



Continual Learning

Scenarios, evaluation, and metrics

Antonio Carta

antonio.carta@unipi.it

- Module Structure
- Definition of Deep Continual Learning
- Catastrophic Forgetting

- Deep Continual Learning and Catastrophic Forgetting (today)
- Methodologies – Baselines and Replay
- Methodologies – EWC and Continual Optimization
- Methodologies – Regularization
- Methodologies – Architectural methods and incremental classifiers

We will also have some labs with practical examples

CL Definition

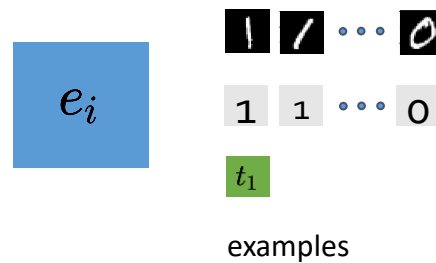
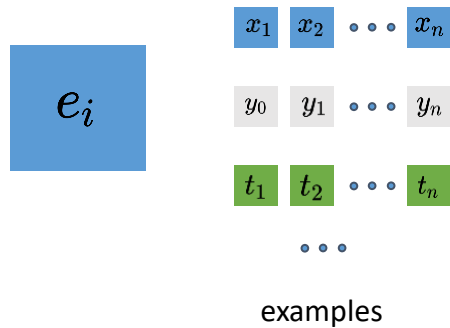
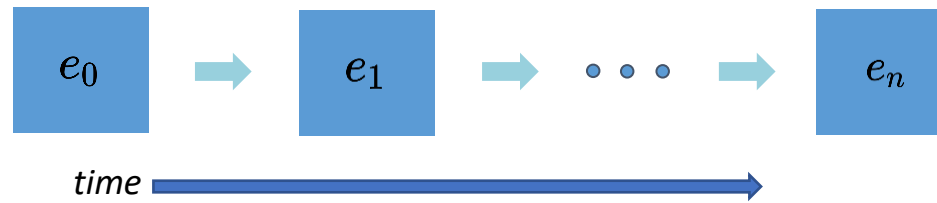
Learning from a non-stationary stream with Deep Neural Networks

- How to learn incrementally with DNN
 - Efficiently (fast adaptation)
 - Without forgetting old tasks (stability)
 - Possibly with some constraints such as privacy of old data or low-data scenarios
- We revisit the previous problems in a more challenging setting:
 - Multi-Task Learning -> Task-Incremental Learning
 - Online with prequential evaluation -> evaluation on the entire stream (i.e. never forget)

- **Lifelong timescales:** learning incrementally on long (lifelong) timescales is a fundamental property of intelligent systems.
- **Limited memory:** No (or limited) access to previously encountered data.
- **Limited computation:** Computational cost is limited over time and it doesn't grow too much.
- **Incremental learning:** the model continually improves over time.

Continual Learning – Data Streams

Stream of Experiences



- Data stream $S = e_1, \dots, e_n$
- Each experience e_i provides samples D_i
- Each sample is a tuple $\langle x, y, t \rangle$ of input, target, and (optionally) task label
- We assume to have parallel train/validation/test streams
- A continual learning algorithm is a method that, given the previous model and the new data, provides a new model
 - The model may have some internal state (e.g. a replay buffer or some proxy for old losses)

Continual Learning - Definition



In continual learning (CL) data arrives in a streaming fashion as a (possibly infinite) sequence of learning experiences $S = e_1, \dots, e_n$. For a supervised classification problem, each experience e_i consists of a batch of samples \mathcal{D}^i , where each sample is a tuple $\langle x_k^i, y_k^i \rangle$ of input and target, respectively, and the labels y_k^i are from the set \mathcal{Y}^i , which is a subset of the entire universe of classes \mathcal{Y} . Usually \mathcal{D}^i is split into a separate train set \mathcal{D}_{train}^i and test set \mathcal{D}_{test}^i .

A continual learning algorithm \mathcal{A}^{CL} is a function with the following signature:

$$\mathcal{A}^{CL} : \langle f_{i-1}^{CL}, \mathcal{D}_{train}^i, \mathcal{M}_{i-1}, t_i \rangle \rightarrow \langle f_i^{CL}, \mathcal{M}_i \rangle \quad (1)$$

where f_i^{CL} is the model learned after training on experience e_i , \mathcal{M}_i a buffer of past knowledge, such as previous samples or activations, stored from the previous experiences and usually of fixed size. The term t_i is a task label which may be used to identify the correct data distribution.

The objective of a CL algorithm is to minimize the loss \mathcal{L}_S over the entire stream of data S :

$$\mathcal{L}_S(f_n^{CL}, n) = \frac{1}{\sum_{i=1}^n |\mathcal{D}_{test}^i|} \sum_{i=1}^n \mathcal{L}_{exp}(f_n^{CL}, \mathcal{D}_{test}^i) \quad (2)$$

$$\mathcal{L}_{exp}(f_n^{CL}, \mathcal{D}_{test}^i) = \sum_{j=1}^{|\mathcal{D}_{test}^i|} \mathcal{L}(f_n^{CL}(\mathbf{x}_j^i), y_j^i), \quad (3)$$

where the loss $\mathcal{L}(f_n^{CL}(\mathbf{x}), y)$ is computed on a single sample $\langle \mathbf{x}, y \rangle$, such as cross-entropy in classification problems.

- The goal is to minimize the average loss on the test stream
- It's different from online learning, where we only evaluate on future data. Here, we want to remember the past

$$\mathcal{L}_S(f_n^{CL}, n) = \frac{1}{\sum_{i=1}^n |\mathcal{D}_{test}^i|} \sum_{i=1}^n \mathcal{L}_{exp}(f_n^{CL}, \mathcal{D}_{test}^i)$$

$$\mathcal{L}_{exp}(f_n^{CL}, \mathcal{D}_{test}^i) = \sum_{j=1}^{|\mathcal{D}_{test}^i|} \mathcal{L}(f_n^{CL}(\mathbf{x}_j^i), y_j^i),$$

Deep CL

- (possibly) Large experiences
- Virtual drift
- Domains: Vision, NLP, speech
- Evaluation: average accuracy on the full stream

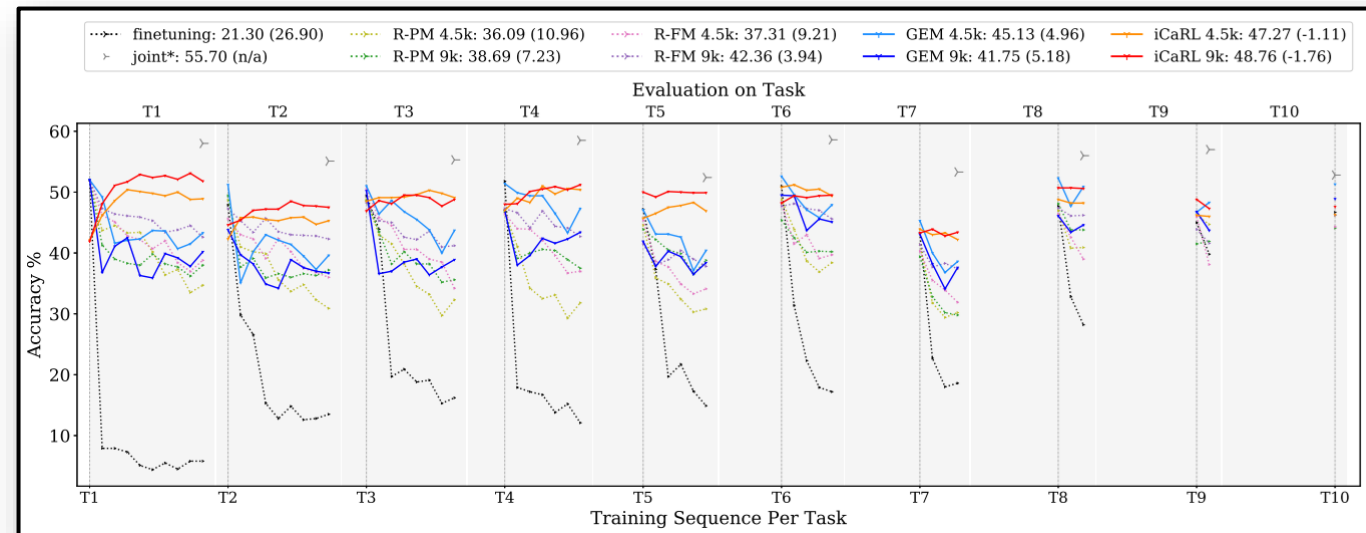
Online ML

- One sample at a time
- Real drift
- Domains: time series data
- Evaluation: prequential accuracy

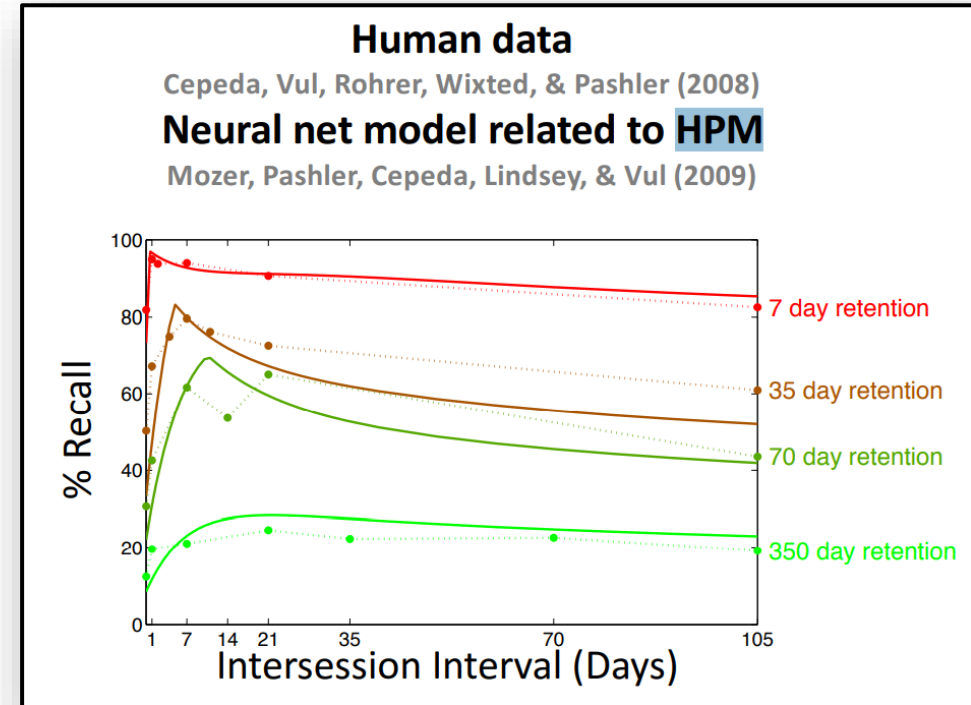
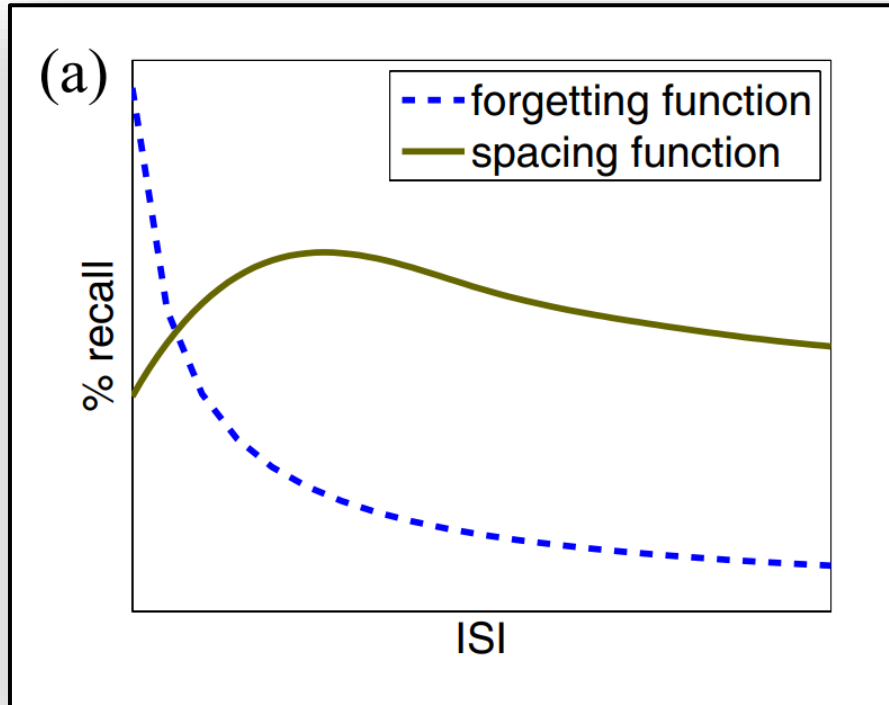
Challenges

