# DATA POISONING ATTACK AIMING THE VULNERABILITY OF CONTINUAL LEARNING

*Gyojin Han[1*], Jaehyun Choi[1*], Hyeong Gwon Hong[2], and Junmo Kim[1]*

[1]School of Electrical Engineering, KAIST, South Korea
[2]Kim Jaechul Graduate School of AI, KAIST, South Korea
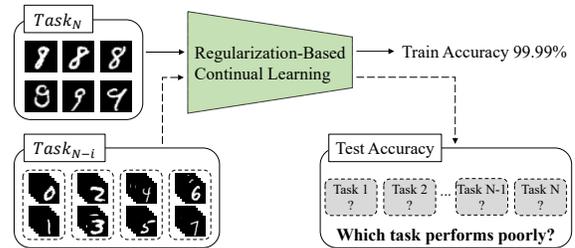
## ABSTRACT

Generally, regularization-based continual learning models limit access to the previous task data to imitate the real-world constraints related to memory and privacy. However, this introduces a problem in these models by not being able to track the performance on each task. In essence, current continual learning methods are susceptible to attacks on previous tasks. We demonstrate the vulnerability of regularization-based continual learning methods by presenting a simple task-specific data poisoning attack that can be used in the learning process of a new task. Training data generated by the proposed attack causes performance degradation on a specific task targeted by the attacker. We experiment with the attack on the two representative regularization-based continual learning methods, Elastic Weight Consolidation (EWC) and Synaptic Intelligence (SI), trained with variants of MNIST dataset. The experiment results justify the vulnerability proposed in this paper and demonstrate the importance of developing continual learning models that are robust to adversarial attacks.

***Index Terms***— Data poisoning, continual learning, catastrophic forgetting

## 1. INTRODUCTION

Humans can continuously learn new concepts throughout their lifetime while retaining previously learned knowledge. In contrast, neural networks trained on a new task that lies in a different distribution from previous tasks, suffer from performance degradation on previous tasks by losing information. This phenomenon, known as *catastrophic forgetting* [1], happens due to one of the limitations of the neural networks: train and test data must be in the same distribution for the network to perform well. To ease *catastrophic forgetting*, continual learning (also termed lifelong or incremental learning) [2] seeks to obtain a single model that works well on all of the learned tasks while incrementally training the model with access only to the current training task data.

Recent continual learning methods have overcome *catastrophic forgetting* showing remarkable performance in both past and current tasks. However, in a real-world situation, it is impossible to verify whether the model still works well on

_____
*Equal contribution.



**Fig. 1**. The regularization-based continual learning methods do not hold the data from the previous data. The model will be validated only by the train accuracy before publishing the model to the users. However, the train accuracy does not show the model's performance on each past task, thereby being unable to detect whether attacks are done to the model or not.
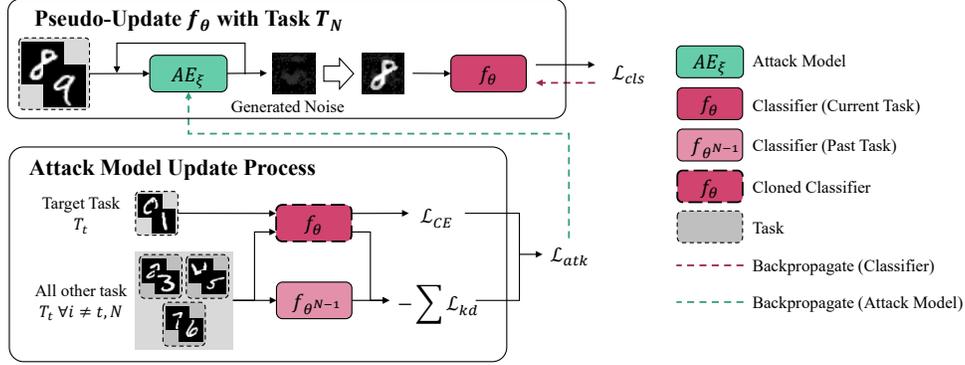
the past tasks since the previous data is not available due to memory and privacy problems. In other words, train accuracy is the only option to validate the performance of the continual learning models, thereby blindly utilizing the model even when the performance on one of the tasks is low as demonstrated in the Fig. 1. The difficulty of tracking the reliability of the model on past tasks is a serious problem especially when the attack is done to a specific task that the model has learned in the past as it would not affect the train accuracy. Despite such a problem, adversarial attacks [3] and defenses are not actively discussed in the field of continual learning.

