

# Continual Learning with Deep Architectures

Tutorial @ ICML 2021

## Part 1: Introduction and State-of-the-Art

Irina Rish

University of Montreal & Mila

*irina.rish@mila.quebec*

# Irina Rish

*Associate Professor @  
University of Montreal*

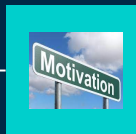
*Core member @ Mila -  
Quebec AI Institute*

*Canada Excellence Research  
Chair (CERC) in Autonomous AI*

*CIFAR chair*



# Outline



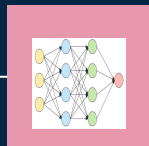
01

Motivation  
& History



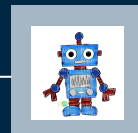
02

Inspirations  
from  
Neuroscience



03

Supervised  
Continual  
Learning



04

Continual  
Reinforcement  
Learning



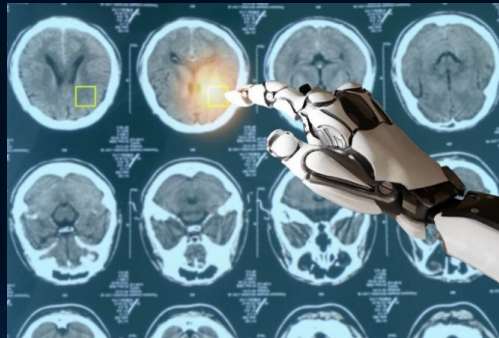
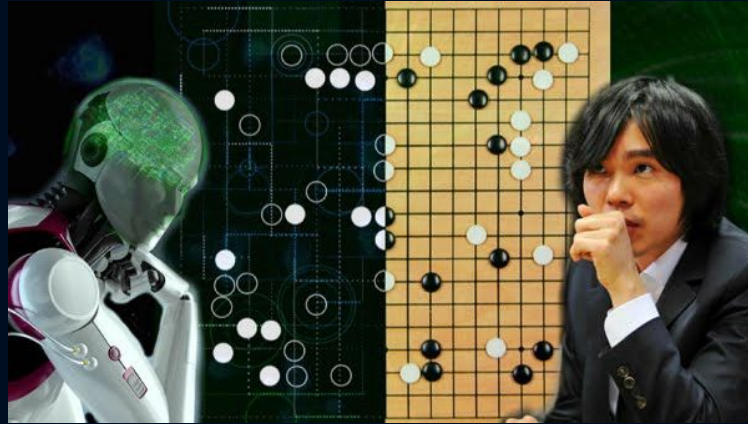
05

Summary



# Motivation & History

# AI Today: Impressive... but (Still) "Narrow"



theguardian.com/commentisfree/2020/sep/08/robot-wrote-this-article-gpt-3

ASO Clients Health Home JS Machine Learning Misc Search Stuff Sharepoint

gence (AI)

## A robot wrote this entire article. Are you scared yet, human?

GPT-3

# OpenAI GPT-3

language generator, to assignment? To

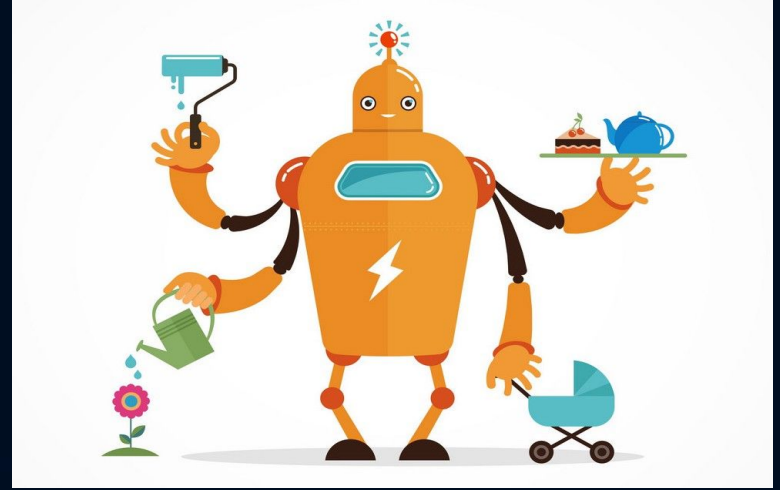
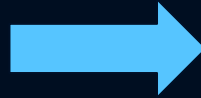
ay was written and

ow

**175 BILLION** Parameters

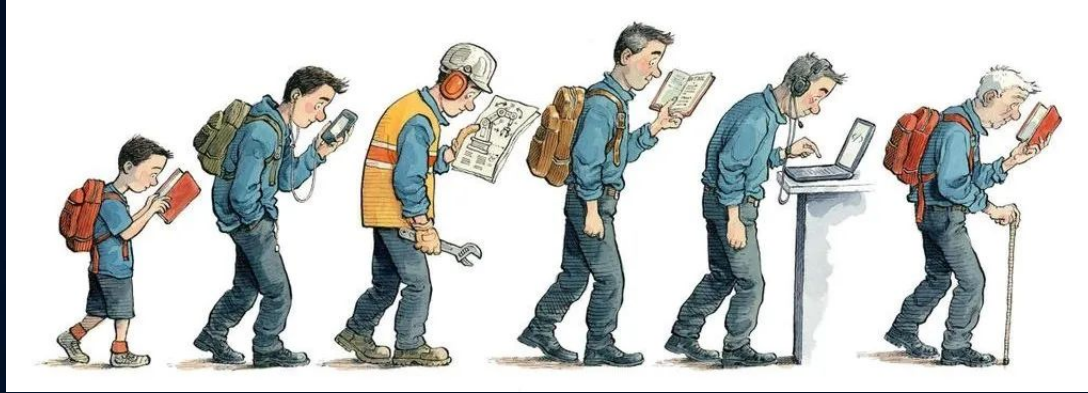
The bar chart shows the number of parameters in GPT-3 models. The y-axis is logarithmic, ranging from 1 to 1000. The x-axis represents different model versions. The bars show a steady increase in parameters, with the final bar (GPT-3) reaching 175 billion parameters, highlighted by a red circle.

# Human-Level AI: “Broad” – Versatile, Multi-Task



**How Do We Achieve This?**

# Lifelong, Continual Learning

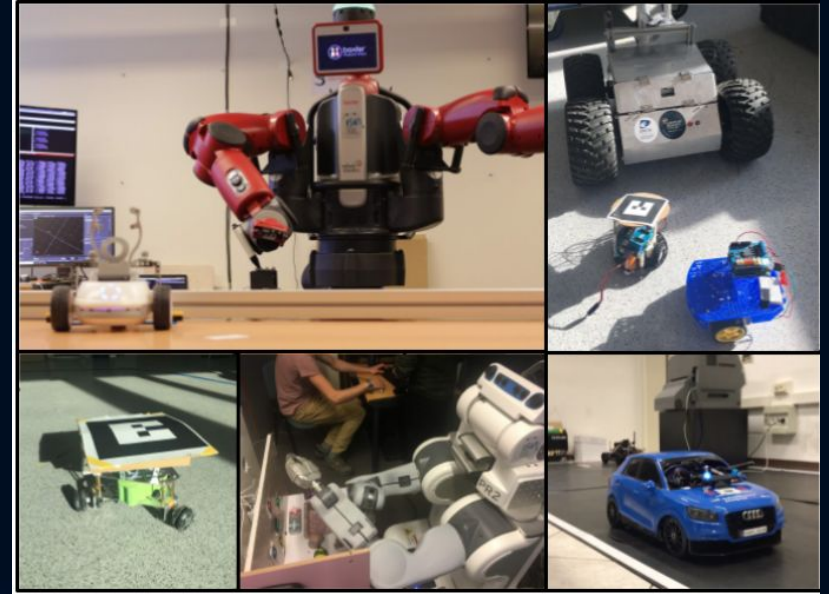


“Continual learning is the constant development of increasingly complex behaviors; the process of building more complicated skills on top of those already developed.”

Ring (1997). CHILD: A First Step Towards Continual Learning.

# Continual Learning Challenges in Practical Applications

A robot acquiring new skills in different environment, adapting to new situations, learning new tasks



S. Thrun and T. Mitchell. **Lifelong robot learning**. *Robotics and Autonomous Systems*, 15:25-46, 1995.

Lesort, T., Lomonaco, V., Stoian, A., Maltoni, D., Filliat, D. and Díaz-Rodríguez, N., **Continual learning for robotics: Definition, framework, learning strategies, opportunities and challenges**. *Information fusion*, 2020.



# Continual Learning Challenges in Practical Applications

A self-driving car adapting to different environments (from a country road to a highway to a city)



# Continual Learning Challenges in Practical Applications

Conversational agents adapting to different users, situations, tasks



# Continual Learning Challenges in Practical Applications

Medical applications:  
adapting to new patients, hos  
conditions



# Continual Learning Challenges in Practical Applications

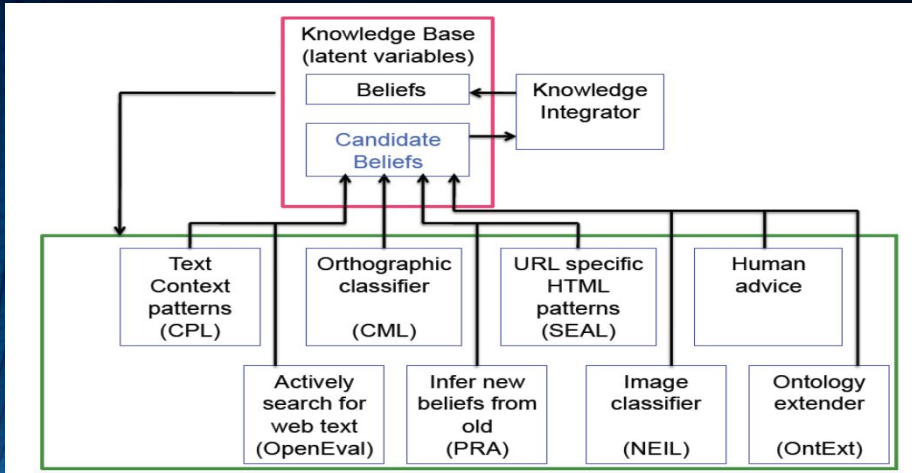
Multi-game environments  
(e.g. OpenAI gym)



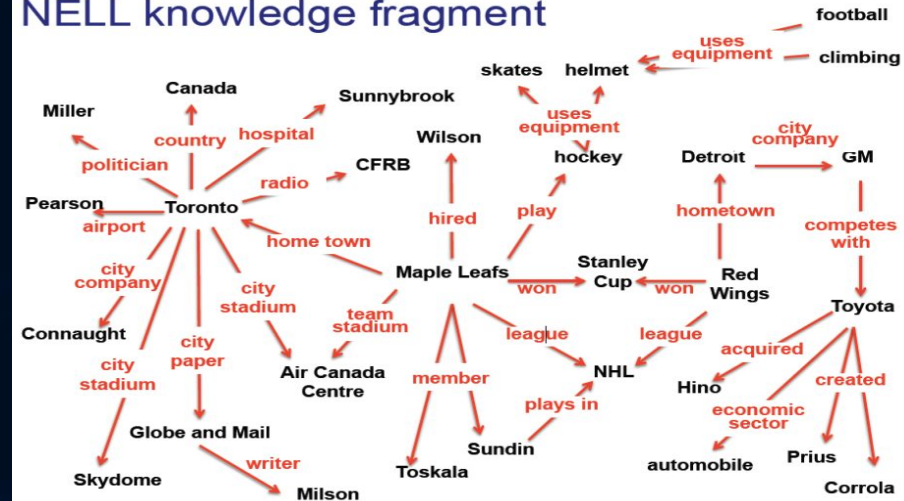
# Example: Never-Ending Language Learner (NELL)

Mitchell et al, Never-Ending Learning, AAAI-2015

<http://rtw.ml.cmu.edu/rtw/>

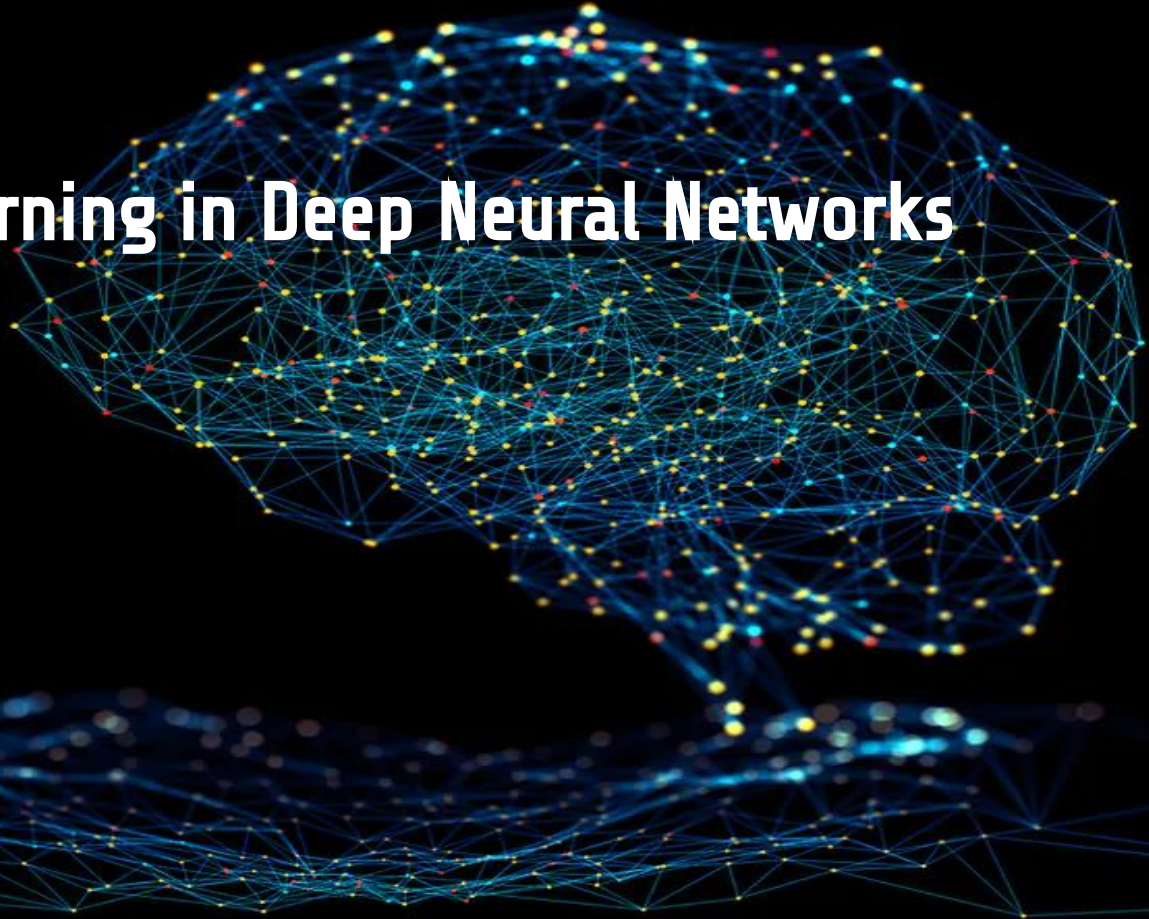


## NELL knowledge fragment



**Our Focus:**

**Continual Learning in Deep Neural Networks**

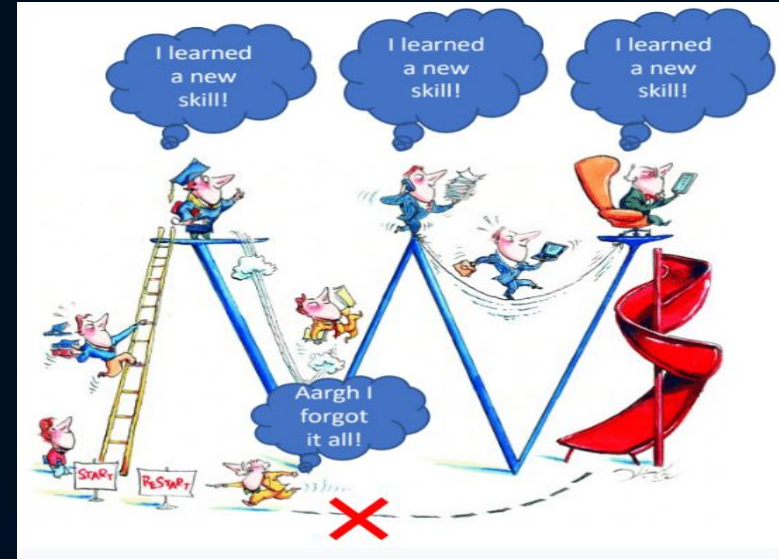


# Challenge: Catastrophic Forgetting (McCloskey and Cohen, 1989)

“...the process of learning a new set of patterns suddenly and completely erased a network’s knowledge of what it had already learned.” (French, 1999)

CF was identified by (McCloskey and Cohen, 1989):

- A neural net trained with backprop learned a set of “one’s addition facts” (i.e., the 17 sums 1+1 through 9+1 and 1+2 through 1+9)
- Then the network learned the 17 “two’s addition facts” (2+1 through 2+9, 1+2 through 9+2).
- Within 1-5 two’s learning trials, accuracy on task 1 had dropped from 100% to 20%, in 5 more trials, to 1%; by 15 trials, to 0%.



McCloskey, M. and Cohen, N.J., 1989. Catastrophic interference in connectionist networks: The sequential learning problem. In *Psychology of learning and motivation* (Vol. 24, pp. 109-165). Academic Press.

French, R.M., 1999. Catastrophic forgetting in connectionist networks. *Trends in cognitive sciences*, 3(4).

# More Generally: Stability vs Plasticity (Carpenter and Grossberg, 1987)

Plasticity  $\Leftrightarrow$  ability to adapt to a new task

Stability  $\Leftrightarrow$  ability to retain the learned skills on the old tasks

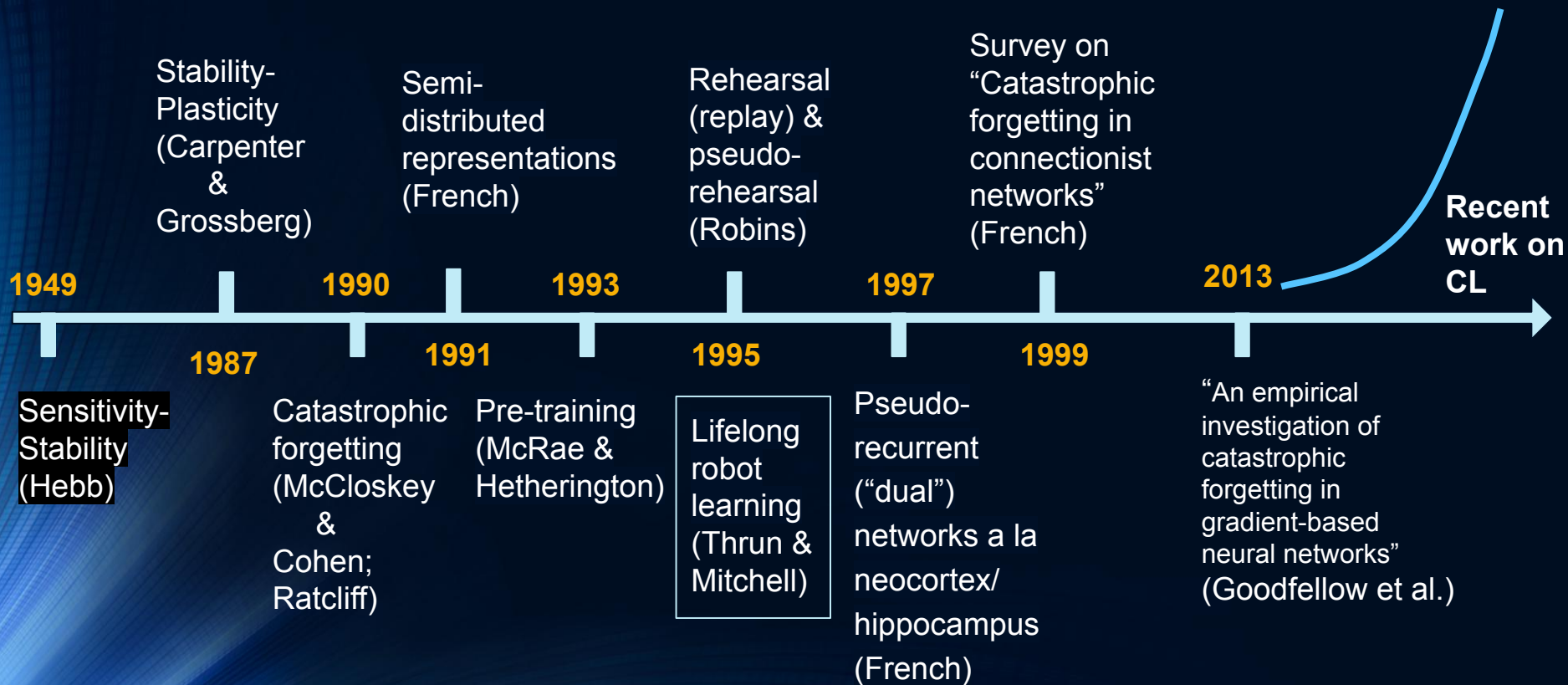
“Catastrophic interference is a radical manifestation of a more general problem for connectionist models of memory — in fact, for any model of memory — the so-called “stability-plasticity” problem [1,2]. **The problem is how to design a system that is simultaneously sensitive to, but not radically disrupted by, new input.**” (French, 1999)

[1] Grossberg, S. (1982) Studies of Mind and Brain: Neural Principles of Learning, Perception, Development, Cognition, and Motor Control.

[2] Carpenter, G. and Grossberg, S. (1987) ART 2: Self-organization of stable category recognition codes for analog input patterns.



# Brief History: CL in Neural Networks





# Inspirations from Neuroscience

# Synaptic Plasticity Regulation for Retaining Knowledge

Learning -> enlargement of (some) dendritic spines -> decreased plasticity of the corresponding synapses.

(Cichon&Gan, 2015; Yang et al., 2009)

Persistent change (months), despite learning new tasks.

If these changes removed via synaptic optogenetics, the task is forgotten (Hayashi-Takagi et al., 2015).



Two-photon data  
(structural imaging)

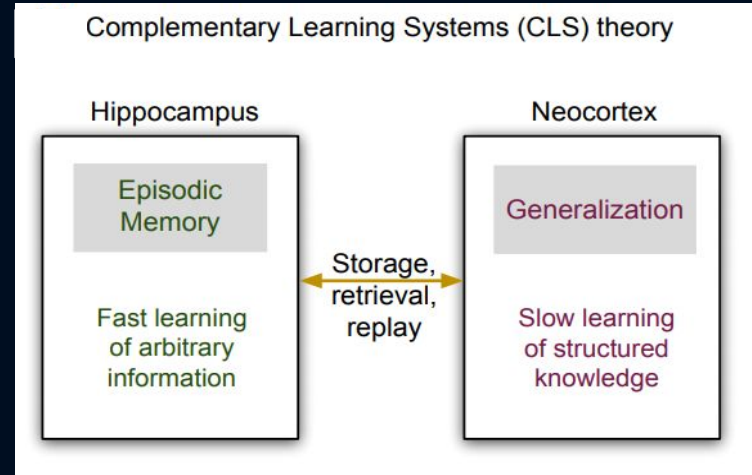
Hassabis et al (2017).  
Neuroscience-inspired  
artificial intelligence.

