

Bayesian continual learning and forgetting in neural networks

Djohan Bonnet, Kellian Cottart, Tifenn Hirtzlin, Tarcisius Januel, Thomas Dalgaty, Elisa Vianello, Damien Querlioz

Nature Communications, Accepted

September 16, 2025



◆ Introduction

Motivations and context

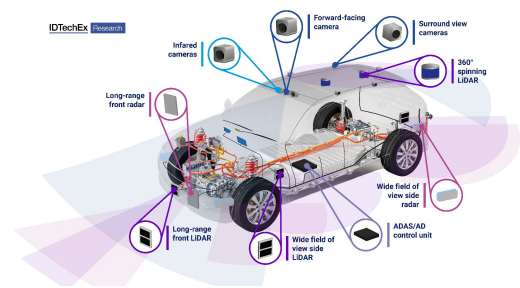
Continual learning

Real-time applications
 (Autonomous Vehicles,
 Healthcare, IoT) receive
unbatched streams of data

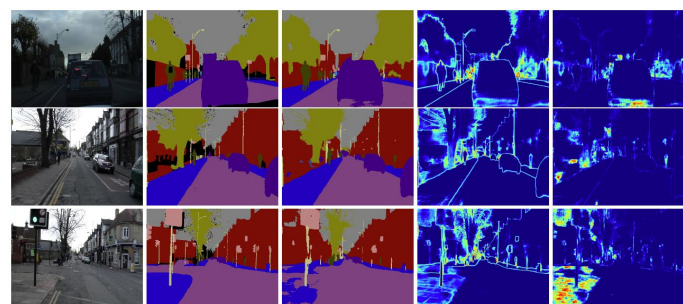
Uncertainties

Safety-critical contexts require
uncertainties to trust the
predictions that are made by
 the model

**This talk: Uncertainties help with
 continual learning**



IDTechEx



◆ Table of contents

1. Catastrophic forgetting & remembering

**2. Regulate forgetting and remembering
through Metaplasticity from Synaptic
Uncertainty**

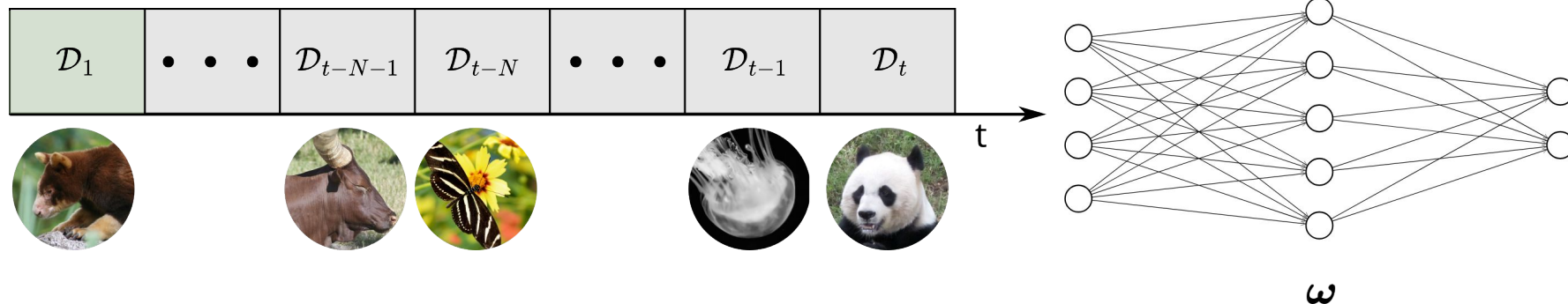
3. Experimental results

Catastrophic forgetting & remembering

◆ Catastrophic forgetting

Forgetting about previously seen data

When deployed in an environment, data distributions may evolve through time.
“Is this animal the one I’m looking for?”

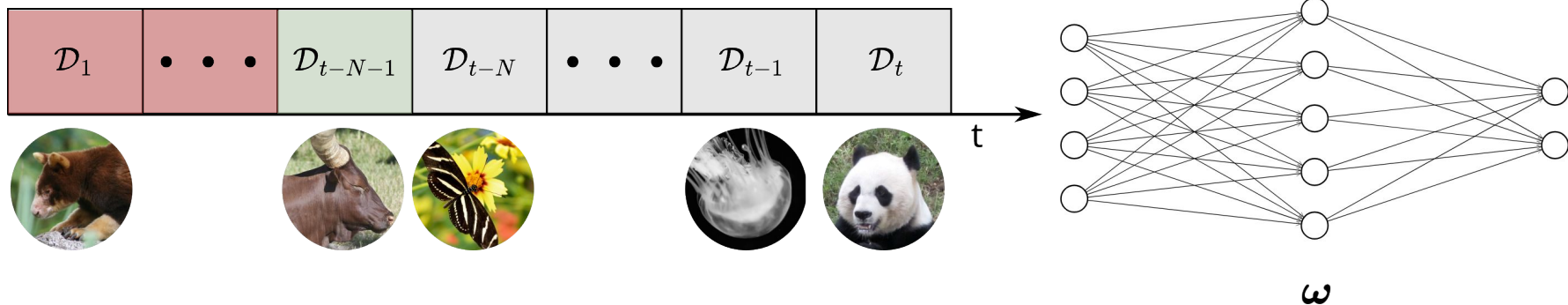


Kirkpatrick, James, et al. "Overcoming catastrophic forgetting in neural networks." *Proceedings of the national academy of sciences* 114.13 (2017): 3521-3526.

◆ Catastrophic forgetting

Forgetting about previously seen data

When deployed in an environment, data distributions may evolve through time.
“Is this animal the one I’m looking for?”

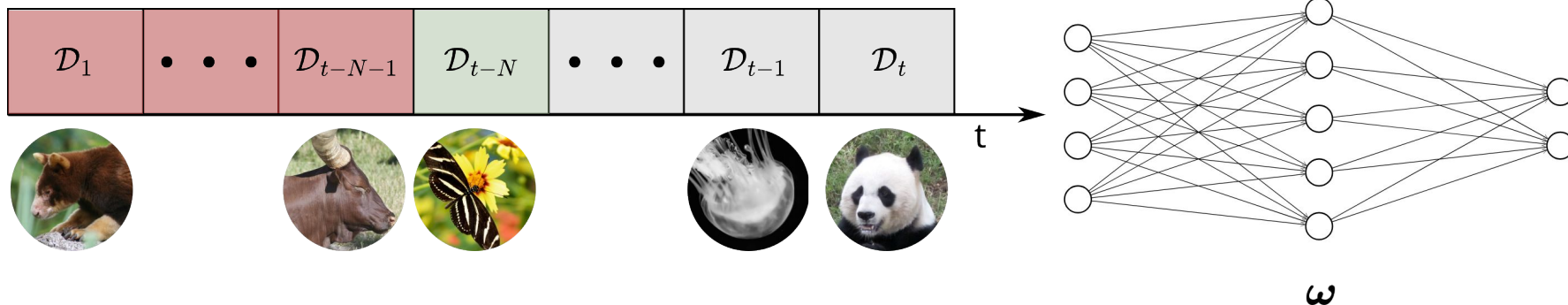


Kirkpatrick, James, et al. "Overcoming catastrophic forgetting in neural networks." *Proceedings of the national academy of sciences* 114.13 (2017): 3521-3526.

◆ Catastrophic forgetting

Forgetting about previously seen data

When deployed in an environment, data distributions may evolve through time.
“Is this animal the one I’m looking for?”

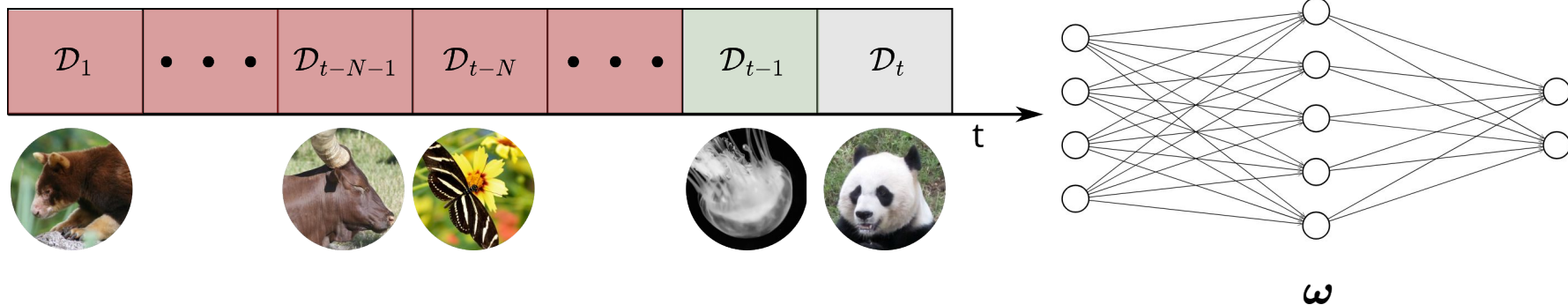


Kirkpatrick, James, et al. "Overcoming catastrophic forgetting in neural networks." *Proceedings of the national academy of sciences* 114.13 (2017): 3521-3526.

◆ Catastrophic forgetting

Forgetting about previously seen data

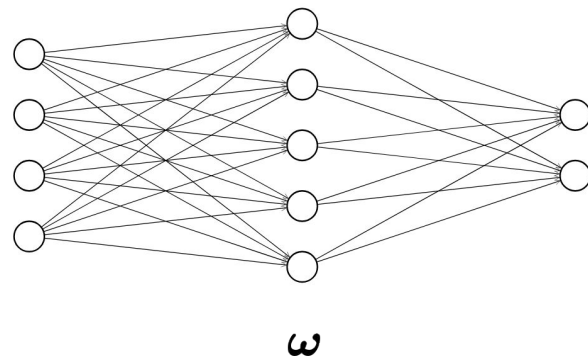
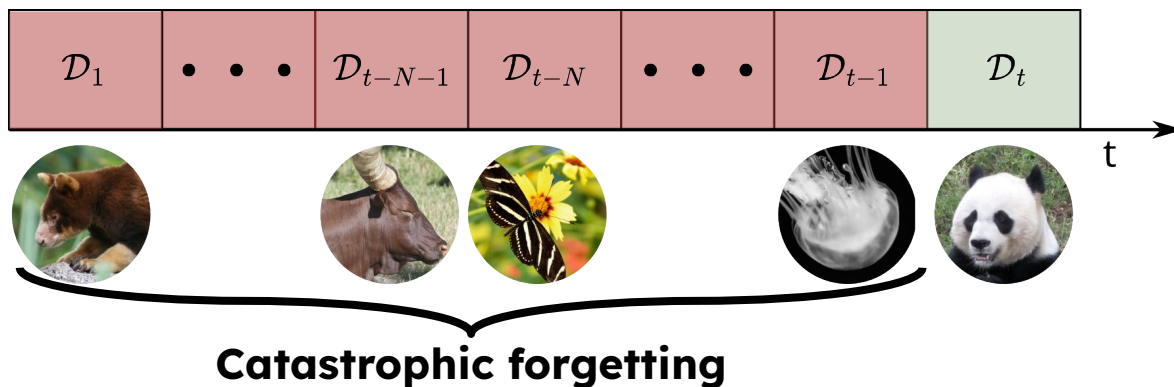
When deployed in an environment, data distributions may evolve through time.
“Is this animal the one I’m looking for?”



◆ Catastrophic forgetting

Forgetting about previously seen data

When deployed in an environment, data distributions may evolve through time.
“Is this animal the one I’m looking for?”



