



Pervasive AI

Tutorial of the 37th AAAI Conference on Artificial Intelligence
February 7, 2023 – Washington DC, USA

Davide Bacciu, Antonio Carta, Patrizio Dazzi, Claudio Gallicchio
University of Pisa, Italy



<http://pai.di.unipi.it/aaai-2023-tutorial-on-pervasive-ai/>

Continual Learning

Davide Bacciu, Antonio Carta, Patrizio Dazzi, Claudio Gallicchio
University of Pisa, Italy

Outline

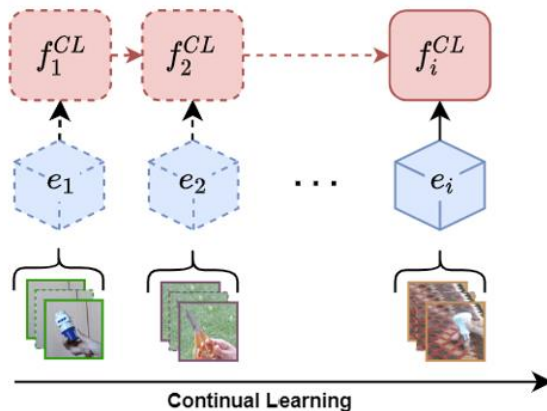
- Continual Learning
 - CL for Edge Devices
 - Distributed CL
-

Introduction to Continual Learning

Continual Learning

CL = Incremental Learning from a non-stationary stream

- + **environment information**: task labels, task boundaries, ...
- + **constraints**: computational/memory limits, privacy, ...
- + **metrics**: minimize forgetting, maximize transfer, ...



Suggested review:

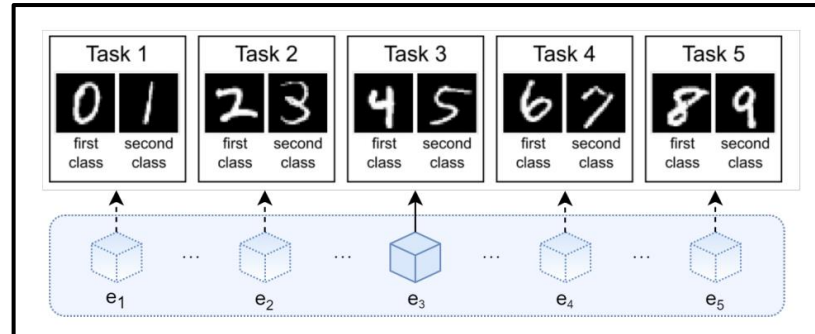
T. Lesort et al. "Continual Learning for Robotics: Definition, Framework, Learning Strategies, Opportunities and Challenges." *Information Fusion*.

<https://doi.org/10.1016/j.inffus.2019.12.004>.

Environment Information – Nomenclature

Different streams require different methods

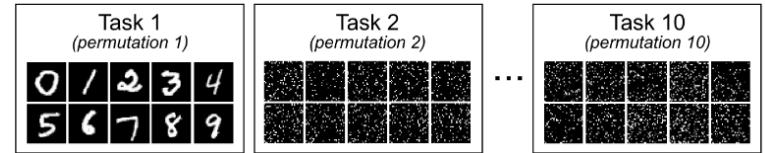
- **Stream:** A list of experiences, each providing a batch of data and some additional information (e.g. task labels)
- **Batch/Online:** How much data do we have in each experience?
- **Class/Domain-incremental:** Do we know the **type of shifts**?
- Do we know **when** the shifts happen?
- Do we have **task labels** at training/inference time?



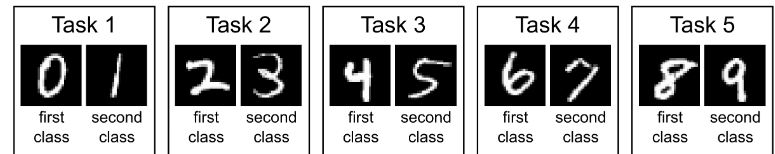
Three Common Scenarios

- **Domain-incremental:** each experience provides new instances for old classes. Old instances are never seen again (in the training stream).
- **Class-Incremental:** each experience provides new classes. Old classes are never revisited (in the training stream).
- **Task-Incremental:** each experience provides task labels
- **task-aware/task-agnostic** to highlight presence of task labels