In this paper, we bring forward the aforementioned problem in continual learning and experimentally justify the vulnerability with a simple task-specific data poisoning attack. The attack is designed to not affect the training accuracy of the model to be indistinguishable when training. More specifically, we generated adversarial data [4] during training that behaves in a new way, to work in continual learning. The generated adversarial data of a new task only severely degrades the performance of the targeted task.

## 2. RELATED WORKS

**Continual Learning.** Continual learning has three streams: rehearsal-based methods [5, 6], architecture-based methods

**Fig. 2**. The attack model loss $\mathcal{L}_{atk}$ is calculated from the pseudo-updated classifier $f_\theta$ and classifier from the past task $f_{\theta^{N-1}}$. The gradient is transmitted to the attack model $AE_\xi$ through partial differentiation. The attack model $AE_\xi$ then generates noise that gets added to the original data from the current task which updates the pseudo-updated classifier $f_\theta$ with cross-entropy loss.

[7, 8, 9, 10, 11], and regularization-based methods [12, 13]. Among them, regularization-based methods add a regularization loss term to the loss function when learning a new task to reduce the amount of change in parameters that are important for classifying the previous tasks. In this paper, we consider the regularization-based method setting, having no access to data from past tasks during the training process for the new task. As there are no previous task data available, the regularization-based continual learning method cannot track the performance of the past task hence relying solely on the training accuracy of the current task when deploying the model. This opens the chance for data poisoning to easily attack the continual learning methods on previous tasks without getting detected during training time. We demonstrate how vulnerable the continual learning methods are with a simple task-specific data poisoning attack.

**Adversarial Attack.** First introduced in [3], adversarial examples refer to samples with a very small perturbation, usually imperceptible by human eyes but noticeable by machine learning models, creating a gap in the inference results between them. Although there are image-agnostic methods [14, 15], this paper deals primarily with image-dependent adversarial attacks which are methods for generating such adversarial examples. It can be categorized into test time adversarial attacks [3, 16, 17] and data poisoning attacks [4]. Test time adversarial attacks generate images during inference and aim for incorrect inference results whereas data poisoning attacks generate images when training to make the model be trained erroneously. In this work, we propose a data poisoning attack for continual learning that adds perturbation to the training data of a new task so that the victim continual learning model loses information about the previous task specified by the attacker while learning the new task.

## 3. PROPOSED METHOD

In this section, we describe how an attacker generates adversarial data that cause a continual learning model to lose information on a previous task. More specifically, the model trainer is provided with adversarial training data that does not affect the train accuracy of a new task to prevent the trainer from detecting the attack. Moreover, as the model user might lose trust in the model if it does not work for all previous tasks, we propose an attack that only affects the performance of a specific task that the attacker intends to target. We assume a regularization-based method setting in which training data $\{D_1, \ldots, D_N\}$ corresponding to $N$ tasks $\{T_1, \ldots, T_N\}$ is provided sequentially. Training on the data of the new task proceeds without access to previous data. The attacker is provided with the training data $\{D_1, \ldots, D_N\}$ and a classifier $f_{\theta^{N-1}}$ trained up to $(N-1)$-th task. The goal of the attacker is to make the victim classifier lose knowledge about a target task $T_t$ while being trained well on the new task, where the $t$-th task is the target task. We emphasize that the attacker only slightly modifies the training data of the new task $T_N$ for flawless training on $T_N$ while losing information of $T_t$.