Deep neural network are weak to evolving data distributions and **forget previous representations**

Kirkpatrick, James, et al. "Overcoming catastrophic forgetting in neural networks." (2017)

◆ Mitigating catastrophic forgetting

Literature review for continual learning

Regularization

Optimization

Replay

Representation

Architecture

◆ Mitigating catastrophic forgetting

Literature review for continual learning

Regularization

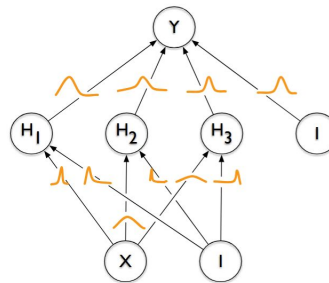
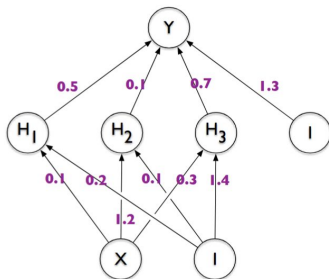
Optimization

Replay

Representation

Architecture

Equipping each weight with a **standard deviation** allows to track uncertainty of each synapse



◆ Mitigating catastrophic forgetting

Literature review for continual learning

Regularization

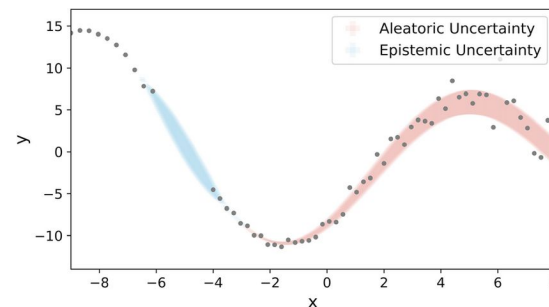
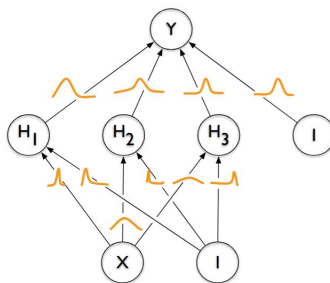
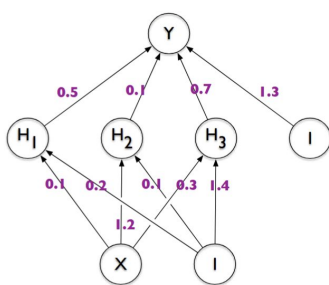
Optimization

Replay

Representation

Architecture

Equipping each weight with a **standard deviation** allows to track **uncertainty of each synapse**



◆ Mitigating catastrophic forgetting

Literature review for continual learning

Regularization

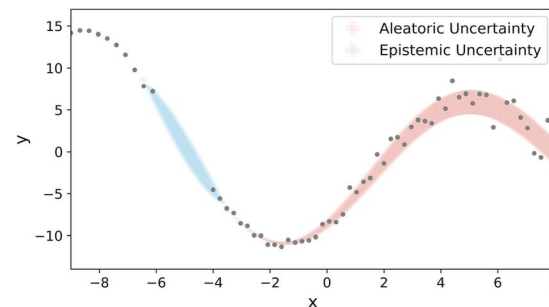
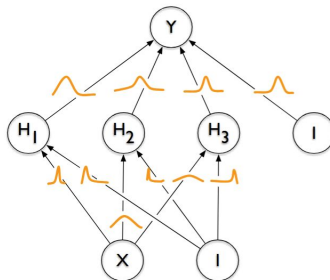
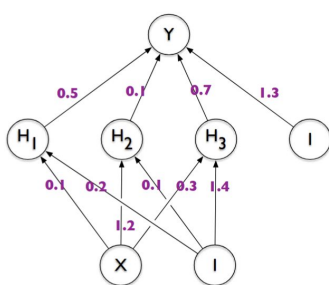
Optimization

Replay

Representation

Architecture

Equipping each weight with a **standard deviation** allows to track **uncertainty of each synapse**



Blundell, Charles, et al. "Weight uncertainty in neural network." 2015.

Conditioning the model on **Bayesian statistics** allows to express the **degree of belief in an event** and **face lifelong learning**

◆ Bayes' Rule against catastrophic forgetting

Updating the neural network based on prior knowledge

$$\underbrace{p(\omega | \mathcal{D}_1, \dots, \mathcal{D}_t)}_{\text{Posterior}} = \frac{\underbrace{p(\mathcal{D}_t | \omega)}_{\text{Likelihood}} \cdot \underbrace{p(\omega | \mathcal{D}_1, \dots, \mathcal{D}_{t-1})}_{\text{Prior}}}{p(\mathcal{D}_t)}.$$

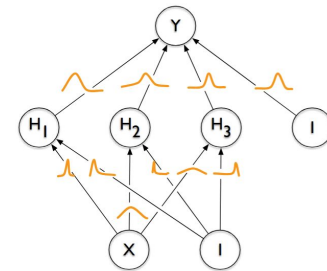
◆ Bayes' Rule against catastrophic forgetting

Updating the neural network based on prior knowledge

$$\underbrace{p(\omega | \mathcal{D}_1, \dots, \mathcal{D}_t)}_{\text{Posterior}} = \frac{\underbrace{p(\mathcal{D}_t | \omega)}_{\text{Likelihood}} \cdot \underbrace{p(\omega | \mathcal{D}_1, \dots, \mathcal{D}_{t-1})}_{\text{Prior}}}{p(\mathcal{D}_t)}.$$

Variational Mean-field Gaussian

$$\underbrace{q_{\theta_t}(\omega)} = \prod_i \mathcal{N}(\omega_i; \mu_{t,i}, \sigma_{t,i}^2)$$



$$\theta_t = (\mu_t, \sigma_t)$$

◆ Bayes' Rule against catastrophic forgetting

Updating the neural network based on prior knowledge

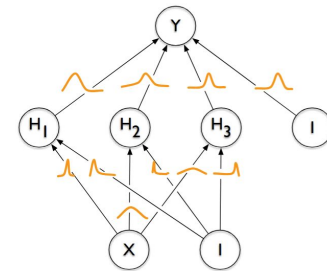
$$\underbrace{p(\omega | \mathcal{D}_1, \dots, \mathcal{D}_t)}_{\text{Posterior}} = \frac{\underbrace{p(\mathcal{D}_t | \omega)}_{\text{Likelihood}} \cdot \underbrace{p(\omega | \mathcal{D}_1, \dots, \mathcal{D}_{t-1})}_{\text{Prior}}}{p(\mathcal{D}_t)}.$$

Minimize using
Kullback-Leibler
Divergence



Variational Mean-field
Gaussian

$$q_{\theta_t}(\omega) = \prod_i \mathcal{N}(\omega_i; \mu_{t,i}, \sigma_{t,i}^2)$$



$$\theta_t = (\mu_t, \sigma_t)$$

◆ Bayes' Rule against catastrophic forgetting

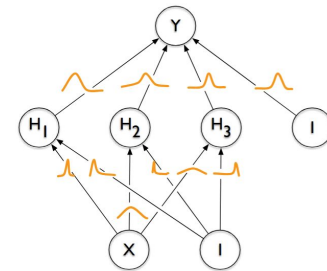
Updating the neural network based on prior knowledge

$$\underbrace{p(\omega | \mathcal{D}_1, \dots, \mathcal{D}_t)}_{\text{Posterior}} = \frac{\underbrace{p(\mathcal{D}_t | \omega)}_{\text{Likelihood}} \cdot \underbrace{p(\omega | \mathcal{D}_1, \dots, \mathcal{D}_{t-1})}_{\text{Prior}}}{p(\mathcal{D}_t)}.$$

Minimize using
Kullback-Leibler
Divergence

Variational Mean-field
Gaussian

$$q_{\theta_t}(\omega) = \prod_i \mathcal{N}(\omega_i; \mu_{t,i}, \sigma_{t,i}^2)$$



$$\theta_t = (\mu_t, \sigma_t)$$

Neural networks are mitigating catastrophic forgetting using Bayes' Rule

Zeno, Chen, et al. "Task Agnostic Continual Learning Using Online Variational Bayes." (2022)

◆ Synaptic metaplasticity through variance

Bayesian neural networks for lifelong learning

FOO-VB Diagonal parameter updates

$$\Delta\mu = -\sigma_{t-1}^2 \frac{\partial \mathcal{C}_t}{\partial \mu_{t-1}} \quad \Delta\sigma = -\frac{\sigma_{t-1}^2}{2} \frac{\partial \mathcal{C}_t}{\partial \sigma_{t-1}} + \sigma_{t-1} \left(\sqrt{1 + \left(\frac{\sigma_{t-1}}{2} \frac{\partial \mathcal{C}_t}{\partial \sigma_{t-1}} \right)^2} - 1 \right)$$

$$\frac{\partial \mathcal{C}_t}{\partial \mu} = \mathbb{E}_\epsilon \left[\frac{\partial \mathcal{L}_t(\omega)}{\partial \omega} \right]$$

$$\frac{\partial \mathcal{C}_t}{\partial \sigma} = \mathbb{E}_\epsilon \left[\frac{\partial \mathcal{L}_t(\omega)}{\partial \omega} \times \epsilon \right]$$

◆ Synaptic metaplasticity through variance

Bayesian neural networks for lifelong learning

FOO-VB Diagonal parameter updates

$$\Delta\mu = -\sigma_{t-1}^2 \frac{\partial \mathcal{C}_t}{\partial \mu_{t-1}} \quad \Delta\sigma = -\frac{\sigma_{t-1}^2}{2} \frac{\partial \mathcal{C}_t}{\partial \sigma_{t-1}} + \sigma_{t-1} \left(\sqrt{1 + \left(\frac{\sigma_{t-1}}{2} \frac{\partial \mathcal{C}_t}{\partial \sigma_{t-1}} \right)^2} - 1 \right)$$

Metaplasticity: each synapse adjusts its own learning rate

$$\frac{\partial \mathcal{C}_t}{\partial \mu} = \mathbb{E}_\epsilon \left[\frac{\partial \mathcal{L}_t(\omega)}{\partial \omega} \right]$$

$$\frac{\partial \mathcal{C}_t}{\partial \sigma} = \mathbb{E}_\epsilon \left[\frac{\partial \mathcal{L}_t(\omega)}{\partial \omega} \times \epsilon \right]$$

◆ Synaptic metaplasticity through variance

Bayesian neural networks for lifelong learning

FOO-VB Diagonal parameter updates

$$\Delta\mu = -\sigma_{t-1}^2 \frac{\partial \mathcal{C}_t}{\partial \mu_{t-1}} \quad \Delta\sigma = -\frac{\sigma_{t-1}^2}{2} \frac{\partial \mathcal{C}_t}{\partial \sigma_{t-1}} + \sigma_{t-1} \left(\sqrt{1 + \left(\frac{\sigma_{t-1}}{2} \frac{\partial \mathcal{C}_t}{\partial \sigma_{t-1}} \right)^2} - 1 \right)$$

Metaplasticity: each synapse adjusts its own learning rate

$$\frac{\partial \mathcal{C}_t}{\partial \mu} = \mathbb{E}_\epsilon \left[\frac{\partial \mathcal{L}_t(\omega)}{\partial \omega} \right]$$

$$\frac{\partial \mathcal{C}_t}{\partial \sigma} = \mathbb{E}_\epsilon \left[\frac{\partial \mathcal{L}_t(\omega)}{\partial \omega} \times \epsilon \right]$$

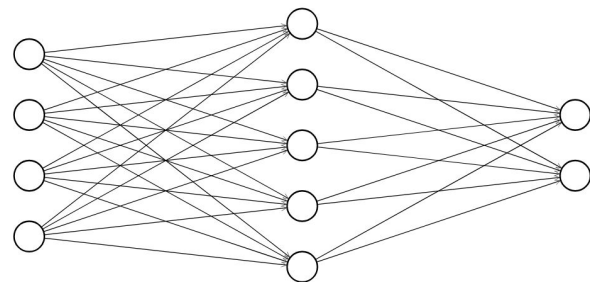
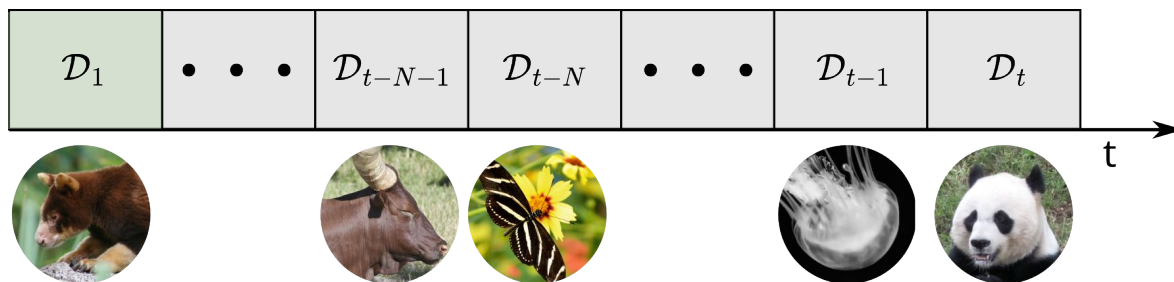
- ◆ Variance tracks **uncertainty**: the **higher** the variance, the **larger** the learning rate
- ◆ When synapses are **very certain**, the mean cannot change: **information is retained**

Bayes' Rule applied to multiple datasets leads inherently metaplastic updates

◆ Catastrophic remembering

Inability to learn induced by too much data

Using Bayesian neural networks adapted to lifelong learning, catastrophic forgetting is prevented.



