

Task-agnostic Continual Hippocampus Segmentation for Smooth Population Shifts^{*}

Camila González¹[0000-0002-4510-7309](✉), Amin Ranem¹, Ahmed Othman²,
and Anirban Mukhopadhyay¹

¹ Darmstadt University of Technology, Karolinenplatz 5, 64289 Darmstadt, Germany
`camila.gonzalez@gris.tu-darmstadt.de`

² University Medical Center Mainz, Langenbeckstraße 1, 55131 Mainz, Germany

Abstract. Most continual learning methods are validated in settings where task boundaries are clearly defined and task identity information is available during training and testing. We explore how such methods perform in a task-agnostic setting that more closely resembles dynamic clinical environments with gradual population shifts. We propose ODEX, a holistic solution that combines out-of-distribution detection with continual learning techniques. Validation on two scenarios of hippocampus segmentation shows that our proposed method reliably maintains performance on earlier tasks without losing plasticity.

Keywords: Continual learning · Lifelong learning · Distribution shift.

1 Introduction

Deep learning methods are mostly validated in stationary environments where the train and test data have been carefully homogenized to preserve the i.i.d. assumption. This does not reflect the reality of clinical deployment, where acquisition conditions and disease patterns evolve over time. *Continual learning* (CL) paradigms are being explored by medical imaging researchers [19,22,27] and regulatory bodies [29] as evaluation settings that are better suited for AI in health-care. Continual methods deal with temporal restrictions on data availability by sequentially accumulating knowledge over a stream of *tasks*, each containing data from a different distribution, without revisiting previous stages.

Yet most CL approaches are validated in settings with *rigid task boundaries* and *known task labels*, which is far from how real dynamic environments behave [7]. When deviating from this simplistic problem formulation, they perform worse than simple baselines [23]. Previous research has established desirable properties for CL methods, illustrated in Fig. 1. These include no reliance on either (1) assumptions on task boundaries during training or (2) access to task identity labels, i.e. the method should be *task-agnostic* [10]. In addition, the model should (3) preserve previous knowledge while (4) maintaining sufficient plasticity to

^{*} Supported by the Bundesministerium für Gesundheit (BMG) with grant [ZMVI1-2520DAT03A].

learn new tasks and (5) not require additional computational resources during training [7,10]. The last three objectives are often deemed to be orthogonal, i.e. most approaches either *catastrophically forget* previous knowledge (too plastic), cannot learn new tasks (too rigid) or the training time and resource requirements grow linearly with the number of tasks.

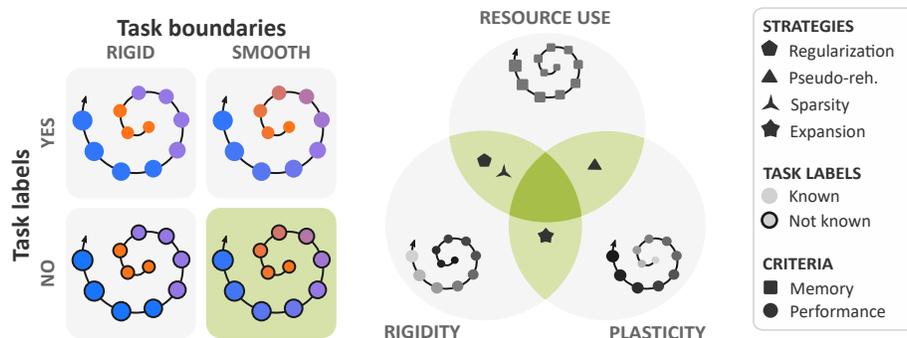


Fig. 1. Desiderata for continual learning [7,10]. Left: methods should not rely on rigid boundaries or task labels. Right: trade-off between plasticity, rigidity and resource use.

Methods for task-agnostic continual learning are overwhelmingly *rehearsal-based* [1,2,12,21,27], i.e. store a subset of past images or features in a memory buffer, which is not admissible in many diagnostic settings due to patient privacy considerations. *Active learning* methods also exist which rely on expert interaction [22].

Other approaches train generative models to identify distribution shifts [24] or only update the shortest sub-path of the network that allows a correct classification [6], but such solutions are computationally expensive and are therefore only evaluated in low-resolution classification settings. The field of continual learning for medical segmentation is still under-studied. Most research follows regularization-based strategies that calculate the importance of parameters and penalize their deviation [19,30]. Approaches have also been proposed for active learning [31], others allow the storage of previous samples [21,28]. Some methods leverage feature disentanglement to alleviate forgetting [16,18] or maintain task-dependent batch normalization layers [13]. To our knowledge, no method has been previously introduced for semantic segmentation that is task-agnostic and does not make use of a rehearsal component.