Catastrophic Forgetting

Catastrophic interference, a.k.a. **catastrophic forgetting**, is the tendency of deep neural networks to completely and abruptly forget previous tasks/domains when learning new ones



Forgetting in Humans



The spacing function (solid line) depicts recall at test following two study sessions separated by a given intersession interval (ISI); the forgetting function (dashed line) depicts recall as a function of the lag between study and test.

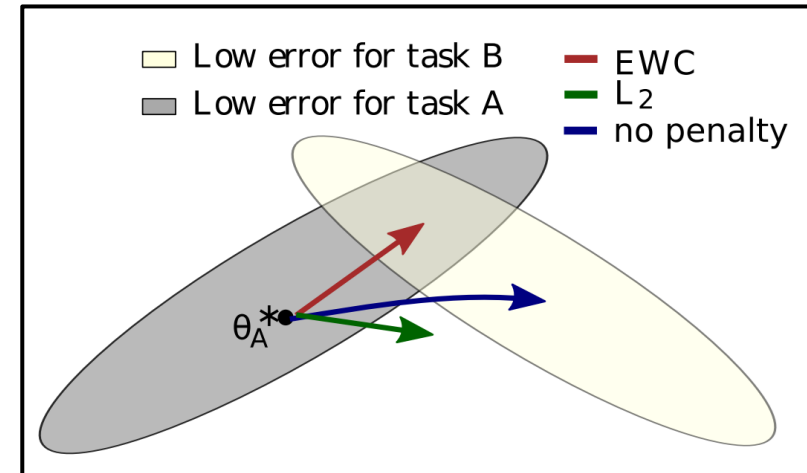
Forgetting – EWC



How can we approximate the loss on the previous task/domain without access to the data?

L_S is made of two parts:

- The loss on current data, which we can compute
- The loss on previous data, which we CANNOT compute because we don't have the data



The objective of a CL algorithm is to minimize the loss \mathcal{L}_S over the entire stream of data S :

$$\mathcal{L}_S(f_n^{CL}, n) = \frac{1}{\sum_{i=1}^n |\mathcal{D}_{test}^i|} \sum_{i=1}^n \mathcal{L}_{exp}(f_n^{CL}, \mathcal{D}_{test}^i) \quad (2)$$

$$\mathcal{L}_{exp}(f_n^{CL}, \mathcal{D}_{test}^i) = \sum_{j=1}^{|\mathcal{D}_{test}^i|} \mathcal{L}(f_n^{CL}(x_j^i), y_j^i), \quad (3)$$

where the loss $\mathcal{L}(f_n^{CL}(x), y)$ is computed on a single sample $\langle x, y \rangle$, such as cross-entropy in classification problems.

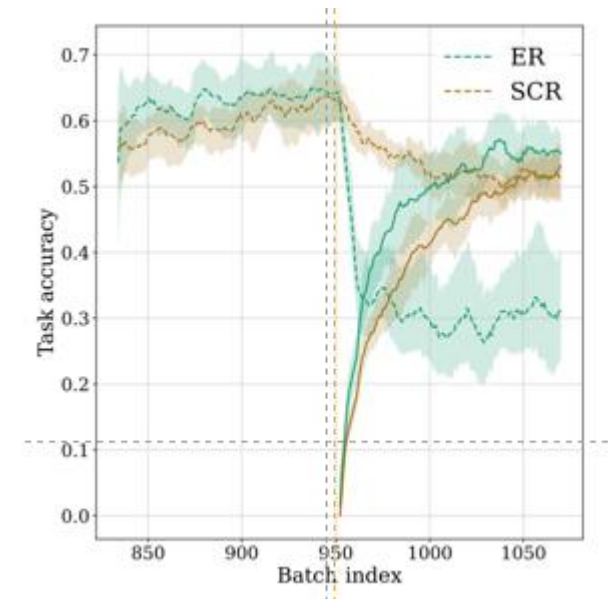
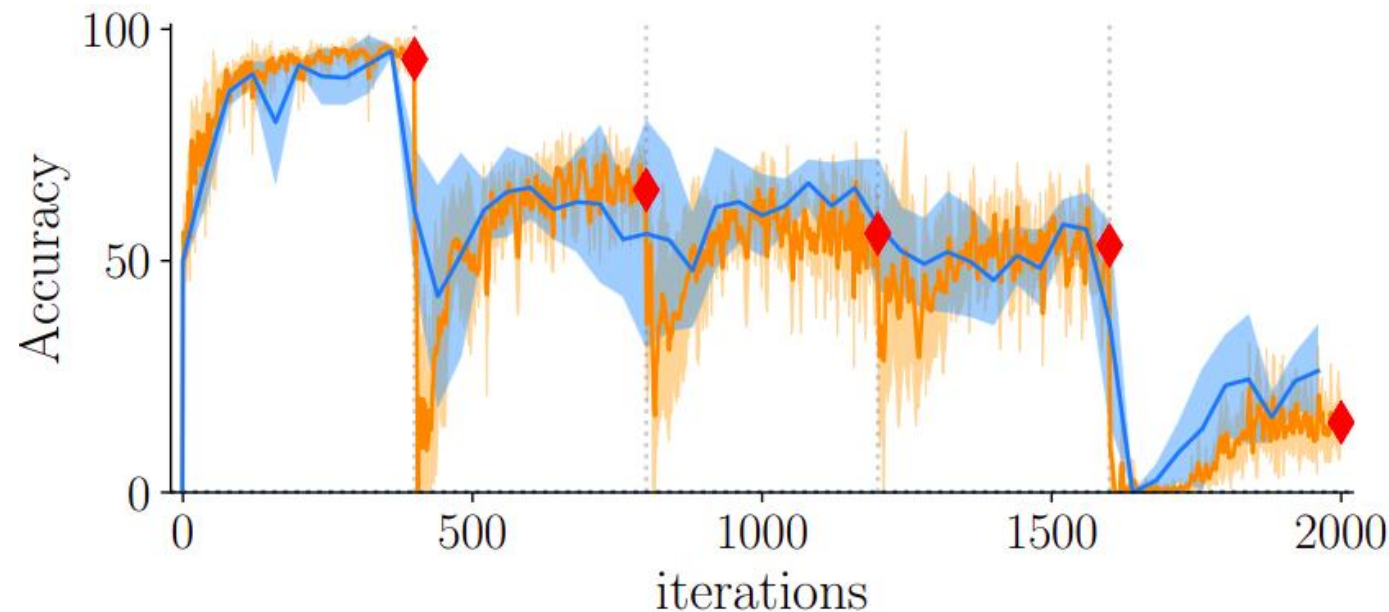
Stability-Plasticity Dilemma



- **Stability:** the ability to avoid forgetting of old task/domains
- **Plasticity:** the ability to learn new tasks
- In general, there is a tradeoff between the two properties
 - I can freeze the network to prevent forgetting
 - I can do a naive finetuning (or even randomly initialize) to have optimal plasticity

Training instability

- Observe the behavior of the accuracy during training (curve from one diamond to the next)
- CL methods forget and re-learn old experiences during training
- This phenomenon is masked with the typical metrics measured only at boundaries (red diamonds)



[1] Mathias Delange et. al, *Continual Evaluation for Lifelong Learning: Identifying the stability gap*, ICLR 2023

[2] Lucas Caccia et. al, *New Insights on Reducing Abrupt Representation Change in Online Continual Learning*, ICLR 2022

[3] A. Soutif et al. "A Comprehensive Empirical Evaluation on Online Continual Learning" 2023

CL - Nomenclature

CL Scenarios

- CL assumptions
- Types of drifts
- Nomenclature
- Benchmarks examples

Evaluation

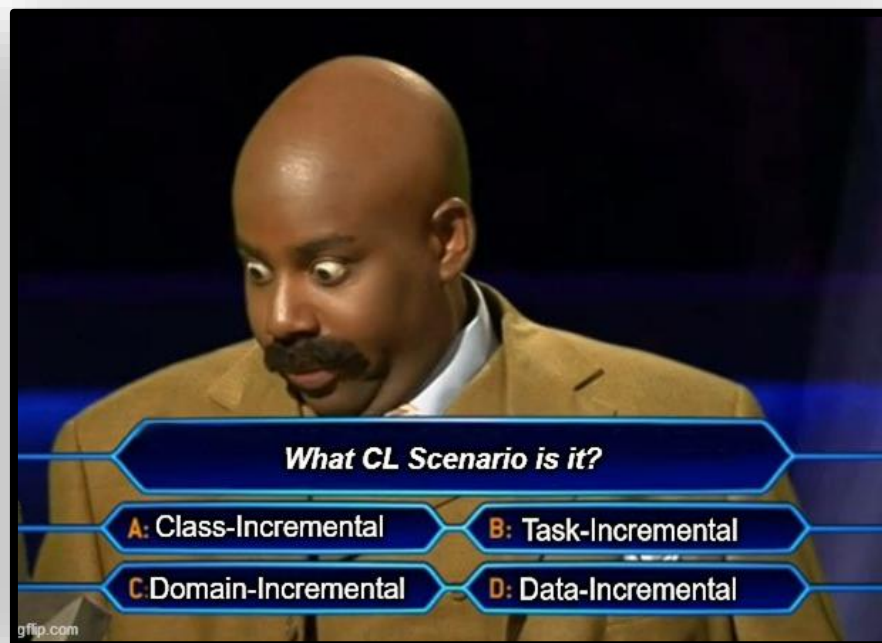
- CL eval vs prequential
- Hyperparameter selection