### 3.1. The attack process

We use an attack model $AE_\xi$ with an encoder-decoder structure to manipulate the training data of a new task into adversarial data. It takes the clean training data of the new task $D_N$ as input and generates noise that is bounded by $(-\epsilon, \epsilon)$. The generated noise is added to $D_N$, and becomes the adversarial training data $D_N{}'$ that can degrade the performance of the continual learning model. To train the attack model $AE_\xi$, we use the optimization method proposed by [4] with modifications in training process. The modified training process repeats the following two steps, (1) recording the trajectories

of a temporary model by updating it with adversarial data, and (2) training the attack model along the trajectories by pseudo-updating the recorded parameters. Fig. 2 illustrates the second step of training process of the attack model $AE_\xi$.

**Recording the trajectories of a temporary model.** For the optimization of $AE_\xi$, we need to approximate the trajectories of $f_{\theta^{N-1}}$ when it learns the adversarial image generated by $AE_\xi$. Therefore, we use a temporary model $f_\theta$ for episodic training. $f_\theta$ is trained with adversarial training data $D_N{}'$ generated by fixed $AE_\xi$. At this time, $f_\theta$ should be trained with a continual learning approach, as in the actual situation. Therefore, the loss for $f_\theta$, $\mathcal{L}_{cls}$ is:

$$\mathcal{L}_{cls} = \mathcal{L}_{CE}\left(f_\theta(x^N + AE_\xi(x^N)), y^N\right) + \Omega_m^{N-1} \quad (1)$$

where $\Omega_m^{N-1}$ is the regularization term of the continual learning method $m$, and $(x^N, y^N)$ is the mini-batch of the training data of the $N$-th task. Through the loss $\mathcal{L}_{cls}$, the parameters $\theta$ are updated and recorded as follows:

$$\theta \leftarrow \theta - \alpha_f \cdot \nabla_\theta \mathcal{L}_{cls} \quad (2)$$

where $\alpha_f$ is the learning rate for $f_\theta$.

**Training the attack model along the trajectories.** We pseudo-update recorded $f_\theta$ using the image generated by $AE_\xi$ with same loss $\mathcal{L}_{cls}$ as in step (1). Cross-entropy loss is calculated from the data of the target task $D_t$ using pseudo-updated $f_\theta$, and gradient values are transmitted to $AE_\xi$ through the loss. Then $AE_\xi$ can be updated with the gradient ascent. However, owing to the relevance between tasks, if an attack on the target task $T_t$ is attempted without any restrictions, the performance of the classifier against the other tasks involved will also be reduced. Therefore, appropriate constraints are required to maintain the classifier performance for other tasks. We want to preserve the outputs of inferences of $f_{\theta^{N-1}}$ for other tasks even if it is trained using adversarial training data. Therefore, we add the knowledge distillation loss term [18] to the loss function for training $AE_\xi$. The knowledge distillation loss between the outputs of $f_{\theta^{N-1}}$ and $f_\theta$, prevents adversarial data from affecting the inferences on other tasks. To calculate knowledge distillation loss, we sampled the data from the training data of each task. The knowledge distillation loss term for the $k$-th task $T_k$ is:

$$\mathcal{L}_{kd} = \alpha_{kd} \cdot T^2 \mathcal{L}_{KLD}\left(\sigma\left(\frac{f_\theta(x^k)}{T}\right), \sigma\left(\frac{f_{\theta^{N-1}}(x^k)}{T}\right)\right) \quad (3)$$

where $\alpha_{kd}$ is the balancing parameter of the knowledge distillation loss, $\mathcal{L}_{KLD}$ is the KL-divergence loss, $T$ is the temperature parameter, and $\sigma(\cdot)$ is the softmax function. Because $AE_\xi$ is trained via gradient ascent, the knowledge distillation loss term for all tasks except $T_t$ and $T_N$ is subtracted from the cross-entropy loss. The loss for $AE_\xi$, including the knowledge distillation loss term is:

$$\mathcal{L}_{atk} = \mathcal{L}_{CE}\left(f_\theta(x^t), y^t\right) - \sum_{i \neq t, N} \mathcal{L}_{kd}(f_\theta(x^i), f_{\theta^{N-1}}(x^i)) \quad (4)$$



(a) Permuted MNIST  (b) Split MNIST

**Fig. 3**. Example of clean samples (left two samples) and adversarial samples (right two samples) for (a) permuted MNIST, and (b) split MNIST. The perturbations on the samples appear differently due to the differences in the target task.

