

# CL-MAE: Curriculum-Learned Masked Autoencoders – Supplementary

Neelu Madan<sup>1,◇</sup>, Nicolae-Cătălin Ristea<sup>2,3,◇</sup>, Kamal Nasrollahi<sup>1,4</sup>,

Thomas B. Moeslund<sup>1</sup>, Radu Tudor Ionescu<sup>3,5,\*</sup>

<sup>1</sup>Aalborg University, Denmark, <sup>2</sup>University Politehnica of Bucharest, Romania,

<sup>3</sup>University of Bucharest, Romania, <sup>4</sup>Milestone Systems, Denmark, <sup>5</sup>SecurifAI, Romania

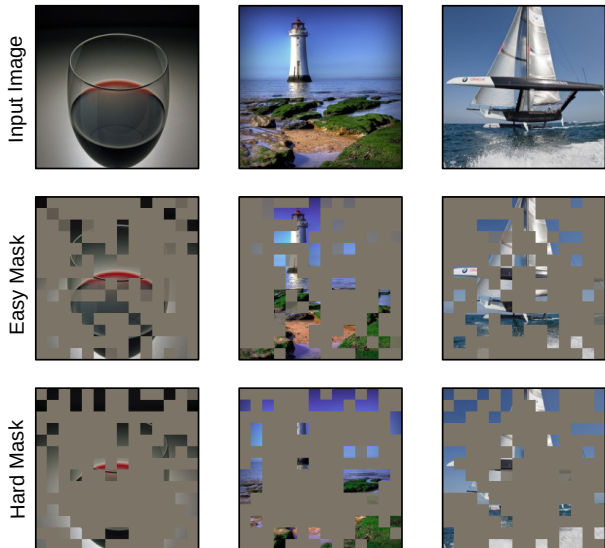


Figure 1. Masks generated by our masking module at two different moments during training, when all losses are in place, for the images on the top row. The masks on the second row are generated halfway during training, when the masking module is still acting as a partner to the MAE. In contrast, the masks on the bottom row are generated in the last epoch, when the masking module is behaving as an adversary to the MAE. Our module shifts its preference from masking non-salient tokens to masking tokens situated on edges and object contours, generating an easy-to-hard curriculum for the MAE.

## 1. Additional Experiments and Discussion

In the supplementary, we analyze the behavior of our loss functions from a qualitative perspective, and present few-shot linear probing results on five downstream tasks.

**Effect of curriculum learning on masking.** In Figure 1, we provide some masks generated by our learnable masking module at two different states during the training process on ImageNet [5]. In the first state (second row), the masking module acts as a partner to the MAE, helping to ease the reconstruction task. In this state, the module seems to mask tokens from non-salient or plain texture regions,

\*corresp. author: raducu.ionescu@gmail.com; ◇equal contribution.

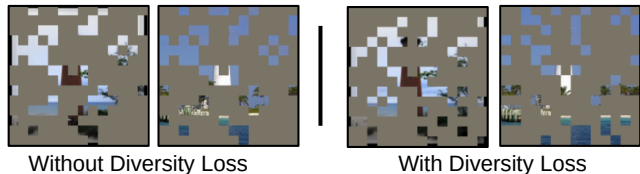


Figure 2. Masks generated by the proposed masking module without (left) and with (right) adding the diversity loss ( $\mathcal{L}_{div}$ ). If the diversity loss ( $\mathcal{L}_{div}$ ) is not included, the masking module can enter mode collapse and produce nearly identical masks. This can lead to overfitting CL-MAE on reconstructing certain patch configurations. The effect is no longer observed when the proposed diversity loss is employed.

