VQGAN(w. Obj Rep)	×	13.2

Table 9. Comparison of training costs in GPU hours with and without pre-processing for 1 epoch training using 500K data and a single A100 GPU.

**Implementation details on downstream tasks.** Following 826  
He et al. [24], we first discard the decoder after pre-training is 827  
complete. For end-to-end FT, we use AdamW [35] optimizer 828  
with base learning rate  $blr = 1.0 \times 10^{-3}$ , weight decay 0.05, 829  
layer decay 0.75 and train for 20 epochs with 5 rounds of 830  
warmup epochs. Additionally, we use drop path 0.1 with 831  
mixup 0.8 and ensure the effective batch size is 1024 by 832  
accumulating SGD iters. For LP, we use base learning rate 833  
 $blr = 1.0 \times 10^{-1}$  and an effective batch size of 16384 while 834  
keeping other settings the same. In our model, each self- 835  
attention layer includes  $\alpha = 16$  attention heads. 836

**Implementation details on pertaining.** For the first stage, 837  
we use AdamW [35] optimizer with a base learning of  $blr =$  838  
 $1.5 \times 10^{-4}$ , weight decay  $wd = 0.05$ , and the cosine learning 839  
rate decay scheduler. We accumulate iterations to emulate 840  
the recommended batch size of 4096 and pre-train the model 841  
for 25 epochs with 5 warmup epochs. During this stage, the 842  
mask ratio is set for  $mr_{patch} = 75\%$ . For the second stage, 843  
we start from the saved checkpoint from stage one. We apply 844  
an object ratio of  $mr_{obj} = 50\%$  which randomly masks 845  
out 25 objects in each image by hiding the patches spatially 846  
covering them. To enable batch processing, we apply an 847  
additional mask ratio constraint of  $mr_{patch} = 60\%$  on all 848  
images. The mask ratio is set 15% lower to accommodate 849  
increased difficulty in the objective. 850

Due to constraints in computing resources, we use pub- 851  
licly available pre-trained checkpoints<sup>4,5</sup> as the starting 852

<sup>4</sup><https://github.com/facebookresearch/mae>

<sup>5</sup>[https://github.com/amirbar/visual\\_prompting](https://github.com/amirbar/visual_prompting)

Model	FT (%)	LP (%)
MIM <sup>†</sup> [24]	83.66	70.80
SemMAE <sup>†</sup> [28]	83.73	71.25
MIM (wo. Obj Rep)	67.72 ↓15.94	58.75 ↓12.05
MIM (w. Obj Rep)	<b>84.43 ↑0.77</b>	<b>71.91 ↑1.11</b>

Table 10. **Linear probing (LP) and finetuning (FT) results on ImageNet-1K.**

853 model for both stages of pre-training, unless otherwise speci-  
 854 fied. Importantly, using pre-trained checkpoints does not  
 855 undermine our objective, as they are trained with a patch-level  
 856 objective, which aligns with the first stage of our framework  
 857 for learning low-level representations (Two Stage Learning  
 858 Section). Essentially, we retrain these models on a different  
 859 dataset with some adaptations.

860 **Loss function for MIM-VQGAN.** MIM-VQGAN was pro-  
 861 posed by Bar et al. [5] to study the effectiveness of visual  
 862 prompting, which effectively shifted the MIM evaluation  
 863 paradigm from fine-tuning on downstream tasks to direct  
 864 output generation via prompting. This can be seen as a  
 865 unified framework for vision tasks. Unlike He et al. [23],  
 866 which computes the MSE loss by directly regressing on pixel  
 867 values, MIM-VQGAN instead computes the cross-entropy  
 868 (CE) loss on the corresponding patch value in the quantized  
 869 codebook. This design effectively alleviates ambiguity in  
 870 generation, as the codebook is discrete, unlike pixel val-  
 871 ues. Notably, the underlying objective—masked autoencod-  
 872 ing—remains unchanged. Hence, MIM-VQGAN provides  
 873 an effective way to directly compare our proposed method.  
 874 In our experiments, we follow the implementation of Bar  
 875 et al. [5].