An inspiration for Elastic Weight Consolidation (EWC), Synaptic Intelligence (SI) and other **regularization methods** for preserving task-important weights.

# Complementary Learning Systems and Experience Replay

CLS theory (McClelland et al. 1995):

- **hippocampus**: fast (one-shot) learning of episodic information, consolidated to the neocortex in sleep (or resting periods) via “replay” of neural activity patterns associated with the episode
- **neocortex**: slow learning of structured knowledge; efficient representation for generalization.



Parisi et al. (2019) Continual lifelong learning with neural networks: A review.

An inspiration for rehearsal/pseudo-rehearsal (Robins, 1995), pseudo-recurrent (“dual”) networks (French, 1997) and many modern **experience replay methods**.

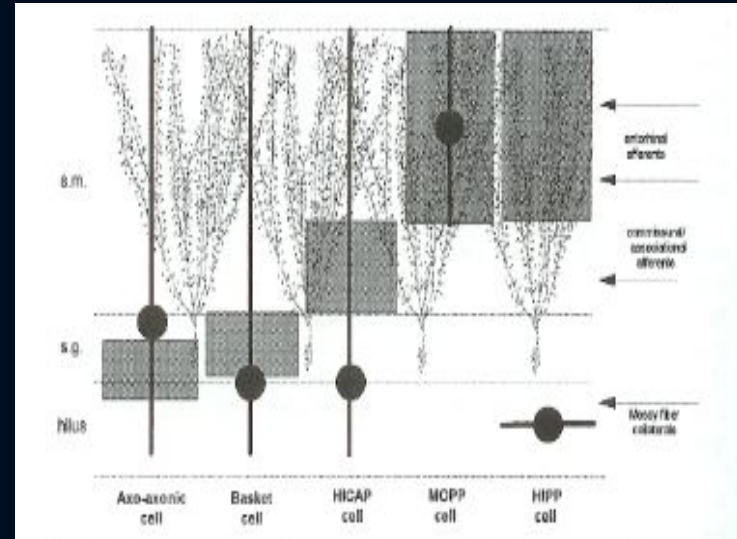
# Structural Plasticity via Neurogenesis


**Adult neurogenesis:** generation of new neurons in adult brains throughout life, balanced by death of unused neurons (“use it or lose it”)

In humans, it occurs primarily in the **dentate gyrus** of the **hippocampus**

Increased neurogenesis is associated with **better** adaptation to new environments.

An inspiration for **adaptive, expanding neural architecture** methods.





# Supervised Continual Learning

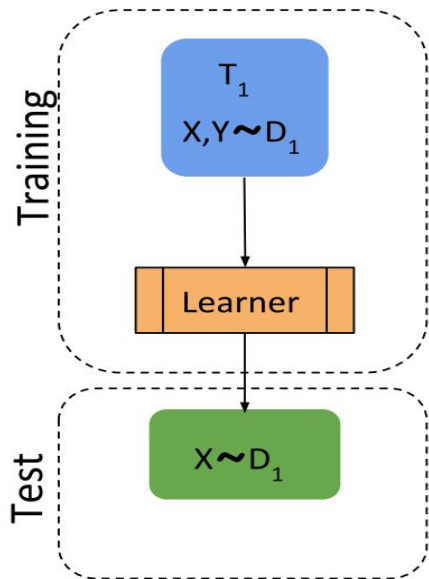
# Supervised Continual Learning

Non-stationary data comes one example at a time in a stream:

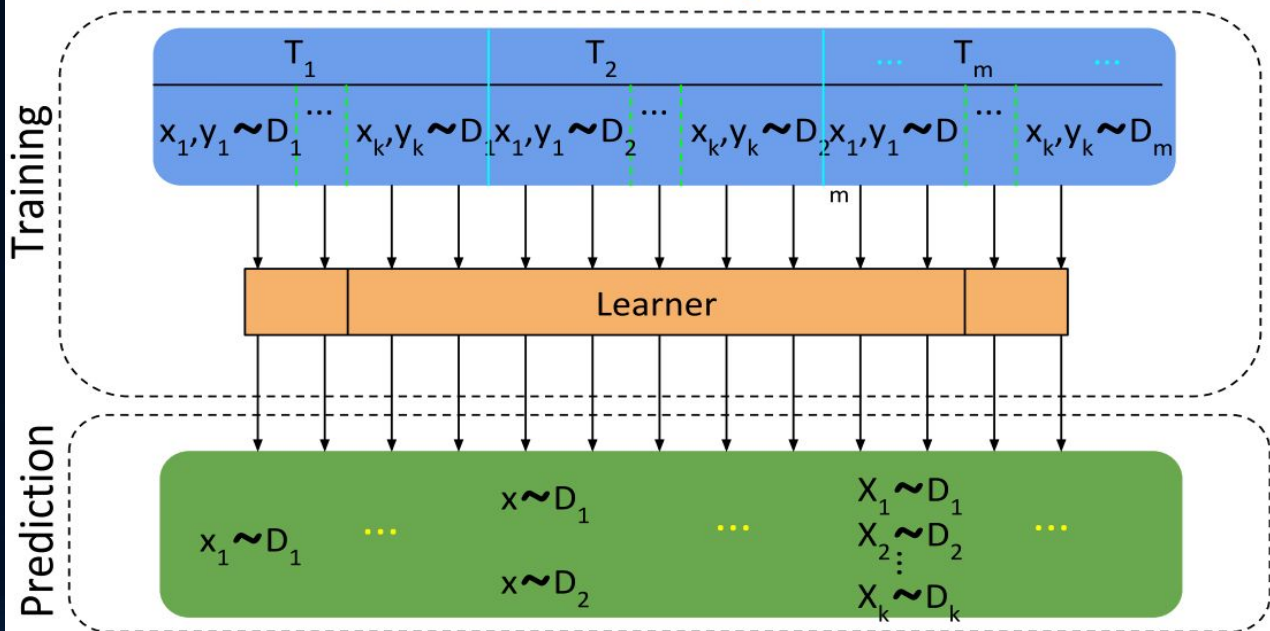
$$(x_1, y_1, t_1), \dots, (x_i, y_i, t_i), \dots, (x_{i+j}, y_{i+j}, t_{i+j})$$

Our data is *locally i.i.d.* – samples for a task are drawn from the same unknown joint probability distribution  $x_i, y_i \sim P_t(x, y)$ .

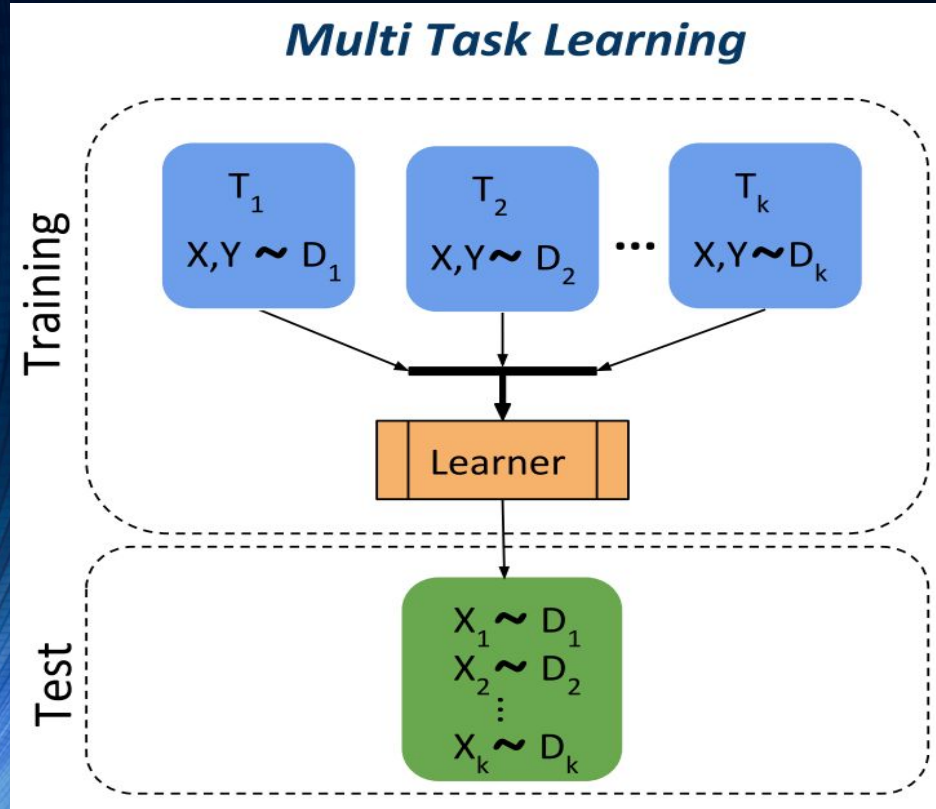
## Standard Supervised Learning



## Continual Learning



# Continual vs Multi-Task Learning



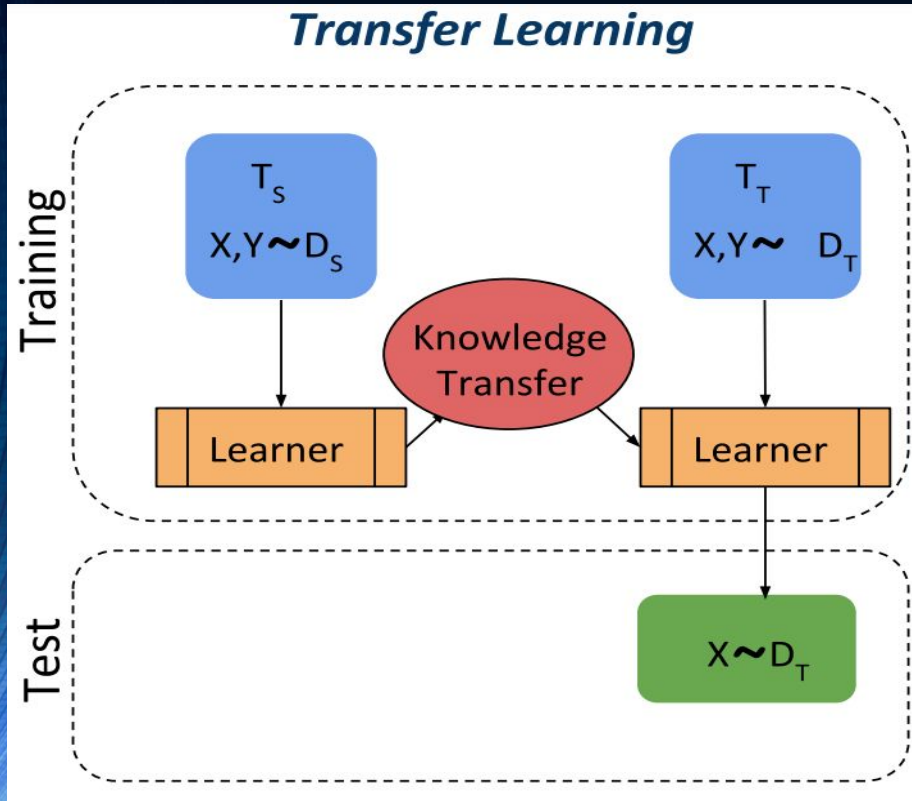
Learning of multiple related tasks **offline**, simultaneously

Using a set or subset of shared parameters

**No continual model adaptation**



# Continual vs Transfer Learning

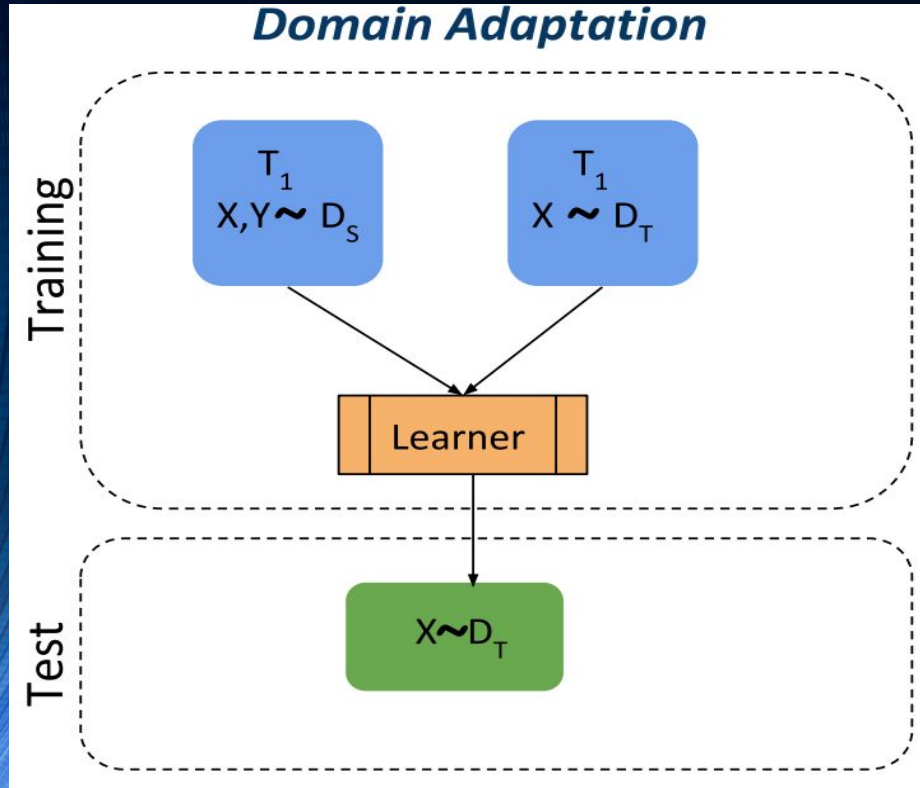


Help learning the target task using model trained on the source task

No continuous adaptation after learning the target task

Performance on the source task(s) is not taken into account

# Continual Learning vs Domain Adaptation



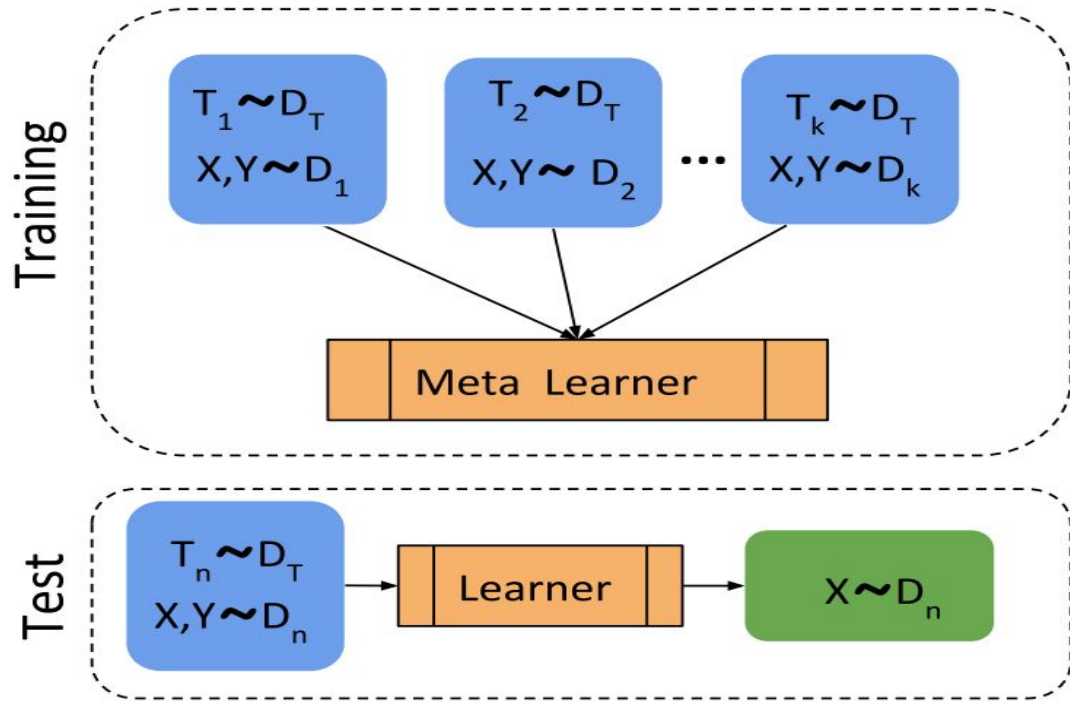
Transfer learning with same source and target tasks, but from different input domains

Trains on the source domain, adapts model to the (with no or only a few labels).

**Unidirectional; does not involve any accumulation of knowledge**

# Continual vs Meta-Learning

## *Meta Learning*



Faster adaptation on a task given a large number of training tasks

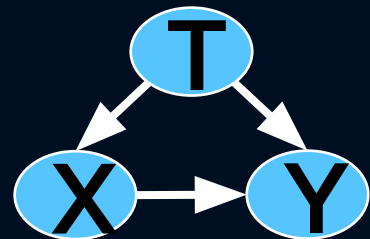
Offline training: a set of training tasks available at the same time

# Continual Learning Settings

X – input vector

Y – class label

T – task (context) defines  $P(X, Y|T)$



Task ID observed at training:

- T observed at test: task-incremental CL
- T not observed at test: class-incremental or domain-incremental CL

Task ID/boundary is not known at training:

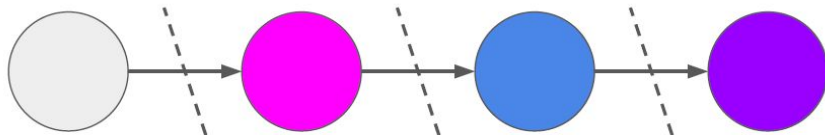
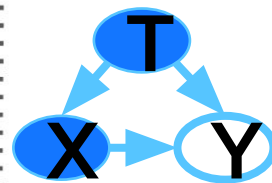
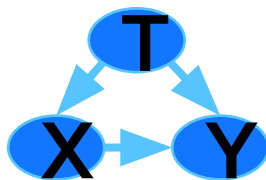
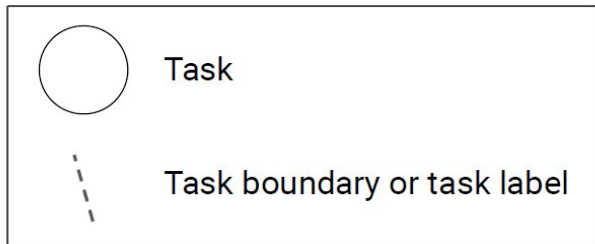
- Task-agnostic CL

# Task-Incremental CL

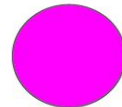
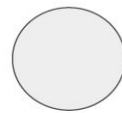
- Models are always informed about which task needs to be performed (both at train and test time)
- The easiest continual learning scenario; possible to train models with task-specific components
- A typical NN architecture: “multihead” output layer - each task has its own output units but the rest of the network may be shared between the tasks
- Assumptions:  $\{\mathcal{Y}^{(t)}\} \neq \{\mathcal{Y}^{(t+1)}\}$

# Task-Incremental CL

Task-Incremental Learning (or multi-head setting)



Training



Test

# Class-Incremental CL

- Models must be able not only to solve each task seen so far, but also to infer which task they are presented with.
- Includes protocols in which new classes need to be learned incrementally.

An example: sequentially learning MNIST digits (split-MNIST)

- Assumptions:

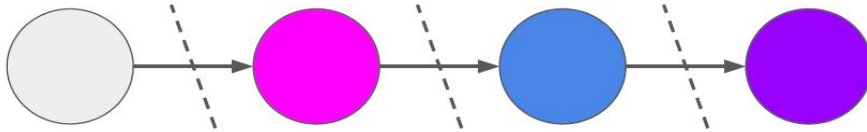
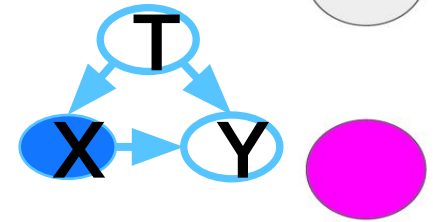
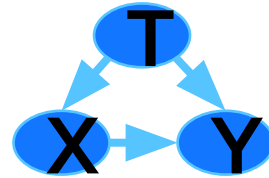
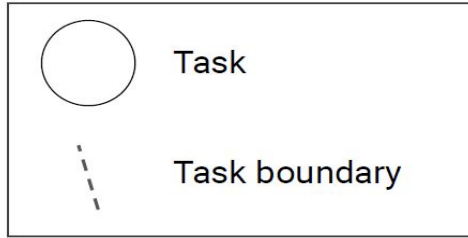
$$P(\mathcal{X}^{(t)}) \neq P(\mathcal{X}^{(t+1)})$$

$$\{\mathcal{Y}^{(t)}\} = \{\mathcal{Y}^{(t+1)}\}$$

$$P(\mathcal{Y}^{(t)}) \neq P(\mathcal{Y}^{(t+1)})$$

# Class-Incremental CL

Class-Incremental Learning (or shared-head setting)



Training

Test

Figure credit: Massimo Caccia



# Domain-Incremental CL

- Task identity is not available at test time
- Models however only need to solve the task at hand; they are not required to infer which task it is
- Typical examples of this scenario: the structure of the tasks is always the same, but the input-distribution is changing (e.g., 'permuted MNIST')
- Assumptions: similar to class-incremental, except for last one:

$$P(\mathcal{X}^{(t)}) \neq P(\mathcal{X}^{(t+1)})$$

$$\{\mathcal{Y}^{(t)}\} = \{\mathcal{Y}^{(t+1)}\}$$

$$P(\mathcal{Y}^{(t)}) = P(\mathcal{Y}^{(t+1)})$$

# Example: Split MNIST

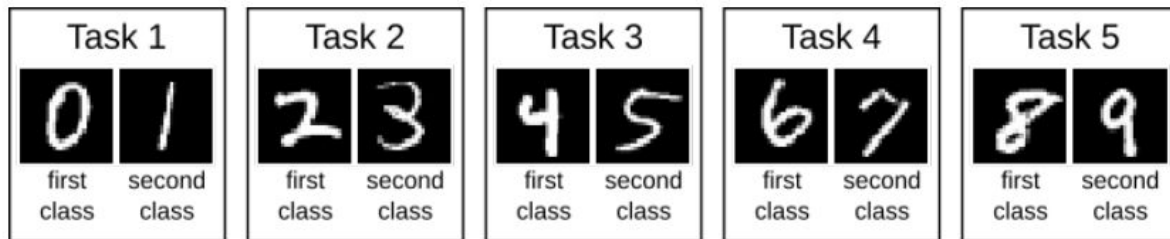


Figure 1: Schematic of the split MNIST task protocol.