Finally, the parameters of the attack model $\xi$ are updated as follows.

$$\xi \leftarrow \xi + \alpha_{AE} \cdot \nabla_\xi \mathcal{L}_{atk} \quad (5)$$

## 4. EXPERIMENTS

To effectively validate the vulnerability of the proposed problem in continual learning, the experiment setting is chosen with the following considerations: 1) utilized continual learning methods must successfully alleviate the catastrophic forgetting, 2) the performance of all the past tasks must be high to show how easily the task-specific data poisoning drops the performance of a specific task. Accordingly, we experiment on two continual learning methods, EWC and SI, and two variants of MNIST [19] dataset, permuted MNIST [20] and split MNIST [13].
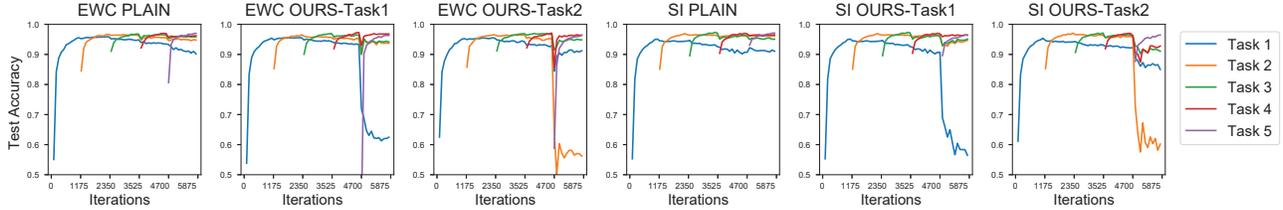
### 4.1. Experiment details.

The attack model consists of an encoder and decoder. The encoder has $3 \times 3$ convolution layers with 16, 64, and 128 channels, and the decoder has a $5 \times 5$ convolution layer with 128 channels and a $2 \times 2$ convolution layer with 64 channels. We train the attack model with Adam optimizer for 10 epochs with learning rate of 0.0001 and batch size of 256. The weight $\epsilon$ which determines the magnitude of the generated noise when adding to the clean sample is set to 0.2.

### 4.2. Results

Adversarial samples made by our task-specific data poisoning attack method can be seen in Fig. 3. Tab. 1 shows the final accuracy of the continual learning model after training for all tasks is completed. SGD in Tab. 1 is the baseline method. The 'Plain' results of EWC and SI show the effectiveness of each method by being higher than baseline results. Additionally, random uniform noise added to the clean sample is denoted by 'Noise'. 'Ours-T1' and 'Ours-T2' denote our task-specific data poisoning attacks done on task 1 and task 2, respectively.

As can be seen in Fig. 4, the noise created by the proposed attack caused the victim classifier to forget the knowledge about $T_1$ as it learns $T_5$. This proves the existence of adversarial data that causes much more severe catastrophic forgetting compared with clean data as the results of 'Noise', 'Ours-T1',

**Fig. 4**. Test accuracy on permuted MNIST for five tasks using EWC and SI. The graphs named 'PLAIN' show the effect of continual learning when no attacks are applied. The graphs named 'OURS-Task1' and 'OURS-Task2' show the test accuracy when $T_1$ and $T_2$ were attacked by our method, respectively. Best viewed zoomed in.

