
Robust semi-supervised segmentation with timestep ensembling diffusion models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Medical image segmentation is challenging due to limited data and annotations.
2 Denoising diffusion probabilistic models (DDPM) show promise in modelling
3 natural image distributions and are successfully applied in medical imaging. Our
4 research focuses on semi-supervised image segmentation using diffusion models'
5 latent representations and addressing domain generalisation. We found that optimal
6 performance depends on choice of diffusion steps and ensembling. Our model out-
7 performed in domain-shifted settings while remaining competitive within domain,
8 highlighting DDPMs' potential for medical image segmentation.¹

9 1 Introduction

10 Denoising diffusion probabilistic models (DDPM) [Sohl-Dickstein et al., 2015, Ho et al., 2020] have
11 recently emerged as a promising approach for modelling the distribution of natural images, outper-
12 forming alternative methods in terms of sample realism and diversity. More recently, DDPM have
13 also been successfully applied to various medical imaging tasks, such as image reconstruction [Xie
14 and Li, 2022], diagnostics [Aviles-Rivero et al., 2022] and segmentation [Wolleb et al., 2022].

15 Image segmentation is crucial in medical settings, where accurate and efficient methods are required
16 to support diagnosis, treatment planning, and disease monitoring. However, limited dataset size
17 and insufficient annotations make it challenging to train accurate models. High variability due to
18 differences in acquisition parameters, scanner types, and patient demographics, known as domain
19 shift, also presents difficulties in generalising segmentation models to new datasets, leading to
20 potential underperformance in clinical settings.

21 Recent research in diffusion models has shown promising results [Baranchuk et al., 2021, Deja et al.,
22 2023] for semi-supervised learning: the bottleneck network tasked to learn the backward process of
23 removing noise from an image also learns an expressive feature representation that can benefit other
24 downstream analysis tasks. However, more research is needed to understand the implications of these
25 models' design choices for generalisation.

26 Our work focuses on optimally leveraging diffusion steps for improving generalisation in semi-
27 supervised image segmentation under domain shift. Based on the analysis of datasets with diverse
28 imaging modalities and domain shifts, our findings demonstrate significant improvements over
29 existing baselines using five different datasets. Our key findings can be summarised as follows:

- 30 • Small diffusion steps are crucial for model generalisation;
- 31 • Concatenating latent representations over steps to predict segmentation maps can hurt
32 generalisation;

¹Demo: <https://huggingface.co/spaces/anonymous2023-21/TEDM-demo>

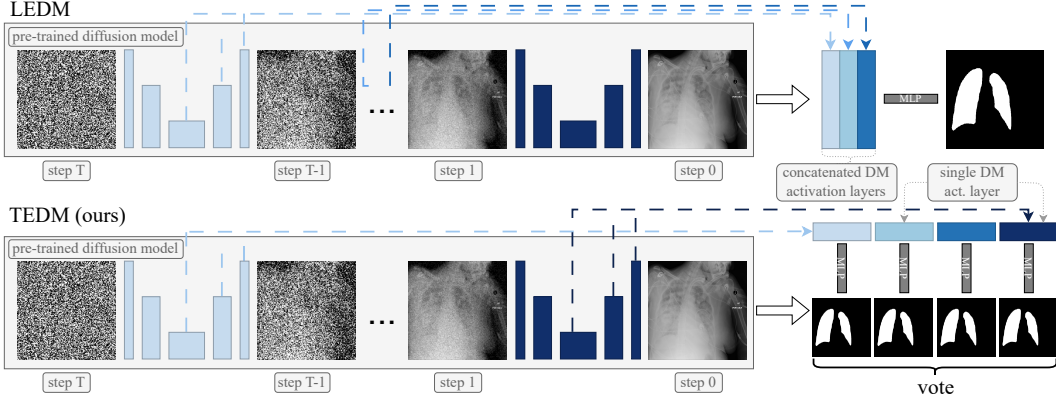


Figure 1: Models diagram. LEDM, the SOTA in semi-supervised segmentation with diffusion models, selects a subset of timesteps and concatenates latent representations extracted from a pretrained diffusion model as features fed to an MLP. Our method (i) selects smaller and more informative timesteps, (ii) predicts through a voting mechanism over our steps selection and (iii) shares the MLP weights across timesteps, resulting in improved segmentation performance.

- Instead, generalisation can be significantly improved by (i) optimising which timesteps to use at test time, (ii) ensembling predictions from individual timesteps using a shared predictor and (iii) using these individual predictions for regularisation during training.

2 Background and related work

DDPM are generative models that use a UNet to iteratively denoise a noise signal over T timesteps and generate samples from a distribution. See Appendix A for more information.

Baranchuk et al. [2021] apply diffusion models to semi-supervised segmentation by using a DDPM pretrained on unlabelled images and extracting latent representation from its UNet’s intermediate layers. Their Label Efficient Diffusion Model (LEDM) selects a set of steps $t \in S \subset \{0, \dots, T\}$ and generates latent representations $\mathbf{z}_t \in \mathbb{R}^{c \times h \times w}$. These are upsampled to the input size and concatenated into a feature map $\mathbf{Z} \in \mathbb{R}^{(|S| \times c) \times H \times W}$. Finally, an ensemble of lightweight multilayer perceptions (MLPs) $C_\phi^n : \mathbf{Z}^{i,j} \rightarrow y^{i,j}$; $n \in \{1, \dots, 10\}$ performs pointwise prediction, trained with a cross-entropy loss. The authors choose diffusion steps $S = \{50, 150, 250\}$ to form the input.

Similarly, Deja et al. [2023] use the latent representations of a pretrained diffusion model for classification. Their proposed method uses all intermediate timesteps to regularise the training process, and only uses the last diffusion step $t = 1$ at test time.

3 Timestep ensembling diffusion models

Preliminary results discussed in Appendix B illustrate that LEDM does not perform well to out-of-distribution (OOD) settings. We hypothesise that using more model regularization and reducing the number of parameters can improve its generalization. Currently, LEDM’s approach of concatenating features from numerous timesteps to feed into the pixel-wise MLP predictor results in an excessively high-dimensional input and a complex predictor. Instead, we propose using a shared MLP trained to generate a prediction map from each latent representation of the steps considered.

