

Continual Learning: On Machines that can Learn Continually

Official Open-Access Course @ University of Pisa, ContinualAI, AIDA

Lecture 2: Understanding Catastrophic Forgetting

Vincenzo Lomonaco

University of Pisa & ContinualAI

vincenzo.lomonaco@unipi.it

TABLE OF CONTENTS



01

What's
Forgetting?



02

Forgetting
with 1
Neuron



03

A Deep
Learning
Example



04

Intro to
Avalanche



What's Forgetting?

A Concrete Example



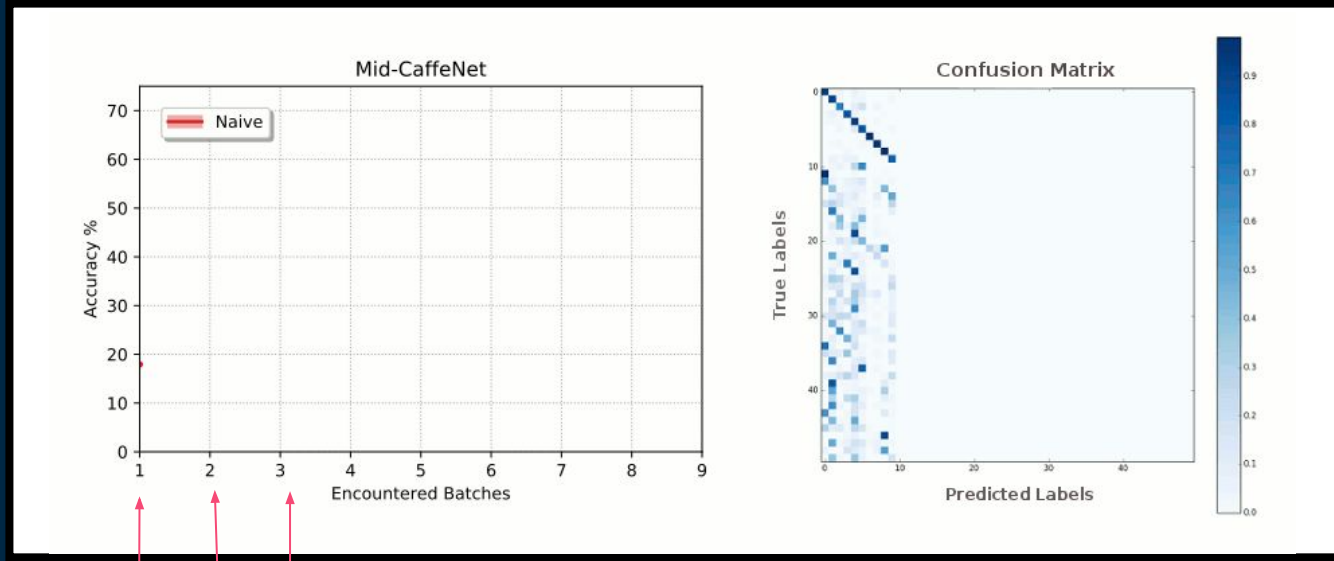
- **50GB/s** streaming data.
- **~30240 TB of data** after only a week.
- **Impossible** to re-train the mini-spot brain from scratch and to **adapt fast**.



R1 Example from *Istituto Italiano di Tecnologia*



Continual Learning: what's the problem?



