

# Rebuttal Response

Anonymous Authors<sup>1</sup>

## 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	<b>0.128<math>\pm</math>0.023</b>	36.25 $\pm$ 2.23	<b>0.083<math>\pm</math>0.019</b>
MEGA	38.94 $\pm$ 2.79	0.389 $\pm$ 0.026	34.03 $\pm$ 2.23	<u>0.134<math>\pm</math>0.023</u>
STREAM	<b>47.92<math>\pm</math>0.45</b>	<u>0.153<math>\pm</math>0.018</u>	<b>43.32<math>\pm</math>0.018</b>	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	<b>0.275<math>\pm</math>0.021</b>	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<math>\pm</math>0.054</u>	21.69 $\pm$ 0.59	<b>0.227<math>\pm</math>0.010</b>
MEGA	<u>33.16<math>\pm</math>3.12</u>	0.471 $\pm$ 0.106	<u>24.84<math>\pm</math>2.26</u>	0.475 $\pm$ 0.091
STREAM	<b>35.35<math>\pm</math>3.52</b>	0.504 $\pm$ 0.098	<b>26.17<math>\pm</math>1.32</b>	<u>0.303<math>\pm</math>0.056</u>

<sup>1</sup>Anonymous 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	<b>0.11</b>	<b>0.79</b>

## 3. Failure Example

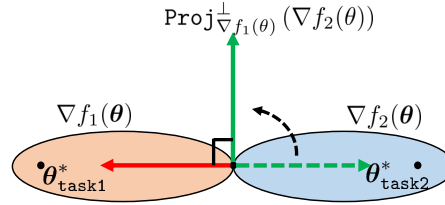


Figure 1. The Failure example for GEM, A-GEM and OGD, where the current gradient  $\nabla f_2(\theta)$  is in the opposite direction to the memory gradient  $\nabla f_1(\theta)$ . The update direction is based on a rotation of the current gradient, which is orthogonal to the memory gradient:  $\text{Proj}_{\nabla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

## 4. Results with Multiple Runs

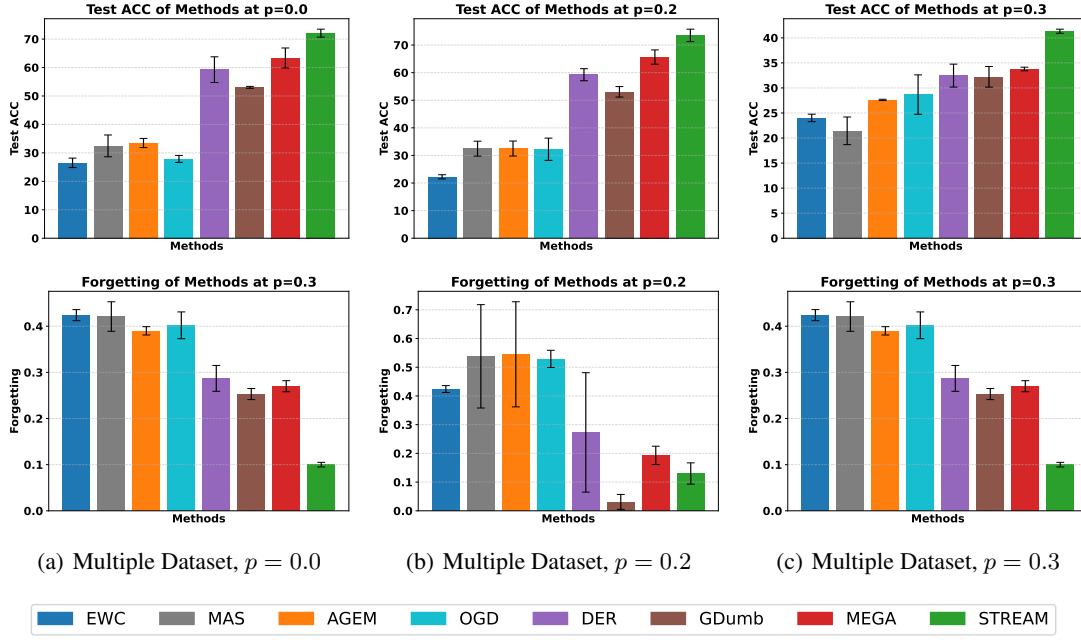


Figure 2. Performance (test accuracy and forgetting) of continual learning methods on Multiple Dataset over 5 runs.