We propose **ODEx**, an expansion-based approach that (1) does not revisit previous stages, (2) is well-suited to a wide array of use cases, including semantic segmentation and (3) is task-agnostic, i.e. requires neither task boundaries nor task labels during training or inference. *ODEx* uses continual out-of-distribution (OOD) detection to signal when to *expand* the model and select the best parameters during inference. Although we maintain multiple parameter states in

persistent memory, each occupies less than 0.2 GB and the continual OOD detection mechanism ensures that this number remains low. Unlike other methods, *ODEx* requires the same GPU memory and training time as regular sequential learning. Our contributions include:

1. proposing a task-agnostic continual learning solution suitable for a wide array of deep learning architectures, and
2. introducing a continual OOD detection mechanism that does not require access to early data for estimating the distance to the training distribution.

We explore the problem of hippocampus segmentation in T1-weighted MRIs, which is crucial for the diagnosis and treatment of neuropsychiatric disorders but highly sensitive to distribution shifts [25], for two non-stationary environments. Our results show that *ODEx* outperforms state-of-the-art approaches while adhering to desirable properties for continual learning.

2 Methodology

We start by defining our problem formulation of task-agnostic continual learning. We then introduce *ODEx*, visualized in Fig. 2 (bottom). During training, we accumulate the mean and covariance of batch normalization layers and detect domain shifts with the Mahalanobis distance. When a domain shift occurs, a new model is initialized with the most appropriate parameters and added to the model pool. During inference, we extract predictions with the best model state.

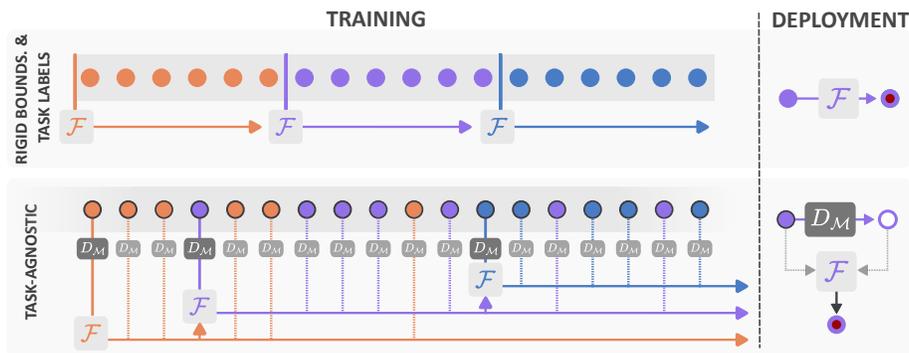


Fig. 2. Top: continual setting with rigid boundaries and task labels. Expansion methods create new parameters at each task boundary. Bottom: the task-agnostic *ODEx* method initializes a new set of parameters when a domain shift is detected.

Task-agnostic continual learning: In continual learning settings, model $\mathcal{F}_\theta : x \rightarrow \hat{y}$ is trained with data samples from an array of N_t different *tasks* or data distributions $\{\mathcal{T}_1 \dots \mathcal{T}_{N_t}\}$, each found at the i_{th} *stage* t_i . The model should be

deployable after finishing the first stage, and evolve over time. For segmentation, each instance has the form (x, y, j) , where x is an image and y the segmentation mask. Additionally, j denotes the *task label*, i.e. that $(x, y) \sim \mathcal{T}_j$. The goal is to find parameters θ that minimize the loss \mathcal{L} over all seen tasks $\{\mathcal{T}_i\}_{i \leq N_t}$ (Eq. 1).

$$\arg \min_{\theta} \sum_{j=1}^{N_t} \mathbb{E}_{(x,y) \sim \mathcal{T}_j} [\mathcal{L}(\mathcal{F}_{\theta}(x), y)] \quad (1)$$

The objective cannot be optimized directly, as at any training stage t_j only data from \mathcal{T}_j is available. The main challenge consists of ensuring enough *rigidity* during training to obtain good performance on $(x, y) \sim \{\mathcal{T}_i\}_{i < j}$ and enough *plasticity* to learn from present and future data $(x, y) \sim \{\mathcal{T}_i\}_{i > j}$.