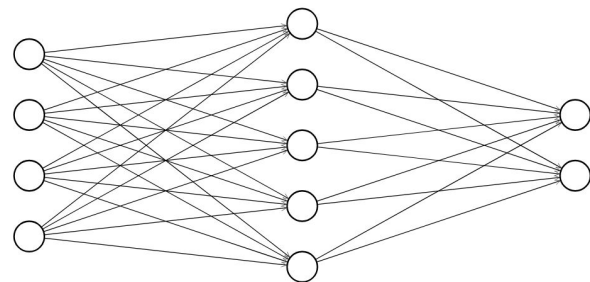
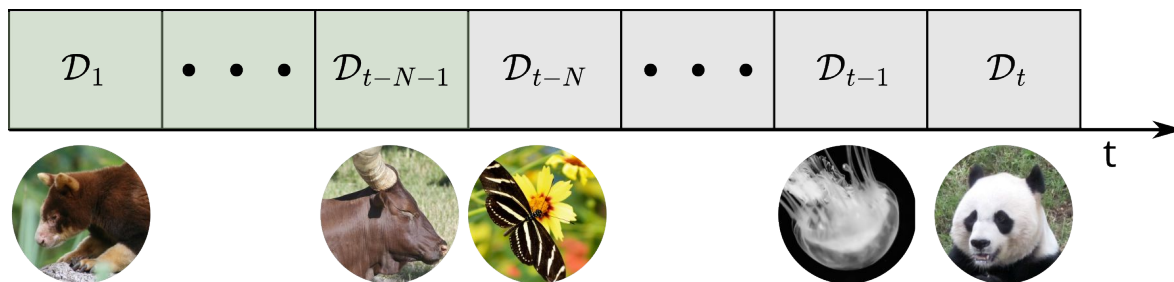
$$q_{\theta_t}(\omega) \approx p(\omega | \mathcal{D}_1, \dots, \mathcal{D}_t)$$

Kaushik, Prakhar, et al.
 "Understanding Catastrophic Forgetting and Remembering in Continual Learning with Optimal Relevance Mapping."

◆ Catastrophic remembering

Inability to learn induced by too much data

Using Bayesian neural networks adapted to lifelong learning, catastrophic forgetting is prevented.



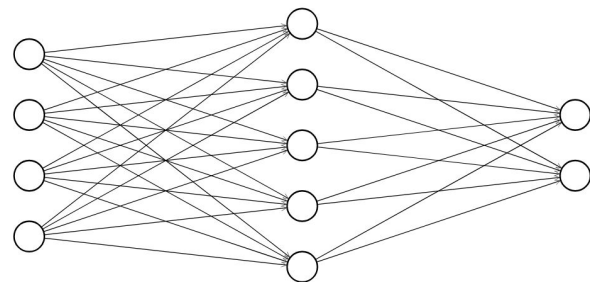
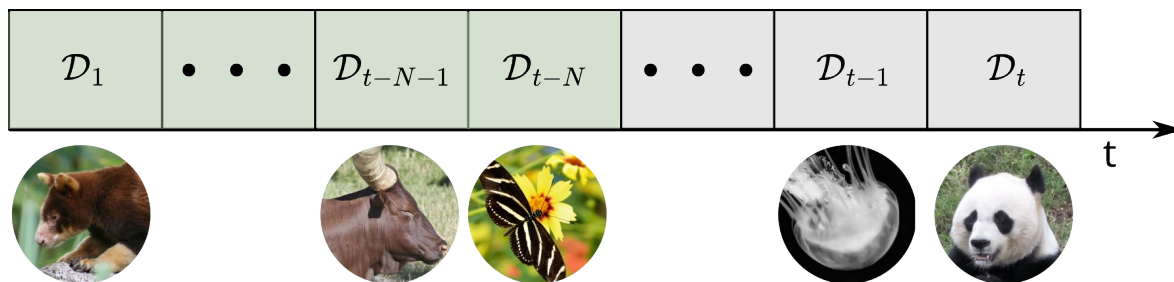
$$q_{\theta_t}(\omega) \approx p(\omega | \mathcal{D}_1, \dots, \mathcal{D}_t)$$

Kaushik, Prakhar, et al.
 "Understanding Catastrophic Forgetting and Remembering in Continual Learning with Optimal Relevance Mapping."

◆ Catastrophic remembering

Inability to learn induced by too much data

Using Bayesian neural networks adapted to lifelong learning, catastrophic forgetting is prevented.



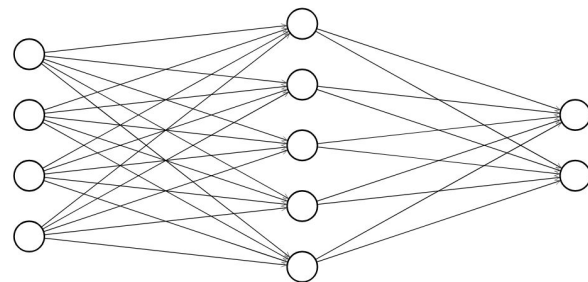
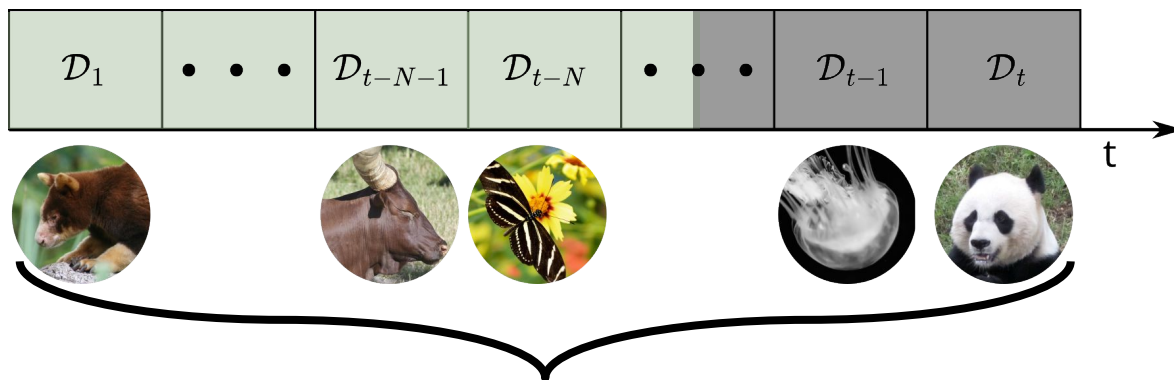
$$q_{\theta_t}(\omega) \approx p(\omega | \mathcal{D}_1, \dots, \mathcal{D}_t)$$

Kaushik, Prakhar, et al.
 "Understanding Catastrophic Forgetting and Remembering in Continual Learning with Optimal Relevance Mapping."

◆ Catastrophic remembering

Inability to learn induced by too much data

Using Bayesian neural networks adapted to lifelong learning, catastrophic forgetting is prevented.



$$q_{\theta_t}(\omega) \approx p(\omega | \mathcal{D}_1, \dots, \mathcal{D}_t)$$

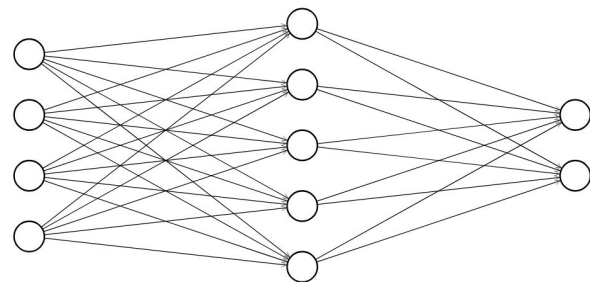
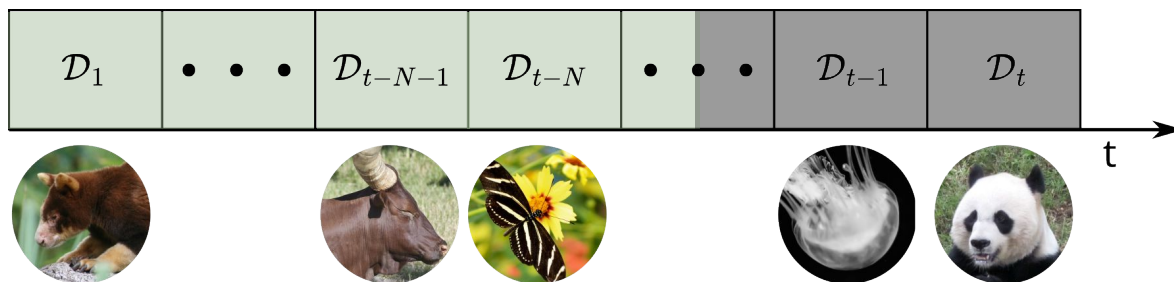
Catastrophic remembering

Kaushik, Prakhar, et al.
 "Understanding Catastrophic Forgetting and Remembering in Continual Learning with Optimal Relevance Mapping."

◆ Catastrophic remembering

Inability to learn induced by too much data

Using Bayesian neural networks adapted to lifelong learning, catastrophic forgetting is prevented.



$$q_{\theta_t}(\omega) \approx p(\omega | \mathcal{D}_1, \dots, \mathcal{D}_t)$$

Kaushik, Prakhar, et al.
 "Understanding Catastrophic Forgetting and Remembering in Continual Learning with Optimal Relevance Mapping."

However, remembering all datasets prevents learning

◆ Research direction

Improving and furthering Bayesian neural networks

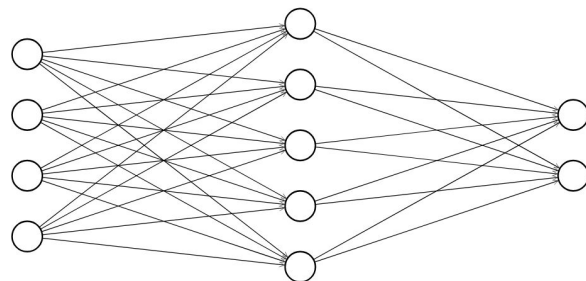
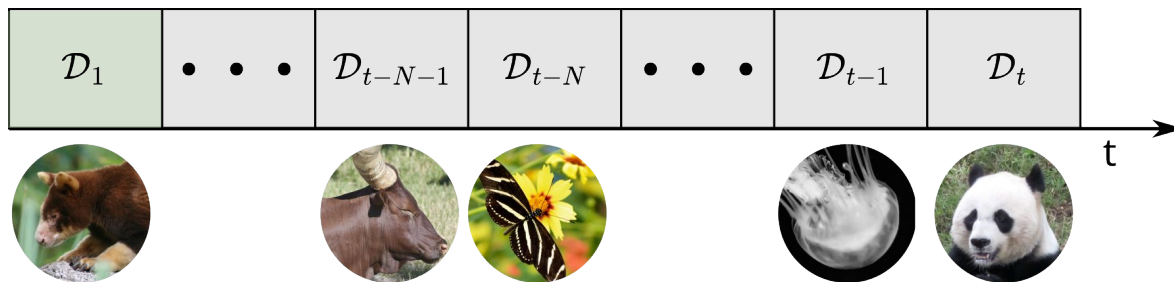
- ◆ **What triggers catastrophic remembering?**
- ◆ **How to maintain optimal discriminative capabilities through time?**
- ◆ **Is it possible to avoid both catastrophic remembering and catastrophic forgetting?**

Regulate forgetting and remembering through Metaplasticity from Synaptic Uncertainty

◆ Bayesian continual learning and forgetting

Avoiding catastrophic remembering and vanishing uncertainties

We consider a memory window of size N , corresponding to the maximum number of datasets presented to retain



