

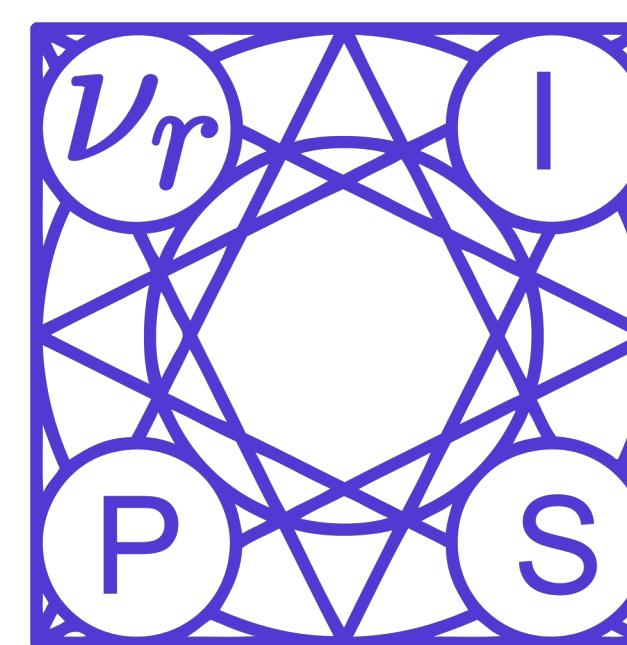
Uncertainty-based Continual Learning with Adaptive Regularization

Hongjoon Ahn^{*1}, Sungmin Cha^{*2}, Donggyu Lee², and Taesup Moon^{1,2}

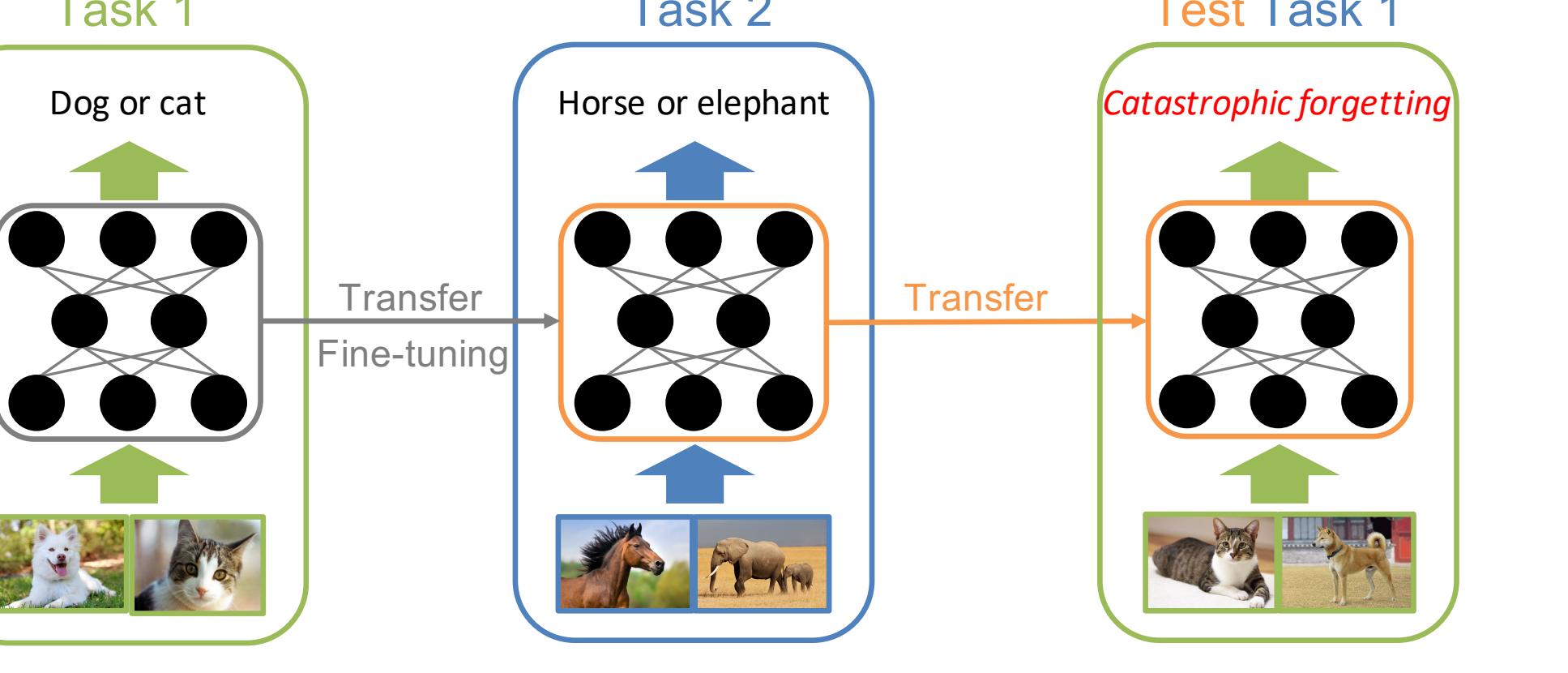
Department of Artificial Intelligence¹, Department of Electrical and Computer Engineering²

