
Class-Incremental Continual Learning for General Purpose Healthcare Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Healthcare clinics regularly encounter dynamic data that changes due to variations
2 in patient populations, treatment policies, medical devices, and emerging disease
3 patterns. Deep learning models can suffer from catastrophic forgetting when fine-
4 tuned in such scenarios, causing poor performance on previously learned tasks.
5 Continual learning allows learning on new tasks without performance drop on
6 previous tasks. In this work, we investigate the performance of continual learning
7 models on four different medical imaging scenarios involving ten classification
8 datasets from diverse modalities, clinical specialties, and hospitals. We implement
9 various continual learning approaches and evaluate their performance in these
10 scenarios. Our results demonstrate that a single model can sequentially learn
11 new tasks from different specialties and achieve comparable performance to naive
12 methods. These findings indicate the feasibility of recycling or sharing models
13 across the same or different medical specialties, offering another step towards the
14 development of general-purpose medical imaging AI that can be shared across
15 institutions.

16 1 Introduction

17 Deep Neural Networks (DNNs) have recently exhibited remarkable achievements in various tasks,
18 surpassing human expertise in some cases (7; 12; 14). However their dependence on fixed, balanced
19 datasets within stable environments presents a significant constraint. The ever-changing nature of the
20 real world requires networks capable of sequential learning over time and adapting to shifting data
21 distributions. This shortcoming is especially pronounced in healthcare and medical imaging. The
22 emergence of new diseases, changes in patient population, treatment policies, disease distribution,
23 imaging hardware, or image acquisition techniques can significantly impact the model’s performance.
24 Fine-tuning exclusively on new data, adapts models to the latest targets, resulting in a rapid loss of
25 previously acquired knowledge. Techniques such as Joint Training (JT) are used to overcome this,
26 where the model is trained on both old and new data. However, healthcare data can’t always be shared
27 due to safety concerns and regulation differences across geographies. On the contrary, the NAIVE
28 approach trains an independent model for each task, increasing computational resources required
29 for training, deployment, and missing performance gain due to shared representation. Continually
30 learning is an active area of research, allowing efficient training and adaptation of algorithms to new
31 data without losing prior knowledge. This can improve the model updating, sharing, and resource
32 optimization for healthcare institutions. Cross-sharing models between multiple hospitals can benefit
33 both institutions. Finally, healthcare professionals will be able to screen and detect patients more
34 effectively by quickly identifying and analyzing new biomarkers as they emerge with disease or
35 population shift.

36 In this work, we show the feasibility of training continually learning models for medical imaging.
37 Our contributions are as follows:

- 38 1. We explore the potential of continual learning for sharing medical imaging AI algorithms
39 across changes in hospitals/geographies, medical specialties, and imaging modalities.
- 40 2. We create 4 continual learning scenarios to assess cross-sharing (inter-hospital scenario and
41 one inter-specialty) and intra-specialty model recycling (pathology, radiology) use cases.
- 42 3. We show that continual learning methods can gain performance on par with naive and joint
43 learning approaches while remembering previous tasks.

44 2 Methods

45 We implement 6 variants of continual learning methods, namely memory aware synapses (MAS),
46 replay using memory indexing (REMINDE), MAS with replay (MAS+r), Neuro-inspired stability-
47 plasticity adaptation (NISPA), dark experience replay (DER), and DER++; with prior two not
48 requiring data access to previous tasks. The final four require data access, generally solved by local
49 storage in replay buffers of fixed sizes. MAS (3) and its replay variant MAS+r are regularization
50 methods that calculate the importance of the model parameters in an online fashion. REMIND (9)
51 stores compressed low-level feature representations instead of actual input data, making it well-suited
52 when past data is not feasible. NISPA (8) uses a rewiring mechanism inspired by the structural
53 plasticity of biological neurons and driven by local activations of units similar to Hebbian learning
54 in the brain. DER (4) and DER++ are replay methods that selectively choose examples with high
55 uncertainty to replay. We train naïve learner, and joint learner for baseline comparison. The naïve
56 learner corresponds to training an independent model for each data, and the joint learner is trained on
57 all data together. We use a 5-layered convolutional neural network (CNN) followed by a linear layer
58 classifier as a backbone model. No task labels are provided during testing. To succeed, the model
59 needs to learn inter-task differences to predict task ID on testing and learn intra-task differences to
60 predict the correct class. We measure task accuracy percentage after every episode, average accuracy
61 on seen classes after completing an episode, and backward transfer.

