

---

# CONTINUAL LEARNING OF LONGITUDINAL HEALTH RECORDS

---

PREPRINT

© **Jacob Armstrong**

Institute of Biomedical Engineering  
Oxford University

`jacob.armstrong@eng.ox.ac.uk`

**David A. Clifton**

Institute of Biomedical Engineering  
Oxford University

`davidc@robots.ox.ac.uk`

October 20, 2021

## ABSTRACT

*Continual learning* denotes machine learning methods which can adapt to new environments while retaining and reusing knowledge gained from past experiences. Such methods address two issues encountered by models in non-stationary environments: ungeneralisability to new data, and the catastrophic forgetting of previous knowledge when retrained. This is a pervasive problem in clinical settings where patient data exhibits covariate shift not only between populations, but also continuously over time. However, while continual learning methods have seen nascent success in the imaging domain, they have been little applied to the multi-variate sequential data characteristic of critical care patient recordings. Here we evaluate a variety of continual learning methods on longitudinal ICU data in a series of representative healthcare scenarios. We find that while several methods mitigate short-term forgetting, domain shift remains a challenging problem over large series of tasks, with only replay based methods achieving stable long-term performance.

Code for reproducing all experiments can be found at <https://github.com/iacobo/continual>

**Keywords** Continual learning · domain adaptation · time series · clinical machine learning · EHR

## I. INTRODUCTION

Clinical and healthcare-related machine learning studies have grown rapidly in recent years, with over a thousand publications annually since 2018 [1]. However many models suffer from ungeneralisability: the distribution of their training data is not representative of the setting in which they are deployed, and hence their real-world performance and utility is overestimated. Further, the distribution of data in a given environment itself continually shifts with time, limiting the use even of models trained on initially representative domains [2, 3].

Unfortunately, naively retraining networks on new data as it becomes available ("fine tuning") commonly results in forgetting of past knowledge. Models can overfit to the specific features of the new dataset, degrading performance on previous tasks in a process known

as *catastrophic forgetting*. This occurs since training on the current task propels updated parameter values far from the previously optimized values (see fig 1). This effectively overwrites learned features pertinent to previous tasks when they are not useful for the current one. While accumulating data and periodically retraining models theoretically alleviates catastrophic forgetting, such approaches are practically encumbered by privacy, storage, and computational hurdles.

Continual learning (CL) has recently emerged as a field to tackle these issues. Models are designed to incrementally update on new datasets while retaining and reusing past knowledge where relevant. Concretely this refers to models which can sequentially train on a series of tasks, while retaining predictive power on previously encountered examples.

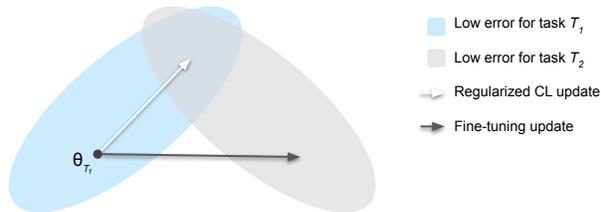


Fig. 1: Under naive transfer learning (grey arrow), there is no guarantee that the parameter values ( $\theta_{T_1}$ ) remain within a region of low error for the previous task  $T_1$  (blue oval) after training on subsequent task  $T_2$ . Regularization techniques like Elastic Weight Consolidation (EWC) enforce such behaviour by penalising the loss, constraining parameter updates to a locus of learned values for previous tasks (figure adapted from [4]).

However a number of state of the art techniques rely on storing past examples and hence may be infeasible in clinical settings due to privacy or data storage limitations. Generative models which create simulated pseudo-examples face further issues of computational limitations.

Further, while a large proportion of Electronic Health Records (EHR) used in patient monitoring and prognostics consists of periodic tabular readings (i.e. multivariate time-series), most current evaluations of continual learning methods are in the image domain [5, 6]. Current benchmarks do not adequately capture the realistic issues faced in a clinical context (e.g. highly imbalanced classes, large multivariate sequences, sparse recordings) [7], and hence the generalisability of their results to these contexts is unclear.

*Contributions:* In this work we present a set of representative continual learning scenarios in the medical domain derived from the open-access eICU-CRD and MIMIC-III ICU datasets [8, 9, 10]. We evaluate a range of methods on these problems, the first (to our knowledge) comprehensive study of Continual Learning methods on medical time-series data. Benchmarks demonstrate common domain shifts encountered by clinical systems in the real world, across geographies, time, and population demographics.

*Related work:* Cossu et al. [11] present a comprehensive evaluation of methods on a set of proposed benchmarks for sequence data. We extend on this work by evaluating such methods on real-world clinical scenarios, over a broader array of model architectures (including Transformers and alternative recurrent networks). Aljundi et al. [12] examine imbalanced classification problems and the effect of dropout regularisation but from a task incremental perspective on imaging data. Kiyasseh et al. [13] investigate domain incremental learning on univariate physiological signals but examine only replay based methods. Churamani et al. [14] investigate domain incremental learning across ethnicity and gender but for facial image data, only evaluating regularization based methods. Guo et al. [15] and Alves et al. [16] investigate temporal and institutional domain shift in ICU data, but

from a domain adaptation perspective, considering only a single source and target dataset.

## II. BACKGROUND

### A. Continual Learning Scenarios

The typical continual learning problem consists of a model encountering a sequence of discrete batches of data, corresponding to different ‘tasks’, where data cannot be stored between tasks.<sup>1</sup> For example a clinical decision model updated annually on new hospital data. The data cannot be retained longer than this due to privacy limitations, but we aspire for the model to generalise to the population with each dataset encountered, and not overfit to the most recent batch as is seen in traditional supervised learning.

Problems are typically split into three scenarios [17]:

- **Task Incremental** Here each task is nominally different. In a classification setting this typically corresponds to each pair of tasks having non-overlapping target sets  $Y_i \cap Y_j = \emptyset \forall i \neq j$ .
- **Class Incremental** Here the set of potential targets expands with each task:  $Y_i \subset Y_j \forall i < j$ .
- **Domain Incremental** Here tasks are nominally the same (i.e. the set of targets is identical for all tasks  $Y_i = Y_j \forall i \neq j$ ), but the distribution of input-features changes with each task.

However, as noted by Cossu et al. [11], this does not capture the full breadth of potential scenarios. For example, newly encountered datasets may introduce a mix of domain shifted instances of old classes, new classes, or novel combinations of classes. Maltoni and Lomonaco [18] divide scenarios into multi-task, single-incremental task, and multiple-incremental tasks, along with a secondary classification for new examples containing new instances of old classes, new classes, or both. However, as discussed by [11], several of the proposed categories are unrealistic or rare. For simplicity we use the terminology of Van de Ven and Tolias [17].

### B. Ontology of methods

A number of methods have been proposed in recent years to mitigate catastrophic forgetting, falling under three general archetypes [19]:

- **Regularization** A regularization constraint is added to the loss function, enforcing updated parameter values to lie within a radius of the current value. This has the benefit of a natural Bayesian interpretation where the posterior values after training on task

<sup>1</sup>More general settings exists in which models encounter a stream of incoming data, sometimes referred to as *online learning*. However, since many CL methods are not designed for such settings we stick to the batch case to allow a broader comparison of methods.

$T_i$  inform the priors for task  $T_{i+1}$ . Methods differ in strategies for choosing which parameters to constrain, and to what degree.

- **Rehearsal** A subset of examples (or generated pseudo-examples) from previous tasks are cached and mixed in with each new task’s training set. Methods differ chiefly in the criteria used for choosing examples. Also known as *replay*.
- **Dynamic architectures** A broad variety of techniques where the network architecture itself adapts with new task presentation. Approaches range from hyper-networks with task-specific subnetworks, to initially small networks which add neurons as resources are required to model new tasks. They are broadly characterised by increasing network complexity with number of tasks.

Such architectural features are not mutually exclusive, and may be hybridised in a number of ways. For example, GEM [20], iCARL [21], and FRoMP [22] employ both rehearsal and regularisation elements. More complex ontologies have been proposed to finer categorise such methods [7].

Replay methods achieve state of the art in many scenarios examined in the literature [23, 24]. However, such techniques are often infeasible in real-world settings, where previous examples cannot be stored or shared due to data privacy constraints [5, 25]. Such a problem is not unique to clinical settings, and while *generative* replay models simulating past examples have been proposed [26], sparse and complex sequential data can prohibit learning of an adequate generative distribution function [27].