Sungkyunkwan University

{hong0805, csm9493, ldk308, tsmoon}@skku.edu



Catastrophic Forgetting



- Algorithms for overcoming Catastrophic forgetting**
 - Dynamic network architecture-based (ex. PNN, DEN)
 - Dual memory system-based (Ex. GEM, iCaRL)
 - Regularization-based** (Ex. EWC, LWF, SI)
- Contributions of our work**
 - Define a notion of “Uncertainty” for each hidden node
 - Much smaller number of additional parameters
 - Add two additional regularization terms
 - Freezing** the weights that are identified to be important
 - The **actively** learning parameters for new tasks

Bayesian Online Learning

- Notation and a review of Bayesian online learning**
 - In **Bayesian online learning**, they apply this standard variational inference method to the continual learning
- $\mathcal{F}(D_t, \theta_t) = \mathbb{E}_{q(\mathcal{W}|\theta_t)}[-\log p(D_t|\mathcal{W})] + D_{KL}(q(\mathcal{W}|\theta_t)||q(\mathcal{W}|\theta_{t-1}))$
- Interpreting KL-divergence and motivation of UCL**

$$\frac{1}{2} \sum_{l=1}^L \left[\underbrace{\left\| \frac{\mu_t^{(l)} - \mu_{t-1}^{(l)}}{\sigma_{t-1}^{(l)}} \right\|_2^2}_{(a)} + \underbrace{1^\top \left\{ \left(\frac{\sigma_t^{(l)}}{\sigma_{t-1}^{(l)}} \right)^2 - \log \left(\frac{\sigma_t^{(l)}}{\sigma_{t-1}^{(l)}} \right)^2 \right\}}_{(b)} \right]$$
 - $\theta_t^{(l)} = (\mu_t^{(l)}, \sigma_t^{(l)})$: the mean and standard deviation of the weight matrix for layer l at task t
 - Term (b) is minimized when $\sigma_t^{(l)} = \sigma_{t-1}^{(l)}$

Uncertainty-based Continual Learning(UCL)

- Contribution 1: A notion of uncertainty for a node**
 - We devise a notion of uncertainty for each **node** of the network
- Contribution 2: Modifying term (a)**
 - Regularization for the important connection**

$$\frac{1}{2} \left(\sum_{l=1}^L \left\| \Lambda^{(l)} \odot (\mu_t^{(l)} - \mu_{t-1}^{(l)}) \right\|_2^2 \right), \Lambda_{ij}^{(l)} \triangleq \max\left(\frac{\sigma_{\text{init}}^{(l)}}{\sigma_{t-1,i}^{(l)}}, \frac{\sigma_{\text{init}}^{(l-1)}}{\sigma_{t-1,j}^{(l-1)}}\right) \quad (4)$$
 - $\sigma_{\text{init}}^{(l)}$: Pivot value for defining regularization strength
 - Constrain the incoming weights to the node to have the same variance parameters
 - Then, set that variance as the uncertainty of the node
 - Regularization for the important weight**

