

Apprentissage continu de représentations visuelles

ENSIMAG
2023-2024



KartEEK Alahari & Diane Larlus

Apprentissage continu

<https://project.inria.fr/bigvisdata/>



Summary: Continual Learning

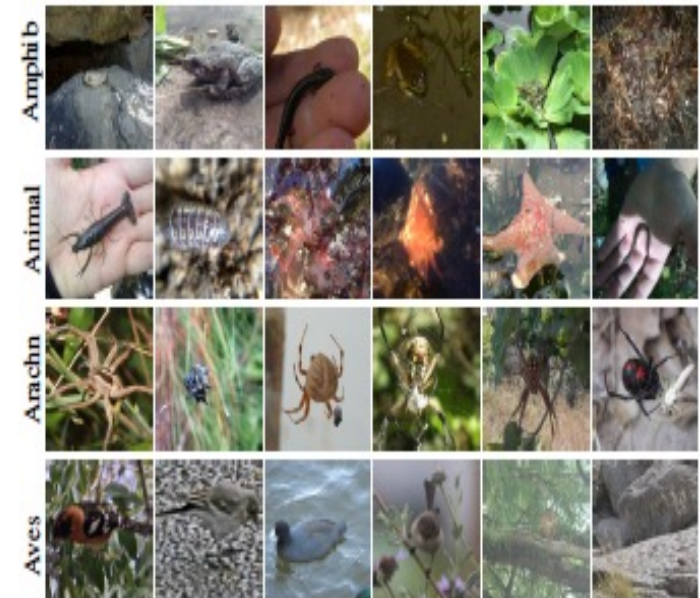
- **Flavour of different approaches:**
 1. Regularization based: LwF, EBLL, EWC, SI, MAS, IMM, ...
 2. Rehearsal / Replay: iCaRL, DGR, GEM, ...
 3. Architecture based: PackNet, progressive nets , HAT, ...
- Other learning frameworks, e.g., self-supervised
- Takeaways

A Comparative Analysis

- TinyImagenet: small, balanced, class-incremental
- iNaturalist: large-scale, unbalanced, task-incremental

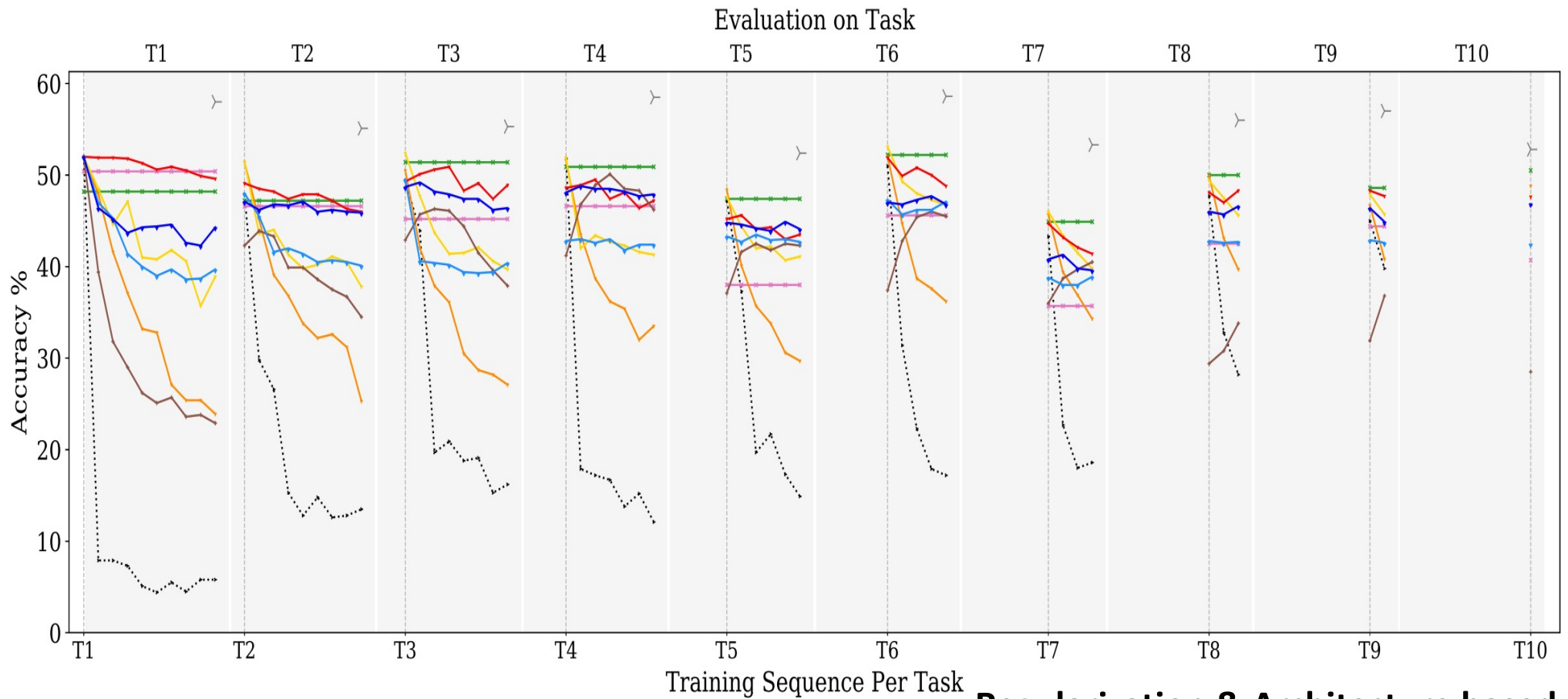
	Tiny Imagenet	iNaturalist
Tasks	10	10
Classes per task	20	5 to 314
Training data per task	8k	0.6k to 66k
Validation data per task	1k	0.1k to 9k
Task Constitution	random class selection	supercategory