Metrics

- What and when to monitor
- Accuracy
- Forgetting
- Backward/forward transfer
- Computational performance

New Instances
& Classes

Task-
Agnostic



Task-
Free

Single-
Incremental-Task

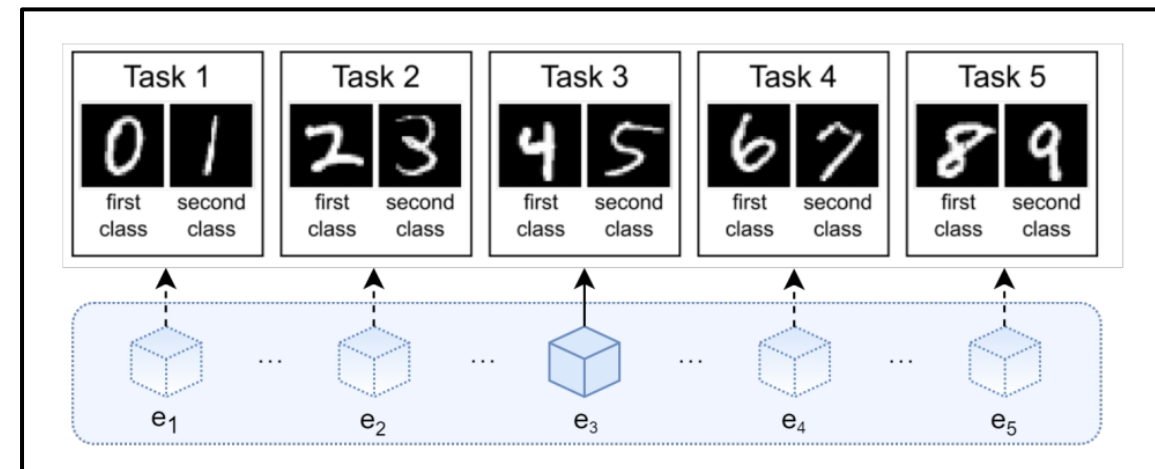
Multi-Task

What is a CL Scenario?

A common nomenclature that:

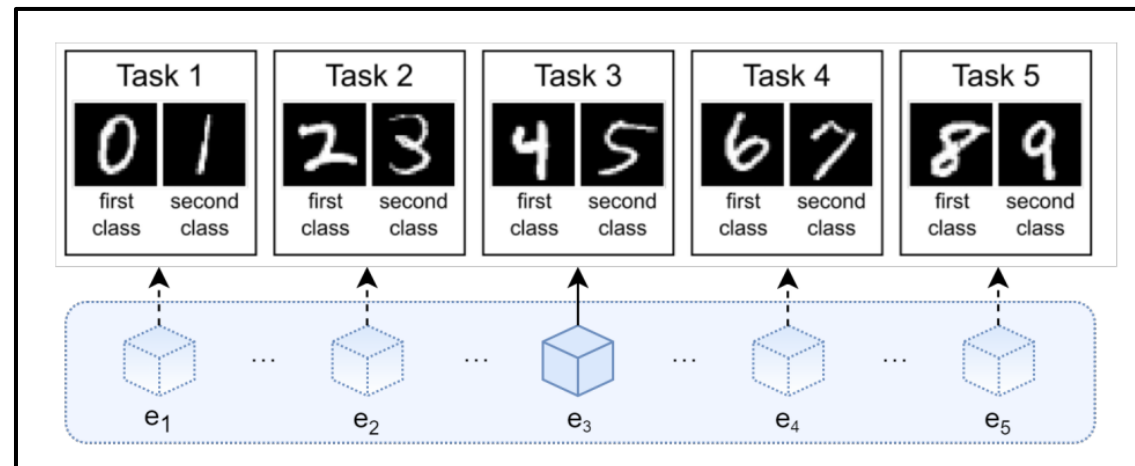
- A set of **CL metrics** that we want to optimize
- A set of constraints that the **learning algorithm** must satisfy
- A restricted form of access to the data through a sequential **data stream**

Given the type of scenario and its constraints we can identify a proper strategy to learn it.



Scenarios Nomenclature: what we need

- **How much data** do we have in each experience?
- Do we know the **type of shifts**?
- Do we know **when** the shifts happen?
- Do we have **task labels** at training/inference time?



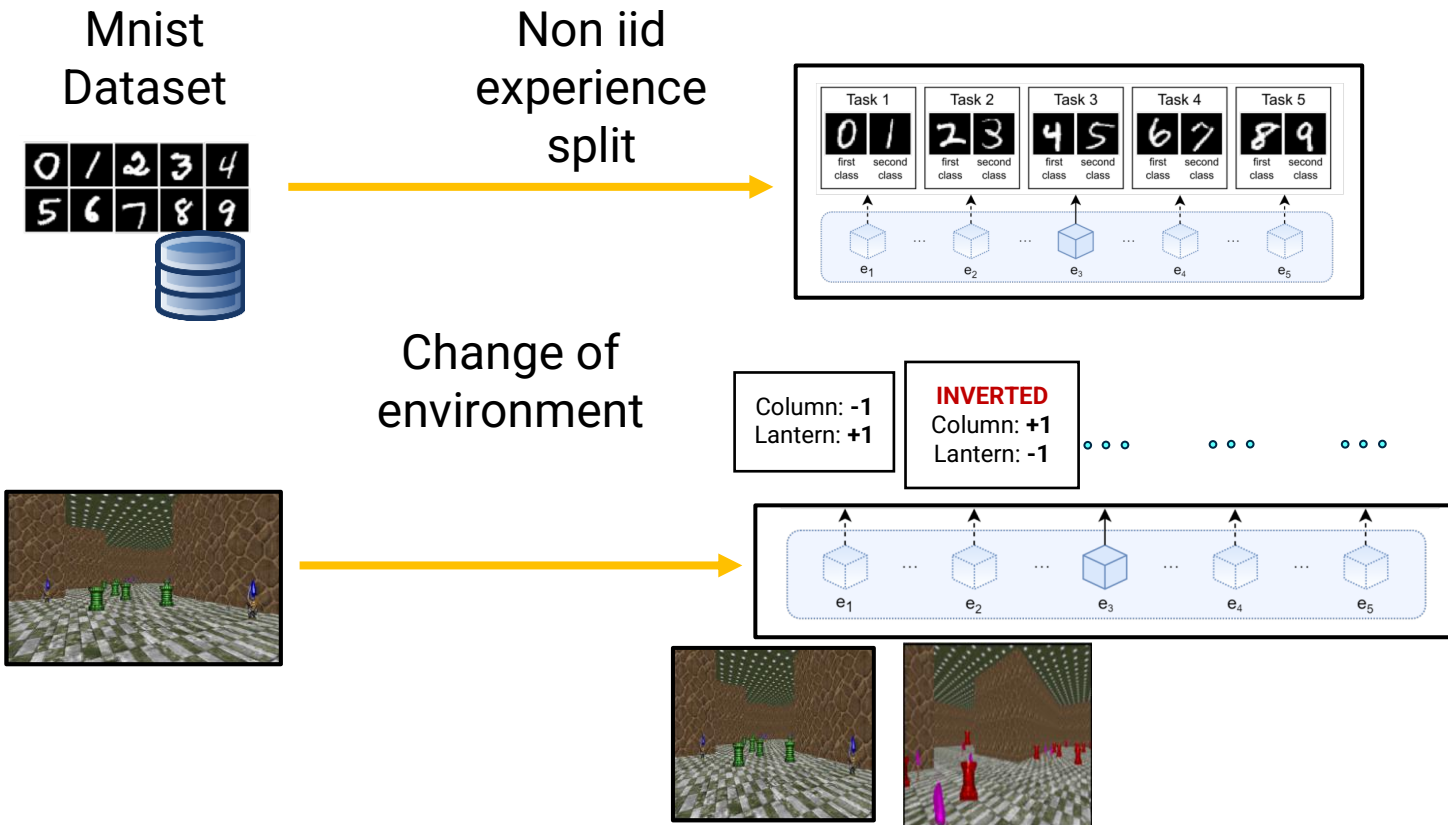
- **Real shift:**

- We have seen it in OML
- The world changes, either abruptly (covid lockdown) or continuously (weather, financial markets). $p(x, y)$ is changed.
- In some application you care only about the present and future and can forget the past

- **Virtual Shift**

- Common CL assumption
- The «world» is fixed. Shifts are «virtual» and due to sample selection bias
- The data changes because the prior $p(x)$ is changing. $p(y|x)$ is fixed.
- You don't want to forget anything because you may encounter the old data again in the future

Real vs Virtual Shift



Non-stationarity Assumptions:

- **Real Shift:** the world is changing (more studied in OML).
- **Virtual Shift:** the world is fixed but there is a sample selection bias.

We will often assume virtual shifts.

REMINDER: Dataset Shift Nomenclature



Dataset Shift: $p_{tra}(x, y) \neq p_{tst}(x, y)$ Informally: any change in the distribution is a shift

Covariate shift: happen in $X \rightarrow Y$ problems when

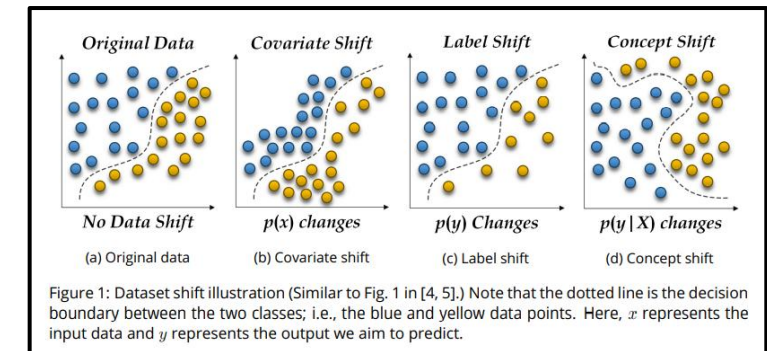
- $p_{tra}(y|x) = p_{tst}(y|x)$ and $p_{tra}(x) \neq p_{tst}(x)$
- informally: the input distribution changes, the input->output relationship does not

Prior probability shift: happen in $Y \rightarrow X$ problems when

- $p_{tra}(x|y) = p_{tst}(x|y)$ and $p_{tra}(y) \neq p_{tst}(y)$
- Informally: output->input relationship is the same but the probability of each class is changed

Concept shift:

- $p_{tra}(y|x) \neq p_{tst}(y|x)$ and $p_{tra}(x) = p_{tst}(x)$ in $X \rightarrow Y$ problems.
- $p_{tra}(x|y) \neq p_{tst}(x|y)$ and $p_{tra}(y) = p_{tst}(y)$ in $Y \rightarrow X$ problems.
- Informally: the «concept» (i.e. the class)