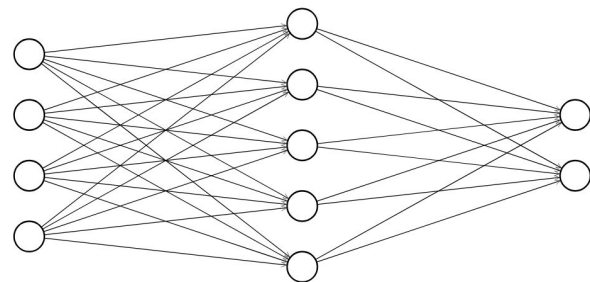
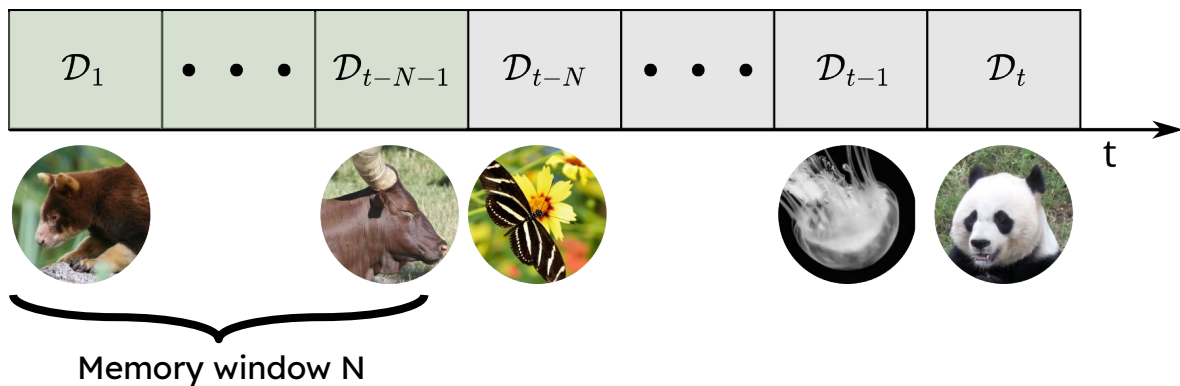
$$q_{\theta_t}(\omega) \approx p(\omega | \mathcal{D}_{t-N}, \dots, \mathcal{D}_t)$$

Memory window N

◆ Bayesian continual learning and forgetting

Avoiding catastrophic remembering and vanishing uncertainties

We consider a memory window of size N , corresponding to the maximum number of datasets presented to retain

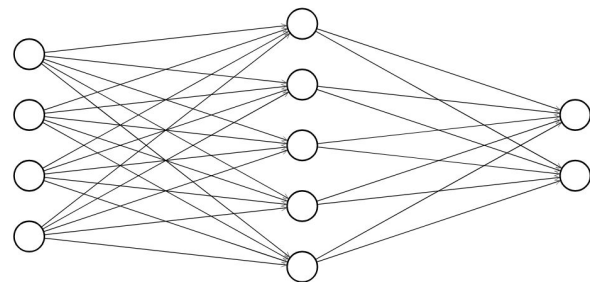
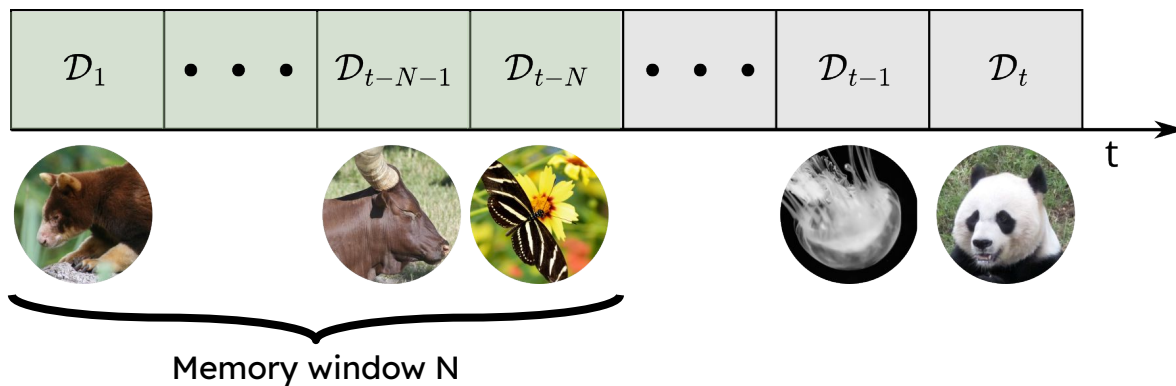


$$q_{\theta_t}(\omega) \approx p(\omega | \mathcal{D}_{t-N}, \dots, \mathcal{D}_t)$$

◆ Bayesian continual learning and forgetting

Avoiding catastrophic remembering and vanishing uncertainties

We consider a memory window of size N , corresponding to the maximum number of datasets presented to retain

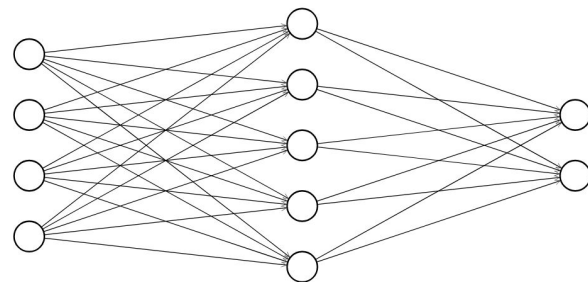
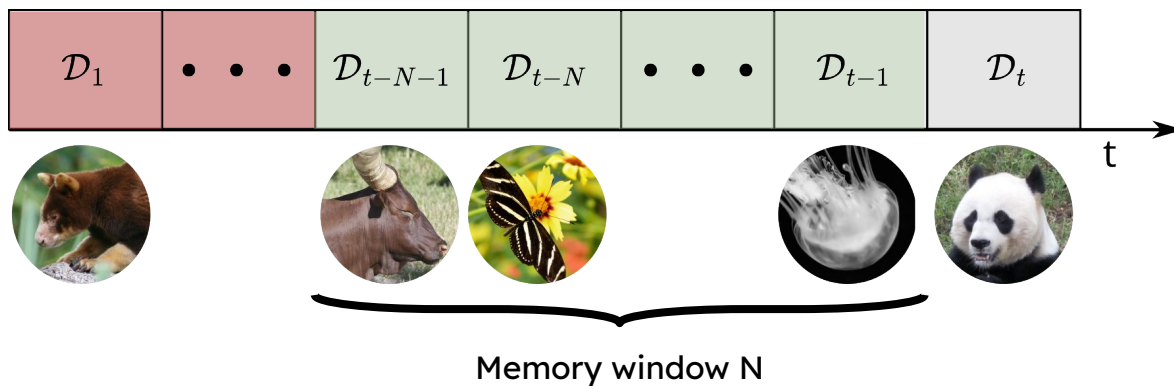


$$q_{\theta_t}(\omega) \approx p(\omega | \mathcal{D}_{t-N}, \dots, \mathcal{D}_t)$$

◆ Bayesian continual learning and forgetting

Avoiding catastrophic remembering and vanishing uncertainties

We consider a memory window of size N , corresponding to the maximum number of datasets presented to retain

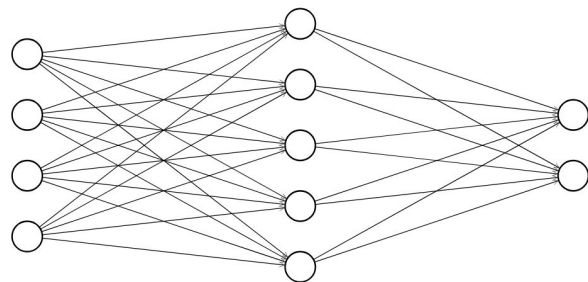
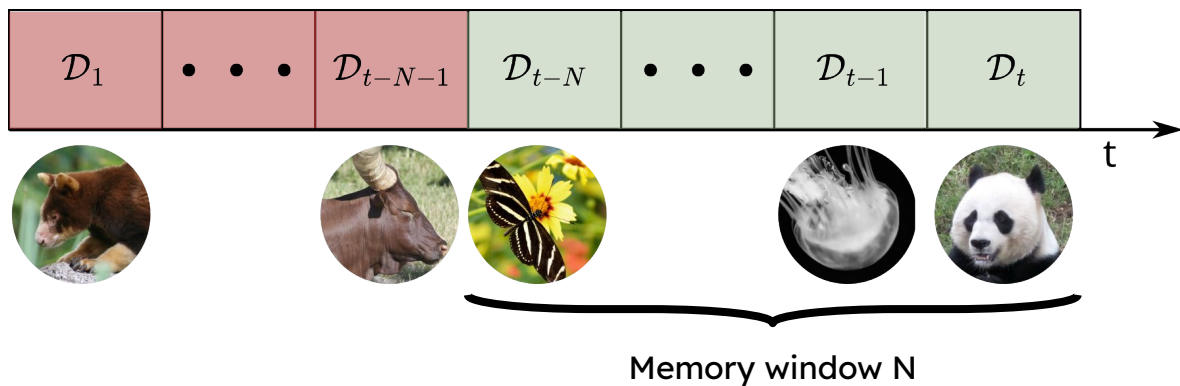


$$q_{\theta_t}(\omega) \approx p(\omega | \mathcal{D}_{t-N}, \dots, \mathcal{D}_t)$$

◆ Bayesian continual learning and forgetting

Avoiding catastrophic remembering and vanishing uncertainties

We consider a memory window of size N , corresponding to the maximum number of datasets presented to retain



$$q_{\theta_t}(\omega) \approx p(\omega | \mathcal{D}_{t-N}, \dots, \mathcal{D}_t)$$

Information out of the memory window is gradually forgotten

◆ Bayesian continual learning and forgetting

Avoiding catastrophic remembering and vanishing uncertainties

If we desire to learn only **N datasets**, Bayes' rule yields

$$\underbrace{p(\omega | \mathcal{D}_{t-N}, \dots, \mathcal{D}_t)}_{\text{Posterior}} = \underbrace{\frac{p(\mathcal{D}_t | \omega) \cdot p(\omega | \mathcal{D}_{t-N-1}, \dots, \mathcal{D}_{t-1})}{p(\mathcal{D}_t)}}_{\text{Learning}} \cdot \underbrace{\frac{p(\mathcal{D}_{t-N-1})}{p(\mathcal{D}_{t-N-1} | \omega)}}_{\text{Forgetting}}.$$

◆ Bayesian continual learning and forgetting

Avoiding catastrophic remembering and vanishing uncertainties

If we desire to learn only **N datasets**, Bayes' rule yields

$$\underbrace{p(\omega | \mathcal{D}_{t-N}, \dots, \mathcal{D}_t)}_{\text{Posterior}} = \underbrace{\frac{p(\mathcal{D}_t | \omega) \cdot p(\omega | \mathcal{D}_{t-N-1}, \dots, \mathcal{D}_{t-1})}{p(\mathcal{D}_t)}}_{\text{Learning}} \cdot \underbrace{\frac{p(\mathcal{D}_{t-N-1})}{p(\mathcal{D}_{t-N-1} | \omega)}}_{\text{Forgetting}}.$$

But we don't have access to that..

◆ Bayesian continual learning and forgetting

Avoiding catastrophic remembering and vanishing uncertainties

If we desire to learn only **N datasets**, Bayes' rule yields

$$\underbrace{p(\omega | \mathcal{D}_{t-N}, \dots, \mathcal{D}_t)}_{\text{Posterior}} = \underbrace{\frac{p(\mathcal{D}_t | \omega) \cdot p(\omega | \mathcal{D}_{t-N-1}, \dots, \mathcal{D}_{t-1})}{p(\mathcal{D}_t)}}_{\text{Learning}} \cdot \underbrace{\frac{p(\mathcal{D}_{t-N-1})}{p(\mathcal{D}_{t-N-1} | \omega)}}_{\text{Forgetting}}. \quad \text{But we don't have access to that..}$$

Assuming each dataset has **equal marginal likelihood** and a **Gaussian prior, posterior and likelihood**, forgetting depends on the prior and the current state

$$p(\mathcal{D}_{t-N-1}, \dots, \mathcal{D}_{t-1} | \omega) = \prod_{i=t-N-1}^{t-1} p(\mathcal{D}_i | \omega) = [p(\mathcal{D}_{t-N-1} | \omega)]^N \propto \mathcal{N}(\omega; \mu_{L_{t-1}}, \text{diag}(\sigma_{L_{t-1}}^2))$$

◆ Bayesian continual learning and forgetting

Avoiding catastrophic remembering and vanishing uncertainties

If we desire to learn only **N datasets**, Bayes' rule yields

$$\underbrace{p(\omega | \mathcal{D}_{t-N}, \dots, \mathcal{D}_t)}_{\text{Posterior}} = \underbrace{\frac{p(\mathcal{D}_t | \omega) \cdot p(\omega | \mathcal{D}_{t-N-1}, \dots, \mathcal{D}_{t-1})}{p(\mathcal{D}_t)}}_{\text{Learning}} \cdot \underbrace{\frac{p(\mathcal{D}_{t-N-1})}{p(\mathcal{D}_{t-N-1} | \omega)}}_{\text{Forgetting}}. \quad \text{But we don't have access to that..}$$

Assuming each dataset has **equal marginal likelihood** and a **Gaussian prior, posterior and likelihood**, forgetting depends on the prior and the current state

$$p(\mathcal{D}_{t-N-1}, \dots, \mathcal{D}_{t-1} | \omega) = \prod_{i=t-N-1}^{t-1} p(\mathcal{D}_i | \omega) = [p(\mathcal{D}_{t-N-1} | \omega)]^N \propto \mathcal{N}(\omega; \mu_{L_{t-1}}, \text{diag}(\sigma_{L_{t-1}}^2))$$