Table 1: The split MNIST task protocol according to each continual learning scenario.

<b>Incremental task learning</b>	With task given, is it the first or second class? (e.g., '0' or '1')
<b>Incremental domain learning</b>	With task unknown, is it a first or second class? (e.g., in ['0', '2', '4', '6', '8'] or in ['1', '3', '5', '7', '9'])
<b>Incremental class learning</b>	With task unknown, which digit is it? (choice from '0' to '9')

# Example: Permuted MNIST

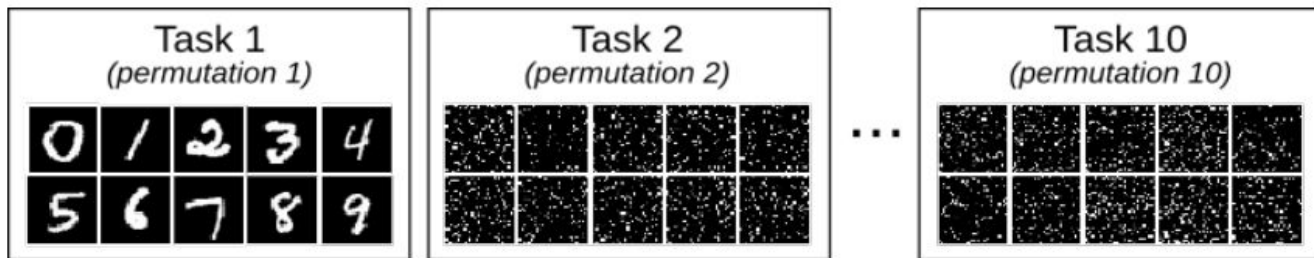


Figure 2: Schematic of the permuted MNIST task protocol.

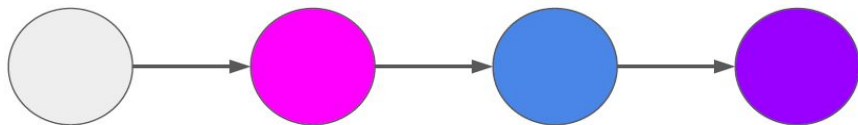
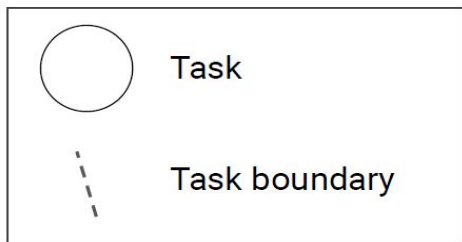
Table 2: The permuted MNIST task protocol according to each continual learning scenario.

<b>Incremental task learning</b>	Given permutation $X$ was applied, which digit is it?
<b>Incremental domain learning</b>	With permutation unknown, which digit is it?
<b>Incremental class learning</b>	Which digit is it <i>and</i> which permutation was applied?

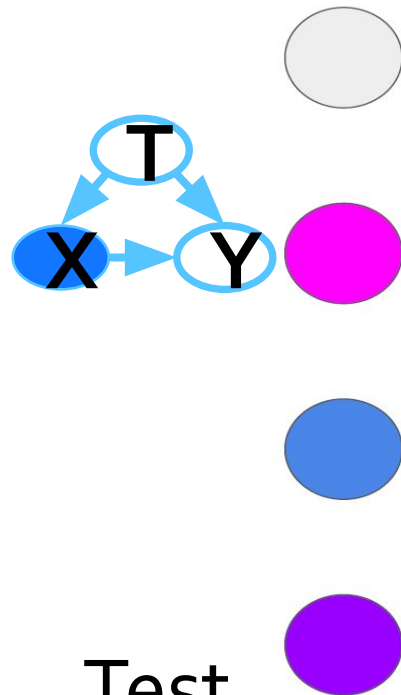
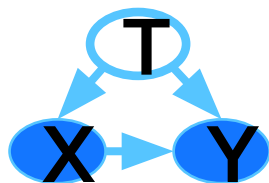
# Task-Agnostic CL: Most Challenging

Task identity is not available even at training time!

## Task Agnostic CL



Training



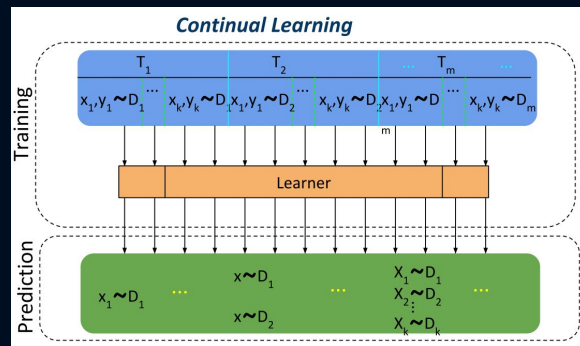
Test



# Objective

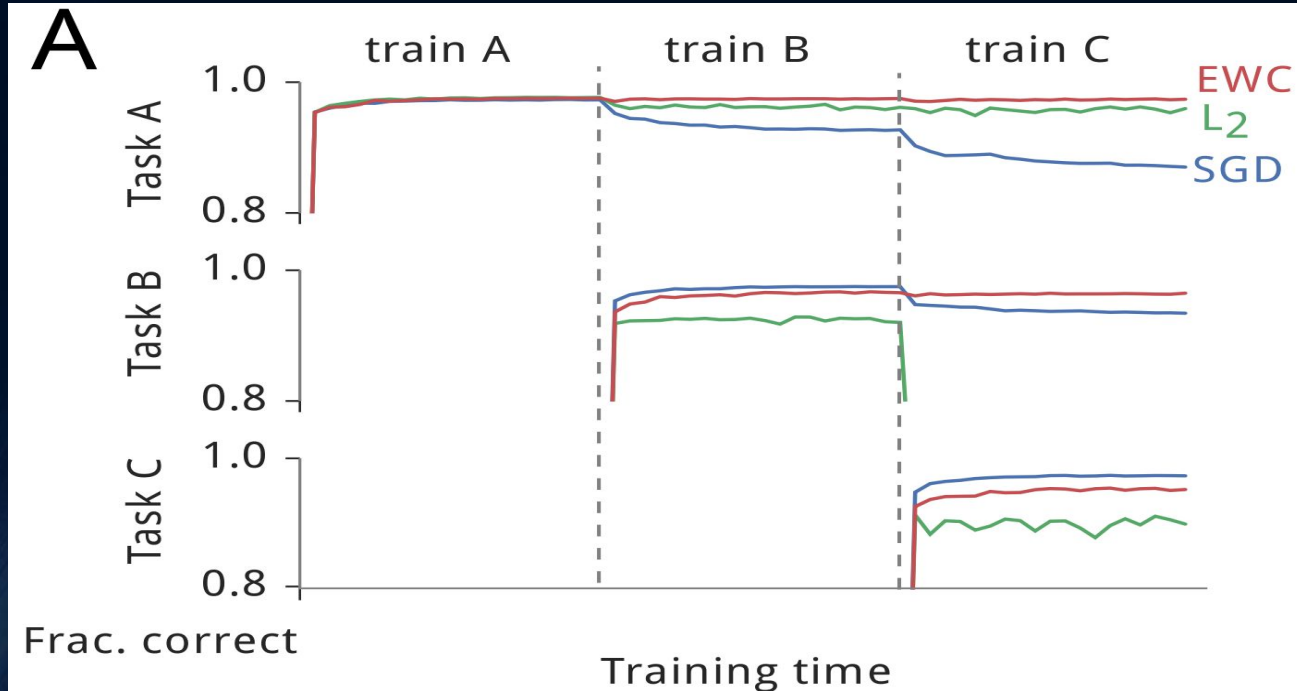
Data  $(\mathcal{X}^{(t)}, \mathcal{Y}^{(t)})$  is randomly drawn from distribution  $D^{(t)}$ , with  $\mathcal{X}^{(t)}$  a set of data samples for task  $t$ , and  $\mathcal{Y}^{(t)}$  the corresponding ground truth labels. The goal is to control the statistical risk of all seen tasks given limited or no access to data  $(\mathcal{X}^{(t)}, \mathcal{Y}^{(t)})$  from previous tasks  $t < \mathcal{T}$ :

$$\sum_{t=1}^{\mathcal{T}} \mathbb{E}_{(\mathcal{X}^{(t)}, \mathcal{Y}^{(t)})} [\ell(f_t(\mathcal{X}^{(t)}; \theta), \mathcal{Y}^{(t)})]$$



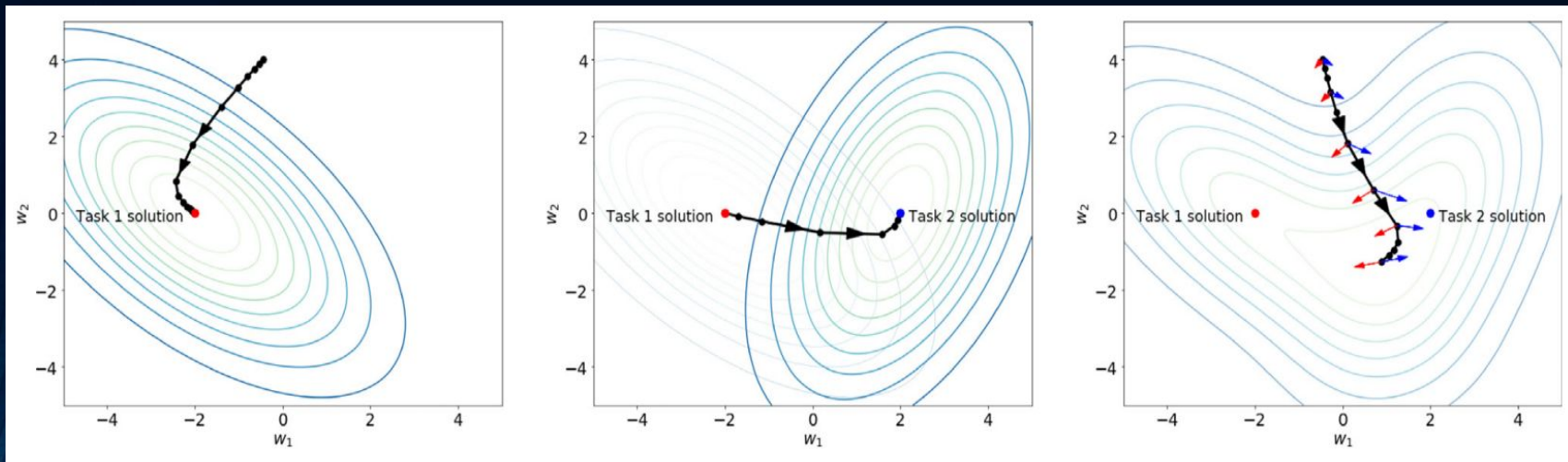
However, we have no access to the previous tasks, and thus cannot compute this empirical risk **exactly**.

# The Curse of Catastrophic Forgetting



Standard learning methods such as SGD quickly “forget” old knowledge (parameters) when the data/task change (i.e., adapt too well).

# Multi-Task Gradient Dynamics: Tug-of-War



Loss(Task1)

Loss(Task2)

Loss(Task1) + Loss(Task2)

However, the tasks are **not available simultaneously** in CL!

Need to use some form of memory, or to modify the gradients, to still take into account what solutions are good for previous tasks

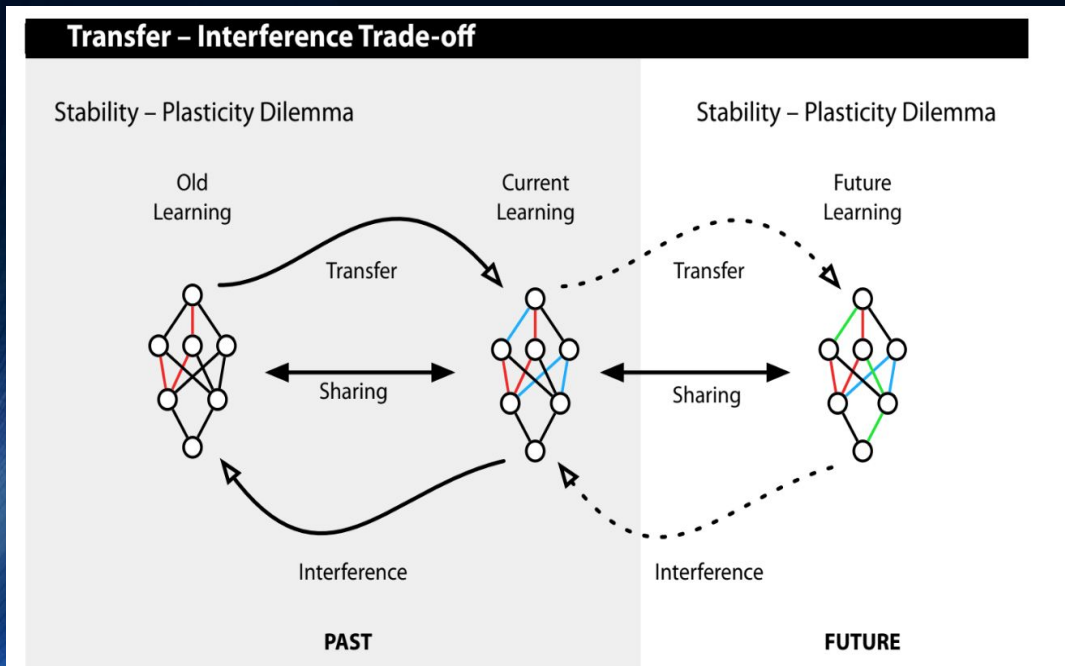


# Transfer vs Interference

**Transfer** from task A to task B  $\Leftrightarrow$  improved performance on B after learning A

**Interference** = negative transfer (i.e., decrease in performance)

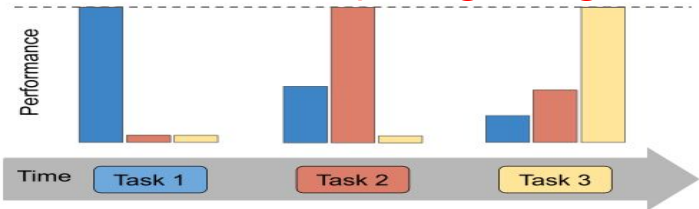
Weight-sharing can cause both  $\Leftrightarrow$  finding a good trade-off is the key!



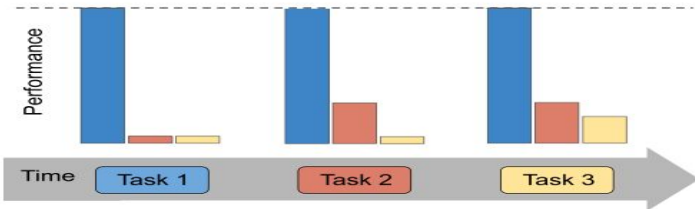
The reason for catastrophic forgetting is “the very thing — **a single set of shared weights** — that gave the networks their remarkable abilities to generalize and degrade gracefully.”  
(French, 1999)

# Possible Scenarios in CL

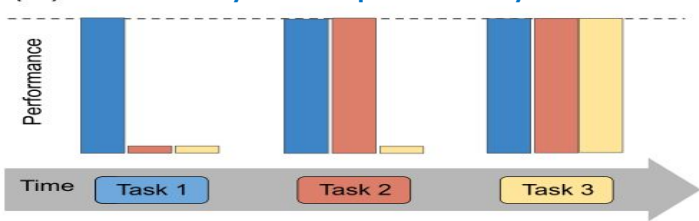
(A) Lack of stability (forgetting)



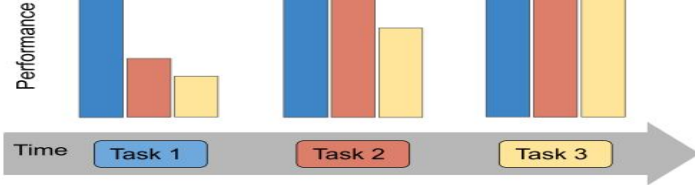
(B) Lack of plasticity



(C) Stability and plasticity



(D) Stability + (positive) forward transfer

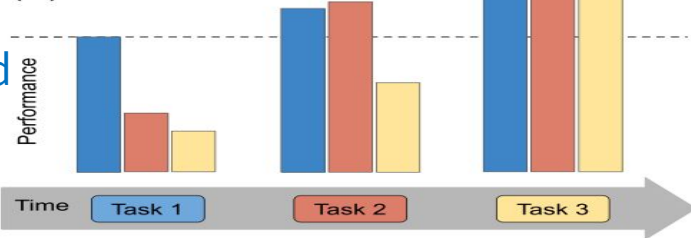


Bad

Good

Both (positive) backward and forward transfer

(E)



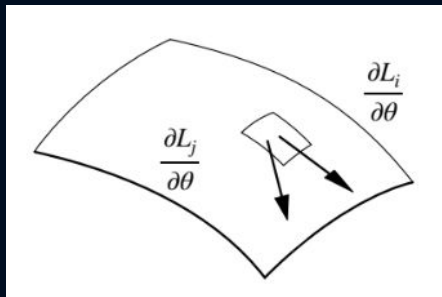
# Transfer vs Interference (“Negative Transfer”)

Riemer et al (2019) Learning to Learn without Forgetting By Maximizing Transfer and Minimizing Interference

## Transfer:

$$\frac{\partial L(x_i, y_i)}{\partial \theta} \cdot \frac{\partial L(x_j, y_j)}{\partial \theta} > 0.$$

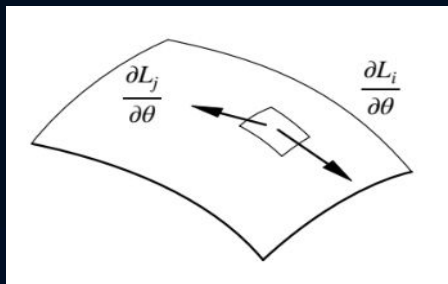
- When we train on one we improve on the other (a form of generalization)
- Analogous to positive transfer in either the forward or backward direction



## Interference:

$$\frac{\partial L(x_i, y_i)}{\partial \theta} \cdot \frac{\partial L(x_j, y_j)}{\partial \theta} < 0.$$

- When we train on one we get worse at the other
- Analogous to negative transfer in either the forward or backward direction



## Continual Learning Methods

```
graph TD; A[Continual Learning Methods] --> B[Replay methods]; A --> C[Regularization-based methods]; A --> D[Parameter isolation methods]
```

Replay  
methods

Regularization-based  
methods

Parameter isolation  
methods

### Memory-based

Maintain a subset of (representative) samples from previous tasks (raw samples or pseudo-samples – e.g., use a generative model).

Replay: reuse these samples as additional inputs in the future, or use them to constrain the new task loss.

### No (explicit) memory

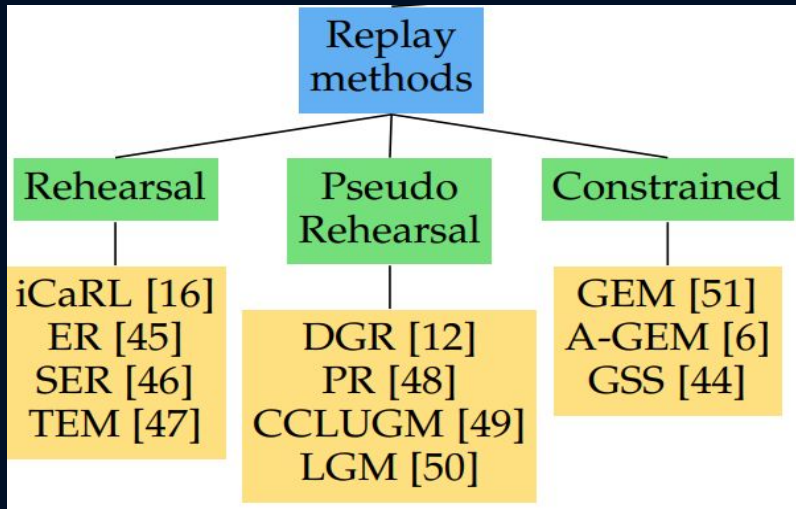
Instead, add regularization term to the loss function, consolidating previous knowledge when learning on new data.

Better for privacy, minimizing memory size and access latency.

### Architecture-based

Each task will use different model parameters (no forgetting!). Can use a static architecture or keep expanding it.

How will the overall performance scale with the model size?



## Constrained optimization:

Minimize interference with old tasks by constraining updates on the new task.

E.g., GEM, in task-incremental setting, projects the estimated gradient direction on the feasible region determined by previous task gradients, etc. More recent work (A-GEM, MER, etc).

## Rehearsal methods:

retrain current model on a subset of stored samples jointly with new tasks (e.g., reservoir sampling [45] )

## Pseudo rehearsal methods:

Feed random input to previous models, use the output as a pseudo-sample [48]. Generative models are also used but add training complexity.

# Example: Replay + Constraints (GEM)

Lopez-Paz and Ranzato (2017). Gradient episodic memory for Continual Learning.

- Idea:
- (1) store small amount of data per task in **memory**
  - (2) when making updates for new tasks, ensure that they don't **unlearn** previous tasks

How do we accomplish (2)?

learning predictor  $y_t = f_\theta(x_t, z_t)$       memory:  $\mathcal{M}_k$  for task  $z_k$

For  $t = 0, \dots, T$

minimize  $\mathcal{L}(f_\theta(\cdot, z_t), (x_t, y_t))$

subject to  $\mathcal{L}(f_\theta, \mathcal{M}_k) \leq \mathcal{L}(f_\theta^{t-1}, \mathcal{M}_k)$  for all  $z_k < z_t$  (i.e. s.t. loss on previous tasks doesn't get worse)

Assume local  
linearity:

$$\langle g_t, g_k \rangle := \left\langle \frac{\partial \mathcal{L}(f_\theta, (x_t, y_t))}{\partial \theta}, \frac{\partial \mathcal{L}(f_\theta, \mathcal{M}_k)}{\partial \theta} \right\rangle \geq 0 \quad \text{for all } z_k < z_t$$

Can formulate & solve as a QP.