- A set of new objects (classes) each day
- 10 the first day, 5 the following

The Stability-Plasticity Dilemma

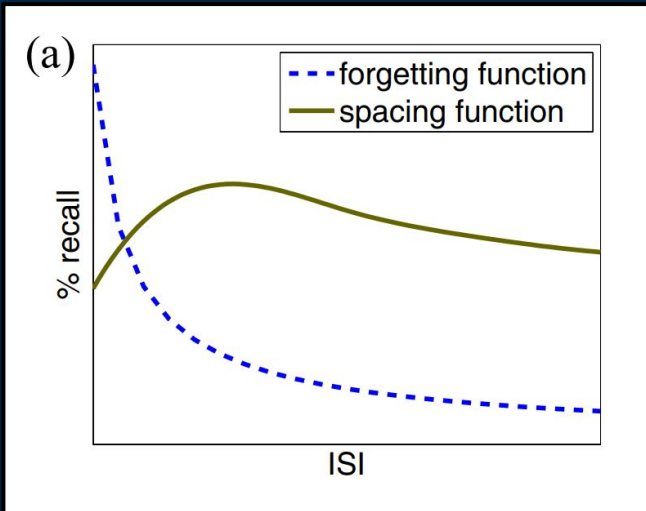
Stability-Plasticity Dilemma:

- Remember past concepts
- Learn new concepts
- Generalize

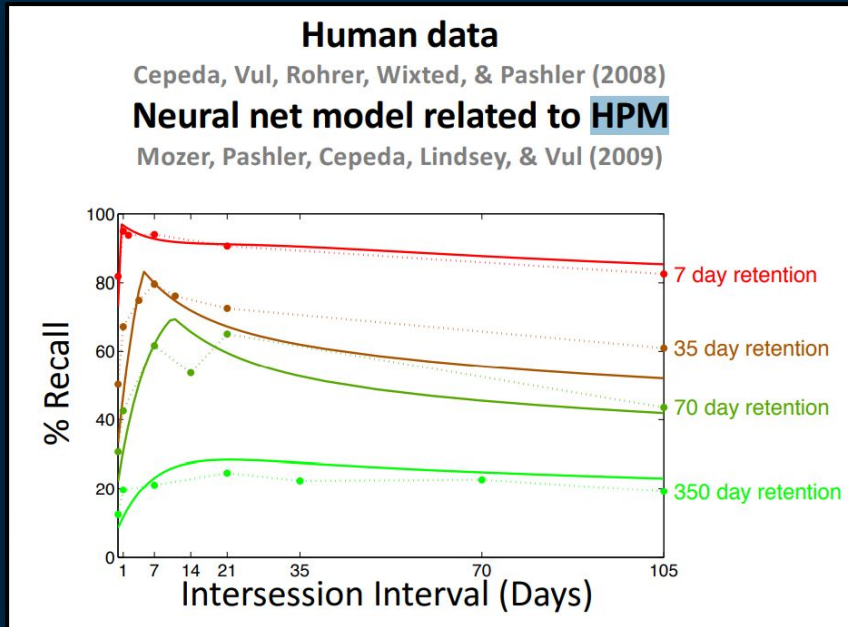
First Problem in Deep Learning:

- Catastrophic Forgetting

Forgetting in Humans

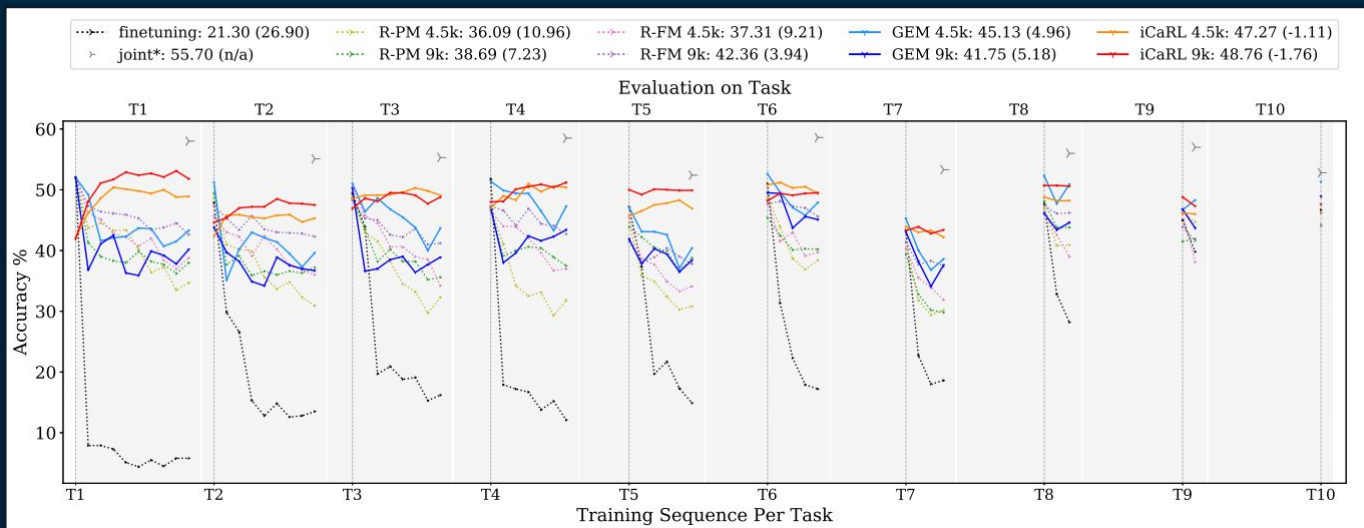


The spacing function (solid line) depicts recall at test following two study sessions separated by a given inter-session interval (ISI); the forgetting function (dashed line) depicts recall as a function of the lag between study and test.

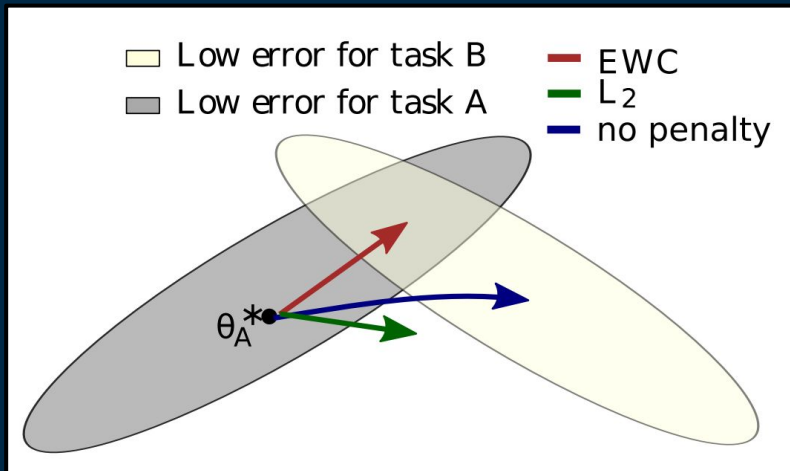


Forgetting in Machines

*Catastrophic interference, also known as **catastrophic forgetting**, is the tendency of an artificial neural networks to completely and abruptly forget previously learned information upon learning new information. -> Mostly due to Gradient Descent.*



Forgetting in Machines



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.

1. We don't have access to previously encountered data.
2. \mathcal{L}_S can only be approximated.

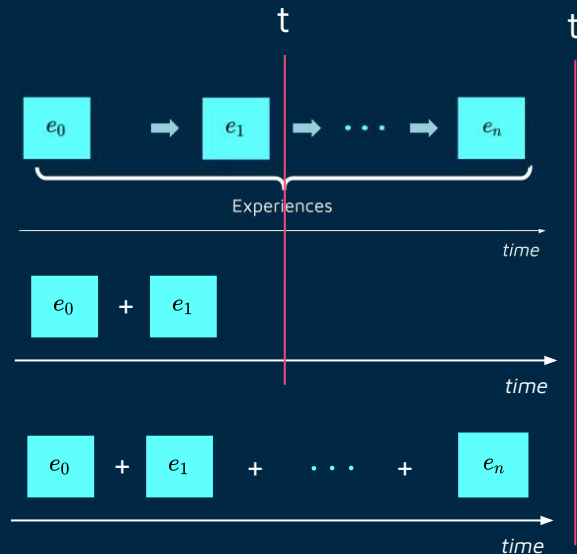
Classic ML vs CL

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.