- Fair way of setting hyperparameters (stability-plasticity tradeoff)



Comparative Evaluation (TinyImagenet)

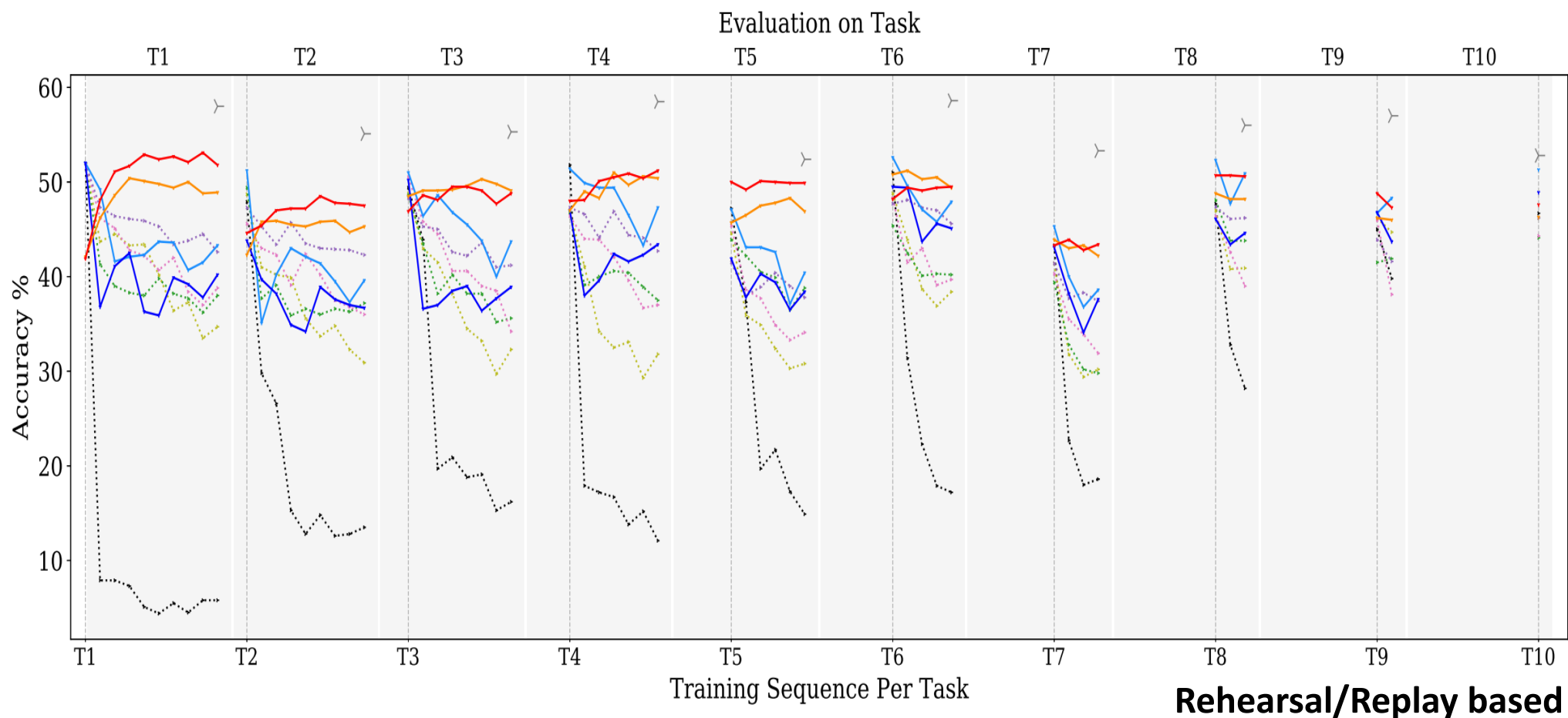
····· finetuning: 21.30 (26.90)	—*— PackNet: 49.13 (0.00)	—*— SI: 33.93 (15.77)	—*— MAS: 46.90 (1.58)	—*— LwF: 41.91 (3.08)
—*— joint*: 55.70 (n/a)	—*— HAT: 43.57 (0.00)	—*— EWC: 42.43 (7.51)	—*— mode-IMM: 36.89 (0.98)	—*— EBLL: 45.34 (1.44)



Regularization & Architecture based

Comparative Evaluation (TinyImagenet)