62 2.1 Datasets and scenarios

63 Figure 1 provides a snapshot of scenarios and datasets used. We devise four continual learning
64 scenarios, divided into 3 major categories, to simulate learning new tasks inside the same specialty,
65 different specialties, and hospitals. **Inter-hospital scenario** simulates model sharing across hospitals
66 from different geographies. We used x-rays for pleural effusion, cardiomegaly, atelectasis, and
67 consolidation from the Chexpert dataset, CXR-14 dataset, and VinBig dataset; in total representing
68 3 hospitals and 2 countries. **Inter-Specialty scenarios** simulate model sharing between different
69 medical specialties inside the same hospital and can help specialties with fewer data to benefit from
70 those with more. We combine three specialties; pathology(10), radiology(11), and dermatology(5; 6).
71 **Intra-Specialty scenarios** simulate model rotation within a specialty, mimicking learning new disease
72 finding that appears later on. The pathology scenario contains three subtasks: histology of colorectal
73 cancer (11), blood cells (1), and kidney cortex cells (13). The radiological scenario contains three
74 subtasks: computed tomography (CT) (15)), ultrasound (2), and chest X-ray (11) images.

75 3 Results and Discussion

76 In this section, we discuss the performance of continual learning methods with comparison at the
77 scenario and algorithm level. Figure 1 shows the average accuracy of methods on test data, with
78 every point representing average accuracy on current and all previous tasks. As expected accuracy
79 for NAIVE method takes a sharp dip. Continual learning methods perform on par or better, with
80 NAIVE on the current task, with little to no drop in their performance on the previous task. Across
81 all continual learning methods, replay methods perform better than regularization methods across all
82 scenarios. MAS with replay and DER++ get high accuracies compared to other methods. MAS+r is
83 consistently the best performer in all scenarios, achieving an average accuracy of 88,82,75,79, on inter-
84 hospital, inter-specialty, pathology, and radiology respectively at the end of each scenario. MAS with
85 replay had the highest backward transfer with a value of -2, -5, +3, and -5 on the respective scenarios.
86 Among non-replay methods, MAS has a huge performance drop. REMIND achieves an average
87 accuracy of 83,77,75,80 on inter-hospital, inter-specialty, pathology, and radiology respectively. It's
88 important to note that replay methods have access to a subset of previous data stored for future

