

# **Learning to Learn without Forgetting by Maximizing Transfer and Minimizing Interference**

---

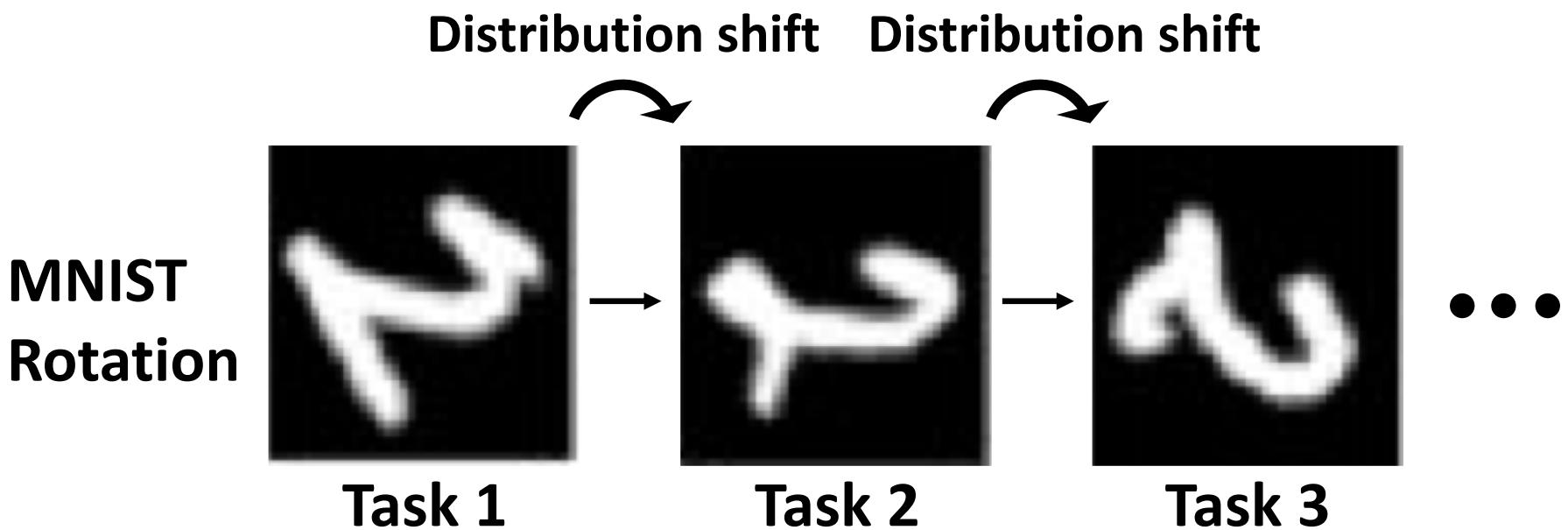
**PRESENTER: BYUNGSOO JEON**

**LLL READING GROUP @ CMU**

**\* THIS PAPER WAS PUBLISHED IN ICLR19**

# Motivation: Challenge in CL

Discovering notions of tasks without supervision while learning incrementally is hard in CL setting:



# Motivation: Challenge in CL

---

What makes incremental learning (or CL) so hard?

- Neural network forgets past tasks while learning new ones

Conceptualization of this challenge

- Catastrophic forgetting (Lack of stability in NN)
  - Focus was on limiting weight sharing
- Stability-plasticity dilemma
  - Focus was on balancing limited weight sharing while ensuring fast learning
- **Transfer-interference trade-off**
  - ?

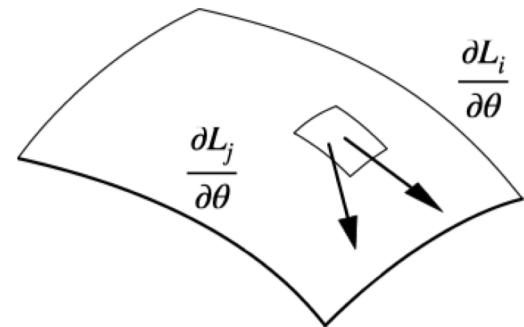
# Transfer and Interference

---

Operational measures of  
Transfer and Interference:

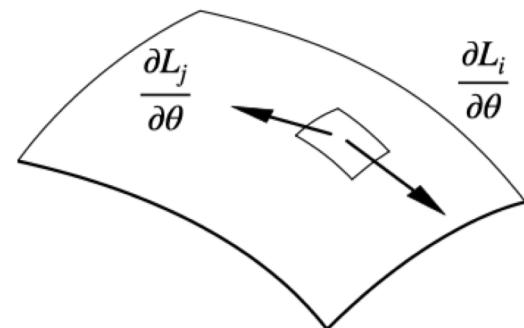
$$\frac{\partial L(x_i, y_i)}{\partial \theta} \cdot \frac{\partial L(x_j, y_j)}{\partial \theta} > 0$$

## B. Transfer



$$\frac{\partial L(x_i, y_i)}{\partial \theta} \cdot \frac{\partial L(x_j, y_j)}{\partial \theta} < 0$$

## C. Interference

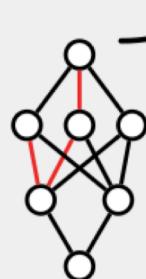


# Transfer-Interference Trade-off

## A. Transfer – Interference Trade-off

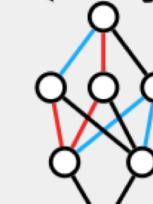
Stability – Plasticity Dilemma

Old  
Learning



Transfer

Current  
Learning



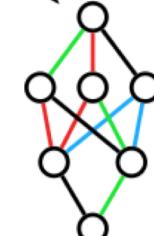
Sharing

Interference

PAST

Stability – Plasticity Dilemma

Future  
Learning



Transfer

Sharing

Interference

FUTURE

# Previous Approach: GEM

---

## Gradient Episodic Memory for Continual Learning

- Goal: Minimize negative backward transfer (catastrophic forgetting) by the efficient use of episodic memory
- Key idea: Use the closest gradient to the original that has positive gradient dot products with gradients from previous tasks

### Limitation

- 1) It does not consider positive forward transfer
- 2) Gradients are adjusted in an ad hoc manner  
(Solving a quadratic equation separately for each example)

# Proposed Method

---

## Meta Experience Replay (a.k.a. MER)

Problem: CL setting without task ID

Goal: Learn the optimal weight sharing scheme to balance between new learning and forgetting

### Advantages

- It learns a generalizable theory about weight sharing
- It enables transfer to improve future performance while not disrupting performance on previous tasks

# Learning to learn w/o forgetting

---

In typical offline supervised learning with **stationary** data distribution,

$$\theta = \arg \min_{\theta} \mathbb{E}_{(x,y) \sim D} [L(x, y)]$$

---

In the same setting with **non-stationary** data distribution, given two samples  $(x_i, y_i), (x_j, y_j)$ ,

$$\theta = \arg \min_{\theta} \mathbb{E}_{[(x_i, y_i), (x_j, y_j)] \sim D} [L(x_i, y_i) + L(x_j, y_j) - \alpha \frac{\partial L(x_i, y_i)}{\partial \theta} \cdot \frac{\partial L(x_j, y_j)}{\partial \theta}]$$

# Challenges

---

In the same setting with **non-stationary** data distribution, given two samples  $(x_i, y_i), (x_j, y_j)$ ,

$$\theta = \arg \min_{\theta} \mathbb{E}_{[(x_i, y_i), (x_j, y_j)] \sim D} [L(x_i, y_i) + L(x_j, y_j) - \alpha \frac{\partial L(x_i, y_i)}{\partial \theta} \cdot \frac{\partial L(x_j, y_j)}{\partial \theta}]$$

## Challenge

- 1) We have non-stationary stream of data
- 2) The gradient of this loss depends on the second derivative of the loss function (EXPENSIVE to compute)

# Learning to learn w/o forgetting

In the same setting with **non-stationary** data distribution, given two samples  $(x_i, y_i), (x_j, y_j)$ ,

$$\theta = \arg \min_{\theta} \mathbb{E}_{[(x_i, y_i), (x_j, y_j)] \sim D} [L(x_i, y_i) + L(x_j, y_j) - \alpha \frac{\partial L(x_i, y_i)}{\partial \theta} \cdot \frac{\partial L(x_j, y_j)}{\partial \theta}]$$