···	finetuning: 21.30 (26.90)	···	R-PM 4.5k: 36.09 (10.96)	···	R-FM 4.5k: 37.31 (9.21)	—	GEM 4.5k: 45.13 (4.96)	—	iCaRL 4.5k: 47.27 (-1.11)
·	joint*: 55.70 (n/a)	···	R-PM 9k: 38.69 (7.23)	···	R-FM 9k: 42.36 (3.94)	—	GEM 9k: 41.75 (5.18)	—	iCaRL 9k: 48.76 (-1.76)

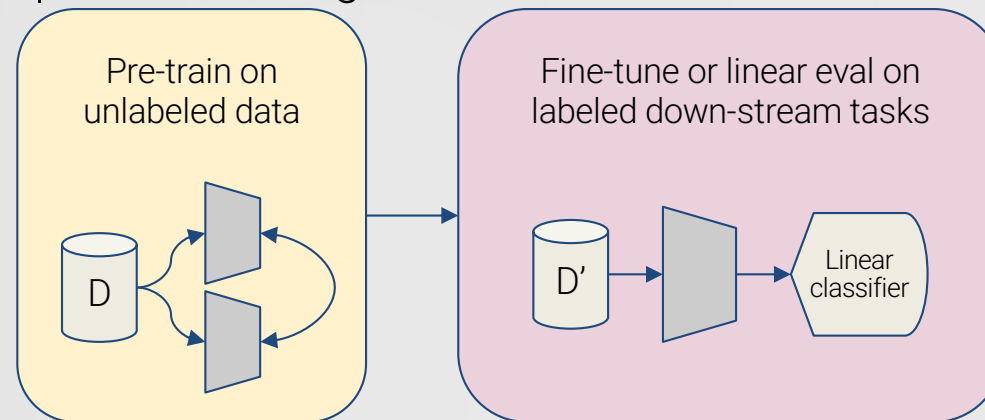


What else will we see in the class?

- Flavour of different approaches:
 1. Regularization based: LwF, EBLL, EWC, SI, MAS, IMM, ...
 2. Rehearsal / Replay: iCaRL, DGR, GEM, ...
 3. Architecture based: PackNet, progressive nets , HAT, ...
- **Other learning frameworks, e.g., self-supervised**
- Takeaways

Self-Supervised Learning (SSL)

- Self-Supervised Learning exploits the intrinsic structure of the data to pretrain strong feature extractors

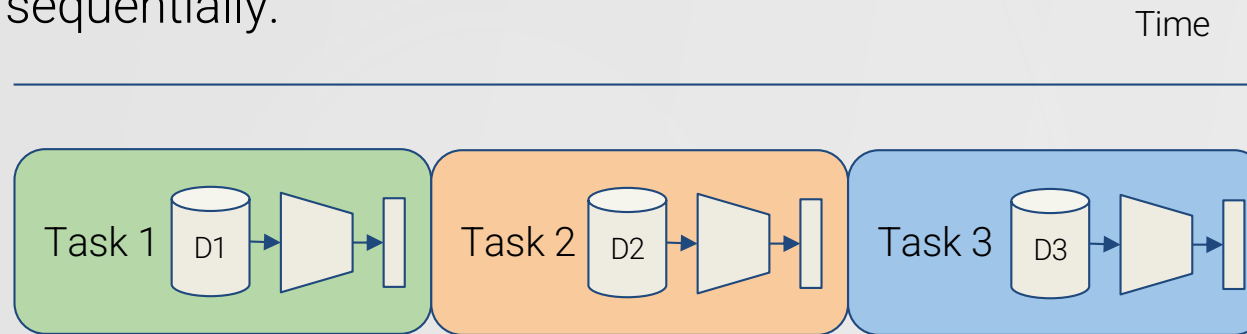


- The pre-training objective is to find parameters such as:

$$\operatorname{argmin}_{\theta} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\mathcal{L}_{SSL}(\mathbf{z}^A, \mathbf{z}^B)]$$

Continual Learning (CL)

- Continual learning tackles the problem of learning tasks sequentially.



- More formally, the objective is to find parameters such as:

$$\operatorname{argmin}_{\theta'} \sum_{t=1}^T \mathbb{E}_{(x, y) \sim \mathcal{D}_t} [\mathcal{L}_{CL}(\mathbf{p}, \mathbf{y})]$$