We define our loss function as follows:

$$\phi = \operatorname{argmin} \mathbb{E}_{\mathcal{D}} \mathbb{E}_{i,j} \mathbb{E}_{s \in S} \operatorname{CE} (C_\phi(\mathbf{z}_s^{i,j}), y^{i,j}) \quad (1)$$

At test time, we use a voting mechanism to ensemble the various prediction maps to obtain a final segmentation map. We call this technique “timestep ensembling” and show that it yields superior performance. Moreover, we leverage the insights from the preliminary results and combine predictions from the diffusion steps $\{1, 10, 25, 50, 200, 400, 600, 800\}$. This approach allows us to benefit from

Table 1: Models performance w.r.t. ground truth segmentations. Reported as mean \pm standard deviation over the dataset. Global CL, Global & Local CL and LEDM are a reproduction of Chen et al. [2020], Chaitanya et al. [2020] and Baranchuk et al. [2021] respectively.

Training size	1 (1%)	3 (2%)	6 (3%)	12 (6%)	197 (100%)
JSRT (in-domain for classifier)					
Sup. Baseline	84.4 \pm 5.4	91.7 \pm 3.7	93.3 \pm 2.9	95.3 \pm 2.3	97.3 \pm 1.2
Global CL	88.8 \pm 5.9	92.7 \pm 1.8	93.6 \pm 1.6	95.3 \pm 1.1	97.1 \pm 1.4
Global & Local CL	89.8 \pm 5.2	93.1 \pm 1.7	92.9 \pm 1.9	94.8 \pm 1.49	97.2 \pm 1.2
LEDM	90.8 \pm 3.5	94.1 \pm 1.6	95.5 \pm 1.4	96.4 \pm 1.4	97.0 \pm 1.3
LEDMe	93.7 \pm 2.6	95.5 \pm 1.5	96.7 \pm 1.5	97.0 \pm 1.1	97.6 \pm 1.2
TEDM (ours)	93.1 \pm 3.4	94.8 \pm 1.4	95.8 \pm 1.2	96.6 \pm 1.1	97.3 \pm 1.2
NIH (in-domain for DDPM, OOD for classifier)					
Sup. Baseline	68.5 \pm 12.8	71.2 \pm 15.1	71.4 \pm 15.9	77.8 \pm 14.0	81.5 \pm 12.7
Global CL	70.7 \pm 14.6	80.3 \pm 12.2	77.1 \pm 16.4	84.6 \pm 10.8	86.9 \pm 10.8
Global & Local CL	71.1 \pm 16.2	79.6 \pm 12.7	81.1 \pm 14.0	82.2 \pm 13.6	87.4 \pm 10.8
LEDM	63.3 \pm 12.2	78.0 \pm 10.1	81.2 \pm 9.3	85.9 \pm 7.4	88.9 \pm 5.9
LEDMe	70.3 \pm 11.4	78.3 \pm 9.8	83.0 \pm 8.6	84.4 \pm 8.1	90.1 \pm 5.3
TEDM (ours)	80.3 \pm 9.0	86.4 \pm 6.2	89.2 \pm 5.5	91.3 \pm 4.1	92.9 \pm 3.2
Montgomery (OOD for DDPM and classifier)					
Sup. Baseline	77.1 \pm 12.0	83.0 \pm 12.2	80.9 \pm 14.7	83.8 \pm 14.9	94.1 \pm 6.6
Global CL	76.1 \pm 15.0	87.6 \pm 9.7	88.8 \pm 11.4	90.4 \pm 10.4	92.9 \pm 10.8
Global & Local CL	77.4 \pm 17.4	88.7 \pm 9.14	89.9 \pm 8.2	90.1 \pm 10.9	92.5 \pm 11.2
LEDM	79.3 \pm 8.1	85.9 \pm 7.4	89.4 \pm 6.7	92.3 \pm 7.2	94.4 \pm 7.2
LEDMe	80.7 \pm 6.6	86.3 \pm 6.5	89.5 \pm 5.9	91.2 \pm 5.6	95.3 \pm 4.0
TEDM (ours)	90.5 \pm 5.3	91.4 \pm 6.1	93.3 \pm 6.0	94.6 \pm 6.0	95.1 \pm 6.9

61 the small steps information content and larger step regularisation effect, unlike LEDM, which only
 62 used timesteps {50, 125, 250}. To better understand the distinctions between our model and LEDM,
 63 please refer to Figure 1.

64 4 Experiments

65 We evaluate our work on the task of chest X-ray lung segmentation, training the DDPM on ChestX-
 66 ray8 [Wang et al., 2017], the MLP on JSRT [Van Ginneken et al., 2006] and testing on JSRT,
 67 NIH [Tang et al., 2019] and the Montgomery [Jaeger et al., 2014] datasets, where NIH is a labelled
 68 subset of ChestX-ray8. For details on the experimental setup, please refer to Appendix B.

69 To test our semi-supervised method, we experiment with various percentages of the JSRT training set
 70 (100%, 12%, 6%, 3%, 2%, and 1%). We compare our timestep ensembling diffusion model (TEDM)
 71 to a fully supervised baseline (described in Appendix B), LEDM, and two other semi-supervised
 72 methods that use contrastive learning (CL): the ‘Global CL’ [Chen et al., 2020] and the ‘Local and
 73 Global CL’ [Chaitanya et al., 2020]. All methods were trained with the same backbone architecture.

74 We perform ablations to analyse the impact of each component in our TEDM model. We compare the
 75 LEDM model with an instance trained using our diffusion steps, henceforth LEDMe. Additionally,
 76 we test the test-time voting mechanism using steps 1, 10, and 25 individually instead of the ensemble.

77 Finally, we test the TEDM method on two additional datasets: the UK Biobank dataset and the BraTS
 78 dataset [Menze et al., 2014, Bakas et al., 2017, 2018]. In the UK Biobank dataset, we segment brain
 79 structures in 2D slices of brain MRI T1 images, while in the BraTS dataset, we segment tumours from
 80 brain MRI of patients. The former dataset is challenging due to the low intensity variation between
 81 structures and background, while the latter is even more difficult as it entails segmenting items of
 82 varied shapes and locations. Further details on the experimental process for these two datasets are
 83 available in Appendix C.

84 5 Results

85 The performance results on chest X-rays and brain MRI are shown quantitatively in Tables 1 and 3,
86 and qualitatively in Figure 4. The ablation results are shown in Table 2. Further results can be found
87 in Appendix D. The best-performing models² are highlighted in bold in all tables.

88 **Using small step sizes improves performance both in- and out-of-domain.** Across all experi-
89 ments, models with small diffusion steps perform the best: LEDMe outperforms LEDM in all but
90 two experiments in Table 1, and for the UK Biobank and BraTS datasets in Table 3 for training sizes
91 larger than 3 and 1, respectively.

92 **Concatenating latent representations hurts generalisability in the low data regime.** TEDM
93 outperforms LEDM and LEDMe, except for $n=197$, in NIH and Montgomery datasets. We deduce
94 that concatenation in LEDM hurts generalisation. In addition, TEDM performs statistically similarly
95 to LEDM for JSRT, indicating that its generalisability comes with no in-domain performance cost.