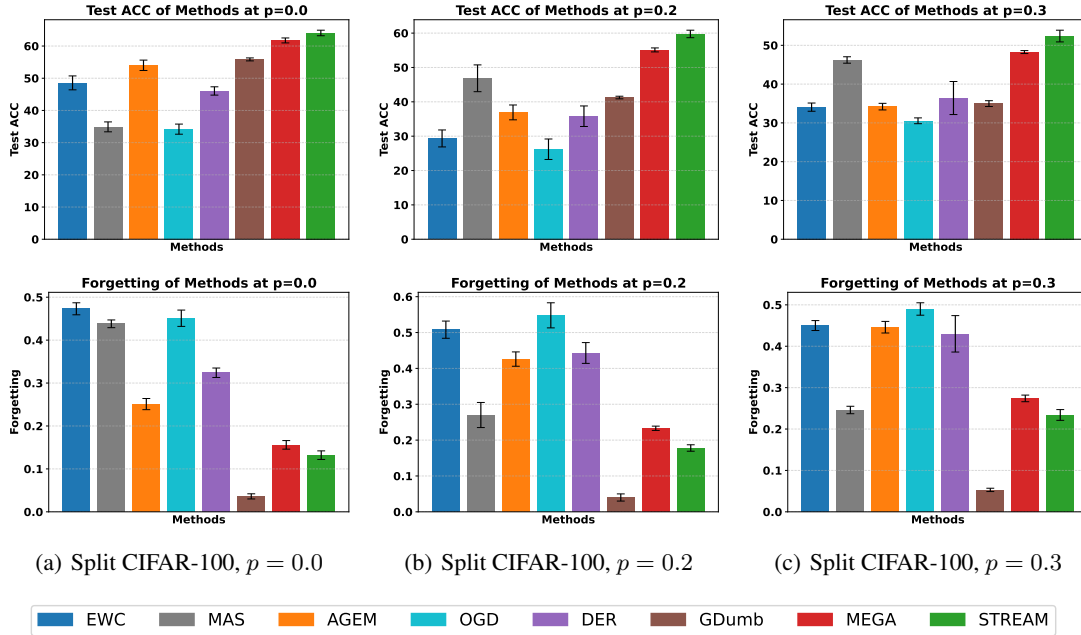


Figure 3. Performance (test accuracy and forgetting) of continual learning methods on Split CIFAR100 over 5 runs.

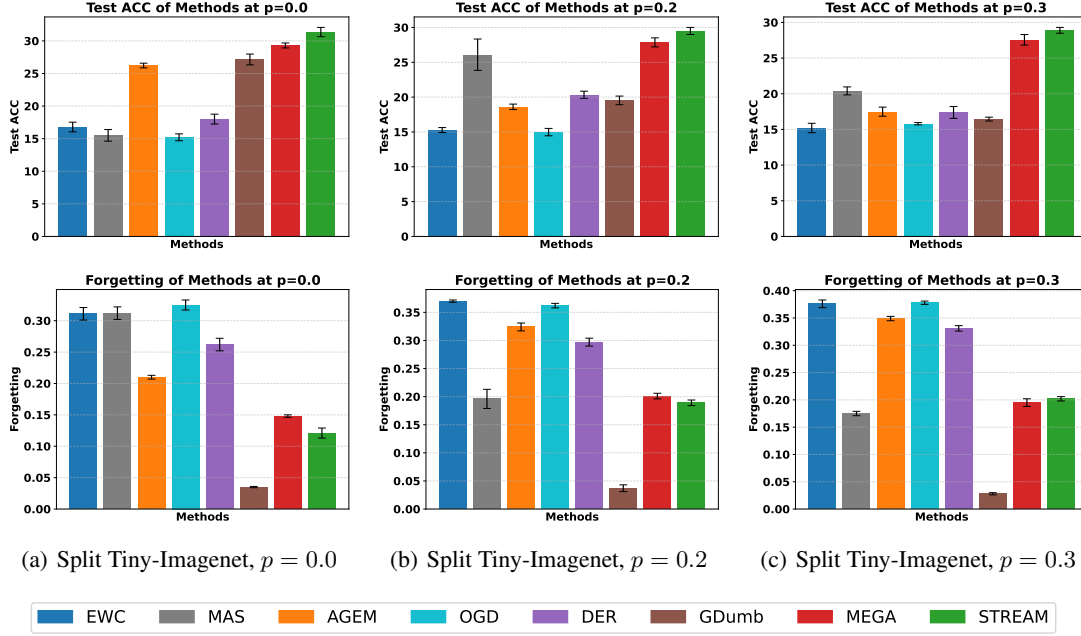


Figure 4. Performance (test accuracy and forgetting) of continual learning methods on Split Tiny-Imagenet over 5 runs.

