

Rebuttal Response

Anonymous Authors¹

1. Experiments of Class-incremental Learning

Table 1. Results of class-incremental experiments on Multiple Dataset (p denotes the noise rate).

Methods	$p = 0.0$		$p = 0.2$	
	ACC (\uparrow)	FGT (\downarrow)	ACC (\uparrow)	FGT (\downarrow)
EWC	13.08 \pm 5.36	0.695 \pm 0.004	10.36 \pm 3.21	0.590 \pm 0.009
MAS	10.06 \pm 1.13	0.613 \pm 0.019	6.04 \pm 0.62	0.487 \pm 0.016
AGEM	13.67 \pm 1.24	0.654 \pm 0.018	12.10 \pm 1.24	0.548 \pm 0.007
OGD	13.00 \pm 0.74	0.707 \pm 0.005	9.75 \pm 0.44	0.586 \pm 0.005
DER	45.67 \pm 2.19	0.260 \pm 0.039	40.88 \pm 2.15	0.248 \pm 0.026
GDumb	40.69 \pm 3.32	0.128\pm0.023	36.25 \pm 2.23	0.083\pm0.019
MEGA	38.94 \pm 2.79	0.389 \pm 0.026	34.03 \pm 2.23	<u>0.134\pm0.023</u>
STREAM	47.92\pm0.45	<u>0.153\pm0.018</u>	43.32\pm0.018	0.367 \pm 0.013

Table 2. Results of class-incremental experiments on Split CIFAR-10 (p denotes the noise rate).

Methods	$p = 0.0$		$p = 0.2$	
	ACC (\uparrow)	FGT (\downarrow)	ACC (\uparrow)	FGT (\downarrow)
EWC	18.36 \pm 0.23	0.714 \pm 0.018	18.01 \pm 0.23	0.717 \pm 0.007
MAS	14.36 \pm 1.13	0.275\pm0.021	11.66 \pm 1.59	0.546 \pm 0.020
AGEM	20.50 \pm 1.92	0.693 \pm 0.035	19.44 \pm 1.14	0.525 \pm 0.071
OGD	19.02 \pm 0.12	0.734 \pm 0.003	18.99 \pm 0.08	0.736 \pm 0.004
DER	29.87 \pm 2.52	0.586 \pm 0.084	22.62 \pm 3.27	0.685 \pm 0.029
GDumb	28.40 \pm 2.13	<u>0.365\pm0.054</u>	21.69 \pm 0.59	0.227\pm0.010
MEGA	<u>33.16\pm3.12</u>	0.471 \pm 0.106	<u>24.84\pm2.26</u>	0.475 \pm 0.091
STREAM	35.35\pm3.52	0.504 \pm 0.098	26.17\pm1.32	<u>0.303\pm0.056</u>

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

2. Running Time Comparison

Table 3. Running time on Multiple Dataset and Split CIFAR-100.

Methods	Multiple Dataset (hours)	Split CIFAR-100 (hours)
EWC	0.16	1.31
MAS	0.17	1.31
AGEM	0.16	1.30
OGD	0.47	3.01
DER	0.18	1.27
GDumb	0.13	0.92
MEGA	0.15	1.05
STREAM	0.11	0.79

3. Failure Example

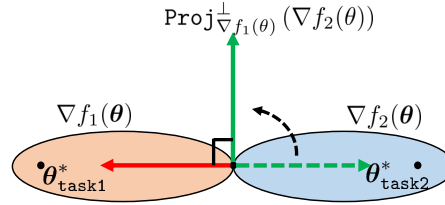


Figure 1. The Failure example for GEM, A-GEM and OGD, where the current gradient $\nabla f_2(\theta)$ is in the opposite direction to the memory gradient $\nabla f_1(\theta)$. The update direction is based on a rotation of the current gradient, which is orthogonal to the memory gradient: $\text{Proj}_{\nabla f_1(\theta)}^\perp(\nabla f_2(\theta))$. Therefore it cannot learn task 2.