96 **Test-time ensembling over timesteps improves generalisation over single-step predictions.** Ta-
97 ble 2 shows that the voting mechanism used in TEDM is more effective than using any individual
98 step, as different steps produce latent representations focused on different aspects of the image.

99 **TEDM performs robustly for increasingly challenging segmentation tasks.** Table 3 shows that
100 TEDM is statistically superior or equal to its competitors for all cases with less than 12 datapoints,
101 demonstrates its competitiveness in challenging in-domain scenarios with low labelled data.

102 **Fully supervised baselines are competitive for in-domain harder segmentation tasks.** Our
103 method TEDM showcases excellent performance on very small dataset sizes (1, 2, 3 and 6 in Table 3).
104 However, for larger datasets (6 patients or more), a well-designed baseline model is more effective
105 than any of the semi-supervised models. This result suggests that although semi-supervised methods
106 with self-supervised pretraining may have their limitations in providing task-specific performance for
107 larger datasets, they present great potential for improving results on small datasets.

108 6 Conclusions

109 This study investigated the impact of different diffusion steps on the performance and generalisation
110 of semi-supervised segmentation models. Our comprehensive experiments across multiple datasets
111 revealed that small diffusion steps are crucial for domain generalisation, requiring only a few
112 training samples to become powerful pixel-wise predictors. Furthermore, we found that ensembling
113 segmentation maps over timesteps significantly improves model generalisation in the low data regime
114 while offering competitive performance in-domain. Conversely, concatenating latent representations
115 can hurt the generalisation of the pixel-wise classifier. These findings were demonstrated by the
116 superior performance of our proposed Timestep Ensembling Diffusion Model on chest X-ray lung
117 segmentation and more challenging tasks such as brain structure and tumour segmentation. Our
118 results indicate that latent representations across different steps share semantics and act as a model
119 regulariser, leading to better generalisation than competing methods. This analysis underscores the
120 importance of thoroughly investigating the design decisions for auxiliary tasks in diffusion models,
121 such as timestep selection and ensembling. These decisions can have a significant impact on the
122 model’s performance.

123 Our findings provide important new insights and may inform the development of new approaches
124 leveraging powerful diffusion models for medical imaging tasks. In future work, the performance of
125 TEDM and similar approaches should be compared to the emerging foundation model techniques,
126 where the pre-training is executed at a larger scale than semi-supervised methods. Here, the ability of
127 diffusion models to efficiently capture the data distribution from extensive, unlabelled data holds a
128 promise to overcome the persistent data scarcity problem in medical image segmentation.

²That is best-performing and statistically equivalent models

References

- 129
130 Fidel Alfaro-Almagro, Mark Jenkinson, Neal K Bangerter, Jesper LR Andersson, Ludovica Grif-
131 fanti, Gwenaëlle Douaud, Stamatios N Sotiropoulos, Saad Jbabdi, Moises Hernandez-Fernandez,
132 Emmanuel Vallee, et al. Image processing and quality control for the first 10,000 brain imaging
133 datasets from uk biobank. *Neuroimage*, 166:400–424, 2018.
- 134 Angelica I Aviles-Rivero, Christina Runkel, Nicolas Papadakis, Zoe Kourtzi, and Carola-Bibiane
135 Schönlieb. Multi-modal hypergraph diffusion network with dual prior for alzheimer classification.
136 In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2022: 25th Inter-*
137 *national Conference, Singapore, September 18–22, 2022, Proceedings, Part III*, pages 717–727.
138 Springer, 2022.
- 139 Spyridon Bakas, Hamed Akbari, Aristeidis Sotiras, Michel Bilello, Martin Rozycki, Justin S Kirby,
140 John B Freymann, Keyvan Farahani, and Christos Davatzikos. Advancing the cancer genome atlas
141 glioma mri collections with expert segmentation labels and radiomic features. *Scientific data*, 4(1):
142 1–13, 2017.
- 143 Spyridon Bakas, Mauricio Reyes, Andras Jakab, Stefan Bauer, Markus Rempfler, Alessandro Crimi,
144 Russell Takeshi Shinohara, Christoph Berger, Sung Min Ha, Martin Rozycki, et al. Identifying
145 the best machine learning algorithms for brain tumor segmentation, progression assessment, and
146 overall survival prediction in the brats challenge. *arXiv preprint arXiv:1811.02629*, 2018.
- 147 Dmitry Baranchuk, Andrey Voynov, Ivan Rubachev, Valentin Khrukov, and Artem Babenko. Label-
148 efficient semantic segmentation with diffusion models. In *International Conference on Learning*
149 *Representations*, 2021.
- 150 Krishna Chaitanya, Ertunc Erdil, Neerav Karani, and Ender Konukoglu. Contrastive learning of
151 global and local features for medical image segmentation with limited annotations. *Advances in*
152 *Neural Information Processing Systems*, 33:12546–12558, 2020.
- 153 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
154 contrastive learning of visual representations. In *International conference on machine learning*,
155 pages 1597–1607. PMLR, 2020.
- 156 Kamil Deja, Tomasz Trzcinski, and Jakub M Tomczak. Learning data representations with joint
157 diffusion models. *arXiv preprint arXiv:2301.13622*, 2023.
- 158 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
159 *Neural Information Processing Systems*, 33:6840–6851, 2020.
- 160 Stefan Jaeger, Sema Candemir, Sameer Antani, Yi-Xiáng J Wáng, Pu-Xuan Lu, and George Thoma.
161 Two public chest X-ray datasets for computer-aided screening of pulmonary diseases. *Quantitative*
162 *imaging in medicine and surgery*, 4(6):475, 2014.
- 163 Bjoern H Menze, Andras Jakab, Stefan Bauer, Jayashree Kalpathy-Cramer, Keyvan Farahani, Justin
164 Kirby, Yuliya Burren, Nicole Porz, Johannes Slotboom, Roland Wiest, et al. The multimodal brain
165 tumor image segmentation benchmark (brats). *IEEE transactions on medical imaging*, 34(10):
166 1993–2024, 2014.
- 167 Brian Patenaude, Stephen M Smith, David N Kennedy, and Mark Jenkinson. A bayesian model of
168 shape and appearance for subcortical brain segmentation. *Neuroimage*, 56(3):907–922, 2011.
- 169 Margherita Rosnati, Fabio De Sousa Ribeiro, Miguel Monteiro, Daniel Coelho de Castro, and Ben
170 Glocker. Analysing the effectiveness of a generative model for semi-supervised medical image
171 segmentation. In *Machine Learning for Health*, pages 290–310. PMLR, 2022.
- 172 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised
173 learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*,
174 pages 2256–2265. PMLR, 2015.
- 175 You-Bao Tang, Yu-Xing Tang, Jing Xiao, and Ronald M Summers. Xlsor: A robust and accurate lung
176 segmentor on chest X-rays using criss-cross attention and customized radiorealistic abnormalities
177 generation. In *International Conference on Medical Imaging with Deep Learning*, pages 457–467.
178 PMLR, 2019.