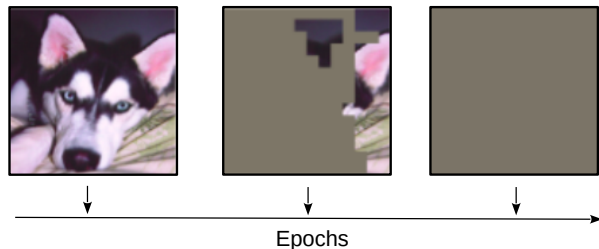


Figure 3. Masks generated by the proposed masking module during training, without the Kullback-Leibler loss. The masks evolve from leaving all patches visible (to reduce the reconstruction error for the MAE) to hiding all patches (to increase the reconstruction error for the MAE). The Kullback-Leibler loss is required to make sure the model always masks the desired number of patches.

which can be easily inferred from surrounding patches. We observe that the object contours are mostly visible in the easy masks. In the second state (third row), our masking module acts as an adversary to the MAE, aiming to make the reconstruction task harder for the MAE. In this state, we observe that the module prefers to mask patches located on edges, object contours and salient regions. These patches are much harder to reconstruct based on the visible information. In summary, the examples presented in Figure 1 confirm that our learnable masking module behaves as expected, creating an easy-to-hard curriculum for the MAE.