Permuted MNIST



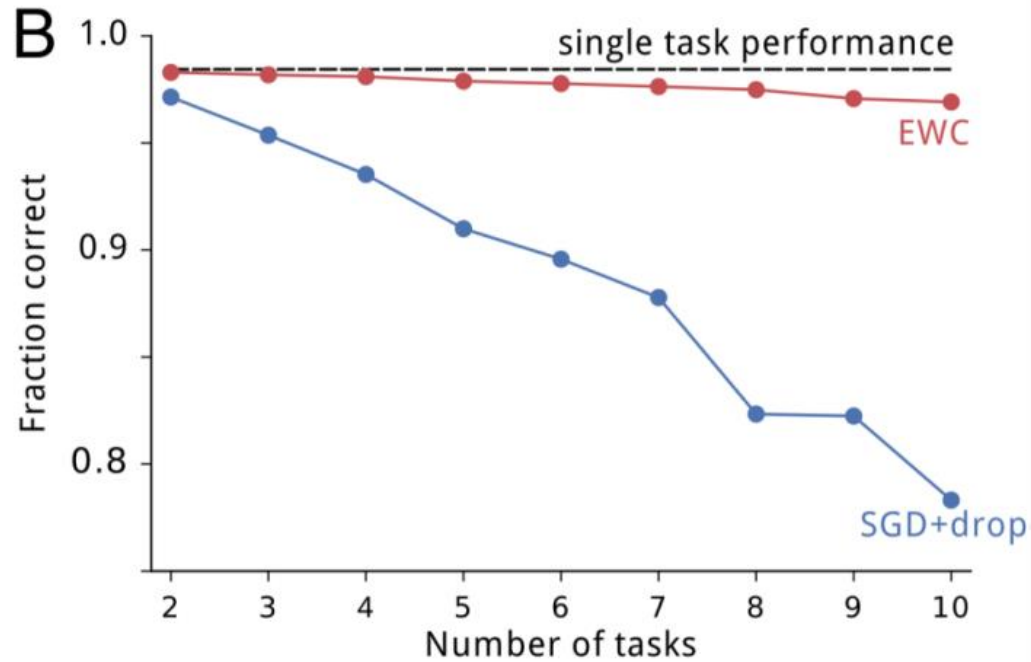
SplitMNIST



“Catastrophic Forgetting”

Deep neural networks completely and abruptly forget previously learned information upon learning new information.

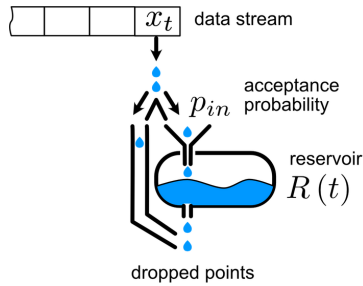
Catastrophic Forgetting



Methodologies

Replay

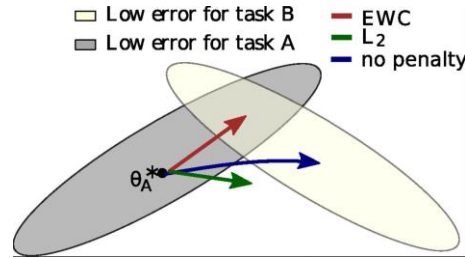
- Keep a buffer of old samples
- Rehearse old samples



Reservoir Sampling

Regularization

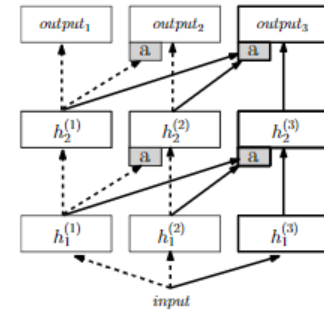
- Regularize the model to balance learning and forgetting



Elastic Weight Consolidation

Architectural

- Expand the model over time with new units/layers



Progressive Neural Networks

CL for Edge Devices

CL for Edge Devices

Applications

- Model personalization
- Wearable sensors
- Streams of video

Edge devices bring their own constraints

- **Privacy:** often we cannot save the user's data
- **Efficiency:** energy, memory, and computational power are limited resources
- **Latency:** often require online training and predictions

Continuous Object Recognition: CORE50

- **Continuous Object Recognition**

- 50 classes
- Short videos of object manipulation with different background
- Temporal coherence from videos