Common Assumptions

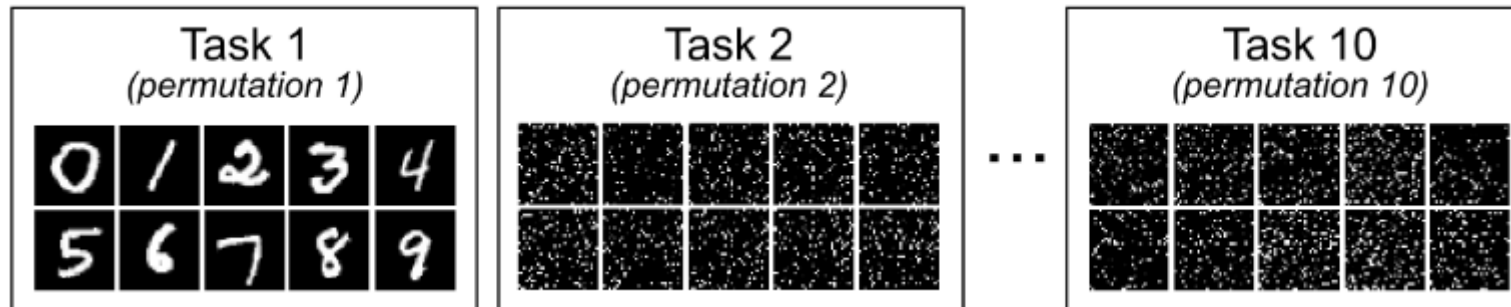


- **Shift is only virtual:** we do not want to forget, we need to accumulate knowledge.
- **No labeling errors/conflicting information:** targets are always correct (but possibly noisy).
- **Unbounded time:** No hard latency requirements. We may have computational constraints.
- **Data in each experience can be freely processed:** you can shuffle them, process them multiple times, etc. like you would do during offline training.

Common Types of Shifts

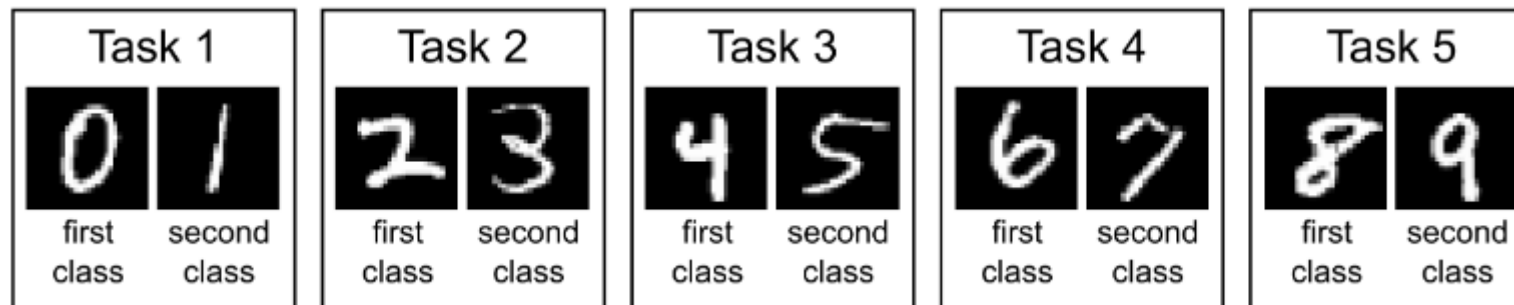
- **Domain-incremental:** each experience provides new instances for old classes. Old instances are never seen again (in the training stream).

Permuted
MNIST



- **Class-incremental:** each experience provides new classes. Usually, old classes are never revisited (in the training stream).

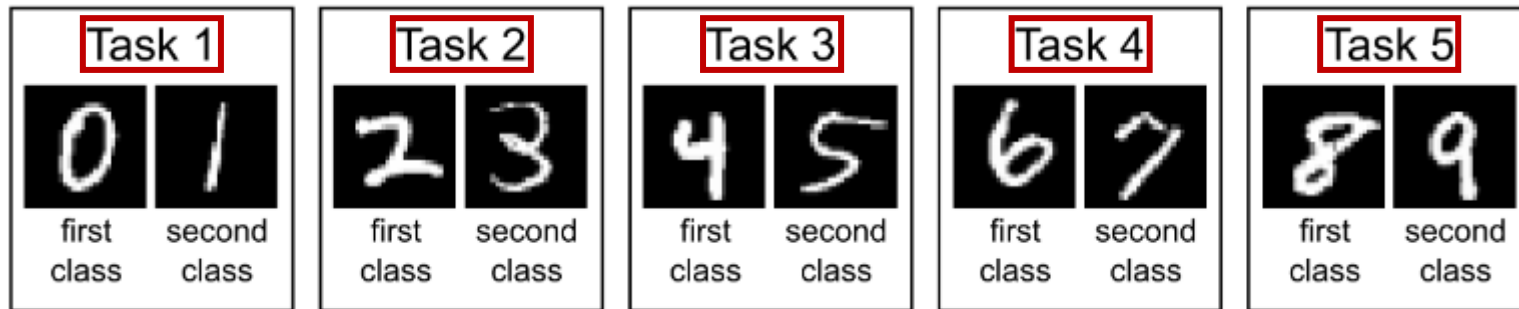
Split
MNIST



Presence of Task Labels



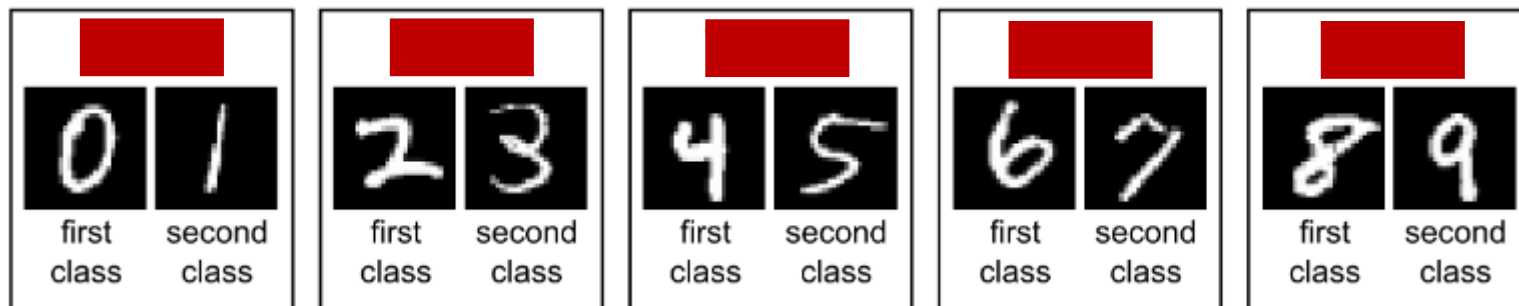
- **Task-aware:** task labels are available during training and inference



NOTE: some references use the term «task» to denote experiences even in task-agnostic scenarios

- **Task-agnostic:** task labels are not available

Task labels can change the output space (single vs multi head)

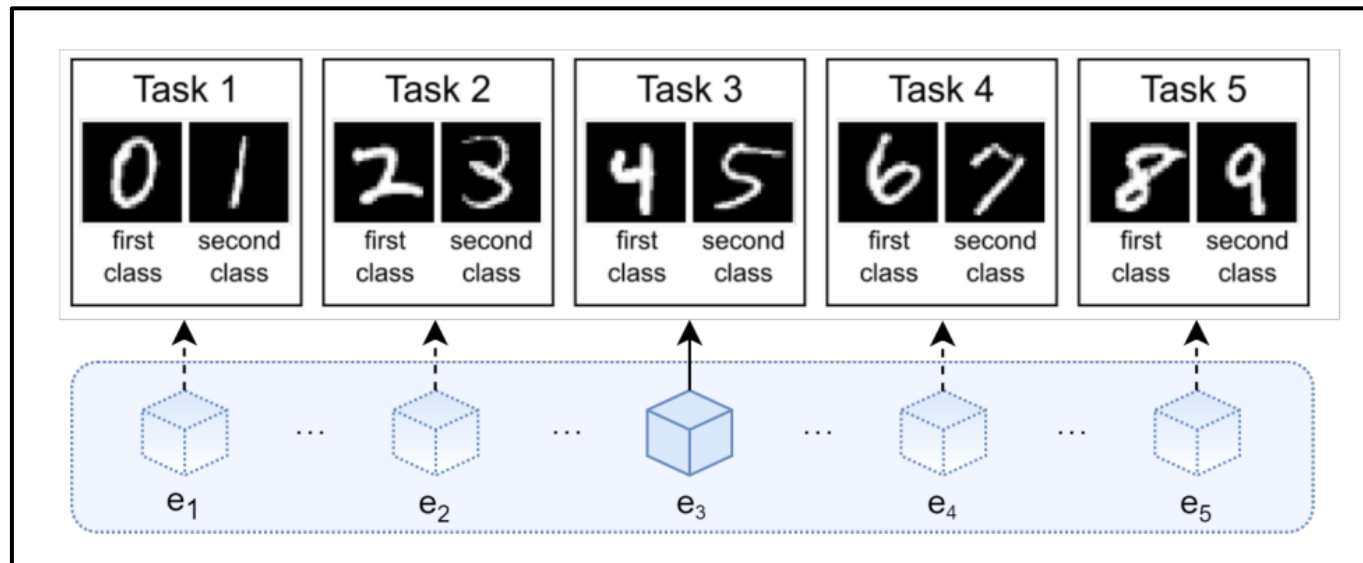


- Task labels simplify the problem
 - We can use multi-task models that take task labels as an explicit argument
 - Modularity also helps to prevent interference (it may limit forward transfer)
 - Output space is smaller:
 - 100 classes divided in 10 tasks -> 10-way classification
 - 100 classes in a single task -> 100-way classification
- Notice: the term task in CL is a bit overloaded
 - Sometimes, experiences are called tasks even when there are no explicit labels or other mechanisms to disambiguate different tasks
 - Often, tasks are actually domains (same problem, different $p(x)$)

Batch vs Online/Streaming CL

How much data for each experience?