Optimize over the new meta-learning loss (Reptile):

$$\theta = \arg \min_{\theta} \mathbb{E}_{B_1, \dots, B_s \sim D} [2 \sum_{i=1}^s [L(B_i) - \sum_{j=1}^{i-1} \alpha \frac{\partial L(B_i)}{\partial \theta} \cdot \frac{\partial L(B_j)}{\partial \theta}]]$$

# Learning to learn w/o forgetting

Optimize over the new meta-learning loss (Reptile):

$$\theta = \arg \min_{\theta} \mathbb{E}_{B_1, \dots, B_s \sim D} [2 \sum_{i=1}^s [L(B_i) - \sum_{j=1}^{i-1} \alpha \frac{\partial L(B_i)}{\partial \theta} \cdot \frac{\partial L(B_j)}{\partial \theta}]]$$

Optimize over the proposed loss (MER):

$$\theta = \arg \min_{\theta} \mathbb{E}_{[(x_{11}, y_{11}), \dots, (x_{sk}, y_{sk})] \sim M} [2 \sum_{i=1}^s \sum_{j=1}^k [L(x_{ij}, y_{ij}) - \sum_{q=1}^{i-1} \sum_{r=1}^{j-1} \alpha \frac{\partial L(x_{ij}, y_{ij})}{\partial \theta} \cdot \frac{\partial L(x_{qr}, y_{qr})}{\partial \theta}]]$$

# MER

1. Prioritizing the current example

2. Combining ER with optimization based meta-learning (Reptile)

3. Reservoir Sampling

---

## Algorithm 1 Meta-Experience Replay (MER)

---

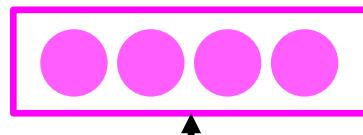
```
procedure TRAIN( $D, \theta, \alpha, \beta, \gamma, s, k$ )
     $M \leftarrow \{\}$ 
    for  $t = 1, \dots, T$  do
        for  $(x, y)$  in  $D_t$  do
            // Draw batches from buffer:
             $B_1, \dots, B_s \leftarrow sample(x, y, s, k, M)$ 
             $\theta_0^A \leftarrow \theta$ 
            for  $i = 1, \dots, s$  do
                 $\theta_{i,0}^W \leftarrow \theta$ 
                for  $j = 1, \dots, k$  do
                     $x_c, y_c \leftarrow B_i[j]$ 
                     $\theta_{i,j}^W \leftarrow SGD(x_c, y_c, \theta_{i,j-1}^W, \alpha)$ 
                end for
                // Within batch Reptile meta-update:
                 $\theta \leftarrow \theta_{i,0}^W + \beta(\theta_{i,k}^W - \theta_{i,0}^W)$ 
                 $\theta_i^A \leftarrow \theta$ 
            end for
            // Across batch Reptile meta-update:
             $\theta \leftarrow \theta_0^A + \gamma(\theta_s^A - \theta_0^A)$ 
            // Reservoir sampling memory update:
             $M \leftarrow M \cup \{(x, y)\}$  (algorithm 3)
        end for
    end for
    return  $\theta, M$ 
end procedure
```

---

# Experiment Setting

**Single-headed**

$$P(y|x)$$



**Hidden  
Layer**

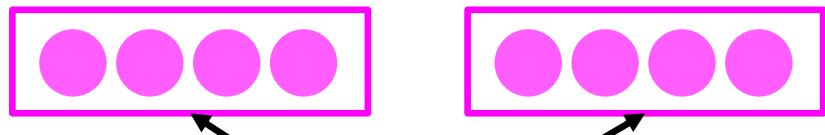
**Input x**



**Task 1   Task 2**

**Multi-headed**

$$P(y|x)$$



**Hidden  
Layer**

**Input x**



**Task 1   Task 2**

# Result

How does MER perform on supervised CL benchmarks?