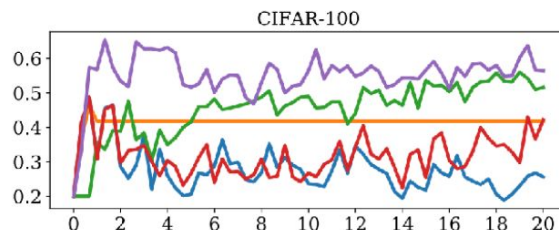
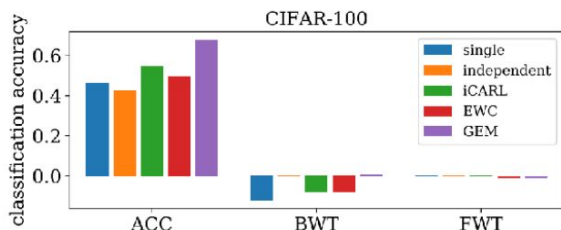
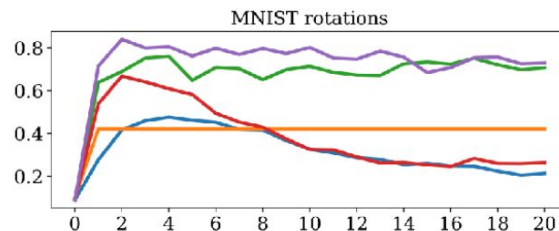
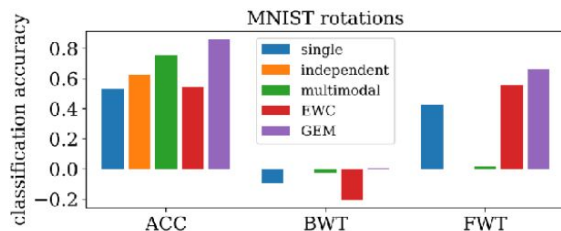
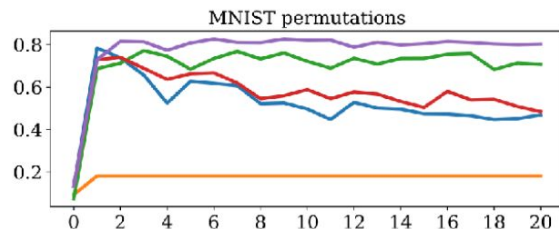
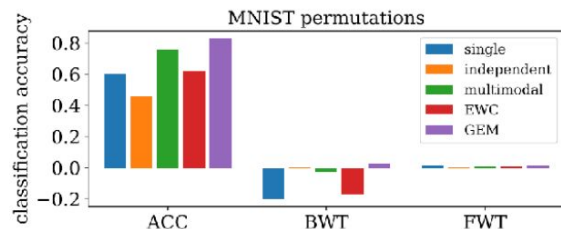
# Experiments

Problems:

- MNIST permutations
- MNIST rotations
- CIFAR-100 (5 new classes/task)

BWT: backward transfer,  
FWT: forward transfer

Total memory size:  
5012 examples



# Meta-Experience Replay (MER) Approach

Riemer et al (2019) Learning to Learn without Forgetting By Maximizing Transfer and Minimizing Interference

Standard offline training objective with dataset D:

$$\theta = \arg \min_{\theta} \mathbb{E}_{(x,y) \in D} [L(x, y)],$$

Modifying it to also learn to maximize transfer and minimize interference in either direction:

$$\theta = \arg \min_{\theta} \mathbb{E}_{(x_i, y_i) \& (x_j, y_j) \in D} [L(x_i, y_i) + L(x_j, y_j) - \alpha \frac{\partial L(x_i, y_i)}{\partial \theta} \cdot \frac{\partial L(x_j, y_j)}{\partial \theta}],$$

*Meta-learning perspective:* we would like to learn to learn each example in a way that generalizes to other examples from the overall distribution.



# Replay + Meta-Learning: Meta-Experience Replay

Reptile [1] is an efficient meta-learning algorithm that approximates the same objective as MAML.

Reptile can be extended to continual learning by integrating with ER! 😊

- ✓ Results from [1] still hold to the extent that our buffer captures the full variation of the distribution of examples seen.
- ✓ We can separate an ER batch into SGD steps over individual examples and apply a Reptile parameter meta-update.
- ✓ We also note that it is important to prioritize the current example before moving on as it may not be added to M.

**Approximate Objective (s batches with k examples each):**

$$\theta = \arg \min_{\theta} \mathbb{E}_{(x_{11}, y_{11}), \dots, (x_{sk}, y_{sk}) \in M} \left[ 2 \sum_{i=1}^s \sum_{j=1}^k [L(x_{ij}, y_{ij}) - \sum_{q=1}^{i-1} \sum_{r=1}^{j-1} \alpha \frac{\partial L(x_{ij}, y_{ij})}{\partial \theta} \cdot \frac{\partial L(x_{qr}, y_{qr})}{\partial \theta}] \right].$$

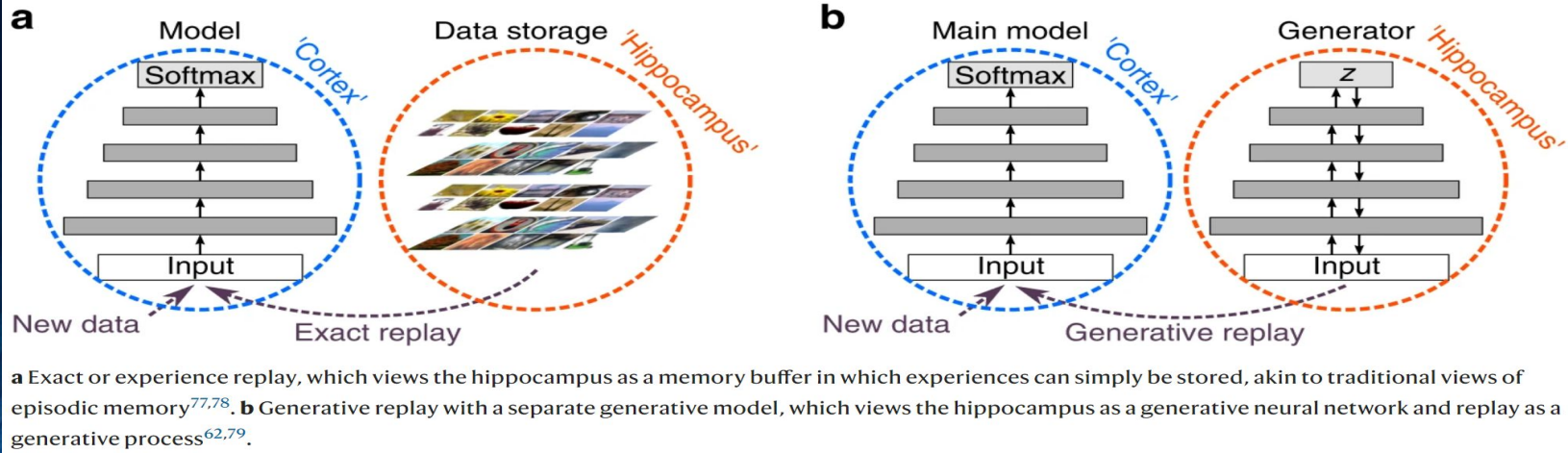
**Retained Accuracy After Training on All Tasks**

Model	Buffer Size	MNIST Rotations	MNIST Permutations	Many Permutations	Incremental Omniglot
Online	N/A	46.40 ± 0.78	55.42 ± 0.65	32.62 ± 0.43	4.36 ± 0.37
EwC	N/A	57.96 ± 1.33	62.32 ± 1.34	33.46 ± 0.46	4.63 ± 0.14
GEM	5120	87.12 ± 0.44	82.50 ± 0.42	56.76 ± 0.29	18.03 ± 0.15
	500	72.08 ± 1.29	69.26 ± 0.66	32.14 ± 0.50	-
	200	66.88 ± 0.72	55.42 ± 1.10	-	-
MER	5120	<b>89.56 ± 0.11</b>	<b>85.50 ± 0.16</b>	<b>61.84 ± 0.25</b>	<b>75.23 ± 0.52</b>
	500	<b>81.82 ± 0.52</b>	<b>77.40 ± 0.38</b>	<b>47.40 ± 0.35</b>	<b>32.05 ± 0.69</b>
	200	<b>77.24 ± 0.47</b>	<b>72.74 ± 0.46</b>	-	-

# Episodic (exact) vs Generative Replay

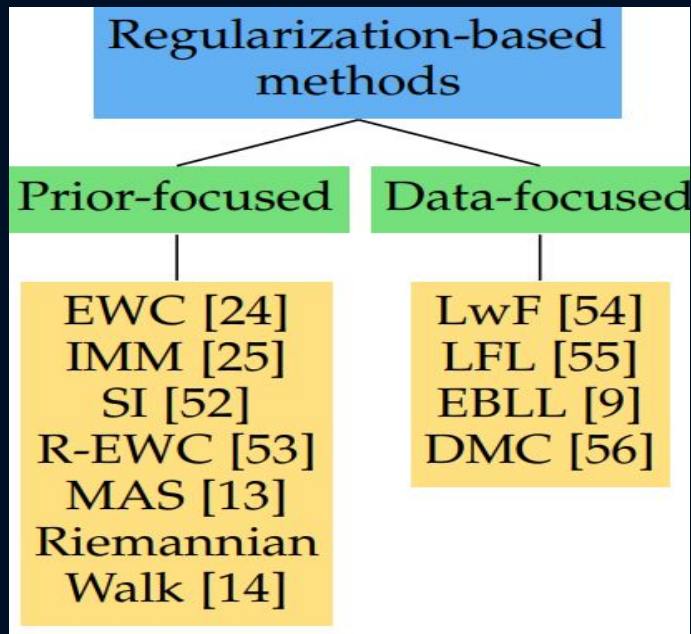
van de Ven et al (2020). Brain-inspired replay for continual learning with artificial neural networks

From: Brain-inspired replay for continual learning with artificial neural networks



Novel GR method: internal or hidden representations are replayed that are generated by the network's own, context-modulated feedback connections.

SOTA performance on challenging CL benchmarks with many tasks ( $\geq 100$ ) or complex inputs (natural images) without storing data



Add regularization term to the loss function, consolidating previous knowledge when learning on new data.

## Data-focused:

Knowledge distillation from a previous model (trained on a previous task) to the model being trained on the new data.

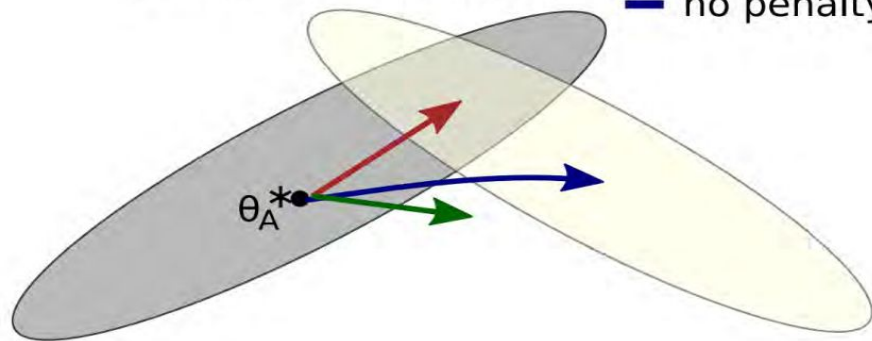
## Prior-focused:

Use an estimated distribution over the model parameters as prior when learning from new data; penalize changes to parameters important for the past tasks (e.g. EWC and later work).

# Elastic Weight Consolidation (EWC)

[Kirkpatrick et al, PNAS 2017]

- Low error for task B
- Low error for task A
- EWC
- L<sub>2</sub>
- no penalty



Idea: Don't let *important* parameters change drastically (reduce plasticity)

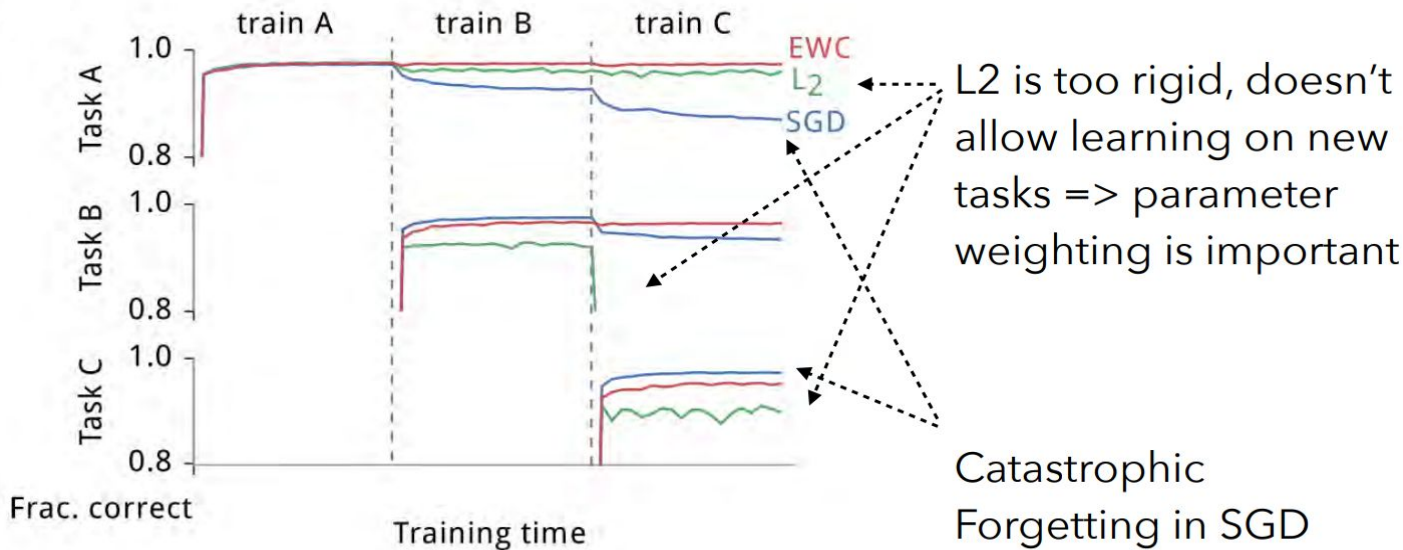
$$\mathcal{L}(\theta) = \mathcal{L}_B(\theta) + \sum_i \frac{\lambda}{2} F_i (\theta_i - \theta_{A,i}^*)^2,$$

Task B Loss

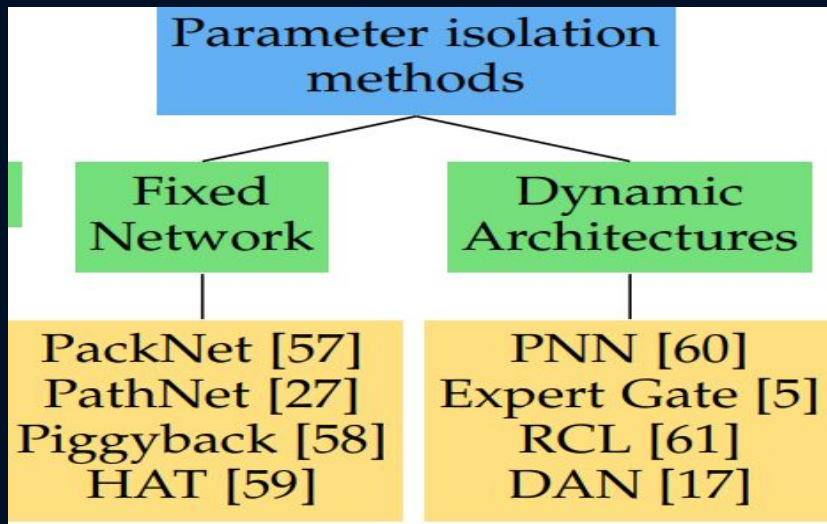
- Inspired by research on synaptic consolidation

# Elastic Weight Consolidation (EWC)

[Kirkpatrick et al., PNAS 2017]



MNIST experiments. New tasks are random pixel permutations.



Idea: avoid forgetting by using different parameters for each task

Best-suited for: task-incremental setting, unconstrained model capacity, performance is the priority.

## Fixed Network Methods:

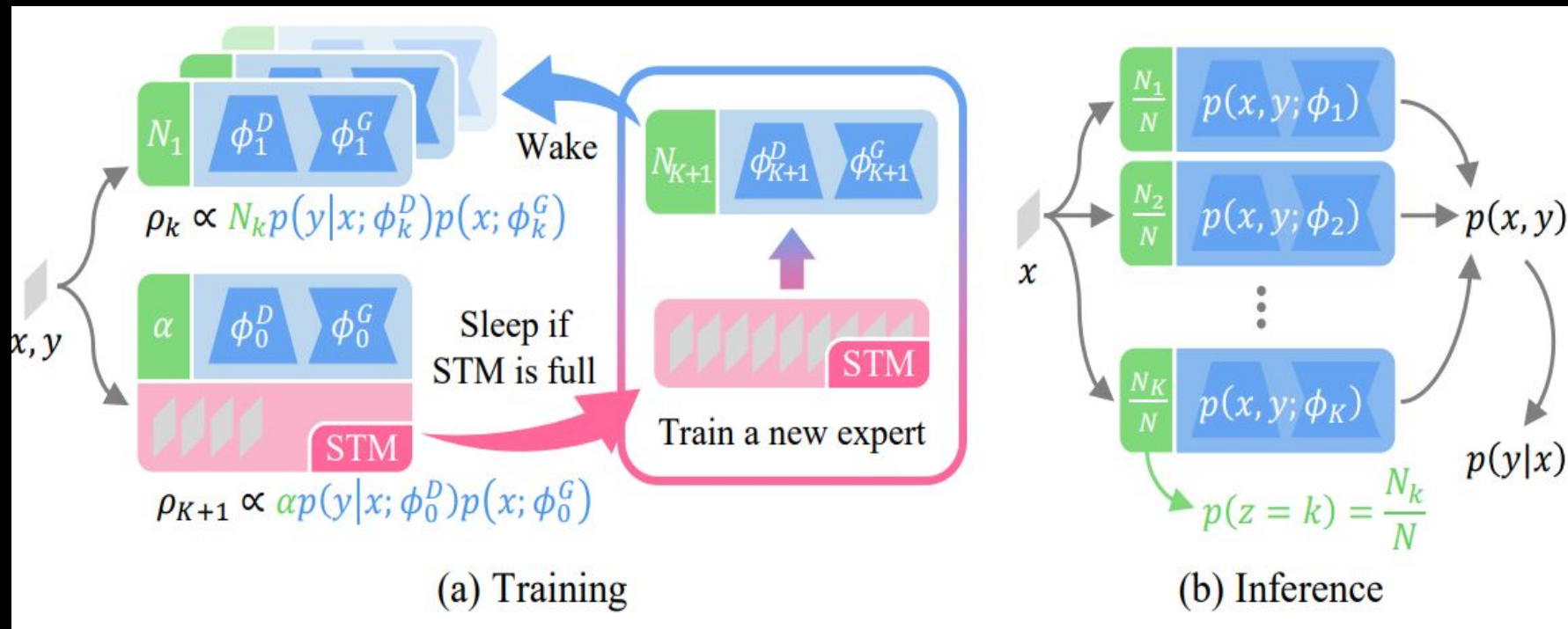
Network parts used for previous tasks are masked out when learning new tasks (e.g., at neuronal level (HAT) or at parameter level (PackNet, PathNet))

## Dynamic Architecture Methods:

When model size is not constrained: grow new branches for new tasks, while freezing previous task parameters (RCL), or dedicate a model copy to each task (Expert Gate), etc.

# Example: Architectural Approaches

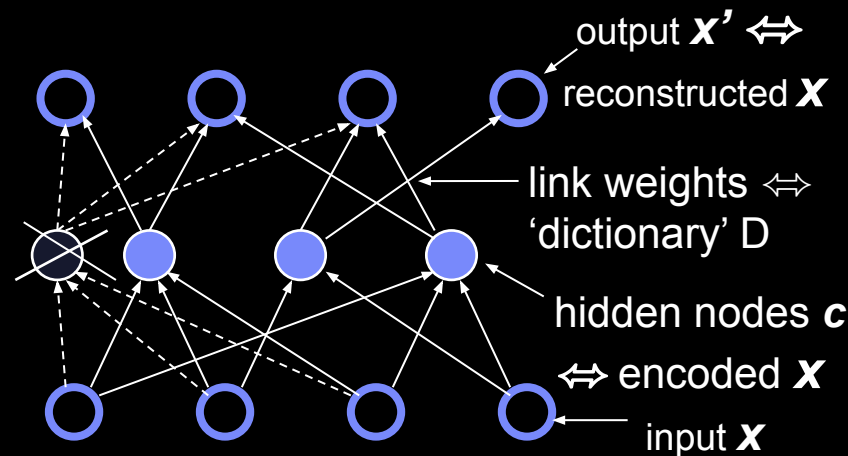
Lee et al (2020) A Neural Dirichlet Process Mixture Model for Task-Free Continual Learning. ICLR2020



# Neurogenetic Autoencoder (Online Dictionary Learner)

Garg et al (2017) Neurogenesis-inspired dictionary learning.

- A sparse autoencoder model (a.k.a. dictionary learning)
- Neuronal “birth”: adding random hidden units
- Neuronal “death”: using  $l_1/l_2$  (group sparsity) regularize



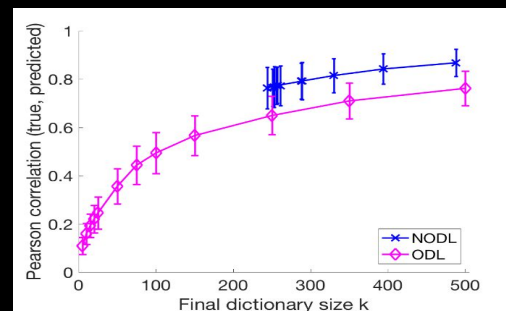
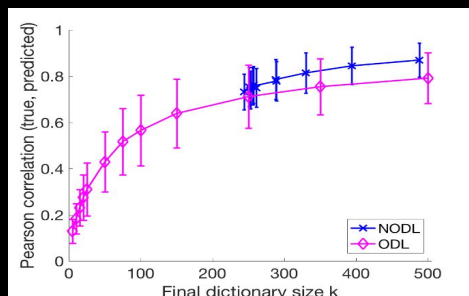
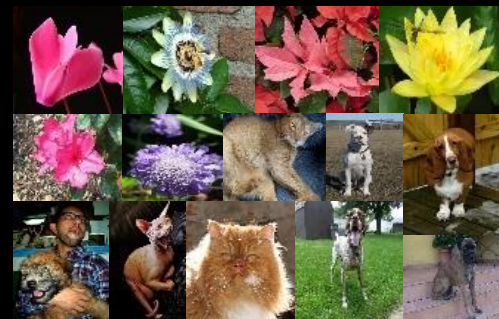
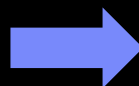
$$\hat{f}_t(D) = \underbrace{\frac{1}{t} \sum_{i=1}^t \frac{1}{2} \|\mathbf{x}_i - D\boldsymbol{\alpha}_i\|_2^2}_{\text{reconstruction error}} + \underbrace{\lambda_c \|\boldsymbol{\alpha}_i\|_1}_{\text{sparsity on codings}} + \underbrace{\lambda_g \sum_j \|\mathbf{d}_j\|_2}_{L_1/L_2 \text{ group sparsity}} + \underbrace{\sum_j \lambda_j \|\mathbf{d}_j\|_1}_{\text{sparse elements}}$$



# “Neurogenesis” Helps CL

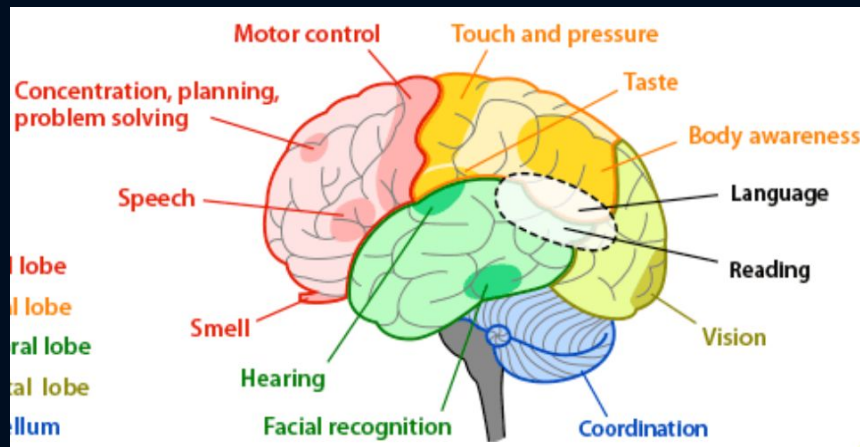
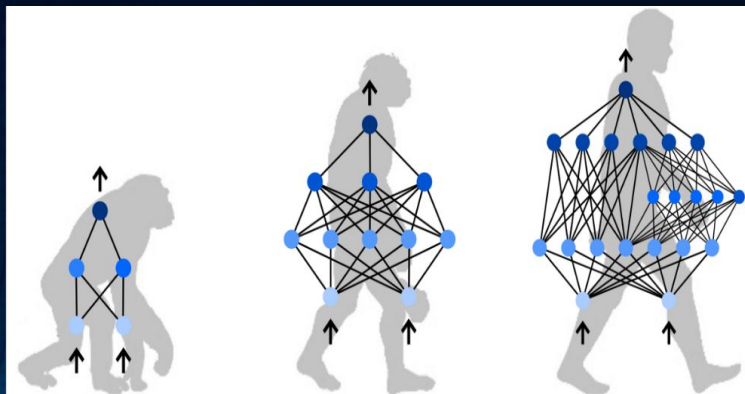
Garg et al (2017) Neurogenesis-inspired dictionary learning.

