

# On the Computation of the Fisher Information in Continual Learning

*Blog post, ICLR 2025*

Gido van de Ven

—

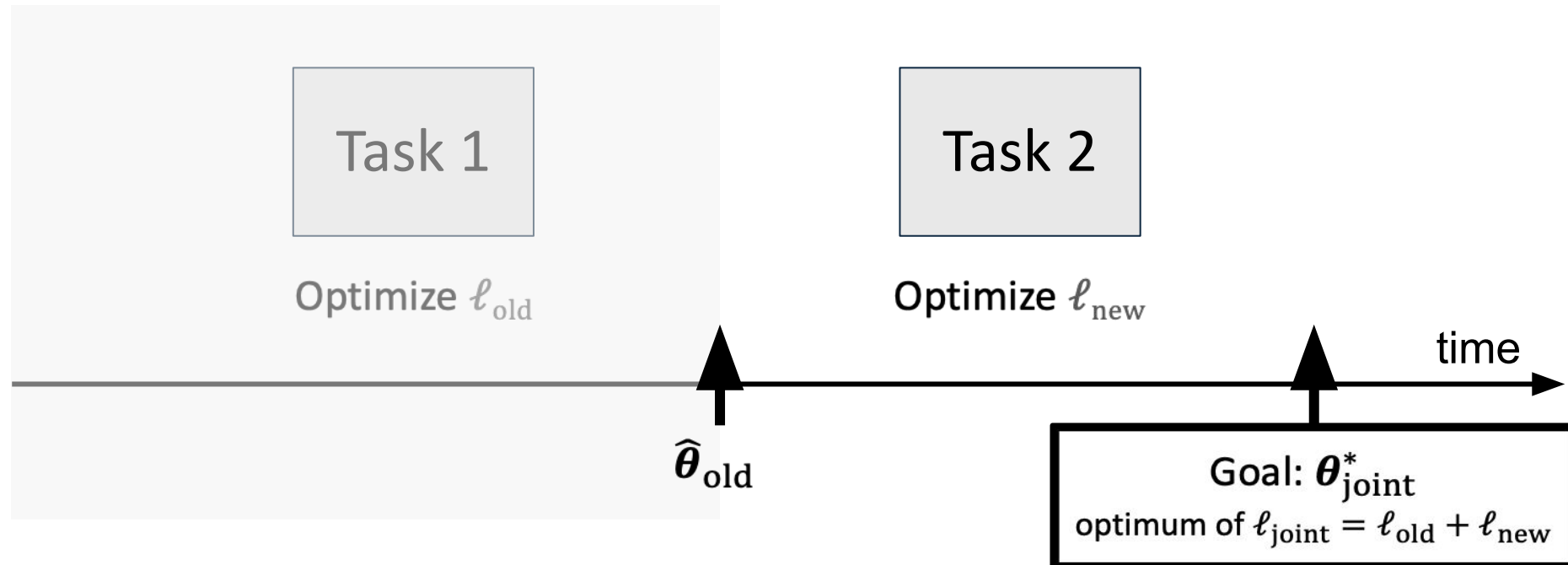
April 2025



Funding acknowledgement: This work has been supported by a senior post-doctoral fellowship from the Research Foundation – Flanders (FWO) under grant number 1266823N.