- **Online CL (OCL) / Streaming CL:** Single example/small minibatch
- **Batch CL:** Large batch, no constraint on the size of the experience








- **Task-free:** the model doesn't know when the shift happens (as in OML)
 - Notice: we don't have task labels AND we don't know WHEN the shift happens
- No common term for the «shift-aware» version
- In a batch scenario the typical assumption is that each experience is the result of a distribution shift
- In OCL knowledge of task boundaries is more useful (because the stream is much longer) but all the methods assume that they don't have access to it (more realistic)

Sharp vs Blurry Shifts

- **Sharp Shifts:** drift happen abruptly
 - **Blurry/Gradual Shifts:** drift happen slowly
- Most CL methods deal with sharp drifts

RotatedMNIST:

				
Task 0	Task 1	Task 2	Task 3	Task 4

Example of
gradual shift:
Rotated MNIST

Remember the
assumption about «no
conflicting
information»? We
may want to remove 6
or 9 here

Nomenclature for Common Scenarios

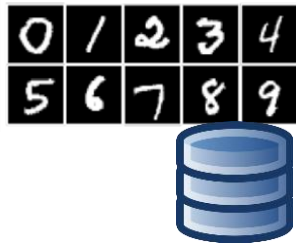


- **Availability of Task/Distribution Labels:** during training and/or testing
- **Task/Shift Boundaries:** during training and/or testing
- **Experience Content:** examples of [same|new] classes
- **Output Space:** [Shared/Separate]
- NOT an exhaustive classification

Name	Task Labels	Boundaries	Classes	Output
Class-Incremental	No	Yes	New	Shared
Task-Incremental	Yes	Yes	New	Separate
Domain-Incremental	No	Yes	Same	Shared
(Online) Task-Free	No	No	Any	Shared

Dataset, Scenarios, Benchmarks

MNIST Dataset



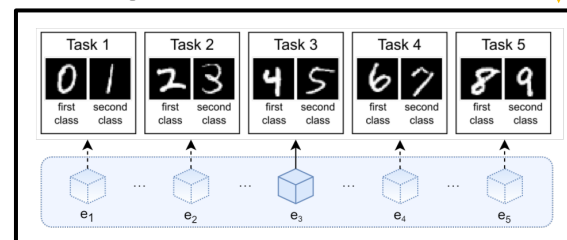
Class-Incremental Learning Scenario

Settings:

1. Each e contains only examples of new classes never seen before (clear boundaries)
2. No t available during train or test.
3. Shared output space



Split MNIST Benchmark



Benchmark Instances

Exact specific sequences and composition of e . For example:

$$S1 = \{\text{class } 0 \ \& \ 1\}, \{\text{class } 2 \ \& \ 3\}$$

$$S2 = \{\text{class } 5 \ \& \ 7\}, \{\text{class } 0 \ \& \ 8\}$$

represents two possible **Benchmark Instances** of *Split MNIST*.

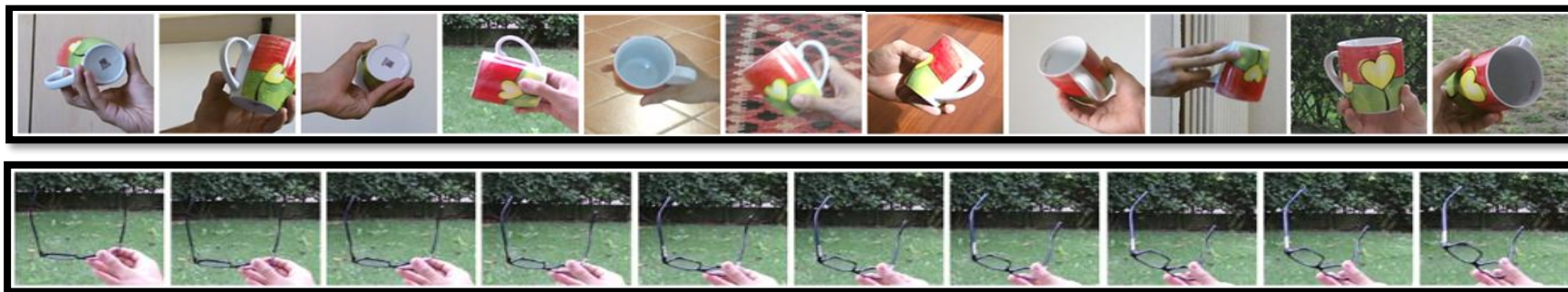
Some examples

Table 3: Benchmarks and environments for continual learning. For each resource, paper use cases in the NI, NC and NIC scenarios are reported.

Benchmark	NI	NC	NIC	Use Cases
Split MNIST/Fashion MNIST		✓		[83, 81, 57, 130]
Rotation MNIST	✓			[92, 83, 127]
Permutation MNIST	✓			[53, 73, 43, 150, 176, 83, 57, 127]
iCIFAR10/100		✓		[125, 97, 70]
SVHN		✓		[71, 145, 130]
CUB200	✓			[80]
CORe50	✓	✓	✓	[91, 115, 97]
iCubWorld28	✓			[116, 90]
iCubWorld-Transformation		✓		[117, 16]
LSUN		✓		[171]
ImageNet		✓		[125, 95]
Omniglot		✓		[77, 144]
Pascal VOC		✓		[104, 151]
Atari	✓			[136, 73, 144]
RNN CL benchmark			✓	[153]
CRLMaze (based on VizDoom)	✓			[89]
DeepMind Lab	✓			[99]

Example: CORe50

- Continuous Object Recognition
 - 50 classes
 - Short videos of object manipulation with different background
 - Temporal coherence from videos
- Many scenarios: batch, online, with repetitions.

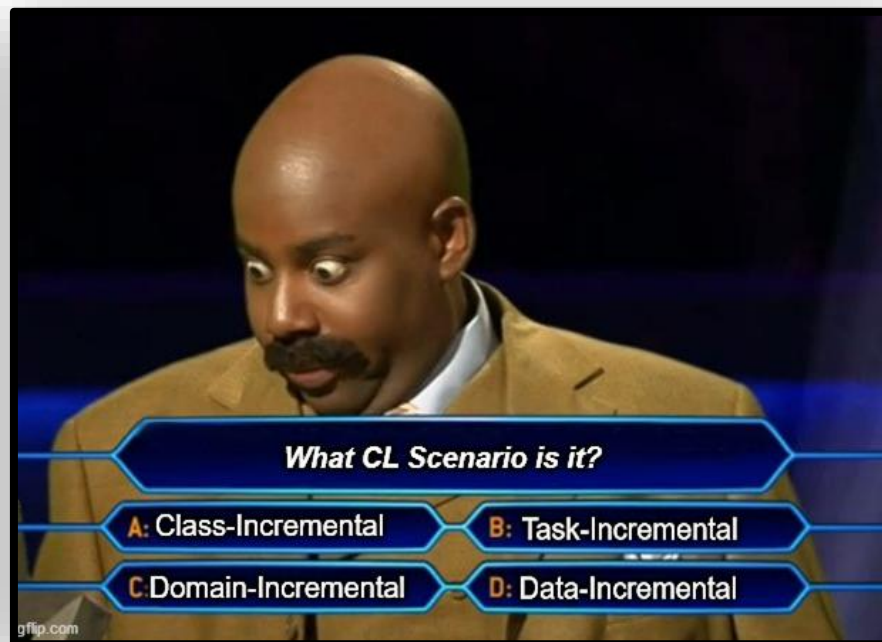




# Images	164,866
Format	RGB-D
Image size	350x350 128x128
# Categories	10
# Obj. x Cat.	5
# Sessions	11
# img. x Sess.	~300
# Outdoor Sess.	3
Acquisition Sett.	Hand held

New Instances
& Classes

Task-
Agnostic



Task-
Free

Single-
Incremental-Task

Multi-Task

CL – Evaluation

Average Stream Accuracy

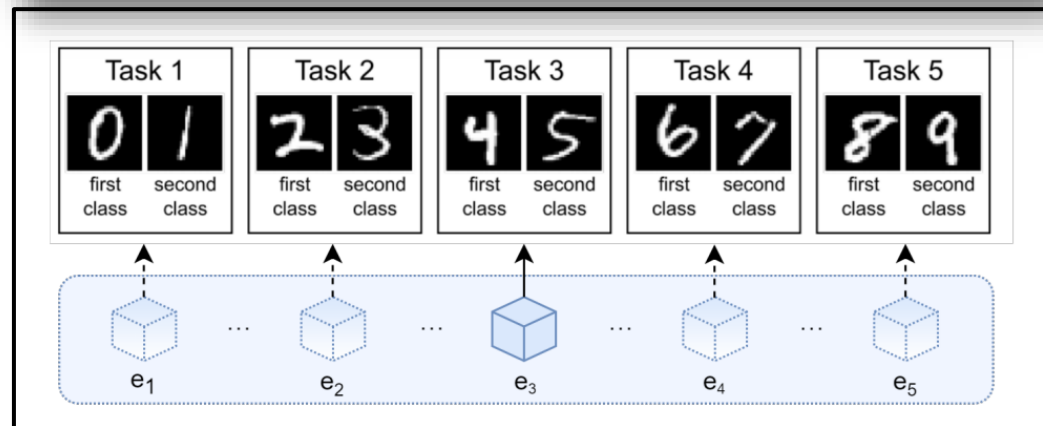
- The model is evaluated on the average accuracy on the entire (test) stream
- It must remember how the classify data from old experiences

The objective of a CL algorithm is to minimize the loss \mathcal{L}_S over the entire stream of data S :

$$\mathcal{L}_S(f_n^{CL}, n) = \frac{1}{\sum_{i=1}^n |\mathcal{D}_{test}^i|} \sum_{i=1}^n \mathcal{L}_{exp}(f_n^{CL}, \mathcal{D}_{test}^i) \quad (2)$$

$$\mathcal{L}_{exp}(f_n^{CL}, \mathcal{D}_{test}^i) = \sum_{j=1}^{|\mathcal{D}_{test}^i|} \mathcal{L}(f_n^{CL}(x_j^i), y_j^i), \quad (3)$$

where the loss $\mathcal{L}(f_n^{CL}(x), y)$ is computed on a single sample $\langle x, y \rangle$, such as cross-entropy in classification problems.



- We assume to have access to three parallel stream
 - Train/Validation/Test streams
 - At time t , they all have data from the same distribution
- Offline Model Selection
 - Train several models on the train stream
 - Select best on the entire validation stream
- Final evaluation on the test stream
- Simple to implement but unrealistic. Assume to have access to the entire stream at the end of training for model selection purposes

Early Model Selection



- Use the first k experiences for model selection
 - Train sequentially on the training splits (until $\text{time}=k$)
 - Evaluate on the validation splits
 - Select the best model on the first k experiences of the validation stream
- Continue training the best model on the rest of the training stream
- Model selection is still offline, but only for the first part of the stream

Continual Hyperparameter Selection



- Can we do the model selection without access to the old data?
- We have two objectives:
 - Maximize plasticity (learning new experiences)
 - Accuracy on current data
 - Is to estimate given the current validation experience
 - Minimize forgetting (of older experiences)
 - Accuracy on past experiences (for the current model)
 - We don't have the data to evaluate this objective
- Let's assume that we have only two hyperparameters:
 - One controls plasticity
 - The other controls stability (forgetting)

Continual Hyperparameter Selection



Example hyperparameters:

- Plasticity: learning rate
- Stability: regularization strength

PSEUDOCODE:

- For each experience:
 - STEP 1: Find optimal plasticity hyperparameters
 - This step will find the max accuracy you can get at the expense of stability
 - fix them
 - STEP 2: Find stability hyperparameters
 - Start with maximal stability
 - Decrease stability hyperparameters until you have a good enough accuracy
 - This a stability-plasticity tradeoff. If you stop too soon you have low plasticity. If you stop too late, you have too much forgetting.

