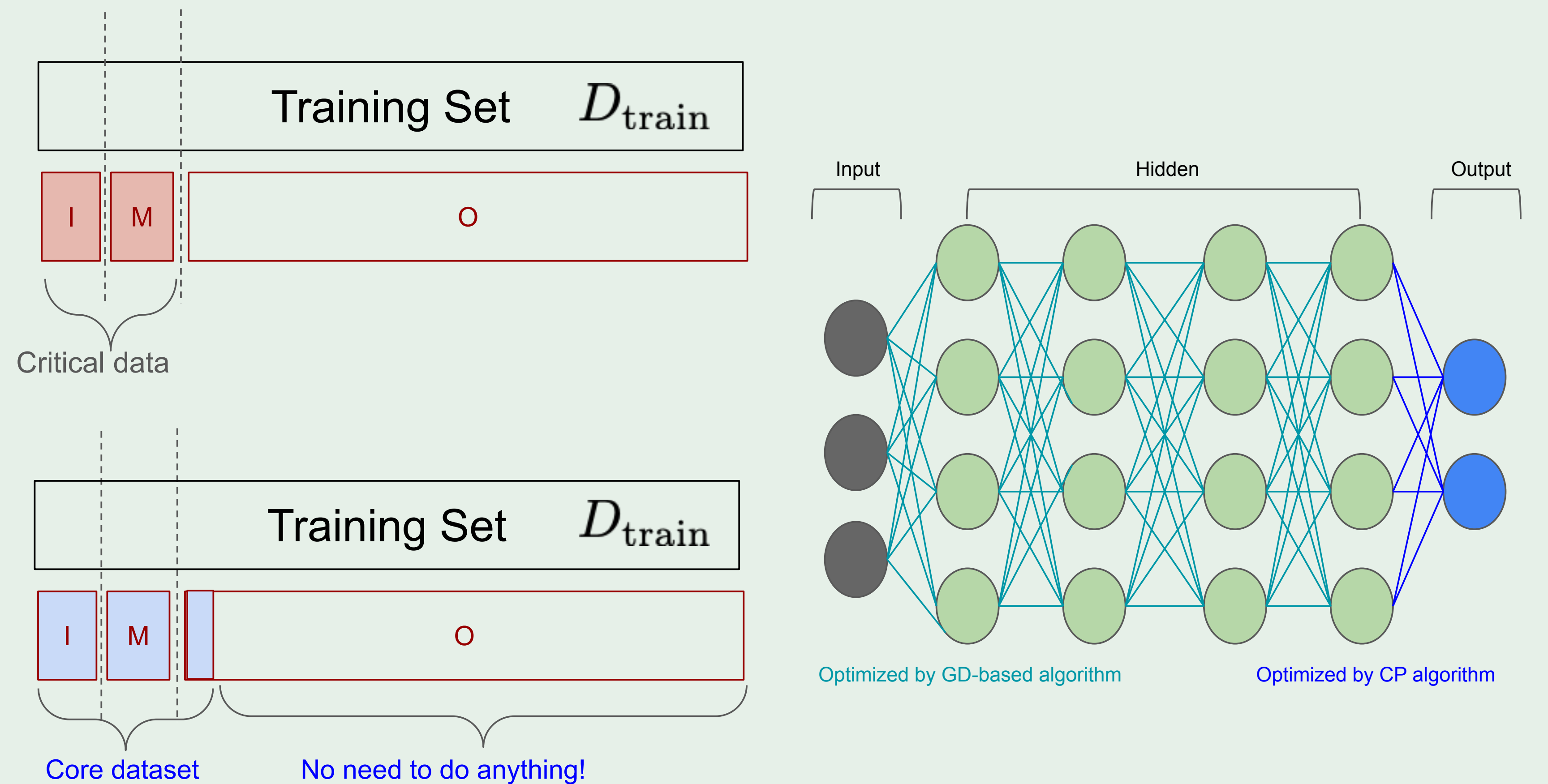


Abstract

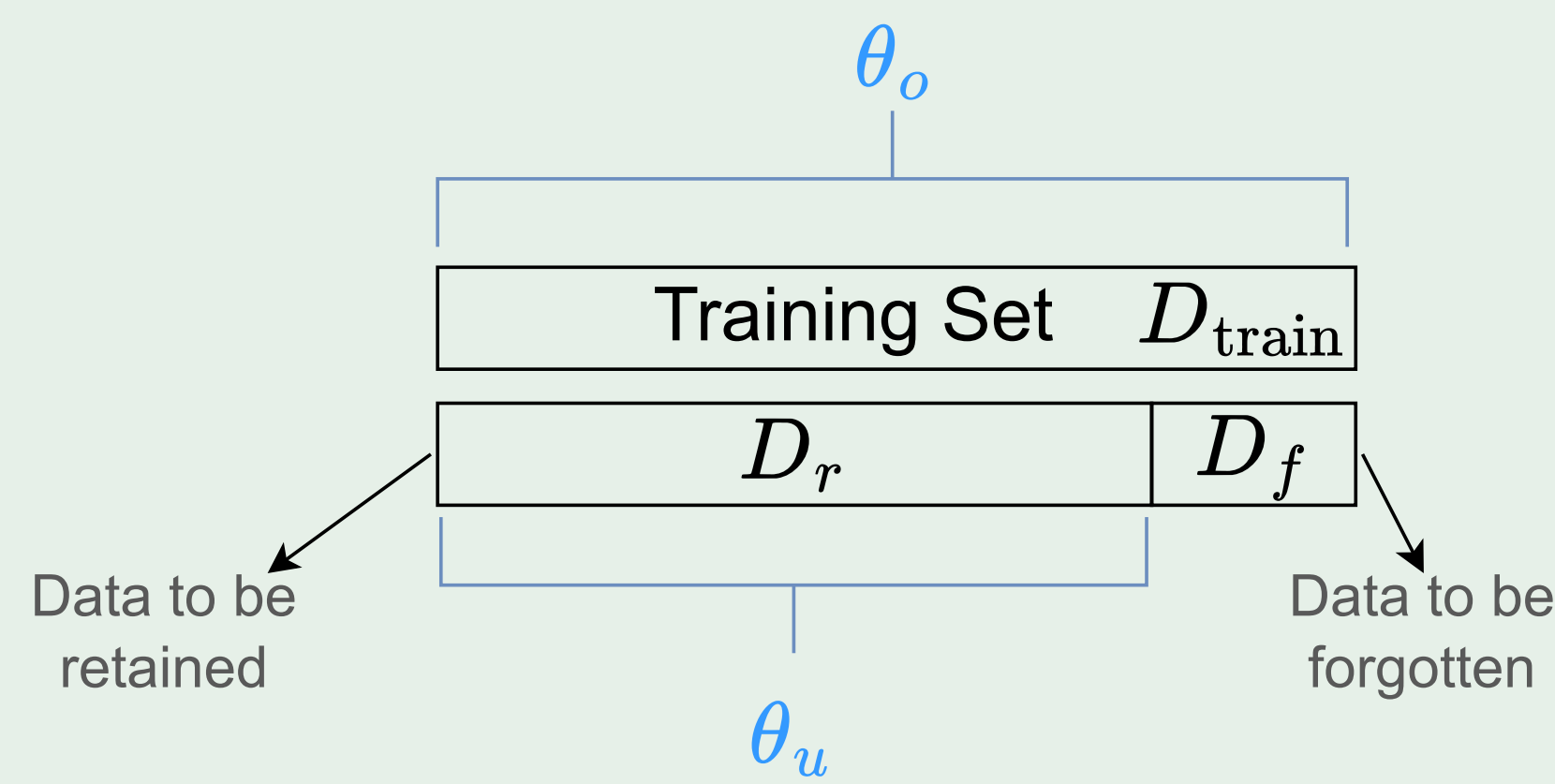
This work introduces ECO, an efficient computational optimization framework that adapts the CP algorithm—originally proposed by Cauwenberghs & Poggio (2000)—for unlearning within deep neural network models.

- ECO simplifies labor-intensive tasks, significantly reducing the workload for service providers compared to previous exact unlearning methods for DNNs.
- We demonstrate that ECO not only boosts efficiency but also maintains the performance of the original base DNN model, and surprisingly, it even surpasses naive retraining (NR) in effectiveness.
- Crucially, we are the first to adapt the CP algorithm’s decremental learning for leave-one-out evaluation to achieve exact unlearning in DNN models. We also open-source a usable base code for the CP algorithm, addressing the previous lack of such resources and encouraging further research and practical applications.