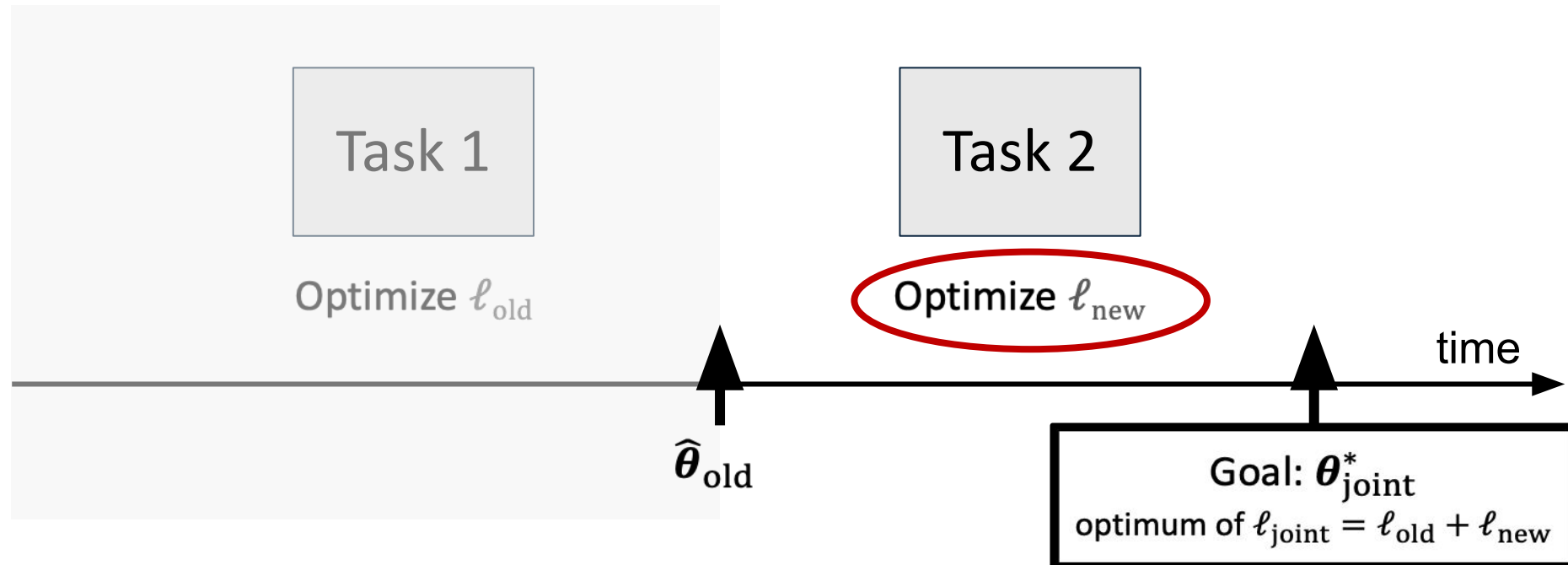
# The Continual Learning Problem

→ Optimize the parameters  $\theta$  of a neural network  $f_\theta$  for two tasks that are observed one after the other