4. Results with Multiple Runs

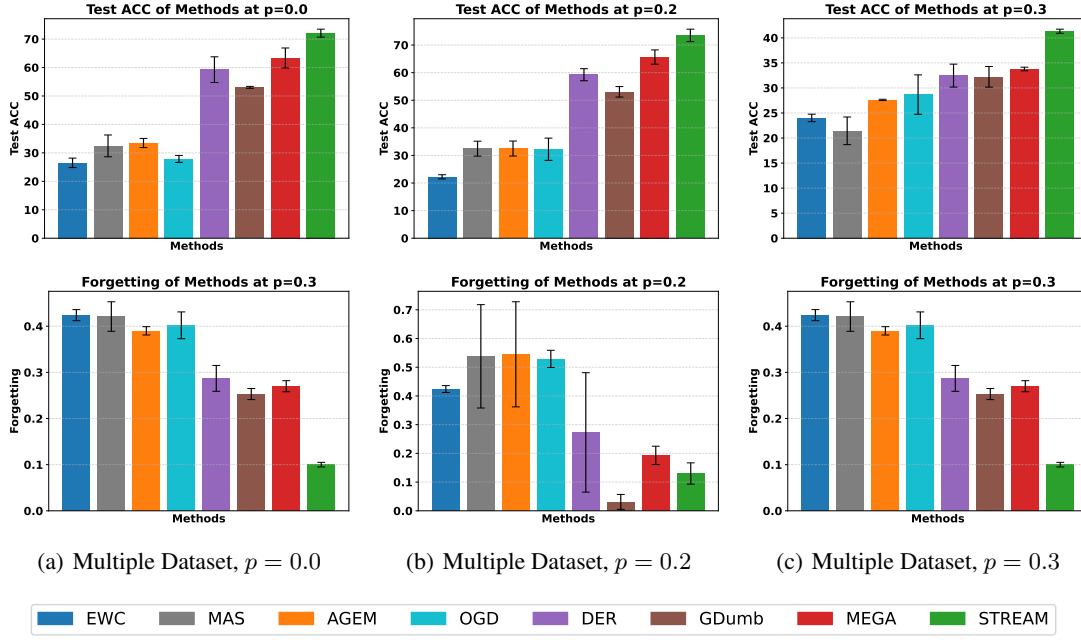


Figure 2. Performance (test accuracy and forgetting) of continual learning methods on Multiple Dataset over 5 runs.