From urban  
 (“Oxford”)  
to nature  
 (flowers,  
 dogs, cats)  
 images



Neurogenetic Online Dictionary Learner (NODL) improves reconstruction accuracy over standard ODL on BOTH old and new data (i.e. avoids forgetting while adapting), and learns more compact representations.

# Role of CL: Evolving a “Library” of “Basis Functions” ?




Network components  $\Leftrightarrow$  finite functional “basis”

$$\{ h_1(x), \dots, h_k(x) \}$$

$$f_1(x) \quad f_2(x) \quad f_3(x) \quad \dots \quad f_n(x) \quad \dots$$

Infinite stream of changing environments and tasks

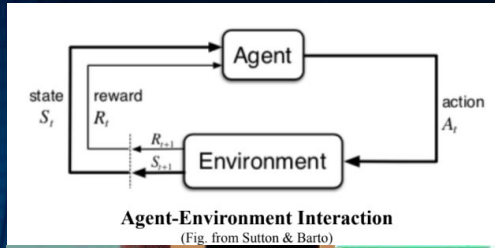


# Continual Reinforcement Learning

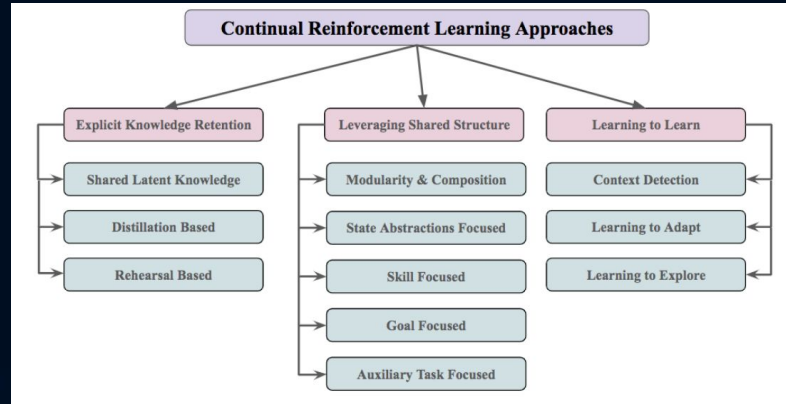
# Continual Reinforcement Learning

Khimya Khetarpal\*, Matthew Riemer\*, Irina Rish, Doina Precup (2020).

Towards Continual Reinforcement Learning: A Review and Perspectives.

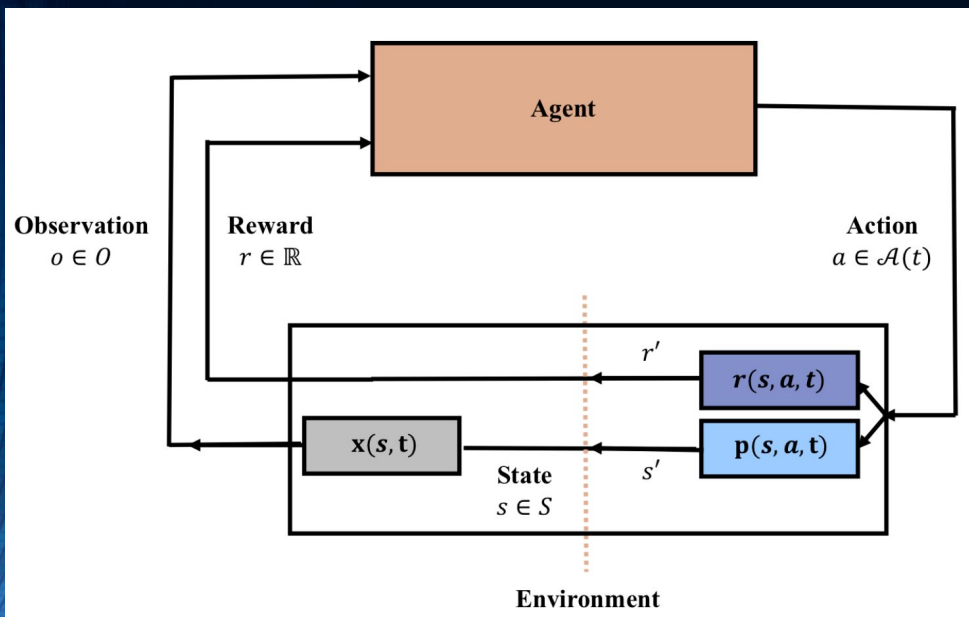


Procgen: A benchmark for procedurally generated set of environments to measure generalization.

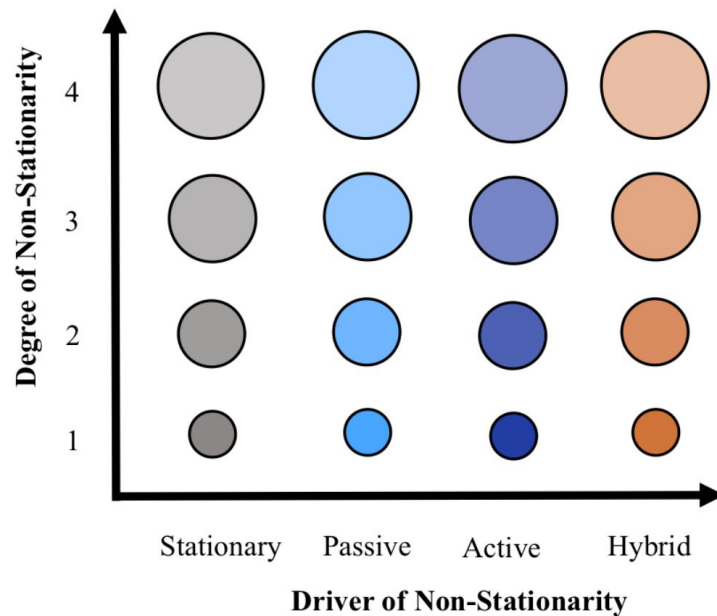


Bsuite: is a collection of carefully-designed experiments that investigate core capabilities of a reinforcement learning (RL) agent.

# Continual Reinforcement Learning



## A Taxonomy of Continual RL Formalisms



# Drivers of Nonstationarity

**Passive Non-stationarity:** In passive non-stationary environments, we assume that the non-stationary behavior (i.e. the evolution of tasks) does not depend on the behavior of the agent itself when interacting with the environment.

- E.g. *“Hidden-Mode Markov Decision Processes for Nonstationary Sequential Decision Making”*

- The evolution of tasks depends on a stochastic function  $P(z'|z)$  as in without having to consider the effects of our own changing policy on this distribution

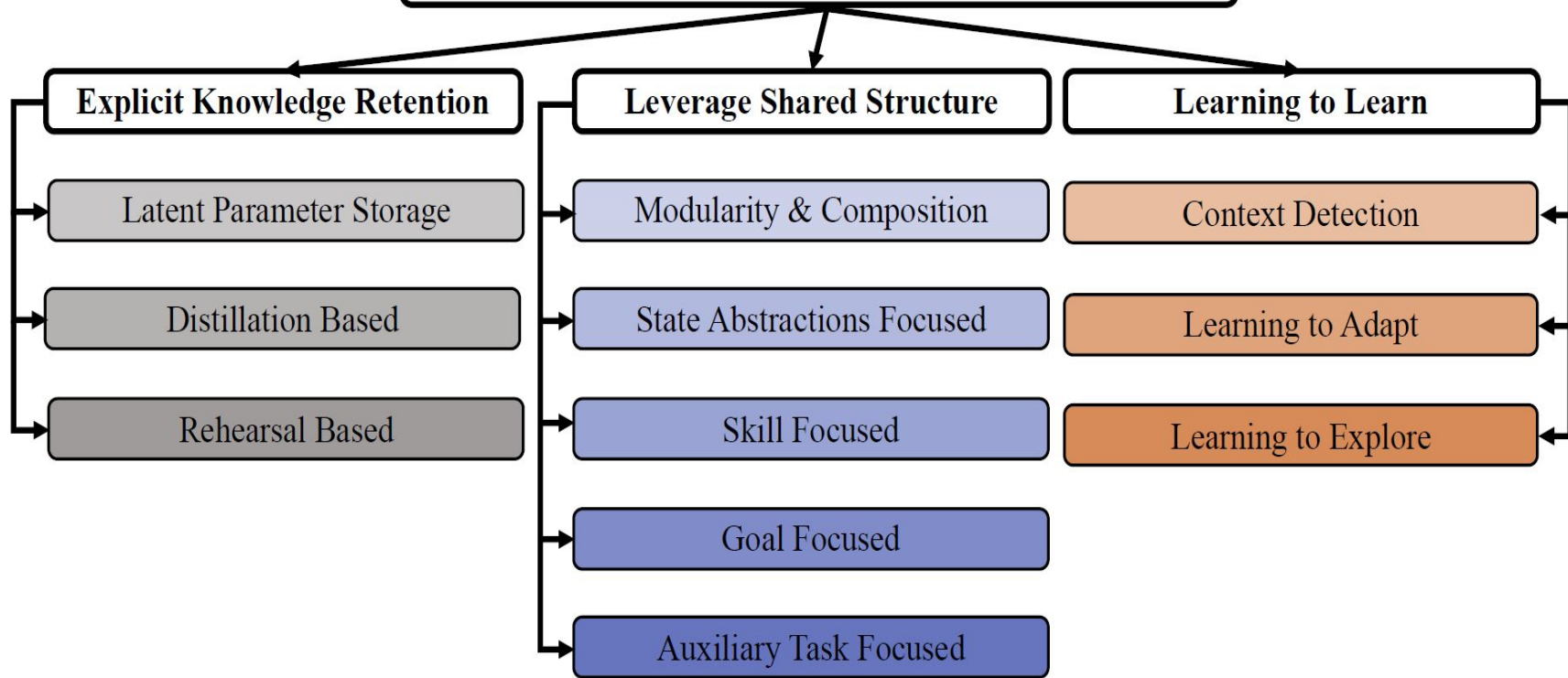
**Active Non-stationarity:** In active non-stationary environments, we consider that the agent's behavior may have an impact on the nature of the non-stationarity in the environment.

- Eg: Intrinsic motivation, curriculum learning

- This setting is foundational to work studying intrinsic motivation or the agent setting its own curriculum.

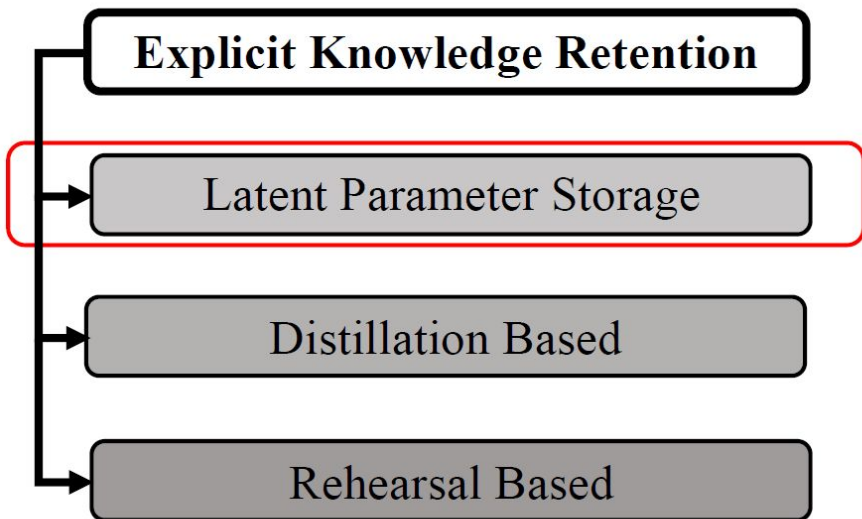
**Hybrid Non-stationarity:** Combining both active and passive sources of non-stationarity.

# Continual Reinforcement Learning Approaches



# Explicit Knowledge Retention

## □ Latent Parameter Storage



*Shared components:* Ammar et al., 2014 : shared latent basis that captures reusable components

Borsa et al., 2016: explicitly model a shared abstraction of state-action space for multi-task setting

*Prior representations:* Rusu et al., 2016: provide representations of networks trained on previous tasks as inputs for subsequent tasks.

Kirkpatrick et al., 2017: store a prior about the extent of past usage of each parameter during learning to preserve important old knowledge

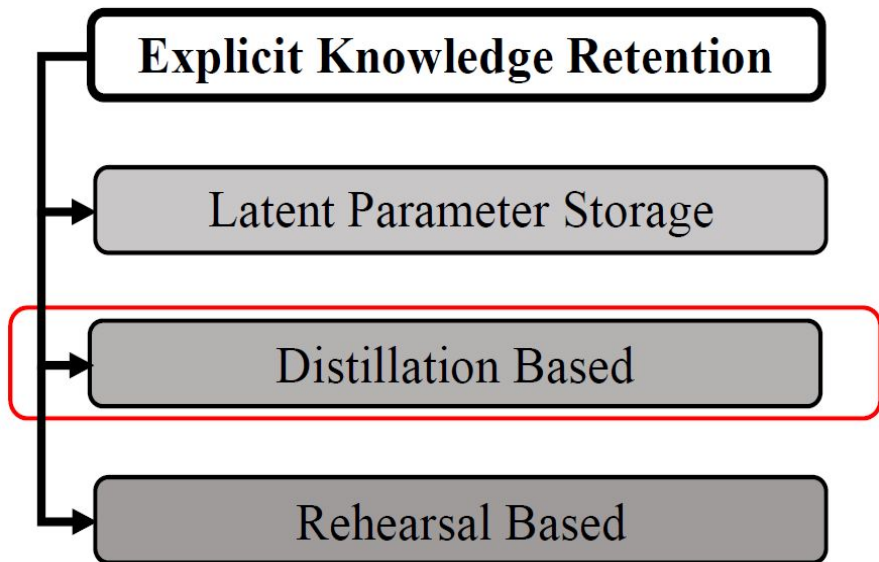
*Single shared representation:* Maurer et al. 2016, extracted features for multiple tasks in a single low-dim shared representation

D'Eramo et al. 2016 derived theoretical bounds for AVI and API showing that learning a shared representation significantly decreases the error propagation.



# Explicit Knowledge Retention

## □ Distillation Based



*Idea:* Leveraging knowledge distillation which dates back to Bucilua et al., 2006 and renewed interest and success in Hinton et al. 2015

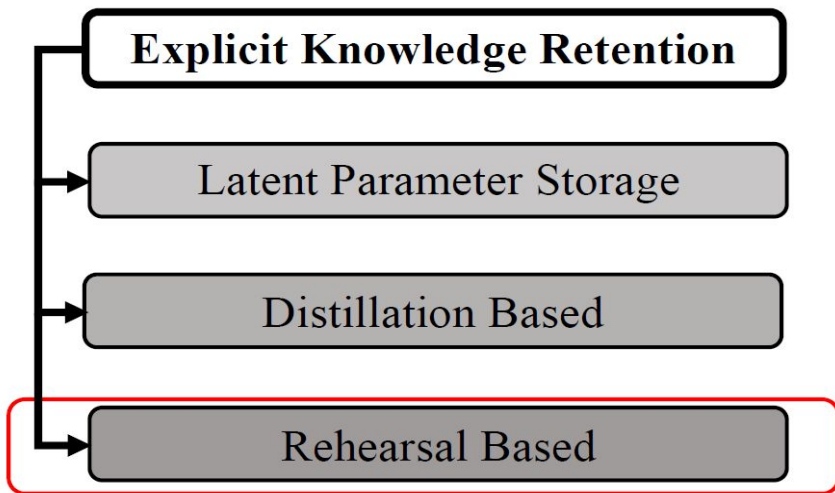
Distill knowledge from past policies when learning a new task as in Rusu et al., 2015; Li and Hoiem, 2016; Riemer et al., 2016; Espeholt et al., 2018; Schwarz et al., 2018; Berseth et al., 2018; Kaplanis et al., 2019; Traoré et al., 2019; Tirumala et al., 2019.

An additional benefit: ultimately learns a separate model for each task.

However, this results in a need for a knowledge compression strategy to scale truly to many many tasks.

# Explicit Knowledge Retention

## ☐ Rehearsal Based



*Idea:* Reinforce the importance of experiences from the past distribution using experience replay (Lin, 1992)

Replay experiences can help **correct the bias in our objective function** towards the short term to the extent that the past is a good proxy for the future. A very successful approach for tackling continual RL as shown in (Isele and Cosgun, 2018; Riemer et al., 2019; Rolnick et al., 2019).

Replay might result in significant storage requirements, and it is not always clear how to prioritize data in replay (the length, the recency, rewards).

Besides they struggle to effectively leverage past data if the shift in distribution is drastic.

*Idea:* replace replay with **pseudo-rehearsals** sampled from a trained generative model of the environment.

Explored in Robins, 1995, Atkinson et al., 2018

# On Evaluation of Continual RL Agents: Benchmarks

A desired CRL Benchmark should allow for

- a range of degrees of non-stationarity (from 1 to 4)
- training in progressive fashion
- discovery and composition of skills
- generate more complex tasks/scenarios with increasing difficulty
- learning casual relationships including affordances associated with objects
- embodied agents
- rich parallel streams of data which are multi modal

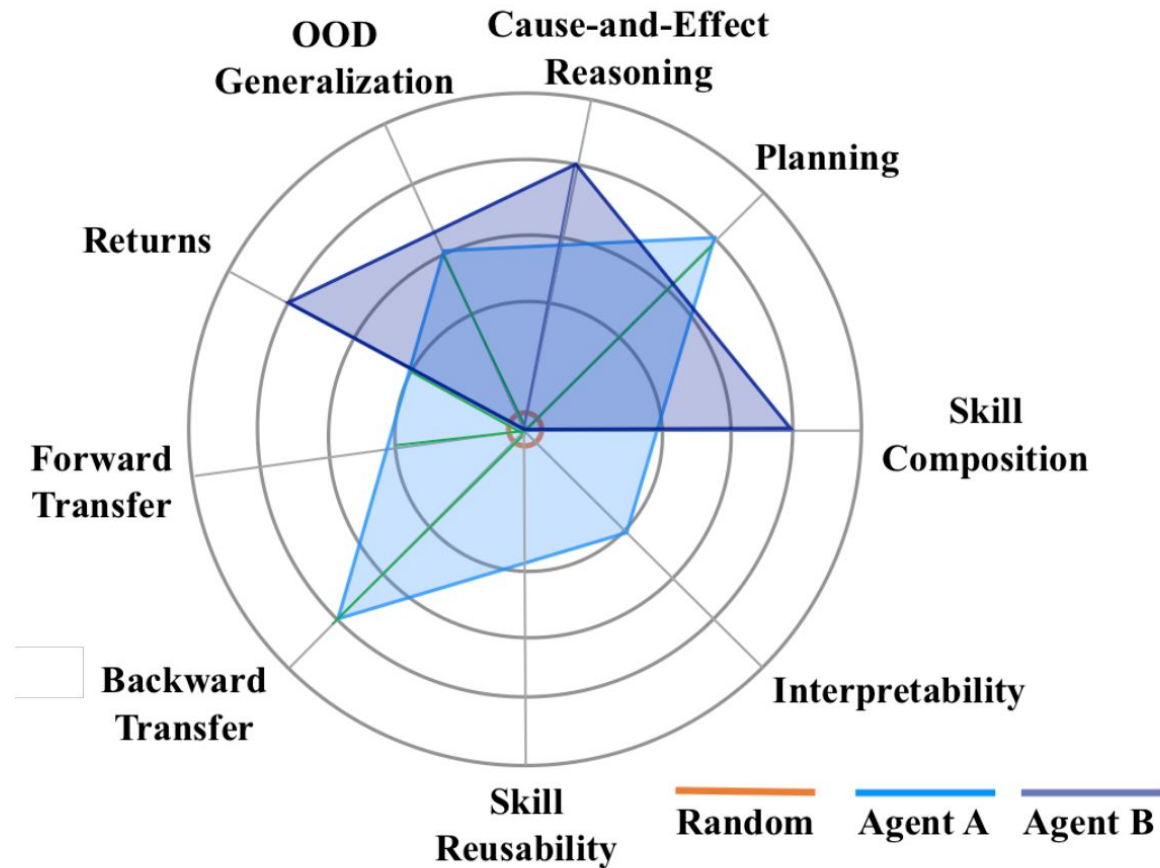


Progen: A benchmark for procedurally generated set of environments to measure generalization.

# Open Problems and Challenges

- ❑ Finding the right inductive biases
- ❑ Task specification and formalism
- ❑ Understanding the agent-environment boundary
- ❑ Experimental design and evaluation i.e. training and testing
- ❑ Interpreting discovered behaviors
- ❑ Learning at scale – scaling laws for Continual RL?