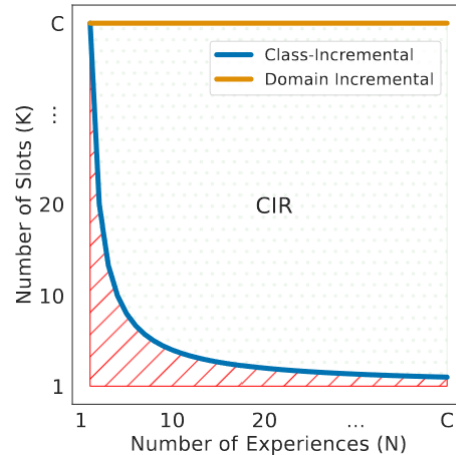
- Many scenarios: batch, online, with repetitions.

- **Advantages:** online, realistic streams



CL with Repetitions

- Real world problems have repetitions of concepts and imbalanced concepts
- Synthetic generator to simulate streams with repetitions and control the difficulty of the benchmark



CL with Repetitions – Results

- Repetitions help, even with naive finetuning
- Models change less and less over time
- Correcting the imbalance in the stream helps

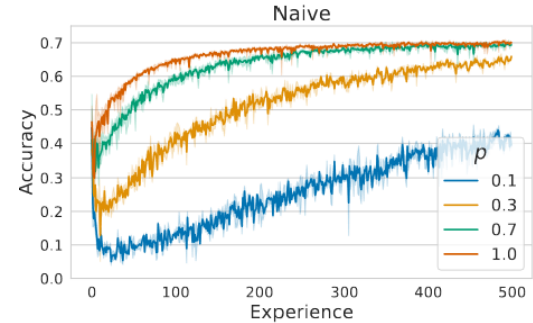


Figure 11: Accuracy of Infrequent Classes.

CL Methods for Edge Devices

- **Rehearsal-free** CL methods avoid storing past data to ensure privacy
- **Online CL** methods learn from small mini-batches (often with replay)
- CL methods can exploit **randomized networks** and **pretrained models**

Efficiency – Replay with Class-Balanced Reservoir Sampling

- **Rehearsal:** keeps a small buffer with the old data
- **Reservoir Sampling:** random sampling
- **Class-Balanced:** The buffer capacity is class-balanced
- **Advantages:** Efficient, suitable for Online CL

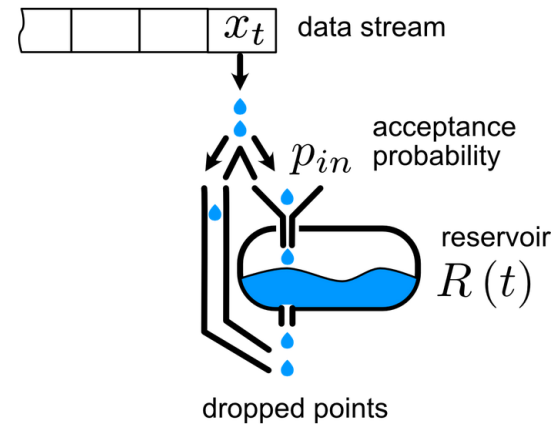


Image from <https://towardsdatascience.com/reservoir-sampling-for-efficient-stream-processing-97f47f85c11b>

Rehearsal-free CL

- **Hard problem:** if you can, always use replay
- **Pretraining** helps
- **Freezing** lower layers helps
- **Classifier bias:** in class-incremental learning, the classifier will be biased towards new classes

CL with Pretrained Models – Deep SLDA

- Fixed pretrained feature extractor
- Classifier trained via Streaming LDA
- No backpropagation: fast training