For an in depth review of continual learning methods generally, we refer to Delange et al. [7], Parisi et al. [19], Luo et al. [28]. For convenience, we briefly outline the methods evaluated in this work below:

#### Regularization approaches:

- **Elastic Weight Consolidation (EWC)** [4] Penalises changes in parameter values relative to the *importance* of parameters to previous task(s). Importances determined via Fisher’s information matrix. Parameters which are important to previous task(s) are highly constrained, and ones of less importance are less constrained during updates.
- **Online EWC** [29] An adaptation of EWC using a running average of task importance penalties, as opposed to distinct penalties for each previous task. Computationally more efficient and tractable for a large number of tasks.
- **Synaptic Intelligence (SI)** [30] Similar to EWC, enforces parameter specific regularization but importances are calculated *online* (i.e. during training) by approximating the effect on loss and gradient update, as opposed to during an additional pass of the network post training.

- **Learning without Forgetting (LwF)** [31] A copy of the model parameters before updating on the current task is stored and compared to the updated version. Parameter values are distilled between both versions for final update. Hence may be categorised as a *functional* regularization strategy.

#### Replay approaches:

- **Replay** Naive storage of a set of random examples per task, which are mixed in with each subsequent task’s training data. May employ more specific storage policies such as class or task-wise balancing of memories.
- **GDumb** [32] A greedy rehearsal method in which the memory buffer is filled with the  $\frac{\text{buffer size}}{n \text{ tasks seen}}$  most recently encountered examples per task. Examples replayed with each new task.
- **Gradient Episodic Memory (GEM)** [20] Stores a set of examples from each task. Selectively updates gradient for a given minibatch on the current task only if the gradient can be projected in a plane which maintains the positivity of the gradient updates for all stored examples.
- **Averaged Gradient Episodic Memory (A-GEM)** [33] Adaptation of GEM considering only the average gradient for a randomly sampled subset of the stored examples.

#### Dynamic approaches:

- **Progressive Neural Network (PNN)** [34] A copy of parameter weights before updating on a new task is stored. If any parameters shift beyond a certain threshold, the previous weights are frozen and cloned to produce a sister neuron with the updated weights. Relies on task identity at inference to ensure shifted ‘sister’ neurons do not interfere with predictions for prior tasks.

## III. EXPERIMENTS

### A. Problem definitions

*Domain Incremental:* We consider 3 natural Domain Incremental experiments, corresponding to  $n$  patient ICU datasets encountered sequentially across time or location. Domain increments correspond to changing:

- time (season) ( $n = 4$ )
- hospital ( $n = 155$ )
- region ( $n = 4$ )

We also consider the following 3 artificial Domain Incremental experiments, simulating imbalanced populations between healthcare environments (due to demographic-specific care in a given institution, or general population imbalance). Domain increments correspond to groups of patients split by:

- age group ( $n = 7$ )

- ethnicity (broad) ( $n = 5$ )
- ICU ward ( $n = 8$ )

The majority of the above domain splits are self explanatory. ICU WARD refers to different types of critical care (i.e. intensive care) unit, which may specialise in cardiac, trauma, neurological etc injuries.

For each task the setting is supervised prediction of a binary outcome (48hr in-hospital mortality). Input data are multivariate time-series, consisting of periodically recorded patient vital signs from an ICU admission. These are sampled at a rate of 1 per hour, and are of duration  $t = 48$  time steps. Static covariates are repeated to the length of the time-varying sequence and concatenated to enable processing by sequential models.

Further experiments on alternative outcomes (acute respiratory failure; shock) and different sequence/prediction window lengths ( $t \in \{4, 12\}$ ) can be found in Appendix A.

Note that REGION and ETHNICITY (BROAD) can be seen as easier, lower resolution versions of the HOSPITAL and ETHNICITY (NARROW) experiments respectively since their domains correspond to non-overlapping supersets of the formers’.

#### Work in progress

**Note:** not all experiments described have been completed (e.g. PNN strategy, Region experiment, transformer architecture for HOSPITAL, class incremental experiments, class incremental experiment, and supplementary experiments on memory size, traditional regularization, alternative outcome definitions, and sequence length).

## B. Experimental setup

*Model architectures:* For each problem, we evaluate 4 basic neural network architectures:

- 1) a dense feedforward network (MLP)
- 2) 1d convolutional neural network (CNN)
- 3) long- short-term memory network (LSTM); and
- 4) transformer

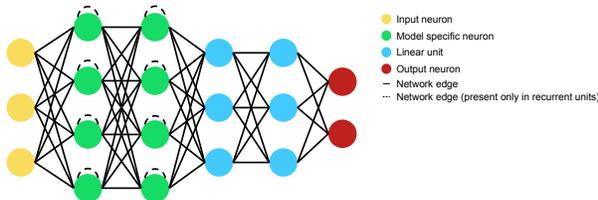


Fig. 2: Generic structure of binary classification model. Example contains 2 hidden ‘feature’ layers of width 4 (green), and 2 fully connected ‘classification’ layers of width 3 (blue).

These were chosen to give a breadth of sequential models, along with a data-structure agnostic model (MLP) for baseline comparison. Further recurrent models are evaluated in Appendix A-B. Models consist of 1 to 4 architecture-specific layers, followed by two dense linear layers (see fig 2).

To enable more fair comparison of methods, model-level parameters (such as number and width of layers) were tuned for the naive baselines and frozen for all other methods (for a given architecture and experiment). Standard regularization features such as dropout were omitted to clearer investigate the effect of the continual learning mechanisms themselves. Additional experiments on the effect of regularization strategies on catastrophic forgetting are found in Appendix C-A. Batch Normalisation was not used in the CNN due to its intensifying effect on catastrophic forgetting [35].

*Strategies:* Each model is equipped with one of the 8 continual learning strategies listed in Table I:

Archetype	Method	Abbreviation	Source
Baseline	Naive fine-tuning	Naive	
	Cumulative multi-task training	Cumulative	
Regularization	Elastic Weight Consolidation	EWC	[4]
	Online EWC	Online EWC	[29]
	Synaptic Intelligence	SI	[30]
	Learning without Forgetting	LwF	[31]
Rehearsal	Naive replay	Replay	
	GDumb	GDumb	[32]
	Gradient Episodic Memory	GEM	[20]
	Averaged Gradient Episodic Memory	AGEM	[33]

TABLE I: Continual Learning methods evaluated.

Rehearsal based methods are given a fixed budget of 256 samples per task, corresponding to approximately 5% and 0.5% of the training data for MIMIC and eICU experiments respectively. See Appendix C-C for experiments on increasing storage capacity.

We further evaluate all models using two baseline methods:

- **Naive:** Naive fine-tuning on each additional task. This is a soft lower bound on performance, equivalent to serial transfer learning with no continual learning mechanism. It is expected to undergo catastrophic forgetting.
- **Cumulative:** Cumulative multi-task training on all tasks seen thus far. This is a soft upper bound on performance, equivalent to transfer learning on a continually expanding dataset, or a rehearsal method with unlimited storage capacity. Note that continual learning methods may outperform this in the instance of strong backwards transfer of information, or on tasks with considerable imbalance in dataset sizes.

*Data:* We use the open-access **eICU-CRD** [9] ICU dataset for all experiments bar seasonal and narrow ethnicity domain increments, for which such information was not available. For these we use the open-access **MIMIC-III** [8, 10] ICU database. For standardisation of

preprocessing and outcome definitions, datasets were preprocessed with the **FIDDLE** pipeline [36]. Data can be accessed at <https://www.physionet.org/content/mimic-eicu-fiddle-feature/1.0.0/>.

Relevant domain shifts identifiable in both datasets is listed in Table II. Full list of domain shifts, along with number of samples in each task, domain, and train/validation/test partition are available in Appendix Table IX.

MIMIC-III	eICU	Domain increment	Number of domains
	✓	Region (US)	4
	✓	Hospital	155
✓	✓	Unit	5-8
✓	✓	Sex	2
✓	✓	Age	6-7
✓	✓	Ethnicity (broad)	5
✓		Ethnicity (narrow)	20
✓		Time (season)	4

TABLE II: Domain shifts annotated in the MIMIC-III and eICU-CRD datasets. When a range of values are given, these correspond to different domains represented across sub-populations with different outcomes (i.e. mortality, Shock, ARF) or datasets (MIMIC-III, eICU).

*Metrics:* We compared the methods using Balanced Accuracy as the main metric.

Since class sizes are highly imbalanced in all experiments (mortality outcome averaging 10% across tasks, see Table IX), and the degree of class imbalance is not constant across domain splits, accuracy is an inappropriate measure of model performance [37]. In minority-event detection, metrics such as sensitivity and specificity (i.e. true positive and true negative rates) are often preferred depending on the relative importance of Type I and Type II errors in the given medical context [38]. To simplify presentation of results, we report the Balanced Accuracy, an average of specificity and sensitivity. Full presentation of sensitivity, specificity, precision, class-accuracy, Area Under the Receiver Operating Curve (AUROC), and Area Under the Precision Recall Curve (AUPRC) can be found in Appendix A.

