



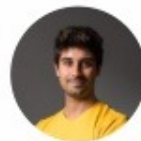
# Catastrophic Forgetting in Deep Graph Networks: an Introductory Benchmark for Graph Classification



**Antonio  
Carta**



**Andrea  
Cossu**

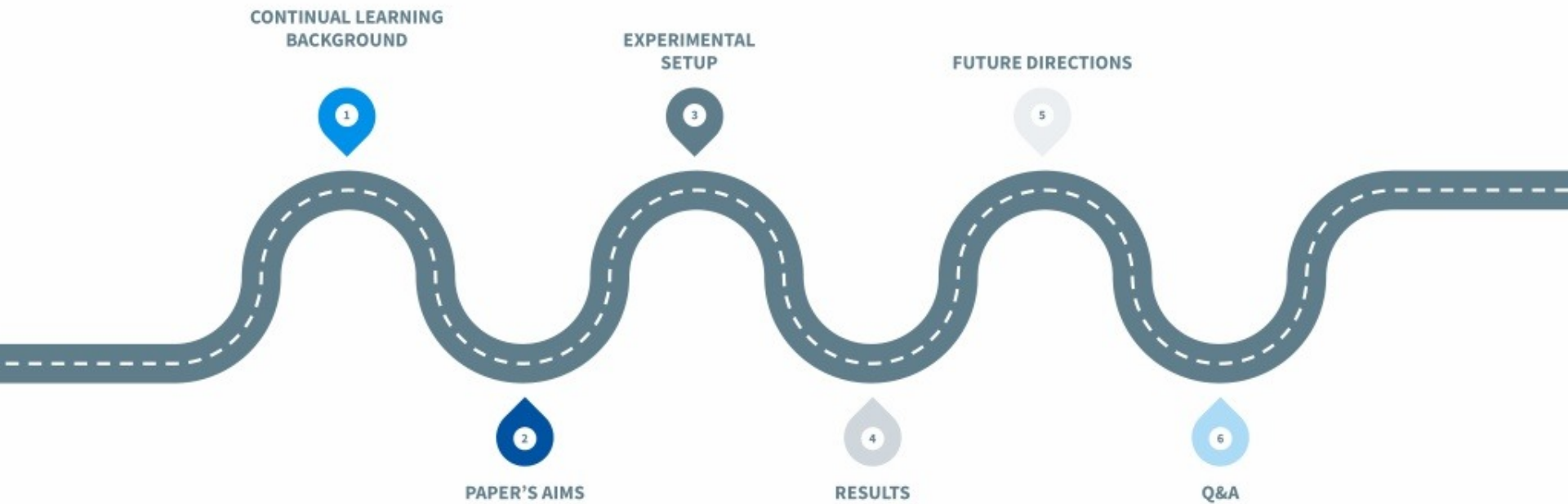


**Federico  
Errica**

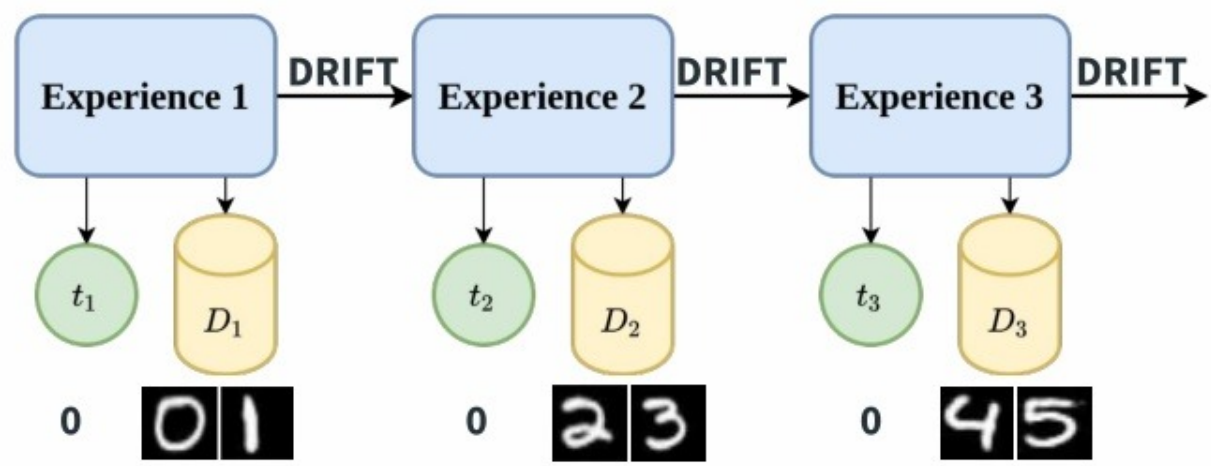


**Davide  
Bacciu**

# Outline



# Continual Learning setting



WE STUDIED CATASTROPHIC FORGETTING

# Why bothering?

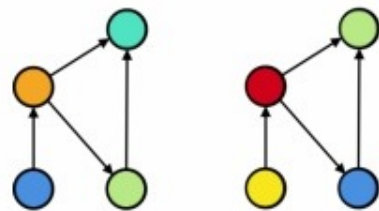
Training on the entire data may even lead to superior results! However...

- ◎ Datasets may be **huge** → **no training on the edge** 😞
- ◎ **New data** after deployment → **retraining**
  - **Expensive** 😞
  - **Inefficient** → most of the **information** is **already in** the **model** 😞

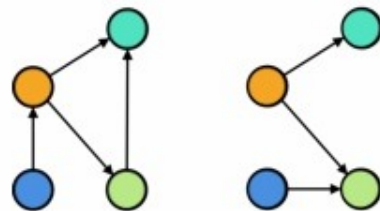
# Paper's Aim (1)

## Does CL work on graphs?

- **Graph** distribution drift
- GRL **regularization** strategies?
- Graph **classification**