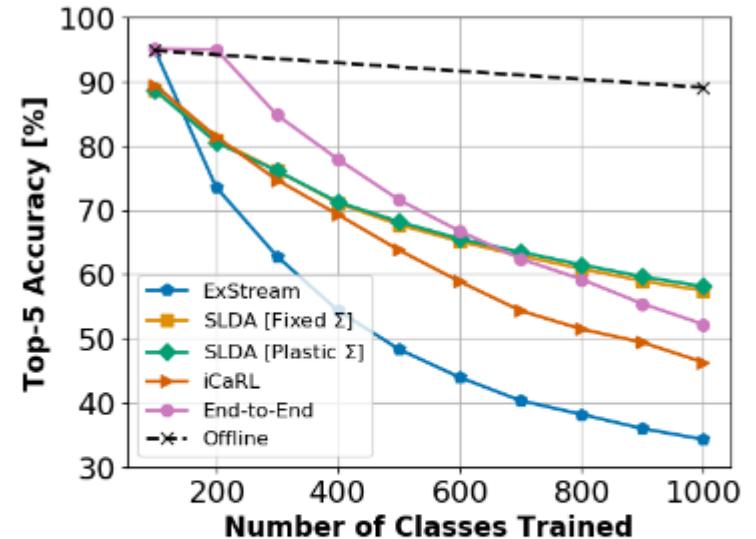


Figure 1: Learning curve for incremental ImageNet. Our Deep SLDA approach achieves the best final top-5 accuracy, while running over 100 times faster and using 1,000 times less memory than the iCaRL and End-to-End models.

CL with Pretrained Models – AR1

- Pretraining
- Architectural+Regularization Method
- Fix classifier bias with Copy-Weight-Reinit
- Regularization with Synaptic Intelligence

Algorithm 3 AR1

- 1: $cw = 0$
- 2: init $\bar{\Theta}$ random or from a pre-trained model (e.g. ImageNet)
- 3: $\Theta = 0$ (Θ are the optimal shared weights resulting from the last training, see Section 2.3)
- 4: $\hat{F} = 0$ (\hat{F} is the weight importance matrix, see Section 2.3).
- 5: for each training batch B_i :
- 6: expand output layer with s_i neurons for the new classes in B_i
- 7: $tw = 0$ (for all neurons in the output layer)
- 8: Train the model with SGD on the s_i classes of B_i by simultaneously:
- 9: learn tw with no regularization
- 10: learn $\bar{\Theta}$ subject to SI regularization according to \hat{F} and Θ
- 11: for each class j among the s_i classes in B_i :
- 12: $cw[j] = tw[j] - \text{avg}(tw)$
- 13: $\Theta = \bar{\Theta}$
- 14: Update \hat{F} according to trajectories computed on B_i (see eq. 10 and 11)
- 15: Test the model by using $\bar{\theta}$ and cw

CL for Time Series – CL with RNNs

- Sequence length increases forgetting
- Replay is still the best method
- Lack of a common benchmark for CL on time series

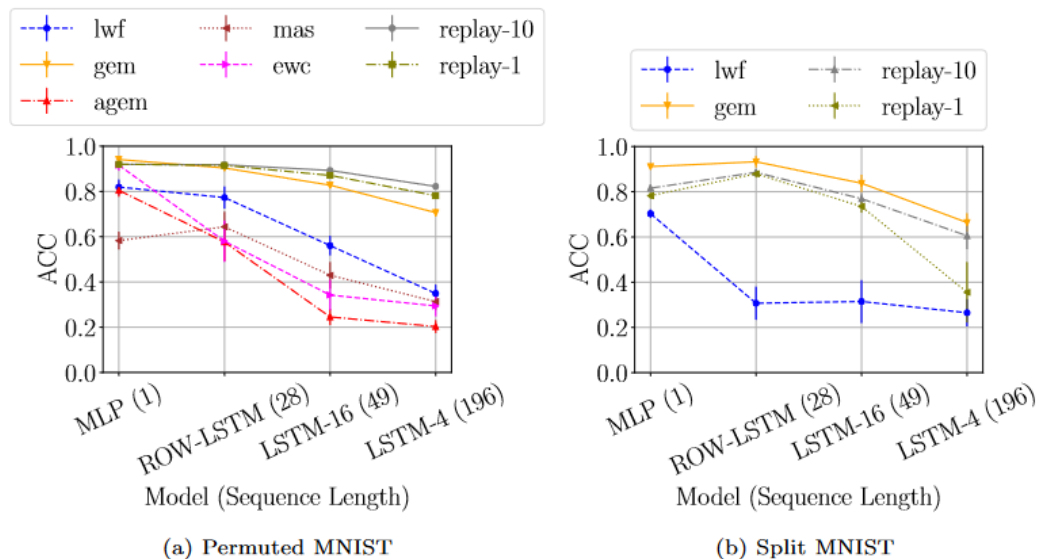


Figure 6: Average ACC on all steps for different sequence lengths and different CL strategies. Sequence length causes a decrease in performances among all strategies. Best viewed in color.