Model	MNIST Rotations			MNIST Permutations		
	RA	LA	BTI	RA	LA	BTI
Online	<b>53.38</b> $\pm$ 1.53	<b>58.82</b> $\pm$ 2.17	<b>-5.44</b> $\pm$ 1.70	<b>55.42</b> $\pm$ 0.65	<b>69.18</b> $\pm$ 0.99	<b>-13.76</b> $\pm$ 1.19
Independent	<b>60.74</b> $\pm$ 4.55	<b>60.74</b> $\pm$ 4.55	-	<b>55.80</b> $\pm$ 4.79	<b>55.80</b> $\pm$ 4.79	-
Task Input	<b>79.98</b> $\pm$ 0.66	<b>81.04</b> $\pm$ 0.34	<b>-1.06</b> $\pm$ 0.59	<b>80.46</b> $\pm$ 0.80	<b>81.20</b> $\pm$ 0.52	<b>-0.74</b> $\pm$ 1.10
EWC	<b>57.96</b> $\pm$ 1.33	<b>78.38</b> $\pm$ 0.72	<b>-20.42</b> $\pm$ 1.60	<b>62.32</b> $\pm$ 1.34	<b>75.64</b> $\pm$ 1.18	<b>-13.32</b> $\pm$ 2.24
GEM	<b>87.58</b> $\pm$ 0.32	<b>86.46</b> $\pm$ 0.46	<b>+1.12</b> $\pm$ 0.35	<b>83.02</b> $\pm$ 0.23	<b>80.46</b> $\pm$ 0.38	<b>+2.56</b> $\pm$ 0.42
MER	<b>89.56</b> $\pm$ 0.11	<b>87.62</b> $\pm$ 0.16	<b>+1.94</b> $\pm$ 0.18	<b>85.50</b> $\pm$ 0.16	<b>86.36</b> $\pm$ 0.21	<b>-0.86</b> $\pm$ 0.21

- RA: Retained Accuracy
- LA: Learning Accuracy
- BTI: Backward Transfer and Interference

# Result

---

How do the performance gains from MER vary as a function of the buffer size?

Model	Buffer	MNIST Rotations			MNIST Permutations		
		RA	LA	BTI	RA	LA	BTI
GEM	5120	<b>87.58</b> $\pm$ 0.32	<b>86.46</b> $\pm$ 0.46	<b>+1.12</b> $\pm$ 0.35	<b>83.02</b> $\pm$ 0.23	<b>80.46</b> $\pm$ 0.38	<b>+2.56</b> $\pm$ 0.42
	500	<b>74.88</b> $\pm$ 0.93	<b>85.90</b> $\pm$ 0.53	<b>-11.02</b> $\pm$ 1.22	<b>69.26</b> $\pm$ 0.66	<b>80.28</b> $\pm$ 0.52	<b>-11.02</b> $\pm$ 0.71
	200	<b>67.38</b> $\pm$ 1.75	<b>85.40</b> $\pm$ 0.53	<b>-18.02</b> $\pm$ 1.99	<b>55.42</b> $\pm$ 1.10	<b>79.84</b> $\pm$ 1.01	<b>-24.42</b> $\pm$ 1.10
MER	5120	<b>89.56</b> $\pm$ 0.11	<b>87.62</b> $\pm$ 0.16	<b>+1.94</b> $\pm$ 0.18	<b>85.50</b> $\pm$ 0.16	<b>86.36</b> $\pm$ 0.21	<b>-0.86</b> $\pm$ 0.21
	500	<b>82.08</b> $\pm$ 0.31	<b>85.66</b> $\pm$ 0.20	<b>-3.58</b> $\pm$ 0.39	<b>77.50</b> $\pm$ 0.46	<b>83.90</b> $\pm$ 0.23	<b>-6.40</b> $\pm$ 0.35
	200	<b>77.42</b> $\pm$ 0.78	<b>83.02</b> $\pm$ 0.23	<b>-5.60</b> $\pm$ 0.70	<b>73.46</b> $\pm$ 0.45	<b>84.42</b> $\pm$ 0.20	<b>-9.96</b> $\pm$ 0.45

- The benefits of MER seem to grow as the buffer becomes smaller.