Learning and forgetting is formalized through a truncated posterior distribution

◆ Metaplasticity from Synaptic Uncertainty

Synapse-wise regularization based on standard deviation

FOO-VB Diagonal

$$\Delta\mu = -\sigma_{t-1}^2 \frac{\partial \mathcal{C}_t}{\partial \mu_{t-1}}$$

$$\Delta\sigma = -\frac{\sigma_{t-1}^2}{2} \frac{\partial \mathcal{C}_t}{\partial \sigma_{t-1}} + \sigma_{t-1} \left(\sqrt{1 + \left(\frac{\sigma_{t-1}}{2} \frac{\partial \mathcal{C}_t}{\partial \sigma_{t-1}} \right)^2} - 1 \right)$$

MESU

$$\Delta\mu = -\sigma_{t-1}^2 \frac{\partial \mathcal{C}_t}{\partial \mu_{t-1}} + \frac{\sigma_{t-1}^2}{N \sigma_{\text{prior}}^2} (\mu_{\text{prior}} - \mu_{t-1})$$

$$\Delta\sigma = -\frac{\sigma_{t-1}^2}{2} \frac{\partial \mathcal{C}_t}{\partial \sigma_{t-1}} + \frac{\sigma_{t-1}}{2 N \sigma_{\text{prior}}^2} (\sigma_{\text{prior}}^2 - \sigma_{t-1}^2)$$

◆ Metaplasticity from Synaptic Uncertainty

Synapse-wise regularization based on standard deviation

FOO-VB Diagonal

$$\Delta\mu = -\sigma_{t-1}^2 \frac{\partial \mathcal{C}_t}{\partial \mu_{t-1}} \quad \longrightarrow \quad \mathbf{0}$$

$$\Delta\sigma = -\frac{\sigma_{t-1}^2}{2} \frac{\partial \mathcal{C}_t}{\partial \sigma_{t-1}} + \sigma_{t-1} \left(\sqrt{1 + \left(\frac{\sigma_{t-1}}{2} \frac{\partial \mathcal{C}_t}{\partial \sigma_{t-1}} \right)^2} - 1 \right)$$

MESU

$$\Delta\mu = -\sigma_{t-1}^2 \frac{\partial \mathcal{C}_t}{\partial \mu_{t-1}} + \frac{\sigma_{t-1}^2}{N \sigma_{\text{prior}}^2} (\mu_{\text{prior}} - \mu_{t-1})$$

$$\Delta\sigma = -\frac{\sigma_{t-1}^2}{2} \frac{\partial \mathcal{C}_t}{\partial \sigma_{t-1}} + \frac{\sigma_{t-1}}{2 N \sigma_{\text{prior}}^2} (\sigma_{\text{prior}}^2 - \sigma_{t-1}^2)$$

◆ Metaplasticity from Synaptic Uncertainty

Synapse-wise regularization based on standard deviation

FOO-VB Diagonal

$$\Delta\mu = -\sigma_{t-1}^2 \frac{\partial \mathcal{C}_t}{\partial \mu_{t-1}} \quad \longrightarrow \quad \mathbf{0}$$

$$\Delta\sigma = -\frac{\sigma_{t-1}^2}{2} \frac{\partial \mathcal{C}_t}{\partial \sigma_{t-1}} + \sigma_{t-1} \left(\sqrt{1 + \left(\frac{\sigma_{t-1}}{2} \frac{\partial \mathcal{C}_t}{\partial \sigma_{t-1}} \right)^2} - 1 \right)$$

MESU

$$\Delta\mu = -\sigma_{t-1}^2 \frac{\partial \mathcal{C}_t}{\partial \mu_{t-1}} + \frac{\sigma_{t-1}^2}{N \sigma_{\text{prior}}^2} (\mu_{\text{prior}} - \mu_{t-1})$$

$$\Delta\sigma = -\frac{\sigma_{t-1}^2}{2} \frac{\partial \mathcal{C}_t}{\partial \sigma_{t-1}} + \frac{\sigma_{t-1}}{2 N \sigma_{\text{prior}}^2} (\sigma_{\text{prior}}^2 - \sigma_{t-1}^2)$$

Forgetting gradually bring synapses that have the highest standard deviation back to their prior

◆ Main claims

What is MESU capable of doing?

- ◆ **MESU prevents catastrophic remembering and rigidity by gradually forgetting information**
- ◆ **MESU conserves OOD detection abilities through time, expressing high uncertainty**
- ◆ **MESU mitigates catastrophic forgetting within the memory window**

Experimental results

◆ Estimating uncertainties

The neural network is a probability distribution

$$\begin{aligned} \mathcal{I}(\omega, y \mid \mathcal{D}, x) &= H[p(y \mid x, \mathcal{D})] - \mathbb{E}_{p(\omega \mid \mathcal{D})} H[p(y \mid x, \omega)], \\ \text{EU} &= \text{TU} - \text{AU}. \end{aligned}$$

◆ Estimating uncertainties

The neural network is a probability distribution

$$\begin{aligned} \mathcal{I}(\omega, y \mid \mathcal{D}, x) &= H[p(y \mid x, \mathcal{D})] - \mathbb{E}_{p(\omega \mid \mathcal{D})} H[p(y \mid x, \omega)], \\ \text{EU} &= \text{TU} - \text{AU}. \end{aligned}$$

Epistemic

- ◆ Computes the mutual information, the **uncertainty of the model given the whole data**

◆ Estimating uncertainties

The neural network is a probability distribution

$$\begin{aligned} \mathcal{I}(\omega, y \mid \mathcal{D}, x) &= H[p(y \mid x, \mathcal{D})] - \mathbb{E}_{p(\omega \mid \mathcal{D})} H[p(y \mid x, \omega)], \\ \text{EU} &= \text{TU} - \text{AU}. \end{aligned}$$

Epistemic

- ◆ Computes the mutual information, the **uncertainty of the model given the whole data**

Aleatoric

- ◆ Computes the expected **uncertainty of the output given the model and the input**

◆ Estimating uncertainties

The neural network is a probability distribution

$$\begin{aligned} \mathcal{I}(\omega, y \mid \mathcal{D}, x) &= H[p(y \mid x, \mathcal{D})] - \mathbb{E}_{p(\omega \mid \mathcal{D})} H[p(y \mid x, \omega)], \\ \text{EU} &= \text{TU} - \text{AU}. \end{aligned}$$

Epistemic

- ◆ Computes the mutual information, the **uncertainty of the model given the whole data**

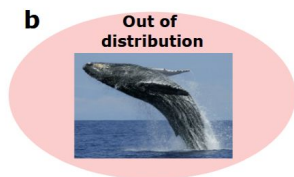
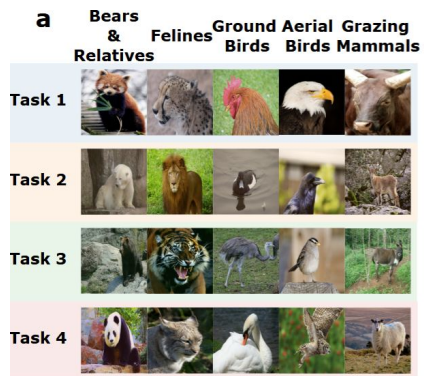
Aleatoric

- ◆ Computes the expected **uncertainty of the output given the model and the input**

Bayesian neural networks generate both aleatoric and epistemic uncertainties to evaluate out of distribution (OOD) data

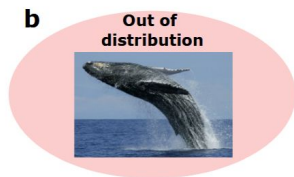
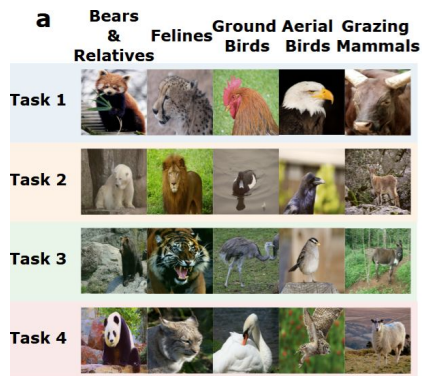
◆ Mitigating catastrophic forgetting

20 epochs of Animals dataset

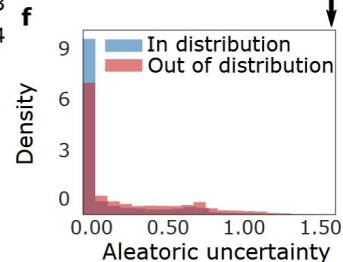
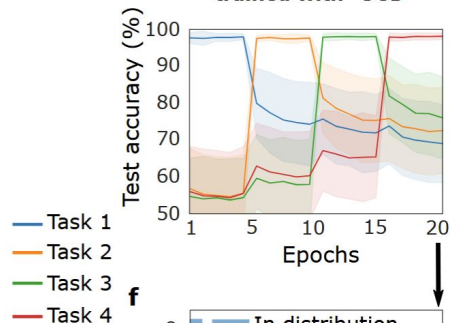


◆ Mitigating catastrophic forgetting

20 epochs of Animals dataset

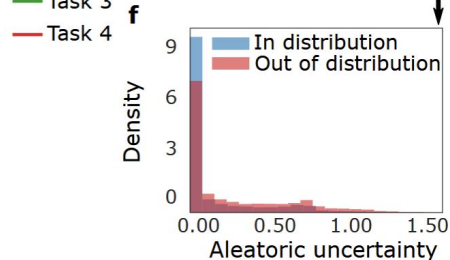
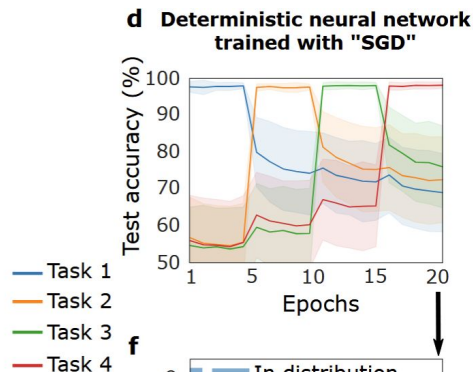
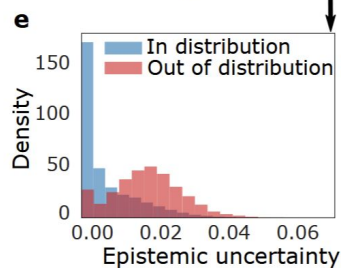
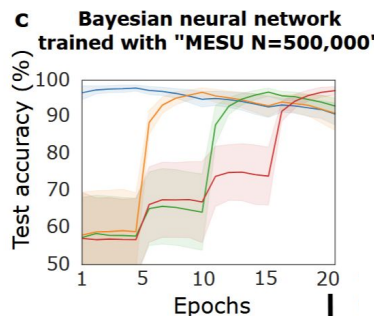
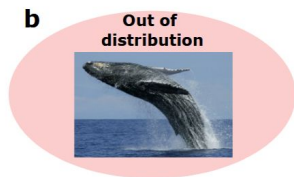
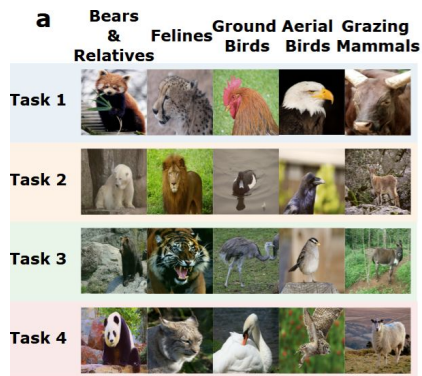


d Deterministic neural network trained with "SGD"



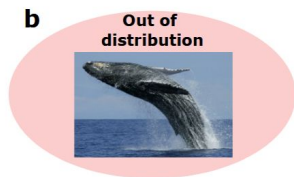
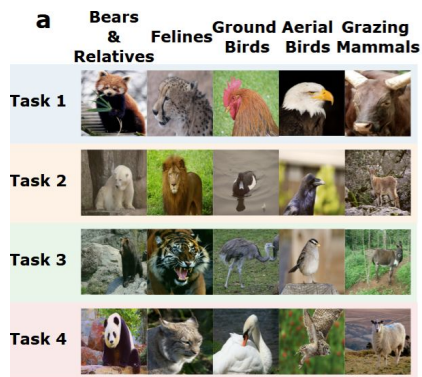
◆ Mitigating catastrophic forgetting