*Pipeline:*

- 1) **Task split** Data is initially split into several tasks via the task identity (i.e. demographic category for Domain Incremental experiments). For some demographics there were no positive outcomes (e.g. some low volume ethnicity groups). These groups were excluded from the dataset for that experiment. Task order was randomized.
- 2) **Train, validation, test split** Data within each task is then split into train, validation, and test subsets for the first two tasks, and into train, test subsets only for all subsequent tasks in proportions 70:15:15 and 70:30 respectively. Since multiple ICU admissions can pertain to the same patient, train/validation/test streams were split along patient identities to avoid data leakage of similar records

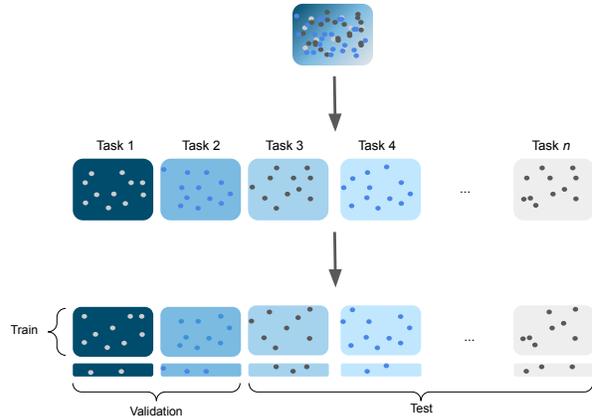


Fig. 3: Data is initially divided into sequential ‘tasks’ split by domain shift. Task order is randomized. The first two tasks are split into training (85%) and validation (15%) sets, the latter used for hyperparameter tuning. Subsequent tasks are assumed to be unavailable for hyperparameter tuning and are split into training (85%) and test data (15%) only. Different colours refer to different domain shifts within the complete dataset.

[39]. Sample counts for each experiment can be found in Table IX.

- 3) **Hyperparameter optimisation** Hyperparameter optimisation requires careful consideration in a continual learning setting, since we should not have access to validation sets from future tasks during the hypothesis generation phase (i.e. model specification). As such, tuning was performed using validation data from the first two tasks only. For setups with a large ( $> 5$ ) number of tasks, these first two tasks are excluded from the final training and testing phase. Otherwise, training and validation data are combined at this stage. This setup is consistent with validation regimes proposed in [33] for Continual Learning setups with a limited number of tasks.

For fairer comparison of methods, generic hyperparameters (i.e. learning rate, batch size, number of layers, hidden depth) were tuned for the Naive baseline run only and frozen for all other methods. Strategy specific hyperparameters were tuned independently for each method.

Hyperparameters were sampled from a range of reasonable values determined from the literature [36, 11]. Where methods shared identical or analogous parameters, the search-space was also shared to ensure fair comparison (for example, regularization strength in EWC, SI, and LwF). A full list of hyperparameter search spaces and the best performing configurations for each model can be found in Appendix B-A.

Hyperparameters were chosen which maximised the average balanced accuracy of the validation predictions for the first two tasks.

- 4) **Training** Once hyper-parameters were selected, each model/strategy combination was trained from scratch on the sequence of tasks' training data. In Appendix C-D we present an extra experiment on the impact of task ordering (cf. *curriculum learning*). The objective function of training was minimising the weighted cross entropy of predictions. Weights are determined by the inverse proportion of class examples in the first two tasks' training data.
- 5) **Evaluation** Models were evaluated on each task's test data, with balanced accuracy, forgetting, and weighted cross entropy loss recorded. Per-task and average metrics were recorded at the end of each training epoch. Training and evaluation was repeated from random initialisation 5 times. Mean performance and bootstrapped 95% confidence intervals are reported.

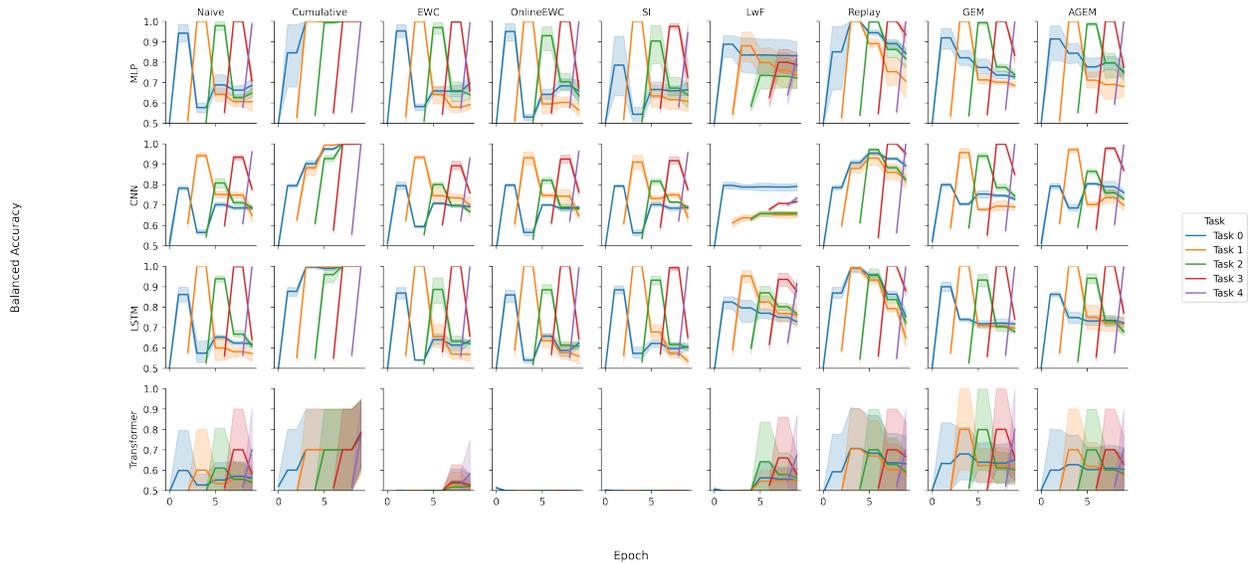


Fig. 4: Domain Incremental results for an outcome of mortality (48h) across domain shift of different ICU WARD. Coloured lines show the training balanced accuracy for each task as it is encountered, for each model and strategy. Shaded regions (left) and black bars (right) refer to bootstrapped 95% confidence intervals. Naive methods (first column) notably undergo catastrophic forgetting as new tasks are introduced. Cumulative training (second column) mitigates this. Regularization methods mitigate this to a degree for the most recently encountered task(s), but do not maintain performance across the entire history of tasks (with the notable exception of LwF). Regularization techniques appear most effective in combination with CNN's. Rehearsal methods achieve greater success across a range of architectures, achieving top performance in most experiments.