Expansion-based methods approach this by keeping *task-dependent* parameters $\{\theta_1 \dots \theta_{N_t}\}$, which in their simplest form comprise the entire model, and perform inference on (x, y, j) with the respective \mathcal{F}_{θ_j} (see Fig. 2, above). In *task-agnostic scenarios*, task labels j are unknown and may not even be clearly defined. The goal is to learn a set of parameters $\Theta = \{\theta_1 \dots \theta_{|\Theta|}\}$ and an inference function $\mathcal{J} : x \rightarrow \theta$ that selects the best parameters during testing (Eq. 2). In the absence of rigid task boundaries, the size of the model pool $|\Theta|$ is unknown. Task-agnostic settings thus signify three additional challenges: (1) detecting when domain shifts occur, (2) keeping $|\Theta|$ low and (3) choosing the best parameters during testing. In the following, we outline how we approach these.

$$\arg \min_{\Theta} \sum_{j=1}^{N_t} \mathbb{E}_{(x,y) \sim \mathcal{T}_j} [\mathcal{L}(\mathcal{F}_{\mathcal{J}(x)}(x), y)] \quad (2)$$

Detecting domain shifts: During training, we extract features z from the first set of *Batch Normalization* layers BN_1 . These normalize inputs and thus contain domain-pertinent information which has been found to play a key role in detecting interference during sequential learning [13]. We estimate a multivariate Gaussian $\mathcal{N}_i(\mu_i, \Sigma_i)$ at the end of training stage t_i as:

$$z_k \leftarrow BN_1(x_k); \quad \mu_i \leftarrow \frac{1}{N} \sum_{k=1}^N z_k; \quad \Sigma_i \leftarrow \frac{1}{N} \sum_{k=1}^N (z_k - \mu_i)(z_k - \mu_i)^T \quad (3)$$

Inspired by previous research on OOD detection for semantic segmentation [9], we detect data shifts by calculating the *Mahalanobis distance* $D_{\mathcal{M}}(z; \mu, \Sigma)$ to the training distribution. In contrast to other methods for assessing similarity, such as the *Gram distance* popular in rehearsal-based continual learning [21,22], the Mahalanobis distance requires storing only μ and Σ .

As we cannot revisit data from previous stages, we cannot estimate \mathcal{N} with all data used to train the model. In a situation with slowly shifting data distributions, if we were to only consider the μ and Σ of the last training batch, then we may never detect a sufficiently large distance signaling the need to expand the model pool. We therefore store μ_i and Σ_i at the end of each training stage t_i

and add this to the *history* \mathcal{B}_i of the model which contains information from all pertinent training stages. At stage t_{i+1} , parameters $\hat{\theta}$ are selected that minimize the summed distance of the present training data to the history of $\hat{\theta}$ (Eq. 4).

$$D_{\mathcal{M}}(z; i) : \min_{\theta_j \in \Theta_i} \sum_{(\mu_j, \Sigma_j) \in \mathcal{B}_j} D_{\mathcal{M}_j}(z; \mu_j, \Sigma_j) \quad (4)$$

Managing the model pool: When data arrives for a new stage t_i , the distance $D_{\mathcal{M}}(z; i)$ is calculated and the best model $\hat{\theta}$ is selected. If $D_{\mathcal{M}}(z; i) < \xi$ (case 1), then $\hat{\theta}$ is updated with the current data. Afterwards, μ_i and Σ_i are calculated and added to the model history $\hat{\mathcal{B}}$. If instead $D_{\mathcal{M}}(z; i) \geq \xi$ (case 2), a domain shift is detected and a new model θ_i is initialized with the parameters of $\hat{\theta}$. After a domain shift, the size of the model pool $|\Theta|$ grows by 1. The history of the new model \mathcal{B}_i is initialized with $\hat{\mathcal{B}}$, so the history of each model contains information pertaining to all data distributions used to train it. Following previous research [9] we normalize the distances between the minimum and doubled maximum in-distribution values, and set $\xi = 2\mu$.

Continuing to train older models instead of initializing a new one for each stage has two advantages: (1) the model pool does not grow linearly with the length of the data stream, which would be prohibiting for deployment over long time periods and (2) models can benefit from further training when the data distributions are compatible, potentially allowing positive backwards transfer.

Performing inference: Inference proceeds as illustrated in Fig. 2 (right). For each image, the summed Mahalanobis distance of the test image to each set of parameters $\theta \in \Theta$ is calculated. Again, the best model $\hat{\theta}$ is selected and, in this case, directly used to extract a segmentation mask $\mathcal{F}_{\hat{\theta}}(x) = \hat{y}$.

3 Experimental Setup

We briefly outline how we build our data base of tasks with smooth distribution shifts from publicly available datasets and report relevant aspects of our experimental setup. For further implementation details, we refer the reader to the supplementary material and our code found under <https://github.com/MECLabTUDA/LifeLong-nnUNet>.

Data: We look at two different scenarios of data streams with slowly shifting distributions for segmentation of the entire hippocampus (head, body and tail) in T1-weighted MRIs. The first is constructed from three public datasets: *HarP* [3] contains 135 healthy and Alzheimer’s disease patients, *Dryad* [15] has 25 healthy adult subjects and *Decathlon* [26] contains 130 healthy and schizophrenia patients. We slowly shift the distribution of cases from each source as illustrated in Appendix A. We refer to this scenario as **shifting source**. For the second scenario, henceforth referred to as **transformed**, we slowly modify the *Decathlon* data using the *TorchIO* library [20]. We apply intensity rescaling up to a contrast stretching of (0.1, 0.9) and affine transformations of up to a (0.8, 1.2) scaling range, 15 degrees rotation and 5 mm translation.