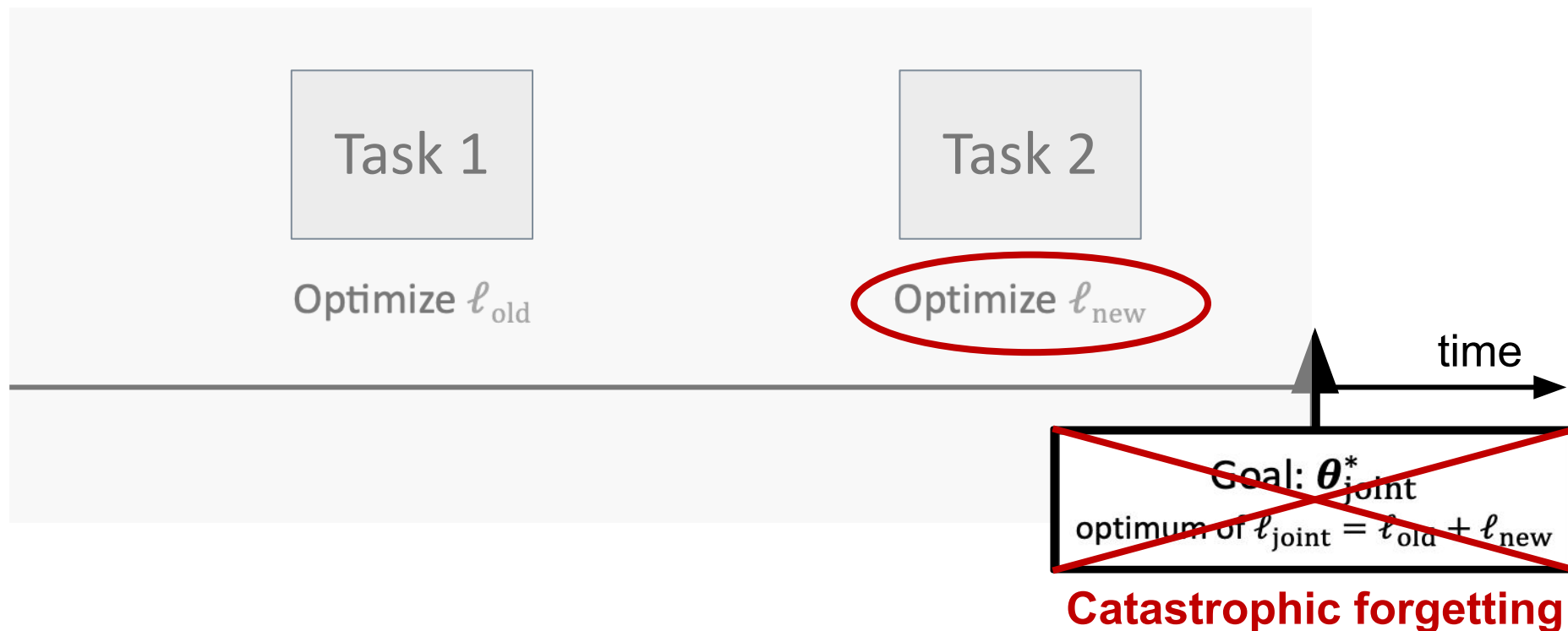
# The Continual Learning Problem

→ Optimize the parameters  $\theta$  of a neural network  $f_\theta$  for two tasks that are observed one after the other



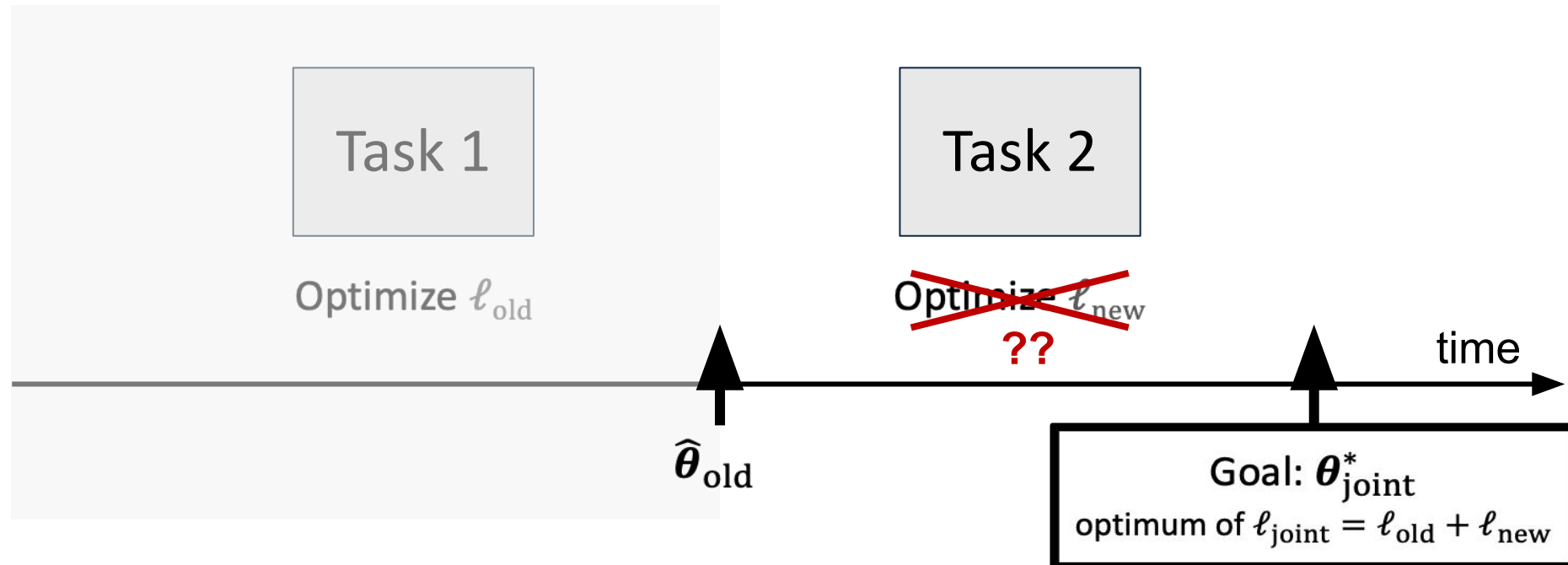
# The Continual Learning Problem

→ Optimize the parameters  $\theta$  of a neural network  $f_\theta$  for two tasks that are observed one after the other