- 179 Bram Van Ginneken, Mikkel B Stegmann, and Marco Loog. Segmentation of anatomical structures
180 in chest radiographs using supervised methods: a comparative study on a public database. *Medical*
181 *image analysis*, 10(1):19–40, 2006.
- 182 Xiaosong Wang, Yifan Peng, Le Lu, Zhiyong Lu, Mohammadhadi Bagheri, and Ronald M Summers.
183 ChestX-ray8: Hospital-scale chest X-ray database and benchmarks on weakly-supervised classi-
184 fication and localization of common thorax diseases. In *Proceedings of the IEEE conference on*
185 *computer vision and pattern recognition*, pages 2097–2106, 2017.
- 186 Julia Wolleb, Florentin Bieder, Robin Sandkühler, and Philippe C Cattin. Diffusion models for
187 medical anomaly detection. In *Medical Image Computing and Computer Assisted Intervention–*
188 *MICCAI 2022: 25th International Conference, Singapore, September 18–22, 2022, Proceedings,*
189 *Part VIII*, pages 35–45. Springer, 2022.
- 190 Yutong Xie and Quanzheng Li. Measurement-conditioned denoising diffusion probabilistic model
191 for under-sampled medical image reconstruction. In *Medical Image Computing and Computer*
192 *Assisted Intervention–MICCAI 2022: 25th International Conference, Singapore, September 18–22,*
193 *2022, Proceedings, Part VI*, pages 655–664. Springer, 2022.

194 A Diffusion models

195 Diffusion models have garnered significant interest in the machine learning community due to their
196 remarkable ability to model complex data distributions efficiently. Diffusion models utilise a series
197 of simple and learnable transformations to diffuse noise iteratively and generate samples from the
198 target distribution. Formally, a DDPM works as follows. Given a data distribution $p(\mathbf{x}_0)$ and forward
199 process:

$$p(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}), \quad (2)$$

200 where $\beta_t \in (0, 1)$ is the variance schedule and $t \in [0, T]$ is the Markov chain time step, a DDPM
201 aims to learn $\mu_\theta(\mathbf{x}_t, t)$ and $\Sigma_\theta(\mathbf{x}_t, t)$ which define the backward process:

$$p(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t)). \quad (3)$$

202 In order to do so, Ho et al. [2020] fix the variance $\Sigma_\theta(\mathbf{x}_t, t)$, reparametrise $\mu_\theta(\mathbf{x}_t, t)$ as a function of
203 the noise $\epsilon_\theta(\mathbf{x}_t, t)$

$$\mu_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right), \quad \text{where } \alpha_t = 1 - \beta_t, \quad \bar{\alpha}_t = \prod_{i=1}^t \alpha_i \quad (4)$$

204 and design a UNet-based neural network architecture

$$G_\theta : (\mathbf{x}_t, t) \rightarrow \epsilon_\theta(\mathbf{x}_t, t)$$

205 for learning to identify the noise. The UNet is trained through cross-entropy between the injected and
206 predicted noise.

207 B On the importance of the diffusion steps for domain generalisation

208 Previous findings suggest that latent representations in larger steps contain coarse information, which
209 becomes more granular as the diffusion steps approach the target data distribution [Baranchuk et al.,
210 2021, Deja et al., 2023]. Here, we are interested in understanding how the wealth of information in
211 each time step $s \in S$ contributes to model generalisation when the training dataset size varies.

212 We train a Ridge logistic regression-based pixel-wise classifier over latent representations extracted
213 from specific timesteps $t = \{1, 10, 25, 50, 200, 400, 600, 800\}$ to isolate the predictive power of each
214 timestep. We compare these timestep-wise predictions to LEDM and a fully supervised baseline
215 using the same UNet backbone as the DDPM backbone.

216 We evaluate our work on the task of chest X-ray lung segmentation. Chest X-rays are among the most
217 frequent radiological examinations in clinical practice, and automatically extracted features from
218 anatomical regions such as the lungs can aid clinical decision-making. Moreover, the availability of

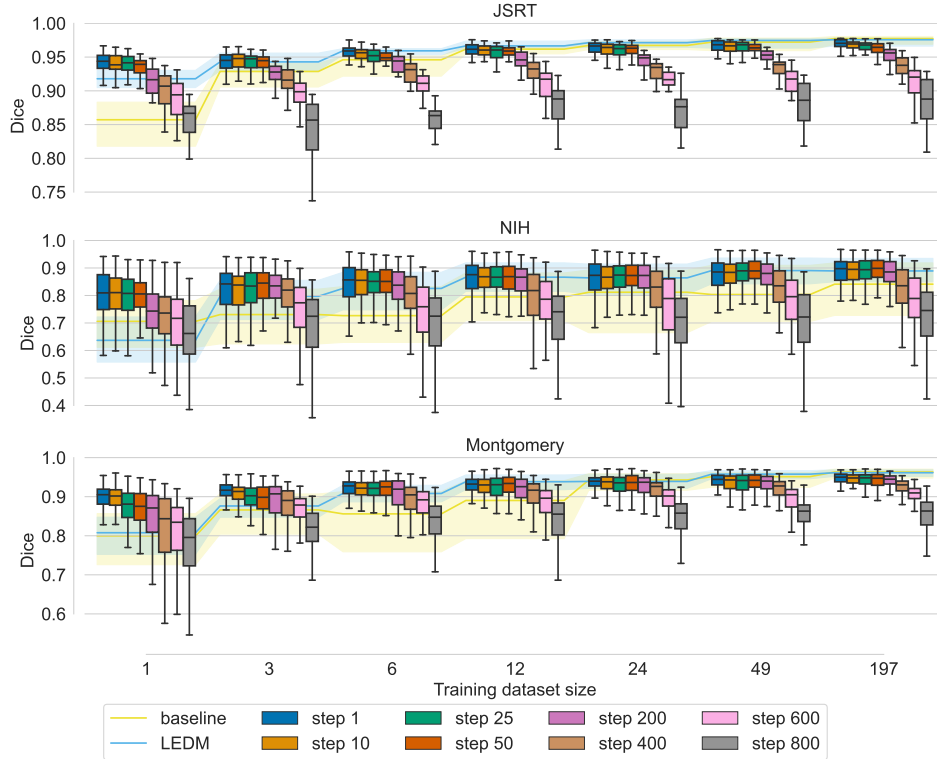


Figure 2: Performance of a logistic regression segmentation model trained on latent features from individual diffusion steps.