## 876 C. Additional Ablation Study.

877 **Influence of different object masking strategies:** As  
 878 shown in Figure 9 and Figure 10, we evaluate reconstruction  
 879 performance using three masking strategies: masking  
 880 strictly based on the object shape, masking the square re-  
 881 gion of the object, and a combination of both. While these  
 882 visualizations demonstrate the superiority of object-based  
 883 masking compared to random masking strategies, they also  
 884 reveal certain limitations. Specifically, relying solely on ob-  
 885 ject shape masking can lead to the model overfitting to the  
 886 mask shape (“cheating”), while using only square masking  
 887 results in sub-optimal performance on details. By combining  
 888 these two strategies, we achieve more realistic and effective  
 889 reconstruction.

890 **Study on how the model captures context:** We investigate  
 891 and visualize if our model has learned to capture the context  
 892 during the pretraining process. Here we focus on learning  
 893 the “shape” and “color”, two of the most important ingredi-  
 894 ents to human learning. As we have addressed learning the  
 895 “shape” in Figure 5 and Discussion Section, we showcase  
 896 the learning of color in Figure 7. In this example, when the

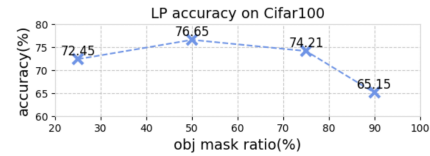


Figure 6. **Effect of object mask ratio:** The number of objects masked out during masked image modeling.

Model	Backbone	Cifar100 Top-1 Acc (%)	
		FT	LP
MIM <sup>†</sup>	ViT-B	89.98	75.01
MIM <sup>†</sup>	ViT-L	92.67	76.20
MIM (w. Obj Rep)	ViT-B	90.08	72.44
MIM (w. Obj Rep)	ViT-L	93.77	76.65

Table 11. Comparison of different model sizes. Results show our approach is able to scale with model size.

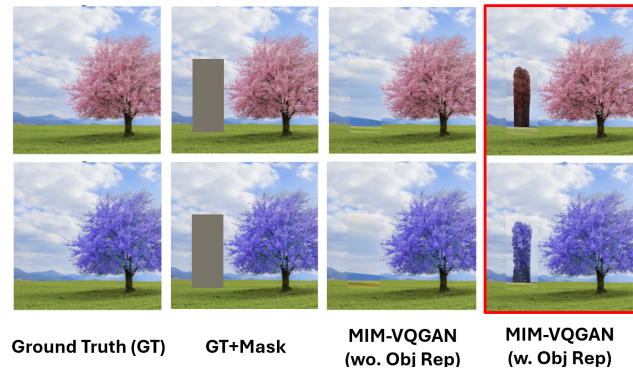


Figure 7. **Extend of color learning example**

897 same pair of examples but with different colors is given to  
 898 the model, it is able to reconstruct objects of colors similar  
 899 to the example, meaning that it does not infer color based on  
 900 memorization but rather from the context that is given.

901 **Study on model sizes:** Table 11 shows the LP and FT re-  
 902 sults on different vit base models, and the result shows our  
 903 observations and findings in Quantitative Evaluation and  
 904 Discussion sections hold for different model sizes.

905 **Obj-Mask Ratio.** To determine the influence of the masking  
 906 strategy, we train our model with different mask ratios, as  
 907 shown in Figure 6. Unlike traditional random patch-level  
 908 masking, as in He et al. [24], object-level masking becomes  
 909 less effective when obj-mask ratios exceed 50%. This de-  
 910 cline occurs because random masking often leaves portions  
 911 of objects visible, which can help guide reconstruction, while  
 912 object-level masking requires the model to learn the seman-  
 913 tic relationships between objects only from other objects.  
 914 We note that a 50% obj-mask ratio effectively masks out  
 915 around 75% of the image.

916 **Loss functions.** We further ablate the effect of object bal-  
 917 ance loss defined in Equation 7. Results in Table 12 shows  
 918 that combining both  $\mathcal{L}_{MIM}$  and  $\mathcal{L}_{obj}$  achieves the best per-  
 919 formance.