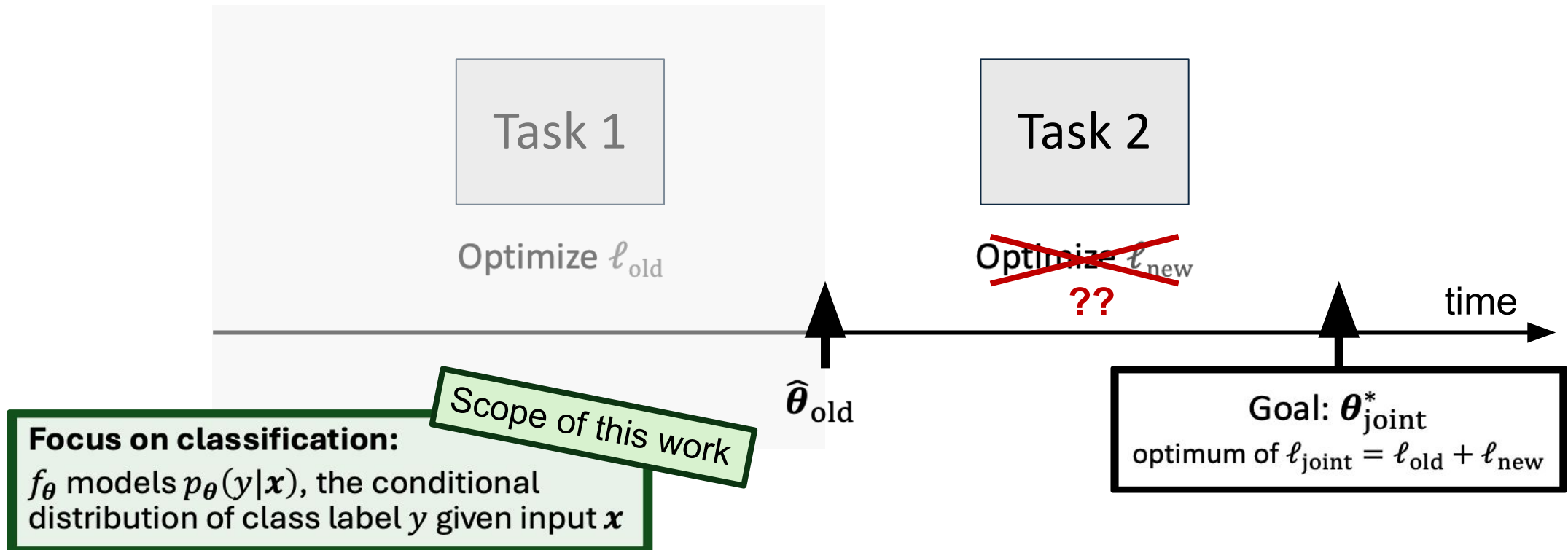
# The Continual Learning Problem

→ Optimize the parameters  $\theta$  of a neural network  $f_\theta$  for two tasks that are observed one after the other



# The Continual Learning Problem

→ Optimize the parameters  $\theta$  of a neural network  $f_\theta$  for two tasks that are observed one after the other



# Elastic Weight Consolidation (EWC)

- One of the most popular continual learning methods, >8000 citations (Google Scholar)
- Used as baseline in large proportion of continual learning studies

## Overcoming catastrophic forgetting in neural networks

**PNAS**

James Kirkpatrick<sup>a,1</sup>, Razvan Pascanu<sup>a</sup>, Neil Rabinowitz<sup>a</sup>, Joel Veness<sup>a</sup>, Guillaume Desjardins<sup>a</sup>, Andrei A. Rusu<sup>a</sup>, Kieran Milan<sup>a</sup>, John Quan<sup>a</sup>, Tiago Ramalho<sup>a</sup>, Agnieszka Grabska-Barwinska<sup>a</sup>, Demis Hassabis<sup>a</sup>, Claudia Clopath<sup>b</sup>, Dharshan Kumaran<sup>a</sup>, and Raia Hadsell<sup>a</sup>

*March, 2017*

# Elastic Weight Consolidation (EWC)

- One of the most popular continual learning methods, >8000 citations (Google Scholar)
- Used as baseline in large proportion of continual learning studies
- When training on a new task, EWC adds an extra term to the loss:

$$\ell_{\text{EWC}}(\boldsymbol{\theta}) = \ell_{\text{new}}(\boldsymbol{\theta}) + \frac{\lambda}{2} \sum_{i=1}^{N_{\text{params}}} F_{\text{old}}^{i,i} (\theta^i - \hat{\theta}_{\text{old}}^i)^2$$

$F_{\text{old}}^{i,i}$   $i^{\text{th}}$  diagonal element of the old network's  
Fisher Information matrix on the old task