219 several public datasets of chest X-ray images allows us to investigate the methods’ generalisation
 220 ability in the presence of changes in dataset characteristics.

221 Following previous work in semi-supervised medical image segmentation [Rosnati et al., 2022],
 222 we use the ChestX-ray8 [Wang et al., 2017] (n=108k) as the unlabelled dataset to train the DDPM
 223 backbone over $T = 1000$ steps and a subset of the JSRT [Van Ginneken et al., 2006] (n=247) labelled
 224 dataset for training (n=197) and validating (n=25) our method. The dataset splits, architecture, and
 225 code are available in our code repository.

226 We reserve the remaining JSRT samples (n=25) along with the NIH [Tang et al., 2019] (n=95), and
 227 Montgomery [Jaeger et al., 2014] (n=138) labelled datasets for final testing. Notably, the NIH dataset
 228 is an annotated subset of the ChestX-ray8 dataset. This setup allows us to test the models on data that
 229 is (i) in-domain for the classifier (JSRT), (ii) out-of-domain for the classifier but in-domain for the
 230 DDPM (ChesX-ray8/NIH) and (iii) out-of-domain for both (Montgomery).

231 Figure 2 shows the Dice coefficients from the step-wise experiment when training our segmen-
 232 tation model, the baseline and LEDM on $n = \{197, 49, 24, 12, 6, 3, 1\}$ JSRT labelled datapoints,
 233 corresponding to $\{100, 50, 25, 12, 6, 3, 2, 1\}$ % of the training dataset. Surprisingly, LEDM does
 234 not significantly³ outperform the baseline in the one-shot setting for domain-shifted datasets (NIH,
 235 Montgomery). This indicates that LEDM may not fully utilise the latent representation information.
 236 Secondly, we find that the predictor trained on a single step $t = 1$ statistically outperforms both
 237 LEDM and the baseline for small training sizes (1, 3, 6 in NIH and Montgomery and for one datapoint
 238 in JSRT). In addition, this predictor remains competitive with both the baseline and LEDM across all
 239 other training dataset sizes.

240 The experiment highlights that latent representations obtained from smaller steps are more powerful
 241 predictors than those obtained from larger steps, particularly for domain generalisation. In particular,
 242 the LEDM steps 50, 125 and 250 are not the optimal choice for segmentation as single-step approaches

³Significance is calculated through a Wilcoxon paired test at level 0.05.

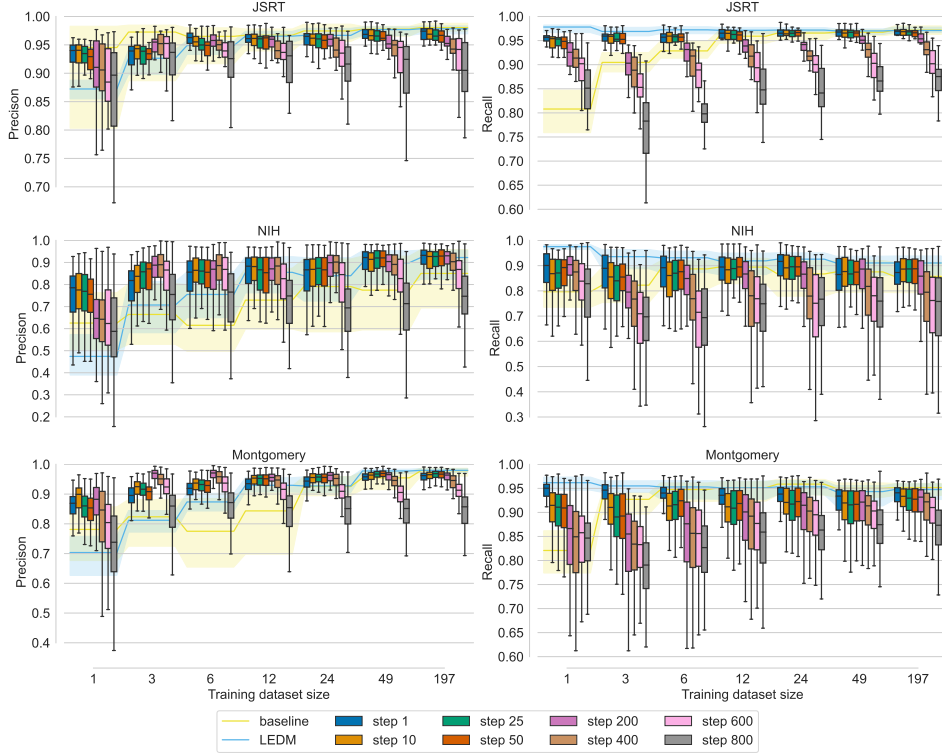


Figure 3: Additional results on the performance of a logistic regression segmentation model trained on latent features from individual diffusion steps.

243 with smaller steps perform better on out-of-distribution datasets. In the next section, we investigate
 244 whether ensembling different steps can still outperform single-step approaches given the right
 245 choice of steps. We investigate several ways of ensembling these steps and their impact on model
 246 generalisation.

247 C Methods details

248 C.1 UK Biobank data preprocessing

249 The UK Biobank brains dataset contains 42 791 patients’ scans. We initially separate the data in three
 250 sets, a training set with $n_{train} = 34\,230$, a validation set with $n_{val} = 4280$ and a test set of $n_{test} =$
 251 4280 patients. After evaluating some methods with $n_{test} = 4280$ and careful consideration of results
 252 variance, we reduced the test set to $n_{test} = 500$ without suffering any drops in metrics accuracy.

253 All scans have voxel size $1mm^3$ and image size $189 \times 233 \times 197$, and are paired with the seg-
 254 mentation of 15 subcortical structures’ volumes from FIRST (FMRIB’s Integrated Registration and
 255 Segmentation Tool Patenaude et al. [2011]) segmentation, and brain masks. For more details on the
 256 scan preprocessing, please refer to Alfaró-Almagro et al. [2018].

257 We preprocess the images by clipping the intensities to $[0, 1500]$ to remove large outliers, then
 258 normalise the brain pixels using the brain masks so that the 1st and 99th quantiles correspond to -1
 259 and 1 respectively:

$$x_{norm}[mask \neq 0] = a \cdot x[mask \neq 0] + b \quad (5)$$

$$\text{such that } a = \frac{2}{x^{99\%} - x^{1\%}} \text{ and } b = 1 - a \cdot x^{99\%} \quad (6)$$

260 where $x^{1\%}$ and $x^{99\%}$ are the 1st and 99th quantiles of $x[mask \neq 0]$.