Table 1. Performance of the joint training upper bound (first row), sequential learning and six continual learning strategies on the two hippocampus segmentation scenarios.

Method	Shifting source			Transformed		
	Dice \uparrow	BWT \uparrow	FWT \uparrow	Dice \uparrow	BWT \uparrow	FWT \uparrow
Joint	.89 \pm .01			.90 \pm .01		
Seq.	.57 \pm .32	-.19 \pm .12	.14 \pm .09	.87 \pm .03	-.02 \pm .02	.09 \pm .05
EWC	.78 \pm .08	.02 \pm .03	.08 \pm .08	.79 \pm .10	.01 \pm .01	.04 \pm .02
MiB	.67 \pm .24	-.10 \pm .07	.14 \pm .10	.87 \pm .04	-.02 \pm .02	.07 \pm .04
RW	.61 \pm .28	-.15 \pm .10	.14 \pm .10	.87 \pm .03	-.03 \pm .03	.09 \pm .05
PLOP	.57 \pm .32	-.22 \pm .14	.13 \pm .09	.86 \pm .02	-.02 \pm .02	.10 \pm .06
LwF	.51 \pm .35	-.23 \pm .13	.10 \pm .07	.86 \pm .04	-.04 \pm .04	.10 \pm .06
ODEx (ours)	.87 \pm .04	-.03 \pm .02	.14 \pm .09	.89 \pm .01	-.01 \pm .01	.09 \pm .05

Network architecture and training: We use a full-resolution *nnUNet* [11] model for all experiments, with the architecture and training settings selected for the first training stage of each data stream. We perform 200 epochs for each stage, with a loss of *Dice* and *Binary Cross Entropy* weighted equally. All experiments were carried out on a *Nvidia Tesla T4* GPU (16 GB).

Metrics: We report the average Dice on test data from all tasks $\{\mathcal{T}_i\}_{i \leq N_t}$ as well as backwards (BWT) and forwards (FWT) transferability [7,10]. BWT is the *inverse forgetting* and displays to what extent the performance on test samples $(x, y) \sim \mathcal{T}_i$ deteriorates with further training in stages $\{t_i\}_{i > N_t}$. FWT instead measures what impact training on each stage $\{t_i\}_{i \leq N_t}$ has on test data $(x, y) \sim \mathcal{T}_i$. Methods that prevent forgetting show high, realistically close to 0, BWT. FWT is high if enough plasticity is maintained to acquire new knowledge. For both metrics, we report the average over test data from all tasks.

Baselines: In Sec. 4.1, we compare our approach against sequential training and five popular continual learning approaches: Elastic Weight Consolidation (EWC) [14], Modelling the Background (MiB) [4], Riemannian Walk (RW) [5], PLOP [8] and Learning without Forgetting (LwF) [17]. We also report the upper bound of joint training. In most cases, we use the hyperparameters suggested in the corresponding publications or code bases (for more details see Appendix B). For MiB, we reduce the *lkd* to prevent loss explosion. In Sec. 4.2 we perform an ablation study and compare the use of the Mahalanobis distance to other methods proposed within task-agnostic learning, namely using the Gram matrix [21] and detecting domain shifts through a fall in training performance [6].

4 Results

We first compare *ODEx* to state-of-the-art continual learning approaches in Sec. 4.1. Afterwards, we take a closer look at the cumulative Mahalanobis distance for identifying domain shifts and selecting the best parameters (Sec. 4.2).

4.1 Continual learning performance

We compare our proposed approach *ODEx* to five continual learning methods in Tab. 1. The first row shows the upper bound of training a model statically with all training data. *Sequential* results show the deterioration of the performance in earlier tasks as training is carried out, and the following rows display how five continual learning strategies alleviate this. From these, only *EWC* maintains performance on earlier tasks, but at the cost of losing model plasticity and being unable to acquire new knowledge. *ODEx* instead reaches a high FWT showing effective learning on later tasks while still performing well on data from the first training stages. This behavior is further illustrated in Fig. 3 (left), where the per-task performance is plotted for *EWC*, which successfully retains old knowledge, *MiB*, which reaches a high Dice on later tasks, and *ODEx* that performs well on data from all stages. This is particularly clear for the more difficult *shifting source* case, but a Wilcoxon one-sided signed-rank test affirms that *ODEx* significantly outperforms all other approaches in terms of Dice score for both scenarios.