# Elastic Weight Consolidation (EWC)

- One of the most popular continual learning methods, >8000 citations (Google Scholar)
- Used as baseline in large proportion of continual learning studies
- When training on a new task, EWC adds an extra term to the loss:

$$\ell_{\text{EWC}}(\boldsymbol{\theta}) = \ell_{\text{new}}(\boldsymbol{\theta}) + \frac{\lambda}{2} \sum_{i=1}^{N_{\text{params}}} F_{\text{old}}^{i,i} (\theta^i - \hat{\theta}_{\text{old}}^i)^2$$

$F_{\text{old}}^{i,i}$   $i^{\text{th}}$  diagonal element of the old network's  
Fisher Information matrix on the old task

My claim: EWC is usually implemented sub-optimally, and most of its currently reported results can likely be improved

# A Closer Look at the Fisher Information

Following [Martens, 2020 – JMLR](#), the  $i^{\text{th}}$  diagonal element of the network's Fisher Information matrix on the data of the old task, is defined as:

$$F_{\text{old}}^{i,i} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{old}}} \left[ \mathbb{E}_{y \sim p_{\hat{\theta}_{\text{old}}}} \left[ \left( \frac{\partial \log p_{\theta}(y|\mathbf{x})}{\partial \theta^i} \bigg|_{\theta = \hat{\theta}_{\text{old}}} \right)^2 \right] \right]$$

# A Closer Look at the Fisher Information

Following [Martens, 2020 – JMLR](#), the  $i^{\text{th}}$  diagonal element of the network's Fisher Information matrix on the data of the old task, is defined as:

$$F_{\text{old}}^{i,i} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{old}}} \left[ \mathbb{E}_{\mathbf{y} \sim p_{\hat{\theta}_{\text{old}}}(\mathbf{y}|\mathbf{x})} \left[ \left( \frac{\partial \log p_{\theta}(\mathbf{y}|\mathbf{x})}{\partial \theta^i} \bigg|_{\theta = \hat{\theta}_{\text{old}}} \right)^2 \right] \right]$$

• There are two expectations:

- (1) An **outer expectation** over  $\mathcal{D}_{\text{old}}$ , the input distribution of the old task
- (2) An **inner expectation** over  $p_{\hat{\theta}_{\text{old}}}(\mathbf{y}|\mathbf{x})$ , the conditional distribution of  $\mathbf{y}$  given  $\mathbf{x}$ , defined by the network after training on the old task

# A Closer Look at the Fisher Information

Following [Martens, 2020 – JMLR](#), the  $i^{\text{th}}$  diagonal element of the network's Fisher Information matrix on the data of the old task, is defined as:

$$F_{\text{old}}^{i,i} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{old}}} \left[ \mathbb{E}_{\mathbf{y} \sim p_{\hat{\theta}_{\text{old}}}(\mathbf{y}|\mathbf{x})} \left[ \left( \frac{\partial \log p_{\theta}(\mathbf{y}|\mathbf{x})}{\partial \theta^i} \bigg|_{\theta = \hat{\theta}_{\text{old}}} \right)^2 \right] \right]$$

• There are two expectations:

- (1) An **outer expectation** over  $\mathcal{D}_{\text{old}}$ , the input distribution of the old task
- (2) An **inner expectation** over  $p_{\hat{\theta}_{\text{old}}}(\mathbf{y}|\mathbf{x})$ , the conditional distribution of  $\mathbf{y}$  given  $\mathbf{x}$ , defined by the network after training on the old task

# A Closer Look at the Fisher Information

Following [Martens, 2020 – JMLR](#), the  $i^{\text{th}}$  diagonal element of the network's Fisher Information matrix on the data of the old task, is defined as:

$$F_{\text{old}}^{i,i} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{old}}} \left[ \mathbb{E}_{y \sim p_{\hat{\theta}_{\text{old}}}} \left[ \left( \frac{\partial \log p_{\theta}(y|\mathbf{x})}{\partial \theta^i} \bigg|_{\theta = \hat{\theta}_{\text{old}}} \right)^2 \right] \right]$$

There are two expectations:

- (1) An **outer expectation** over  $\mathcal{D}_{\text{old}}$ , the input distribution of the old task
- (2) An **inner expectation** over  $p_{\hat{\theta}_{\text{old}}}(y|\mathbf{x})$ , the conditional distribution of  $y$  given  $\mathbf{x}$ , defined by the network after training on the old task