261 We then split the image and segmentation in 189 2D slices, and discard all slices where no brain
 262 structures are present in the segmentation, resulting in roughly 100 2D slices per brain image.

Table 2: Ablation study on test-time ensembling over timesteps. Each ‘Step i ’ experiment only uses predictions from timestep i at test time.

Training size	1 (1%)	3 (2%)	6 (3%)	12 (6%)	197 (100%)
JSRT (in-domain for classifier)					
Step 1	91.1 \pm 5.0	94.5 \pm 2.1	96.0 \pm 1.4	96.8 \pm 1.1	97.4 \pm 1.3
Step 10	91.6 \pm 4.6	94.6 \pm 1.8	96.0 \pm 1.3	96.9 \pm 1.0	97.4 \pm 1.2
Step 25	91.7 \pm 4.2	94.5 \pm 1.6	95.8 \pm 1.2	96.8 \pm 1.0	97.3 \pm 1.2
TEDM	93.1 \pm 3.4	94.8 \pm 1.4	95.8 \pm 1.2	96.6 \pm 1.1	97.3 \pm 1.2
NIH (in-domain for DDPM, OOD for classifier)					
Step 1	70.4 \pm 10.9	78.9 \pm 9.4	84.2 \pm 8.3	87.5 \pm 6.5	91.9 \pm 3.3
Step 10	73.2 \pm 10.3	81.1 \pm 8.3	85.8 \pm 7.3	88.8 \pm 5.6	91.8 \pm 3.3
Step 25	75.1 \pm 9.8	82.6 \pm 7.7	86.5 \pm 6.7	89.4 \pm 5.2	91.9 \pm 3.3
TEDM	80.3 \pm 9.0	86.4 \pm 6.2	89.2 \pm 5.5	91.3 \pm 4.1	92.9 \pm 3.2
Montgomery (OOD for DDPM and classifier)					
Step 1	85.9 \pm 4.0	89.3 \pm 4.2	92.2 \pm 4.2	93.9 \pm 3.9	94.9 \pm 5.3
Step 10	87.1 \pm 4.5	89.3 \pm 4.8	92.1 \pm 5.2	94.1 \pm 5.0	94.8 \pm 6.5
Step 25	87.4 \pm 5.3	89.1 \pm 5.5	91.7 \pm 6.2	93.7 \pm 6.3	94.6 \pm 7.0
TEDM	90.5 \pm 5.3	91.4 \pm 6.1	93.3 \pm 6.0	94.6 \pm 6.0	95.1 \pm 6.9

263 C.2 BraTS data preprocessing

264 The BraTS dataset consists of 338 patients’ scans. For each patient, four scanner modalities are
 265 available, “native T1, post-contrast T1-weighted (T1Gd), T2-weighted (T2), and T2 Fluid Attenuated
 266 Inversion Recovery (T2-FLAIR) volumes”⁴. Segmentation maps for GD-enhancing tumour, the
 267 peritumoural oedema, and the necrotic and non-enhancing tumour core are provided. In addition, the
 268 scans are co-registered, resampled to $1mm^3$ resolution as skull stripped. For more information about
 269 the BraTS dataset preprocessing, please refer to Bakas et al. [2018], Menze et al. [2014]. We separate
 270 the data in three sets, a training set with $n_{train} = 269$, a validation set with $n_{val} = 36$ and a test set
 271 of $n_{test} = 33$. For each scan modality, we calculate the mean and variance of the brain pixels across
 272 the training set, excluding the background. We use the calculated mean and variance to normalise the
 273 data distribution to mean 0 and standard deviation 1.

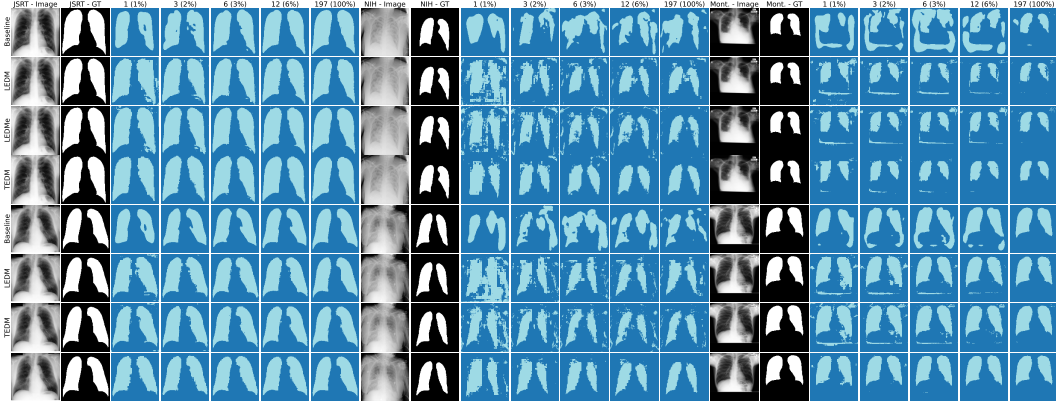
274 We then split the images and segmentation in 155 2D slices. For each slice, concatenate the four
 275 modalities, and take a centre crop of 176×176 .

276 C.3 Training hyperparameters

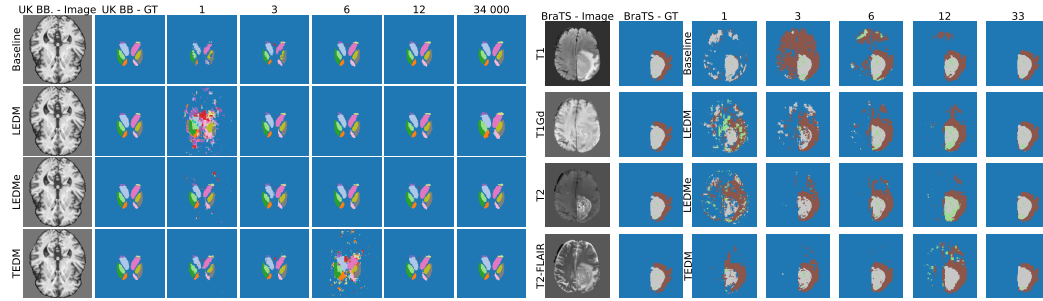
277 We train the DDPM for 100 000 steps with batch size 4 and learning rate $\eta = 0.0001$ on a single
 278 NVIDIA TITAN X GPU with 12GB capacity. Similarly, we train the Global CL and Global & Local
 279 CL models for 100 000 steps. All downstream models - the supervised baseline, Global CL and
 280 Global & Local CL fine-tuning, LEDM, LEDMe and TEDM - are trained for 20 000 steps, with the
 281 same learning rate.