**Effect of diversity loss on masking.** Figure 2 shows the

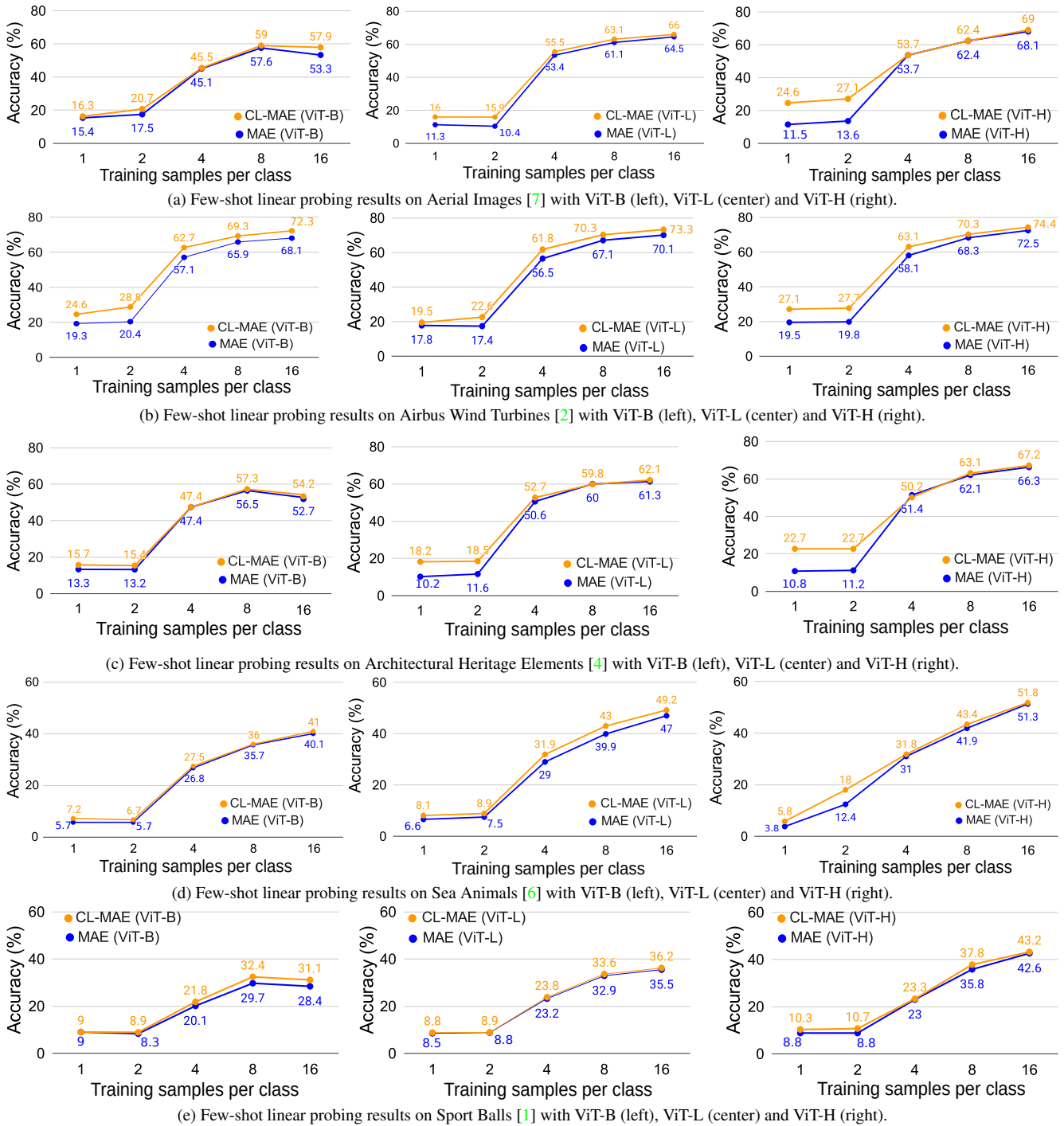


Figure 4. Few-shot linear probing results for MAE [3] and CL-MAE (ours) based on various backbones (ViT-B, ViT-L, ViT-H). The number of training samples per class varies between 1 and 16. The reported accuracy rates are averaged over three runs. Best viewed in color.

results before and after introducing the diversity loss ( $\mathcal{L}_{div}$ ) as an objective of our learnable masking module. When the masking module is trained without adding the diversity loss, the module enters mode collapse, producing nearly

identical masks, regardless of the input sample. This can be problematic, as CL-MAE might overfit to the task and learn to reconstruct only a certain configuration of patches. In contrast, integrating the diversity loss helps our module