# Different Ways to Compute the Fisher:

## (1) *Exact*

- **Outer expectation:** estimate using all training data  $D_{\text{old}}$
- **Inner expectation:** compute exactly

$$F_{\text{old}}^{i,i} = \frac{1}{|D_{\text{old}}|} \sum_{\mathbf{x} \in D_{\text{old}}} \left( \sum_{y=1}^{N_{\text{classes}}} p_{\hat{\boldsymbol{\theta}}_{\text{old}}}(y|\mathbf{x}) \left( \frac{\partial \log p_{\boldsymbol{\theta}}(y|\mathbf{x})}{\partial \theta^i} \bigg|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_{\text{old}}} \right)^2 \right)$$

Definition

I will refer to this option as **EXACT**

$$F_{\text{old}}^{i,i} = \mathbb{E}_{\mathbf{x} \sim D_{\text{old}}} \left[ \mathbb{E}_{y \sim p_{\hat{\boldsymbol{\theta}}_{\text{old}}}} \left[ \left( \frac{\partial \log p_{\boldsymbol{\theta}}(y|\mathbf{x})}{\partial \theta^i} \bigg|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_{\text{old}}} \right)^2 \right] \right]$$

# Different Ways to Compute the Fisher:

## *(2) Sampling data points*

- **Outer expectation:** estimate using  $n$  random samples from  $D_{\text{old}}$
- **Inner expectation:** compute exactly

$$F_{\text{old}}^{i,i} = \frac{1}{n} \sum_{\mathbf{x} \in S_{D_{\text{old}}}^{(n)}} \left( \sum_{y=1}^{N_{\text{classes}}} p_{\hat{\theta}_{\text{old}}}(y|\mathbf{x}) \left( \frac{\partial \log p_{\theta}(y|\mathbf{x})}{\partial \theta^i} \bigg|_{\theta=\hat{\theta}_{\text{old}}} \right)^2 \right)$$

Definition

I will explore this option using  $n = 500$ ,  
referring to it as **EXACT ( $n=500$ )**

$$F_{\text{old}}^{i,i} = \mathbb{E}_{\mathbf{x} \sim D_{\text{old}}} \left[ \mathbb{E}_{y \sim p_{\hat{\theta}_{\text{old}}}} \left[ \left( \frac{\partial \log p_{\theta}(y|\mathbf{x})}{\partial \theta^i} \bigg|_{\theta=\hat{\theta}_{\text{old}}} \right)^2 \right] \right]$$

# Different Ways to Compute the Fisher:

## *(3) Sampling labels*

- **Outer expectation:** estimate using over all training data  $D_{\text{old}}$
- **Inner expectation:** estimate using a single Monte Carlo sample

$$F_{\text{old}}^{i,i} = \frac{1}{|D_{\text{old}}|} \sum_{\mathbf{x} \in D_{\text{old}}} \left( \left. \frac{\partial \log p_{\boldsymbol{\theta}}(c_{\mathbf{x}}|\mathbf{x})}{\partial \theta^i} \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_{\text{old}}} \right)^2 \text{ with } c_{\mathbf{x}} \text{ randomly sampled from } p_{\hat{\boldsymbol{\theta}}_{\text{old}}}(\cdot|\mathbf{x})$$

Definition

I will refer to this option as **SAMPLE**

$$F_{\text{old}}^{i,i} = \mathbb{E}_{\mathbf{x} \sim D_{\text{old}}} \left[ \mathbb{E}_{y \sim p_{\hat{\boldsymbol{\theta}}_{\text{old}}}} \left[ \left( \left. \frac{\partial \log p_{\boldsymbol{\theta}}(y|\mathbf{x})}{\partial \theta^i} \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_{\text{old}}} \right)^2 \right] \right]$$



# Different Ways to Compute the Fisher:

## *(4) Empirical Fisher*

- **Outer expectation:** estimate using all training data  $D_{\text{old}}$
- **Inner expectation:** approximate by computing the squared gradient only for the ground-truth label

$$F_{\text{old}}^{i,i} = \frac{1}{|D_{\text{old}}|} \sum_{(\mathbf{x}, y) \in D_{\text{old}}} \left( \left. \frac{\partial \log p_{\boldsymbol{\theta}}(y|\mathbf{x})}{\partial \theta^i} \right|_{\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}_{\text{old}}} \right)^2$$

Definition

I will refer to this option as **EMPIRICAL**

$$F_{\text{old}}^{i,i} = \mathbb{E}_{\mathbf{x} \sim D_{\text{old}}} \left[ \mathbb{E}_{y \sim p_{\hat{\boldsymbol{\theta}}_{\text{old}}}} \left[ \left( \left. \frac{\partial \log p_{\boldsymbol{\theta}}(y|\mathbf{x})}{\partial \theta^i} \right|_{\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}_{\text{old}}} \right)^2 \right] \right]$$

# Different Ways to Compute the Fisher:

## *(5) Batched approximation of Empirical Fisher*

- **Outer expectation:** estimate by averaging over batched version of  $D_{\text{old}}$
- **Inner expectation:** approximate using the square of mini-batch averaged gradients w.r.t. ground-truth labels

Most used one!