As for resource utilization, *ODEx* requires no more GPU memory than sequential training, as we update one model at a time. The estimation of Σ and the calculation of $D_{\mathcal{M}}$ can be carried out in the CPU given the low resolution of z . Fig. 3 (right) shows that *ODEx* takes only marginally longer than training without any method for forgetting prevention. Though several models are stored (two for *shifting source* and four for *transformed*, see Tab. 2) each weights less than 200 MB, being far from a limiting factor in practice.

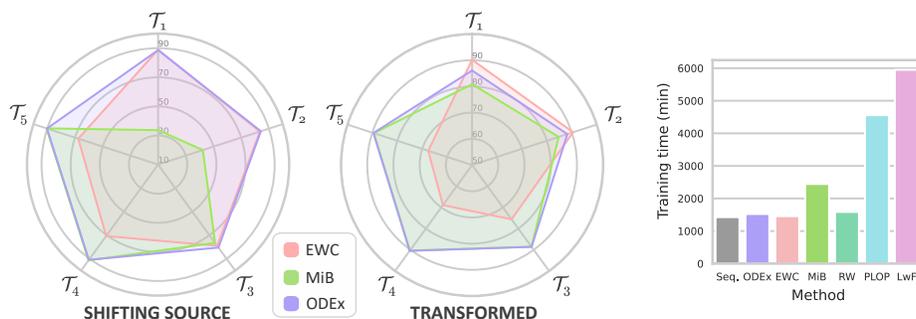


Fig. 3. Left: Per-task Dice. EWC and MiB are at opposite ends of the plasticity/rigidity spectrum, whereas ODEx allows for further training without compromising performance on previous tasks. Right: training times for the shifting source scenario.

Fig. 4 qualitatively shows in the upper row the sequential deterioration of the segmentation for a test subject $(x, y) \sim \mathcal{T}_1$. The lower row displays the segmentation masks produced by each continual learning method. Though the head is mostly segmented well by several methods, only *EWC* and *ODEx* properly segment the body and tail and maintain the integrity of the shape.

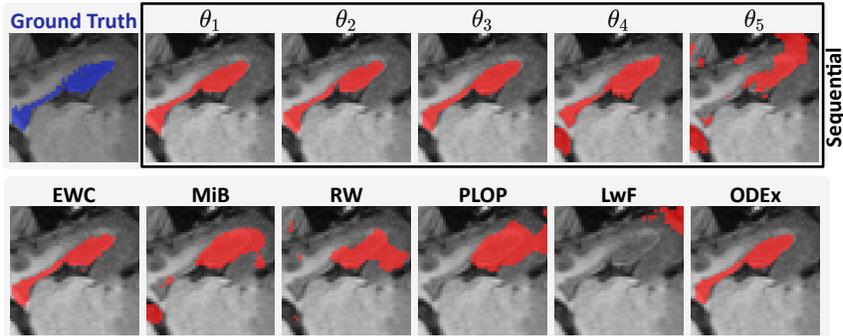


Fig. 4. Crops with overlaid segmentations for axial slice 25 of a subject from \mathcal{T}_1 (shifting source). Top: ground truth (blue) and performance deterioration with regular SGD. Bottom: six continual learning methods, after finishing training on last stage.

4.2 Ablation study

In Tab. 2 we compare our strategy for detecting when to grow the model pool to previous work in the field of task-agnostic learning. The performance of all methods is very similar for the easier *transformed* scenario, but we see clear differences in *shifting source*. We first explore two versions of *ODEx* that use our proposed strategy for selecting the best model but detect domain shifts in a different fashion. *ODEx* $-\infty \xi$ creates a new model for every stage. The lower Dice suggests that the models suffer from the lack of training data, and $|\Theta|$ grows linearly with the number of training stages. *DiceEx* initializes a new model when the training Dice falls more than 10%, which results in higher forgetting. *ODEx* $-\mathcal{B}$ shows the situation where we do not keep a history for the training distributions of previous stages and only calculate the distance to the last stage. For this version, no new model is initialized for *shifting source* and the single available model significantly forgets previous knowledge. Finally, we test the use of the Gram distance instead of Mahalanobis for both training and testing, and find that it does not properly detect distribution shifts for *shifting source*.

Table 2. Performance of different strategies for detecting domain boundaries and/or selecting a model state during inference.

