

Learning the Heterogeneous Representation of Brain's structure from Serial SEM Images Using a Masked Autoencoder

Ao Cheng^{1,2}, Lirong Wang^{1,2,*}, and Ruobing Zhang^{2,3,*}

¹ School of Electronic and Information Engineering, Soochow University, Suzhou 215009, China

² Jiangsu Key Laboratory of Medical Optics, Suzhou Institute of Biomedical Engineering and Technology, Chinese Academy of Sciences, Suzhou 215163, China

³ Institute of Artificial Intelligence, Hefei Comprehensive National Science Center, Hefei 230088, China

Correspondence*:

Lirong Wang
wanglirong@suda.edu.cn

Ruobing Zhang
zhangrb@sibet.ac.cn

2 ABSTRACT

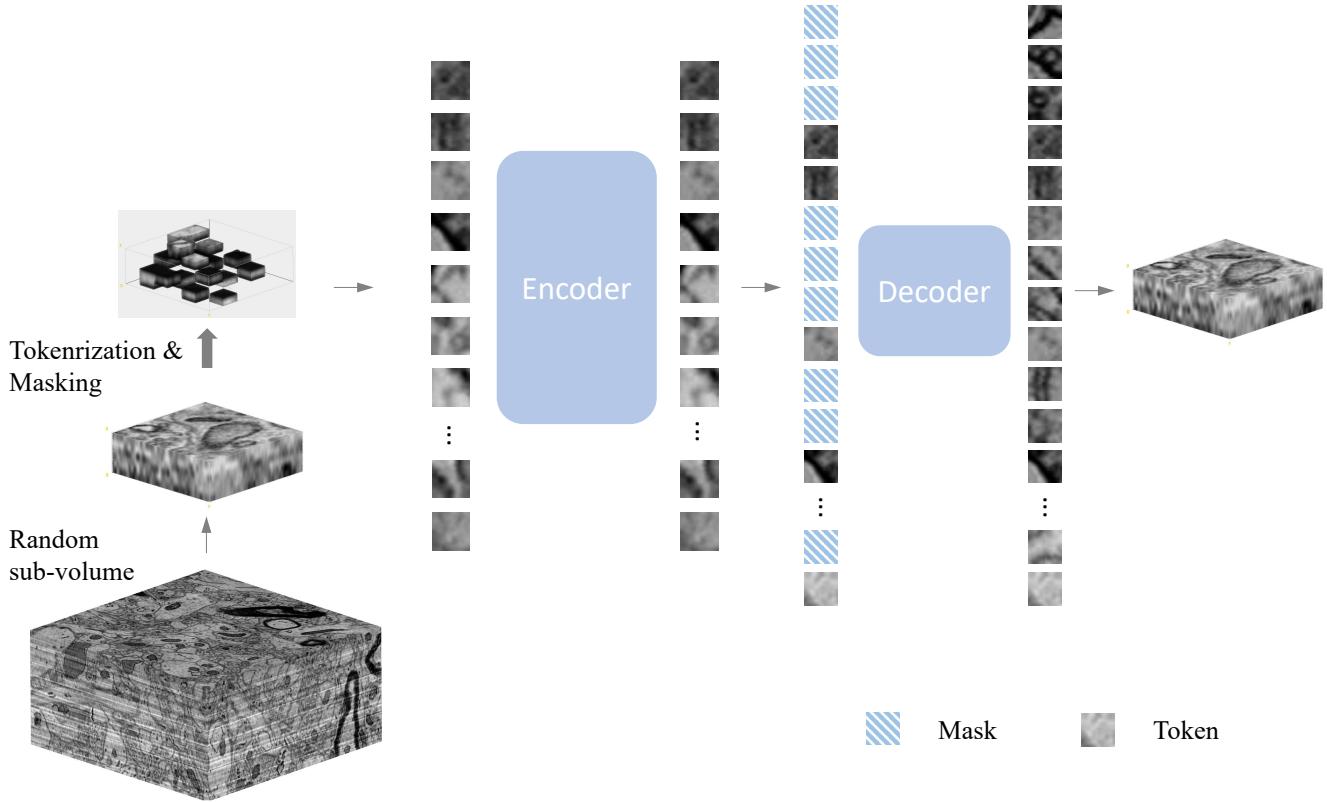
The exorbitant cost of accurately annotating the large-scale serial scanning electron microscope (SEM) images as the ground truth for training has always been a great challenge for brain map reconstruction by deep learning methods in neural connectome studies. The representation ability of the model is strongly correlated with the number of such high-quality labels. Recently, the masked autoencoder (MAE) has been shown to effectively pre-train Vision Transformers (ViT) to improve their representational capabilities. In this paper, we investigated a self-pre-training paradigm for serial SEM images with MAE to implement downstream segmentation tasks. We randomly masked voxels in three-dimensional brain image patches and trained an autoencoder to reconstruct the neuronal structures. We tested different pre-training and fine-tuning configurations on three different serial SEM datasets of mouse brains, including two public ones, SNEMI3D(Lee et al., 2017) and MitoEM-R, and one acquired in our lab. A series of masking ratios were examined and the optimal ratio for pre-training efficiency was spotted for 3D segmentation. The MAE pre-training strategy significantly outperformed the supervised learning from scratch. Our work shows that the general framework of MAE(He et al., 2021) can be a unified approach for effective learning of the representation of heterogeneous neural structural features in serial SEM images to greatly facilitate brain connectome reconstruction.

Keywords: Image segmentation, Neural segmentation, Masked autoencoder, Self-supervised learning, SEM image

1 INTRODUCTION

Three-dimensional segmentation of neural structures in serial scanning electron microscope (SEM) images is one of the core tasks of brain connectomic studies.(Kasthuri et al., 2015; Eberle et al., 2018). Supervised deep learning with annotations as ground truth (e.g., U-Net(Ronneberger et al., 2015)) has become the

Figure 1. Illustration of the masked autoencoder. Firstly, we randomly sample a sub-volume from the volume dataset as training data. Next, we randomly mask voxels with a specific masking ratio after the patch embedding operation. Afterward, the encoder processes all the unmasked voxels. A smaller decoder operates the full set of voxels, which includes the masked voxels as well, outputs the final reconstruction result. Note that the value of masked voxels here is given the value 0 because our goal is to predict those masked voxels by encoded voxels.

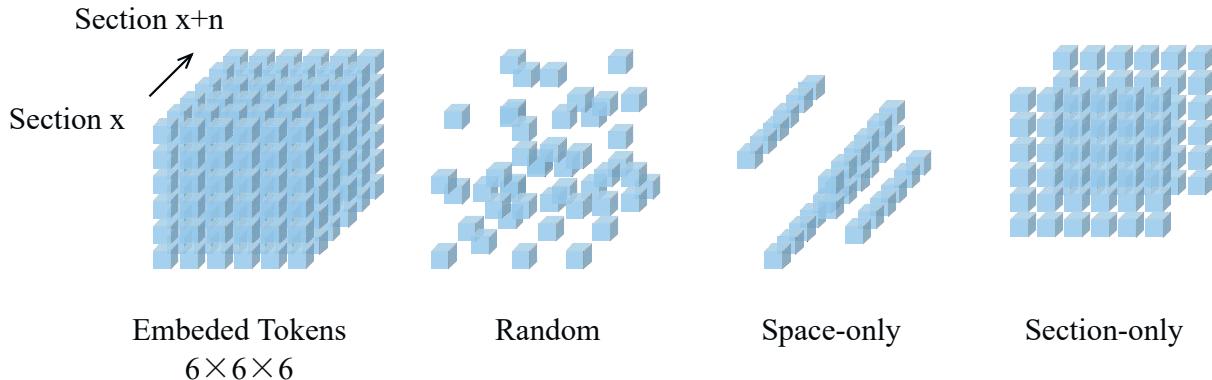


23 dominant approach for reconstruction process. However, due to the exorbitant cost of annotation, this
24 method is not friendly to large-scale serial image segmentation tasks.