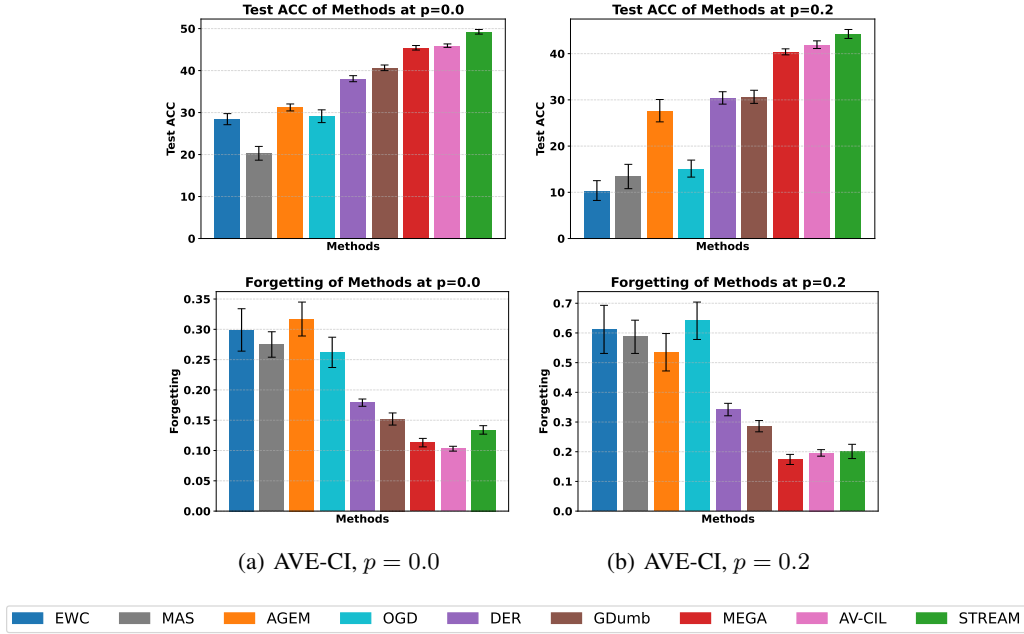
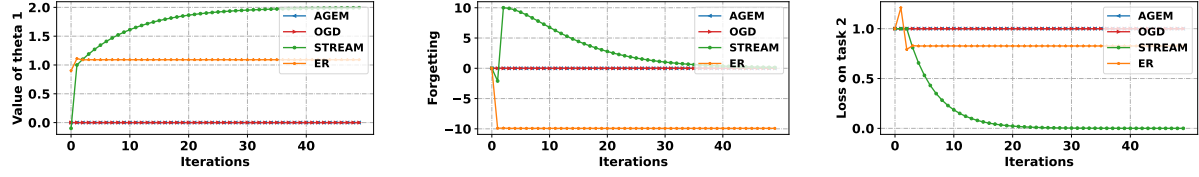


Figure 5. Performance (test accuracy and forgetting) of continual learning methods on AVE-CI over 5 runs