A new perspective

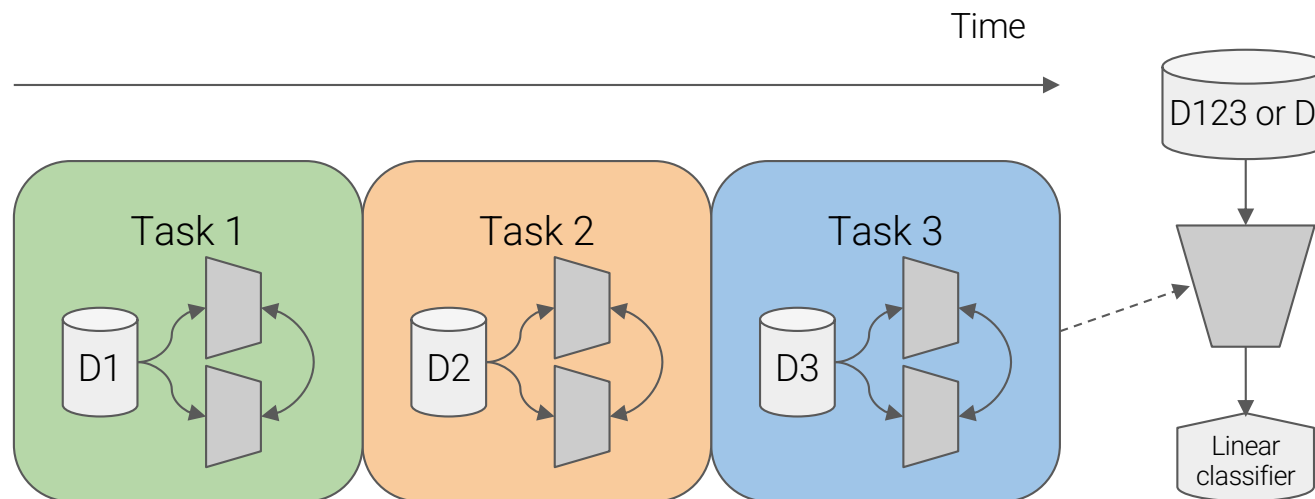
- Most of the literature studies CL with the following assumptions:
 - ◆ Availability of **Labels** (supervised learning)
 - ◆ Focus on learning a **classifier** that solves all tasks

- What if we looked at CL from a **different perspective**?

Assumption	Motivation
No labels	Availability of labels when learning online or sequentially is unlikely
Focus on learning representations	Training a linear classifier or fine-tuning a pre-trained feature extractor on a down-stream task is very simple and inexpensive with modern hardware

Continual Self-Supervised Learning (CSSL)

- Continual Self-Supervised Learning is the problem of learning strong feature extractors from streams of unlabeled data



- The continual pre-training objective is to find parameters such as:

$$\operatorname{argmin}_{\theta} \sum_{t=1}^T \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_t} [\mathcal{L}_{SSL}(\mathbf{z}^A, \mathbf{z}^B)]$$

Continual Self-Supervised Learning (CSSL)

- **Self-Supervised** fine-tuning sometimes **outperforms supervised** fine-tuning
- **Plasticity** of representations is fundamental in CSSL
- SSL benefits from **longer** training, the evolution of representations should not be overly **constrained**
- SSL methods exhibit different losses and feature normalizations that **interfere** with CL regularization losses and vice-versa

Methods	Loss	Equation
SimCLR [13] MoCo [28] NNCLR [19]	InfoNCE	$-\log \frac{\exp(\text{sim}(\mathbf{z}_i^A, \mathbf{z}_i^B)/\tau)}{\sum_{\mathbf{z}_j \in \eta(i)} \exp(\text{sim}(\mathbf{z}_i^A, \mathbf{z}_j)/\tau)} \quad (6)$
BYOL [26] SimSiam [15] VICReg [3]	MSE	$-\ \mathbf{q}^A - \mathbf{z}^B\ _2^2 \quad (7)$
SwAV [7] DCV2 [7] DINO [8]	Cross-entropy	$-\sum_d \mathbf{a}_d^B \log \frac{\exp(\text{sim}(\mathbf{z}^A, \mathbf{c}_d)/\tau)}{\sum_k \exp(\text{sim}(\mathbf{z}^A, \mathbf{c}_k)/\tau)} \quad (8)$
Barlow Twins [58] VICReg [3]	Cross-correlation	$\sum_u (1 - \mathcal{C}_{uv})^2 + \lambda \sum_u \sum_{v \neq u} \mathcal{C}_{uv}^2 \quad (9)$