$$\sum_{l=1}^L (\sigma_{\text{init}}^{(l)})^2 \left\| \left(\frac{\mu_{t-1}^{(l)}}{\sigma_{t-1}^{(l)}} \right)^2 \odot (\mu_t^{(l)} - \mu_{t-1}^{(l)}) \right\|_1 \quad (5)$$
 - A weight is important if the ratio μ/σ is high
 - The ℓ_1 -norm will promote sparsity and $\mu_{t,ij}^{(l)}$ will tend to **freeze** to $\mu_{t-1,ij}^{(l)}$
- Contribution 3: Modifying term (b)**
 - Regularization for the actively learning**

$$\frac{1}{2} 1^\top \left((\sigma_t^{(l)})^2 - \log(\sigma_t^{(l)})^2 \right) \quad (6)$$
 - which forces $\sigma_t^{(l)}$ to get close to $\sqrt{2}\sigma_{t-1}^{(l)}$ when minimized with term (b)
 - Increase the number of “**actively**” learning nodes
 - Result in **gracefully forgetting** the past tasks

- Final loss function for UCL**

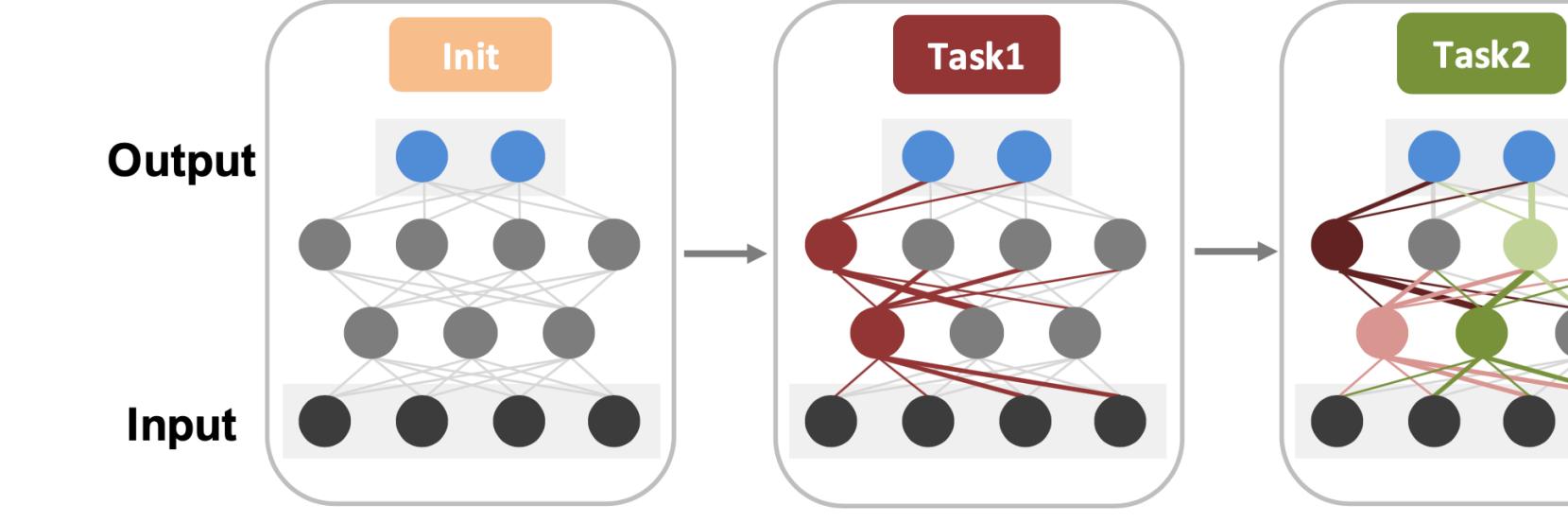
$$-\log p(D_t|\mathcal{W}) + \sum_{l=1}^L \left[\left(\frac{1}{2} \left\| \Lambda^{(l)} \odot (\mu_t^{(l)} - \mu_{t-1}^{(l)}) \right\|_2^2 + (\sigma_{\text{init}}^{(l)})^2 \left\| \left(\frac{\mu_{t-1}^{(l)}}{\sigma_{t-1}^{(l)}} \right)^2 \odot (\mu_t^{(l)} - \mu_{t-1}^{(l)}) \right\|_1 \right) + \frac{\beta}{2} 1^\top \left\{ \left(\frac{\sigma_t^{(l)}}{\sigma_{t-1}^{(l)}} \right)^2 - \log \left(\frac{\sigma_t^{(l)}}{\sigma_{t-1}^{(l)}} \right)^2 + (\sigma_t^{(l)})^2 - \log(\sigma_t^{(l)})^2 \right\} \right], \quad (7)$$
 - β : the hyperparameter controls the increasing or decreasing speed of $\sigma_t^{(l)}$
 - The number of sampling is 1 for each iteration
- Illustration of the regularization mechanism**