Method	Shifting source				Transformed			
	Dice \uparrow	BWT \uparrow	FWT \uparrow	$ \Theta \downarrow$	Dice \uparrow	BWT \uparrow	FWT \uparrow	$ \Theta \downarrow$
ODEx (ours)	.87 \pm .04	-.03 \pm .02	.14 \pm .09	2	.89 \pm .01	-.01 \pm .01	.09 \pm .05	4
ODEx $-\infty \xi$.83 \pm .04	.00 \pm .00	.11 \pm .09	5	.89 \pm .01	.00 \pm .00	.09 \pm .05	5
DiceEx [6]	.84 \pm .08	-.07 \pm .03	.14 \pm .10	2	.89 \pm .02	-.01 \pm .01	.09 \pm .05	2
ODEx $-\mathcal{B}$ [9]	.57 \pm .32	-.19 \pm .12	.14 \pm .09	1	.89 \pm .01	.00 \pm .00	.09 \pm .05	3
Gram [21]	.57 \pm .32	-.19 \pm .12	.14 \pm .09	1	.90 \pm .01	.00 \pm .00	.09 \pm .05	3

5 Conclusion

We introduce *ODEx*, an expansion-based continual learning strategy suitable for real clinical environments with smooth acquisition and population shifts. We evaluate our approach on two hippocampus segmentation scenarios and show that it outperforms state-of-the-art methods by maintaining good performance on data from early stages without compromising model plasticity. *ODEx* requires only marginally higher training times than regular sequential learning, and the same amount of GPU memory. While additional persistent storage is needed to store different sets of parameters, the OOD detection strategy keeps this number low. Each explored scenario required less than 0.8 GB, rendering this limitation insignificant in practice. Future work should explore whether it suffices to maintain only a subset of domain-specific parameters, such as the last decoder blocks or batch normalization layers. By releasing our code and models, we hope to boost continual learning research in task-agnostic medical settings.

References

1. Aljundi, R., Belilovsky, E., Tuytelaars, T., Charlin, L., Caccia, M., Lin, M., Page-Caccia, L.: Online continual learning with maximal interfered retrieval. *NeurIPS* **32** (2019)
2. Aljundi, R., Lin, M., Goujaud, B., Bengio, Y.: Gradient based sample selection for online continual learning. *NeurIPS* **32** (2019)
3. Boccardi, M., Bocchetta, M., Morency, F.C., Collins, D.L., Nishikawa, M., Ganzola, R., Grothe, M.J., Wolf, D., Redolfi, A., Pievani, M., et al.: Training labels for hippocampal segmentation based on the eadc-adni harmonized hippocampal protocol. *Alzheimer’s & Dementia* **11**(2), 175–183 (2015)
4. Cermelli, F., Mancini, M., Bulò, S.R., Ricci, E., Caputo, B.: Modeling the background for incremental learning in semantic segmentation. In: *CVPR*. pp. 9233–9242 (2020)
5. Chaudhry, A., Dokania, P.K., Ajanthan, T., Torr, P.H.: Riemannian walk for incremental learning: Understanding forgetting and intransigence. In: *ECCV*. pp. 532–547 (2018)
6. Chen, H.J., Cheng, A.C., Juan, D.C., Wei, W., Sun, M.: Mitigating forgetting in online continual learning via instance-aware parameterization. *NeurIPS* **33**, 17466–17477 (2020)
7. Delange, M., Aljundi, R., Masana, M., Parisot, S., Jia, X., Leonardis, A., Slabaugh, G., Tuytelaars, T.: A continual learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021)
8. Douillard, A., Chen, Y., Dapogny, A., Cord, M.: Plop: Learning without forgetting for continual semantic segmentation. In: *CVPR*. pp. 4040–4050 (2021)
9. Gonzalez, C., Gotkowski, K., Bucher, A., Fischbach, R., Kaltenborn, I., Mukhopadhyay, A.: Detecting when pre-trained nnu-net models fail silently for covid-19 lung lesion segmentation. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. pp. 304–314. Springer (2021)
10. Hadsell, R., Rao, D., Rusu, A.A., Pascanu, R.: Embracing change: Continual learning in deep neural networks. *Trends in cognitive sciences* **24**(12), 1028–1040 (2020)