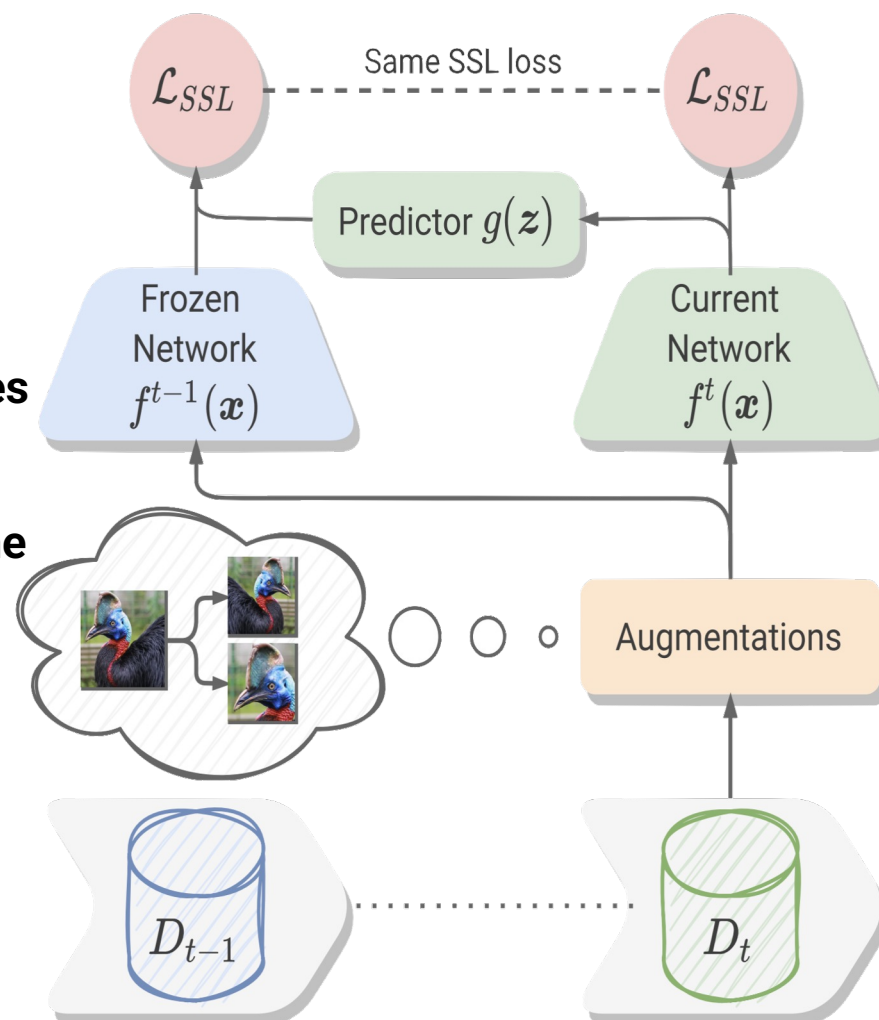
Proposed Method (CaSSLe)

Basically, two simple ideas:

- a **predictor** network that maps the current state of the representations to their past state
- a **family** of adaptable **distillation losses** inherited from the SSL literature

... and an important key insight: use the **same loss** for distillation and representation learning!

$$\begin{aligned}\mathcal{L} &= \mathcal{L}_{SSL}(z^A, z^B) + \mathcal{L}_D(z^A, \bar{z}^A) \\ &= \mathcal{L}_{SSL}(z^A, z^B) + \mathcal{L}_{SSL}(g(z^A), \bar{z}^A)\end{aligned}$$



Experiments

→ We train six SSL models:

- ◆ **Barlow Twins**
- ◆ **SwAV**
- ◆ **BYOL**
- ◆ **VICReg**
- ◆ **MoCoV2+**
- ◆ **SimCLR**

→ We evaluate three CL settings:

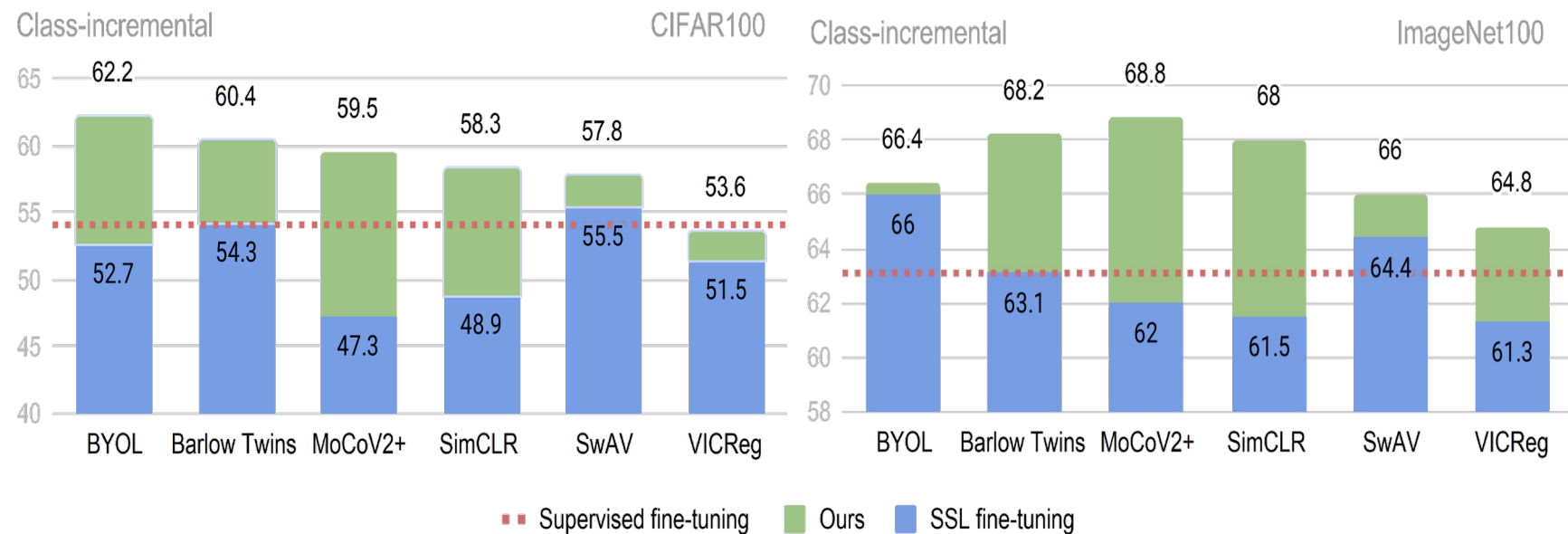
- ◆ **Class-incremental:**
each task contains a new set of classes
- ◆ **Data-incremental:**
each task contains new samples of the same classes
- ◆ **Domain-incremental:**
each task contains new domain