25 As a feasible alternative, self-supervised learning acquires supervised information from the data itself and
26 has recently been shown to successfully address the need for data and be able to learn dense representations
27 of the input (Hung et al., 2018; Mittal et al., 2019; Lin et al., 2020; He et al., 2021). For the pretext
28 tasks, masked image modeling is such a pre-training learning task to enhance the representation capability:
29 mask part of the input information and try to predict the masked information. This paradigm has been
30 very successful in NLP, as self-supervised learning algorithms based on masked language modeling tasks
31 have revolutionized the discipline. Methods such as BERT(Devlin et al., 2018) and GPT(Radford et al.,
32 2018, 2019) have demonstrated that they can learn on unlabeled text data and are suitable for a variety of
33 applications. With the introduction of Vision Transformers(ViT) (Vaswani et al., 2017), Masked autoencoder
34 (MAE)(He et al., 2021) is also used to enhance the representation ability of self-attention mechanism
35 models(Wei et al., 2022; Xie et al., 2022; He et al., 2021). Following this philosophy, state-of-the-art
36 methods based on MAE have demonstrated their effectiveness in developing vision models.

37 Other common self-supervised methods on downstream tasks aim to exploit existing labels for unlabeled
38 domains. One approach is to use discriminator constraints on the spatial distribution of predictions on

Figure 2. Demonstration of mask sampling strategies. In this illustration, our embedded dimension ($Z \times X \times Y$) is $6 \times 6 \times 6$. Blue cubes represent embedded tokens. Random sampling is an agnostic spatial-wise sampling strategy. Space-only random sampling masks the tokens to all sections, and section-only random sampling masks random sections.



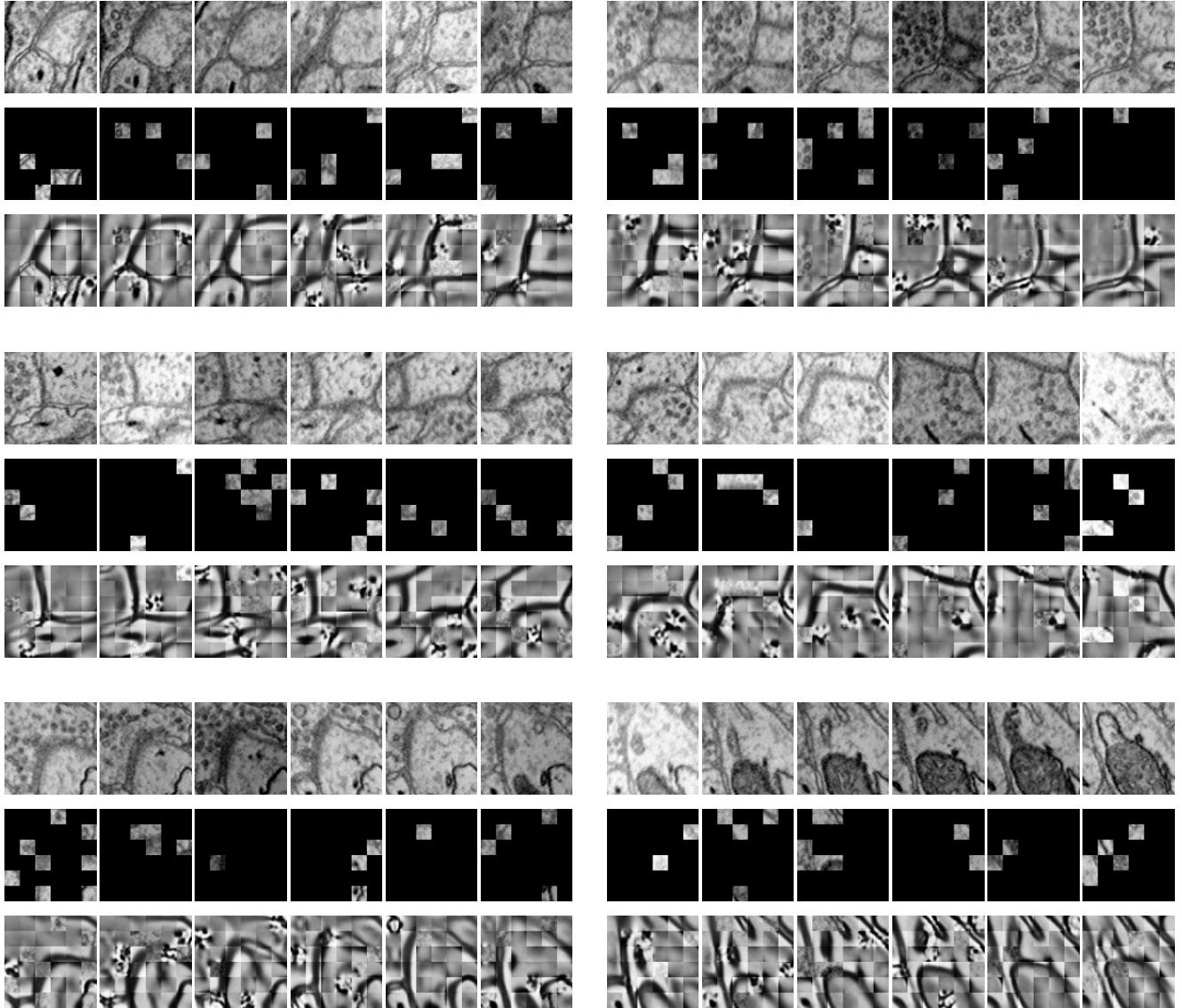
39 unlabeled images to improve model accuracy(Hung et al., 2018). Another approach uses unpaired image-
 40 to-image translation models, such as CycleGAN(You et al., 2020) on MRI images, CySGAN(Lauenburg
 41 et al., 2022) on SEM images, to domain-shift the dataset. However, regardless of the downstream task, the
 42 segmentation relies on an optimized translation model, and several independent modules also increase the
 43 complexity of the pipeline.

44 To address these challenges, we implement MAE(He et al., 2021) as a model pre-text task that the
 45 pre-trained backbone can be utilized for downstream tasks. Currently, since MAE(He et al., 2021) method
 46 has not demonstrated feasibility in 3D electron microscope images, and its applicability has not been
 47 thoroughly investigated. Therefore, we want to investigate the unified pre-training paradigm for different
 48 tasks on the 3D electron microscope images. Naturally, we wondered whether MAE(He et al., 2021) would
 49 also advance 3D electron microscopy image analysis. In this work, we aim to address this problem through
 50 the following attempts: the improvements in accuracy compared with training from scratch; the evaluation
 51 results on the public datasets.