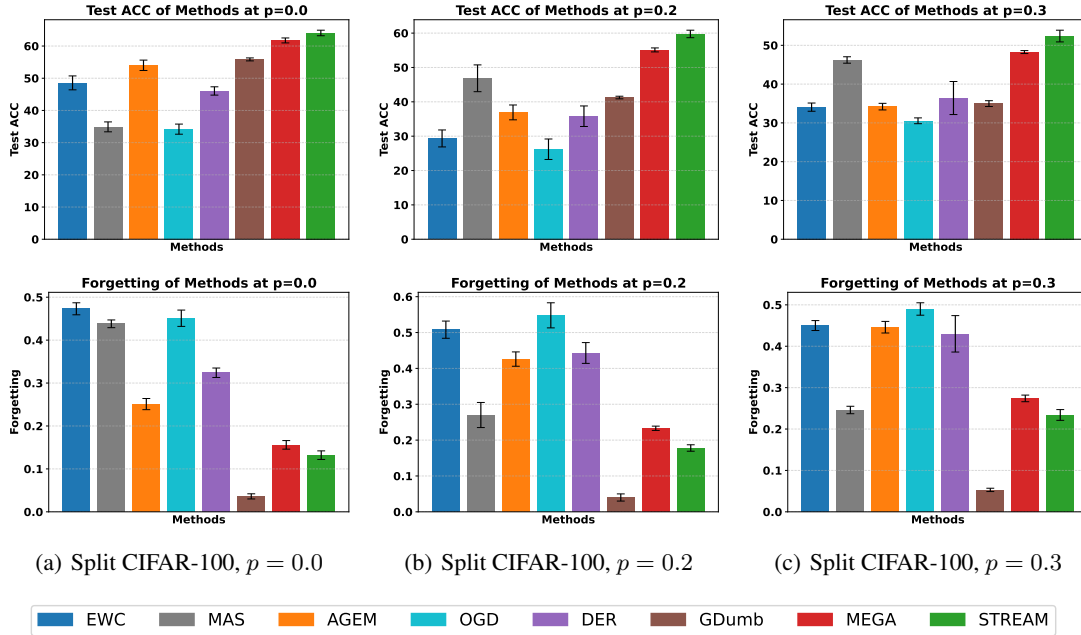


Figure 3. Performance (test accuracy and forgetting) of continual learning methods on Split CIFAR100 over 5 runs.

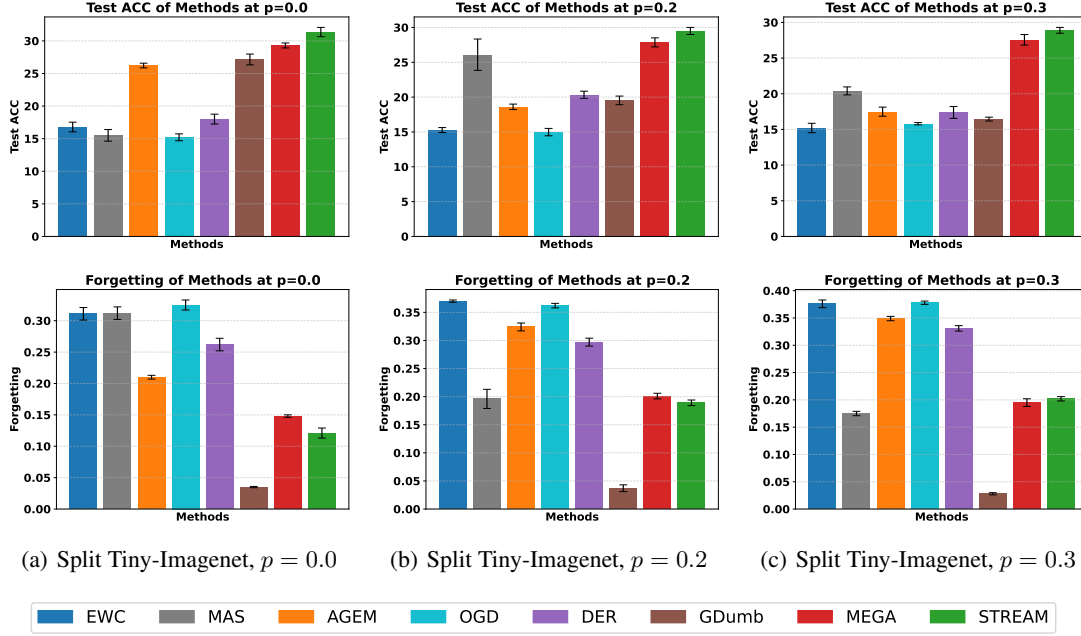


Figure 4. Performance (test accuracy and forgetting) of continual learning methods on Split Tiny-Imagenet over 5 runs.