# Result

---

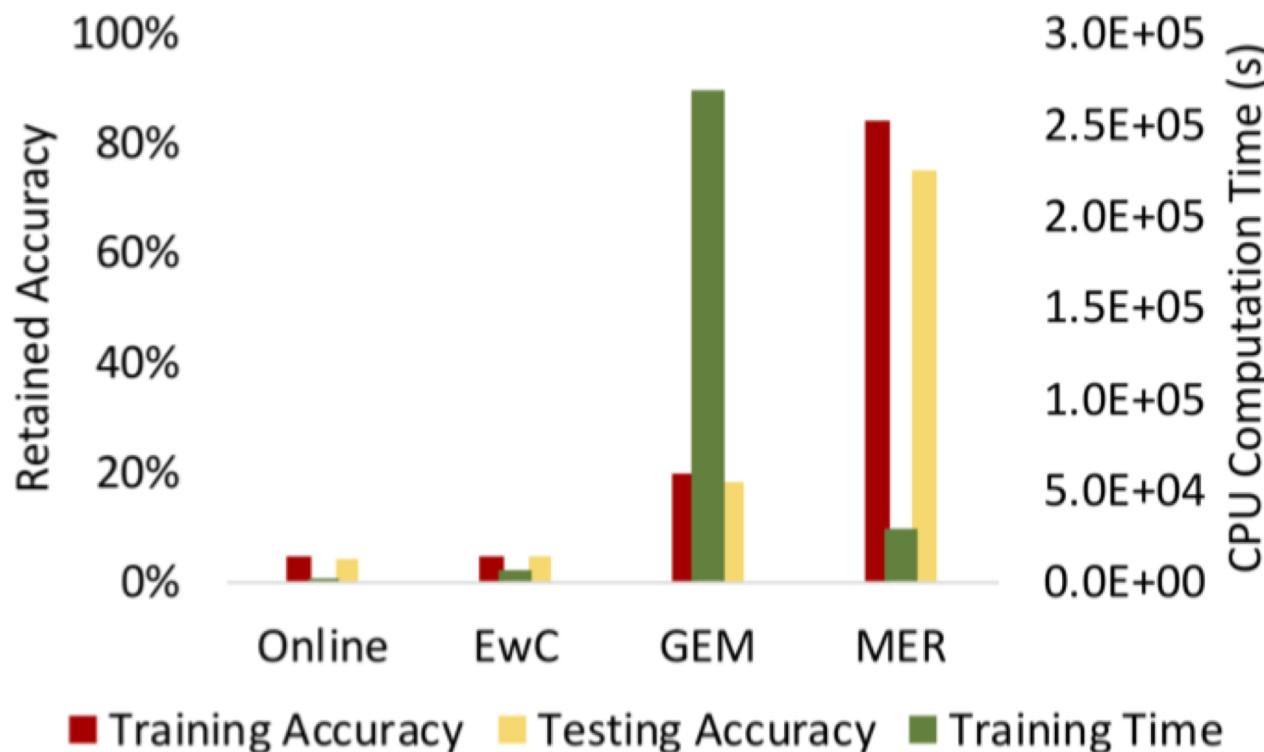
How effective is MER at dealing with increasingly non-stationary settings?

Model	Buffer	Many Permutations			Omniglot		
		RA	LA	BTI	RA	LA	BTI
Online	0	32.62 ± 0.43	51.68 ± 0.65	-19.06 ± 0.86	4.36 ± 0.37	5.38 ± 0.18	-1.02 ± 0.33
EWC	0	33.46 ± 0.46	51.30 ± 0.81	-17.84 ± 1.15	4.63 ± 0.14	9.43 ± 0.63	-4.80 ± 0.68
GEM	5120	56.76 ± 0.29	<b>59.66</b> ± 0.46	-2.92 ± 0.52	18.03 ± 0.15	3.86 ± 0.09	<b>+14.19</b> ± 0.19
	500	32.14 ± 0.50	55.66 ± 0.53	-23.52 ± 0.87	-	-	-
MER	5120	<b>62.52</b> ± 0.32	<b>59.44</b> ± 0.23	<b>+3.08</b> ± 0.31	<b>75.23</b> ± 0.52	<b>69.12</b> ± 0.83	<b>+6.11</b> ± 0.62
	500	<b>47.40</b> ± 0.35	<b>65.18</b> ± 0.20	<b>-17.78</b> ± 0.39	<b>32.05</b> ± 0.69	<b>28.78</b> ± 0.91	<b>+3.27</b> ± 1.04

- Many Permutations is a variant of MNIST Permutations that has 5 times more tasks (100 total) and 5 times less training examples per task (200 each)

# Result

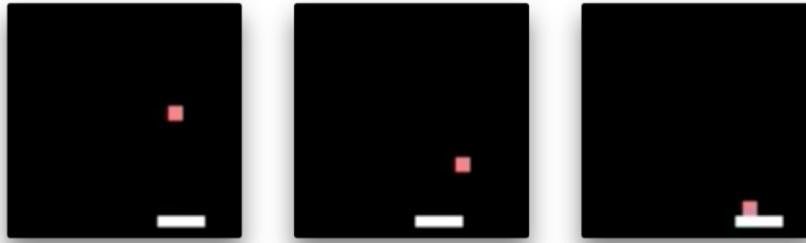
How effective is MER at dealing with increasingly non-stationary settings?