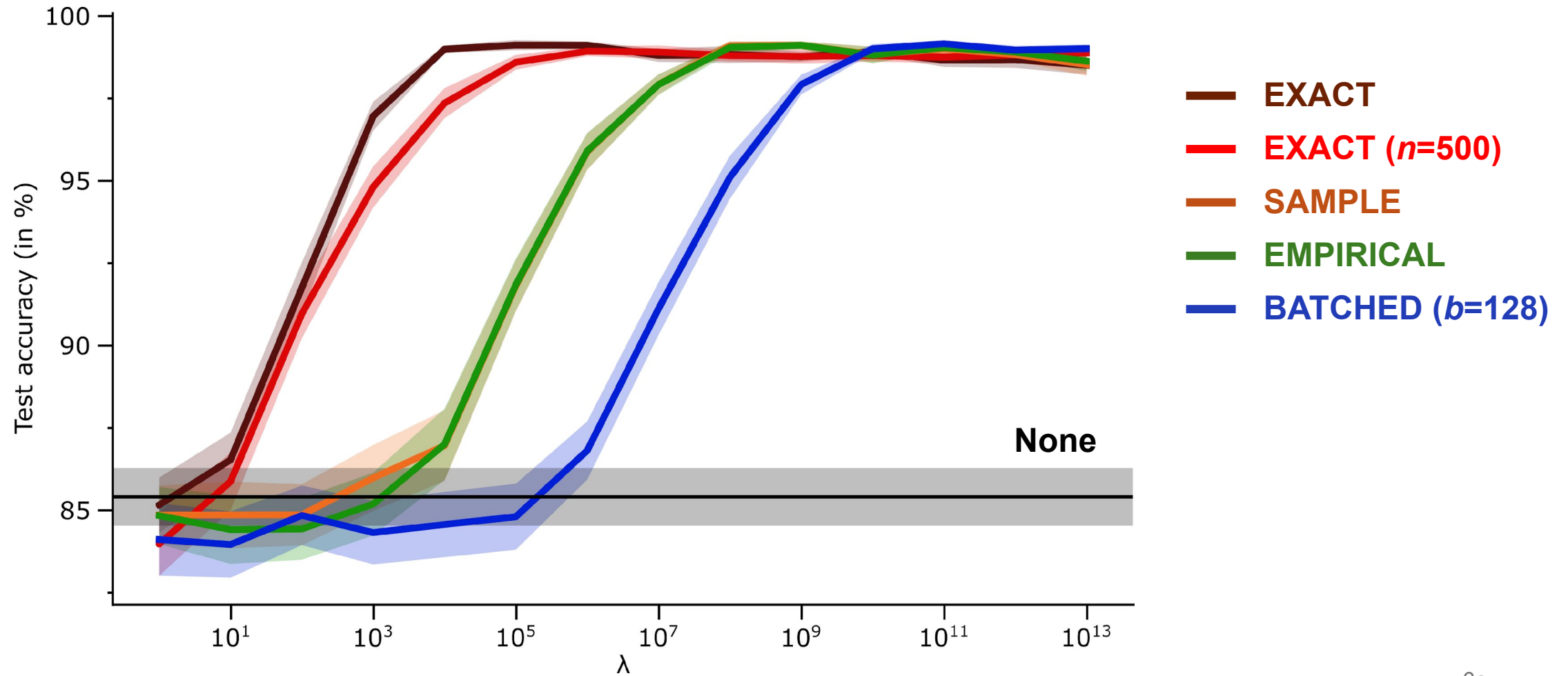
$$F_{\text{old}}^{i,i} = \frac{1}{|D_{\text{old}}^{(b)}|} \sum_{\mathcal{B} \in D_{\text{old}}^{(b)}} \left( \sum_{(\mathbf{x}, y) \in \mathcal{B}} \frac{\partial \log p_{\theta}(y|\mathbf{x})}{\partial \theta^i} \bigg|_{\theta = \hat{\theta}_{\text{old}}} \right)^2$$