## 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

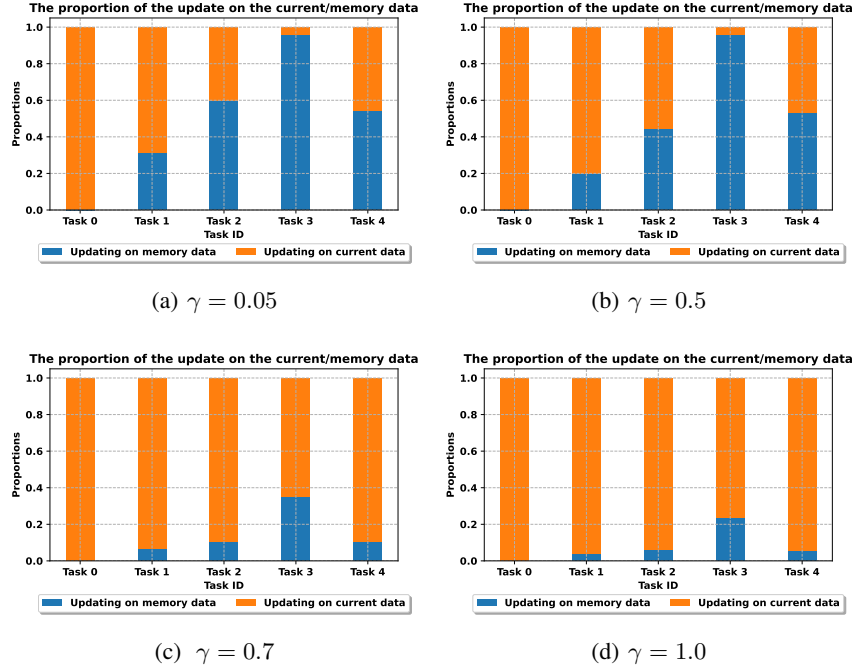


Figure 7. The frequency statistics of model update on the current/memory data.

## 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	<b>0.032<math>\pm</math>0.026</b>	58.64 $\pm$ 0.42	<b>0.000<math>\pm</math>0.006</b>	69.11 $\pm$ 2.83	<b>0.000<math>\pm</math>0.018</b>
MEGA	<u>63.36<math>\pm</math>3.51</u>	0.210 $\pm$ 0.109	<u>66.44<math>\pm</math>2.03</u>	0.292 $\pm$ 0.015	<u>67.24<math>\pm</math>2.38</u>	0.214 $\pm$ 0.003
STREAM	<b>72.08<math>\pm</math>1.40</b>	<u>0.152<math>\pm</math>0.035</u>	<b>74.14<math>\pm</math>2.01</b>	<u>0.145<math>\pm</math>0.016</u>	<b>74.07<math>\pm</math>0.95</b>	<u>0.025<math>\pm</math>0.048</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	<b>0.035<math>\pm</math>0.001</b>	32.81 $\pm$ 0.25	<b>0.040<math>\pm</math>0.004</b>	<u>33.51<math>\pm</math>0.14</u>	<b>0.021<math>\pm</math>0.002</b>
MEGA	<u>29.30<math>\pm</math>0.38</u>	0.148 $\pm$ 0.002	<u>32.87<math>\pm</math>0.63</u>	0.137 $\pm$ 0.007	33.20 $\pm$ 0.44	<u>0.104<math>\pm</math>0.001</u>
STREAM	<b>31.36<math>\pm</math>0.71</b>	<u>0.121<math>\pm</math>0.008</u>	<b>33.02<math>\pm</math>0.53</b>	<u>0.115<math>\pm</math>0.003</u>	<b>34.67<math>\pm</math>0.14</b>	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	<b>0.058<math>\pm</math>0.002</b>	55.85 $\pm$ 0.46	<b>0.036<math>\pm</math>0.006</b>	67.56 $\pm$ 0.53	<b>0.008<math>\pm</math>0.009</b>
MEGA	<u>48.03<math>\pm</math>1.44</u>	0.176 $\pm$ 0.003	<u>61.74<math>\pm</math>0.77</u>	0.156 $\pm$ 0.010	<u>68.55<math>\pm</math>0.57</u>	0.094 $\pm$ 0.007
STREAM	<b>50.33<math>\pm</math>0.66</b>	<u>0.167<math>\pm</math>0.008</u>	<b>64.06<math>\pm</math>0.86</b>	<u>0.132<math>\pm</math>0.010</u>	<b>69.70<math>\pm</math>0.37</b>	<u>0.080<math>\pm</math>0.010</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	<b>0.062<math>\pm</math>0.001</b>	45.03 $\pm$ 4.16	<b>0.049<math>\pm</math>0.012</b>
	MetaSP	57.14 $\pm$ 1.10	0.113 $\pm$ 0.042	47.14 $\pm$ 1.66	0.081 $\pm$ 0.027
	STREAM	<b>72.08<math>\pm</math>1.40</b>	0.152 $\pm$ 0.035	<b>73.50<math>\pm</math>2.25</b>	0.130 $\pm$ 0.037
Split Tiny-Imagenet	OCS	28.16 $\pm$ 0.09	<b>0.112<math>\pm</math>0.001</b>	20.42 $\pm$ 0.94	<b>0.061<math>\pm</math>0.005</b>
	MetaSP	27.66 $\pm$ 0.51	0.127 $\pm$ 0.002	22.43 $\pm$ 0.34	0.068 $\pm$ 0.007
	STREAM	<b>31.36<math>\pm</math>0.71</b>	0.121 $\pm$ 0.008	<b>29.50<math>\pm</math>0.51</b>	0.189 $\pm$ 0.005