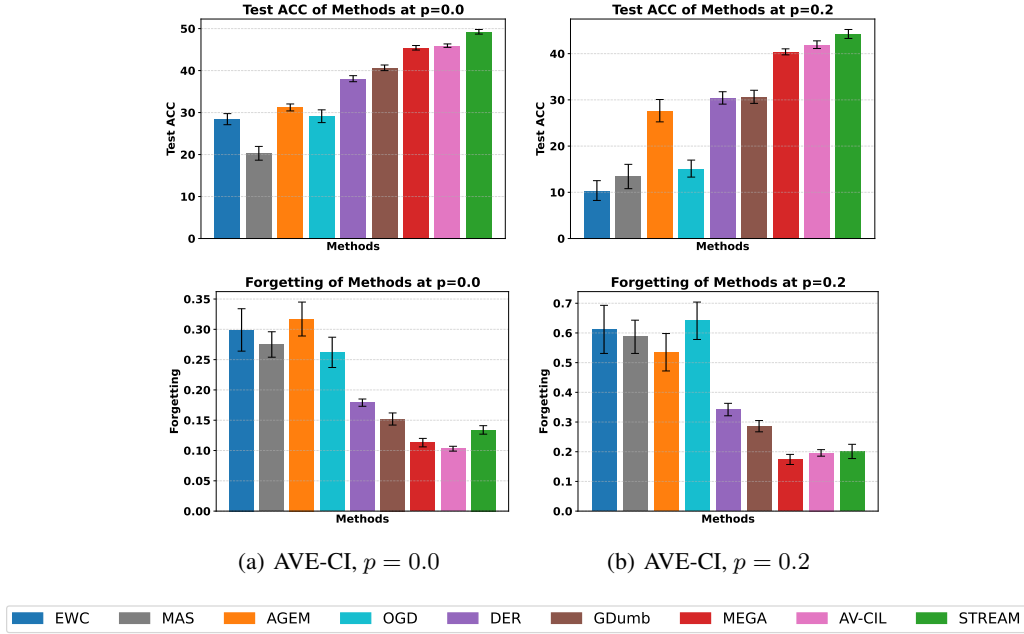
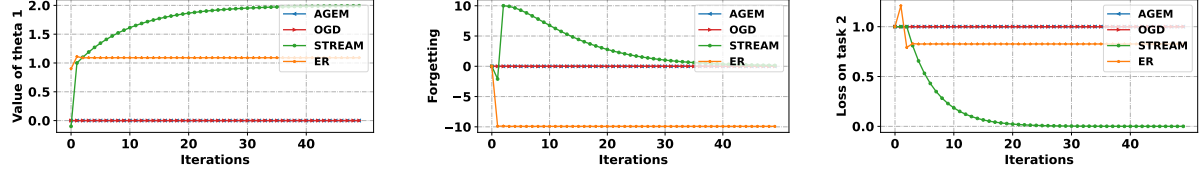


Figure 5. Performance (test accuracy and forgetting) of continual learning methods on AVE-CI over 5 runs

5. Synthetic Experiment for the counterexample



(a) The value of the $\theta^{(1)}$ vs. iterations. (b) The forgetting vs. iterations. (c) Loss on task 2 vs. iterations.

Figure 6. Synthetic experiment for the counterexample. (a), (b), (c) show the evolution of the value of $\theta^{(1)}$, the forgetting on task 1 and the loss on task 2. STREAM can find the optimal $\theta^{(1)}$ and achieve minimal forgetting and loss on the new task. But A-GEM and OGD fail to update their parameter $\theta^{(1)}$ throughout the training process, thus cannot minimize the loss on the new task. ER minimizes the loss functions f_1 and f_2 jointly (Assume that memory is large enough, it can visit f_1 as it needs). ER cannot find the optimal value of $\theta^{(1)}$, thus can not achieve the low loss on task2. ER performs well in terms of forgetting, but STREAM still exhibits good forgetting metric (forgetting = 0).