# Key Metrics for Continual RL



# CL and Neural Scaling Laws?

<https://youtu.be/V8FEFw50lg4>

## Explaining Neural Scaling Laws

Physics n ML 6.16.2021

Ethan Dyer

Based on 2102.06701 w/ Yasaman Bahri, Jared Kaplan, Jaehoon Lee, Utkarsh Sharma

BLUESHIFT >>

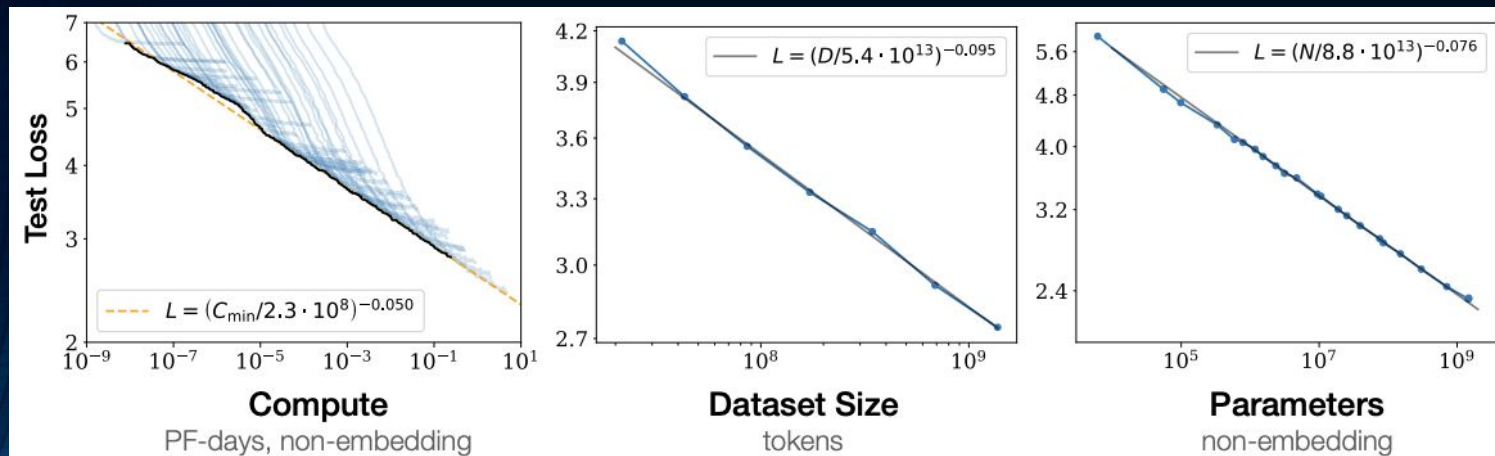


0:26 / 1:17:51



# Neural Scaling Laws

Kaplan et al (2020). Scaling laws for neural language models.



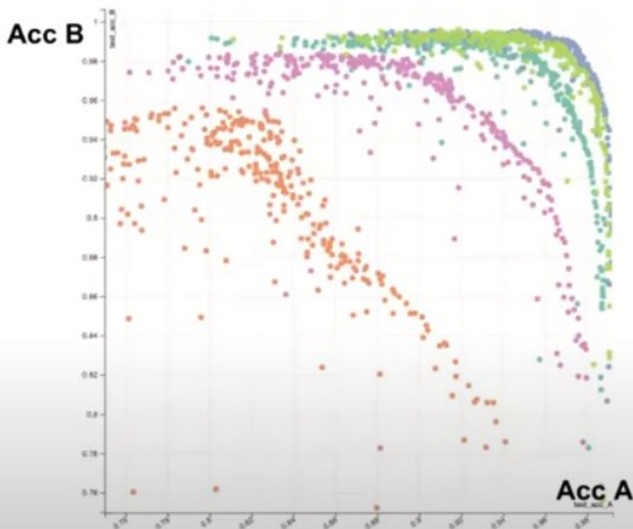
**Figure 1** Language modeling performance improves smoothly as we increase the model size, dataset size, and amount of compute<sup>2</sup> used for training. For optimal performance all three factors must be scaled up in tandem. Empirical performance has a power-law relationship with each individual factor when not bottlenecked by the other two.

What aspects of models, algorithms and data have the largest effect on scaling laws?  
How can we design learning approaches with better scaling (e.g., scaling exponent)?

# Can Scaling Solve Catastrophic Forgetting?



What is solved by scale, what is not?



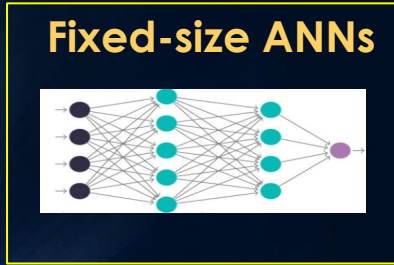
ResNet pre-trained on ImageNet 21k: 26 50 101 152 200



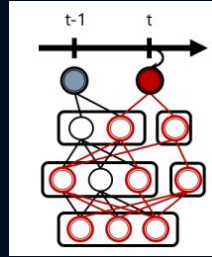


# Environment Complexity vs Model Capacity

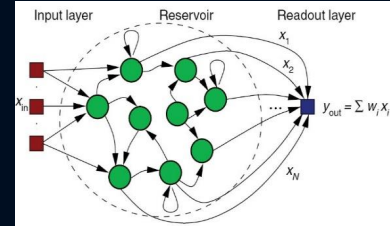
Representation  
Capacity



**Expanding ANNs:**



**3rd-generation ANNs?**



Space-time coding?

Environment Complexity

# Summary: Desirable Properties of CL Systems

- Constant memory (infinite data/task stream)
- No task boundary info
- Online learning (without offline training on large batches/tasks)
- Forward transfer (e.g., OOD generalization)
- Backward transfer (beyond not forgetting)
- Problem agnostic (e.g., not limited to classification)
- Adaptively learning from any partial data (e.g., semi-supervised)
- No test time oracle
- Task revisiting to strengthen prior knowledge
- Graceful forgetting (compression) to balance stability and plasticity

# Recent Surveys on Continual Learning

Hadsell et al. (2020) Embracing Change: Continual Learning in Deep Neural Networks.

Khetarpal et al. (2020) Towards Continual Reinforcement Learning: A Review and Perspectives.

Mundt et al. (2020) A wholistic view of continual learning with deep neural networks: Forgotten lessons and the bridge to active and open world learning.

De Lange et al. (2019) Continual learning: A comparative study on how to defy forgetting in classification tasks.

Parisi et al. (2019) Continual lifelong learning with neural networks: A review.

Chen & Liu (2018). Lifelong Machine Learning.

Soltoggio et al. (2017) Born to learn: the inspiration, progress, and future of evolved plastic artificial neural networks.

Thank you!

The background is a dark blue gradient. It features several vertical white lines of varying lengths. Scattered throughout are small squares in three colors: light blue, light orange, and light pink. Some squares are solid, while others are hollow outlines. The text '10 minutes Break' is centered in the middle of the image.

10 minutes  
Break

# Future Directions and Open Challenges

Continual Learning with Deep Architectures  
Tutorial @ ICML 2021 - Part 2

Vincenzo Lomonaco  
University of Pisa & ContinualAI  
*vincenzo.lomonaco@unipi.it*

# Vincenzo Lomonaco

*Assistant Professor @  
University of Pisa*

*Co-founding President and Lab  
Director @ [ContinualAI.org](https://ContinualAI.org)*

*Co-founder & Board Member @  
[AlforPeople.org](https://AlforPeople.org)*





Focus on Deep  
Neural Networks



# TABLE OF CONTENTS



01

Continual  
Unsupervised  
Learning



02

Continual  
Learning  
Applications



03

Impact on  
Sustainable AI



04

Open  
Questions



05

Conclusion



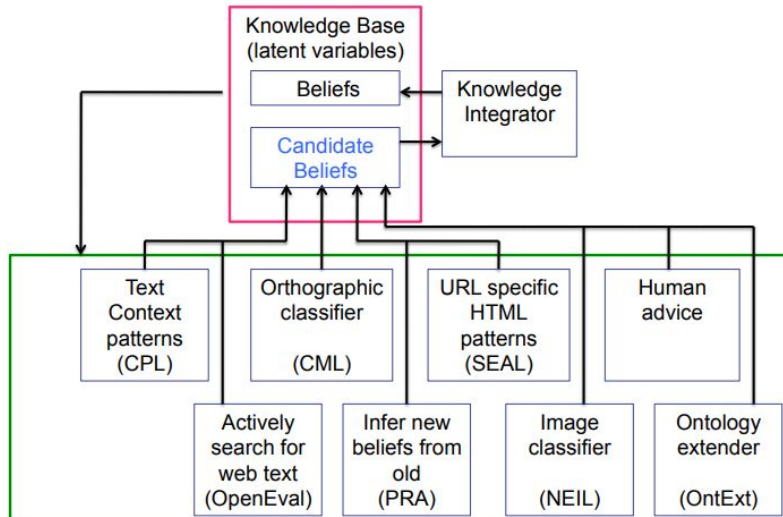
# Continual Unsupervised Learning

# NELL: Never-Ending Language Learning

## Key Ideas

- *Semi-Supervised Learning System*
- Ran 24x7, from January, 2010 to September 2018
- Combination of many learning algorithms (CPL, CML, SEAL, OpenEval, PRA, NEIL)
- Intended as a case-study for a **never-ending agent**

## NELL Architecture



# NELL: Never-Ending Language Learning

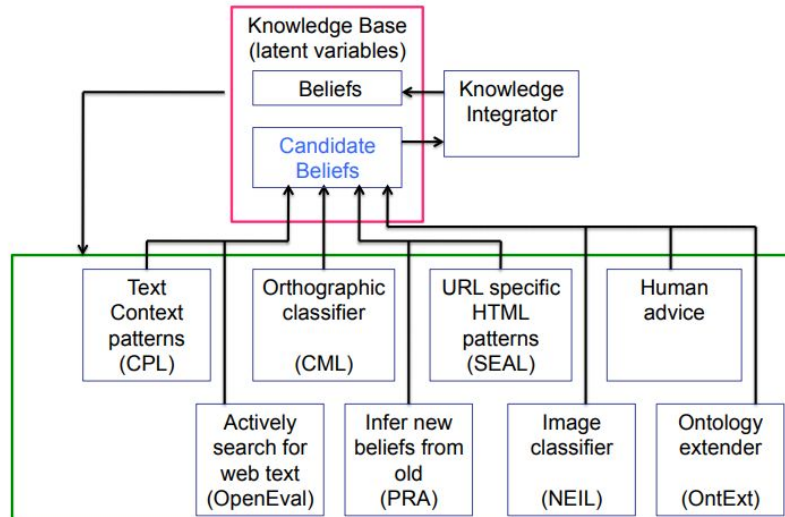
## Key Ideas

- *Semi-Supervised Learning System*

In his 2019 book "Human Compatible", Stuart Russell commented that "Unfortunately NELL has confidence in only 3% of its beliefs and relies on human experts to clean out false or meaningless beliefs on a regular basis."

- Intended as a case-study for a **never-ending agent**

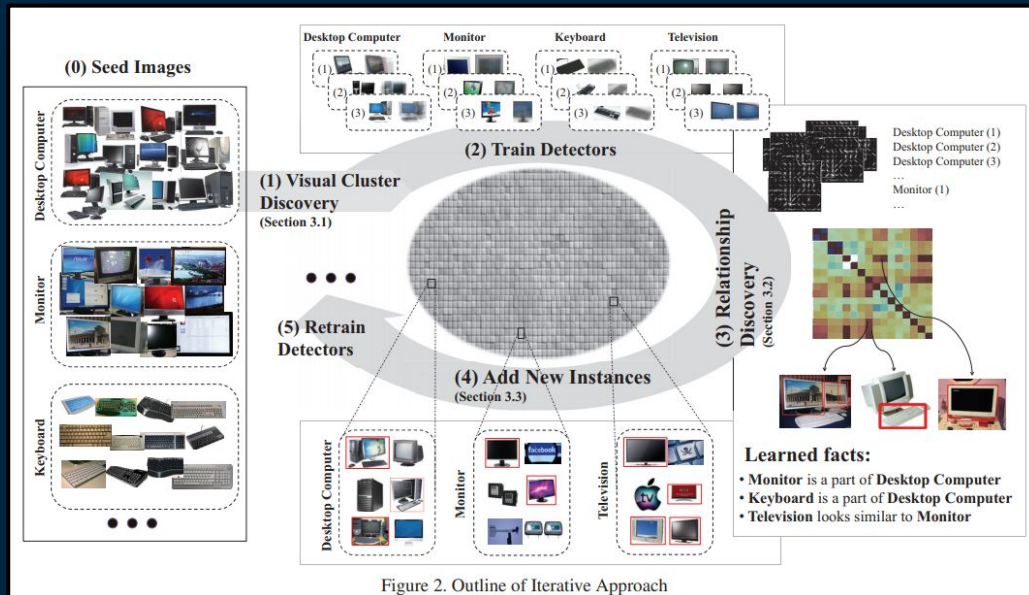
## NELL Architecture



# NEIL: Extracting Visual Knowledge from Web Data

## Key Ideas

- *Semi-Supervised Learning System*
- **Cumulative approach:** not incremental, SVM as main learning algorithm
- Feature: GIST, SIFT, HOG, Lab color space, and Texton
- 2.5 months on 200 core cluster: 16 iterations, 400K self-labeled instances, 1152 object categories, 1034 scene categories



# Lifelong Topic Modeling

## Key Ideas

- *Traditionally an unsupervised learning task.*
- **The "topics" produced by topic modeling techniques are clusters of similar words.**
- Set of shared words among some topics generated from multiple domains are **more likely to be coherent** for a particular topic.
- **Focus:** knowledge accumulation rather than learning an incremental function

---

### Algorithm 1 PriorTopicsGeneration( $D$ )

---

```
1: for  $r = 0$  to  $R$  do
2:   for each domain corpus  $D_i \in D$  do
3:     if  $r = 0$  then
4:        $S_i \leftarrow \text{LDA}(D_i)$ ;
5:     else
6:        $S_i \leftarrow \text{LTM}(D_i, S)$ ;
7:     end if
8:   end for
9:    $S \leftarrow \cup_i S_i$ ;
10: end for
```

---

---

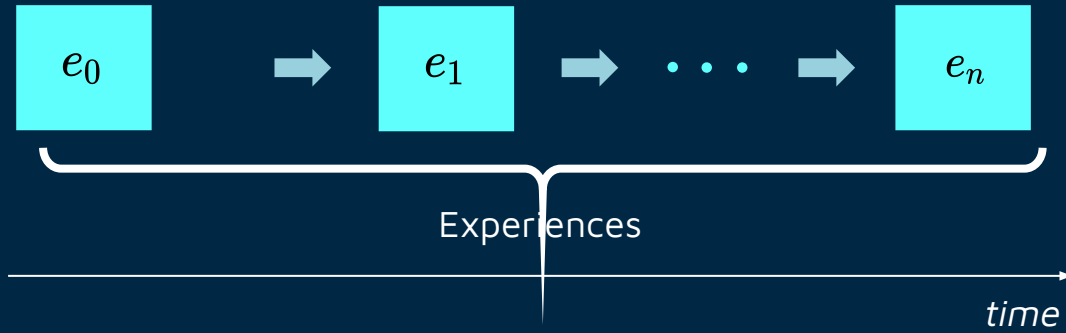
### Algorithm 2 LTM( $D^t, S$ )

---

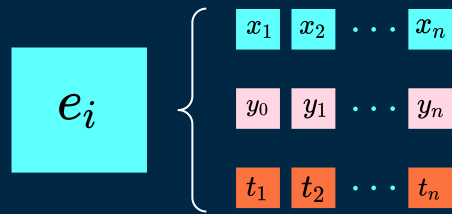
```
1:  $A^t \leftarrow \text{GibbsSampling}(D^t, \emptyset, N)$ ; // Run  $N$  Gibbs iterations with no knowledge (equivalent to LDA).
2: for  $i = 1$  to  $N$  do
3:    $K^t \leftarrow \text{KnowledgeMining}(A^t, S)$ ;
4:    $A^t \leftarrow \text{GibbsSampling}(D^t, K^t, 1)$ ; // Run with knowledge  $K^t$ .
5: end for
```

---

# Continual Unsupervised Learning

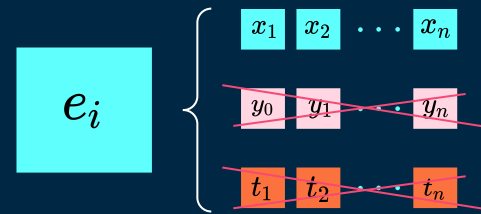


Supervised Learning



examples

Unsupervised Learning



examples

# Continual Unsupervised Learning

## Ideal Paradigm to Combine with CL

- **No Continual Labeling**
- **Less Bias**
- **Why this is still not the case?**
  - *Changing the paradigm:*  
More Data, Less Supervision
  - Less impactful applications (for now)

### ■ “Pure” Reinforcement Learning (cherry)

- ▶ The machine predicts a scalar reward given once in a while.
- ▶ **A few bits for some samples**

### ■ Supervised Learning (icing)

- ▶ The machine predicts a category or a few numbers for each input
- ▶ Predicting human-supplied data
- ▶ **10→10,000 bits per sample**

### ■ Unsupervised/Predictive Learning (cake)

- ▶ The machine predicts any part of its input for any observed part.
- ▶ Predicts future frames in videos
- ▶ **Millions of bits per sample**

■ (Yes, I know, this picture is slightly offensive to RL folks. But I’ll make it up)





# Continual Unsupervised Representation Learning

## Key Ideas

- **Fully Generative Approach**
- $\mathbf{y}$  can be interpreted as representing some discrete clusters in the data
- **Mixture of Gaussian with Dynamic Expansion**
- *Difficult to scale*: tested only on MNIST and Omniglot

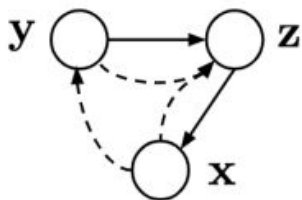


Figure 1: Graphical model for CURL. The categorical task variable  $\mathbf{y}$  is used to instantiate a latent mixture-of-Gaussians  $\mathbf{z}$ , which is then decoded to  $\mathbf{x}$ .

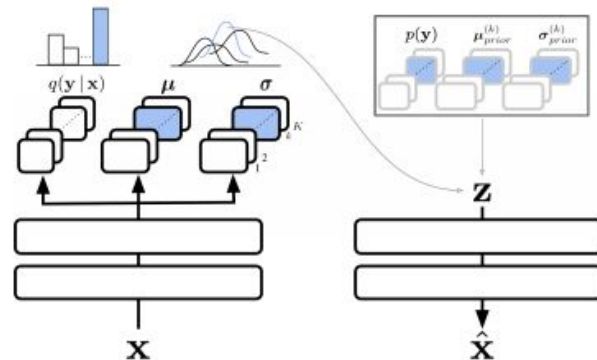


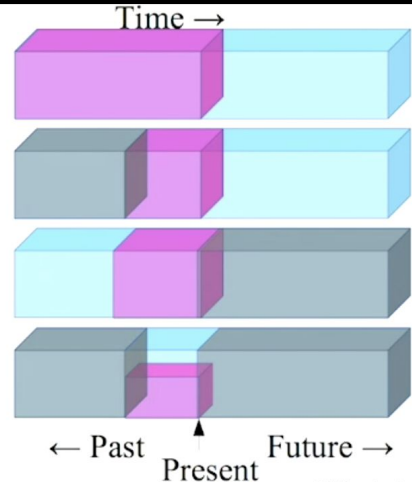
Figure 2: Diagram of the proposed approach, showing the inference procedure and architectural components used.

# Continual Unsupervised Learning

## Huge Exploration Opportunities

- **Self-Supervised Learning**
- Sequence Learning
- Contrastive Learning
- Hebbian-like Learning
- Active Learning
- Weakly/Semi-Supervised Learning
- Randomized Networks

- ▶ Predict any part of the input from any other part.
- ▶ Predict the **future** from the **past**.
- ▶ Predict the **future** from the **recent past**.
- ▶ Predict the **past** from the **present**.
- ▶ Predict the **top** from the **bottom**.
- ▶ Predict the **occluded** from the **visible**
- ▶ **Pretend there is a part of the input you don't know and predict that.**

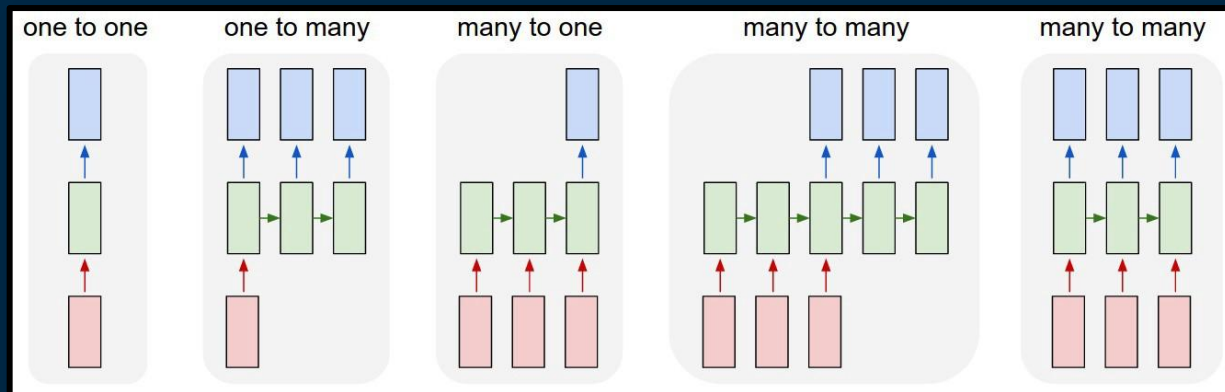


Slide: LeCun

# Continual Unsupervised Learning

## Huge Exploration Opportunities

- Self-Supervised Learning
- **Sequence Learning**
- Contrastive Learning
- Hebbian-like Learning
- Active Learning
- Weakly/Semi-Supervised Learning
- Randomized Networks



Y. Cui et al. *Continuous online sequence learning with an unsupervised neural network model*. *Neural Computation*, 2016.

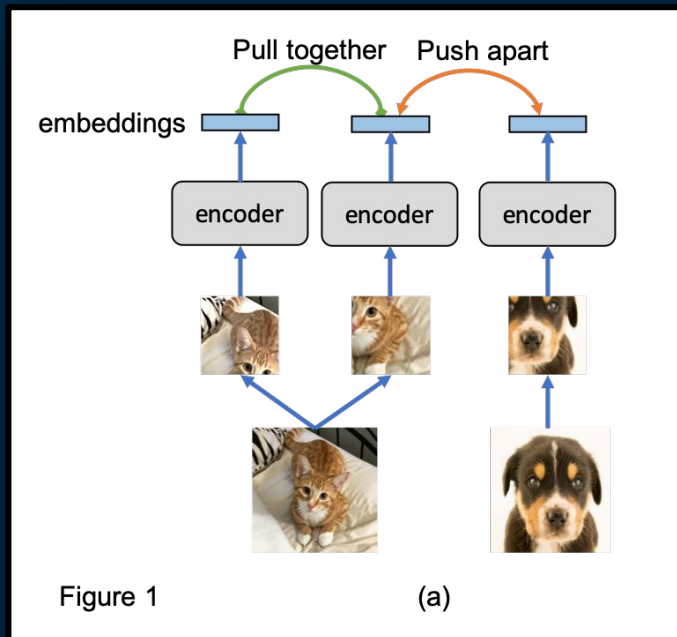
A. Cossu et al. *Continual Learning for Recurrent Neural Networks: an Empirical Evaluation*. Elsevier *Neural Networks*, 2021.