2 RELATED WORK

52 Self-supervised learning approaches focus on learning representations from unlabeled data to achieve
 53 high precision, high accuracy, and rich representations. Transfer learning from natural images is used for
 54 medical image analysis regardless of differences in image statistics, scale, and task-related features.
 55 Raghu et al.(Raghu et al., 2019) showed that transfer learning from imageNet can accelerate the
 56 convergence of medical images, which is particularly useful when medical image training data is limited.
 57 In electron microscope images, transfer learning using domain-specific data can also help address domain
 58 differences and reduce labeling costs(Januszewski and Jain, 2019; Lauenburg et al., 2022). januszewski et
 59 al.(Januszewski and Jain, 2019) migrated the pre-trained segmentation model to the target data without
 60 labels so that the more accurate pseudo labels of the new dataset can be obtained directly. Lauenburg
 61 et al. (Lauenburg et al., 2022) proposed additional self-supervised and segmentation-based adversarial
 62 objectives in addition to the two steps of domain translation and image segmentation. Although this strategy
 63 effectively improves the representation ability of the model, it requires a part of the label as a constraint,
 64 and this data is expensive and time-consuming to collect. Besides, these domain-based self-supervised

Figure 3. Demonstration of masked results on the SNEMI3D dataset(Lee et al., 2017). Each row represents the image volume (top), masked volume (middle), and the final prediction (bottom). The image volume has a size of $6 \times 96 \times 96$ and is embedded with patch size $1 \times 16 \times 16$. After patch embedding, we obtain 216 tokens in total with the shape of $6 \times 6 \times 6$. Since the masking ratio is up to 90%, there are only 21 visible tokens for further encoding.



65 learning are difficult to combine with each other. Recent improvements in self-supervised learning offer a
66 feasible alternative, allowing specific representations to be learned to use unlabeled data, which is massive
67 and often more accessible.

68 The masked autoencoder is a self-supervised learning method that learns representations from the image
69 itself. DAE(Vincent et al., 2008, 2010) is a pioneering work in this field that presents masking as a type of
70 noise. It develops with the MLM task in NLP, the most representative is BERT(Devlin et al., 2018). In
71 the field of CV, such methods continue to develop and have proven effective (Dosovitskiy et al., 2021;
72 He et al., 2021; Wei et al., 2022; Xie et al., 2022; Pathak et al., 2016). Recent methods are based on the

73 transformer(Vaswani et al., 2017) structure, which is a self-attention-based model capable of solving image
74 and language tasks.

3 PROPOSED METHOD

75 As shown in Figure.1, our method is an extension of MAE(He et al., 2021) to 3D electron microscopy image
76 data. Our objective is to develop methods that are applicable to electron microscopy images under a general
77 and unified framework. Masked Image Modeling typically masks parts of the input image or encoded
78 image tokens and promotes the model to reconstruct the masked regions. Many existing Masked Image
79 Modeling methods employ an encoder-decoder design followed by a prediction head, such as BEiT(Bao
80 et al., 2021) and MAE(He et al., 2021). The encoder helps to pattern the latent feature representation, while
81 the decoder helps to process the latent features to the original image. Moreover, designing the decoder
82 components in a lightweight size minimizes training time. In our experience, lightweight decoders not
83 only reduce computational complexity, but also maximize the ability of the encoders to learn more general
84 representations. In this work, we thoroughly investigate the effectiveness of different MAE models on 3D
85 SEM image data. The following components provide more details:

86 3.1 Patch Embedding

87 Following the original ViT (Dosovitskiy et al., 2021), given a patch, we divide it into a regular grid of non-
88 overlapping blocks in space. These patches are flattened and embedded by linear projection(Dosovitskiy
89 et al., 2021). The positional embedding (Vaswani et al., 2017) is added to the embedded token. The token
90 and position embedding process is the only voxel-wise aware process. Unlike the 2D MAE(He et al.,
91 2021) design, due to the different spatial resolutions during imaging, we do not use down-sampling in the
92 z-direction, which ensures the 3D resolution of the voxel is close to a cube.

93 3.2 Masking

94 We randomly sample patches from the embedded patch set without replacement. This random sampling
95 is independent of spatial structure. As shown in Figure.2, the structure-independent random sampling
96 strategy is similar to the one-dimensional (Devlin et al., 2018) and two-dimensional (Wei et al., 2022; He
97 et al., 2021) methods. In (He et al., 2021), it is assumed that the optimal masking ratio is related to the
98 information redundancy of the data. For unstructured random masks, BERT (Devlin et al., 2018) uses a
99 masking ratio of 15% for languages, while MAE (He et al., 2021) uses a masking ratio of 75% for images,
100 indicating that images are more information redundancy. Our experimental results on patch data support
101 this hypothesis. The best masking ratio we observed for 3D MAE(He et al., 2021) on SEM images can
102 reach 90%. This is consistent with the general assumption that the 3D SEM data are spatially coherent and
103 more informative.

104 Figure.3 shows the results of our MAE reconstructing the masked data, with a masking ratio of 90%.
105 Spatial random sampling may be more efficient than structure-aware sampling strategies. Since voxels are
106 coherent, with a very high masking ratio, space-only or slice-only sampling may retain less information
107 and produce an overly difficult pre-training task. For example, 83.3% masking ratio with the slice-only
108 sampling of embedded dimension $6 \times 6 \times 6$ means that only one slice is maintained, which presents an
109 extremely challenging task of predicting other sections. We observe that the optimal masking ratio for
110 structure-aware sampling is generally lower. In contrast, spatial random sampling has higher efficiency on
111 the limited number of visible patches, thus allowing the use of a higher masking ratio.

112 3.3 Autoencoding

113 Our encoder is a vanilla ViT(Vaswani et al., 2017), applied only to visible embedded patches, following
 114 (He et al., 2021). This design greatly reduces time and memory complexity and leads to a more practical
 115 solution. A masking ratio of 90% reduces the encoder complexity to 1/10. Unlike SimmIM(Xie et al.,
 116 2022), MAE’s decoder is an encoded patch set and a set of masked tokens(He et al., 2021) concatenated
 117 with another set of vanilla ViT. Decoder-specific position embeddings are added to this set(He et al., 2021).
 118 Although both are Vit structures, the size of the decoder is designed to be smaller than the encoder (He
 119 et al., 2021). Moreover, the decoder handles the complete set, but it is originally less complex than the
 120 encoder. In addition, unlike the 2D MAE, we utilize 3D sin-cos similarity as our 3D MAE’s positional
 121 embedding to provide information about the spatial location.

122 We use the decoder to predict patches in the voxel space. We follow (He et al., 2021), predicting full
 123 spatial voxels (e.g., $Z \times 16 \times 16$) and the normalized value of each block of the original voxel. The training
 124 loss function is the mean-squared error (MSE) between the prediction and its target, averaged over unknown
 125 blocks(Devlin et al., 2018). This method relies on global self-attention to learn useful knowledge from the
 126 data, following (Dosovitskiy et al., 2021).

4 EXPERIMENT RESULTS

127 4.1 Implementation

128 Our encoder and decoder are the vanilla ViT architectures (Vaswani et al., 2017). We use a patch size of 1
 129 for the z-direction’s patch embedding, which follows the features of the SEM dataset. And we implement a
 130 space patch size of 16×16 (Dosovitskiy et al., 2021), denoted as $1 \times 16 \times 16$. We use the same patch size for
 131 ViT-B/L (Dosovitskiy et al., 2021) for simplicity. For a $6 \times 96 \times 96$ input, this patch size produces $6 \times 6 \times 6$
 132 tokens and is embedded with 3D positional embeddings for further encoding.