		AGE				ETHNICITY (BROAD)				ICU WARD				TIME (SEASON)			
		CNN	LSTM	MLP	Transformer												
Baseline	Cumulative	64.9 $\pm$ 1.6	64.3 $\pm$ 0.8	63.0 $\pm$ 0.4	59.0 $\pm$ 4.9	61.6 $\pm$ 1.1	60.6 $\pm$ 1.2	60.6 $\pm$ 0.8	53.7 $\pm$ 4.8	59.2 $\pm$ 0.9	58.2 $\pm$ 1.8	58.8 $\pm$ 1.0	57.3 $\pm$ 4.1	64.1 $\pm$ 1.1	65.9 $\pm$ 1.4	65.5 $\pm$ 2.0	53.5 $\pm$ 6.8
	Naive	64.0 $\pm$ 0.7	62.8 $\pm$ 1.2	50.0 $\pm$ 0.0	57.8 $\pm$ 3.8	67.6 $\pm$ 0.9	67.6 $\pm$ 0.7	68.8 $\pm$ 0.9	50.0 $\pm$ 0.0	64.2 $\pm$ 1.7	56.6 $\pm$ 2.0	58.7 $\pm$ 1.2	53.7 $\pm$ 4.4	67.2 $\pm$ 2.3	67.8 $\pm$ 1.5	67.6 $\pm$ 1.0	50.0 $\pm$ 0.0
Regularization	EWC	63.8 $\pm$ 1.3	63.6 $\pm$ 0.9	50.0 $\pm$ 0.0	55.7 $\pm$ 5.7	67.3 $\pm$ 1.1	66.8 $\pm$ 1.6	69.1 $\pm$ 0.4	53.5 $\pm$ 6.8	62.9 $\pm$ 1.1	58.8 $\pm$ 3.5	58.6 $\pm$ 1.3	51.9 $\pm$ 3.7	66.4 $\pm$ 1.2	66.9 $\pm$ 0.9	68.0 $\pm$ 0.6	50.0 $\pm$ 0.0
	LwF	<b>64.7</b> $\pm$ 0.9	<b>64.4</b> $\pm$ 0.8	64.3 $\pm$ 0.7	58.4 $\pm$ 4.2	67.5 $\pm$ 0.2	<b>67.8</b> $\pm$ 1.1	69.2 $\pm$ 0.6	52.6 $\pm$ 6.0	64.3 $\pm$ 0.5	<b>61.8</b> $\pm$ 1.0	<b>62.4</b> $\pm$ 2.3	<b>54.5</b> $\pm$ 5.5	67.1 $\pm$ 2.0	67.0 $\pm$ 1.4	67.8 $\pm$ 0.8	52.4 $\pm$ 4.7
	OnlineEWC	63.7 $\pm$ 0.8	63.1 $\pm$ 0.6	<b>64.6</b> $\pm$ 0.5	<b>61.1</b> $\pm$ 0.7	<b>67.8</b> $\pm$ 0.5	66.7 $\pm$ 1.7	<b>70.0</b> $\pm$ 0.8	53.2 $\pm$ 6.3	64.2 $\pm$ 0.8	59.5 $\pm$ 2.4	58.9 $\pm$ 1.9	50.0 $\pm$ 0.0	67.7 $\pm$ 1.1	<b>67.8</b> $\pm$ 0.5	68.1 $\pm$ 0.4	50.0 $\pm$ 0.0
	SI	63.9 $\pm$ 1.4	62.8 $\pm$ 1.9	63.7 $\pm$ 0.3	57.3 $\pm$ 4.3	67.5 $\pm$ 1.0	67.3 $\pm$ 1.6	69.9 $\pm$ 0.4	54.8 $\pm$ 6.5	<b>64.5</b> $\pm$ 0.5	58.9 $\pm$ 1.1	60.6 $\pm$ 1.8	50.0 $\pm$ 0.0	66.1 $\pm$ 0.6	67.6 $\pm$ 0.6	67.6 $\pm$ 0.6	52.7 $\pm$ 5.4
Rehearsal	AGEM	64.5 $\pm$ 1.0	62.2 $\pm$ 0.9	64.1 $\pm$ 0.6	58.0 $\pm$ 4.1	64.8 $\pm$ 2.2	67.3 $\pm$ 1.5	68.7 $\pm$ 0.2	<b>56.1</b> $\pm$ 7.3	63.9 $\pm$ 1.3	59.2 $\pm$ 1.5	60.8 $\pm$ 0.9	53.9 $\pm$ 4.7	<b>68.4</b> $\pm$ 1.6	67.2 $\pm$ 2.1	<b>68.6</b> $\pm$ 0.9	50.0 $\pm$ 0.0
	GEM	63.1 $\pm$ 0.8	60.6 $\pm$ 1.1	61.7 $\pm$ 0.6	58.5 $\pm$ 1.4	58.2 $\pm$ 1.1	57.8 $\pm$ 1.1	60.2 $\pm$ 0.4	50.8 $\pm$ 1.6	60.3 $\pm$ 1.6	57.4 $\pm$ 1.5	57.3 $\pm$ 1.3	53.8 $\pm$ 3.2	60.1 $\pm$ 1.1	60.1 $\pm$ 2.4	63.7 $\pm$ 0.9	54.4 $\pm$ 5.2
	Replay	60.0 $\pm$ 1.2	58.1 $\pm$ 1.8	51.1 $\pm$ 2.2	59.0 $\pm$ 1.6	61.6 $\pm$ 3.7	60.3 $\pm$ 3.6	61.6 $\pm$ 2.1	51.5 $\pm$ 3.0	59.0 $\pm$ 1.7	55.7 $\pm$ 1.6	58.7 $\pm$ 1.5	53.2 $\pm$ 3.8	65.9 $\pm$ 3.0	61.4 $\pm$ 2.3	65.2 $\pm$ 1.8	<b>55.6</b> $\pm$ 4.7

		HOSPITAL (7)			HOSPITAL (14)			HOSPITAL (21)			HOSPITAL (28)			HOSPITAL (35)		
		CNN	LSTM	MLP												
Baseline	Cumulative	57.3 $\pm$ 1.2	55.2 $\pm$ 0.8	56.5 $\pm$ 0.3	62.2 $\pm$ 2.5	61.6 $\pm$ 0.8	61.5 $\pm$ 0.3	57.9 $\pm$ 1.0	60.3 $\pm$ 1.2	60.9 $\pm$ 0.8	54.6 $\pm$ 0.6	55.5 $\pm$ 0.9	56.1 $\pm$ 0.7	56.0 $\pm$ 1.7	56.9 $\pm$ 1.5	56.1 $\pm$ 1.6
	Naive	52.6 $\pm$ 0.1	52.4 $\pm$ 0.3	55.0 $\pm$ 0.1	57.4 $\pm$ 1.4	57.9 $\pm$ 1.8	61.9 $\pm$ 0.8	58.3 $\pm$ 1.8	57.0 $\pm$ 0.9	61.1 $\pm$ 0.9	52.0 $\pm$ 0.5	52.6 $\pm$ 0.7	54.1 $\pm$ 0.4	52.2 $\pm$ 0.4	52.0 $\pm$ 0.4	52.5 $\pm$ 0.1
Regularization	EWC	52.6 $\pm$ 0.0	52.5 $\pm$ 0.1	54.5 $\pm$ 1.1	57.9 $\pm$ 1.4	58.9 $\pm$ 0.5	61.2 $\pm$ 1.1	58.8 $\pm$ 1.7	57.4 $\pm$ 1.6	61.8 $\pm$ 1.0	52.4 $\pm$ 0.6	<b>54.7</b> $\pm$ 1.8	54.2 $\pm$ 0.4	51.9 $\pm$ 0.1	52.5 $\pm$ 0.8	52.5 $\pm$ 0.1
	LwF	52.6 $\pm$ 0.1	52.6 $\pm$ 0.1	55.0 $\pm$ 0.1	56.6 $\pm$ 0.4	57.4 $\pm$ 1.0	61.1 $\pm$ 1.0	58.8 $\pm$ 1.2	57.8 $\pm$ 1.0	61.8 $\pm$ 0.9	51.8 $\pm$ 0.5	53.9 $\pm$ 0.9	54.1 $\pm$ 0.6	51.9 $\pm$ 0.1	51.8 $\pm$ 0.3	52.4 $\pm$ 0.0
	OnlineEWC	52.6 $\pm$ 0.0	52.5 $\pm$ 0.1	55.0 $\pm$ 0.1	57.1 $\pm$ 0.7	58.6 $\pm$ 1.0	61.5 $\pm$ 1.1	58.1 $\pm$ 1.1	57.5 $\pm$ 1.9	61.1 $\pm$ 1.1	51.6 $\pm$ 0.5	53.6 $\pm$ 0.9	54.1 $\pm$ 0.5	52.2 $\pm$ 0.4	52.6 $\pm$ 0.9	52.6 $\pm$ 0.3
	SI	52.6 $\pm$ 0.0	<b>53.7</b> $\pm$ 1.3	54.5 $\pm$ 1.0	58.3 $\pm$ 2.1	58.1 $\pm$ 1.1	61.6 $\pm$ 1.1	57.6 $\pm$ 0.7	57.9 $\pm$ 1.7	61.6 $\pm$ 0.7	51.7 $\pm$ 0.1	51.9 $\pm$ 0.7	53.6 $\pm$ 0.8	52.1 $\pm$ 0.4	52.4 $\pm$ 0.8	52.7 $\pm$ 0.2
Rehearsal	AGEM	52.3 $\pm$ 0.3	52.5 $\pm$ 0.1	56.1 $\pm$ 1.7	57.3 $\pm$ 1.7	57.6 $\pm$ 0.9	<b>62.9</b> $\pm$ 1.5	<b>59.5</b> $\pm$ 1.8	57.3 $\pm$ 3.4	<b>63.3</b> $\pm$ 0.6	51.9 $\pm$ 0.4	53.8 $\pm$ 1.5	<b>56.1</b> $\pm$ 0.9	52.2 $\pm$ 0.4	52.5 $\pm$ 0.8	52.9 $\pm$ 0.4
	GEM	<b>54.8</b> $\pm$ 1.5	50.5 $\pm$ 1.0	<b>56.9</b> $\pm$ 1.3	58.2 $\pm$ 1.4	58.9 $\pm$ 1.6	61.0 $\pm$ 0.8	57.9 $\pm$ 1.8	<b>58.9</b> $\pm$ 0.9	59.3 $\pm$ 1.0	<b>53.0</b> $\pm$ 0.3	54.3 $\pm$ 0.6	55.4 $\pm$ 0.7	<b>54.0</b> $\pm$ 1.1	<b>55.5</b> $\pm$ 1.3	<b>58.1</b> $\pm$ 1.1
	Replay	54.4 $\pm$ 1.3	53.2 $\pm$ 1.2	55.8 $\pm$ 1.9	<b>58.8</b> $\pm$ 1.4	<b>59.7</b> $\pm$ 0.4	62.1 $\pm$ 2.3	57.5 $\pm$ 1.2	56.9 $\pm$ 1.1	59.9 $\pm$ 1.8	52.5 $\pm$ 0.3	53.0 $\pm$ 0.6	53.8 $\pm$ 0.9	52.8 $\pm$ 1.0	52.9 $\pm$ 0.6	52.7 $\pm$ 0.1