**node distribution drift**



**neighborhood distribution drift**



## Paper's Aim (2)

### Reproducible Research + Framework

- **Foster** future research
- Avoid **common mistakes**
- Handle boilerplate code

PYDGN ♥ CL

[https://github.com/diningphil/continual\\_learning\\_for\\_graphs](https://github.com/diningphil/continual_learning_for_graphs)



## CL + GRL techniques

**Elastic Weight Consolidation (EWC)**  $\longrightarrow \mathcal{R}(\Theta, \Omega) = \lambda \sum_{i=1}^{n-1} \Omega_i \|\Theta_i - \Theta_n\|_2^2$

**Learning without Forgetting (LwF)**  $\longrightarrow \mathcal{R}(\Theta_n, \Theta_{n-1}; \mathbf{x}, \mathbf{y}) = \alpha \text{KL}[p_{\Theta_n}(\mathbf{y}|\mathbf{x}) \parallel p_{\Theta_{n-1}}(\mathbf{y}|\mathbf{x})]$

**Structure Preserving Regularization (REG)**  $\longrightarrow p(\mathbf{A} | \mathbf{Z}) = \prod_{i=1}^N \prod_{j=1}^N p(A_{ij} | \mathbf{z}_i, \mathbf{z}_j)$  [Ref. 3]  
with  $p(A_{ij} = 1 | \mathbf{z}_i, \mathbf{z}_j) = \sigma(\mathbf{z}_i^\top \mathbf{z}_j)$