Continual Hyperparameter Selection



ALGORITHM:

- for each experience:
 - finetune on new data, coarse grid search on lr, get acc A
 - train on new data with CL method (lr from prev step), get acc A^*
 - while performance on new data is too low ($A - A^*$ too big)
 - Decrease stability hparams
 - train on new data, get acc A^*

INTUITIVELY: after finding optimal plasticity hparams, decrease forgetting hparams until the performance is close enough to the optimal plasticity (tradeoff)

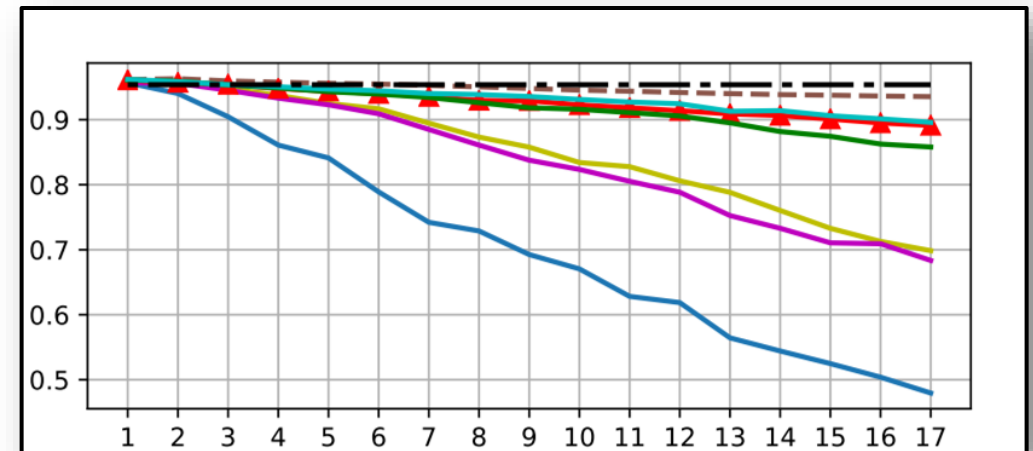
Algorithm 1. Continual Hyperparameter Selection Framework

```
input  $\mathcal{H}$  hyperparameter set,  $\alpha \in [0, 1]$  decaying factor,  $p \in [0, 1]$  accuracy drop margin,  $D^{t+1}$  new task data,  $\Psi$  coarse learning rate grid
require  $\theta^t$  previous task model parameters
require CLM continual learning method
  //Maximal Plasticity Search
1:  $A^* = 0$ 
2: for  $\eta \in \Psi$  do
3:    $A \leftarrow \text{Finetune}(D^{t+1}, \eta; \theta^t) \triangleright$  Finetuning accuracy
4:   if  $A > A^*$  then
5:      $A^*, \eta^* \leftarrow A, \eta \triangleright$  Update best values
  //Stability Decay
6: do
7:    $A \leftarrow \text{CLM}(D^{t+1}, \eta^*; \theta^t)$ 
8:   if  $A < (1 - p)A^*$  then
9:      $\mathcal{H} \leftarrow \alpha \cdot \mathcal{H} \triangleright$  Hyperparameter decay
10: while  $A < (1 - p)A^*$ 
```

CL – Metrics

What to monitor?

- Performance on current/past/future experience
- Resource consumption: Memory, CPU, Disk usage
- Model size growth
- Execution time and latency
- Data efficiency
- ...



Accuracy



Let N be the stream length, T the current timestep, $R_{t,i}$ be the accuracy on the experience i at time t

Time – which model to use:

- Last model: $R_{T,T}$
- Averaged over time: $\frac{1}{T+1} \sum_{t=0}^T \sum_{i=0}^t R_{t,i}$

Data – which data to use:

- Current experience: $R_{t,t}$
- Data seen up to now: $\frac{1}{T+1} \sum_{i=0}^T R_{T,i}$
- Full stream: $\frac{1}{N} \sum_{i=0}^{N-1} R_{T,i}$

R	Te_1	Te_2	Te_3
Tr_1	$R_{1,1}$	$R_{1,2}$	$R_{1,3}$
Tr_2	$R_{2,1}$	$R_{2,2}$	$R_{2,3}$
Tr_3	$R_{3,1}$	$R_{3,2}$	$R_{3,3}$

$$\text{Average Accuracy: ACC} = \frac{1}{T} \sum_{i=1}^T R_{T,i}$$

$$A = \frac{\sum_{i=1}^N \sum_{j=1}^i R_{i,j}}{\frac{N(N+1)}{2}}$$

Zero-shot Forward Transfer



Q: Is continual learning improving the performance on future experiences?

- FWT Metric compares
 - Accuracy on future experience $i+k$ after training on experience i
- against
 - \bar{b}_i Accuracy on experience i of a model trained with a random initialization
 - Averaged over $i=2, \dots, T$

R	Te_1	Te_2	Te_3
Tr_1	$R_{1,1}$	$R_{1,2}$	$R_{1,3}$
Tr_2	$R_{2,1}$	$R_{2,2}$	$R_{2,3}$
Tr_3	$R_{3,1}$	$R_{3,2}$	$R_{3,3}$

$$\text{FWT} = \frac{1}{T-1} \sum_{i=2}^T R_{i-1,i} - \bar{b}_i.$$

- FWT assumes that the model can predict future experiences
 - Most models can't predict unseen classes
 - Only makes sense for new instances, not new classes
- Alternative solution: Evaluate whether the latent representation helps learning unseen tasks
- **Linear Probing:**
 - Learn a linear classifier on top of the learned representation
 - Compare against random feature extractor and previous models
- Measures whether the learned features transfer to the new data

Backward Transfer



Q: Is continual learning improving the performance on **OLD** experiences?

BWT Metric

- Accuracy on experience i after training on experience T

Minus:

- Accuracy on experience i after training on experience i

Averaged over $i=1, \dots, T-1$

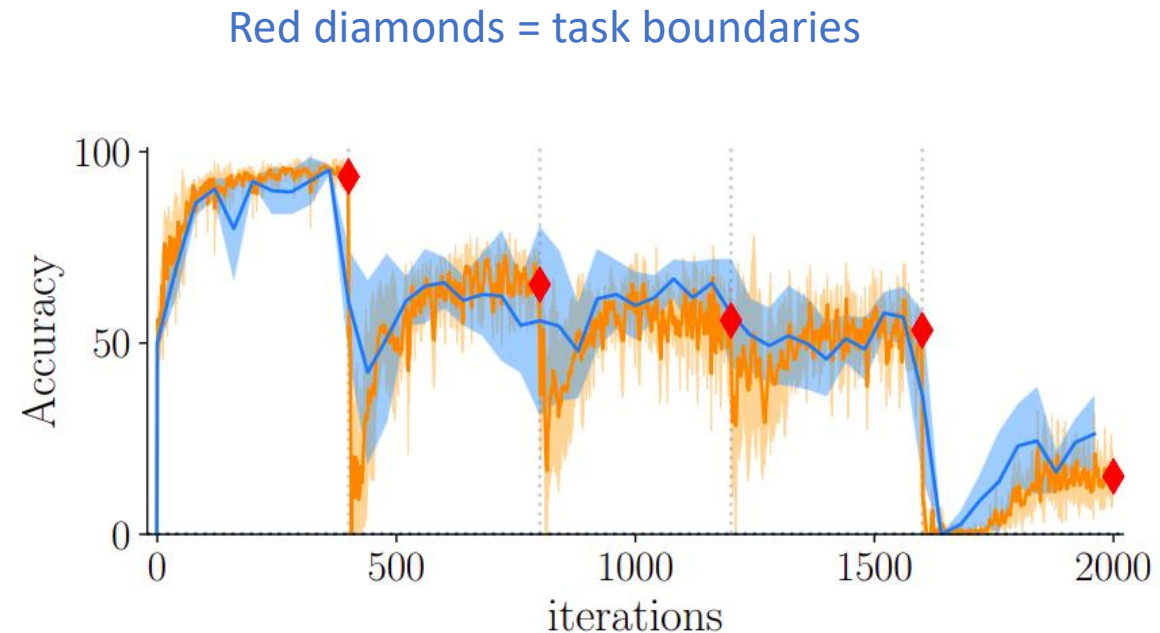
FORGETTING = - BWT

R	Te_1	Te_2	Te_3
Tr_1	$R_{1,1}$	$R_{1,2}$	$R_{1,3}$
Tr_2	$R_{2,1}$	$R_{2,2}$	$R_{2,3}$
Tr_3	$R_{3,1}$	$R_{3,2}$	$R_{3,3}$

$$\text{BWT} = \frac{1}{T-1} \sum_{i=1}^{T-1} R_{T,i} - R_{i,i}$$

Evaluation in Online CL is more difficult because the stream is longer and task boundaries are unknown.

- **Knowledge Accumulation:** the model should improve over time
 - At any point in time
 - High average accuracy but also fast adaptation
- **Continual Stability:** the model should not forget previous knowledge
 - At any point in time
 - We often assume virtual drifts when measuring stability
- **Representation Quality:** the latent representations should improve over time
 - A weaker form of knowledge accumulation/stability
 - Can be evaluated on out-of-distribution data or self-supervised models



Knowledge Accumulation

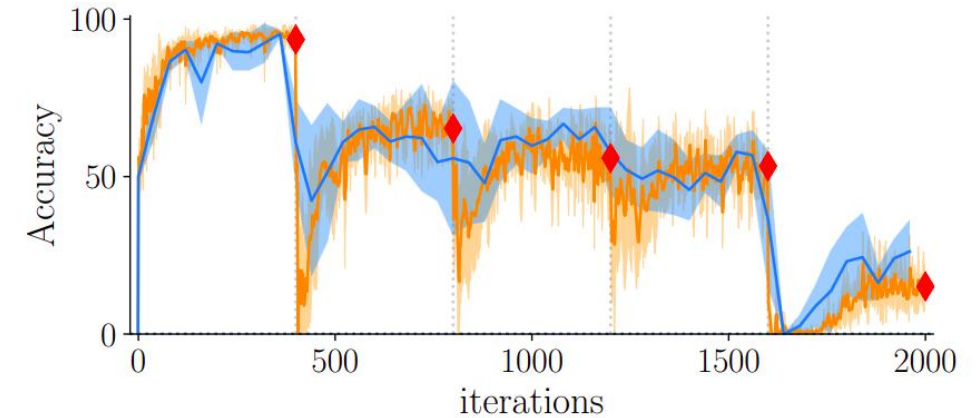


- **Average Anytime Accuracy:** accuracy along the entire curve.
- Do not confuse with
 - Avg accuracy at the end of training (final diamond)
 - Avg at task boundaries (avg of diamonds)

Notation:

- f_i model at time i
- E_i experience i
- $A(E_i, f_i)$ accuracy of model f_i for experience E_i

Red diamonds = task boundaries



$$AAA_t = \frac{1}{t} \sum_{j=1}^t \left(\frac{1}{k} \sum_{i=1}^k A(E_i, f_j) \right)$$

Average along the training curve

Average accuracy of data seen up to now

- **Model Size (MS):** How much space does your model occupy? (MB, # of params, etc.)
- Scalability over time: What is the increment in space required for each new experience?
- **Samples Storage Size (SSS):** How much space do you require for additional information (replay buffer, past models...)?

$$MS = \min\left(1, \frac{\sum_{i=1}^N \frac{Mem(\theta_1)}{Mem(\theta_i)}}{N}\right)$$

$$SSS = 1 - \min\left(1, \frac{\sum_{i=1}^N \frac{Mem(M_i)}{Mem(D)}}{N}\right)$$

Q: What is the computational overhead during training and inference?