TABLE III: Final average balanced accuracy for 48hr mortality prediction across demographic domain shift (AGE, ETHNICITY, WARD) and TIME (top), and HOSPITAL shift (bottom). Average performance over 5 runs are presented with bootstrapped 95% confidence intervals. Bold values refer to the best average performance for each model and experiment. For the hospital experiment we report the current performance after training on  $n$  hospitals for  $n \in \{7, 14, 21, 28, 35\}$  in addition to final performance (i.e. after all hospitals). Bracketed numbers refer to the number of different hospitals sequentially trained on thus far.

## IV. RESULTS

We present the results of the Domain Incremental experiments in Table III. For brevity we show only the results on outcome of 48hr mortality, see Appendix A for results on other outcomes (ARF, shock). Results show the final average test balanced accuracy across all tasks for each method. Reported values are means over 5 runs from random initialisation, with bootstrapped 95% confidence intervals.

For the HOSPITAL domain shift experiments we present the average performance on all tasks thus-seen as the number of tasks increases (i.e. as the models encounter an increasing number of hospitals). Figure 5 displays this performance over time graphically (for the training data).

### A. Model Architectures

Models are generally comparable over a small but constant number (40) of training epochs per domain shift, with the exception of Transformers which demonstrated much more volatile performance over repeated runs.

Highest training efficiency (measured by number of training epochs required to saturate the current task’s loss) was achieved by MLP, followed by LSTM. However a higher training efficiency was correlated with faster (and greater) forgetting upon introduction of new tasks (see for example, MLP EWC vs CNN EWC in Figure 4). We are currently working on introducing an early stopping mechanism to terminate training on each task only once saturation of a given metric has been achieved (as opposed to a fixed number of epochs) to enable a fairer comparison of methods.

### B. Continual Learning strategies

**Regularization** methods showed superior or comparable performance with replay based methods across limited number of domain shifts (AGE, WARD, and ETHNICITY (BROAD), Table III top), but decreasing performance as the number of tasks grew large. LwF achieved superior performance on the largest amount of experiments, achieving the lowest degrees of forgetting (note the ‘flatline’ shape of the LwF task curves in Figure 4). For the HOSPITAL domain shift experiments, regularization methods failed to mitigate catastrophic forgetting for  $n$  tasks  $\geq 5$ , performing on par with Naive fine tuning (no statistically significant difference in final performances). Such performance is expected of regularization methods on domain incremental problems, having been observed in toy problems generally [40, 4], and in recurrent networks specifically [11]. This is likely due to regularization methods only ‘delaying the inevitable’ when faced with a large number of tasks, as models are ‘walled off’ into shrinking locally optimised regions for parameters.

**Rehearsal** methods outperformed all other strategies for a large number of domain shifts. This is consistent with class- and domain-incremental results in other benchmarks [20]. Rehearsal methods all improved with larger storage capacity (Appendix C-C).

As shown in Figure 5, regularization methods were generally volatile across a large number of domain shifts, likely corresponding to sets of hospitals more or less similar to the first few encountered. Contrary to this, the rehearsal methods A-GEM and GEM showed relatively stable performance as more hospitals were encountered. This stability in performance over domain shifts demonstrates sustained generalisation as the task population becomes more heterogeneous.

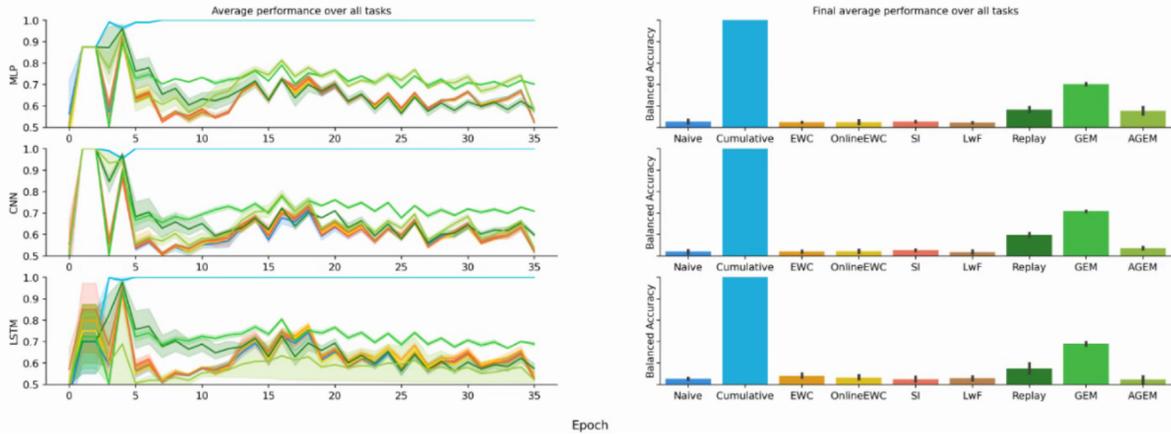


Fig. 5: Domain Incremental results for an outcome of mortality (48h) across domain shift of different HOSPITAL. Results show the average training balanced accuracy over all tasks thus encountered, for each model and strategy. Shaded regions (left) and black bars (right) refer to bootstrapped 95% confidence intervals over 5 runs. Regularization strategies (orange-reds) mitigate catastrophic forgetting to an extent for the first few tasks (hospitals) encountered, but quickly drop to the same performance as the NAIVE baseline (dark blue). A-GEM (lime green) suffers similar behaviour due to averaging of past memory gradients being insufficient to capture the variability in domains. Rehearsal style methods achieve superior performance across the entire range of tasks, with explicit REPLAY achieving the highest performance in all but one instance. No method achieves comparative performance with CUMULATIVE upper bound (light blue) for  $n$  tasks  $\geq 5$ .

## V. DISCUSSION

Our experiments show that simple deep neural networks trained on rich multi-variate sequential data are also prone to catastrophic forgetting in a domain incremental setting.

We observe that regularization methods are prone to more forgetting than rehearsal based methods across a large sequence of tasks, but for few tasks achieve superior or comparable performance to replay based methods (given a fixed small replay buffer).

In the case of patient health records, data may comprise sensitive patient data and hence sharing between institutions or storage over time may require data sharing agreements and ethical approval. This may be prohibitively time-consuming or infeasible, making rehearsal based methods inapplicable. Data-free rehearsal methods such as generative models overcome this issue, but there is a high computational burden to the learning of accurate generative models for such time-series data.

Future work to be performed:

- Implement early stopping mechanism to allow all model/strategies to saturate in current performance before training on new task(s).
- Complete supplementary experiments.
- Investigate domain shift across different countries / healthcare systems / datasets i.e. :
  - MIMIC
  - eICU
  - HIRID [41]
  - AmsterdamUMCDB [42]

- investigate MIMIC-IV. The seasonal information of MIMIC-III appears to be too obfuscated, since models do not seriously undergo catastrophic forgetting in this domain. Use annual information preserved in MIMIC-IV for more realistic experiment.
- Explore continual learning as a means of bias mitigation (compare CL methods on demographic splits with traditional bias mitigation strategies).

## VI. ACKNOWLEDGEMENTS

Jacob Armstrong is supported by the EPSRC Center for Doctoral Training in Health Data Science (EP/S02428X/1).