Figure 2: Colored hidden nodes and edges denote important nodes and highly regularized weights due to (4), respectively. The width of colored edge denotes the regularization strength of (5). Note as new task comes the uncertainty level of a node can vary due to (6), represented with color changes.

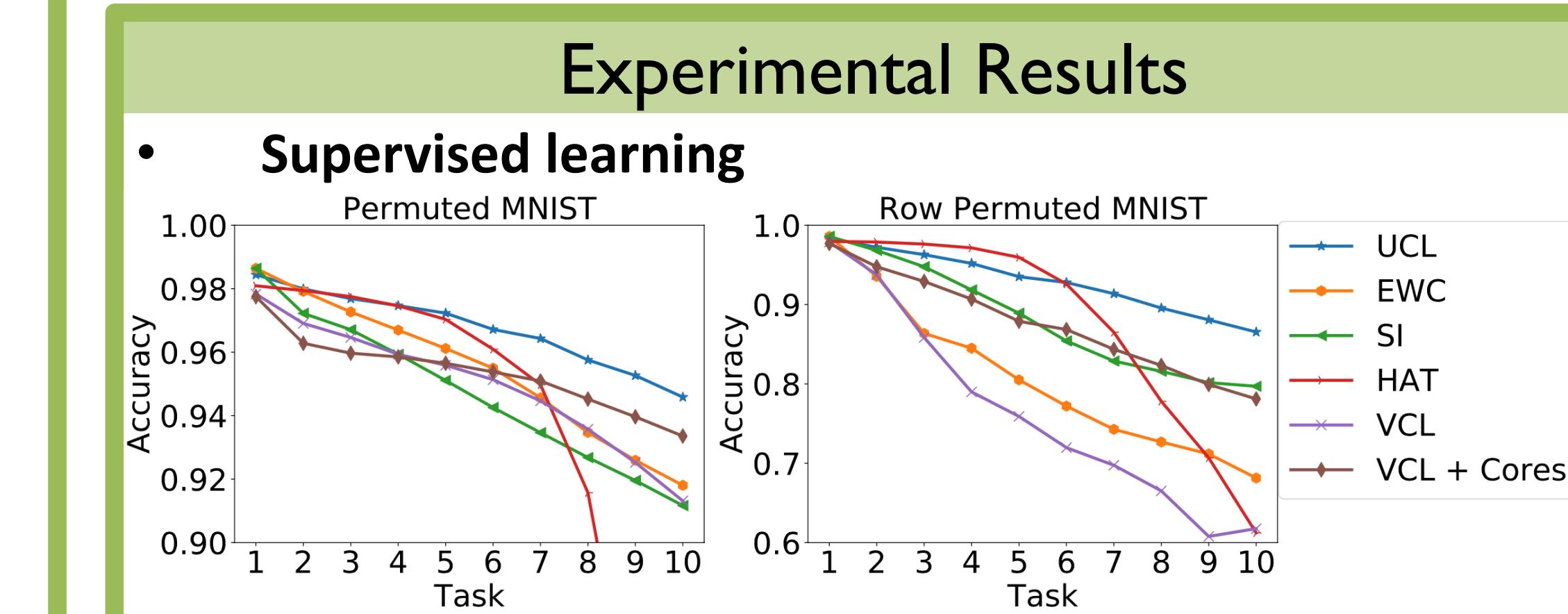


Figure 3: Experimental results on various Permutated MNIST with single head.