- #MAC – Multiply and Accumulate
- Running Time, CPU/GPU time
- scalability over time

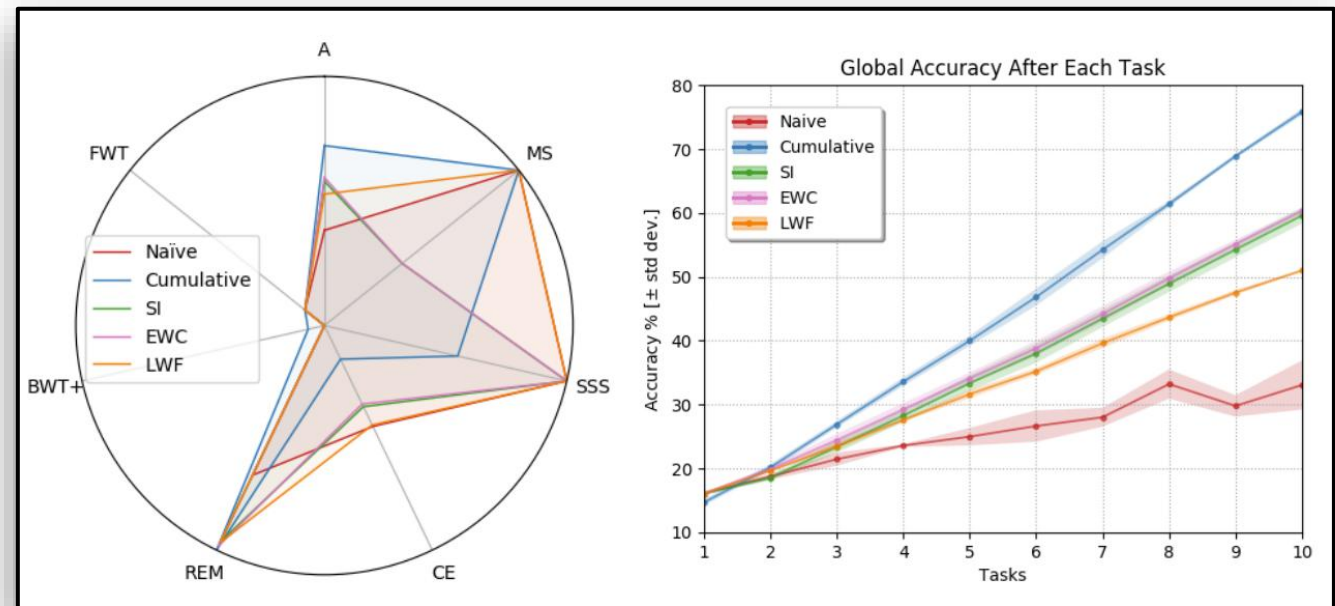
$$CE = \min\left(1, \frac{\sum_{i=1}^N \frac{Ops_{\uparrow\downarrow}(Tr_i) \cdot \epsilon}{1 + Ops(Tr_i)}}{N}\right)$$

NOTE: evaluating GPU performance is tricky because you have to consider CPU \leftrightarrow GPU communication, synchronization costs, parallelization. #MAC can be misleading.

In general, multiple objectives!

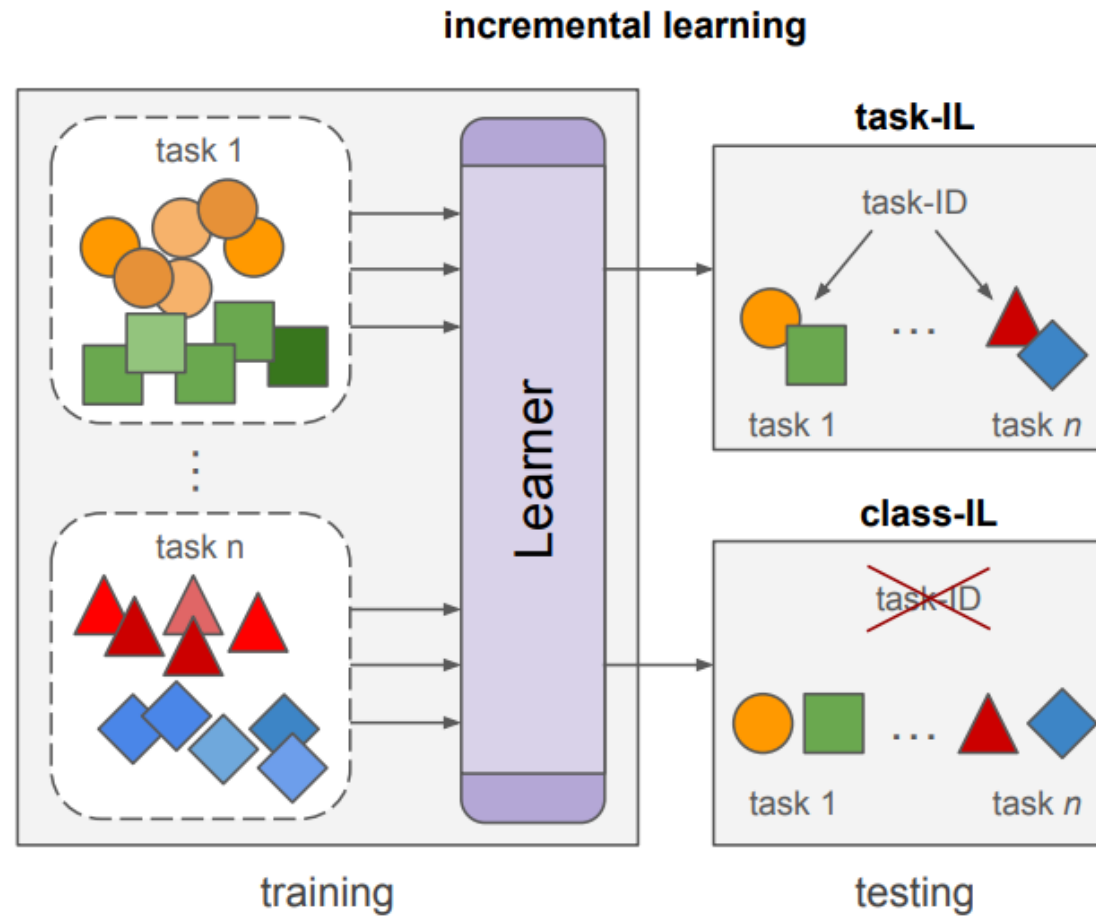
Each application is different and multiple objectives may interfere with each other:

- Computational constraints
- Privacy Constraints
- Accuracy and Forgetting



Forgetting in CL Scenarios

Task-Incremental vs Class-Incremental



Task-Incremental vs Class-Incremental

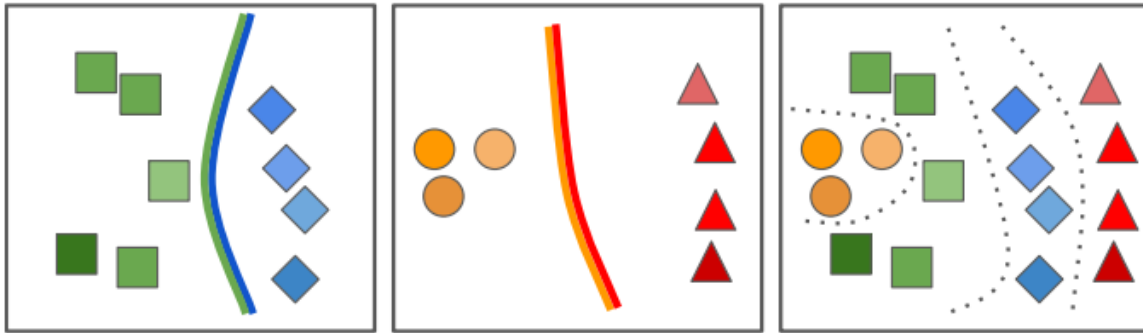


Fig. 2: A network trained continually to discriminate between task 1 (left) and task 2 (middle) is unlikely to have learned features to discriminate between the four classes (right). We call this problem *inter-task confusion*.

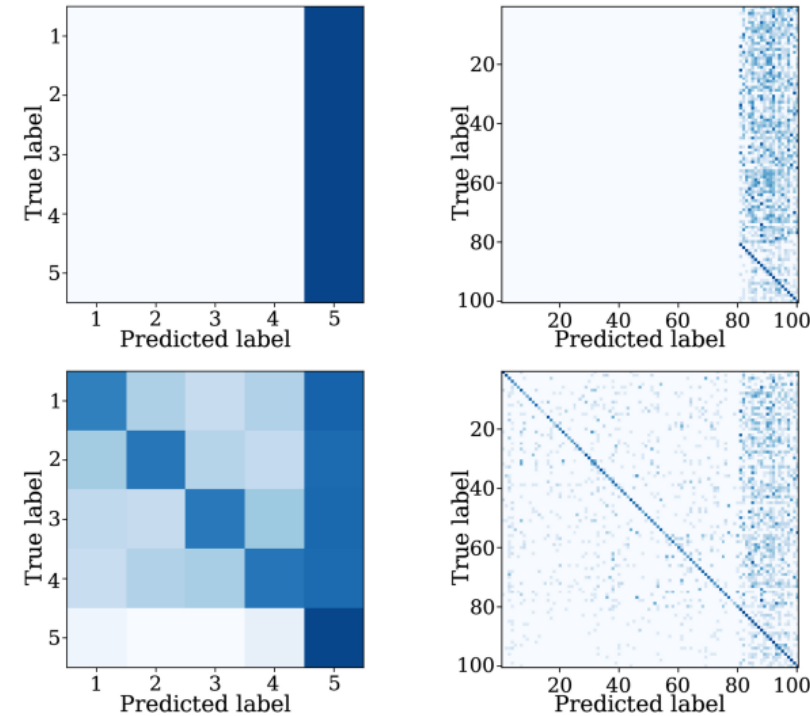


Fig. 3: Examples of task and class confusion matrices for Finetuning (top row) and Finetuning with 2,000 exemplars (bottom row) on CIFAR-100. Note the large bias towards the classes of the last task for Finetuning. By exploiting exemplars, the resulting classifier is clearly less biased.