I will explore this option using  $b = 128$ ,  
referring to it as **BATCHED ( $b=128$ )**

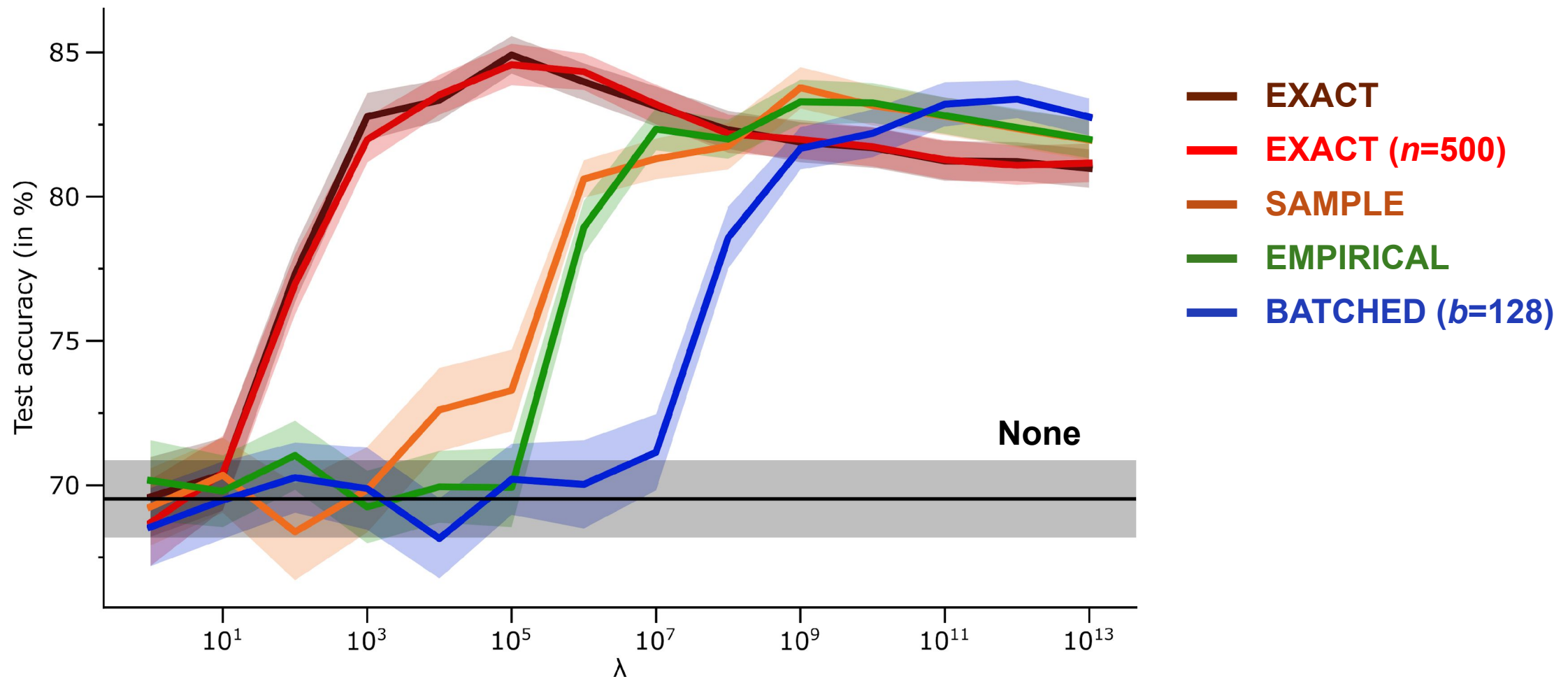
Definition

$$F_{\text{old}}^{i,i} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{old}}} \left[ \mathbb{E}_{y \sim p_{\hat{\theta}_{\text{old}}}} \left[ \left( \frac{\partial \log p_{\theta}(y|\mathbf{x})}{\partial \theta^i} \bigg|_{\theta = \hat{\theta}_{\text{old}}} \right)^2 \right] \right]$$

# Empirical Comparisons – Split MNIST



# Empirical Comparisons – Split CIFAR-10



# Conclusion

- The way in which the Fisher Information is computed can have substantial impact on the performance of EWC

# Recommendations

- (1) When using EWC, give details of how the Fisher is computed
- (2) Do not simply “*use the best performing hyperparameters from another paper*”, if you cannot guarantee that the Fisher is computed in the same way
- (3) It might be better to estimate the Fisher with fewer training samples, than to cut corners in some other way