282 D Further results and visualisations

⁴<https://www.med.upenn.edu/cbica/brats2020/data.html>



(a) JSRT (LHS), NIH (middle) and Montgomery (RHS), where NIH and Montgomery are OOD for the classifier, and for the classifier and backbone respectively. Please zoom in for better visibility of details.



(b) UK Biobank

(c) BraTS

Figure 4: Segmentation examples. Col. 1 and 2 are the image and ground truth segmentation. Subsequent columns correspond to models trained with n training datapoints (see title). Row 1 corresponds to the baseline outcomes, and row 2, 3 and 4 to LEDM, LEDMe and TEDM (our method) respectively.

Table 3: Dice scores on the UK Biobank and BraTS datasets. For both datasets, the model was trained on 2D slices, the results are reported on the 3D images. The training size refers to the number of patients in the labelled training set. The number of 2D slices is roughly 100x larger. Here, statistical equivalence is calculated with Bonferroni correction to account for multiple classes per patient.

UK Biobank ($n_{train}^{unlabelled} = 34000, n_{test} = 500$)					
Training size	1	3	6	12	34 000
Sup. Baseline	54.6 ± 18.6	76.8 ± 12.3	83.1 ± 8.5	85.1 ± 7.6	89.6 ± 5.2
Global CL	42.7 ± 20.4	77.3 ± 11.0	82.0 ± 8.7	85.2 ± 7.4	88.7 ± 5.6
Global & Local CL	44.3 ± 20.3	74.0 ± 11.8	80.6 ± 9.4	82.0 ± 8.9	87.4 ± 6.8
LEDM	60.8 ± 17.1	81.3 ± 7.9	82.3 ± 8.9	83.0 ± 9.2	87.7 ± 5.8
LEDMe	54.7 ± 17.8	79.4 ± 10.8	82.5 ± 9.1	83.8 ± 8.6	86.6 ± 7.0
TEDM (ours)	71.0 ± 14.8	81.0 ± 9.0	82.8 ± 8.8	83.2 ± 9.3	85.1 ± 7.4

BraTS ($n_{train}^{unlabelled} = 268, n_{test} = 33$)					
Training size	1	3	6	12	33
Sup. Baseline	12.5 ± 18.9	30.9 ± 31.2	40.7 ± 33.1	47.1 ± 33.8	69.5 ± 25.7
Global CL	4.7 ± 13.6	25.5 ± 29.4	32.3 ± 32.1	40.5 ± 32.0	56.9 ± 28.6
Global & Local CL	11.7 ± 19.1	27.3 ± 30.5	34.1 ± 31.5	38.3 ± 32.2	55.4 ± 30.0
LEDM	24.0 ± 22.9	31.0 ± 31.4	40.8 ± 31.9	48.0 ± 31.2	62.6 ± 26.7
LEDMe	21.2 ± 22.7	33.1 ± 31.4	42.8 ± 32.7	49.5 ± 31.7	63.2 ± 27.6
TEDM (ours)	27.3 ± 26.1	35.6 ± 31.7	41.9 ± 32.3	47.5 ± 31.7	59.8 ± 29.0

Table 4: Models precision and recall w.r.t. ground truth segmentations, as per Table 1.

Training size	1	3	6	12	197
Precision - JSRT (in-domain for classifier)					
Sup. Baseline	89.2 ± 12.1	93.2 ± 7.2	93.8 ± 5.9	95.3 ± 3.7	97.9 ± 1.1
Global CL	86.8 ± 10.5	95.5 ± 3.0	97.3 ± 2.6	97.0 ± 2.0	97.7 ± 1.5
Global & Local CL	90.2 ± 9.2	97.1 ± 2.2	96.8 ± 2.0	96.2 ± 2.0	97.1 ± 1.6
LEDM	85.2 ± 5.8	91.7 ± 2.9	94.1 ± 1.9	96.3 ± 1.5	97.5 ± 1.3
LEDMe	90.1 ± 4.6	93.6 ± 2.3	96.2 ± 2.0	96.6 ± 1.6	97.9 ± 0.9
TEDM (ours)	91.3 ± 7.2	95.4 ± 2.9	95.6 ± 2.0	96.4 ± 1.7	97.5 ± 1.2
Recall - JSRT (in-domain for classifier)					
Sup. Baseline	81.5 ± 6.6	90.6 ± 3.4	93.2 ± 2.6	95.4 ± 2.5	96.8 ± 2.2
Global CL	91.8 ± 3.2	90.2 ± 3.0	90.3 ± 3.8	93.8 ± 1.9	96.6 ± 2.3
Global & Local CL	90.0 ± 3.2	89.5 ± 3.3	89.5 ± 3.7	93.4 ± 2.4	97.2 ± 1.8
LEDM	97.4 ± 1.1	96.7 ± 1.2	97.0 ± 1.6	96.6 ± 2.1	96.6 ± 2.1
LEDMe	97.7 ± 1.0	97.5 ± 1.4	97.1 ± 1.5	97.4 ± 1.3	97.3 ± 1.9
TEDM (ours)	95.4 ± 2.1	94.3 ± 1.6	96.2 ± 1.5	96.9 ± 1.4	97.2 ± 1.9
Precision - NIH (in-domain for DDPM, OOD for classifier)					
Sup. Baseline	63.0 ± 17.0	65.6 ± 18.3	63.6 ± 20.3	72.0 ± 18.3	80.5 ± 17.4
Global CL	60.8 ± 17.9	78.7 ± 15.9	76.0 ± 20.4	83.2 ± 14.4	89.4 ± 13.6
Global & Local CL	65.1 ± 19.1	81.7 ± 15.2	84.5 ± 15.4	81.7 ± 16.8	88.0 ± 13.9
LEDM	48.4 ± 13.6	69.4 ± 14.8	74.7 ± 14.0	83.0 ± 11.4	88.4 ± 9.2
LEDMe	56.8 ± 14.1	69.3 ± 13.7	77.0 ± 12.9	79.8 ± 12.0	90.8 ± 7.8
TEDM (ours)	70.5 ± 13.3	82.0 ± 10.6	86.3 ± 9.3	90.4 ± 6.9	95.3 ± 3.6
Recall - NIH (in-domain for DDPM, OOD for classifier)					
Sup. Baseline	77.7 ± 10.3	80.5 ± 12.0	85.4 ± 10.1	87.4 ± 8.0	84.2 ± 9.9
Global CL	88.6 ± 9.7	83.6 ± 8.1	80.1 ± 13.6	87.4 ± 7.7	85.3 ± 8.9
Global & Local CL	80.9 ± 14.5	78.6 ± 11.7	78.9 ± 14.0	84.0 ± 11.5	87.6 ± 8.5
LEDM	96.4 ± 4.2	91.8 ± 5.5	91.1 ± 6.2	90.2 ± 6.5	89.9 ± 5.5
LEDMe	96.3 ± 3.2	92.5 ± 6.2	91.8 ± 6.7	90.9 ± 7.2	89.9 ± 5.7
TEDM (ours)	95.7 ± 4.0	92.4 ± 4.2	92.9 ± 4.1	92.7 ± 4.4	90.8 ± 5.0
Precision - Montgomery (OOD for DDPM and classifier)					
Sup. Baseline	75.1 ± 16.4	77.6 ± 16.1	73.5 ± 18.6	78.1 ± 19.0	94.9 ± 8.9
Global CL	68.3 ± 18.3	86.7 ± 13.7	88.8 ± 15.8	89.2 ± 13.8	93.7 ± 14.1
Global & Local CL	72.2 ± 20.9	90.1 ± 12.7	92.2 ± 11.0	89.2 ± 14.4	92.9 ± 14.7
LEDM	68.7 ± 10.5	79.4 ± 9.7	85.9 ± 8.8	92.0 ± 6.8	97.5 ± 2.7
LEDMe	69.7 ± 9.2	78.8 ± 9.1	84.8 ± 8.5	88.5 ± 7.3	96.4 ± 3.7
TEDM (ours)	88.7 ± 5.3	90.9 ± 5.9	93.5 ± 4.9	96.9 ± 2.4	98.5 ± 1.0
Recall - Montgomery (OOD for DDPM and classifier)					
Sup. Baseline	80.9 ± 7.2	90.9 ± 5.9	93.0 ± 5.6	93.0 ± 5.8	93.6 ± 4.8
Global CL	88.7 ± 7.2	89.9 ± 4.8	90.1 ± 5.5	92.8 ± 5.7	93.0 ± 6.5
Global & Local CL	86.1 ± 10.9	88.3 ± 5.5	88.4 ± 6.4	92.2 ± 5.9	93.2 ± 6.0
LEDM	94.9 ± 4.7	94.5 ± 4.2	93.9 ± 4.8	92.9 ± 8.3	92.0 ± 9.4
LEDMe	97.0 ± 3.5	96.3 ± 3.7	95.3 ± 4.3	94.3 ± 5.1	94.4 ± 5.1
TEDM (ours)	92.9 ± 6.7	92.4 ± 6.9	93.3 ± 7.1	92.8 ± 7.9	92.6 ± 9.1