B. Ehret et al. *Continual learning in Recurrent Neural Networks*. ICLR 2021.

# Continual Unsupervised Learning

## Huge Exploration Opportunities

- Self-Supervised Learning
- Sequence Learning
- **Contrastive Learning**
- Hebbian-like Learning
- Active Learning
- Weakly/Semi-Supervised Learning
- Randomized Networks



Junnan Li, 2020

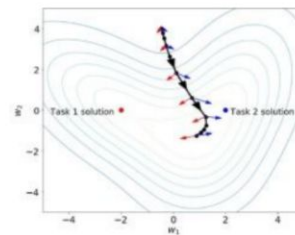
# Continual Unsupervised Learning

## Huge Exploration Opportunities

- Self-Supervised Learning
- Sequence Learning
- Contrastive Learning
- **Hebbian-like Learning**
- Active Learning
- Weakly/Semi-Supervised Learning
- Randomized Networks

## Gradient-based optimization and tug-of-war dynamics

- Continual Learning is a huge challenge for deep learning models because of **gradient-based optimization**.
  - Gradient-based learning is effective and cheap, the de rigeur method for training neural networks for close to 4 decades.
  - However, a close look at the learning dynamics reveals a problem.
  - Each training sample produces a gradient for each parameter in the network that votes to make the parameter bigger or smaller.
  - In a mini-batch, a gradient is produced by each sample in parallel and they are summed to decide the winning direction.
- ➔ The result is a **tug-of-war** over the direction of change of each parameter.

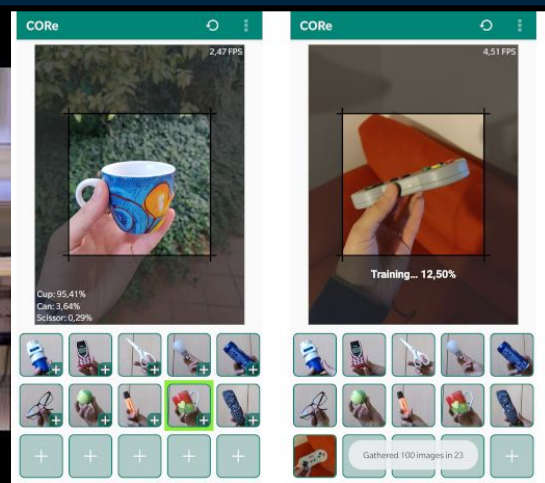


R. Pascanu, 2021

# Continual Unsupervised Learning

## Huge Exploration Opportunities

- Self-Supervised Learning
- Sequence Learning
- Contrastive Learning
- Hebbian-like Learning
- **Active Learning**
- Weakly/Semi-Supervised Learning
- Randomized Networks

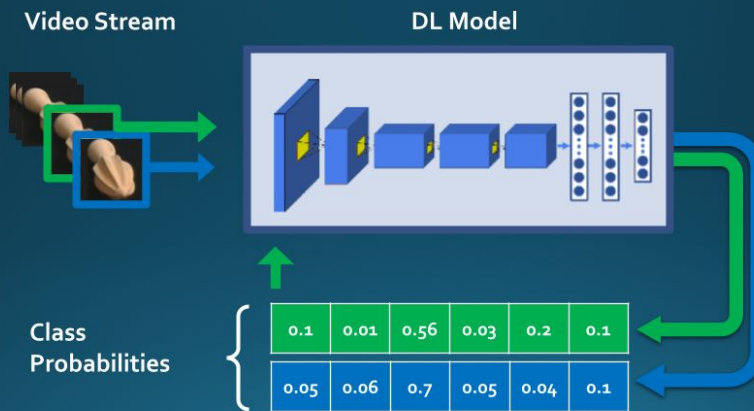


# Continual Unsupervised Learning

## Huge Exploration Opportunities

- Self-Supervised Learning
- Sequence Learning
- Contrastive Learning
- Hebbian-like Learning
- Active Learning
- **Weakly/Semi-Supervised Learning**
- Randomized Networks

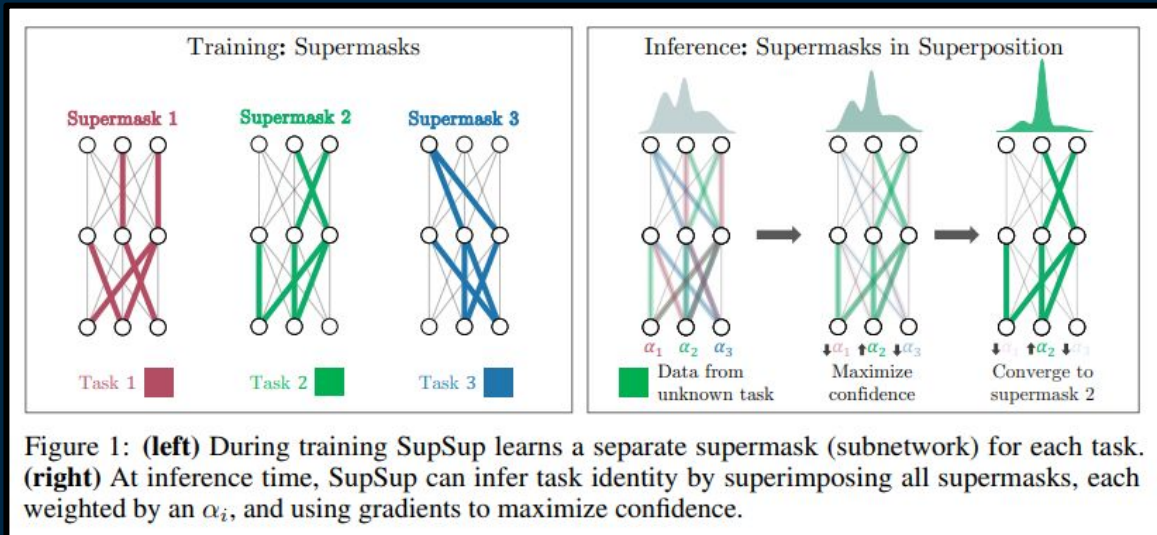
## Semi-Supervised Tuning from Temporal Coherence



# Continual Unsupervised Learning

## Huge Exploration Opportunities

- Self-Supervised Learning
- Sequence Learning
- Contrastive Learning
- Hebbian-like Learning
- Active Learning
- Weakly/Semi-Supervised Learning
- **Randomized Networks**



M. Wortsman, 2020



# Continual Unsupervised Learning

## Other relevant works in this area

- A. Bertugli et al. ***Few-Shot Unsupervised Continual Learning through Meta-Examples***. Workshop on Meta-Learning at NeurIPS 2020.
- I. Muñoz-Martín et al. ***Unsupervised learning to overcome catastrophic forgetting in neural networks***. IEEE Journal on Exploratory Solid-State Computational Devices and Circuits, 2019.
- L. Caccia et al. ***SPeCiaL: Self-Supervised Pretraining for Continual Learning***, arXiv 2021.
- W. Sun et al. ***ILCOC: An Incremental Learning Framework Based on Contrastive One-Class Classifiers***. CLVision Workshop at CVPR 2021.
- J. He et al. ***Unsupervised Continual Learning Via Pseudo Labels***. arXiv 2020.
- S. Khar et al. ***Unsupervised Class-Incremental Learning through Confusion***. arXiv 2021.



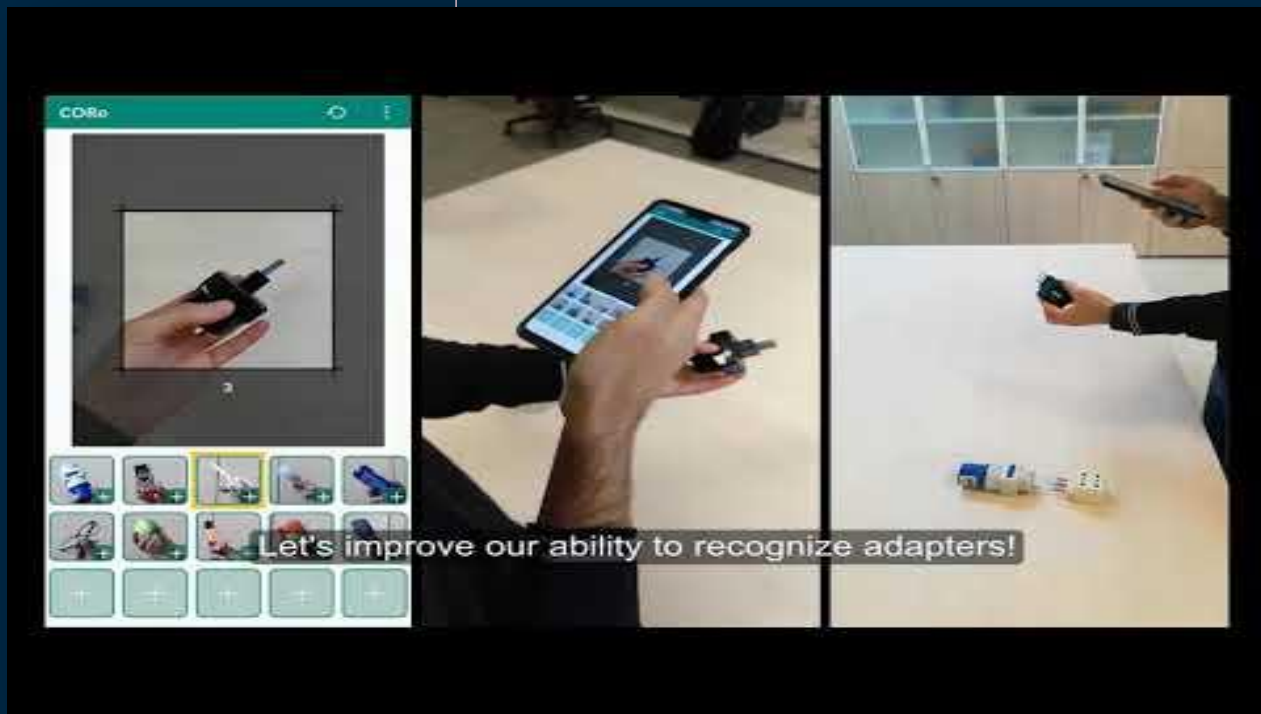
# Continual Learning Applications

# Continual Learning Applications

## Main Possibilities

- Edge
  - **Embedded systems and Robotics:** +privacy, +efficiency, +fast adaptation, +on the edge, -Internet connection (e.g. Autonomous Cars, Robotics Arms/Hands)
- Cloud
  - **AutoML and CI systems for AI models:** +scalability, +efficiency, +fast adaptation, -energy consumption, -\$\$\$ (e.g. Recommendation Systems)
- Continuum Edge-Cloud
  - **Pervasive AI systems:** Efficient Communication, fluid & dynamic computation
  - **Neural Patches:** +security patches, +fairness patches, +fast update
  - **Continual Distributed Learning:** understudied relationship with parallel and federated learning

# On-Device Personalization without Forgetting



L. Pellegrini et al. *Latent Replay for Real-Time Continual Learning*, IROS 2020.

L. Pellegrini et al. *Continual Learning at the Edge: Real-Time Training on Smartphone Devices*. ESANN, 2021.

G. Demosthenous et al. *Continual Learning on the Edge with TensorFlow Lite*. arXiv 2021.

L. Ravaglia et al. *Memory-Latency-Accuracy Trade-offs for Continual Learning on a RISC-V Extreme-Edge Node*. SiPS 2020.

# AR1: a Flexible Hybrid Strategy for CL

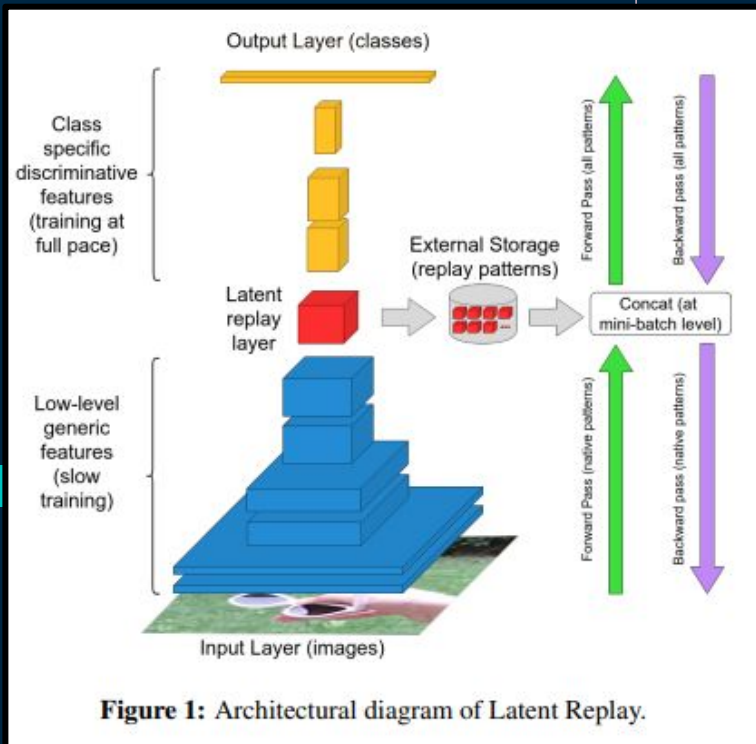


Figure 1: Architectural diagram of Latent Replay.

## Key Ideas

- **Architectural, Regularization and Replay components:**
  - *CWR\** for the output layer (arch)
  - *Online Synaptic Intelligence* (reg)
  - *Latent Replay* (replay)

$$\tilde{L}_\mu = L_\mu + \lambda \sum_k \Omega_k^\mu (\bar{\theta}_k - \theta_k)^2$$

$$w_k^v = \int_{t^{\mu-1}}^{t^\mu} \frac{\partial L}{\partial \theta_k} \cdot \frac{\partial \theta_k}{\partial t}$$

D. Maltoni et al. *Continuous Learning in Single-Incremental-Task Scenarios*, Neural Networks, 2019.

L. Pellegrini et al. *Latent Replay for Real-Time Continual Learning*, IROS 2020.

V. Lomonaco et al. *Rehearsal-Free Continual Learning over small I.I.D Batches*. CLVision at CVPR 2020.

# Continual Learning in Production

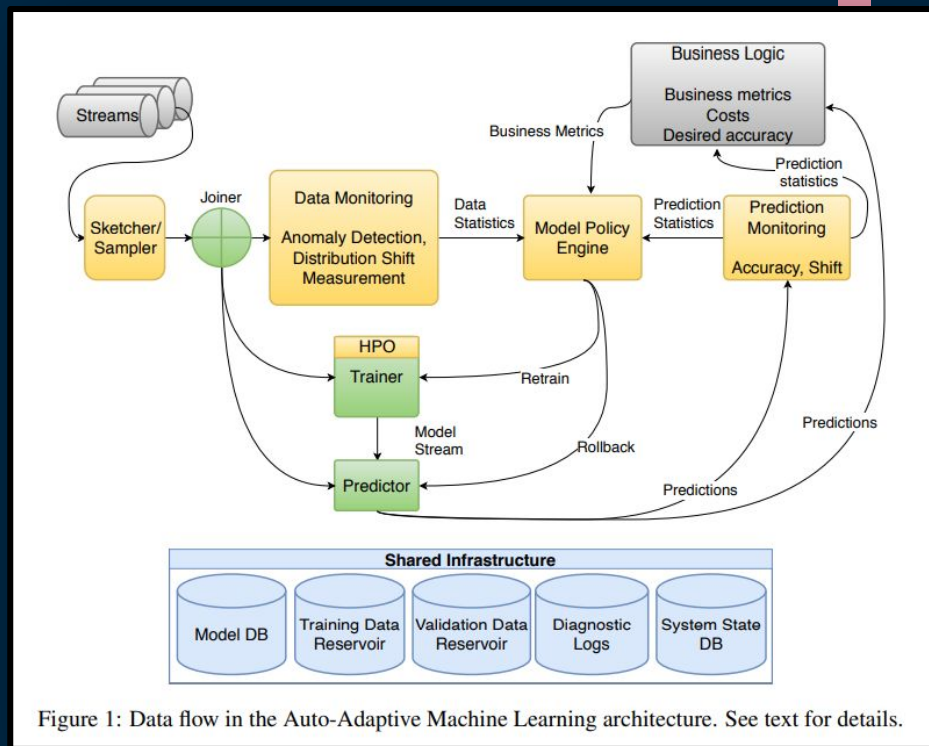


Figure 1: Data flow in the Auto-Adaptive Machine Learning architecture. See text for details.

# Use-Cases: Google Play and Tesla

The Anatomy of a Production-Scale Continuously-Trained Model

Watch later Share

### Warm Starting

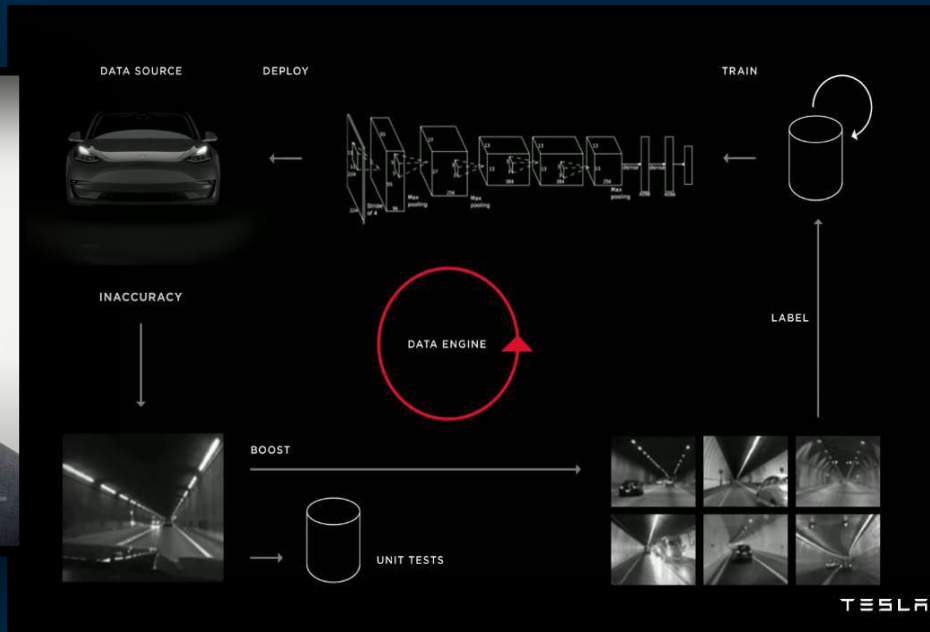
Time t-1 Time t

MORE VIDEOS

2:08 / 4:02

YouTube

The diagram illustrates the concept of 'Warm Starting' in a neural network. It shows two sequential neural network architectures, one at Time t-1 and one at Time t. The network at Time t-1 has a set of weights (represented by blue circles) that are passed to the network at Time t, allowing it to start with a pre-trained state. The video player interface includes a play button, a volume icon, and a progress bar.



- V. Lomonaco. **Continual Learning for Production Systems: The new "Agile" in the Machine Learning Era.** ContinualAI Publication, 2019.
- D. Baylor et al. **FX: A TensorFlow-Based Production-Scale Machine Learning Platform.** KDD, 2017.
- A.Karpathy. **Building the Software 2.0 Stack.** Spark+AI Summit, 2018.

# Some Startups: Cogitai, Neurala, Gantry



Data evolves. Build ML systems that adapt.

Gantry gives you full  
retrain, wh



Product Solutions Technology Resources About Partners

## Improve Quality Inspections with Vision AI Software

Reduce product defects. Increase inspection rates. Prevent production downtime.

TALK WITH OUR EXPERTS

### Vision AI for Industrial Inspections

Neurala is dedicated to helping manufacturers enhance their vision inspection process. AI-powered visual quality inspection



Cogitai is happy to announce that we are now part of Sony AI

Visit Sony AI

Sony AI still supports Continua, our Reinforcement Learning platform. For more information, contact us at [info@coigitai.com](mailto:info@coigitai.com)

Try Continua

History:

Cogitai was founded in 2015 by [Mark Ring](#), [Satinder Singh](#), and [Peter Stone](#), leading AI innovators with a combined total of over 60 years of active research in designing AI algorithms to learn knowledge and actions from experience. The company's founding purpose was to make Reinforcement Learning (RL), and eventually Continual Learning, accessible to a wide range of industrial applications. Cogitai's Continua platform represented the first step towards the company's vision by incorporating the best publicly available and proprietary RL algorithms into a scalable, easy-to-use SaaS platform. Cogitai is excited to be able to continue its research and development efforts as a part of Sony AI.