## 9. Ablation Study for Large Batch Size

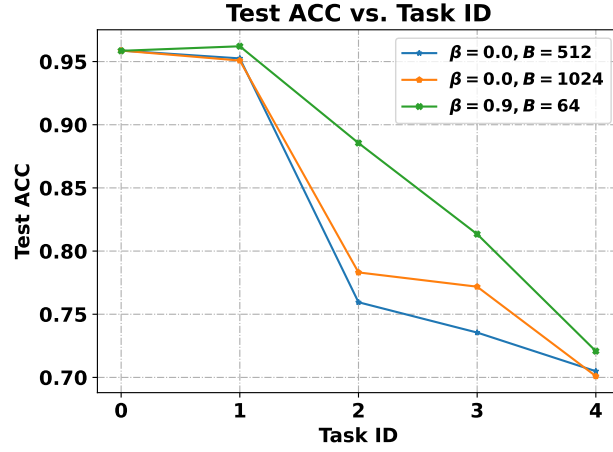


Figure 8. Ablation study for large batch size, where  $\beta$  is momentum parameter,  $B$  is the batch size, and " $\beta = 0.9, B = 64$ " denotes STREAM algorithm, and " $\beta = 0.0, B = 512$ ", " $\beta = 0.0, B = 1024$ " denote SSG with different batch size.

## 10. Results on Split Tiny-Imagenet with DER Settings

Table 8. Hyperparameter settings for Split Tiny-Imagenet.

Methods	Hyperparameter settings
EWC	lr: 0.03, batch_size: 32, $\lambda$ : 25
MAS	lr: 0.03, batch_size: 32, $\lambda$ : 1.0
AGEM	lr: 0.01, batch_size: 32
OGD	lr: 0.01, batch_size: 32
DER	lr: 0.03, batch_size: 32, softmax_temp: 2.0, $\alpha$ : 0.1
GDumb	lr: 0.10, batch_size: 32
MEGA	lr: 0.10, batch_size: 32
STREAM	lr: 0.05, batch_size: 32, $\gamma$ : 0.05

Table 9. Results on Tiny-Imagenet (with DER settings).

Methods	Class-incremental		Task-incremental	
	ACC ( $\uparrow$ )	FGT ( $\downarrow$ )	ACC ( $\uparrow$ )	FGT ( $\downarrow$ )
EWC	7.41 $\pm$ 0.19	0.737 $\pm$ 0.007	15.67 $\pm$ 0.99	0.646 $\pm$ 0.010
MAS	6.91 $\pm$ 0.23	0.676 $\pm$ 0.010	23.99 $\pm$ 0.33	0.143 $\pm$ 0.007
AGEM	7.69 $\pm$ 0.02	<u>0.335<math>\pm</math>0.010</u>	24.57 $\pm$ 0.94	0.480 $\pm$ 0.005
OGD	7.90 $\pm$ 0.04	0.778 $\pm$ 0.003	17.02 $\pm$ 0.82	0.676 $\pm$ 0.010
DER	8.26 $\pm$ 0.84	0.698 $\pm$ 0.570	40.56 $\pm$ 0.82	0.302 $\pm$ 0.033
GDumb	7.34 $\pm$ 0.67	<b>0.032<math>\pm</math>0.026</b>	39.34 $\pm$ 0.67	<b>0.042<math>\pm</math>0.004</b>
MEGA	<u>8.39<math>\pm</math>0.20</u>	0.743 $\pm$ 0.016	<u>41.45<math>\pm</math>0.97</u>	0.465 $\pm$ 0.013
STREAM	<b>8.57<math>\pm</math>0.13</b>	0.726 $\pm$ 0.002	<b>42.04<math>\pm</math>0.82</b>	0.280 $\pm$ 0.016