133 The 3D MAE pre-training configuration on SNEMI3D(Lee et al., 2017) is shown in Table.?.?. We use
 134 the AdamW optimizer (Kingma and Ba, 2014) with a batch size of 128 on 6 NVIDIA RTX3090 GPUs.
 135 We evaluate the pre-training quality by end-to-end fine-tuning. Furthermore, we remove the pre-trained
 136 decoder and implement UNETR(Hatamizadeh et al., 2022) as our decoding method. In the experiments
 137 of fine-tuning, we compare the different predicting targets, affinity maps(Lee et al., 2017), and multi-
 138 task predictions(Wei et al., 2020). The loss function of predicting affinity map follows the proposed
 139 configurations from (Lin et al., 2021; Lee et al., 2017). And the loss function for multi-task predictions
 140 following the configurations from(Lin et al., 2021; Wei et al., 2020). In addition, the following post-
 141 processing step for affinity maps and multi-task predictions are using the default configuration(Lin et al.,
 142 2021).

143 The SNEMI3D(Lee et al., 2017) leaderboard use adapted Rand F-score (A-Rand) (Nunez-Iglesias et al.,
 144 2013; Rand, 1971) as evaluation metrics. To show the significance of different methods, we demonstrated
 145 segmentation accuracy through variation of information (VI)(Bogovic et al., 2013) and adapted Rand
 146 F-score (Nunez-Iglesias et al., 2013; Rand, 1971). VI is defined as:

$$VI(S, T) = H(S | T) + H(T | S) \quad (1)$$

Figure 4. Illustration of the segmentation results. The first row is consecutive EM images from SNEMI3D(Lee et al., 2017). The second and fourth row represents the segmentation results of the model that is training from scratch, and the pre-training method results are shown in the third and fifth row. We use zwatershed as the post-processing step to generate segmentation results from the predicted affinity map. The post-processing algorithm of BCD (Binary maps, contours, distance) predictions are following the configuration from (Lin et al., 2021).

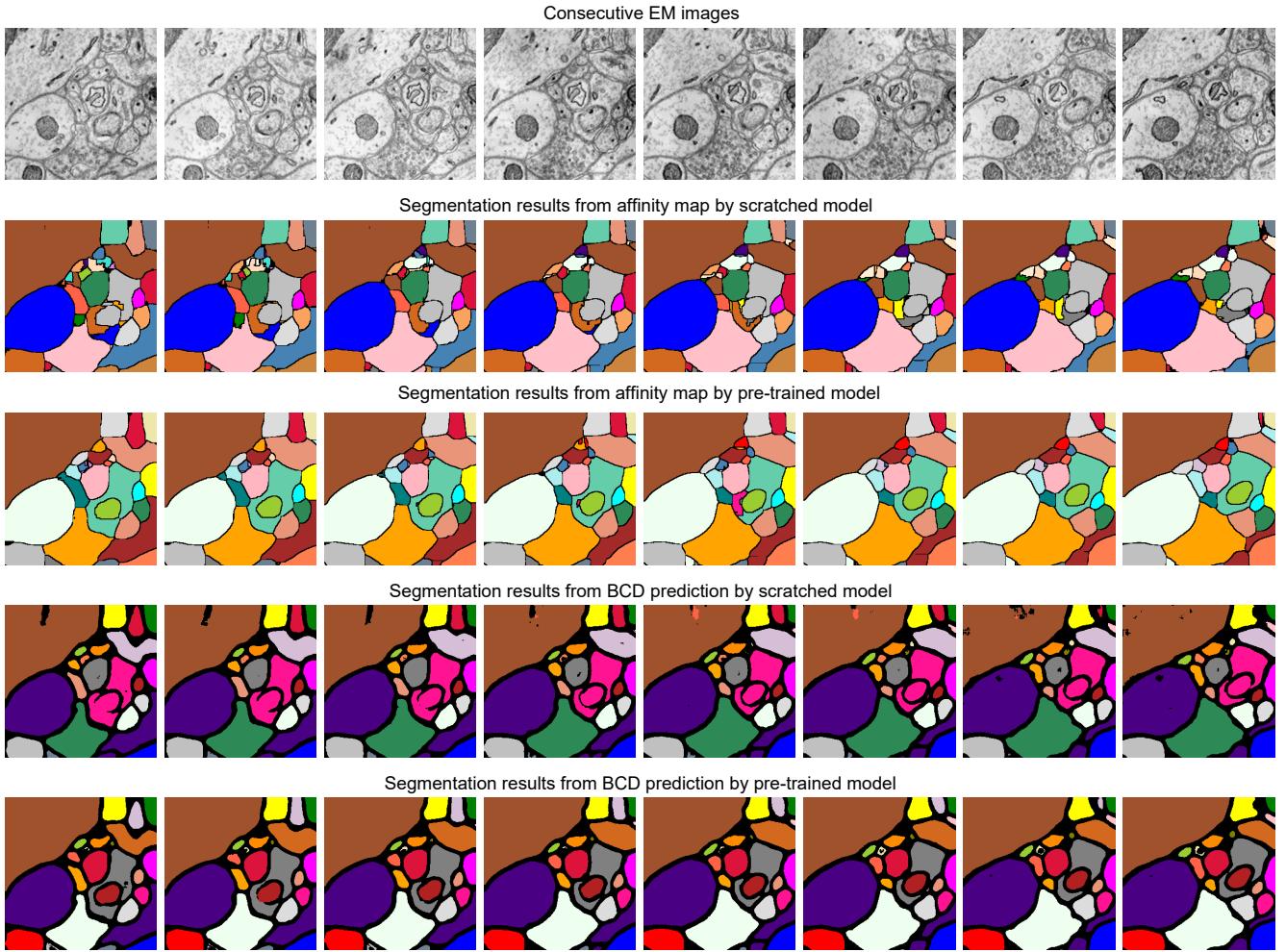


Table 1. Demonstration of mask Sampling. Figure.2 shows the random mask sampling results as well.

Mask	Ratio	Voi-S	Voi-M	A-Rand
Random	90%	0.234	0.129	0.071
Space	89%	0.271	0.140	0.079
Section	83%	0.313	0.156	0.103

147 Where S and T represents segmentation results and its related ground truth. Then the conditional entropy
 148 $H(S|T)$ measures oversegmentation errors (split error), and $H(T|S)$ measures undersegmentation errors
 149 (merger error). We defined split error and merge error as Voi-S and Voi-M respectively.

150 4.2 Ablation study

151 In this section, we assessed the model's pre-training performance across four aspects: sampling strategy,
 152 masking ratio, decoding depth, and decoding dimension. The ultimate fine-tuned models with different

Table 2. Demonstration of masking ratios with random mask sampling. In this table, the masking ratio is increased from 50% up to 90%, and the parameters of the decoder remain unchanged.

Mask	Ratio	Voi-S	Voi-M	A-Rand
Random	50%	0.253	0.131	0.081
Random	75%	0.243	0.116	0.073
Random	90%	0.234	0.129	0.071

Table 3. Demonstration of decoder depth. The masking ratio remains unchanged at 90%. An overlay decoder depth degrades the accuracy.