20 epochs of Animals dataset

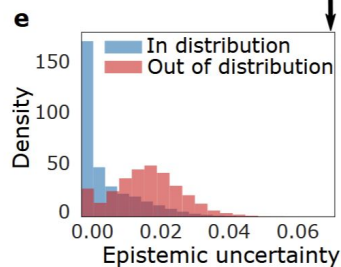
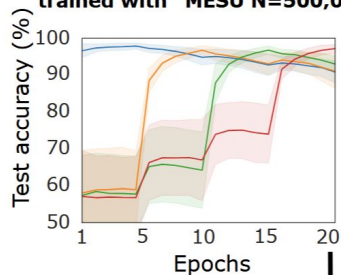


◆ Mitigating catastrophic forgetting

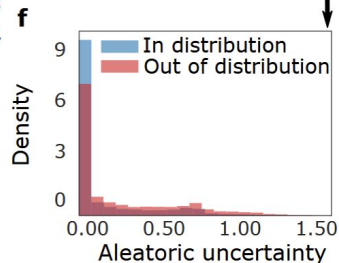
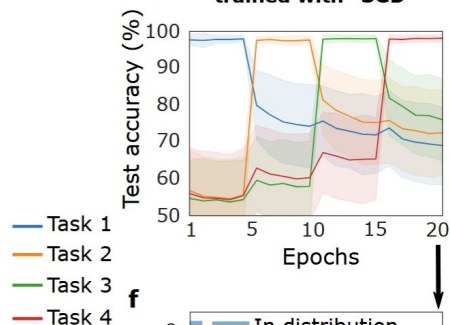
20 epochs of Animals dataset



c Bayesian neural network trained with "MESU N=500,000"



d Deterministic neural network trained with "SGD"



MESU mitigates catastrophic forgetting and allows to maintain multiple domains at the same time

◆ State-of-the-art

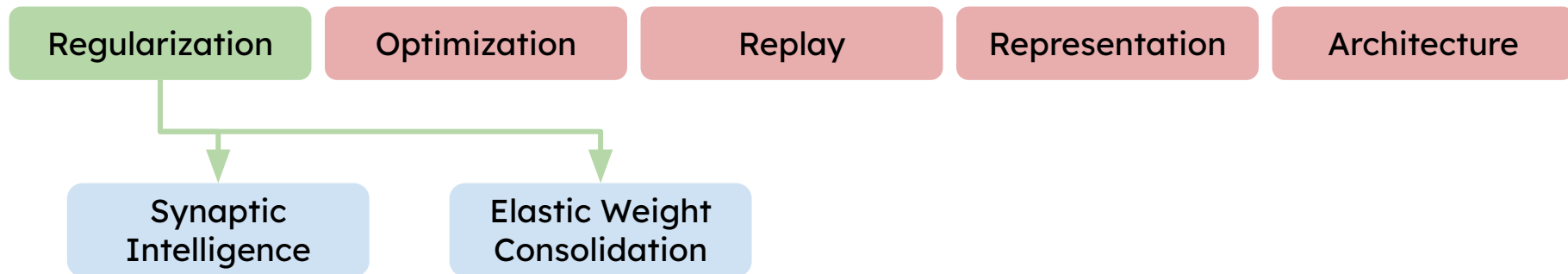
About task boundaries



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

◆ State-of-the-art

About task boundaries



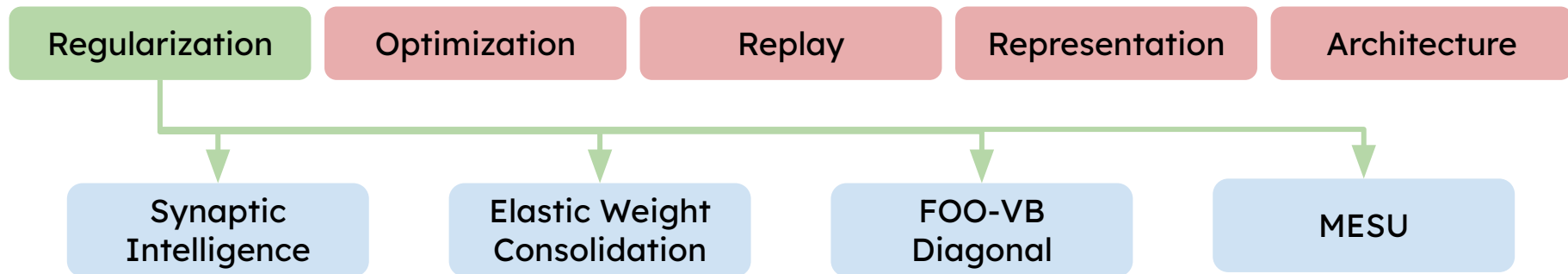
$$\tilde{L}_\mu = L_\mu + c \sum_k \Omega_k^\mu (\tilde{\theta}_k - \theta_k)^2 \quad \mathcal{L}(\theta) = \mathcal{L}_B(\theta) + \sum_i \frac{\lambda}{2} F_i(\theta_i - \theta_{A,i}^*)^2$$

Both SI and EWC methods require tasks boundaries to work, whereas FOO-VB Diagonal and MESU do not

Task boundaries: the algorithm knows ahead of time when the task changes, and **requires the previous parameters** in memory

◆ State-of-the-art

About task boundaries



$$\tilde{L}_\mu = L_\mu + c \sum_k \Omega_k^\mu (\tilde{\theta}_k - \theta_k)^2 \quad \mathcal{L}(\theta) = \mathcal{L}_B(\theta) + \sum_i \frac{\lambda}{2} F_i(\theta_i - \theta_{A,i}^*)^2$$

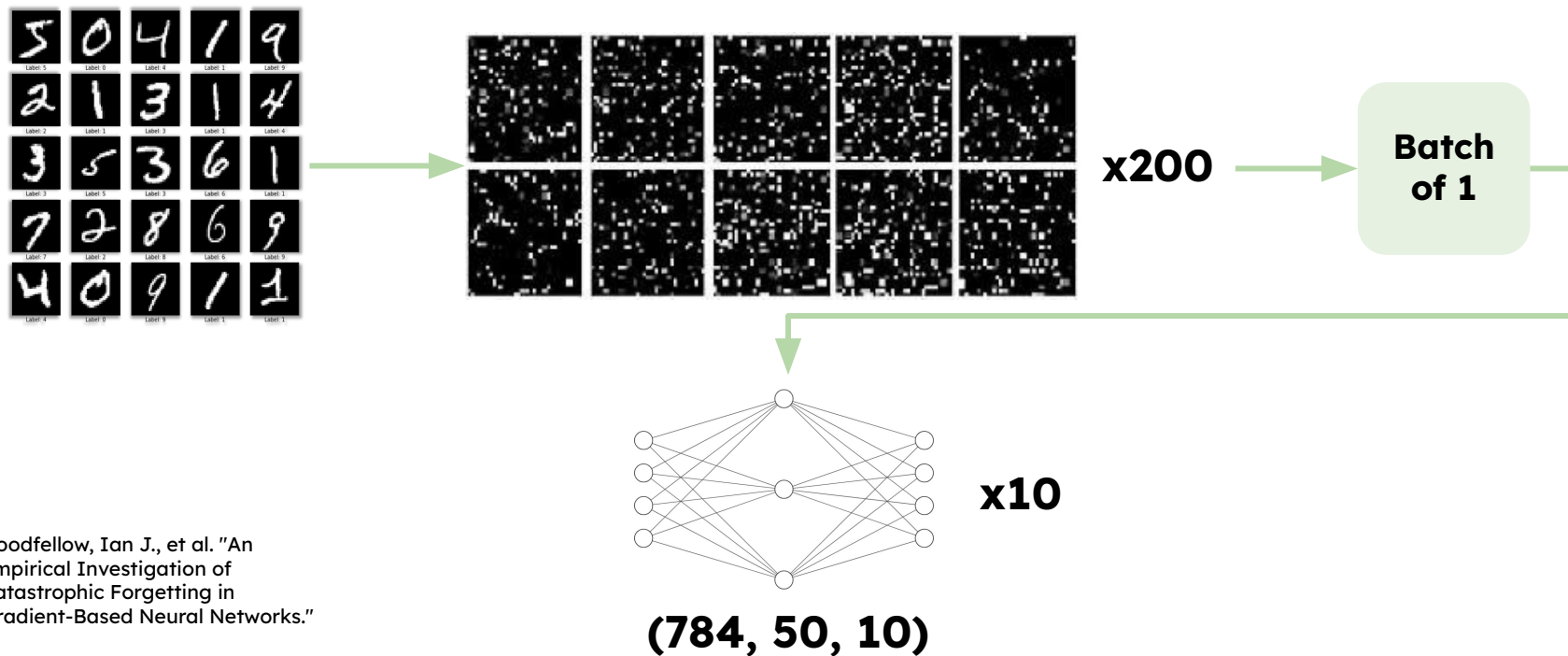
Both SI and EWC methods require task boundaries to work, whereas FOO-VB Diagonal and MESU do not

Task boundaries: the algorithm knows ahead of time when the task changes, and **requires the previous parameters** in memory

◆ Permuted MNIST Benchmark

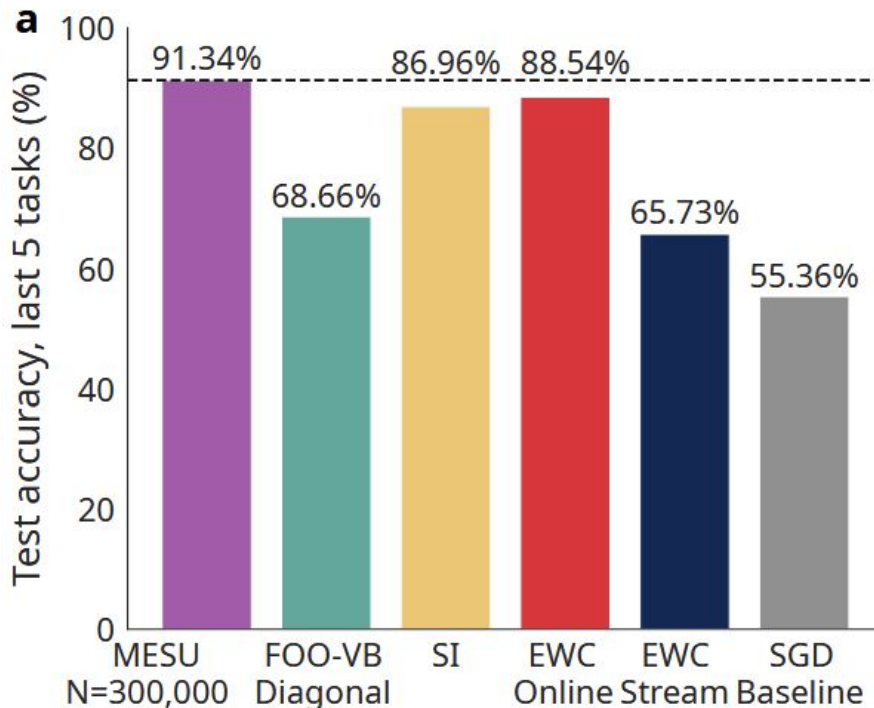
Permutation of pixels from MNIST

Simulation setup with a low amount of synapses in Online Continual Learning strategy