## REFERENCES

- [1] E Hope Weissler, Tristan Naumann, Tomas Andersson, Rajesh Ranganath, Olivier Elemento, Yuan Luo, Daniel F Freitag, James Benoit, Michael C Hughes, Faisal Khan, et al. The role of machine learning in clinical research: transforming the future of evidence generation. *Trials*, 22(1):1–15, 2021.
- [2] Christopher J Kelly, Alan Karthikesalingam, Mustafa Suleyman, Greg Corrado, and Dominic King. Key challenges for delivering clinical impact with artificial intelligence. *BMC medicine*, 17(1):1–9, 2019.
- [3] Joseph Futoma, Morgan Simons, Trishan Panch, Finale Doshi-Velez, and Leo Anthony Celi. The myth of generalisability in clinical research and machine learning in health care. *The Lancet Digital Health*, 2(9):e489–e492, 2020.
- [4] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- [5] Cecilia S Lee and Aaron Y Lee. Clinical applications of continual learning machine learning. *The Lancet Digital Health*, 2(6):e279–e281, 2020.
- [6] Chaitanya Baweja, Ben Glocker, and Konstantinos Kamnitsas. Towards continual learning in medical imaging. *arXiv preprint arXiv:1811.02496*, 2018.
- [7] Matthias Delange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Greg Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [8] Alistair EW Johnson, Tom J Pollard, Lu Shen, H Lehman Li-Wei, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016.
- [9] Tom Pollard, Alistair Johnson, Jesse Raffa, Leo Anthony Celi, Omar Badawi, and Roger Mark. eICU collaborative research database. *PhysioNet*, 2009. doi:10.13026/C2WM1R. Version 2.0.
- [10] Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. Physiobank, physiotookit, and physionet: components of a new research resource for complex physiologic signals. *circulation*, 101(23):e215–e220, 2000.
- [11] Andrea Cossu, Antonio Carta, Vincenzo Lomonaco, and Davide Bacciu. Continual learning for recurrent neural networks: an empirical evaluation. *arXiv preprint arXiv:2103.07492*, 2021.
- [12] Rahaf Aljundi, Klaas Kelchtermans, and Tinne Tuytelaars. Task-free continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11254–11263, 2019.
- [13] Dani Kiyasseh, Tingting Zhu, and David A. Clifton. CLOPS: continual learning of physiological signals. *CoRR*, abs/2004.09578, 2020.
- [14] Nikhil Churamani, Ozgur Kara, and Hatice Gunes. Domain-incremental continual learning for mitigating bias in facial expression and action unit recognition. *arXiv preprint arXiv:2103.08637*, 2021.
- [15] Lin Lawrence Guo, Stephen R Pfohl, Jason Fries, Alistair Johnson, Jose Posada, Catherine Aftandilian, Nigam Shah, and Lillian Sung. Evaluation of domain generalization and adaptation on improving model robustness to temporal dataset shift in clinical medicine. *medRxiv*, 2021.
- [16] Tiago Alves, Alberto Laender, Adriano Veloso, and Nivio Ziviani. Dynamic prediction of icu mortality risk using domain adaptation. In *2018 IEEE International Conference on Big Data (Big Data)*, pages 1328–1336. IEEE, 2018.
- [17] Gido M Van de Ven and Andreas S Tolias. Three scenarios for continual learning. *arXiv preprint arXiv:1904.07734*, 2019.
- [18] Davide Maltoni and Vincenzo Lomonaco. Continuous learning in single-incremental-task scenarios. *Neural Networks*, 116:56–73, 2019.
- [19] German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *Neural Networks*, 113:54–71, 2019.
- [20] David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. *Advances in neural information processing systems*, 30:6467–6476, 2017.
- [21] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 2001–2010, 2017.
- [22] Pingbo Pan, Siddharth Swaroop, Alexander Immer, Runa Eschenhagen, Richard E Turner, and Mohammad Emtiyaz Khan. Continual deep learning by functional regularisation of memorable past. *arXiv preprint arXiv:2004.14070*, 2020.
- [23] Matthew Riemer, Ignacio Cases, Robert Ajemian, Miao Liu, Irina Rish, Yuhai Tu, and Gerald Tesauero. Learning to learn without forgetting by maximizing transfer and minimizing interference. *arXiv preprint arXiv:1810.11910*, 2018.

- [24] Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaisyasingam Ajanthan, Puneet K Dokania, Philip HS Torr, and Marc’ Aurelio Ranzato. On tiny episodic memories in continual learning. *arXiv preprint arXiv:1902.10486*, 2019.
- [25] Sebastian Farquhar and Yarin Gal. Towards robust evaluations of continual learning. *arXiv preprint arXiv:1805.09733*, 2018.
- [26] Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. *arXiv preprint arXiv:1705.08690*, 2017.
- [27] Benjamin Ehret, Christian Henning, Maria R Cervera, Alexander Meulemans, Johannes von Oswald, and Benjamin F Grewe. Continual learning in recurrent neural networks with hypernetworks. *arXiv preprint arXiv:2006.12109*, 2020.
- [28] Yong Luo, Liancheng Yin, Wenchao Bai, and Keming Mao. An appraisal of incremental learning methods. *Entropy*, 22(11):1190, 2020.
- [29] Jonathan Schwarz, Wojciech Czarnecki, Jelena Luketina, Agnieszka Grabska-Barwinska, Yee Whye Teh, Razvan Pascanu, and Raia Hadsell. Progress & compress: A scalable framework for continual learning. In *International Conference on Machine Learning*, pages 4528–4537. PMLR, 2018.
- [30] Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *International Conference on Machine Learning*, pages 3987–3995. PMLR, 2017.
- [31] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
- [32] Ameya Prabhu, Philip HS Torr, and Puneet K Dokania. Gdumb: A simple approach that questions our progress in continual learning. In *European conference on computer vision*, pages 524–540. Springer, 2020.
- [33] Arslan Chaudhry, Marc’ Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient life-long learning with A-GEM. *CoRR*, abs/1812.00420, 2018. URL <http://arxiv.org/abs/1812.00420>.
- [34] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *arXiv preprint arXiv:1606.04671*, 2016.
- [35] Vincenzo Lomonaco, Davide Maltoni, and Lorenzo Pellegrini. Rehearsal-free continual learning over small non-iid batches. In *CVPR Workshops*, pages 989–998, 2020.
- [36] S. Tang, P. Davarmanesh, Y. Song, D. Koutra, M. Sjoding, and J. Wiens. Mimic-iii and eicu-crd: Feature representation by fiddle preprocessing (version 1.0.0). *Journal of the American Medical Informatics Association*, 27(12):1921–1934, 10 2020.
- [37] Subhrajit Roy, Diana Mincu, Eric Loreaux, Anne Mottram, Ivan Protsyuk, Natalie Harris, Emily Xue, Jessica Schrouff, Hugh Montgomery, Ali Connell, et al. Multi-task prediction of organ dysfunction in the icu using sequential sub-network routing. *Journal of the American Medical Informatics Association*, 28(9):1936–1946, 06 2021.
- [38] Steven Hicks, Inga Strüke, Vajira Thambawita, Malek Hammou, Pål Halvorsen, Michael Riegler, and Sravanthi Parasa. On evaluation metrics for medical applications of artificial intelligence. *medRxiv*, 2021.
- [39] Shachar Kaufman, Saharon Rosset, Claudia Perlich, and Ori Stitelman. Leakage in data mining: Formulation, detection, and avoidance. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 6(4): 1–21, 2012.
- [40] Xuejun Han and Yuhong Guo. Continual learning with dual regularizations. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 619–634. Springer, 2021.
- [41] Hugo Yèche, Rita Kuznetsova, Marc Zimmermann, Matthias Hüser, Xinrui Lyu, Martin Faltys, and Gunnar Ratsch. Hirid-icu-benchmark—a comprehensive machine learning benchmark on high-resolution icu data. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021.
- [42] Patrick J Thoral, Jan M Peppink, Ronald H Driessen, Eric JG Sijbrands, Erwin JO Kompanje, Lewis Kaplan, Heatherlee Bailey, Jozef Kesecioglu, Maurizio Cecconi, Matthew Churpek, et al. Sharing icu patient data responsibly under the society of critical care medicine/european society of intensive care medicine joint data science collaboration: The amsterdam university medical centers database (amsterdamumcdb) example. *Critical care medicine*, 49(6):e563, 2021.
- [43] Vincenzo Lomonaco, Lorenzo Pellegrini, Andrea Cossu, Antonio Carta, Gabriele Graffieti, Tyler L Hayes, Matthias De Lange, Marc Masana, Jary Pomponi, Gido M van de Ven, et al. Avalanche: an end-to-end library for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3600–3610, 2021.
- [44] Idar Johan Brekke, Lars Håland Puntervoll, Peter Bank Pedersen, John Kellest, and Mikkel Brabrand. The value of vital sign trends in predicting and monitoring clinical deterioration: A systematic review. *PloS one*, 14(1):e0210875, 2019.