Mask	Depth	Voi-S	Voi-M	A-Rand
Random	2	0.264	0.157	0.083
Random	4	0.234	0.129	0.071
Random	8	0.256	0.143	0.079

Table 4. Demonstration of decoder dimension. The masking ratio remains unchanged at 90%. An overlay decoder dimension degrades the accuracy.

Mask	Dim.	Voi-S	Voi-M	A-Rand
Random	128	0.325	0.173	0.089
Random	256	0.317	0.166	0.085
Random	512	0.234	0.129	0.071

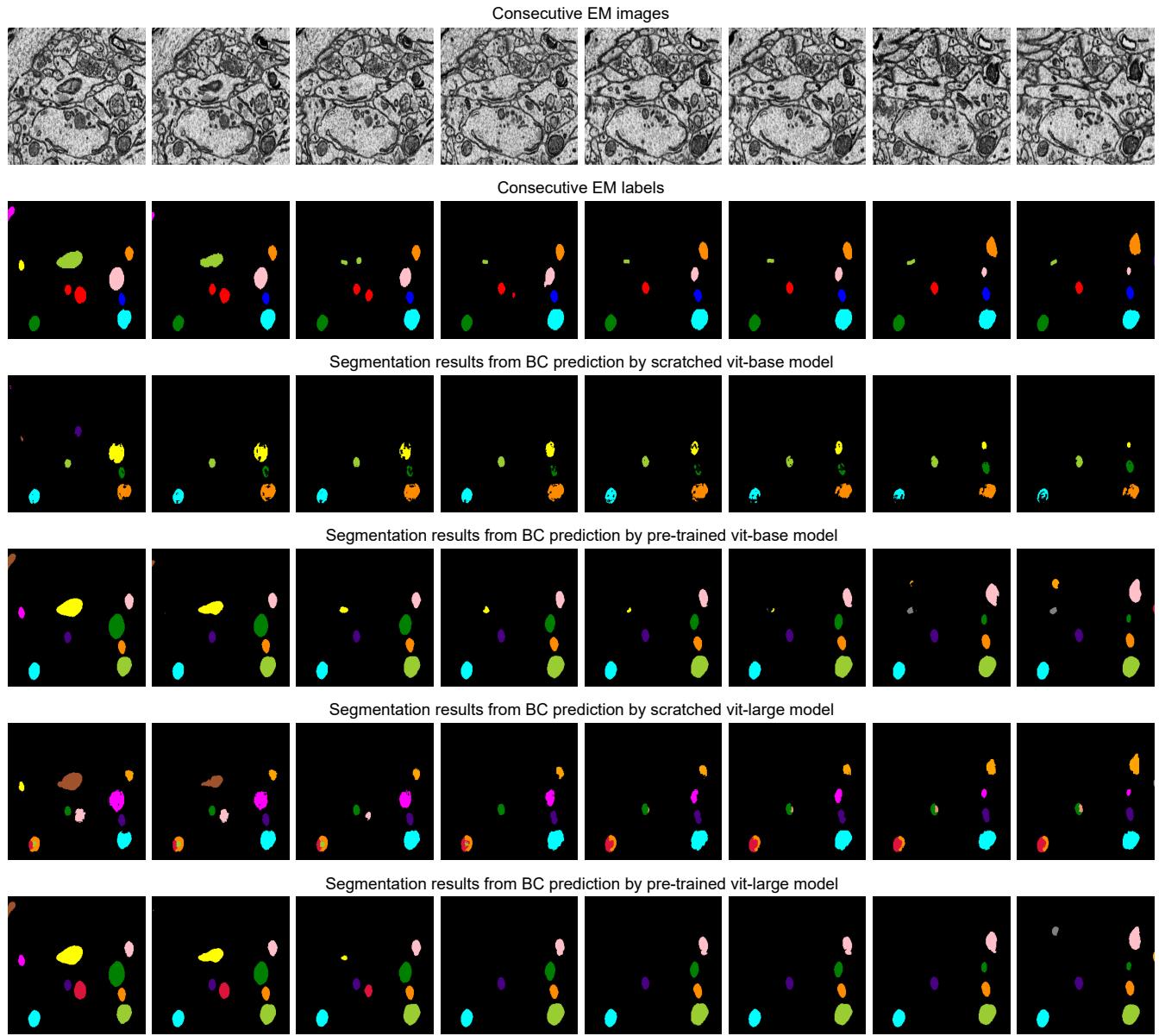
153 sampling methods are evaluated on the SNEMI3D dataset(Lee et al., 2017). Furthermore, we exclusively
 154 used Vit-Base for all ablation experiments and the pre-training dataset is the training dataset from
 155 SNEMI3D.

156 Table.1 shows the masking strategy between random, space-only, and section-only sampling. For a fair
 157 comparison of the masking strategy and masking ratio, we set the decoder depth and decoder dimension
 158 to 4 and 512 respectively. Moreover, in the experiment of the masking strategy, we decided to make
 159 the masking ratio as close as possible. We demonstrate different sampling strategies in Table.1. For the
 160 space-only sampling, we reserved $4 \times 6 = 24$ tokens, leading to a masking ratio up to 89%. This strategy
 161 performs close to random sampling. The masking ratio for the section-only sampling strategy is 83%, it
 162 processes one section's voxels. As shown in Table.1, the section-only sampling has the worst performance.
 163 Since this sampling strategy needs to predict the other 5 sections, It is hard to learn a general representation
 164 with the lack of information. Random sampling strategy has the best performance with the highest masking
 165 ratio, which gains 0.071 of A-Rand.

166 To have a more comprehensive look at the random sampling strategy, we analyzed the impact on the
 167 mask ratio. In this part, we only change the masking ratio and keep the decoder depth and dimension into 4
 168 and 512. Table.2 shows the influence of the masking ratio jointly with the pre-training length. The ratio of
 169 90% works the best. Because of the information redundancy of the data, the masking ratio of the random
 170 sampling strategy can increase to 90%. Furthermore, a higher masking ratio conducts in fewer tokens
 171 encoded by the encoder, which means the training speed is faster.

172 Table.3 and 4 report the influence of the decoder depth and dimension. In Table 3, the best decoder
 173 dimension was set to 512. In Table.4, the best decoder depth was determined to 4. The accuracy is degraded
 174 by large margins when using an overly decoding architecture. In 2D MAE(He et al., 2021), the proposed
 175 decoding depth from the ablation study is 8. In our 3D task, the optimal decoder depth is 4 which is

Figure 5. Illustration of the segmentation results on the MitoEM-R dataset(Wei et al., 2020). The first row is consecutive EM images, and the second row is its related labels. The third and fifth row represents the segmentation results of the model that is training from scratch, and the pre-training method results are shown in the fourth and sixth row. We use the watershed algorithm as the post-processing step to generate segmentation results from the predicted binary maps and contours, following the configuration from (Lin et al., 2021).



176 lower than the proposed depth on 2D MAE(He et al., 2021). This part is also related to the differences in
177 information redundancy between the 2D and 3D data.

178 4.3 Evaluation results

179 4.3.1 SNEMI3D

180 Table. 5 studies the differences between the pre-training strategy and training from scratch on the
181 SNEMI3D dataset(Lee et al., 2017). Moreover, it shows the difference between predicting targets. In Table.

Table 5. Evaluation results of SNEMI3D(Lee et al., 2017). Time and params are measured in millisecond (ms) and million (m).