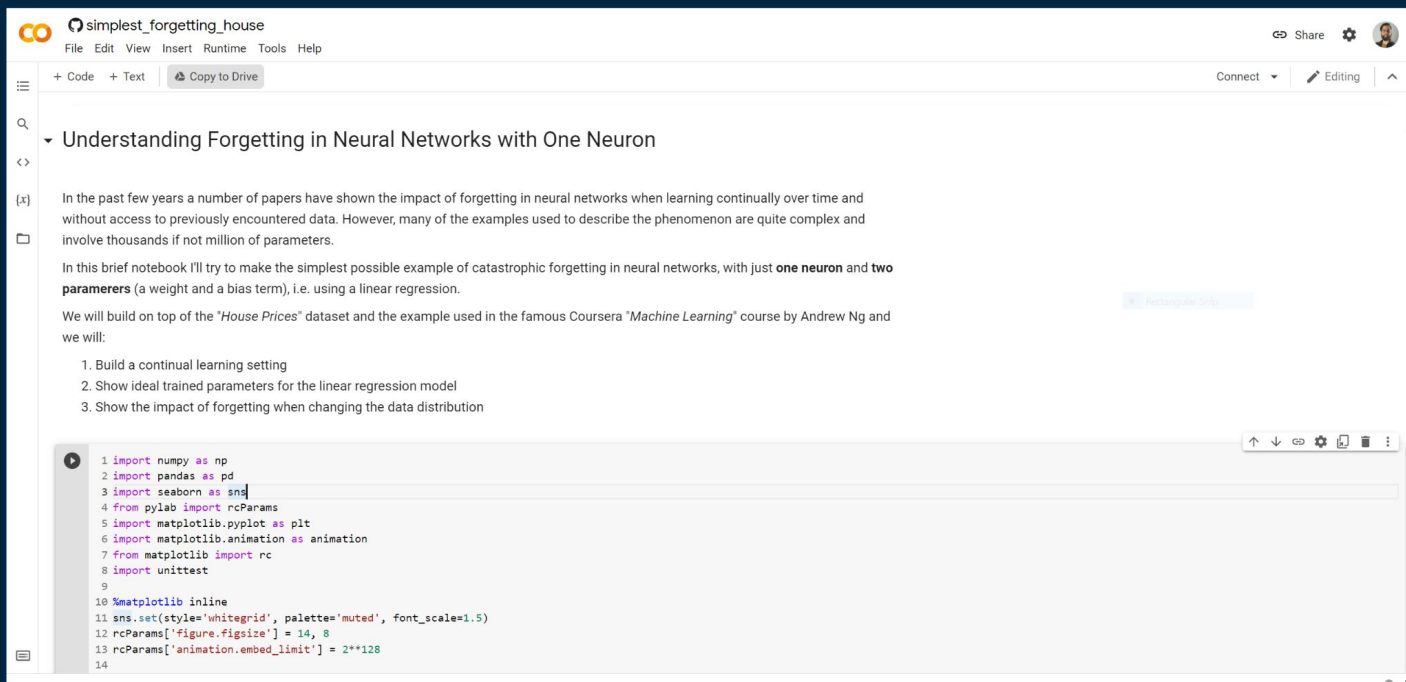
- Given a time t , you can always transform the continual learning problem in a static one.
- Continual Learning is mostly about leveraging previously learned knowledge, since t is unbounded.
- More efficiency (low & bounded memory and computation) -> more forgetting.

The background is a dark blue gradient. It features several thin, vertical white lines of varying lengths scattered across the frame. Interspersed among these lines are small squares in three colors: light blue, light orange, and light pink. Some squares are solid, while others are outlined. The overall aesthetic is modern and minimalist.

■ Forgetting with 1 Neuron

House Prices Estimation Example

Demo Session!



The screenshot shows a Google Colab notebook titled "simplest_forgetting_house". The interface includes a menu bar (File, Edit, View, Insert, Runtime, Tools, Help), a toolbar with options like "+ Code", "+ Text", and "Copy to Drive", and a "Share" button. The notebook content is titled "Understanding Forgetting in Neural Networks with One Neuron".

The text in the notebook reads:

In the past few years a number of papers have shown the impact of forgetting in neural networks when learning continually over time and without access to previously encountered data. However, many of the examples used to describe the phenomenon are quite complex and involve thousands if not million of parameters.