Classifier Bias in CIL

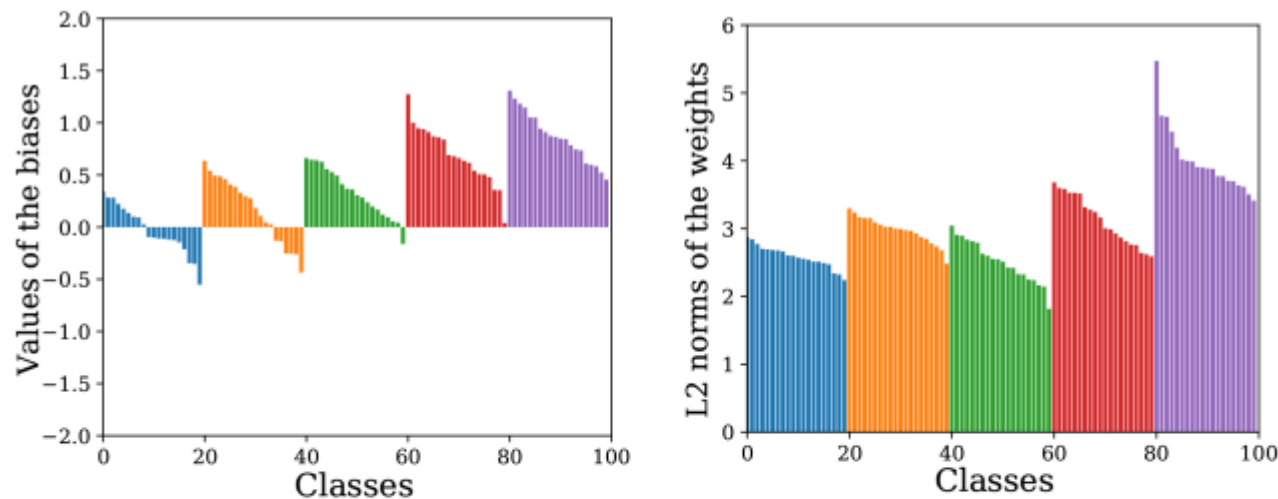


Fig. 4: Bias and weight analysis for iCaRL with **2,000 exemplars** on CIFAR-100. We show the ordered biases and norm of the last classification layer of the network for different tasks. Note how the bias and the norm of the weights are higher for the last tasks. This results in a *task-recency bias*.

Replay does not fix the task-recency bias

Self-Supervised Models



- SSL methods are more robust to continual training
- We evaluate the representation (linear probing)
- More realistic because we don't need labels during training
- We still need labels to finetune the classifier (linear probing), otherwise we can't use the model
- SSL methods are slow to converge and require large amounts of data

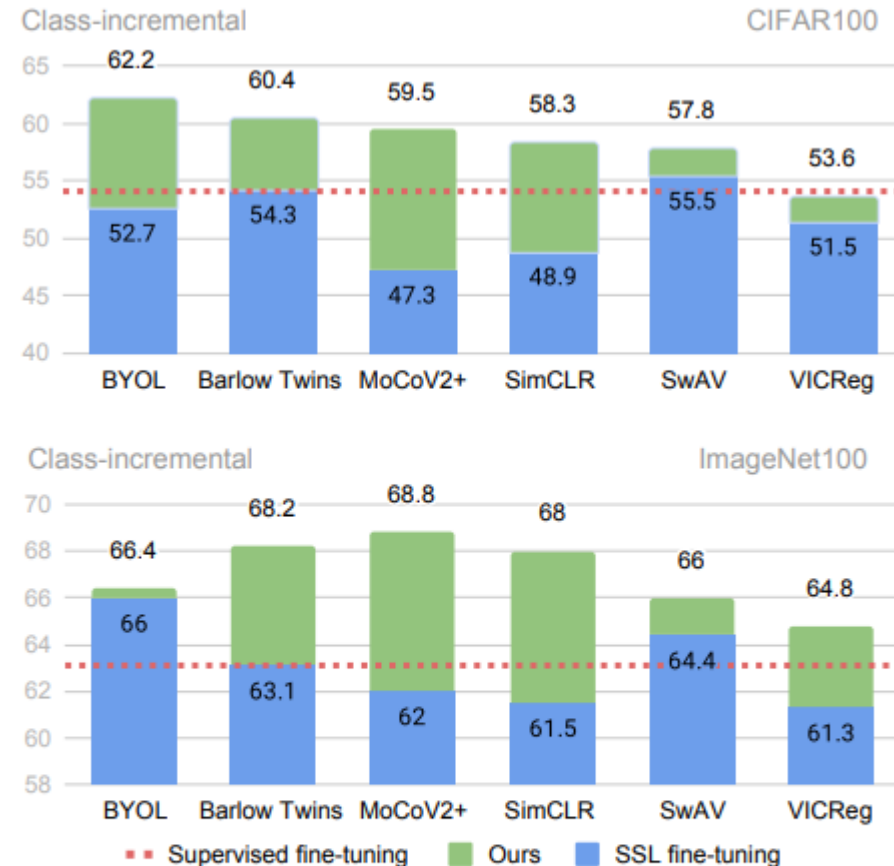


Figure 1. Linear evaluation accuracy of representations learned with different self-supervised methods on class-incremental CIFAR100 and ImageNet100. In blue the accuracy of SSL fine-tuning, in green the improvement brought by CaSSLe. The red dashed line is the accuracy attained by supervised fine-tuning.

Continual Pretraining

- SSL pre-training, incrementally over time
- Finetuning on the downstream tasks using the last model

Again, SSL > supervised, both in vision and text

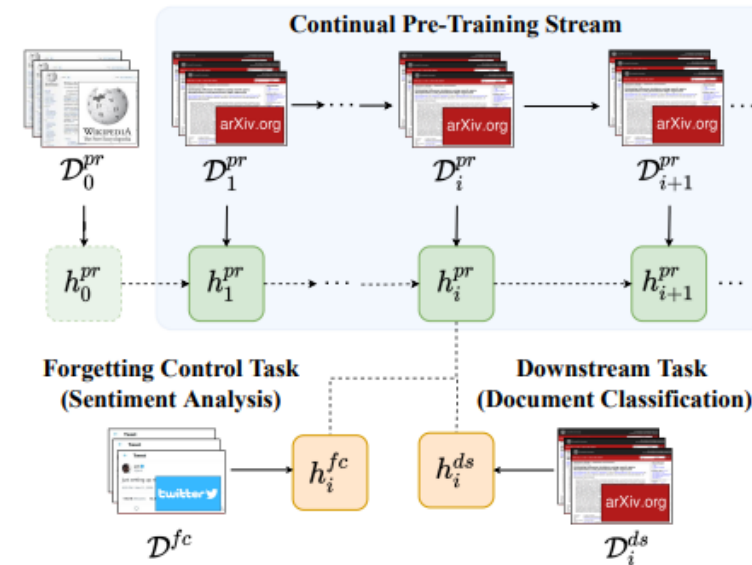
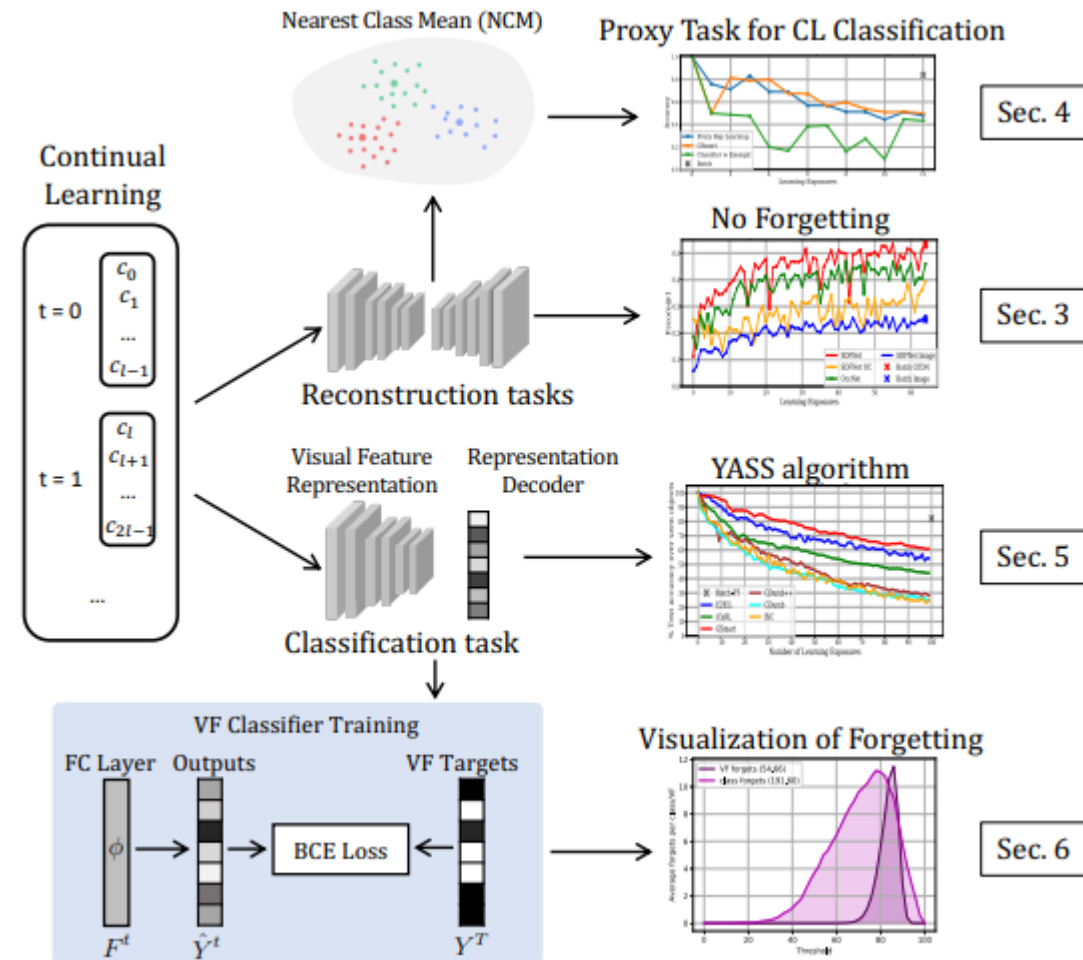


Figure 1: The Continual Pre-training scenario. During each stage (experience) i of continual pre-training (top), the model h_i^{pr} is pre-trained (center) on the dataset \mathcal{D}_i^{pr} (e.g., *scientific abstracts*). Subsequently (bottom), the model is fine-tuned against one (or more) downstream task \mathcal{D}_i^{ds} (e.g. *scientific abstracts* classification). Forgetting is measure by fine-tuning on \mathcal{D}^{fc} (e.g. *sentiment analysis*). At each stage, only the current pre-trained and downstream datasets/models are available.

Forgetting in Different Scenarios

- Some tasks are much more robust to CL than others
- Incremental classification results in catastrophic forgetting
- SSL methods are more robust
- Tasks such as reconstruction are very robust
- Forgetting will also depend on the drifts (iid vs class vs domain vs gradual...)



Take-Home Messages



While DCL is similar to OCL, there are many differences:

- Domains: vision, NLP, ...
- Drift: real vs virtual
- Evaluation: avg. accuracy vs prequential
 - Forgetting is a major concern which we don't have in the prequential setting
 - We can have forward transfer with deep learning models
- DCL methods are much more expensive than OML. The main challenge is avoiding forgetting, not fast adaptation

General References



- CL Course: <https://course.continualai.org/>
 - Ph.D. level but it's recorded and it has most (but not all) of the topics covered here
- Avalanche tutorials: https://avalanche.continualai.org/from-zero-to-hero-tutorial/01_introduction

- Notebooks:
 - Avalanche e2e
 - Forgetting
 - Avalanche standalone

We start looking at the algorithms:

- Baselines for CL
- Replay Methods