APPENDIX A  
FULL RESULTS

A. *Additional outcomes*

Here we present the results for predicting additional outcomes omitted in the main results for brevity, namely the domain incremental experiments on outcomes of ARF (4h) and Shock (4h):

		ARF (4h)				ETHNICITY			
		AGE							
		CNN	LSTM	MLP	Transformer	CNN	LSTM	MLP	Transformer
Baseline	Cumulative	67.9 $\pm$ 0.7	64.8 $\pm$ 1.2	65.2 $\pm$ 1.6	67.2 $\pm$ 0.7	66.1 $\pm$ 0.3	66.0 $\pm$ 1.2	69.5 $\pm$ 0.3	67.8 $\pm$ 1.3
	Naive	67.2 $\pm$ 0.3	66.7 $\pm$ 0.6	62.1 $\pm$ 0.8	65.8 $\pm$ 1.0	69.0 $\pm$ 0.8	68.5 $\pm$ 1.2	68.3 $\pm$ 0.5	66.8 $\pm$ 1.8
Regularization	EWC	66.7 $\pm$ 0.5	<b>66.7</b> $\pm$ 1.3	62.7 $\pm$ 0.7	64.2 $\pm$ 0.9	68.7 $\pm$ 0.3	69.1 $\pm$ 0.7	68.3 $\pm$ 0.7	66.6 $\pm$ 1.2
	LwF	67.3 $\pm$ 0.1	66.2 $\pm$ 0.8	<b>65.4</b> $\pm$ 1.4	65.4 $\pm$ 1.0	69.0 $\pm$ 0.1	<b>69.4</b> $\pm$ 0.6	67.6 $\pm$ 1.5	<b>67.4</b> $\pm$ 1.5
	OnlineEWC	67.1 $\pm$ 0.6	65.2 $\pm$ 1.1	62.8 $\pm$ 0.9	65.1 $\pm$ 1.7	68.4 $\pm$ 0.8	68.7 $\pm$ 0.5	<b>68.5</b> $\pm$ 0.4	65.4 $\pm$ 1.0
	SI	<b>67.4</b> $\pm$ 0.2	66.5 $\pm$ 0.9	63.1 $\pm$ 0.9	65.2 $\pm$ 1.0	<b>69.1</b> $\pm$ 0.2	69.1 $\pm$ 0.7	68.1 $\pm$ 1.1	66.7 $\pm$ 2.0
Rehearsal	AGEM	66.4 $\pm$ 0.6	66.2 $\pm$ 0.6	59.3 $\pm$ 0.6	<b>65.6</b> $\pm$ 0.9	68.8 $\pm$ 0.6	68.9 $\pm$ 0.1	68.2 $\pm$ 0.7	64.4 $\pm$ 5.1
	GEM	61.3 $\pm$ 0.3	59.9 $\pm$ 0.7	57.7 $\pm$ 1.1	61.9 $\pm$ 0.4	68.3 $\pm$ 0.2	65.8 $\pm$ 0.8	68.3 $\pm$ 0.3	66.8 $\pm$ 1.0
	Replay	61.3 $\pm$ 2.0	62.8 $\pm$ 0.5	58.5 $\pm$ 0.9	64.1 $\pm$ 0.5	67.1 $\pm$ 1.2	66.8 $\pm$ 1.0	66.9 $\pm$ 0.8	65.3 $\pm$ 1.5

TABLE IV: Results for outcome of 4h Acute Respiratory Failure. Similar to the main results on mortality, regularization methods achieve best performance over a limited number of domain shifts. Transformers achieve much more stable performance over the shorter sequence experiments.

		Shock (4h)				ETHNICITY			
		AGE							
		CNN	LSTM	MLP	Transformer	CNN	LSTM	MLP	Transformer
Baseline	Cumulative	62.3 $\pm$ 0.4	64.3 $\pm$ 1.4	65.0 $\pm$ 0.8	67.3 $\pm$ 0.6				
	Naive	65.0 $\pm$ 0.6	65.9 $\pm$ 0.9	65.3 $\pm$ 0.3	64.1 $\pm$ 1.8				
Regularization	EWC	64.8 $\pm$ 0.5	<b>67.5</b> $\pm$ 0.7	<b>66.4</b> $\pm$ 0.5	63.3 $\pm$ 2.5				
	LwF	65.1 $\pm$ 0.5	66.8 $\pm$ 0.3	65.3 $\pm$ 1.0	<b>65.6</b> $\pm$ 0.5				
	OnlineEWC	64.8 $\pm$ 0.6	67.1 $\pm$ 0.7	65.8 $\pm$ 0.8	63.1 $\pm$ 1.5				
	SI	<b>65.3</b> $\pm$ 0.4	67.5 $\pm$ 0.6	65.1 $\pm$ 0.9	63.5 $\pm$ 1.4				
Rehearsal	AGEM	62.3 $\pm$ 0.7	64.7 $\pm$ 1.3	65.1 $\pm$ 0.7	63.2 $\pm$ 0.4				
	GEM	61.7 $\pm$ 0.5	61.8 $\pm$ 1.1	62.5 $\pm$ 0.6	63.0 $\pm$ 1.1				
	Replay	61.0 $\pm$ 0.5	60.7 $\pm$ 0.7	62.2 $\pm$ 0.9	63.2 $\pm$ 0.9				

TABLE V: Results for outcome of 4h Shock.

In contrast to other outcomes, prediction of shock (4h) shows little variation between the naive baseline, continual learning methods, and cumulative upper bound. This may be due to shock presenting similarly across domain shifts

B. *Additional sequential models*

Work in progress 

Here we evaluate a number of other sequential model architectures omitted from the main results for brevity, namely, we evaluate an RNN and GRU (in addition to the LSTM of the main results).

APPENDIX B  
MODEL AND DATA SPECIFICATIONS

A. Hyperparameters

Hyperparameter tuning was performed via grid search over the following discrete space (parameter names refer to their `kwarg` names in the Avalanche implementations [43]):

Hyperparameter	Values	Hyperparameter	Values	MLP	CNN	LSTM	Transformer
<code>mem_size</code>	{256}	<code>hidden_dim</code>	[64, 128, 256]	✓	✓	✓	✓
<code>patterns_per_exp</code>	{256}	<code>n_layers</code>	[3, 4]	✓	✓	✓	✓
<code>sample_size</code>	{256, 512}	<code>nonlinearity</code>	[relu, tanh*]	✓	✓	✗	✓
<code>ewc_lambda</code>	{0.001, 0.01, 0.1, 1, 10, 100}	<code>n_heads</code>	[12, 16, 24]	✗	✗	✗	✓
<code>si_lambda</code>	{0.001, 0.01, 0.1, 1, 10, 100}	<code>bidirectional</code>	[True, False]	✗	✗	✓	✗
<code>lambda_e</code>	{0.001, 0.01, 0.1, 1, 10, 100}						
<code>alpha</code>	{0.001, 0.01, 0.1, 1, 10, 100}						
<code>temperature</code>	{0.5, 1.0, 1.5, 2.0, 2.5, 3.0}						
<code>decay_factor</code>	{0.2, 0.4, 0.6, 0.8, 0.9, 1}						
<code>memory_strength</code>	{0.2, 0.4, 0.6, 0.8, 0.9, 1}						

TABLE VI: Grid for method hyperparameter search for all experiments. Left table refers to strategy specific hyperparameters. Right table refers to model specific hyperparameters. Check marks and crosses detail whether hyperparameters are included in the respective model. \*gelu nonlinearity used instead of tanh for the Transformer model.