Table 5: Precision and recall scores on the UK Biobank and BraTS datasets, as per Table 3

UK Biobank ($n_{train}^{unlabelled} = 34\,000, n_{test} = 500$)					
Training size	1	3	6	12	34 000
Precision					
Sup. Baseline	67.3 ± 18.9	84.5 ± 11.4	84.0 ± 10.7	85.8 ± 9.5	88.7 ± 9.0
Global CL	59.3 ± 23.3	83.1 ± 11.5	82.9 ± 11.1	85.2 ± 9.7	89.4 ± 8.6
Global & Local CL	52.3 ± 22.5	75.1 ± 15.0	80.3 ± 11.5	81.7 ± 10.7	88.6 ± 9.2
LEDM	64.9 ± 21.3	83.2 ± 9.6	84.0 ± 9.6	85.5 ± 9.1	86.9 ± 8.8
LEDMe	51.3 ± 19.5	86.0 ± 8.9	86.4 ± 9.2	85.9 ± 9.0	88.5 ± 8.9
TEDM	85.9 ± 11.7	88.8 ± 8.3	86.8 ± 9.1	87.8 ± 9.0	87.7 ± 9.2
Recall					
Sup. Baseline	41.3 ± 20.5	67.8 ± 16.4	79.7 ± 11.4	82.5 ± 11.2	88.6 ± 6.4
Global CL	30.6 ± 19.6	70.2 ± 14.9	78.8 ± 11.3	82.8 ± 10.4	85.8 ± 9.5
Global & Local CL	39.6 ± 19.4	73.6 ± 11.0	81.1 ± 9.9	82.7 ± 10.1	86.6 ± 7.6
LEDM	64.4 ± 17.7	76.2 ± 13.2	81.4 ± 10.2	81.5 ± 11.2	86.2 ± 8.0
LEDMe	66.0 ± 18.1	75.2 ± 12.7	79.4 ± 11.1	82.5 ± 10.7	85.0 ± 8.0
TEDM	58.6 ± 20.3	73.2 ± 13.3	79.7 ± 11.1	80.0 ± 11.9	83.0 ± 8.6
BraTS ($n_{train}^{unlabelled} = 268, n_{test} = 33$)					
Training size	1	3	6	12	33
Precision					
Sup. Baseline	25.7 ± 30.0	45.1 ± 37.4	54.6 ± 37.6	62.2 ± 35.1	74.1 ± 26.8
Global CL	12.0 ± 25.3	38.6 ± 34.9	48.3 ± 37.1	57.1 ± 34.9	66.6 ± 29.9
Global & Local CL	31.6 ± 35.7	40.5 ± 37.2	49.5 ± 36.3	60.7 ± 35.1	66.3 ± 29.2
LEDM	26.4 ± 28.5	44.5 ± 37.9	56.7 ± 35.8	61.6 ± 35.0	70.6 ± 27.4
LEDMe	27.9 ± 29.4	51.2 ± 37.6	60.8 ± 35.2	61.4 ± 34.8	70.4 ± 27.5
TEDM	46.2 ± 34.2	61.4 ± 35.8	67.2 ± 33.6	67.4 ± 33.4	72.4 ± 27.0
Recall					
Sup. Baseline	18.9 ± 28.4	43.7 ± 36.4	48.1 ± 35.8	49.5 ± 35.5	71.1 ± 26.2
Global CL	13.6 ± 29.1	38.9 ± 36.2	33.1 ± 33.6	45.2 ± 33.5	56.9 ± 30.9
Global & Local CL	21.0 ± 31.3	38.3 ± 35.8	40.6 ± 34.9	39.4 ± 33.3	56.8 ± 31.8
LEDM	35.8 ± 26.7	37.0 ± 34.3	45.8 ± 33.6	51.0 ± 32.0	63.8 ± 26.9
LEDMe	26.8 ± 26.8	36.0 ± 32.9	46.7 ± 34.6	53.1 ± 32.5	64.7 ± 27.7
TEDM	27.6 ± 28.2	37.3 ± 33.3	42.4 ± 33.3	47.9 ± 32.9	59.3 ± 30.2