◆ Average accuracies over five tasks

Increased performance due to the memory window

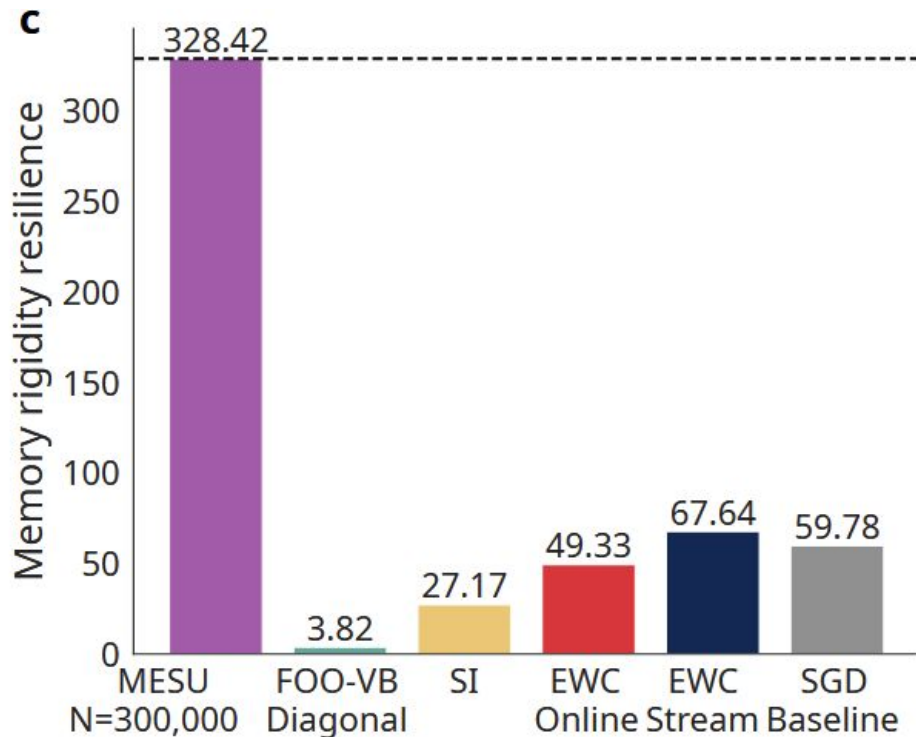


Stream: previous input
Online: previous task

Test accuracy on the
memory window
N = 300,000 = 5 tasks

◆ Resilience to loss of plasticity

Increased performance due to the memory window



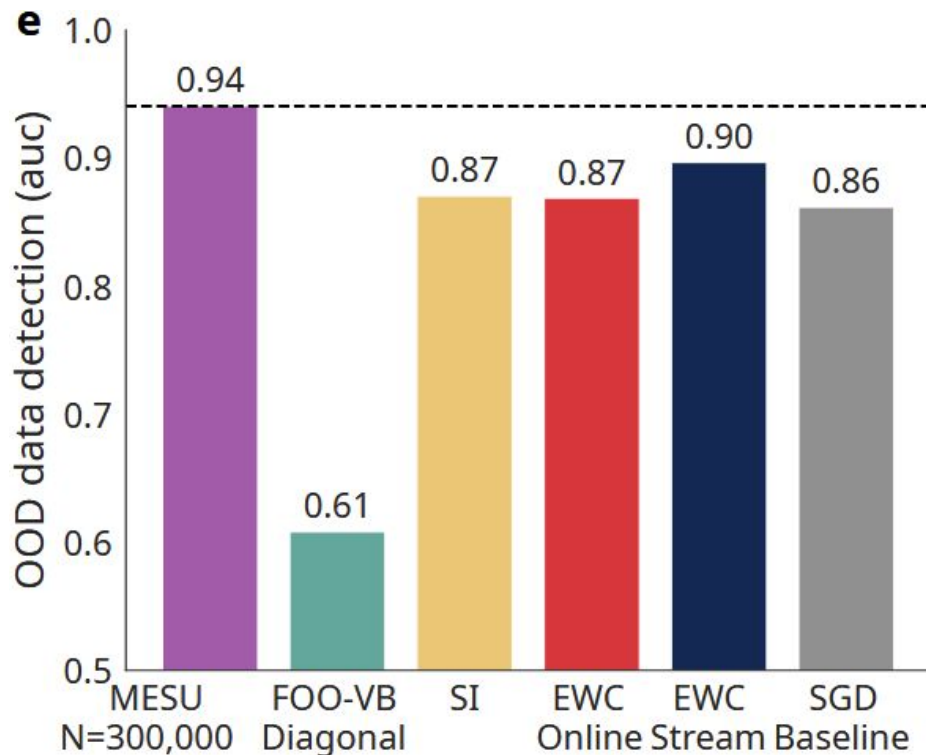
Stream: previous input
Online: previous task

Memory rigidity
 resilience between the
 last and the first task

$$\mathcal{R}_t = \frac{1}{|\mathcal{A}_0 - \mathcal{A}_t|}$$

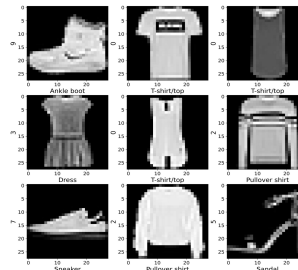
◆ Out-of-distribution data detection

Increased performance due to the memory window



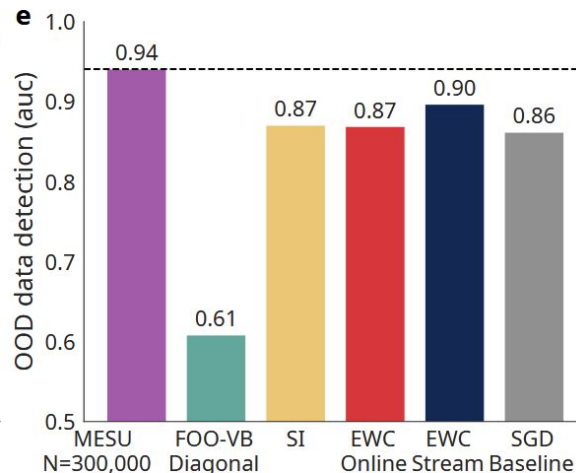
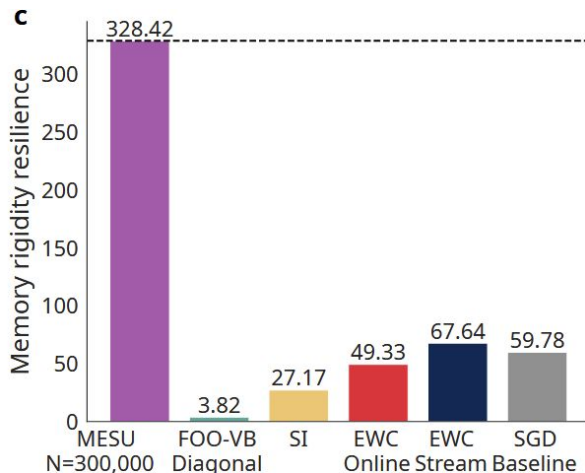
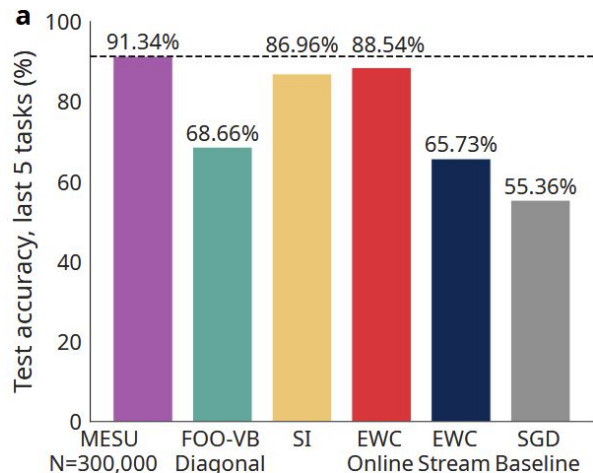
Stream: previous input
Online: previous task

OOD detection between
ID distribution
Permuted MNIST and
OOD distribution
Fashion-MNIST



◆ Summary

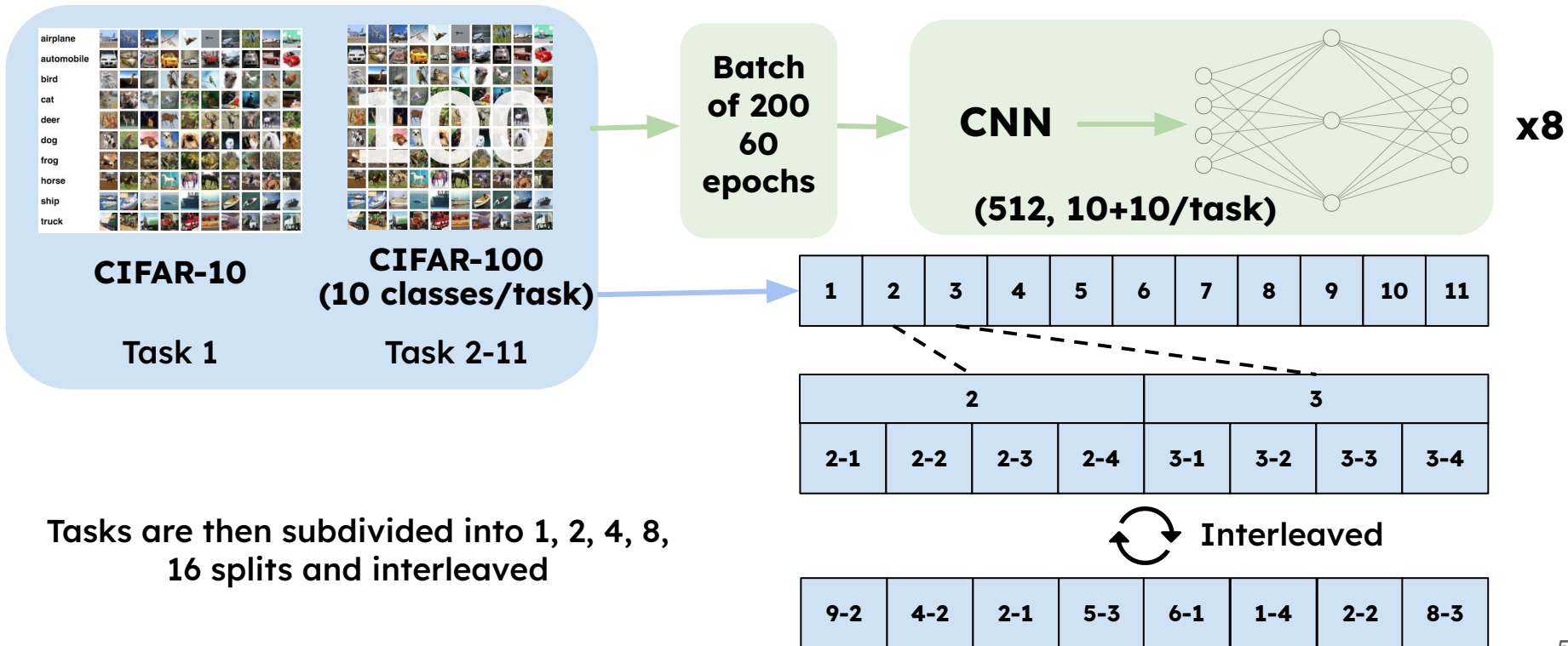
Increased performance due to the memory window



MESU is resilient against loss of plasticity, even after more than 200 tasks and outperforms methods with task-boundaries

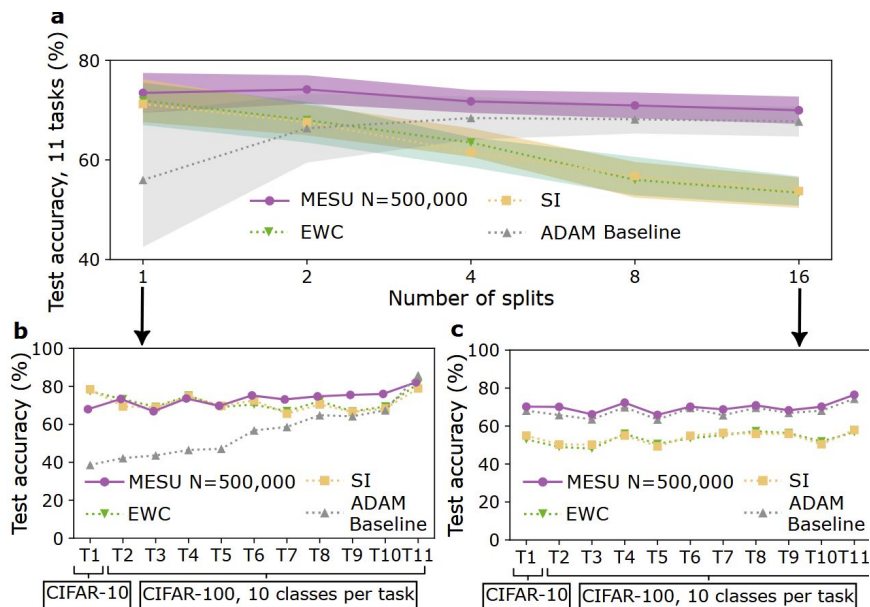
◆ CIFAR-10 / CIFAR-100 Benchmark

Domain-incremental learning



◆ Beyond Multi-layer perceptrons

CIFAR-10 / CIFAR-100 splits



MESU outperforms methods with task-boundaries on end-to-end training with domain-incremental tasks

Bayesian continual learning and forgetting in neural networks

Djohan Bonnet, Kellian Cottart, Tifenn Hirtzlin, Tarcisius Januel, Thomas Dalgaty, Elisa Vianello, Damien Querlioz

Nature Communications, Accepted

Thank you for listening!