Target	Method	Backbone	Voi-S	Voi-M	A-Rand	Time	Params
Affinity	Scratch	Vit-B	0.431	0.334	0.109	138	154
	Pre-train	Vit-B	0.234	0.129	0.071		
	Scratch	Vit-L	0.379	0.318	0.092	195	455
	Pre-train	Vit-L	0.211	0.106	0.063		
BCD	Scratch	Vit-B	0.482	0.341	0.116	142	154
	Pre-train	Vit-B	0.331	0.215	0.084		
	Scratch	Vit-L	0.415	0.301	0.095	200	455
	Pre-train	Vit-L	0.281	0.185	0.079		

Table 6. MitoEM-R(Wei et al., 2020) evaluation results.

Method	Backbone	AP-50	AP-75
Scratch	Vit-B	0.549	0.174
Pre-train	Vit-B	0.895	0.514
Scratch	Vit-L	0.797	0.431
Pre-train	Vit-L	0.923	0.679

182 5, for the prediction of the target, BCD prediction represents the multi-task learning method (Wei et al.,
183 2020), which includes predicting binary maps, contours, and distance (BCD). As shown in Figure.4, we
184 find that predicting the affinity maps presents a more accurate result on small objects, regardless of whether
185 it is pre-trained. Moreover, evaluate metric also shows that predicting affinity maps performs better than
186 predicting BCD. For the training method, as shown in Table 5, the pre-trained model gains comprehensive
187 improvement on both predicting targets compared with the scratched model. The pre-trained Vit-Base
188 and Vit-Large gain 0.071 and 0.063 of the A-Rand value, respectively. Moreover, as shown in Figure.3,
189 regardless of the failure of high-frequency information reconstruction in MAE(He et al., 2021) pre-training,
190 the pre-trained models outperform the scratched models.

191 We also profile the parameters of the models and inferencing times in Table 5. Time and parameters are
192 measured in millisecond (ms) and million (m). We measure the inference time of a single batch with batch
193 size 1. Moreover, we observe that predicting the affinity map has the fastest inference time.

194 4.3.2 MitoEM-R

195 The task of this dataset(Wei et al., 2020) is the instance segmentation of mitochondria. Following the
196 same experiment settings from (Wei et al., 2020), we use the binary maps and instance contours as our
197 targets to fine-tune the models. The configurations are shown in Table.9. The post-processing steps for all
198 the models are following the default configuration from (Wei et al., 2020; Lin et al., 2021). Moreover, we
199 calculate the value of mAP on the validation dataset of MitoEM-R. Table. 6 demonstrates the differences
200 between the pre-training strategy and the training from scratch on the MitoEM-R dataset(Lee et al., 2017).
201 As shown in Table. 6, pre-trained vit-large obtains best results on both AP-75 and AP-50. Moreover, the
202 vit-large with training from scratch performs worse than the pre-trained vit-base. It proves the MAE's(He
203 et al., 2021) capability of representation learning on small objects such as mitochondria. Furthermore, we
204 notice the value of AP-75 from vit-large has a giant improvement compared with vit-base. Higher AP-75
205 means the accurate shape and contour predictions from the model, see Figure. 5, pre-trained vit-large
206 present the best segmentation results compared with other methods.

207 4.3.3 A white matter dataset

208 In this part, we demonstrate the generalization of the model that was pre-trained on SNEMI3D(Lee
209 et al., 2017), a grey matter dataset, to a white matter dataset containing very different structural patterns.
210 The fine-tuning dataset is on the region of the corpus callosum, which contains amounts of myelinated
211 axons and some blurry sections. We manually annotated two different volumes from this dataset for further
212 segmentation experiments. The shape of training and testing volume is $50 \times 3000 \times 0$ and $59 \times 3000 \times 3000$
213 with the resolution of 4 nanometers per pixel, respectively. The fine-tuning process is following the same
214 configuration for SNEMI3D(Lee et al., 2017) fine-tuning8. As shown in table 7, the pre-trained model
215 outperforms the model that trains from scratch in terms of predicted targets. The pre-trained vit-large gains
216 0.197 of A-Rand. The visual results are shown in Figure.6. It proves the representation learning from 3D
217 MAE(He et al., 2021) can promote model performance even when the pre-training dataset and fine-tuning
218 dataset are enormously different. Moreover, in Figure.6, we notice the model that predicting BCD performs
219 better than affinity prediction. Empirically, because of the additional constraining of contour prediction, it
220 allows the model overcomes the impact of blur affectations. In addition, the contours of the myelin sheath
221 are thicker than the cell's membrane, which degrades the challenge of predicting boundaries.

5 DISCUSSION

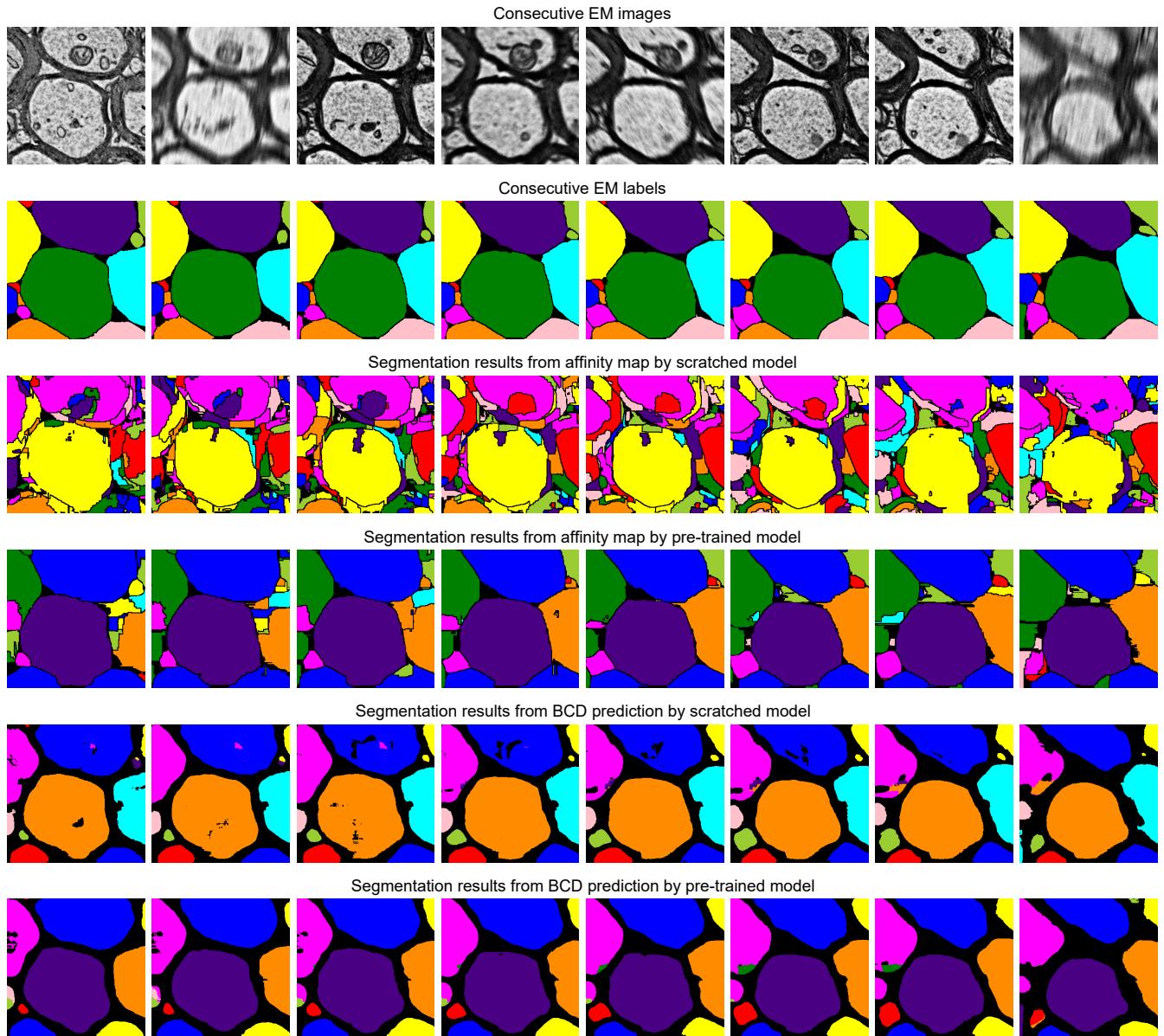
222 This paper proposes the paradigm of implementing MAE(He et al., 2021) to the SEM dataset, which the
223 pre-trained Vit(Vaswani et al., 2017) can be implemented as the backbone of the UNETR(Hatamizadeh
224 et al., 2022) for downstream tasks. We optimized the best configuration of three-dimensional MAE(He
225 et al., 2021) pre-training for the SEM dataset. In the aconfigurationsTable.1 and Table.2, empirically show
226 the efficient representation learning with the 90% random mask sampling strategy. Moreover, Table.3 and
227 Table.4 proves the overlay decoder design can cause the degradation of accuracy.