Problem Definition

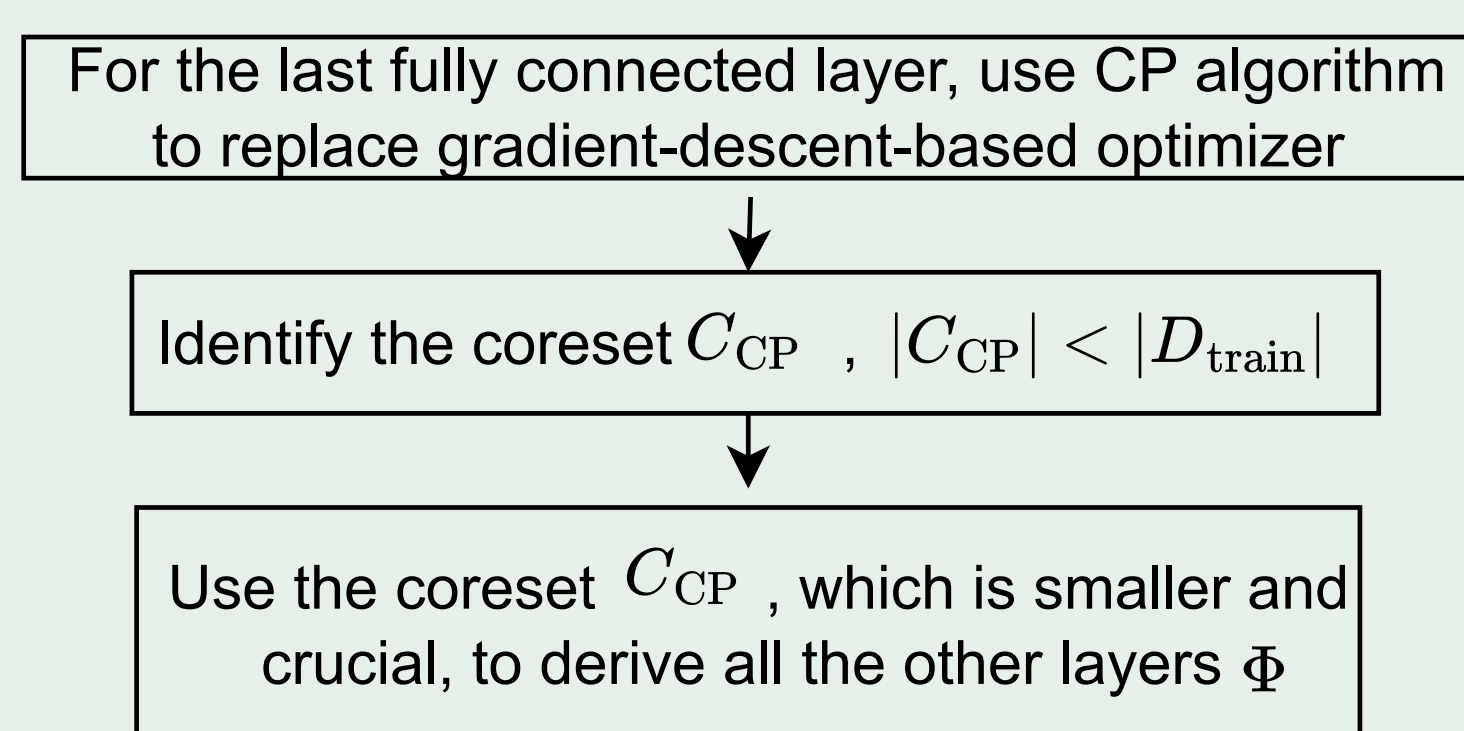
Machine unlearning aims to eliminate the influence of specific training data from an already-trained machine learning model θ_o .



The primary challenge of machine unlearning is to develop an efficient and effective method to transition from θ_o to θ_u .