Few-shot scenario	Method	Aerial Images		Airbus Wind Turbines		Architectural Heritage Elements		Sea Animals		Sport Balls	
		Acc@1	Acc@5	Acc@1	Acc@5	Acc@1	Acc@5	Acc@1	Acc@5	Acc@1	Acc@5
1-shot	MAE (ViT-B) [3]	15.4	48.6	19.3	54.3	13.3	47.2	5.7	20.6	<b>9.0</b>	32.6
	CL-MAE (ViT-B)	<b>16.3</b>	<b>58.9</b>	<b>24.6</b>	<b>63.9</b>	<b>15.7</b>	<b>61.7</b>	<b>7.2</b>	<b>26.7</b>	<b>9.0</b>	<b>34.5</b>
	MAE (ViT-L) [3]	11.3	39.8	17.8	56.2	10.2	40.5	6.6	27.7	8.5	34.9
	CL-MAE (ViT-L)	<b>16.0</b>	<b>50.7</b>	<b>19.5</b>	<b>63.8</b>	<b>18.2</b>	<b>55.8</b>	<b>8.1</b>	<b>31.6</b>	<b>8.8</b>	<b>36.9</b>
	MAE (ViT-H) [3]	11.5	54.1	19.5	57.2	10.8	51.6	3.8	<b>22.8</b>	8.8	34.8
	CL-MAE (ViT-H)	<b>24.6</b>	<b>58.2</b>	<b>27.1</b>	<b>59.0</b>	<b>22.7</b>	<b>56.0</b>	<b>5.8</b>	22.4	<b>10.3</b>	<b>35.7</b>
2-shot	MAE (ViT-B) [3]	17.5	53.1	20.4	65.3	13.2	51.4	5.7	24.6	8.3	34.9
	CL-MAE (ViT-B)	<b>20.7</b>	<b>60.3</b>	<b>28.8</b>	<b>69.9</b>	<b>15.4</b>	<b>60.2</b>	<b>6.7</b>	<b>25.8</b>	<b>8.9</b>	<b>37.2</b>
	MAE (ViT-L) [3]	10.4	39.6	17.4	60.1	11.6	42.9	7.5	<b>33.6</b>	8.8	32.9
	CL-MAE (ViT-L)	<b>15.9</b>	<b>51.3</b>	<b>22.6</b>	<b>66.2</b>	<b>18.5</b>	<b>53.5</b>	<b>8.9</b>	31.9	<b>8.9</b>	<b>34.7</b>
	MAE (ViT-H) [3]	13.6	50.7	19.8	61.7	11.2	51.5	12.4	43.0	8.8	34.9
	CL-MAE (ViT-H)	<b>27.1</b>	<b>58.8</b>	<b>27.7</b>	<b>69.2</b>	<b>22.7</b>	<b>55.8</b>	<b>18.0</b>	<b>52.2</b>	<b>10.7</b>	<b>35.0</b>
4-shot	MAE (ViT-B) [3]	45.1	83.3	57.1	89.8	<b>47.4</b>	87.3	26.8	62.5	20.1	58.0
	CL-MAE (ViT-B)	<b>45.5</b>	<b>89.9</b>	<b>62.7</b>	<b>91.5</b>	<b>47.4</b>	<b>91.6</b>	<b>27.5</b>	<b>63.2</b>	<b>21.8</b>	<b>58.8</b>
	MAE (ViT-L) [3]	53.4	87.2	56.5	88.1	50.6	88.7	29.0	65.3	23.2	58.4
	CL-MAE (ViT-L)	<b>55.5</b>	<b>91.8</b>	<b>61.8</b>	<b>90.4</b>	<b>52.7</b>	<b>89.9</b>	<b>31.9</b>	<b>66.4</b>	<b>23.8</b>	<b>59.4</b>
	MAE (ViT-H) [3]	<b>53.7</b>	90.8	58.1	89.9	<b>51.4</b>	89.3	31.0	65.7	23.0	<b>58.5</b>
	CL-MAE (ViT-H)	<b>53.7</b>	<b>90.9</b>	<b>63.1</b>	<b>92.0</b>	50.2	<b>91.4</b>	<b>31.8</b>	<b>67.0</b>	<b>23.3</b>	<b>58.5</b>
8-shot	MAE (ViT-B) [3]	57.6	95.9	65.9	97.8	56.5	<b>96.6</b>	35.7	70.6	29.7	65.5
	CL-MAE (ViT-B)	<b>59.0</b>	<b>96.5</b>	<b>69.3</b>	<b>98.2</b>	<b>57.3</b>	96.3	<b>36.0</b>	<b>71.5</b>	<b>32.4</b>	<b>69.4</b>
	MAE (ViT-L) [3]	61.1	97.2	67.1	97.9	<b>60.0</b>	97.4	39.9	74.7	32.9	68.1
	CL-MAE (ViT-L)	<b>63.1</b>	<b>97.4</b>	<b>70.3</b>	<b>98.7</b>	59.8	<b>97.7</b>	<b>43.0</b>	<b>76.3</b>	<b>33.6</b>	<b>69.4</b>
	MAE (ViT-H) [3]	<b>62.4</b>	<b>96.3</b>	68.3	97.9	62.1	95.7	41.9	76.2	35.8	70.1
	CL-MAE (ViT-H)	<b>62.4</b>	96.2	<b>70.3</b>	<b>98.6</b>	<b>63.1</b>	<b>96.6</b>	<b>43.3</b>	<b>76.5</b>	<b>37.8</b>	<b>72.1</b>
16-shot	MAE (ViT-B) [3]	53.3	93.7	68.1	97.4	52.7	<b>93.1</b>	40.1	73.9	28.4	67.6
	CL-MAE (ViT-B)	<b>57.9</b>	<b>96.8</b>	<b>72.3</b>	<b>99.0</b>	<b>54.2</b>	92.7	<b>41.0</b>	<b>76.4</b>	<b>31.1</b>	<b>69.0</b>
	MAE (ViT-L) [3]	64.5	95.1	70.1	97.8	61.3	95.4	47.0	81.1	35.5	73.7
	CL-MAE (ViT-L)	<b>66.0</b>	<b>97.3</b>	<b>73.3</b>	<b>99.1</b>	<b>62.1</b>	<b>95.5</b>	<b>49.2</b>	<b>81.9</b>	<b>36.2</b>	<b>75.3</b>
	MAE (ViT-H) [3]	68.1	96.7	72.5	98.2	66.3	<b>96.7</b>	51.3	82.9	42.6	76.9
	CL-MAE (ViT-H)	<b>69.0</b>	<b>98.1</b>	<b>74.4</b>	<b>99.1</b>	<b>67.2</b>	96.6	<b>51.8</b>	<b>82.7</b>	<b>43.2</b>	<b>78.1</b>