<https://www.neurala.com>

<https://gantry.io>

<https://www.cogitai.com>

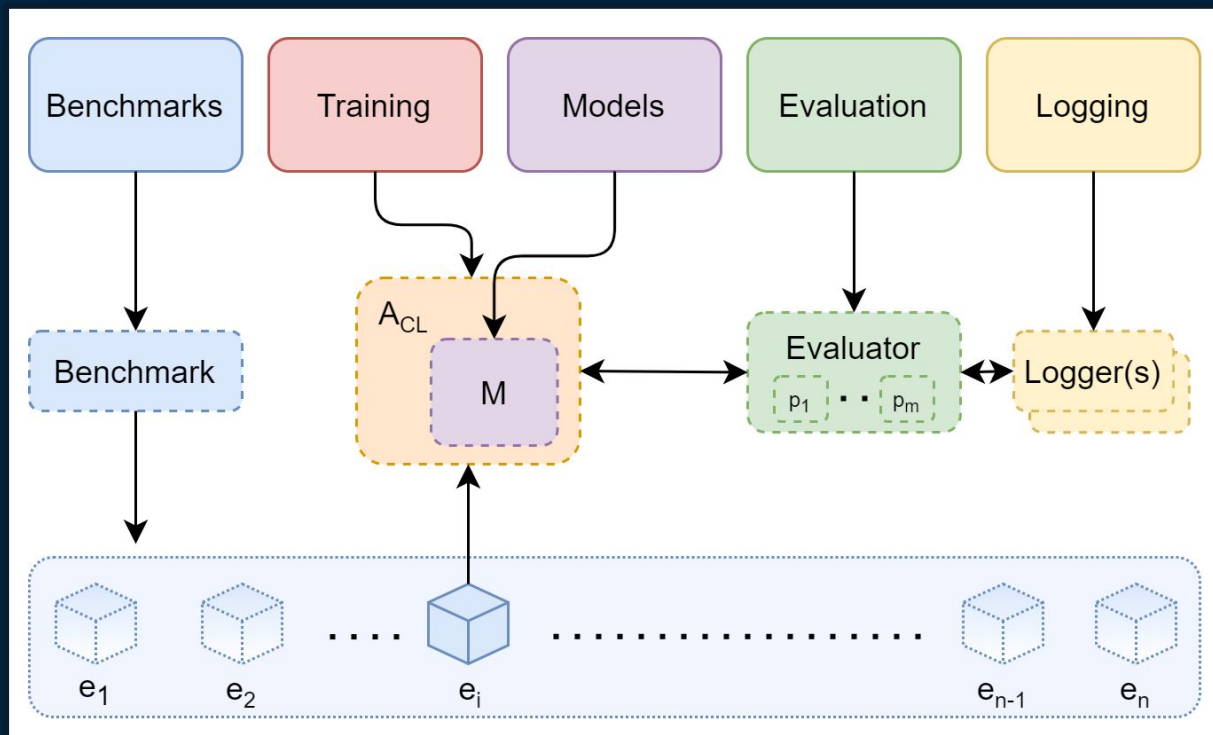


# Continual Learning Tools

## Research & Development Tools

- V. Lomonaco et al. ***Avalanche: an End-to-End Library for Continual Learning***. CLVision Workshop at CVPR 2021.
- A. Douillard et al. ***Continuum: Simple management of complex continual learning scenarios***. CLVision Workshop at CVPR 2021.
- S.I. Mirzadeh et al. ***CL-Gym: Full-Featured PyTorch Library for Continual Learning***. CLVision Workshop at CVPR 2021.
- F. Normandin et al. ***Sequoia: A Software Framework to Unify Continual Learning Research***. CLVision Workshop at CVPR 2021.

# Avalanche for R&D



# Avalanche for R&D

## Avalanche Key Links

- **Avalanche Official Website:**  
<https://avalanche.continualai.org>
- **Avalanche GitHub:**  
<https://github.com/ContinualAI/avalanche>
- **Avalanche API-DOC:**  
<https://avalanche-api.continualai.org>
- **Avalanche ContinualAI Slack:** #avalanche channel

```
With Avalanche | Without Avalanche
1 import torch
2 from torch.nn import CrossEntropyLoss
3 from torch.optim import SGD
4
5 from avalanche.benchmarks.classic import PermutedMNIST
6 from avalanche.extras.models import SimpleMLP
7 from avalanche.training.strategies import Naive
8
9 # Config
10 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
11
12 # model
13 model = SimpleMLP(num_classes=10)
14
15 # CL Benchmark Creation
16 perm_mnist = PermutedMNIST(n_experiences=3)
17 train_stream = perm_mnist.train_stream
18 test_stream = perm_mnist.test_stream
19
20 # Prepare for training & testing
21 optimizer = SGD(model.parameters(), lr=0.001, momentum=0.9)
22 criterion = CrossEntropyLoss()
23
24 # Continual learning strategy
25 cl_strategy = Naive(
26     model, optimizer, criterion, train_mb_size=32, train_epochs=2,
27     eval_mb_size=32, device=device)
28
29 # train and test loop
30 results = []
31 for train_task in train_stream:
32     cl_strategy.train(train_task, num_workers=4)
33     results.append(cl_strategy.eval(test_stream))
```

# Avalanche for R&D

```
replay = ReplayPlugin(mem_size)
ewc = EWCPlugin(ewc_lambda)
strategy = BaseStrategy(
    model, optimizer,
    criterion, mem_size,
    plugins=[replay, ewc])
```



# Impact on Sustainable AI

# Sustainable AI Principles

## General Principles

- **Accuracy & Robustness**
- **Explainability, Transparency & Accountability**
- **Bias**
- **Fairness**
- **Privacy & Security**
- **Human, Social and Environmental Wellbeing**

L. Royakkers et al. ***Societal and ethical issues of digitization***. Ethics and Information Technology, 2018.  
B.D. Mittelstadt et al. ***The ethics of algorithms: Mapping the debate***. Big Data & Society, 2016.  
A. Jobin et al. ***The global landscape of AI ethics guidelines***. Nature Machine Intelligence, 2019.  
<https://www.aiforpeople.org/ethical-ai/>

# Continual Learning Impact

## ...On each Principle:

- **Accuracy & Robustness** → Robustness & Autonomy, Continual & Fast Improvement
- **Bias** → CL as the new Agile: Bias Patches
- **Fairness** → Efficient Fairness Patches
- **Privacy & Security** → Security Patches
- **Human, Social and Environmental Wellbeing** → improved efficiency & scalability: less energy consumption, CO2 emission; sustainable & “progressive” by design
- **Explainability, Transparency & Accountability** → Neuroscience-grounded, Human-centered AI

L. Royakkers et al. *Societal and ethical issues of digitization*. Ethics and Information Technology, 2018.

B.D. Mittelstadt et al. *The ethics of algorithms: Mapping the debate*. Big Data & Society, 2016.

A. Jobin et al. *The global landscape of AI ethics guidelines*. Nature Machine Intelligence, 2019.

<https://www.aiforpeople.org/ethical-ai/>

The background is a dark teal color. It features several vertical white lines of varying lengths. Scattered throughout are small squares in teal, orange, and pink. Some squares are solid, while others are hollow. The text 'Open Questions' is centered in the middle of the page. 'Open' is in white, and 'Questions' is in orange. The 'Q' in 'Questions' has a small hollow square above it.

# Open Questions



# Open Questions (1/2)

1. Is it possible to learn **robust, deep representations continually**?
2. Are currently addressed **scenarios** and **eval metrics enough**?
3. What is the right **level of supervision**?
4. How to know **what to forget** and **what to remember**?
5. What's the relationship with **concept drift**?
6. Is **replay** a research direction worth pursuing?
7. Is **computation** more important than **memory**?
8. Is **gradient descent** the right algorithm to learn continually?
9. **Continual Meta-Learning & Meta-Continual Learning**: what's the right relationship?
10. What is the relationship with **Sequence** and **Continual Learning**?

# Open Questions (2/2)

1. Is **curiosity** important for **continual learning**?
2. What about **Curriculum Learning**?
3. **Compositionality** is a key aspect of human intelligence: what to expect for CL Systems?
4. **Self-Reflection**\*: accuracy of learned functions, given only unlabeled data?
5. **Self-reflection** that can detect every possible shortcoming (called impasse) of the agent\*
6. (External) **Knowledge and Reasoning**\*

*...and much more!*

# On the Future of CL (Short-Medium Term)

## 1. More Natural Scenarios

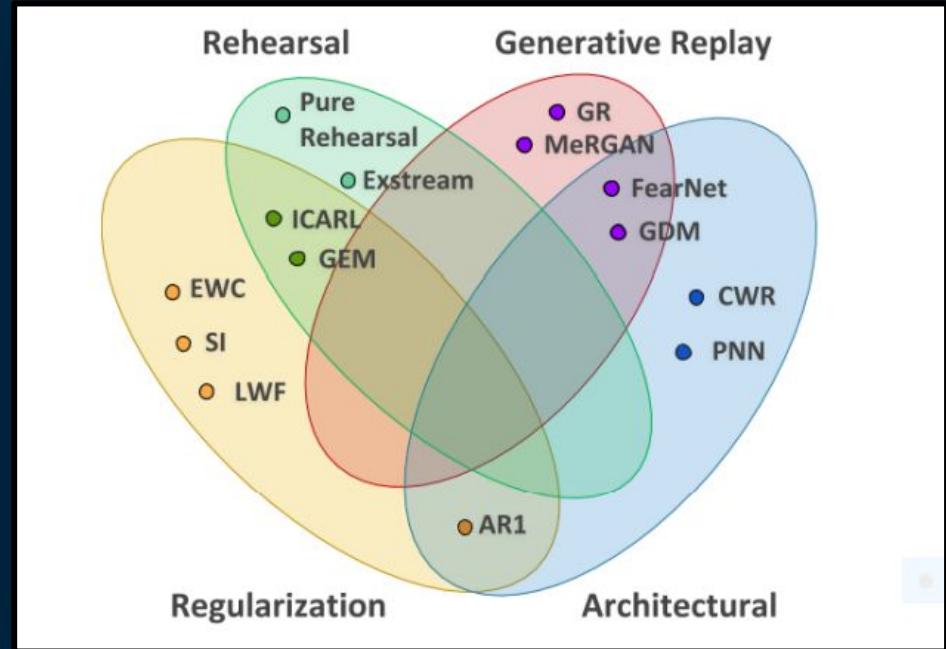
- **Domain, Task** and **Class-Incremental** are not enough.
- Longer streams of “*experiences*”.
- More metrics, focus on scalability.

## 2. Move towards unsupervised training

- Mostly **Semi-Supervised**, **Self-Supervised** and **Sequence Learning**.

## 3. Hybrid Continual Learning Strategies

## 4. Continual Learning Applications



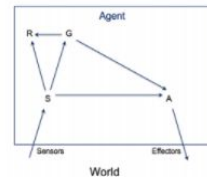
# On the Future of CL (Long-Term)

1. Fundamentally a question of **agent architecture**\*
2. **Two main paths for (deep) CL**
  - a. **Neuroscience-Inspired**
  - b. **Distributed Continual Learning**

## What should a theory of Learning Agents answer?

might model learning agent A as tuple  $\langle S, E, M, F, G, L \rangle$

- S = sensors
- E = effectors
- F = set of functions
- M = set of memory units
- G = graph specifying data flow among F, M, S, E
- L = learning mechanism



might model L as another agent  $L = \langle S_L, E_L, M_L, F_L, G_L \rangle$

- where  $S_L, E_L$  sense and act on Agent, especially its F, M, G

The background is a dark teal color. It features several vertical white lines of varying lengths. Scattered throughout are small squares in teal, orange, and pink. Some squares are solid, while others are hollow. The word "Conclusions" is centered in a large, orange, sans-serif font.

# Conclusions

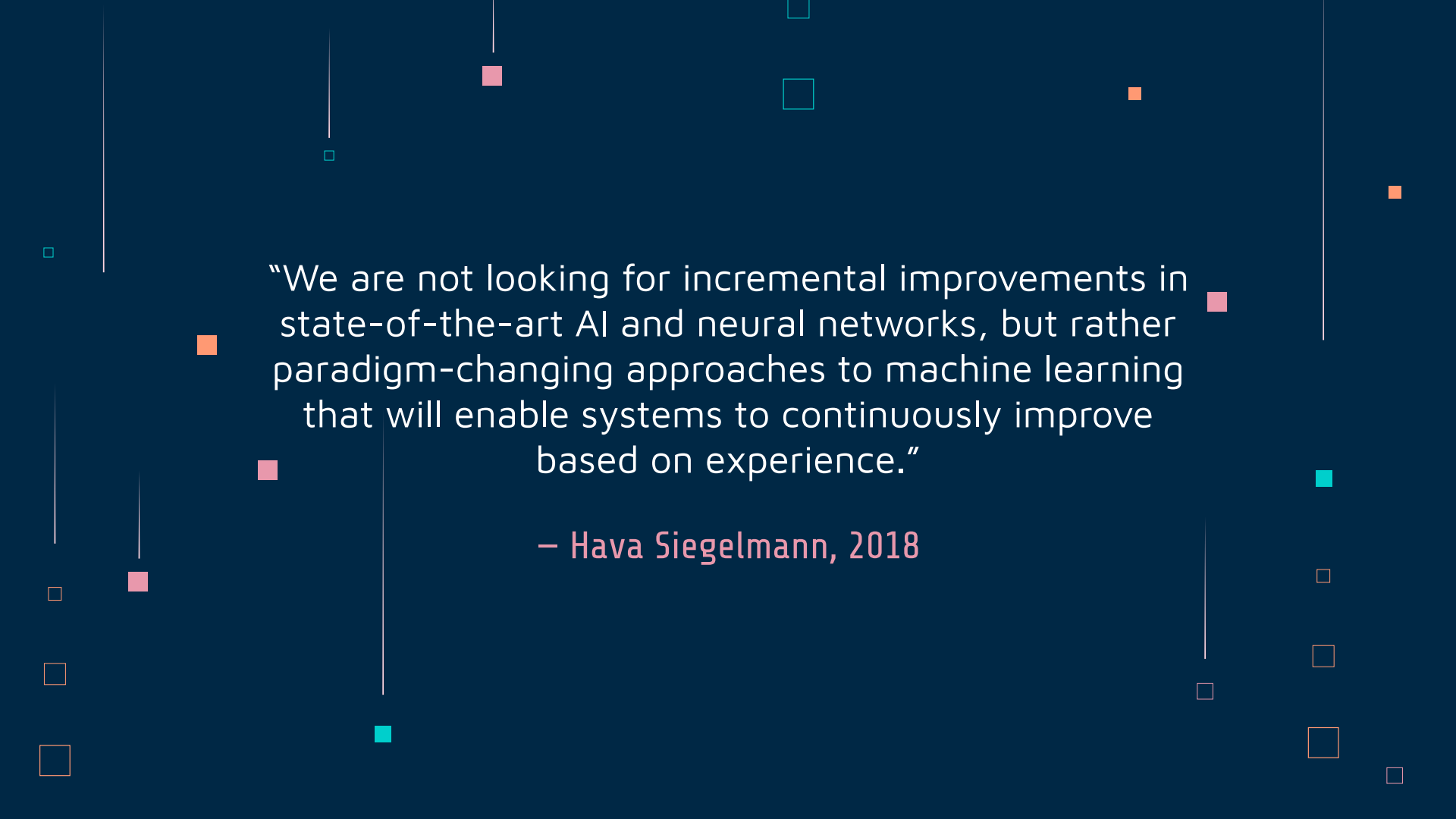
# Conclusions

## *What we have seen*

- Significant and **growing Interest** in the last few years on Continual learning within Deep Learning
- Significant **improvements over standard benchmark** but **focus still mostly on simplified scenarios** and forgetting centered metrics.
- **Huge space of possible and significant explorations.**

## *Take-Home Messages*

1. Continual Learning is a **paradigm-changing approach** trying to break the fundamental i.i.d. assumption in statistical learning.
2. CL pushes for the **next step in Neuroscience-grounded approaches** to learning
3. CL pushes for the next generation of truly intelligent robust and autonomous AI systems: **efficient, effective, scalable, hence sustainable.**

The background is a dark blue gradient. It features several vertical white lines of varying lengths. Scattered throughout are small squares in various colors: light blue, pink, orange, and teal. Some squares are solid, while others are hollow outlines. The overall aesthetic is clean and modern, typical of a presentation slide.

“We are not looking for incremental improvements in state-of-the-art AI and neural networks, but rather paradigm-changing approaches to machine learning that will enable systems to continuously improve based on experience.”

– Hava Siegelmann, 2018

The background features a dark blue field with scattered geometric elements. These include solid squares in teal, orange, and pink, as well as hollow squares in the same color palette. Thin white vertical lines of varying lengths are also present, some extending from the top or bottom edges towards the center. The overall aesthetic is clean and modern.

# Resources



# Additional Resources (1/3)

## Continual Learning with Deep Architectures

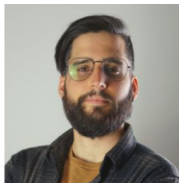
Vincenzo Lomonaco (University of Pisa & ContinualAI), Irina Rish (University of Montreal & MILA)

Tutorial @ ICML 2021

Mon Jul 19 08:00 AM -- 11:00 AM (PDT)

**Official Tutorial Website: Slides, Q&As, Recordings, etc.**

### Authors



Vincenzo Lomonaco  
University of Pisa &  
ContinualAI



Irina Rish  
University of  
Montreal & MILA



# ContinualAI

A Non-profit Research Organization and Open Community on  
**Continual Learning for AI**

[Home](#)

[News & Events](#)

[Research](#)

[Lab](#)

[Forum](#)

[Supporters](#)

[About us](#)

[Join us!](#)



ContinualAI.org

# Additional Resources (2/3)

- **ContinualAI Wiki**: a shared and collaboratively maintained *knowledge base* for Continual Learning: tutorials, workshops, demos, tutorials, courses, etc.
- **Continual Learning Papers**: curated list of CL papers & books with meta-data by ContinualAI
- **ContinualAI Forum** + Slack: discussions / Q&As about Continual Learning
- **ContinualAI Research Consortium**: networks of Top CL Labs across the world.

## Publications


In this section we maintain an updated list of publications related to Continual Learning. This references list is automatically generated by a single bibtex file maintained by the ContinualAI community through an open Mendeley group! Join our group here to add a reference to your paper! Please, remember to follow the (very simple) contributions guidelines when adding new papers.

Search among 262 papers!

Filter list by keyword:

Filter list by regex:

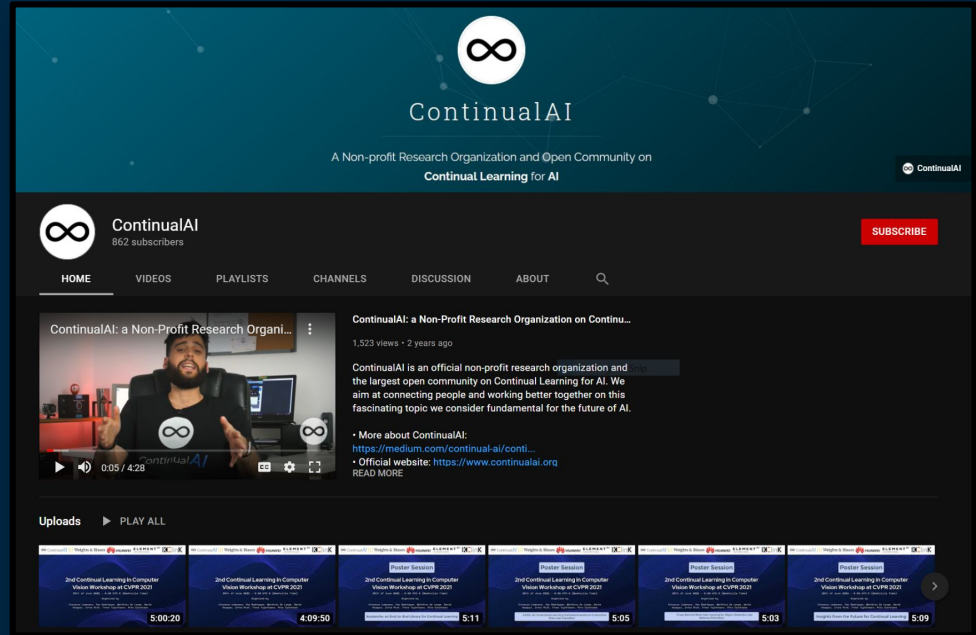
Filter list by year:

  
ContinualAI

[framework] [som] [sparsity] [dual] [spiking] [rnn] [ltp] [graph] [vision] [hebbian] [audio] [bayes]  
[generative] [mnist] [fashion] [cifar] [coresc] [imagenet] [omniglot] [cubs] [experimental]  
[theoretical]

# Additional Resources (3/3)

- **ContinualAI Publication**: a curated list of original blog posts on CL.
- **Continual Learning & AI Mailing List+**: curated list of CL papers & books with meta-data by ContinualAI.
- **ContinualAI Newsletter**: news from the ContinualAI community and the CL World in one place.
- **ContinualAI Seminars**: weekly invited talks on CL.
- **ContinualAI YouTube**: collection of videos about CL.



You can find more at: [www.continualai.org](http://www.continualai.org)

# Continual Learning: On Machines that Can Learn Continually

1st Official Open-Access Course  
on CL Offered by **University of  
Pisa & ContinualAI**

01

# Join our Pervasive AI Lab!

**PERVASIVE AI LAB**

HOME RESEARCH TEAM EVENTS NEWS

**ARTIFICIAL INTELLIGENCE AS A PERVASIVE TECHNOLOGY**

...t in IC  
...ies an  
...nd trus

 Massimo Coppola CONSIGLIO	 Andrea Cosbu SCUOLA NORMALE SUPERIORE	 Valerio De Caro UNIVERSITÀ DI PISA	 Marco Di Benedetto CONSIGLIO
 Daniele Di Sarli UNIVERSITÀ DI PISA	 Federico Erica UNIVERSITÀ DI PISA	 Fabrizio Falchi CONSIGLIO NAZIONALE DELLE RICERCHE	 Stefano Forti UNIVERSITÀ DI PISA
 Claudio Gallicchio UNIVERSITÀ DI PISA	 Alberto Gotta CONSIGLIO NAZIONALE DELLE RICERCHE	 Alessio Gravina UNIVERSITÀ DI PISA	 Alexander Kocian UNIVERSITÀ DI PISA
 Giacomo Landiano SCUOLA NORMALE SUPERIORE	 Francesco Landolfi UNIVERSITÀ DI PISA	 Vincenzo Lomonaco CONSIGLIO NAZIONALE DELLE RICERCHE	 Paolo Manghi UNIVERSITÀ DI PISA
 Chiara Renzo CONSIGLIO NAZIONALE DELLE RICERCHE	 Jacopo Soldani UNIVERSITÀ DI PISA	 Salvatore Trani CONSIGLIO NAZIONALE DELLE RICERCHE	
 Daniele Mazzei UNIVERSITÀ DI PISA	 Gabriele Mancagli UNIVERSITÀ DI PISA	 Davide Maroni CONSIGLIO NAZIONALE DELLE RICERCHE	 Alessio Micheli UNIVERSITÀ DI PISA
 Ugo Montanari UNIVERSITÀ DI PISA	 Orsina Maresca CONSIGLIO	 Franco Maria Martelli CONSIGLIO	 Dario Nameroff UNIVERSITÀ DI PISA
 Marlio Racco CONSIGLIO NAZIONALE DELLE RICERCHE	 Paolo Saracchi CONSIGLIO NAZIONALE DELLE RICERCHE	 Antonio Brogi UNIVERSITÀ DI PISA	 Emanuele Carlini CONSIGLIO NAZIONALE DELLE RICERCHE
 Antonio Carta UNIVERSITÀ DI PISA	 Pietro Cassarà CONSIGLIO NAZIONALE DELLE RICERCHE	 Daniele Castellana UNIVERSITÀ DI PISA	 Stefano Chessa UNIVERSITÀ DI PISA

Teaching & Supervision



Research



Spin-off



Consultancy



The lab is in Pisa, Italy! Feel free to visit and get in touch with us anytime! Official website: [Pervasive AI Lab \(unipi.it\)](http://Pervasive AI Lab (unipi.it))

Do you have any questions?

[vincenzo.lomonaco@unipi.it](mailto:vincenzo.lomonaco@unipi.it)

[vincenzolomonaco.com](http://vincenzolomonaco.com)

University of Pisa

# THANKS



CREDITS: This presentation template was created by [Slidesgo](#), including icons by [Flaticon](#), and infographics & images by [Freepik](#)