6. The Frequency Statistics of Model Update on Current/Memory Data

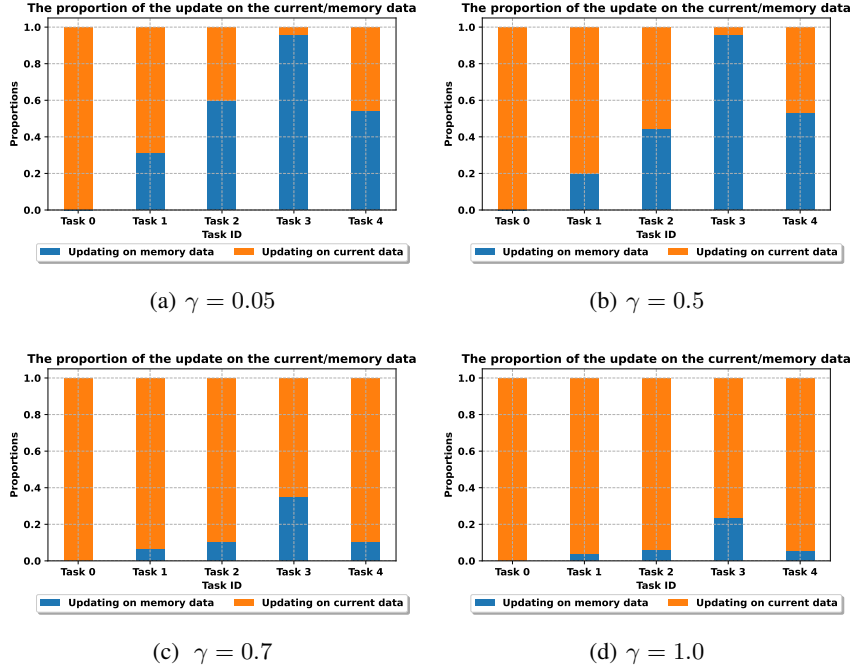


Figure 7. The frequency statistics of model update on the current/memory data.

7. Hyperparameter Tuning

Table 4. Results vs. memory size (m is the memory size) on Multiple Dataset.

Methods	$m = 64$		$m = 128$		$m = 256$	
	ACC (\uparrow)	FGT (\downarrow)	ACC (\uparrow)	FGT (\downarrow)	ACC (\uparrow)	FGT (\downarrow)
EWC	26.50 \pm 1.66	0.377 \pm 0.013	26.67 \pm 1.40	0.245 \pm 0.017	26.59 \pm 0.11	0.223 \pm 0.008
MAS	32.46 \pm 3.83	0.538 \pm 0.181	33.54 \pm 1.47	0.538 \pm 0.041	34.15 \pm 1.63	0.536 \pm 0.039
AGEM	33.47 \pm 1.61	0.541 \pm 0.179	36.04 \pm 0.94	0.480 \pm 0.005	41.90 \pm 1.52	0.424 \pm 0.026
OGD	27.88 \pm 1.23	0.375 \pm 0.015	32.24 \pm 1.34	0.544 \pm 0.051	33.63 \pm 1.24	0.303 \pm 0.022
DER	59.25 \pm 4.52	0.273 \pm 0.106	59.78 \pm 0.94	0.225 \pm 0.015	60.53 \pm 2.67	0.240 \pm 0.028
GDumb	53.03 \pm 0.34	0.032\pm0.026	58.64 \pm 0.42	0.000\pm0.006	69.11 \pm 2.83	0.000\pm0.018
MEGA	<u>63.36\pm3.51</u>	0.210 \pm 0.109	<u>66.44\pm2.03</u>	0.292 \pm 0.015	<u>67.24\pm2.38</u>	0.214 \pm 0.003
STREAM	72.08\pm1.40	<u>0.152\pm0.035</u>	74.14\pm2.01	<u>0.145\pm0.016</u>	74.07\pm0.95	<u>0.025\pm0.048</u>

Table 5. Results vs. memory size (m is the memory size) Split Tiny ImageNet.