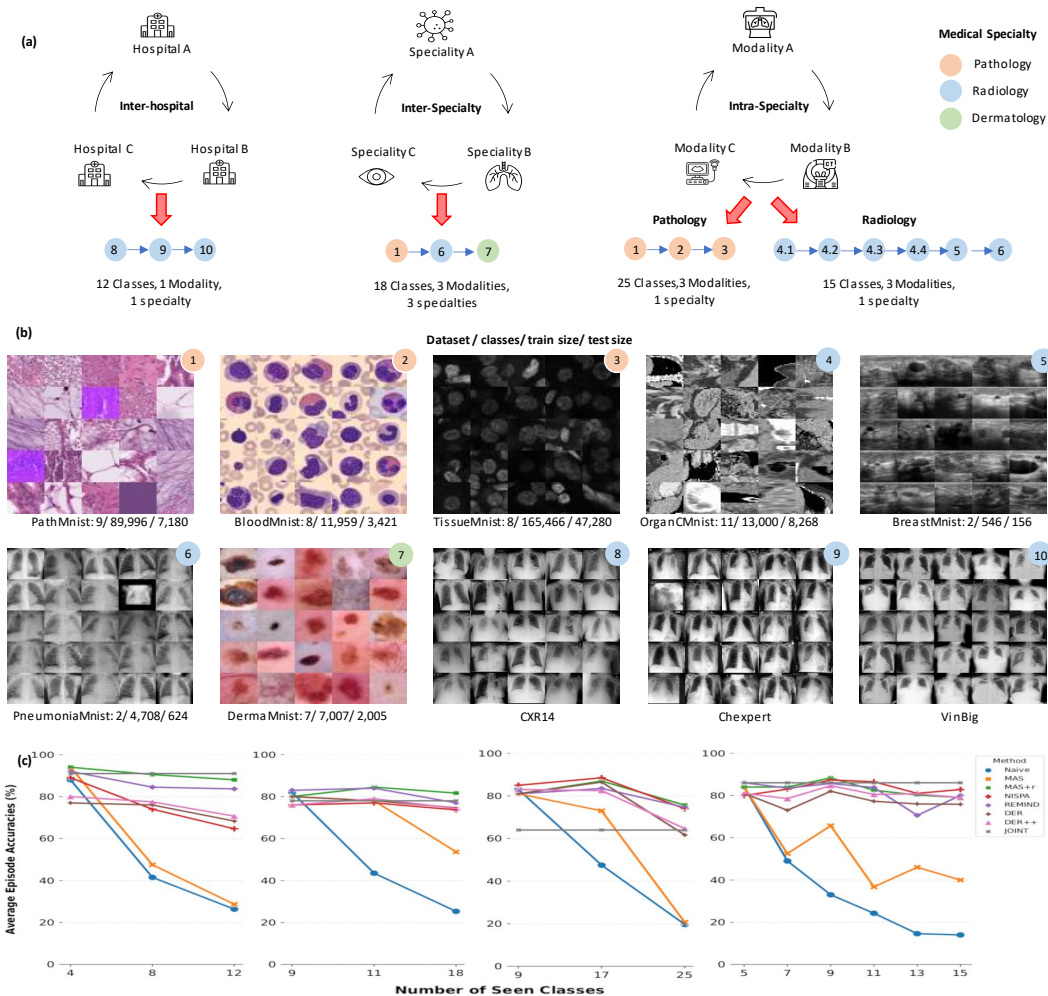


Figure 1: Different learning scenarios and dataset. (a) shows scenarios and sequences of learning tasks used to train the algorithms. (b) shows snapshots of images from different datasets, color-coded (as per their specialty) sequence index on the right top of images. Each dataset has two or more classification labels. Note the visual similarities inside each dataset and the diversity among the different tasks in the scenario. Dataset 4 is split into 4 sub-tasks with unique classes, named 4.1,4.2,4.3 and 4.4

89 retraining. On the other hand, REMIND stores compressed representations, instead of the actual
 90 data, allowing memory efficiency and data privacy. Another critical point to note is REMIND’s
 91 dependence on its feature extractor. Since the initial feature extractor is frozen after its initialization,
 92 the model is less flexible for learning tasks unrelated to initial studies. Using pre-trained weights
 93 from bigger datasets can provide a boost in model performance.

94 **Limitations and Future work** We don’t explore the effect of the sequence of tasks on learning, as
 95 this can impact quality of features learned. We used 32*32*3 image size on a small CNN architecture.
 96 Using higher resolution medical imaging and larger pre-trained models can help boost performance
 97 further. This can be an important future direction.

98 4 Conclusion

99 In this work, we explore the performance of continuous learning methods in varying specialties,
 100 modalities, and geographies. We show that continual learning algorithms can learn new tasks while
 101 maintaining performance on previous tasks, even while changing modalities specialties, and hospitals.

102 This shows potential in developing general-purpose medical imaging AI that can be shared across
103 institutions, with the ability to adapt to new tasks.

104 **5 Potential Negative Societal Impacts**

105 Continual learning models may inherit biases present in the data on which they are trained. If the
106 training data is not representative, these biases can lead to disparities in medical diagnoses and
107 treatment recommendations. While these risks are inherent in deep learning models, automatic
108 unsupervised training of continual learning can exacerbate these biases when deployed, often going
109 unnoticed. Furthermore, due to concerns related to quality control and disparities in deployment
110 regions, continual learning models may inadvertently generate incorrect or misleading medical
111 images or interpretations, which can have detrimental consequences for patients if not adequately
112 monitored and controlled. AI systems are highly sensitive, potentially leading to overdiagnosis and
113 overtreatment if not properly calibrated, thereby resulting in unnecessary medical interventions and
114 increased healthcare costs. Lastly, as with any technological advancement, AI systems are susceptible
115 to hacking and cybersecurity threats. Breaches of medical AI systems could result in unauthorized
116 access to sensitive patient data or manipulation of diagnostic results, posing significant privacy and
117 security risks.