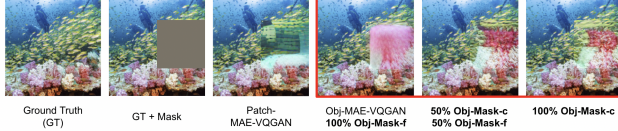


Figure 8. **Failure Cases:** (4): Failure case of reconstruction with fine-grained object masking (Obj-Mask-f). (5)-(6): Remedy by using coarse object masking (Obj-Mask-c)

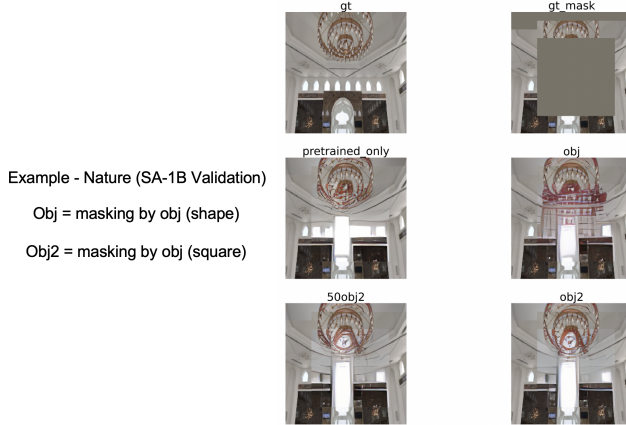


Figure 10. **Ablation Study of Masking Strategies (B)**

Model Variant	VQA (v2.0) Acc. (%)
MIM (w. Obj Rep)	53.02
+ $\mathbb{L}_{MIM}$ only	55.44
+ $\mathbb{L}_{obj}$ only	52.48
+ $\mathbb{L}_{MIM}$ + $\mathbb{L}_{obj}$ (Eq. 7)	<b>56.89</b>

Table 12. Effect of adding different loss terms in Eq. 7 on VQA (v2.0). Combining both  $\mathbb{L}_{MIM}$  and  $\mathbb{L}_{obj}$  achieves the best performance.

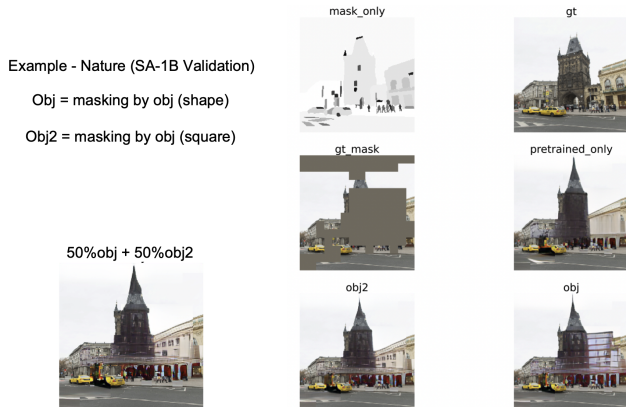


Figure 9. **Ablation Study of Masking Strategies (A)**

920

## D. Additional Discussions.

921

**Model size.** Here we show LP results on Cifar-100 classification with ViT-B and ViT-L. Table 11 indicates that our approach is scalable with respect to increasing model sizes.

922

923

**Additional motivation for using object-level representation.** Besides computer vision research, neuroscience studies

924

925



Figure 11. **GAN loss** can further help with better details.

have also found that the human brain uses an object-centric approach for visual recognition [6, 7, 38]. Within computer vision research, object segmentations have also been found to be helpful for tasks such as instance segmentation [19] and weakly supervised learning [55]. Hence, we conjecture “object” as a plausible candidate and explore it as the masking unit in MAE by simply masking out random objects and inpainting them instead of random patches.

926

927

928

929

930

931

932

933

**Generalizability of object-centric objective.** The surprising result is that while Patch-MAE severely degrades downstream fine-tuning performance, Obj-MIM can recover such gap in a short GPU-hour, demonstrating that object-centric learning objective enables the learning of highly semantic and generalizable features where the original Patch-MIM cannot, especially given the underlying semantic difference (domain gap) between the datasets.

934

935

936

937

938

939

940

941

**Further enhancing visual details with Gan loss.** Generative adversarial networks (GAN) [20] learn representation through the competition of a generator and a discriminator. Recent studies show that adding GAN losses can enhance visual details [17, 24, 37, 47]. Following this intuition, we add an auxiliary GAN loss to our objective in Equation 7:

942

943

944

945

946

947

$$\mathbb{L}_{OBJ-MAE} = \mathbb{L}_{MAE} + \lambda_1 \cdot \mathbb{L}_{obj} + \lambda_2 \cdot \mathbb{L}_{GAN} \quad (8)$$

948

This can be achieved by adding a simple discriminator and using the original network as the generator; details can be found in the Appendix. Results in (Figure 11) confirm that GAN loss can help produce more detailed images.

949

950

951

952