CL for Time Series – CL with ESNs

- Alternative to pretraining for RNNs
- Methods that train only the classifier can be used (SLDA)
- Efficient
- rehearsal-free
- Applicable for Online CL

	SMNIST	LSTM [†]	ESN		SSC	LSTM [†]	ESN
EWC		0.21 \pm 0.02	0.20 \pm 0.00	EWC		0.10 \pm 0.00	0.09 \pm 0.02
LWF		0.31 \pm 0.07	0.47 \pm 0.07	LWF		0.12 \pm 0.01	0.12 \pm 0.02
REPLAY		0.85\pm0.03	0.74 \pm 0.03	REPLAY		0.74\pm0.07	0.36 \pm 0.07
SLDA		—	0.88\pm0.01	SLDA		—	0.57\pm0.03
NAIVE		0.20 \pm 0.00	0.20 \pm 0.00	NAIVE		0.10 \pm 0.00	0.10 \pm 0.00
JOINT		0.97 \pm 0.00	0.97 \pm 0.01	JOINT		0.89 \pm 0.02	0.91 \pm 0.02

Table 1: Mean ACC and standard deviation over 5 runs on SMNIST and SSC benchmarks. SLDA is applied only to ESN since it assumes a fixed feature extractor. SMNIST contains 5 experiences, while SSC contains 10 experiences. [†] results are taken from [3], except for replay which has been recomputed to guarantee the use of the same replay policy (200 patterns in memory).

Towards Distributed Continual Learning

CL for Pervasive Learning

Personalized Models

- Properties:
 - Trained on edge devices
 - Limited computational resources
 - Small and private datasets
- Objectives:
 - Personalization
 - Low-latency
 - Privacy

How do we learn personalized models? Can we improve their performance via transfer between different models?

Large Pretrained Models

- Properties:
 - Trained on a dedicated server/cluster
 - High computational resources
 - Large and diverse datasets
- Objectives:
 - Learning general knowledge
 - Provide forward transfer for downstream tasks

How do we learn large pretrained models with continual learning? How do we use them in CL?

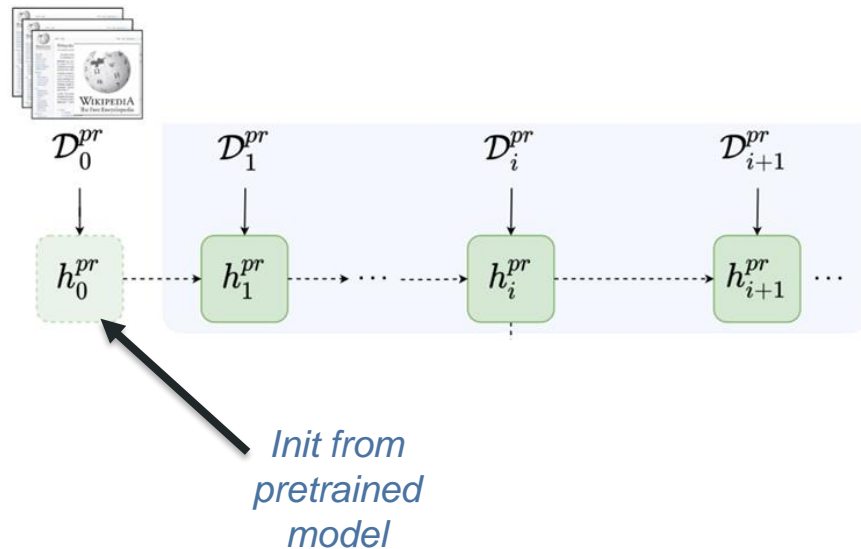
Pretraining in CL

Pretrained Models in the CL literature: used as a base initialization.

PROBLEM: pretraining helps only before training. Once we start learning, we cannot use pretrained models anymore.