Tuned parameters for each model (base model and CL strategy) and experiment are listed below:

		AGE					ETHNICITY (BROAD)				
		lambda	decay_factor	temperature	sample_size	patterns_per_exp					
MLP	EWC	0.1					MLP	EWC	0.1		
	OnlineEWC	0.01	0.9					OnlineEWC	0.001	0.2	
	LwF	0.1		1.5				LwF	1.0		0.5
	SI	0.01						SI	100.0		
	Replay				640.0			Replay			
CNN	AGEM				128.0	128.0	AGEM				128.0
	GEM			0.2		128.0	GEM		0.8		128.0
	EWC	0.01					EWC	0.001			
	OnlineEWC	100.0	0.5				OnlineEWC	100.0	0.9		
	LwF	0.001		2.0			LwF	100.0		1.5	
LSTM	SI	0.1					SI	0.001			
	Replay				128.0		Replay				128.0
	AGEM				128.0	128.0	AGEM				128.0
	GEM			0.4		128.0	GEM		0.8		128.0
	EWC	0.01					EWC	0.001			
Transformer	OnlineEWC	0.01	0.5				OnlineEWC	1.0	0.2		
	LwF	1.0		3.0			LwF	1.0		1.0	
	SI	0.01					SI	10.0			
	Replay				128.0		Replay				128.0
	AGEM				128.0	128.0	AGEM				128.0
Transformer	GEM			0.6		128.0	GEM		0.8		128.0
	EWC	10.0					EWC	100.0			
	OnlineEWC	0.01	0.6				OnlineEWC	10.0	0.8		
	LwF	0.001		0.5			LwF	10.0		3.0	
	SI	10.0					SI	100.0			
Transformer	Replay				128.0		Replay				128.0
	AGEM				128.0	128.0	AGEM				128.0
	GEM			0.7		128.0	GEM		0.2		128.0

		TIME (SEASON)					ICU WARD				
		lambda	decay_factor	temperature	sample_size	patterns_per_exp					
MLP	EWC	0.1					MLP	EWC	0.1		
	OnlineEWC	0.001	0.9					OnlineEWC	0.001	0.9	
	LwF	0.001		2.5				LwF	10.0		1.5
	SI	0.01						SI	0.001		
	Replay				128.0			Replay			
CNN	AGEM				128.0	128.0	AGEM				512.0
	GEM			0.6		128.0	GEM		0.6		256.0
	EWC	10.0					EWC	100.0			
	OnlineEWC	0.001	0.4				OnlineEWC	1.0	0.8		
	LwF	0.001		2.0			LwF	100.0		2.0	
LSTM	SI	0.1					SI	100.0			
	Replay				128.0		Replay				1280.0
	AGEM				128.0	128.0	AGEM				512.0
	GEM			0.4		128.0	GEM		0.8		256.0
	EWC	100.0					EWC	0.01			
Transformer	OnlineEWC	100.0	0.9				OnlineEWC	1.0	0.4		
	LwF	0.001		1.0			LwF	10.0		3.0	
	SI	0.01					SI	0.01			
	Replay				128.0		Replay				1280.0
	AGEM				128.0	128.0	AGEM				512.0
Transformer	GEM			0.4		128.0	GEM		0.8		256.0
	EWC	0.001					EWC	0.1			
	OnlineEWC	0.001	0.2				OnlineEWC	0.01	0.8		
	LwF	0.1		1.0			LwF	0.1		0.5	
	SI	0.001					SI	0.001			
Transformer	Replay				128.0		Replay				1280.0
	AGEM				128.0	128.0	AGEM				256.0
	GEM			0.4		128.0	GEM		1.0		256.0

TABLE VII: Tuned hyperparameters for main experiments (outcome of MORTALITY (48H)).

### B. Training partitions

Total number and number of positive samples in each train/validation/test split for each experiment:

AGE													
task	0		1		2		3		4		5		Outcome
	Total	Outcome	Total	Outcome	Total	Outcome	Total	Outcome	Total	Outcome	Total	Outcome	
train	10794	760 	10528	998 	10365	1191 	11245	1407 	9385	1476 	1798	309 	
val	2308	185 	2248	236 	2217	221 	2402	300 	2012	331 	372	61 	
test	2273	173 	2214	229 	2179	259 	2348	329 	1964	329 	395	65 	

ETHNICITY												
task	0		1		2		3		4		Outcome	
	Total	Outcome	Total	Outcome	Total	Outcome	Total	Outcome	Total	Outcome		
train	6394	705 	912	101 	41380	4722 	2014	270 	2433	301 		
val	1386	149 	172	26 	8880	1019 	405	55 	500	45 		
test	1303	139 	186	25 	8931	1060 	434	61 	521	57 		

WARD											
task	0		1		2		3		4		Outcome
	Total	Outcome	Total	Outcome	Total	Outcome	Total	Outcome	Total	Outcome	
train	771	92 	762	23 	2654	400 	1108	143 	740	97 	
val	166	21 	169	5 	553	87 	230	20 	160	11 	
test	151	16 	189	8 	540	79 	225	19 	159	10 	

TIME (SEASON)									
task	0		1		2		3		Outcome
	Total	Outcome	Total	Outcome	Total	Outcome	Total	Outcome	
train	1428	186 	1551	189 	1510	178 	1546	202 	
val	322	31 	309	43 	342	35 	305	35 	
test	319	35 	293	37 	338	27 	314	33 	

TABLE VIII: Train, test, validation and outcome breakdowns for 48h mortality. Red and green bars represent proportion of positive and negative outcomes respectively per partition per task. Hospital splits have been omitted due to space constraints.

## C. Domain splits

Dataset	Outcome	MICU	SICU	CSRU	TSICU	CCU	Neuro ICU	Med-Surg ICU	CSICU	CTICU	Cardiac ICU	CCU-CTICU
mimic3	mortality (48h)	✓	✓	✓	✓	✓						
	ARF (4h)	✓	✓	✓	✓	✓						
	Shock (4h)	✓	✓	✓	✓	✓						
	ARF (12h)	✓	✓	✓	✓	✓						
	Shock (12h)	✓	✓	✓	✓	✓						
eicu	mortality (48h)	✓	✓		✓		✓	✓	✓	✓	✓	✓
	ARF (4h)	✓	✓				✓	✓	✓	✓	✓	✓
	Shock (4h)	✓	✓				✓	✓	✓	✓	✓	✓
	ARF (12h)	✓	✓				✓	✓	✓	✓	✓	✓
	Shock (12h)	✓	✓				✓	✓	✓	✓	✓	✓

TABLE IX: Domain shifts exhibited for the subset of patients in each outcome dataset.

Dataset	Outcome	Hispanic	Asian	Other/Unknown	Caucasian	African American	Native American
eicu	mortality 48h	✓	✓	✓	✓	✓	
	ARF 4h	✓	✓	✓	✓	✓	✓
	Shock 4h	✓	✓	✓	✓	✓	✓
	ARF 12h	✓	✓	✓	✓	✓	
	Shock 12h	✓	✓	✓	✓	✓	

TABLE X: Domain shifts exhibited for the subset of patients in each outcome dataset.

Dataset	Outcome	UNKNOWN/NOT SPECIFIED	ASIAN - CHINESE	WHITE	HISPANIC OR LATINO	WHITE - OTHER EUROPEAN	ASIAN - VIETNAMESE	PORTUGUESE	HISPANIC/LATINO - DOMINICAN	BLACK/AFRICAN	UNABLE TO OBTAIN	PATIENT DECLINED TO ANSWER	ASIAN - ASIAN INDIAN	HISPANIC/LATINO - PUERTO RICAN	WHITE - RUSSIAN	ASIAN	MULTI RACE ETHNICITY	BLACK/CAPE VERDEAN	BLACK/HAITIAN	OTHER	WHITE - BRAZILIAN	BLACK/AFRICAN AMERICAN	MIDDLE EASTERN	ASIAN - FILIPINO	HISPANIC/LATINO - GUATEMALAN	WHITE - EASTERN EUROPEAN	
mimic3	mortality 48h	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	
	ARF 4h	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	
	Shock 4h	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	
	ARF 12h	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
	Shock 12h	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

TABLE XI: Domain shifts exhibited for the subset of patients in each outcome dataset.

Work in progress 

## APPENDIX C ADDITIONAL EXPERIMENTS

### A. *Generic regularization*

Here we investigate the effect of architecture-agnostic regularization methods on forgetting. We investigate:

- dropout  $p \in 0.2, 0.4, 0.6, 0.8$
- L2 regularization
- SGD momentum  $\in 0.9, 0.8, 0.7, 0.6, 0.4, 0.2$
- training batch size  $n \in 16, 32, 64, 128$

Results:

[...]

### B. *Sequence length*

Here we evaluate the effect of sequence length on forgetting. We subsample the data stream at a more granular level of 1-hourly, 2-hourly, 4-hourly, 8-hourly, 12-hourly, and 24-hourly. We opted for sub sampling since truncating the datastreams may unfairly bias performance towards larger streams as most pertinent information (in terms of clinically observed changes to vital signs) to patient deterioration occurs closer to the deterioration event [44].

### C. *Replay buffer size*

Here we evaluate the rehearsal methods with an increasing storage buffer (from 10% of the training examples incrementing to full memory i.e. Cumulative strategy).

### D. *Curricula*

Here we evaluate the comparative performance of models given different curriculum orderings of their tasks. We consider random, correlated, reverse-correlated orderings for the (i) age, (ii) region, and (iii) time experiments.

### E. *Reduced feature set*

Here we evaluate the models on a reduced feature set consisting only of routinely recorded vital signs (as well as static demographic information).