11. Isensee, F., Jaeger, P.F., Kohl, S.A., Petersen, J., Maier-Hein, K.H.: nnu-net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature methods* **18**(2), 203–211 (2021)
12. Jin, X., Sadhu, A., Du, J., Ren, X.: Gradient-based editing of memory examples for online task-free continual learning. *NeurIPS* **34** (2021)
13. Karani, N., Chaitanya, K., Baumgartner, C., Konukoglu, E.: A lifelong learning approach to brain mr segmentation across scanners and protocols. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. pp. 476–484. Springer (2018)
14. Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A.A., Milan, K., Quan, J., Ramalho, T., Grabska-Barwinska, A., et al.: Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences* **114**(13), 3521–3526 (2017)
15. Kulaga-Yoskovitz, J., Bernhardt, B.C., Hong, S.J., Mansi, T., Liang, K.E., Van Der Kouwe, A.J., Smallwood, J., Bernasconi, A., Bernasconi, N.: Multi-contrast submillimetric 3 tesla hippocampal subfield segmentation protocol and dataset. *Scientific Data* **2**(1), 1–9 (2015)
16. Lao, Q., Jiang, X., Havaei, M., Bengio, Y.: Continuous domain adaptation with variational domain-agnostic feature replay. *arXiv preprint arXiv:2003.04382* (2020)
17. Li, Z., Hoiem, D.: Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence* **40**(12), 2935–2947 (2017)
18. Memmel, M., Gonzalez, C., Mukhopadhyay, A.: Adversarial continual learning for multi-domain hippocampal segmentation. In: *Domain Adaptation and Representation Transfer, and Affordable Healthcare and AI for Resource Diverse Global Health*, pp. 35–45. Springer (2021)
19. Özgün, S., Rickmann, A.M., Roy, A.G., Wachinger, C.: Importance driven continual learning for segmentation across domains. In: *International Workshop on Machine Learning in Medical Imaging*. pp. 423–433. Springer (2020)
20. Pérez-García, F., Sparks, R., Ourselin, S.: Torchio: a python library for efficient loading, preprocessing, augmentation and patch-based sampling of medical images in deep learning. *Computer Methods and Programs in Biomedicine* p. 106236 (2021). <https://doi.org/https://doi.org/10.1016/j.cmpb.2021.106236>, <https://www.sciencedirect.com/science/article/pii/S0169260721003102>
21. Perkonigg, M., Hofmanninger, J., Herold, C.J., Brink, J.A., Pianykh, O., Prosch, H., Langs, G.: Dynamic memory to alleviate catastrophic forgetting in continual learning with medical imaging. *Nature Communications* **12**(1), 1–12 (2021)
22. Perkonigg, M., Hofmanninger, J., Langs, G.: Continual active learning for efficient adaptation of machine learning models to changing image acquisition. In: *International Conference on Information Processing in Medical Imaging*. pp. 649–660. Springer (2021)
23. Prabhu, A., Torr, P.H., Dokania, P.K.: Gdumb: A simple approach that questions our progress in continual learning. In: *European conference on computer vision*. pp. 524–540. Springer (2020)
24. Rao, D., Visin, F., Rusu, A., Pascanu, R., Teh, Y.W., Hadsell, R.: Continual unsupervised representation learning. *NeurIPS* **32** (2019)
25. Sanner, A., Gonzalez, C., Mukhopadhyay, A.: How reliable are out-of-distribution generalization methods for medical image segmentation? In: *DAGM German Conference on Pattern Recognition*. pp. 604–617. Springer (2021)
26. Simpson, A.L., Antonelli, M., Bakas, S., Bilello, M., Farahani, K., van Ginneken, B., Kopp-Schneider, A., Landman, B.A., Litjens, G., Menze, B.H., Ronneberger,

- O., Summers, R.M., Bilic, P., Christ, P.F., Do, R.K.G., Gollub, M., Golia-Pernicka, J., Heckers, S., Jarnagin, W.R., McHugo, M., Napel, S., Vorontsov, E., Maier-Hein, L., Cardoso, M.J.: A large annotated medical image dataset for the development and evaluation of segmentation algorithms. CoRR **abs/1902.09063** (2019)
27. Srivastava, S., Yaqub, M., Nandakumar, K., Ge, Z., Mahapatra, D.: Continual domain incremental learning for chest x-ray classification in low-resource clinical settings. In: Domain Adaptation and Representation Transfer, and Affordable Healthcare and AI for Resource Diverse Global Health, pp. 226–238. Springer (2021)
 28. Venkataramani, R., Ravishankar, H., Anamandra, S.: Towards continuous domain adaptation for medical imaging. In: IEEE 16th ISBI. pp. 443–446. IEEE (2019)
 29. Vokinger, K.N., Gasser, U.: Regulating ai in medicine in the united states and europe. *Nature machine intelligence* **3**(9), 738–739 (2021)
 30. Zhang, J., Gu, R., Wang, G., Gu, L.: Comprehensive importance-based selective regularization for continual segmentation across multiple sites. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 389–399. Springer (2021)
 31. Zheng, E., Yu, Q., Li, R., Shi, P., Haake, A.: A continual learning framework for uncertainty-aware interactive image segmentation. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 35, pp. 6030–6038 (2021)