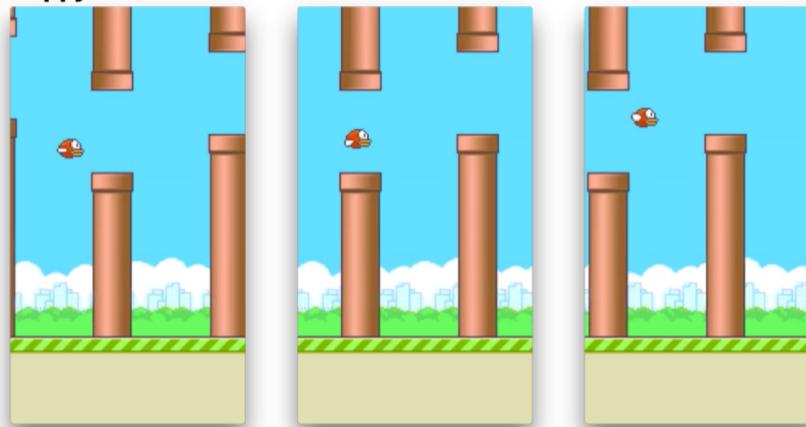
# Result

Can MER improve a DQN with ER in CRL settings?

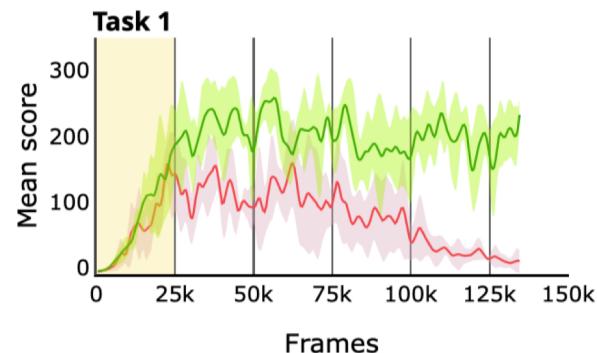
Catcher



Flappy Bird

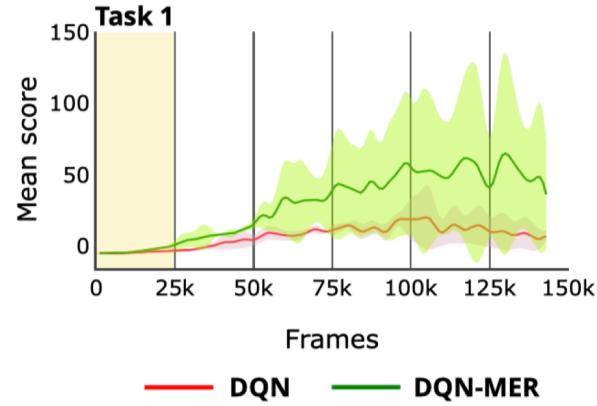


Task 1



Frames

Task 1

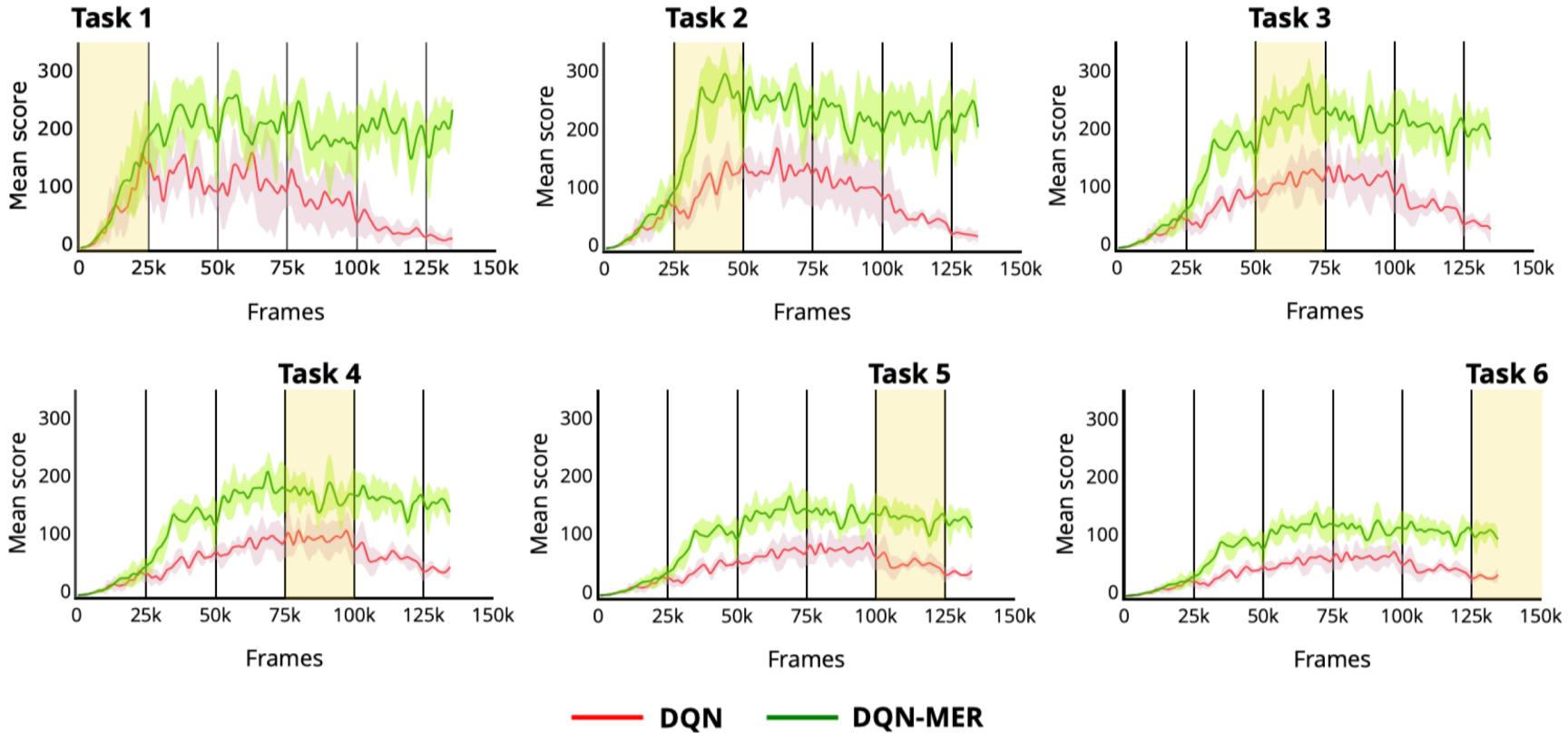


Frames

— DQN — DQN-MER

# Result

Can MER improve a DQN with ER in CRL settings?



# Result

---

Does MER lead to a shift in the distribution of gradient dot products?

Model	MNIST Permutations	MNIST Rotations	Many Permutations
ER	-0.569 $\pm$ 0.077	-1.652 $\pm$ 0.082	-1.280 $\pm$ 0.078
MER	+0.042 $\pm$ 0.017	+0.017 $\pm$ 0.007	+0.131 $\pm$ 0.027

- It verifies that a very significant and reproducible difference in the mean gradient encountered is seen for MER in comparison to ER alone

# Conclusion

---

This work cast a new perspective on CL problem: trade-off btw transfer and interference.

They also propose a new algorithm MER that

- regularizes the objective of experience replay so that gradients on incoming examples are more likely to have transfer and less likely to have interference with respect to past examples

The experimental results show that it outperforms baselines in non-stationary supervised and reinforcement learning setting

# Future Work

---

Combine MER with mixtures of experts model to amplify the ability to control weight sharing dynamics

Combine MER with option frameworks allowing for orthogonalization of weights relating to different skills

# References

---

[MIR19] Matthew Riemer et al., “Learning to Learn without Forgetting by Maximizing Transfer and Minimizing Interference”, ICLR 2019

# Question or Discussion

---