→ On three widely used datasets:

- ◆ **CIFAR100**
- ◆ **ImageNet100**
- ◆ **DomainNet**

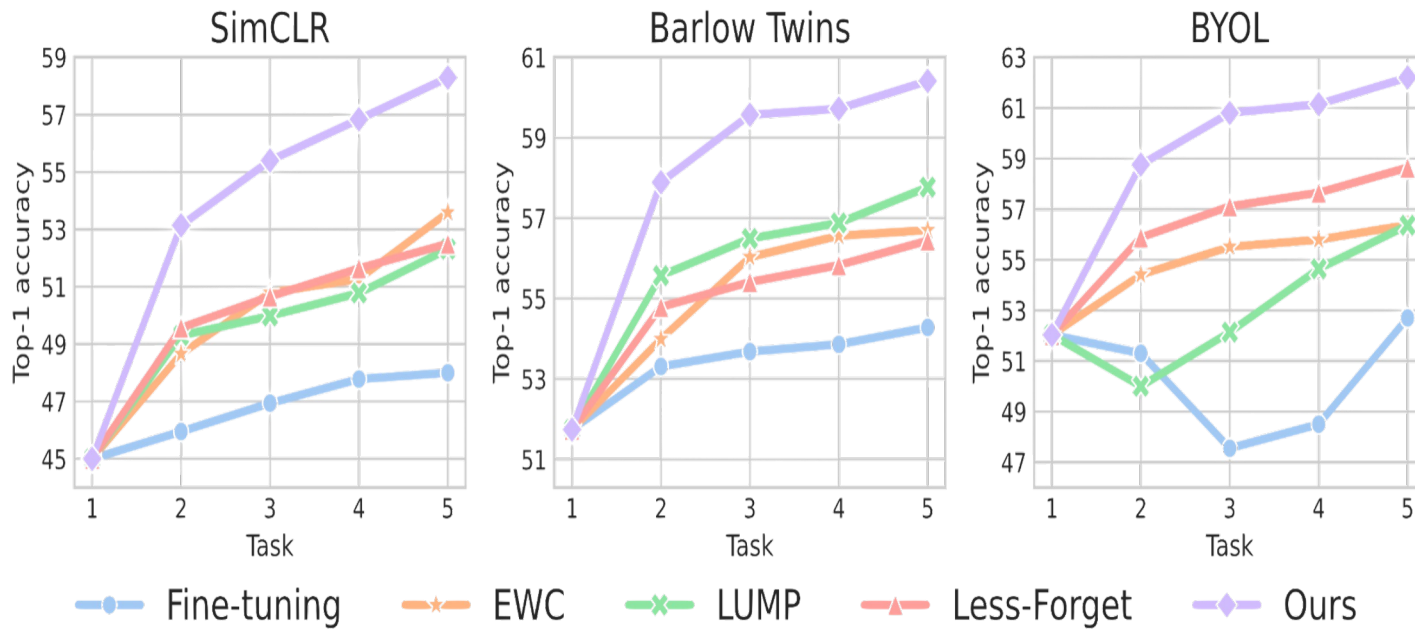
Experiments

→ It turns out that by just using this simple technique you can obtain significant improvements:



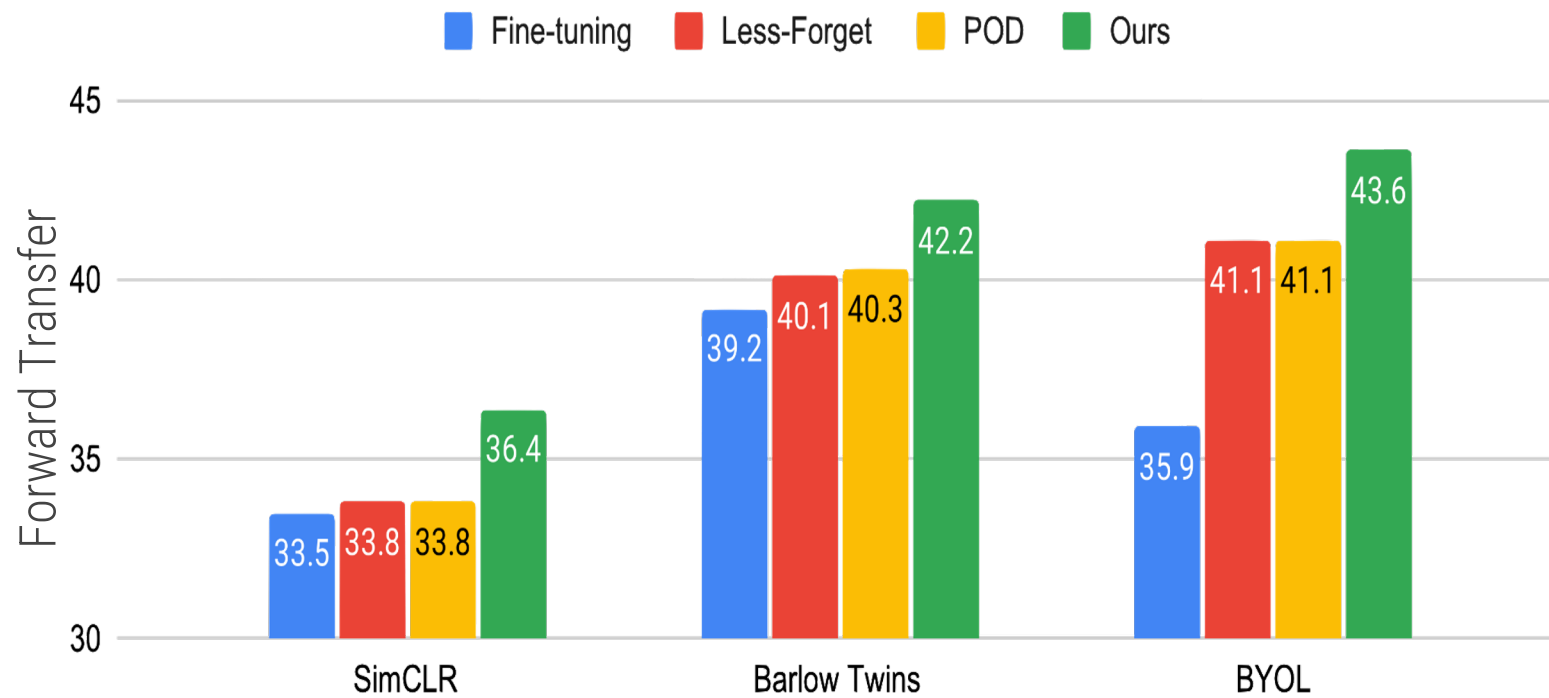
Experiments

→ ...and outperform other CL methods by large margins throughout the whole training trajectory:



Experiments

→ ...and yield better forward transfer:





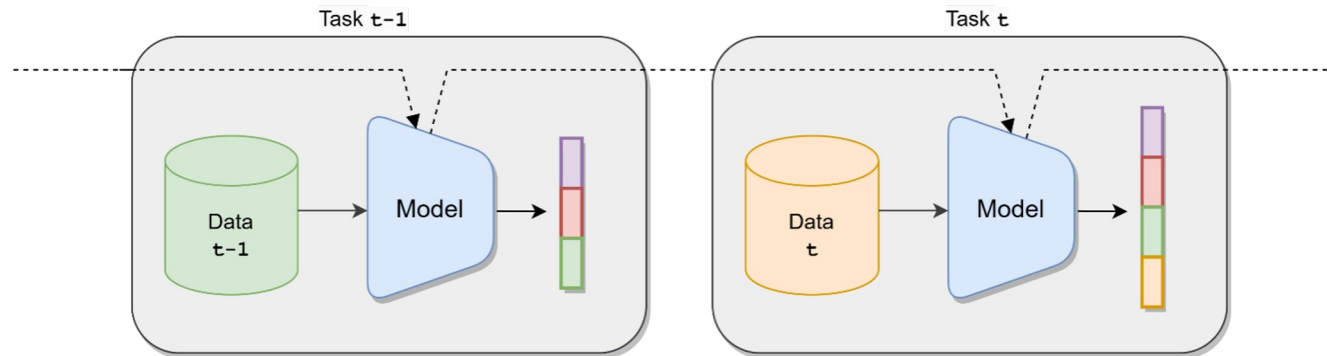
A soft nearest-neighbor framework for continual semi-supervised learning

Zhiqi Kang*, Enrico Fini*, Moin Nabi, Elisa Ricci and Karteek Alahari



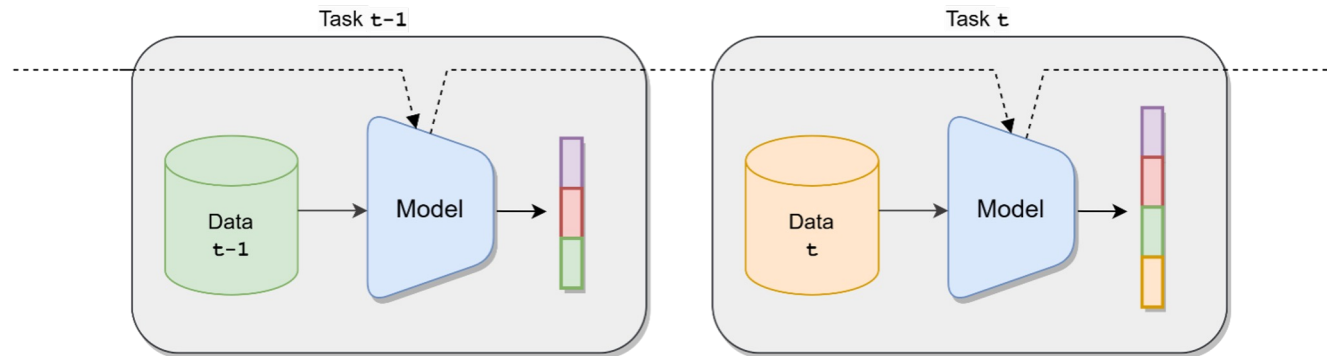
Continual semi-supervised learning

Continual Learning (CL) deals with a dynamic learning scenario.