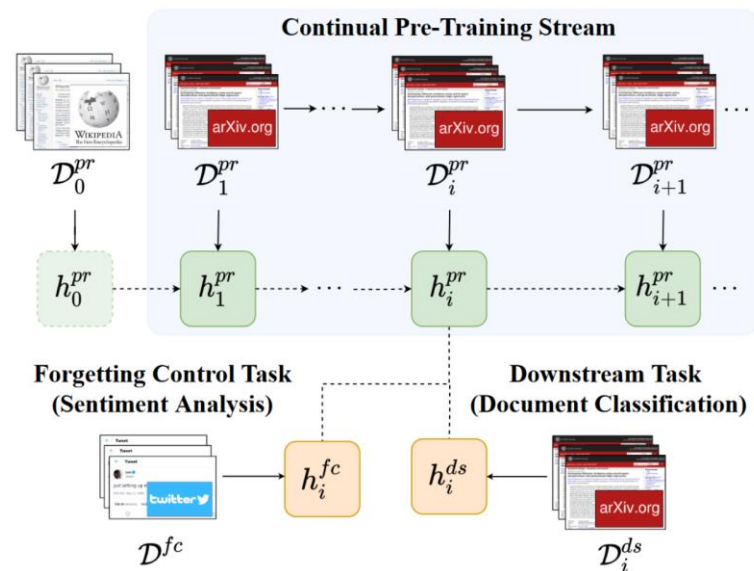
Proposed solutions:

- *Federated CL*: centralized approach, no difference between large pretrained model and personalized models
- *Continual Pretraining*: CL of the pretrained model + Continual finetuning on downstream tasks
- *Ex-Model CL*: Devices send and receive “expert model” messages



Continual Pretraining

- **Continual Pretraining:** the large pretrained model is periodically updated via CL
- **Continual Finetuning:** whenever a new pretrained model becomes available, finetune the downstream model
 - Can be trained from scratch efficiently with linear probing.



Continual Pretraining: Results

Evaluation on the Forgetting Control Task

Table 2: Accuracy on the entire dataset of `sentiment analysis` with RoBERTa model. Continual pre-training has been performed sequentially over each experience of `scientific abstracts`. Base refers to the model pre-trained on Wikipedia, while NT refers to the model with vocabulary expansion.

RoBERTa		Accuracy					1-epoch Accuracy				
Base	93.40					92.40					
Exp.	e1	e2	e3	e4	e5	e1	e2	e3	e4	e5	
Pretr	93.40	93.15	93.35	93.20	92.90	92.40	91.80	92.30	91.85	92.20	
Pretr. NT	93.75	93.70	93.75	93.60	94.10	91.75	91.15	92.00	92.30	92.45	

← fast adaptation

Forgetting is limited even with naive finetuning.
Dynamic vocabulary expansion (NT) slightly improves the performance.

Federated Continual Learning

- Combines Federated and Continual Learning
- Example: **FedWelt**
 - **Client's** model parameters: base + task-specific + weighted sum of other tasks
 - **Server** keeps all the parameters
 - **Forward transfer**: Clients learn which parameters from the other tasks are useful

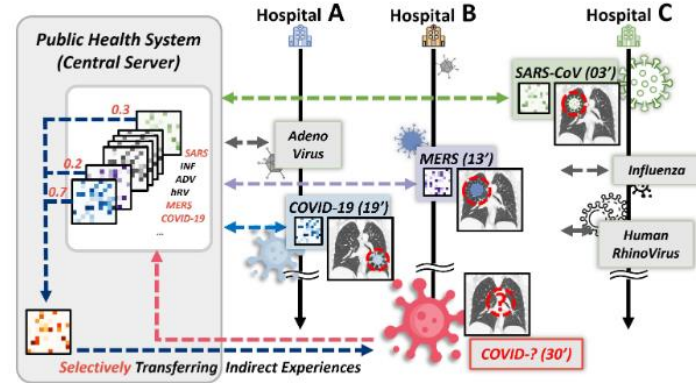
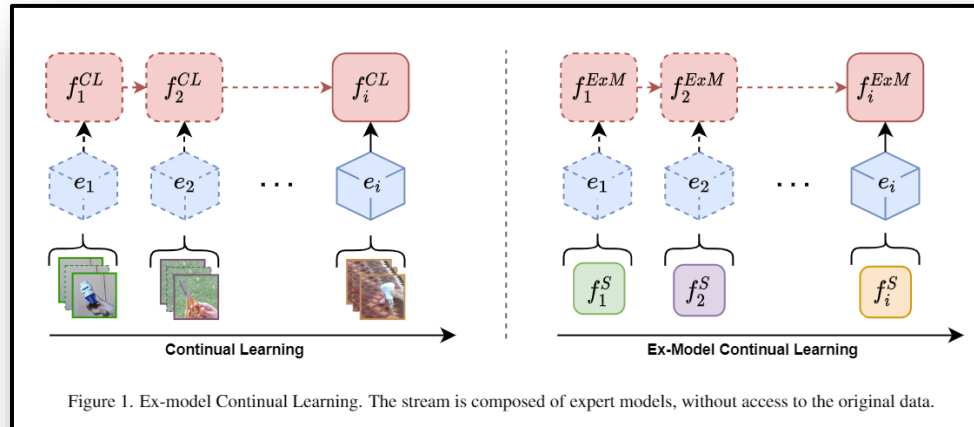