Replay Memory

Naive Strategy

# Which Models?

## Structure-agnostic Baseline

- Impact of structure in **CL**
- MLP + Mean Global Pooling

$$\mathbf{h}_v = \psi(\mathbf{x}_v), \quad x_v \in \mathcal{X}_g,$$

$$\psi(x_v) = \mathbf{W}_L^T (\sigma(\dots (\sigma(\mathbf{W}_1^T x_v + \mathbf{b}_1) \dots) + \mathbf{b}_L)$$

$$\mathbf{h}_g = \Psi_g(\{\mathbf{h}_v \mid v \in \mathcal{V}_g\})$$

## Generic and simple DGN

- Based on GraphSAGE
- Mean aggregator & Global Pooling

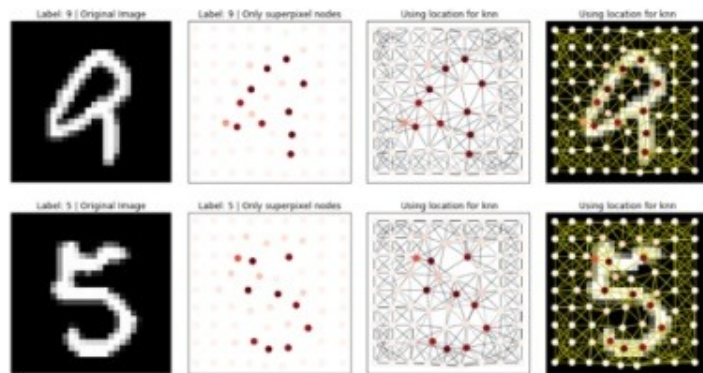
$$\mathbf{h}_v^{\ell+1} = \phi^{\ell+1}(\mathbf{h}_v^\ell, \Psi_n(\{\psi^{\ell+1}(\mathbf{h}_u^\ell) \mid u \in \mathcal{N}_v\}))$$

$$\mathbf{h}_g = \Psi_g(\{\mathbf{h}_v \mid v \in \mathcal{V}_g\})$$

**MLP**  
**GRAPH CONVOLUTION(S)**  
**GLOBAL POOLING**

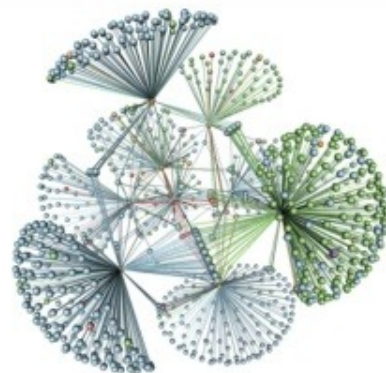


# Datasets



(a) MNIST [Ref. 5]

**CIFAR10 (same preprocessing of MNIST)**



**OGBG-PPA**

# Evaluation Setup

- Class-incremental Scenario
- Monitor  $ACC = \frac{1}{T} \sum_{t=1}^T R_{T,t}$
- Hold-out + Hyper-param optimization **for all models**

	MNIST	CIFAR10	OGBG-PPA
Size	70000	60000	158100
Node Attrs.	3	5	0
Edge Attrs.	0	0	7
Classes	10	10	37
Avg $ \mathcal{V}_g $	70,57	117,63	243,4
Avg $ \mathcal{E}_g $	564,63	941,07	2266,1
Data Split	55K/5K/15K	45K/5K/15K	49%/29%/22%
Class Split	2+2+2+2+2	2+2+2+2+2	17+5+5+5+5

