

Appendix

This appendix contains tables for evaluation on the edit images (i.e. each of the 6 training images). We include these evaluations to check if the edited models perform better than the original model on the image used for the corresponding edit. The evaluations on the unseen images are included in the main paper. In general, we see that the performance of each mapping on the edit images loosely corresponds to its performance on the unseen images.

Image	Precision		Recall		IoU	
	Original	Edited	Original	Edited	Original	Edited
A	0.82 ± 0.19	0.88 ± 0.19	0.78 ± 0.24	0.63 ± 0.28	0.66 ± 0.25	0.57 ± 0.29
B	0.53 ± 0.28	0.54 ± 0.27	0.78 ± 0.24	0.67 ± 0.24	0.66 ± 0.27	0.61 ± 0.26
C	0.83 ± 0.20	0.95 ± 0.09	0.77 ± 0.24	0.70 ± 0.23	0.64 ± 0.30	0.67 ± 0.24
D	0.79 ± 0.23	0.91 ± 0.13	0.78 ± 0.28	0.74 ± 0.23	0.64 ± 0.30	0.66 ± 0.26
E	0.85 ± 0.15	0.92 ± 0.10	0.70 ± 0.26	0.63 ± 0.28	0.55 ± 0.32	0.58 ± 0.29
F	0.83 ± 0.16	0.90 ± 0.10	0.82 ± 0.22	0.72 ± 0.21	0.67 ± 0.26	0.64 ± 0.23

Table 9: **Mapping non-archaeocyathids to red mud: Performance on edit image.** This table contains the mean instance-level metrics \pm one standard deviation run on each image used for mapping non-archaeocyathids to red mud. For example, the first row contains the mean instance-level precision, recall, and IoU across the identified archaeocyathids in image *A* alone. All the precisions increase, and some of the IoU scores increase as well. The recall scores decrease.

Image	Precision		Recall		IoU	
	Original	Edited	Original	Edited	Original	Edited
A	0.82 ± 0.19	0.82 ± 0.18	0.78 ± 0.24	0.81 ± 0.22	0.66 ± 0.25	0.67 ± 0.26
B	0.53 ± 0.28	0.53 ± 0.25	0.78 ± 0.24	0.77 ± 0.24	0.66 ± 0.27	0.62 ± 0.29
C	0.83 ± 0.20	0.84 ± 0.15	0.77 ± 0.24	0.75 ± 0.30	0.64 ± 0.30	0.60 ± 0.34
D	0.79 ± 0.23	0.85 ± 0.16	0.78 ± 0.28	0.82 ± 0.26	0.64 ± 0.30	0.66 ± 0.28
E	0.85 ± 0.15	0.80 ± 0.16	0.70 ± 0.26	0.73 ± 0.30	0.55 ± 0.32	0.48 ± 0.35
F	0.83 ± 0.16	0.82 ± 0.14	0.82 ± 0.22	0.81 ± 0.26	0.67 ± 0.26	0.65 ± 0.28

Table 10: **Mapping all archaeocyathids to pitted texture: Performance on edit image.** This table contains the mean instance-level metrics \pm one standard deviation run on each image used for mapping archaeocyathids to the pitted texture. For example, the first row contains the mean instance-level precision, recall, and IoU across the identified archaeocyathids in image *A* alone. There does not appear to be a clear trend in any of the metrics. For example, image *D* produces an improvement across all metrics, while image *F* does not.

Image	Precision		Recall		IoU	
	Original	Edited	Original	Edited	Original	Edited
A	0.82 ± 0.19	0.80 ± 0.21	0.78 ± 0.24	0.60 ± 0.32	0.66 ± 0.25	0.44 ± 0.35
B	0.53 ± 0.28	0.51 ± 0.32	0.78 ± 0.24	0.64 ± 0.28	0.66 ± 0.27	0.51 ± 0.32
C	0.83 ± 0.20	0.83 ± 0.16	0.77 ± 0.24	0.64 ± 0.32	0.64 ± 0.30	0.44 ± 0.36
D	0.79 ± 0.23	0.80 ± 0.21	0.78 ± 0.28	0.77 ± 0.24	0.64 ± 0.30	0.53 ± 0.36
E	0.85 ± 0.15	0.85 ± 0.14	0.70 ± 0.26	0.60 ± 0.32	0.55 ± 0.32	0.46 ± 0.34
F	0.83 ± 0.16	0.87 ± 0.11	0.82 ± 0.22	0.74 ± 0.24	0.67 ± 0.26	0.64 ± 0.26

Table 11: **Mapping simultaneously: Performance on edit image.** This table contains the mean instance-level metrics \pm one standard deviation run on each image used for simultaneously mapping non-archaeocyathids to red mud and archaeocyathids to the pitted texture. For example, the first row contains the mean instance-level precision, recall, and IoU across the identified archaeocyathids in image *A* alone. Only image *F* produces an increase in precision. None of the edited models have higher recall or IoU scores than the original model.