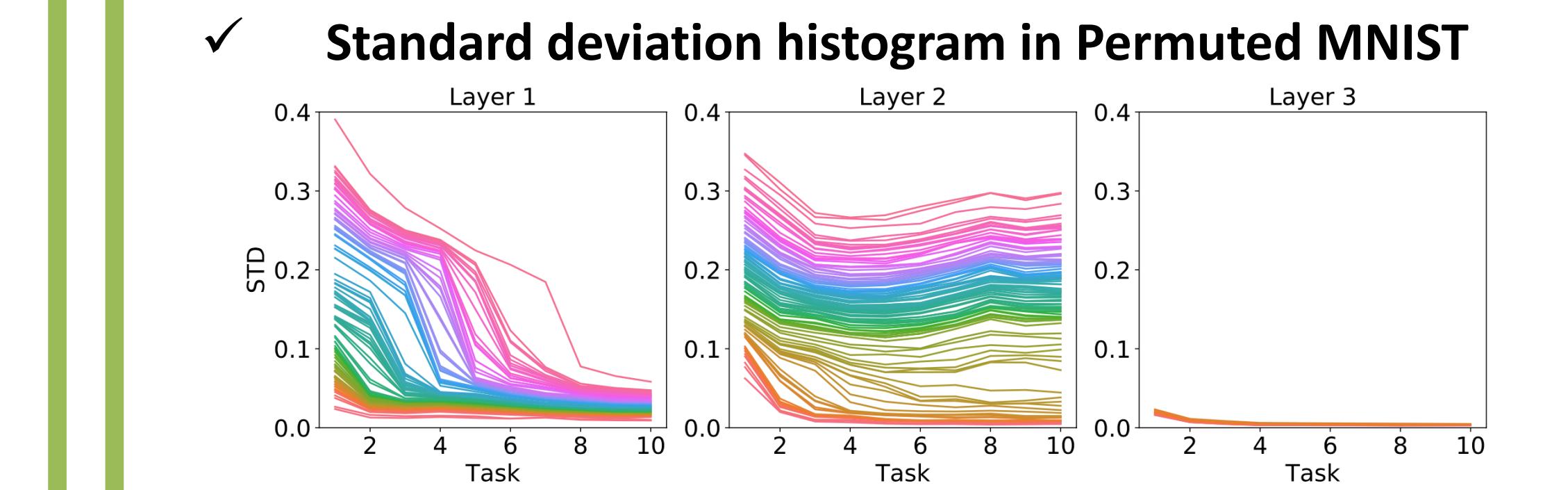


Figure 4: Standard deviation histogram in Permutated MNIST experiment. We randomly selected 100 standard deviations for layer 1 and 2. In layer 3, all 10 nodes are shown.