In this brief notebook I'll try to make the simplest possible example of catastrophic forgetting in neural networks, with just **one neuron** and **two parameters** (a weight and a bias term), i.e. using a linear regression.

We will build on top of the "House Prices" dataset and the example used in the famous Coursera "Machine Learning" course by Andrew Ng and we will:

1. Build a continual learning setting
2. Show ideal trained parameters for the linear regression model
3. Show the impact of forgetting when changing the data distribution

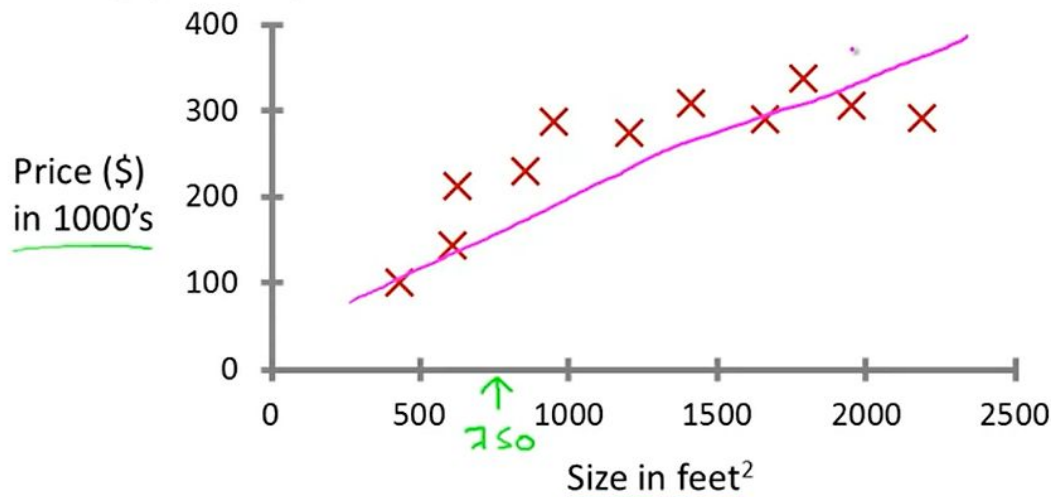
Below the text is a code cell containing the following Python code:

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 from pylab import rcParams
5 import matplotlib.pyplot as plt
6 import matplotlib.animation as animation
7 from matplotlib import rc
8 import unittest
9
10 %matplotlib inline
11 sns.set(style='whitegrid', palette='muted', font_scale=1.5)
12 rcParams['figure.figsize'] = 14, 8
13 rcParams['animation.embed_limit'] = 2**128
14
```

House Prices Estimation Example

Demo Session!

Housing price prediction.



Andrew Ng

Demo Session!

House Prices Estimation Example

$$y = \sum_j w_j x_j + b.$$

$$\mathcal{L}(y, t) = \frac{1}{2}(y - t)^2.$$

$$\begin{aligned}\mathcal{E}(w_1, \dots, w_D, b) &= \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y^{(i)}, t^{(i)}) \\ &= \frac{1}{2N} \sum_{i=1}^N \left(y^{(i)} - t^{(i)} \right)^2 \\ &= \frac{1}{2N} \sum_{i=1}^N \left(\sum_j w_j x_j^{(i)} + b - t^{(i)} \right)^2\end{aligned}$$

$$\frac{\partial \mathcal{E}}{\partial w_j} = \frac{1}{N} \sum_{i=1}^N x_j^{(i)} (y^{(i)} - t^{(i)})$$

$$\frac{\partial \mathcal{E}}{\partial b} = \frac{1}{N} \sum_{i=1}^N y^{(i)} - t^{(i)}.$$

Just add a “dummy” input x_0 which always takes the value 1; then the weight w_0 acts as a bias.

$$w_j \leftarrow w_j - \alpha \frac{\partial \mathcal{E}}{\partial w_j}.$$

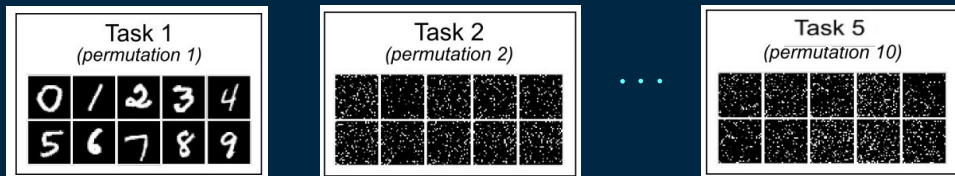
$$\mathbf{w} \leftarrow \mathbf{w} - \frac{\alpha}{N} \mathbf{X}^\top (\mathbf{y} - \mathbf{t}),$$

The background is a dark blue field decorated with an abstract pattern of geometric elements. It includes numerous small squares in teal, pink, and orange, as well as thin white vertical lines of varying lengths, creating a modern, digital aesthetic.

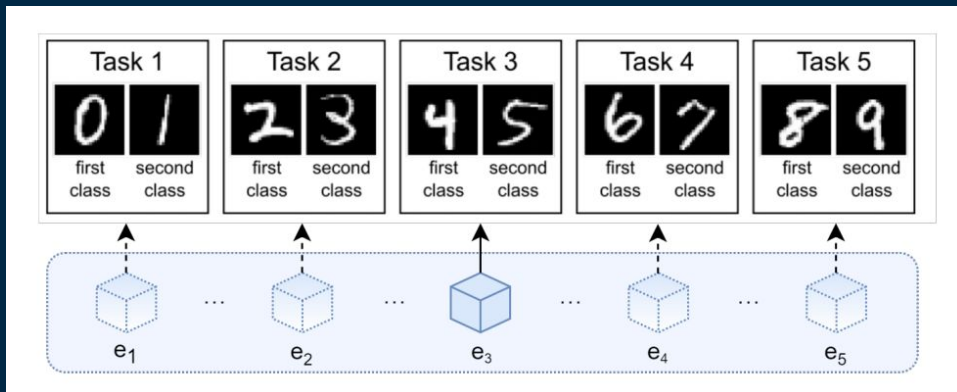
MNIST: a Classic Computer Vision Example

Permuted MNIST and SplitMINIST

Demo Session!



or



Permuted MNIST

Split MNIST

The background is a dark blue field decorated with a pattern of small squares and thin vertical lines. The squares are in three colors: light blue, orange, and pink. Some squares are solid, while others are hollow. The vertical lines are thin and white, extending from the top or bottom of the frame. The overall aesthetic is modern and minimalist.

How to Solve Forgetting?

How to Solve Forgetting?

**Brainstorming
Session!**

- Can we store just a portion of the previously encountered data?
- Can we make an ensemble of models?
- **Can we freeze some parameters of the network?**