# Results

**Baseline:  
Competitive!**

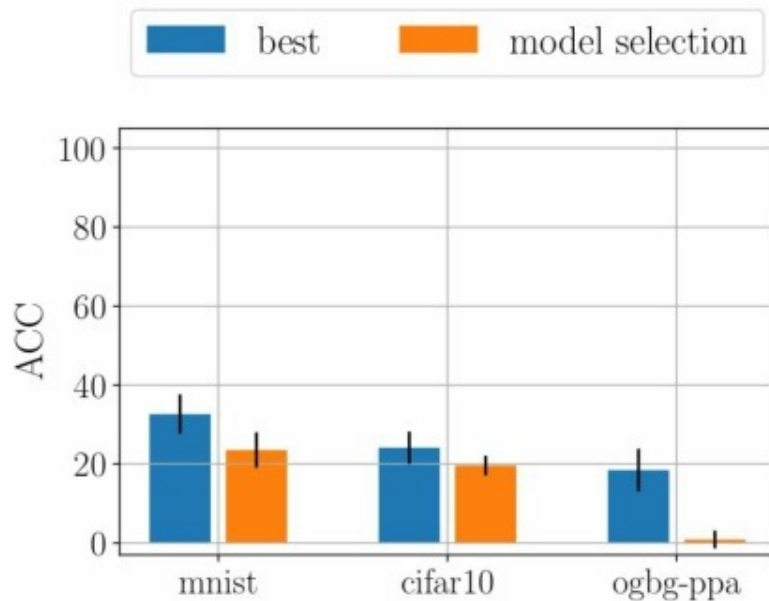


**Are DGNs ignoring  
the structure?**

	Model	Strategy			
		Naïve	EWC	Replay	LWF
MNIST	Baseline	19.56 $\pm$ 0.1	19.39 $\pm$ 0.1	86.13 $\pm$ 4.5	33.16 $\pm$ 13.1
	DGN	19.19 $\pm$ 0.1	18.95 $\pm$ 0.3	79.52 $\pm$ 1.9	32.64 $\pm$ 5.0
	DGN+reg	19.31 $\pm$ 0.1	—	81.42 $\pm$ 2.4	—
CIFAR10	Baseline	17.49 $\pm$ 0.1	17.49 $\pm$ 0.1	42.87 $\pm$ 3.7	26.77 $\pm$ 5.1
	DGN	17.11 $\pm$ 0.2	17.10 $\pm$ 0.2	39.55 $\pm$ 2.3	24.13 $\pm$ 4.1
	DGN+reg	17.13 $\pm$ 0.1	—	46.61 $\pm$ 3.5	—
OGBG-PPA	Baseline	14.53 $\pm$ 0.5	13.90 $\pm$ 0.8	55.96 $\pm$ 3.0	20.83 $\pm$ 6.1
	DGN	14.47 $\pm$ 0.3	14.15 $\pm$ 0.5	56.34 $\pm$ 2.5	18.46 $\pm$ 5.4
	DGN+reg	15.18 $\pm$ 0.8	—	57.27 $\pm$ 3.2	—

**Table 2: Mean accuracy and mean standard deviation (in parenthesis) among all steps. Replay results are related to memory size of 1000. Results are averaged over 5 final runs. We treat the regularization loss as a separate strategy.**

## Results: LwF sensitivity to hyper-parameters



## The road ahead

- ◎ Better understand the **role of graph distribution drift** on forgetting
  - Design ad-hoc regularization
- ◎ **Need more benchmarks!**
- ◎ Study **node classification** and other tasks

# References

1. Cossu A., Carta A., Errica F., Bacciu D., "Catastrophic Forgetting in Deep Graph Networks: an Introductory Benchmark for Graph Classification", **Graph Learning Benchmarks Workshop, WWW 2021**.
2. Bacciu D., Errica F., Micheli A., Podda M., "A Gentle Introduction to Deep Learning for Graphs", **Neural Networks, 2021**
3. Kipf, T. N., Welling, M., "Variational Graph Auto-Encoders", **Bayesian Deep Learning Workshop, NIPS, 2016**
4. Hamilton, W., Ying, Z., Leskovec, J., "Inductive representation learning on large graphs", **NIPS, 2017**
5. Dwivedi, V. P., Joshi, C. K., Laurent, T., Bengio, Y., & Bresson, X., "Benchmarking Graph Neural Networks", *arXiv 2020*.



# Questions?

## Thank you!

You can reach me at:

[andrea.cossu@sns.it](mailto:andrea.cossu@sns.it)

[andreacossu.github.io](https://andreacossu.github.io)

CIML homepage: [ciml.di.unipi.it](https://ciml.di.unipi.it)

PAI Lab: <http://pai.di.unipi.it/>