228 As shown in Table.4, the experiment results on public dataset SNEMI3D(Lee et al., 2017) illustrate
229 the performance, parameters, and inference time of using 3D MAEHe et al. (2021) in downstream tasks.
230 By predicting affinity map, the pre-trained vit-large gains 0.063 of the A-Rand, while the training from
231 scratch method achieves 0.092 of the A-Rand. Moreover, from Table.6, the experiments on MitoEM-R(Wei
232 et al., 2020) demonstrate that pre-training on SNEMI3D(Lee et al., 2017) can also significantly enhance
233 the performance of the tasks on MitoEM-R dataset, while pre-trained vit-large obtains 0.679 of AP-75.
234 Note that the pre-training of the backbone was on SNEMI3D(Lee et al., 2017). We also discovered that
235 the pre-trained backbone has a positive impact on the corpus callosum dataset, which is the region of the
236 grey matter. As shown in Table.7, the pre-trained vit-large gains 0.197 of A-Rand compared with 0.316 of
237 A-Rand by the training from scratch. Such enormous improvement proves the difference across datasets
238 does not constrain the representation learning from 3D MAE(He et al., 2021). In addition, the experiments
239 demonstrate the potential of implementing pre-trained vit(Vaswani et al., 2017) as the backbone to solve
240 the downstream tasks.

6 CONCLUSION

241 We explored the paradigm of implementing MAE(He et al., 2021) to the SEM dataset. We found that
242 representation learning for neural structure heterogeneity is possible with minimal domain knowledge.
243 Similar to the MAE(He et al., 2021) and BERT(Devlin et al., 2018), the masking ratio is strongly related
244 to the information redundancy of the data. Therefore, we found the time cost of the MAE pre-training
245 paradigm for the SEM volume dataset can be tremendously reduced. We reported encouraging results of

Figure 6. Illustration of the segmentation results on the corpus callosum dataset. The first row is consecutive EM images, and the second row is its related labels. The third and fifth row represents the segmentation results of the model that is training from scratch, and the pre-training method results are shown in the fourth and sixth row. We use zwatershed as the post-processing step to generate segmentation results from the predicted affinity map. The post-processing algorithm of BCD (binary maps, contours, distances) predictions are following the configuration from (Lin et al., 2021).



246 using pre-trained vit(Vaswani et al., 2017) as the backbone on two public white matter datasets, and a grey
 247 matter dataset. The pre-training method achieves strong performance and shows the capability of efficient
 248 representation learning across different structure patterns.

CONFLICT OF INTEREST STATEMENT

249 The authors declare that the research was conducted in the absence of any commercial or financial
 250 relationships that could be construed as a potential conflict of interest.

Table 7. Evaluation results on the dataset of corpus callosum, which is the region of the grey matter.

Target	Method	Backbone	Voi-S	Voi-M	A-Rand
Affinity	Scratch	Vit-B	4.122	0.884	0.600
	Pre-train	Vit-B	1.407	0.378	0.214
	Scratch	Vit-L	3.426	0.678	0.316
	Pre-train	Vit-L	0.974	0.205	0.197
BCD	Scratch	Vit-B	1.270	1.116	0.388
	Pre-train	Vit-B	1.098	1.320	0.326
	Scratch	Vit-L	0.898	0.921	0.278
	Pre-train	Vit-L	0.775	0.826	0.241

Table 8. SNEMI3D(Lee et al., 2017) pre-training and fine-tuning configuration

Config (pre-training)	Value
Optimizer	AdamW(Loshchilov and Hutter, 2017)
Optimizer Momentum	$\beta_1, \beta_2=0.9, 0.95$ (Mark et al., 2020)
Weight Decay	0.005
Base Learning Rate	1e-4
Learning Rate Schedule(Loshchilov and Hutter, 2016)	Cosine decay
Warmup Iteration(Goyal et al., 2017)	50,000
Total Iteration	400,000
Batch Size	128
Input Size	6 * 96 * 96

Config (fine-tuning)	Value
Optimizer	AdamW(Loshchilov and Hutter, 2017)
Optimizer Momentum	$\beta_1, \beta_2=0.9, 0.95$ (Mark et al., 2020)
Weight Decay	0.05
Base Learning Rate	1e-4
Learning Rate Schedule(Loshchilov and Hutter, 2016)	Cosine decay
Warmup Iteration(Goyal et al., 2017)	5,000
Dropout(Srivastava et al., 2014)	0.3
Dropout path(Huang et al., 2016)	0.1
Total Iteration	200,000
Augmentation	Default by (Lin et al., 2021)
Batch Size	8
Input Size	6 * 96 * 96

AUTHOR CONTRIBUTIONS

251 AC designed the research and participated in the entire research including data processing, model
 252 construction, result interpretation, and manuscript drafting. LW and RZ designed the research and
 253 participated in the data collection and revisions of the manuscripts.

FUNDING

254 This work was supported by the Hundred Talents Program of the Chinese Academy of Sciences, and the
 255 Leader in Innovation and Entrepreneurship Program of the Province of Jiangsu.

Table 9. MitoEM-R(Wei et al., 2020) fine-tuning configuration

Config	Value
Optimizer	SGD
Weight Decay	0.0001
Base Learning Rate	4e-3
Learning Rate Schedule(Loshchilov and Hutter, 2016)	Cosine decay
Warmup Iteration(Goyal et al., 2017)	10,000
Dropout(Srivastava et al., 2014)	0.3
Dropout path(Huang et al., 2016)	0.1
Total Iteration	300,000
Augmentation	Default by (Lin et al., 2021)
Scales	[1, 0.5, 0.5]
Batch Size	8
Input Size	6 * 96 * 96

ACKNOWLEDGMENTS

256 We would like to thank the Hanhua Lab from the Institute of Automation, Chinese Academy of Sciences
 257 for the corpus callosum SEM data.