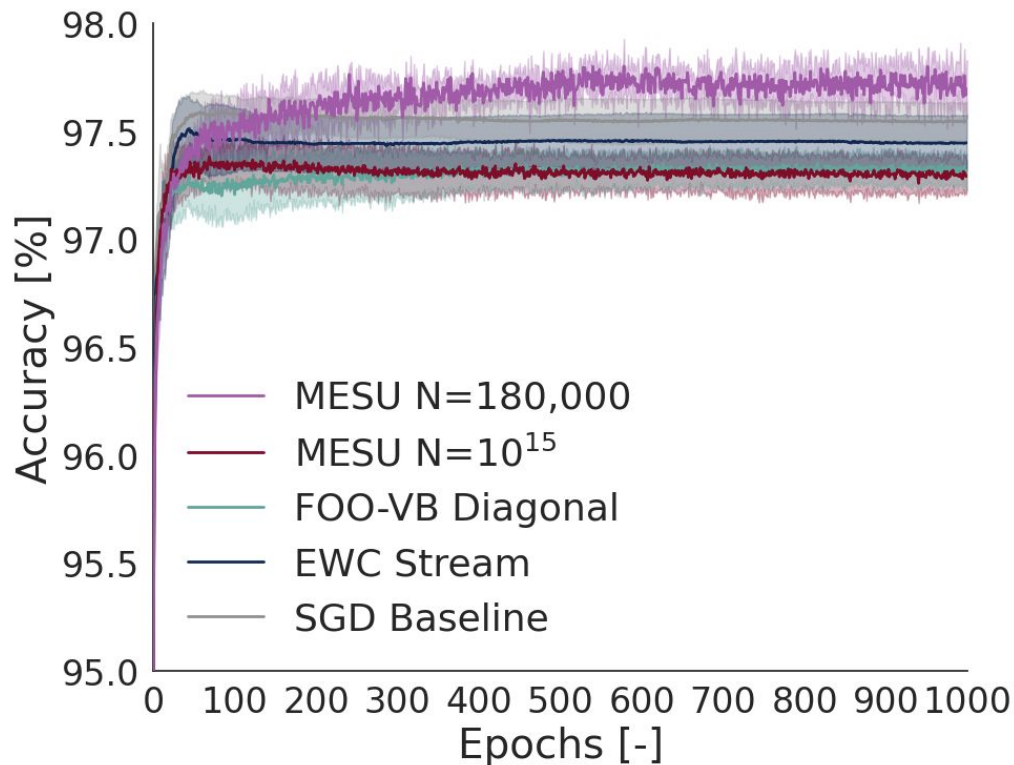
Q&A (5 minutes)



Appendix

◆ Rigidity and uncertainties - Accuracy

1000 epochs of MNIST (784 - 50 - 10) with a permutation of MNIST as OOD



◆ CNN architecture

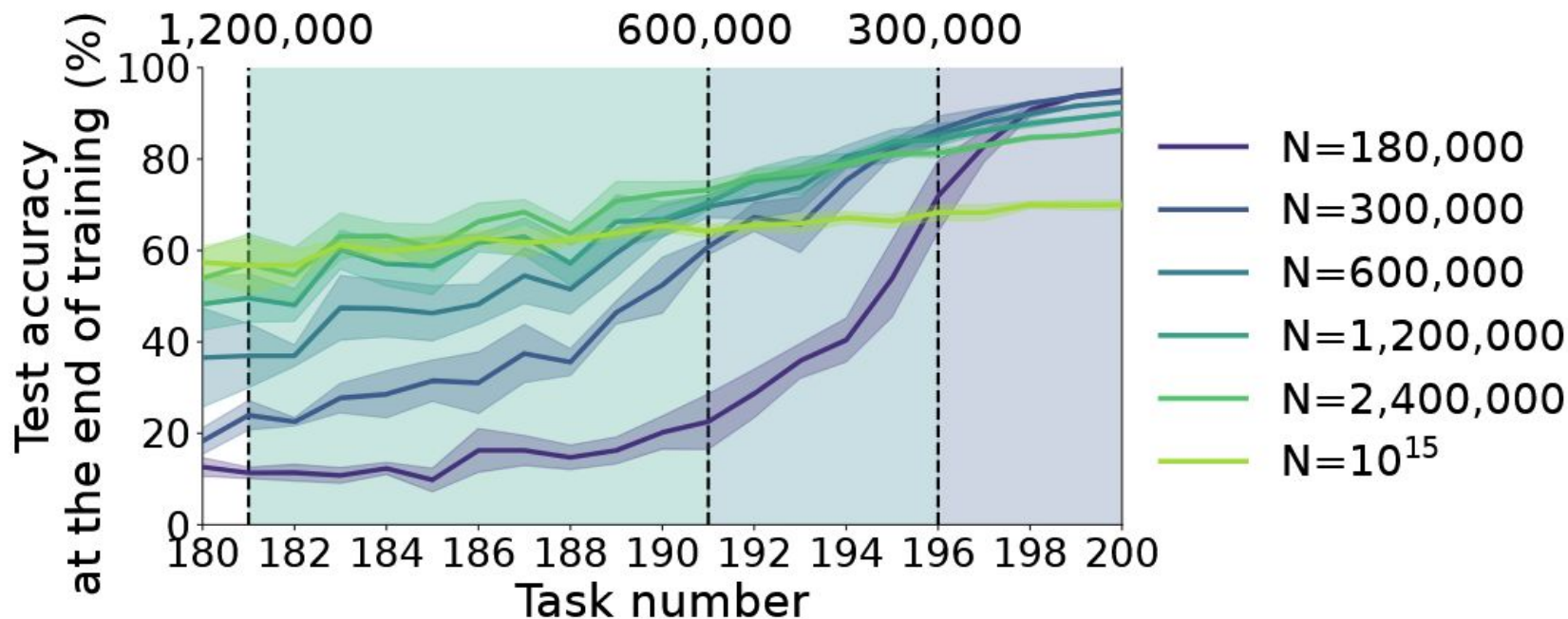
CIFAR10-100 CNN architecture

Operation	Kernel	Stride	Filters	Dropout	Nonlin.
Input	-	-	-	-	-
Convolution	3×3	1×1	32	-	ReLU
Convolution	3×3	1×1	32	-	ReLU
MaxPool	2×2	-	-	0.25	-
Convolution	3×3	1×1	64	-	ReLU
Convolution	3×3	1×1	64	-	ReLU
MaxPool	2×2	-	-	0.25	-
Fully connected	-	-	512	0.5	ReLU
Task 1: Fully connected	-	-	10	-	-
⋮	-	-	⋮	-	-
Task m : Fully connected	-	-	10	-	-

Table 2. CIFAR-10/100 model architecture, and dropout parameters. m : number of tasks.

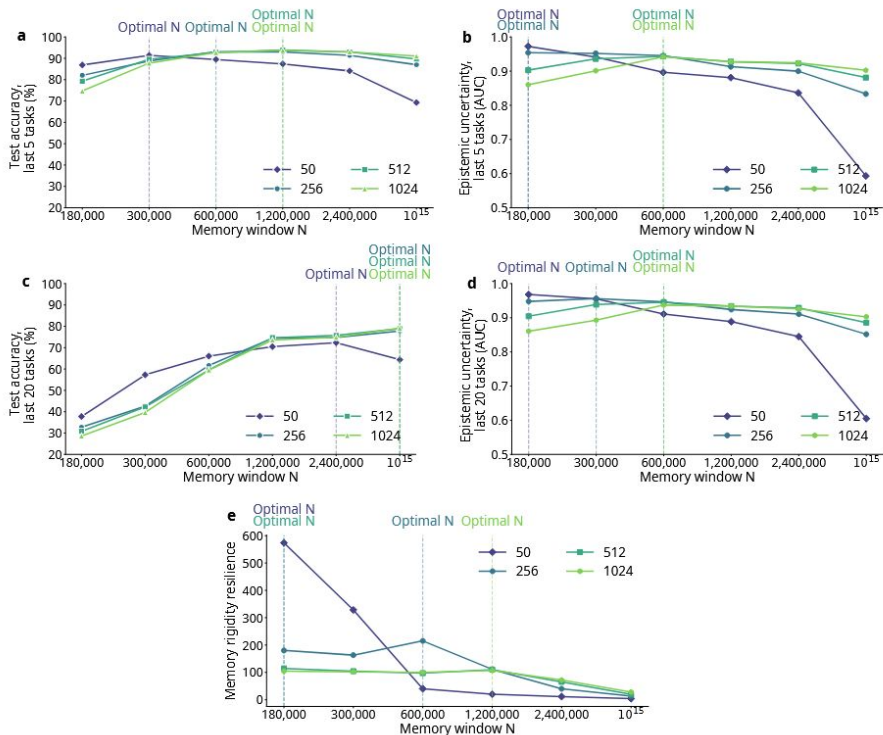
◆ N ablation study

Effect of the memory window on performance



◆ N ablation study

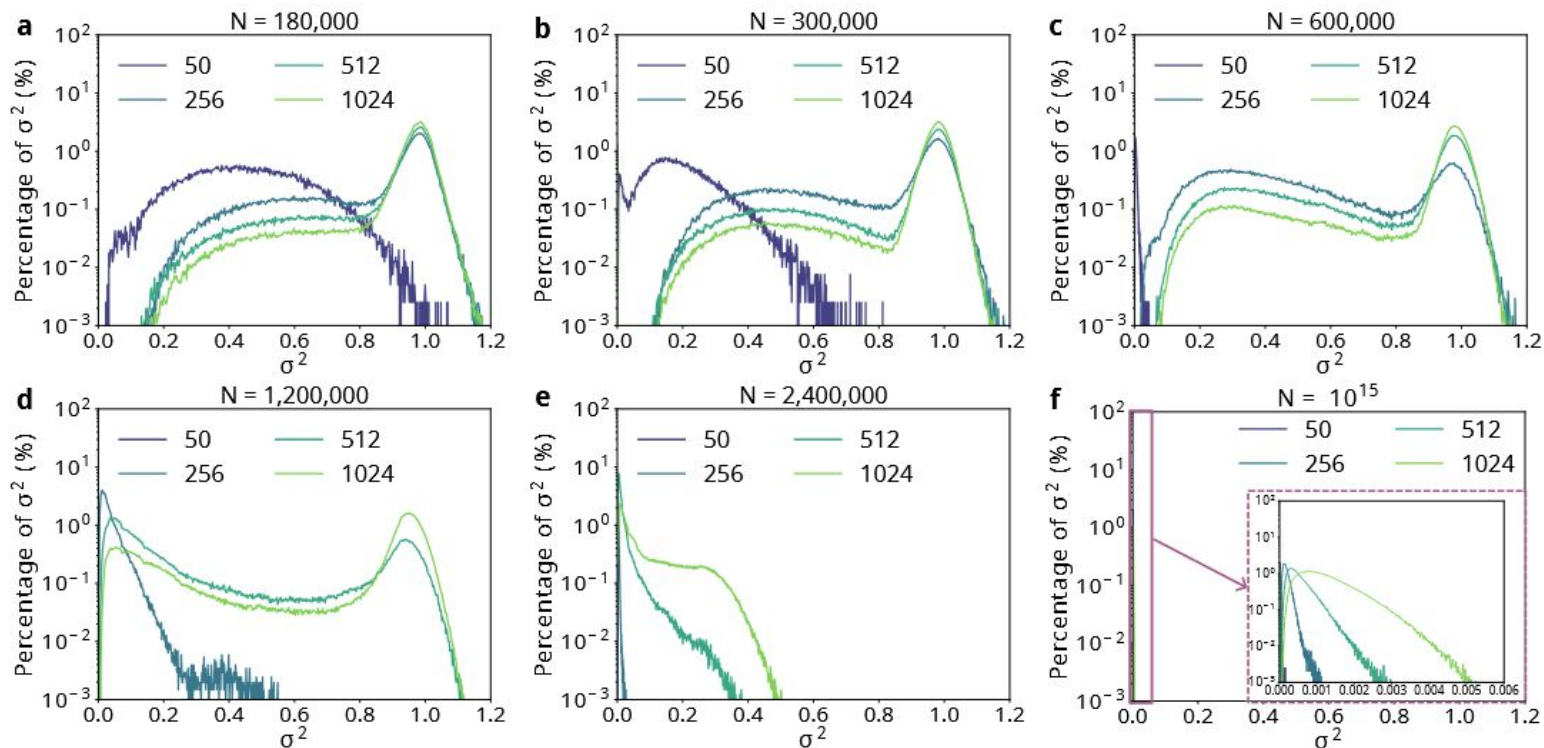
What is the optimal N value?



Optimal N value depends on what we are trying to maximize: amount of tasks, accuracy, OOD detection

◆ N ablation study

Effect of the memory window on variance



◆ Permuted MNIST detail

Effect of the memory window on variance