| Dataset | Method | | $T_1$ | $T_2$ | $T_3$ | $T_4$ | $T_5$ |
|---|---|---|---|---|---|---|---|
| Permuted MNIST | | SGD | 0.8059 | 0.8817 | 0.9306 | 0.9558 | 0.9520 |
| | EWC | Plain | 0.9029 | 0.9493 | 0.9570 | 0.9623 | 0.9700 |
| | | Noise | 0.8657 | 0.9356 | 0.9273 | 0.9687 | 0.9617 |
| | | Ours-$T_1$ | **0.6005** | 0.9344 | 0.9360 | 0.9666 | 0.9620 |
| | | Ours-$T_2$ | 0.9054 | **0.5625** | 0.9486 | 0.9639 | 0.9619 |
| | SI | Plain | 0.9102 | 0.9605 | 0.9529 | 0.9658 | 0.9717 |
| | | Noise | 0.9152 | 0.9561 | 0.9324 | 0.9589 | 0.9646 |
| | | Ours-$T_1$ | **0.5190** | 0.9364 | 0.9386 | 0.9592 | 0.9589 |
| | | Ours-$T_2$ | 0.8584 | **0.5618** | 0.9013 | 0.9183 | 0.962 |
| Split MNIST | | SGD | 0.4019 | 0.5901 | 0.1441 | 0.9084 | 0.9844 |
| | EWC | Plain | 0.4317 | 0.7424 | 0.1254 | 0.9305 | 0.9813 |
| | | Noise | 0.4132 | 0.6459 | 0.1660 | 0.8676 | 0.9803 |
| | | Ours-$T_1$ | **0.3825** | 0.5843 | 0.1596 | 0.9592 | 0.9773 |
| | | Ours-$T_2$ | 0.4463 | **0.5563** | 0.2006 | 0.9350 | 0.9692 |
| | SI | Plain | 0.4790 | 0.8242 | 0.3010 | 0.9728 | 0.9531 |
| | | Noise | 0.4643 | 0.7919 | 0.3116 | 0.9733 | 0.9531 |
| | | Ours-$T_1$ | **0.3939** | 0.8095 | 0.4242 | 0.8454 | 0.9576 |
| | | Ours-$T_2$ | 0.4577 | **0.7767** | 0.4248 | 0.8348 | 0.9551 |

**Table 1**. Final accuracy of the victim classifier $f_\theta$ to evaluate the performance of our attack.

| | Plain | Ours w/o $\mathcal{L}_{kd}$ | Ours with $\mathcal{L}_{kd}$ |
|---|---|---|---|
| EWC | -0.0129 | -0.0717 | -0.0195 |
| SI | -0.0087 | -0.0600 | -0.0237 |

**Table 2**. Backward transfer of $T_5$ for $T_2$, $T_3$, and $T_4$ on permuted MNIST.

distillation loss term. The backward transfer B was calculated as follows:

$$\text{B} = \frac{1}{N-2} \sum_{k \neq t, N} R_{N,k} - R_{k,k} \qquad (6)$$

where $R_{i,j}$ is the test accuracy of the classifier on $T_j$ just after trained with $T_i$.

As shown in Tab. 2, by placing a constraint on using the knowledge distillation loss when training the attack model, the increase in negative backward transfer owing to the attack is significantly reduced.

## 5. CONCLUSION

We reported weakness in continual learning caused by not having access to the data of previous tasks. This hinders performance tracking of previous tasks, which might reduce the reliability of the continual learning models and pose a serious problem for detecting adversarial attacks. In this regard, we propose a task-specific data poisoning attack scenario that this vulnerability could cause. The proposed attack degrades the performance of the continual learning model on the targeted task by adding perturbations to the training data of a new task. We highlight the importance of developing robust continual learning models by demonstrating the existence of adversarial data that causes the loss of knowledge of past tasks and suggest a simple attack scenario.

and 'Ours-T2' in Tab. 1 show. For split MNIST, the results of non-targeted tasks are not stable. This is due to the digits in MNIST dataset sharing many morphological features at the image patch level (e.g. 1&7, 3&8). Therefore, attacking a task inevitably affects the parameters of the other tasks decreasing or increasing the performance of non-targeted tasks in split MNIST. The targeted task result of 'Ours-T1' and 'Ours-T2' being lower than the 'Noise' shows that our attack method successfully attacks the targeted task. More importantly, the performance of non-targeted tasks stays in the reasonable range in line with the continual learning methods. This demonstrates that the attack on the target task cannot be noticed until inference on the target task occurs even for the deployed models. Furthermore, this shows that highly covert attacks on past tasks are possible because of the untraceable accuracy problem for past tasks in continual learning.