DATA AVAILABILITY STATEMENT

258 The codes and datasets of this study can be found in the [https://github.com/aoc777/
 259 connectomics/tree/master/projects/EmMAE](https://github.com/aoc777/connectomics/tree/master/projects/EmMAE).

REFERENCES

- 260 Bao, H., Dong, L., and Wei, F. (2021). Beit: Bert pre-training of image transformers. *arXiv preprint
 261 arXiv:2106.08254*
- 262 Bogovic, J. A., Huang, G. B., and Jain, V. (2013). Learned versus hand-designed feature representations
 263 for 3d agglomeration. In *CVPR*
- 264 Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). Bert: Pre-training of deep bidirectional
 265 transformers for language understanding. *arXiv preprint arXiv:1810.04805*
- 266 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., et al. (2021). An
 267 image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*
- 268 Eberle, A., Mikula, S., Schalek, R., Lichtman, J., Tate, M., and Zeidler, D. (2018). High-resolution,
 269 high-throughput imaging with a multibeam scanning electron microscope. *Journal of Microscopy*
- 270 Goyal, P., Dollár, P., Girshick, R., Noordhuis, P., Wesolowski, L., Kyrola, A., et al. (2017). Accurate, large
 271 minibatch sgd: Training imagenet in 1 hour. In *CVPR*
- 272 Hatamizadeh, A., Tang, Y., Nath, V., Yang, D., Myronenko, A., Landman, B., et al. (2022). Unetr:
 273 Transformers for 3d medical image segmentation. In *Proceedings of the IEEE/CVF Winter Conference
 274 on Applications of Computer Vision*. 574–584
- 275 He, K., Chen, X., Xie, S., Li, Y., Dollár, P., and Girshick, R. (2021). Masked autoencoders are scalable
 276 vision learners. In *arXiv*. doi:10.48550/ARXIV.2111.06377
- 277 Huang, G., Sun, Y., Liu, Z., Sedra, D., and Weinberger, K. Q. (2016). Deep networks with stochastic depth.
 278 In *European conference on computer vision* (Springer), 646–661
- 279 Hung, W., Tsai, Y., Liou, Y., Lin, Y., and Yang, M. (2018). Adversarial learning for semi-supervised
 280 semantic segmentation. *IEEE TPAMI*

- 281 Januszewski, M. and Jain, V. (2019). Segmentation-enhanced cyclegan. *bioRxiv*, 548081
- 282 Kasthuri, N., Hayworth, K., Berger, D., Schalek, R., Conchello, J., Knowles-Barley, S., et al. (2015).
283 Saturated reconstruction of a volume of neocortex. *Cell*
- 284 Kingma, D. and Ba, J. (2014). Adam: A method for stochastic optimization. In *ICLR*
- 285 Lauenburg, L., Lin, Z., Zhang, R., Santos, M. d., Huang, S., Arganda-Carreras, I., et al. (2022). Instance
286 segmentation of unlabeled modalities via cyclic segmentation gan. *arXiv preprint arXiv:2204.03082*
- 287 Lee, K., Zung, J., Li, P., Jain, V., and Seung, H. S. (2017). Superhuman accuracy on the snemi3d
288 connectomics challenge. *arXiv preprint arXiv:1706.00120*
- 289 Lin, Z., Wei, D., Jang, W.-D., Zhou, S., Chen, X., Wang, X., et al. (2020). Two stream active query
290 suggestion for active learning in connectomics. In *European Conference on Computer Vision* (Springer),
291 103–120
- 292 Lin, Z., Wei, D., Lichtman, J., and Pfister, H. (2021). Pytorch connectomics: a scalable and flexible
293 segmentation framework for em connectomics. *arXiv preprint arXiv:2112.05754*
- 294 Loshchilov, I. and Hutter, F. (2016). Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint*
295 *arXiv:1608.03983*
- 296 Loshchilov, I. and Hutter, F. (2017). Decoupled weight decay regularization. *arXiv preprint*
297 *arXiv:1711.05101*
- 298 Mark, C., Alec, R., Rewon, C., Jeffrey, W., Heewoo, J., David, L., et al. (2020). Generative pretraining
299 from pixels. In *ICML*
- 300 Mittal, S., Tatarchenko, M., and Brox, T. (2019). Semi-supervised semantic segmentation with high- and
301 low-level consistency. *IEEE TPAMI*
- 302 Nunez-Iglesias, J., Kennedy, R., Parag, T., Shi, J., and B. Chklovskii, D. (2013). Machine learning of
303 hierarchical clustering to segment 2d and 3d images. *PLOS ONE*
- 304 Pathak, D., Krahenbuhl, P., Donahue, J., Darrell, T., and Efros, A. A. (2016). Context encoders: Feature
305 learning by inpainting. In *Proceedings of the IEEE conference on computer vision and pattern recognition*.
306 2536–2544
- 307 Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., et al. (2018). Improving language understanding
308 by generative pre-training
- 309 Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. (2019). Language models are
310 unsupervised multitask learners. *OpenAI blog* 1, 9
- 311 Raghu, M., Zhang, C., Kleinberg, J., and Bengio, S. (2019). Transfusion: Understanding transfer learning
312 for medical imaging. *Advances in neural information processing systems* 32
- 313 Rand, W. M. (1971). Objective criteria for the evaluation of clustering methods. *Journal of the American
314 Statistical Association*
- 315 Ronneberger, O., Fischer, P., and Brox, T. (2015). U-net: Convolutional networks for biomedical
316 image segmentation. In *International Conference on Medical image computing and computer-assisted
317 intervention* (Springer), 234–241
- 318 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout: a
319 simple way to prevent neural networks from overfitting. *The journal of machine learning research* 15,
320 1929–1958
- 321 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., et al. (2017). Attention is all
322 you need. In *NeurIPS*
- 323 Vincent, P., Larochelle, H., Bengio, Y., and Manzagol, P.-A. (2008). Extracting and composing robust
324 features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine
325 learning*. 1096–1103

- 326 Vincent, P., Larochelle, H., Lajoie, I., Bengio, Y., Manzagol, P.-A., and Bottou, L. (2010). Stacked
327 denoising autoencoders: Learning useful representations in a deep network with a local denoising
328 criterion. *Journal of machine learning research* 11
- 329 Wei, C., Fan, H., Xie, S., Wu, C.-Y., Yuille, A., and Feichtenhofer, C. (2022). Masked feature prediction
330 for self-supervised visual pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision
331 and Pattern Recognition (CVPR)*. 14668–14678
- 332 Wei, D., Lin, Z., Franco-Barranco, D., Wendt, N., Liu, X., Yin, W., et al. (2020). Mitoem dataset:
333 large-scale 3d mitochondria instance segmentation from em images. In *International Conference on
334 Medical Image Computing and Computer-Assisted Intervention* (Springer), 66–76
- 335 Xie, Z., Zhang, Z., Cao, Y., Lin, Y., Bao, J., Yao, Z., et al. (2022). Simmim: A simple framework for
336 masked image modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
337 Recognition*. 9653–9663
- 338 You, C., Li, G., Zhang, Y., Zhang, X., Shan, H., Li, M., et al. (2020). Ct super-resolution gan constrained
339 by the identical, residual, and cycle learning ensemble (gan-circle). *IEEE Transactions on Medical
340 Imaging* 39, 188–203. doi:10.1109/TMI.2019.2922960