118 **References**

- 119 [1] Acevedo, A., Merino, A., Alférez, S., Molina, , Boldú, L., Rodellar, J.: A dataset of microscopic peripheral
120 blood cell images for development of automatic recognition systems. Data in brief **30**, 105474 (jun 2020).
121 <https://doi.org/10.1016/j.dib.2020.105474>, <http://dx.doi.org/10.1016/j.dib.2020.105474>
- 122 [2] Al-Dhabyani, W., Gomaa, M., Khaled, H., Fahmy, A.: Dataset of breast ultrasound images. Data in brief **28**,
123 104863 (feb 2020). <https://doi.org/10.1016/j.dib.2019.104863>, [http://dx.doi.org/10.1016/j.dib.](http://dx.doi.org/10.1016/j.dib.2019.104863)
124 2019.104863
- 125 [3] Aljundi, R., Babiloni, F., Elhoseiny, M., Rohrbach, M., Tuytelaars, T.: Memory aware synapses: Learning
126 what (not) to forget. In: The European Conference on Computer Vision (ECCV) (2018)
- 127 [4] Buzzega, P., Boschini, M., Porrello, A., Abati, D., Calderara, S.: Dark experience for general continual
128 learning: a strong, simple baseline. In: Advances in Neural Information Processing Systems. vol. 33 (2020)
- 129 [5] Chen, E.Z., Dong, X., Wu, J., Jiang, H., Li, X., Rong, R.: Lesion attributes segmentation for melanoma
130 detection with deep learning. <https://doi.org/10.1101/381855>, [https://www.biorxiv.org/content/](https://www.biorxiv.org/content/10.1101/381855v4)
131 10.1101/381855v4, pages: 381855 Section: New Results
- 132 [6] Codella, N., Rotemberg, V., Tschandl, P., Celebi, M.E., Dusza, S., Gutman, D., Helba, B.,
133 Kalloo, A., Liopyris, K., Marchetti, M., Kittler, H., Halpern, A.: Skin lesion analysis toward
134 melanoma detection 2018: A challenge hosted by the international skin imaging collaboration (ISIC).
135 <https://doi.org/10.48550/arXiv.1902.03368>, <http://arxiv.org/abs/1902.03368>
- 136 [7] Guo, Y., Liu, Y., Oerlemans, A., Lao, S., Wu, S., Lew, M.S.: Deep learning for visual understanding: A
137 review. *Neurocomputing* **187** (2016)
- 138 [8] Gurbuz, M.B., Dovrolis, C.: Nispa: Neuro-inspired stability-plasticity adaptation for continual learning in
139 sparse networks. arXiv preprint arXiv:2206.09117 (2022)
- 140 [9] Hayes, T.L., Kafle, K., Shrestha, R., Acharya, M., Kanan, C.: Remind your neural network to prevent
141 catastrophic forgetting. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK,
142 August 23–28, 2020, Proceedings, Part VIII 16. pp. 466–483. Springer (2020)
- 143 [10] Kather, J.N., Krisam, J., Charoentong, P., Luedde, T., Herpel, E., Weis, C.A., Gaiser, T., Marx,
144 A., Valous, N.A., Ferber, D., Jansen, L., Reyes-Aldasoro, C.C., Zörnig, I., Jäger, D., Brenner, H.,
145 Chang-Claude, J., Hoffmeister, M., Halama, N.: Predicting survival from colorectal cancer histology
146 slides using deep learning: A retrospective multicenter study. *PLoS Medicine* **16**(1), e1002730 (jan
147 2019). <https://doi.org/10.1371/journal.pmed.1002730>, [http://dx.doi.org/10.1371/journal.pmed.](http://dx.doi.org/10.1371/journal.pmed.1002730)
148 1002730
- 149 [11] Kermany, D.S., Goldbaum, M., Cai, W., Valentim, C.C.S., Liang, H., Baxter, S.L., McKeown, A.,
150 Yang, G., Wu, X., Yan, F., Dong, J., Prasadha, M.K., Pei, J., Ting, M.Y.L., Zhu, J., Li, C., Hewett, S.,
151 Dong, J., Ziyar, I., Shi, A., Zhang, R., Zheng, L., Hou, R., Shi, W., Fu, X., Duan, Y., Huu, V.A.N.,

- 152 Wen, C., Zhang, E.D., Zhang, C.L., Li, O., Wang, X., Singer, M.A., Sun, X., Xu, J., Tafreshi, A.,
153 Lewis, M.A., Xia, H., Zhang, K.: Identifying medical diagnoses and treatable diseases by image-based
154 deep learning. *Cell* **172**(5), 1122–1131.e9 (feb 2018). <https://doi.org/10.1016/j.cell.2018.02.010>, <http://dx.doi.org/10.1016/j.cell.2018.02.010>
- 156 [12] LeCun, Y., Bengio, Y., Hinton, G.: Deep learning. *Nature* **521** (2015)
- 157 [13] Ljosa, V., Sokolnicki, K.L., Carpenter, A.E.: Annotated high-throughput microscopy image sets for
158 validation. *Nature Methods* **9**(7), 637 (jun 2012). <https://doi.org/10.1038/nmeth.2083>, <http://www.nature.com/doi/10.1038/nmeth.2083>
- 160 [14] Silver, D., Huang, A., Maddison, C., Guez, A., Sifre, L., Driessche, G., Schrittwieser, J., Antonoglou,
161 I., Panneershelvam, V., Lanctot, M., Dieleman, S., Grewe, D., Nham, J., Kalchbrenner, N., Sutskever, I.,
162 Lillicrap, T., Leach, M., Kavukcuoglu, K., Graepel, T., Hassabis, D.: Mastering the game of go with deep
163 neural networks and tree search. *Nature* **529** (2016)
- 164 [15] Xu, X., Zhou, F., Liu, B., Fu, D., Bai, X.: Efficient multiple organ localization in CT
165 image using 3D region proposal network. *IEEE transactions on medical imaging* (jan 2019).
166 <https://doi.org/10.1109/TMI.2019.2894854>, <http://dx.doi.org/10.1109/TMI.2019.2894854>