Our Proposed Method: ECO

Model Preparation:



Model Serving:

Input $D_f, f_{C_{CP}}(x), C_{CP}$

Output $f_{C_{CP}}(x), C_{CP}$

- if $D_f \cap C_{CP} \neq \emptyset$ then
- if $D_f \cap (M \cup I) \neq \emptyset$ then
- $\alpha, b \leftarrow$ Employ Algorithm 3 in our paper to unlearn D_f
- $C_{CP} \leftarrow$ Construct new C_{CP} via (9) in our paper
- else
- $C_{CP} \leftarrow C_{CP} \setminus D_f$
- end if
- $\Phi_{C_{CP}} \leftarrow$ Learn a new feature transformation function $\Phi_{C_{CP}}(x)$ with (10) in our paper
- else
- Remain the input model $f_{C_{CP}}(x)$ and the set C_{CP}
- end if

Main Results

We compare the proposed method, ECO, with the gold standard naive retraining (NR), which involves retraining the DNN models upon receiving an unlearning request.

Table 1. In-time unlearning

	NR	ECO _i	ECO
Acc _r ↑	0.988 ± 0.005	<u>0.996 ± 0.001</u>	0.997 ± 0.001
Acc _f ↑	1.0 ± 0.0	1.0 ± 0.0	1.0 ± 0.0
Acc _{test} ↑	0.986 ± 0.002	0.989 ± 0.001	0.989 ± 0.001
Acc _{all} ↑	0.974 ± 0.007	<u>0.986 ± 0.002</u>	0.987 ± 0.002
Time cost (sec) ↓	20.968 ± 7.992	1.493 ± 4.476	<u>1.842 ± 5.524</u>

Table 2. Off-time batch unlearning

p	Acc _{all} (= Acc _r × Acc _f × Acc _{test}) ↑		
	NR	ECO _i	ECO
1	0.965 ± 0.006	<u>0.974 ± 0.007</u>	0.979 ± 0.003
5	0.964 ± 0.005	<u>0.973 ± 0.004</u>	0.979 ± 0.002
10	0.963 ± 0.006	<u>0.974 ± 0.005</u>	0.978 ± 0.004
20	0.960 ± 0.005	<u>0.973 ± 0.002</u>	0.978 ± 0.004
30	0.961 ± 0.005	<u>0.973 ± 0.002</u>	0.978 ± 0.003
40	0.961 ± 0.005	<u>0.971 ± 0.003</u>	0.979 ± 0.003
50	0.961 ± 0.006	<u>0.966 ± 0.004</u>	0.979 ± 0.004

All the metrics mentioned above are expressed as $a \pm b$, where ‘a’ represents the mean and ‘b’ denotes the standard deviation across 10 independent trials with different random seeds. The symbol ‘↑’ indicates that higher values are better, and ‘↓’ indicates that lower values are preferable. The best result is highlighted in bold, and the second-best result is underlined.

The Forgetfulness Quality: MIA

Table 3. The forgetfulness quality (MIA)

p	NR		ECO _i		ECO	
	D_f	D_{test}	D_f	D_{test}	D_f	D_{test}
1	0.06 ± 0.01	0.07 ± 0.0	0.15 ± 0.25	0.16 ± 0.24	0.03 ± 0.01	0.04 ± 0.01
5	0.14 ± 0.26	0.15 ± 0.26	0.07 ± 0.02	0.08 ± 0.02	0.03 ± 0.01	0.04 ± 0.01
10	0.15 ± 0.26	0.15 ± 0.26	0.07 ± 0.02	0.08 ± 0.01	0.04 ± 0.01	0.05 ± 0.01
20	0.32 ± 0.39	0.32 ± 0.39	0.07 ± 0.02	0.07 ± 0.01	0.04 ± 0.01	0.04 ± 0.01
30	0.15 ± 0.26	0.15 ± 0.26	0.06 ± 0.01	0.07 ± 0.01	0.04 ± 0.01	0.05 ± 0.01
40	0.15 ± 0.26	0.15 ± 0.26	0.06 ± 0.02	0.07 ± 0.01	0.04 ± 0.01	0.04 ± 0.01
50	0.24 ± 0.34	0.24 ± 0.35	0.08 ± 0.02	0.08 ± 0.02	0.04 ± 0.01	0.05 ± 0.01

Table 3 compares the MIA scores between D_f and D_{test} for each unlearning approach. A lower discrepancy between these two values indicates better forgetfulness quality of the model. For the NR method, there is hardly any disparity between D_f and D_{test} . This pattern is also consistent for ECO_i and ECO. As expected, this alignment occurs because all three methods adhere to an exact unlearning approach.

For more details, please refer to:

<https://openreview.net/pdf?id=SeBVP0zxKp>