Continual semi-supervised learning

Continual Learning (CL) deals with a dynamic learning scenario.



Why is the **semi-supervised setting** realistic and interesting?

Annotating the entire dataset (**fully supervised**) can be:

- Expensive
- Impractical

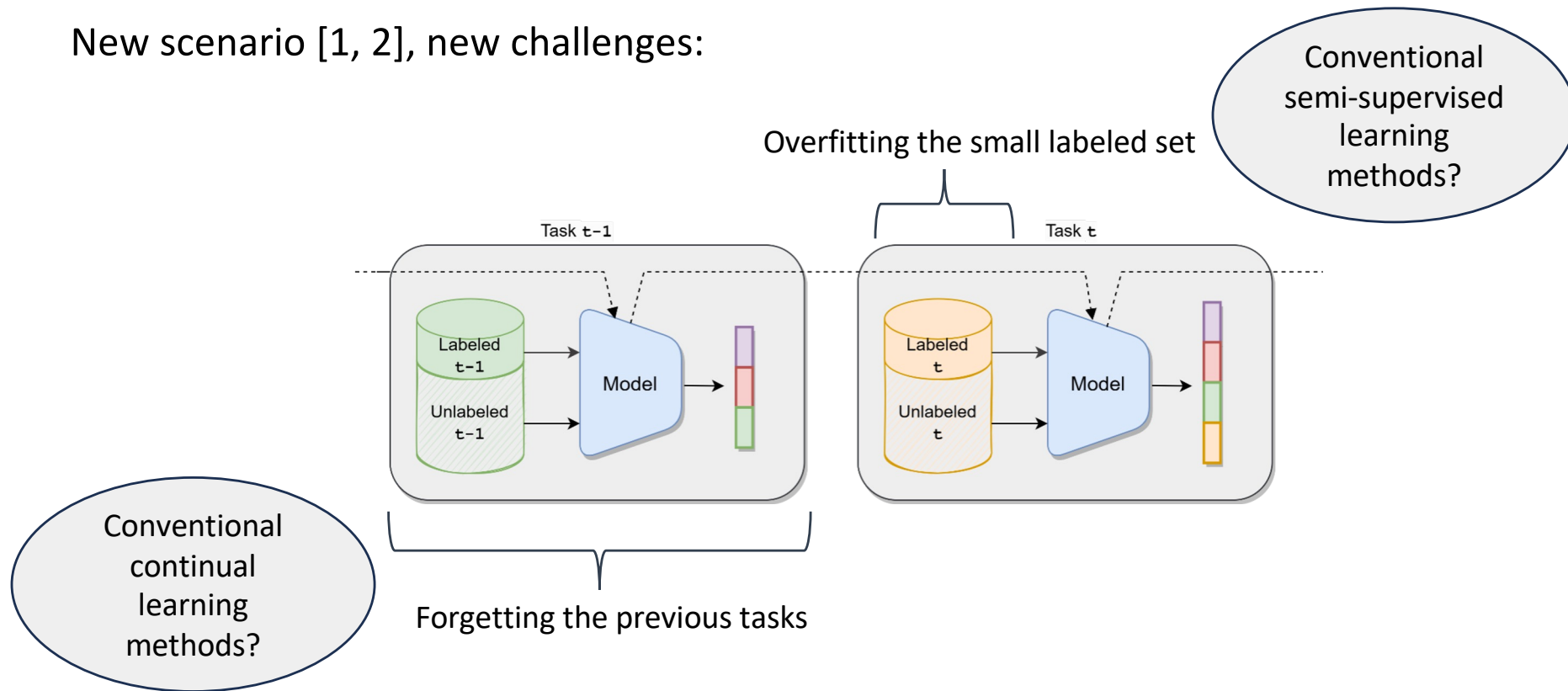
Training on unannotated datasets (**self-supervised**) can be:

- Expensive
- Complex

Semi-supervised setting allows for a good trade-off

Continual semi-supervised learning

New scenario [1, 2], new challenges:



[1] Wang, L., Yang, K., Li, C., Hong, L., Li, Z., and Zhu, J. (2021). Ordisco: Effective and efficient usage of incremental unlabeled data for semi-supervised continual learning. In CVPR.
 [2] Boschini, M., Buzzega, P., Bonicelli, L., Porrello, A., and Calderara, S. (2022). Continual semi-supervised learning through contrastive interpolation consistency. Pattern Recognition Letters.