- Have the model parameters the same importance? Can we leverage this to avoid forgetting?
- How model topology effects forgetting? What about pre-trained networks?
- ...?

CL: Not Only Forgetting

- Catastrophic forgetting is just a self-evident prominent failure of current ML systems
- **Continual Learning is much more!**
- **Efficiency of Learning (memory, computation)**
- **Backward and Forward transfer**
- **Compositionality**
- **Robustness**
- **Learning to Learn**
- ...



Intro to Avalanche

Avalanche: an End-to-End Library for CL

Avalanche is an **End-to-End Continual Learning Library** based on PyTorch, born within **ContinualAI** with the unique goal of providing a collaborative and community-driven open-source (MIT licensed) codebase for **fast prototyping, training** and **reproducible evaluation** of continual learning algorithms.

Avalanche can help Continual Learning researchers and practitioners in several ways:

- *Write less code, prototype faster & reduce errors.*
- *Improve reproducibility.*
- *Improve modularity and reusability.*
- *Increase code efficiency, scalability & portability.*
- *Augment impact and usability of your research products.*

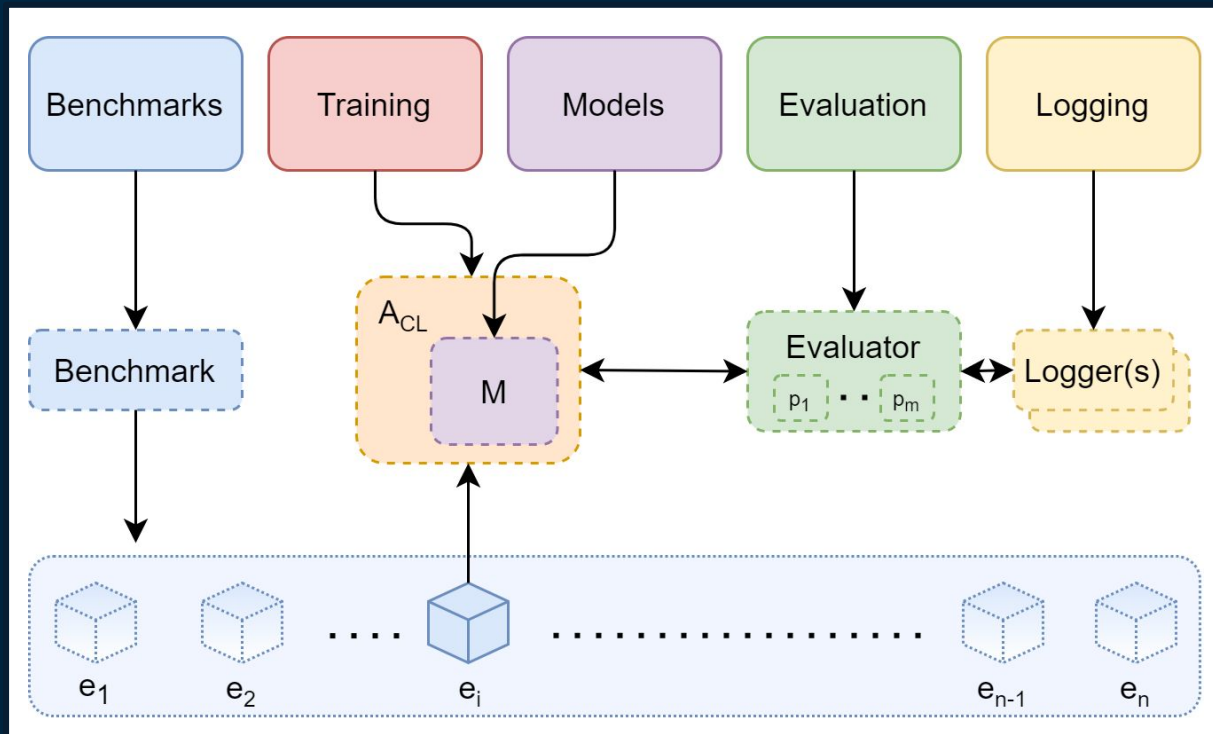
Avalanche Design Principles

1. Comprehensiveness and Consistency
2. Ease-of-Use
3. Low Constraints / Assumptions
4. Reproducibility and Portability
5. Modularity and Independence
6. Efficiency and Scalability
7. Community-driven

Introduction to Avalanche

- 5 main maintainers (Vincenzo Lomonaco, Lorenzo Pellegrini, Andrea Cossu, Antonio Carta, Gabriele Graffieti)
- 40+ contributors and top CL players around the world (more on the people here)
- 15+ organizations and leading CL labs from Europe, USA, China.
- Integrating different codebases (like FACIL)
- Huge Sponsorships on the way...

Avalanche in a Nutshell



Avalanche: an End-to-End Library for CL

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))
```

The background is a dark blue gradient. It features several thin, vertical white lines of varying lengths scattered across the frame. Interspersed among these lines are small squares in three colors: light blue, pink, and orange. Some squares are solid, while others are outlined. The overall aesthetic is modern and minimalist.

Next: Scenarios & Benchmarks

Do you have any questions?

vincenzo.lomonaco@unipi.it

vincenzolomonaco.com

University of Pisa

THANKS



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