Figure 1. Concept. A continual learner at a hospital which learns on sequence of disease prediction tasks may want to utilize relevant task parameters from other hospitals. FCL allows such inter-client knowledge transfer via the communication of task-decomposed parameters.

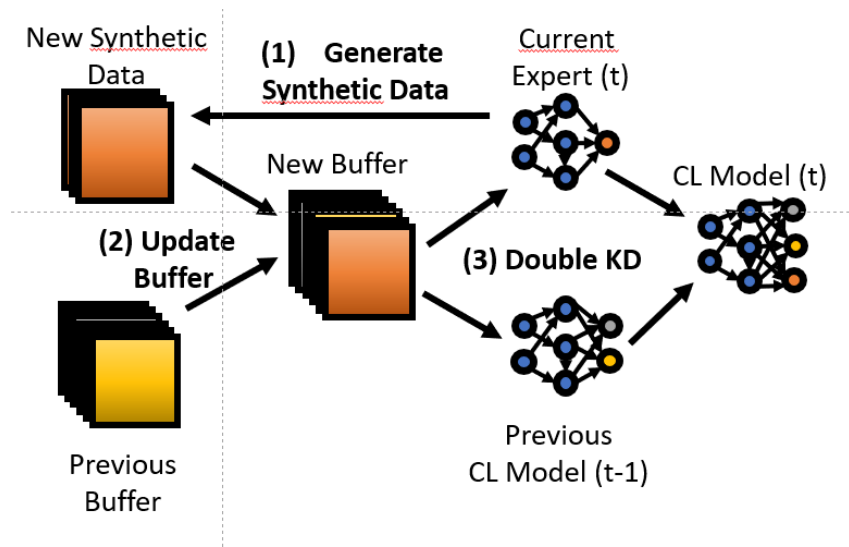
Ex-Model Continual Learning

- **Continual Learning:** one batch of data at each experience
- **Ex-Model CL:** one pretrained model at each experience
- **Consolidation:** learn to combine two different models



Ex-Model Continual Learning: Double Distillation

- **Double Distillation:** learn to consolidate the new expert with the previous CL model
- **Data-free:** use out-of-distribution data or synthetic samples



Avalanche

- CL library built on top of Pytorch
- Currently the most extensive collection of CL benchmarks and algorithms
- Used by the CL community for research, new benchmarks, challenges and courses
- **Tutorial on Avalanche tomorrow @ 9:00AM at the Continual-Causality Bridge!**

Website: avalanche.continualai.org/

CL-baselines:

<https://github.com/continualAI/continual-learning-baselines/>

Avalanche-demo:

<https://github.com/AntonioCarta/avalanche-demo>



powered by



ContinualAI

```
○○○
1 # model
2 model = SimpleMLP(num_classes=10)
3
4 # CL Benchmark Creation
5 perm_mnist = PermutedMNIST(n_experiences=3)
6 train_stream = perm_mnist.train_stream
7 test_stream = perm_mnist.test_stream
8
9 # Prepare for training & testing
10 optimizer = SGD(model.parameters(), lr=0.001, momentum=0.9)
11 criterion = CrossEntropyLoss()
12
13 # Continual learning strategy
14 cl_strategy = Naive(
15     model, optimizer, criterion, train_mb_size=32, train_epochs=2,
16     eval_mb_size=32, device=device)
17
18 # train and test loop over the stream of experiences
19 results = []
20 for train_exp in train_stream:
21     cl_strategy.train(train_exp)
22     results.append(cl_strategy.eval(test_stream))
```

Conclusion

Continual Learning provides efficient solutions to learn in constrained environments

- **Efficient solutions:** Replay with CBRS
- **Online CL:** fix classifier bias + use of pretrained models + freezing/regularization
- promising avenues for **future research:**
 - Interaction between CL and pretraining
 - Federated CL
 - Ex-Model CL