Table 1. Few-shot linear probing results on five benchmarks: Aerial Images, Airbus Wind Turbines, Architectural Heritage Elements, Sea Animals, and Sport Balls. The results are reported for MAE [3] and CL-MAE (ours) based on various backbones (ViT-B, ViT-L, ViT-H). The reported accuracy rates are averaged over three runs. The top scores for each backbone on each data set are in bold.

to escape mode collapse and generate diverse masks, which solves the problem of overfitting to certain mask configurations.

**Effect of Kullback-Leibler loss on masking.** Figure 3 illustrates three masks generated by our masking module during training, illustrating the effect of removing the Kullback-Leibler loss. At the beginning of the training process, when the module aims to make the reconstruction task easy for the MAE, it learns to leave all patches visible. As the training progresses, the module starts behaving like an adversary, aiming to make the reconstruction task hard for the MAE. From this point on, our module starts masking a number of tokens until it converges to masking all tokens. At this point, the MAE has literally no chance at reconstructing the input, being unable to further learn any useful information. To mitigate this issue, we add the Kullback-Leibler loss, which ensures the number of masked tokens complies with the desired ratio given as hyperparameter, ir-

respective of the complexity of the task.

**Few-shot linear probing results.** In Figure 4, we present few-shot linear probing results on Aerial Images [7], Airbus Wind Turbines [2], Architectural Heritage Elements [4], Sea Animals [6], and Sport Balls [1] data sets. We complement the graphs illustrated in Figure 4 with the results reported in Table 1. These experiments are meant to quantify the transferability of the self-supervised representations learned by MAE [3] and CL-MAE in the few-shot learning scenario. In 138 out of 150 cases, CL-MAE outperforms MAE, with absolute gains varying between +0.1% and +13.5%. We generally observe that CL-MAE tends to bring higher gains for the 1-shot and 2-shot scenarios, especially when the ViT-L and ViT-H backbones are applied on the Aerial Images [7], Airbus Wind Turbines [2], Architectural Heritage Elements [4], and Sea Animals [6] data sets. Overall, the few-shot experiments confirm the observations on the nearest neighbor and linear probing experiments pre-

sented in the main article. We thus conclude that our curriculum learning approach represents a useful addition to the MAE framework.

## References

- [1] Samuel Cortinhas. Sport Balls - Multiclass Image Classification. <https://www.kaggle.com/datasets/samuelcortinhas/sports-balls-multiclass-image-classification>, 2022. Accessed: 2023-08-28. [2](#), [3](#)
- [2] Airbus DS GEO. Airbus Wind Turbines Patches. <https://www.kaggle.com/datasets/airbusgeo/airbus-wind-turbines-patches>, 2022. Accessed: 2023-08-28. [2](#), [3](#)
- [3] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of CVPR*, pages 16000–16009, 2022. [2](#), [3](#)
- [4] Ivan Kobzev and Vasiliev Roman. Architectural Heritage Elements Image Dataset. <https://www.kaggle.com/datasets/ikobzev/architectural-heritage-elements-image64-dataset>, 2021. Accessed: 2023-08-28. [2](#), [3](#)
- [5] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. [1](#)
- [6] Lanz Vencer. Sea animals image dataset. <https://www.kaggle.com/datasets/vencerlanz09/sea-animals-image-dataste>, 2023. Accessed: 2023-08-28. [2](#), [3](#)
- [7] Gui-Song Xia, Jingwen Hu, Fan Hu, Baoguang Shi, Xiang Bai, Yanfei Zhong, Liangpei Zhang, and Xiaoqiang Lu. AID: A Benchmark Dataset for Performance Evaluation of Aerial Scene Classification. *IEEE Transactions on Geoscience and Remote Sensing*, 55(7):3965–3981, 2017. [2](#), [3](#)