### 4.3. Ablation study

**Knowledge distillation loss.** We calculated the backward transfer [21] of the new task for the remaining tasks, except the target task to confirm the effectiveness of the knowledge

# 6. REFERENCES

[1] Robert M French, "Catastrophic forgetting in connectionist networks," *Trends in cognitive sciences*, vol. 3, no. 4, pp. 128–135, 1999.

[2] Zhiyuan Chen and Bing Liu, "Lifelong machine learning," *Synthesis Lectures on Artificial Intelligence and Machine Learning*, vol. 12, no. 3, pp. 1–207, 2018.

[3] Ian Goodfellow, Jonathon Shlens, and Christian Szegedy, "Explaining and harnessing adversarial examples," in *International Conference on Learning Representations*, 2015.

[4] Ji Feng, Qi-Zhi Cai, and Zhi-Hua Zhou, "Learning to confuse: Generating training time adversarial data with auto-encoder," in *Advances in Neural Information Processing Systems*, 2019, pp. 11971–11981.

[5] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert, "icarl: Incremental classifier and representation learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, July 2017.

[6] David Lopez-Paz and Marc' Aurelio Ranzato, "Gradient episodic memory for continual learning," in *Advances in Neural Information Processing Systems*, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds. 2017, vol. 30, Curran Associates, Inc.

[7] Andrei A Rusu, Neil C Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, and Raia Hadsell, "Progressive neural networks," *arXiv preprint arXiv:1606.04671*, 2016.

[8] Ju Xu and Zhanxing Zhu, "Reinforced continual learning," in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds. 2018, vol. 31, Curran Associates, Inc.

[9] Chrisantha Fernando, Dylan Banarse, Charles Blundell, Yori Zwols, David Ha, Andrei A Rusu, Alexander Pritzel, and Daan Wierstra, "Pathnet: Evolution channels gradient descent in super neural networks," *arXiv preprint arXiv:1701.08734*, 2017.

[10] Arun Mallya and Svetlana Lazebnik, "Packnet: Adding multiple tasks to a single network by iterative pruning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, June 2018.

[11] Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou, "Overcoming catastrophic forgetting with hard attention to the task," in *Proceedings of International Conference on Machine Learning*, Jennifer Dy and Andreas Krause, Eds. 10–15 Jul 2018, vol. 80 of *Proceedings of Machine Learning Research*, pp. 4548–4557, PMLR.

[12] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell, "Overcoming catastrophic forgetting in neural networks," *Proceedings of the National Academy of Sciences*, vol. 114, no. 13, pp. 3521–3526, 2017.

[13] Friedemann Zenke, Ben Poole, and Surya Ganguli, "Continual learning through synaptic intelligence," in *International Conference on Machine Learning*. PMLR, 2017, pp. 3987–3995.

[14] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard, "Universal adversarial perturbations," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1765–1773.

[15] Tejas Borkar, Felix Heide, and Lina Karam, "Defending against universal attacks through selective feature regeneration," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 709–719.

[16] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu, "Towards deep learning models resistant to adversarial attacks," *arXiv preprint arXiv:1706.06083*, 2017.

[17] Nicholas Carlini and David Wagner, "Towards evaluating the robustness of neural networks," in *2017 IEEE Symposium on Security and Privacy*, 2017, pp. 39–57.

[18] Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean, "Distilling the knowledge in a neural network," in *NIPS Deep Learning and Representation Learning Workshop*, 2015.

[19] Yann LeCun, Corinna Cortes, and Christopher J.C. Burges, "The mnist database of handwritten digits," 1998.

[20] Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou, "Overcoming catastrophic forgetting with hard attention to the task," in *International Conference on Machine Learning*. PMLR, 2018, pp. 4548–4557.

[21] David Lopez-Paz and Marc'Aurelio Ranzato, "Gradient episodic memory for continual learning," *Advances in Neural Information Processing Systems*, vol. 30, pp. 6467–6476, 2017.