A Data scenarios

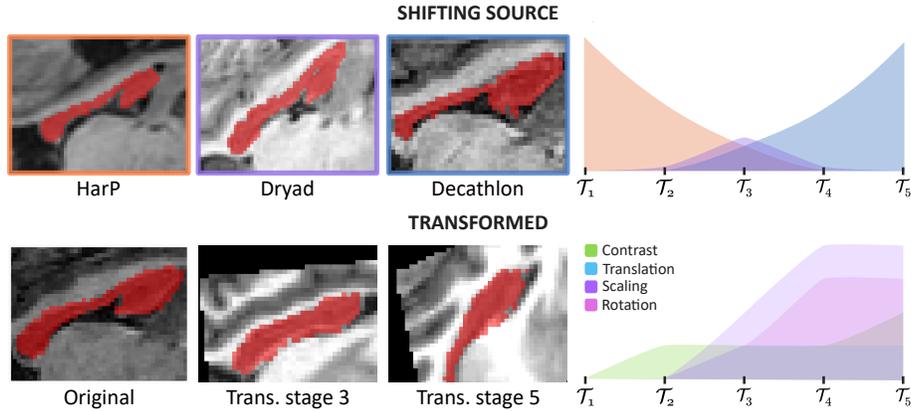


Fig. 1. The two scenarios of data streams with distribution shifts explored in this work. Top: number of cases from three datasets is slowly shifted. Bottom: the *Decathlon* dataset is artificially transformed. We used the first 80/20 split generated by the nnUNet framework for *HarP*, *Dryad* and *Decathlon* and ensured that test cases remained as such across both scenarios.

B Architecture and training parameters

Table 1. Hyperparameters for training continual learning methods. The settings specified in the first row were used for all experiments.

Method	Setting
All	optimizer = SGD, lr = 0.01, weight decay = $3e - 5$, momentum = .99, nr. blocks = 4 for <i>shifting source</i> , nr. blocks = 3 for <i>transformed</i>
EWC	$\lambda = 0.4$
MiB	$\alpha = 0.9$, lkd = 1 for <i>shifting source</i> , lkd = 0.1 for <i>transformed</i>
RW	$\alpha = 0.9$, $\lambda = 0.4$, update after = 10
PLOP	$\lambda = 0.01$, scales = 3, resampling to (48, 48, 48) for <i>transformed</i> , no resampling for <i>shifting source</i>
LwF	$T = 2$

C Calculation of evaluation metrics

Considering $\mathcal{F}_i(x) = \hat{y}_i$ as the prediction made at stage t_i , backwards transfer (BWT) is the change in performance after training with each subsequent task $\{\mathcal{T}\}_{j>i}$, averaged over the number of samples in \mathcal{T}_i and the number of tasks (Eq. 1). BWT is not defined for the last task \mathcal{T}_{N_t} , as $\{t_j\}_{j>N_t} = \emptyset$.

$$BWT = \frac{1}{N_t} \sum_{i=1}^{N_t} \left[\frac{1}{|\{t_j\}_{j>i}|} \sum_{j>i} \left[\frac{1}{|\mathcal{T}_i|} \sum_{k=1}^{|\mathcal{T}_i|} \text{Dice}(\mathcal{F}_j(x_k), y_k) - \text{Dice}(\mathcal{F}_i(x_k), y_k) \right] \right] \quad (1)$$

Forwards transfer (FWT) is, for each task \mathcal{T}_i , the change in performance in each stage before and up to t_i , averaged over the number of samples and tasks. FWT is not defined for the first task \mathcal{T}_1 , as $\{t_j\}_{j<1} = \emptyset$.

$$FWT = \frac{1}{N_t} \sum_{i=1}^{N_t} \left[\frac{1}{|\{t_j\}_{j\leq i}|} \sum_{j\leq i} \left[\frac{1}{|\mathcal{T}_i|} \sum_{k=1}^{|\mathcal{T}_i|} \text{Dice}(\mathcal{F}_j(x_k), y_k) - \text{Dice}(\mathcal{F}_{j-1}(x_k), y_k) \right] \right] \quad (2)$$

D Static learning results

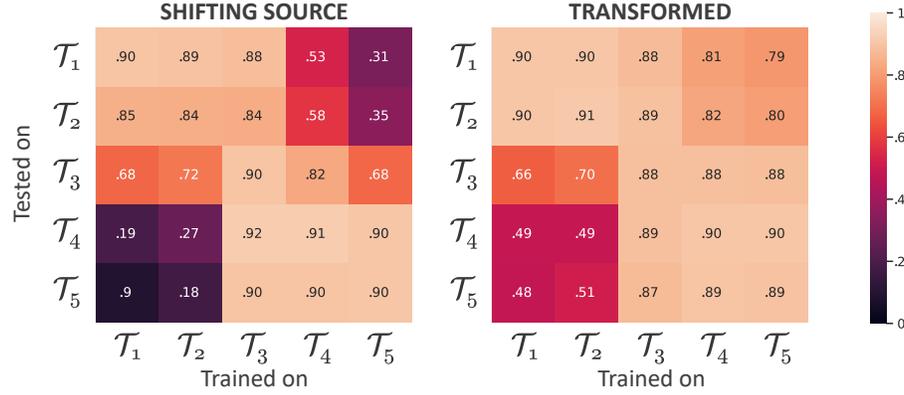


Fig. 2. Base transferability in terms of Dice score of training separate models statically with each task on test data from each task.