Methods	$m = 256$		$m = 512$		$m = 1024$	
	ACC (\uparrow)	FGT (\downarrow)	ACC (\uparrow)	FGT (\downarrow)	ACC (\uparrow)	FGT (\downarrow)
EWC	16.79 \pm 0.74	0.311 \pm 0.010	17.15 \pm 0.44	0.306 \pm 0.002	17.05 \pm 0.10	0.324 \pm 0.003
MAS	15.51 \pm 0.89	0.312 \pm 0.010	17.03 \pm 0.08	0.332 \pm 0.008	18.56 \pm 0.52	0.316 \pm 0.000
AGEM	26.22 \pm 0.36	0.210 \pm 0.003	27.95 \pm 0.58	0.217 \pm 0.004	32.97 \pm 0.15	0.181 \pm 0.004
OGD	15.21 \pm 0.53	0.325 \pm 0.008	16.53 \pm 0.67	0.357 \pm 0.007	17.15 \pm 0.52	0.347 \pm 0.005
DER	18.00 \pm 0.76	0.262 \pm 0.010	21.46 \pm 1.13	0.236 \pm 0.005	23.10 \pm 0.51	0.214 \pm 0.007
GDumb	27.15 \pm 0.83	0.035\pm0.001	32.81 \pm 0.25	0.040\pm0.004	<u>33.51\pm0.14</u>	0.021\pm0.002
MEGA	<u>29.30\pm0.38</u>	0.148 \pm 0.002	<u>32.87\pm0.63</u>	0.137 \pm 0.007	33.20 \pm 0.44	<u>0.104\pm0.001</u>
STREAM	31.36\pm0.71	<u>0.121\pm0.008</u>	33.02\pm0.53	<u>0.115\pm0.003</u>	34.67\pm0.14	0.104 \pm 0.002

Table 6. Results vs. the number of tasks (T denotes the number of tasks) on Split CIFAR-100.

Methods	$T = 10$		$T = 20$		$T = 25$	
	ACC (\uparrow)	FGT (\downarrow)	ACC (\uparrow)	FGT (\downarrow)	ACC (\uparrow)	FGT (\downarrow)
EWC	27.98 \pm 1.57	0.403 \pm 0.013	48.56 \pm 2.16	0.473 \pm 0.014	47.43 \pm 3.23,	0.530 \pm 0.033
MAS	31.40 \pm 2.34	0.380 \pm 0.021	34.90 \pm 1.54	0.438 \pm 0.009	31.68 \pm 2.16	0.496 \pm 0.022
AGEM	39.72 \pm 2.43	0.292 \pm 0.019	54.02 \pm 1.61	0.251 \pm 0.013	48.19 \pm 1.65	0.345 \pm 0.019
OGD	29.13 \pm 1.52	0.262 \pm 0.025	34.19 \pm 1.57	0.451 \pm 0.019	36.35 \pm 1.43	0.353 \pm 0.023
DER	41.99 \pm 1.80	0.264 \pm 0.018	46.05 \pm 1.29	0.324 \pm 0.011	49.65 \pm 0.90	0.334 \pm 0.008
GDumb	38.80 \pm 1.08	0.058\pm0.002	55.85 \pm 0.46	0.036\pm0.006	67.56 \pm 0.53	0.008\pm0.009
MEGA	<u>48.03\pm1.44</u>	0.176 \pm 0.003	<u>61.74\pm0.77</u>	0.156 \pm 0.010	<u>68.55\pm0.57</u>	0.094 \pm 0.007
STREAM	50.33\pm0.66	<u>0.167\pm0.008</u>	64.06\pm0.86	<u>0.132\pm0.010</u>	69.70\pm0.37	<u>0.080\pm0.010</u>

8. New Advanced Baselines

Table 7. Results of task-incremental experiments on Multiple Dataset and Split Tiny-Imagenet (p denotes the noise rate).

	Methods	$p = 0.0$		$p = 0.2$	
		ACC (\uparrow)	FGT (\downarrow)	ACC (\uparrow)	FGT (\downarrow)
Multiple Dataset	OCS	55.65 \pm 2.26	0.062\pm0.001	45.03 \pm 4.16	0.049\pm0.012
	MetaSP	57.14 \pm 1.10	0.113 \pm 0.042	47.14 \pm 1.66	0.081 \pm 0.027
	STREAM	72.08\pm1.40	0.152 \pm 0.035	73.50\pm2.25	0.130 \pm 0.037
Split Tiny-Imagenet	OCS	41.29 \pm 0.09	0.112\pm0.001	35.36 \pm 0.94	0.061\pm0.005
	MetaSP	43.33 \pm 0.32	0.127 \pm 0.002	37.18 \pm 0.76	0.068 \pm 0.007
	STREAM	47.92\pm0.45	0.153 \pm 0.018	43.32\pm0.02	0.367 \pm 0.013