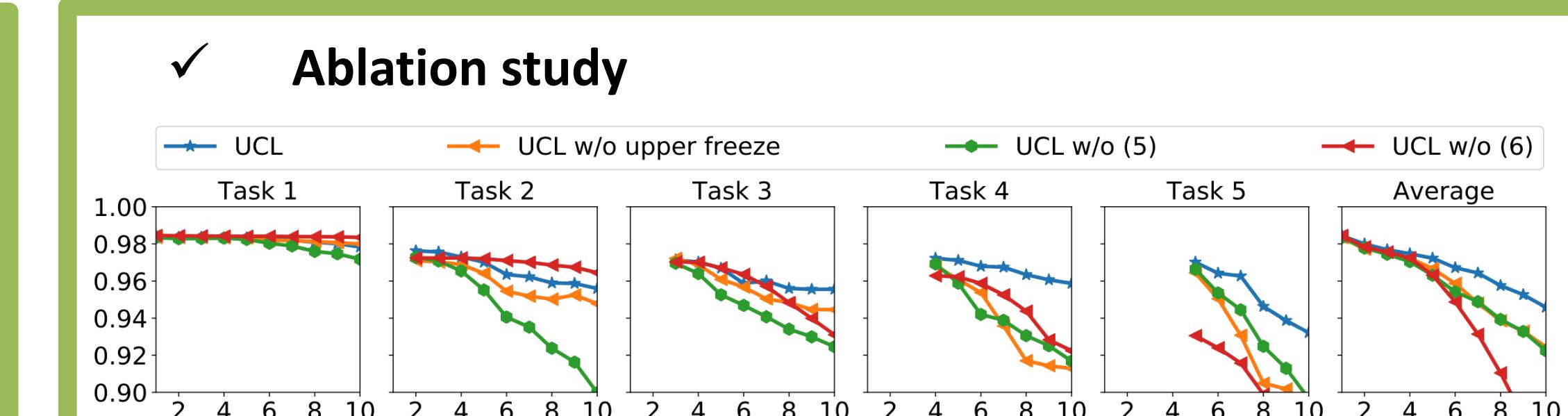


Figure 6: Ablation study in Permutated MNIST. Each line denote the test accuracy.

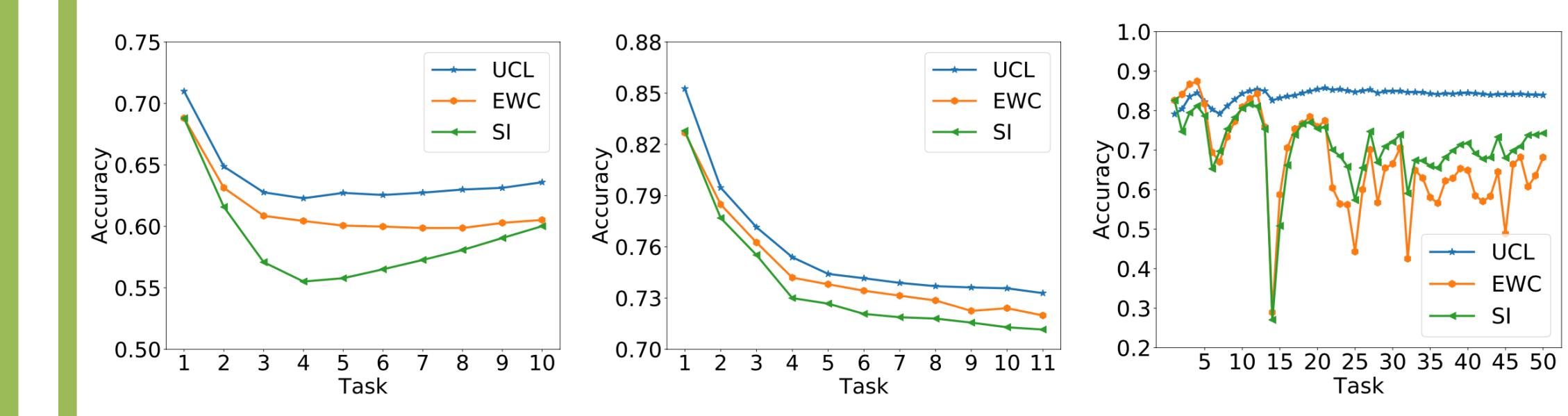


Figure 7: Experiments on supervised learning using convolutional neural network

The number of parameters

Table 1: The number of parameters used for each benchmark.

Dataset	Methods	Vanilla	UCL	EWC	SI	HAT	VCL
Permuted MNIST	478K	960K	1435K	1435K	486K	1914K	
Split MNIST	270K	538K	808K	808K	272K	1077K	
Split notMNIST	187K	375K	559K	559K	190K	749K	
Split CIFAR10/100	839K	1655K	2467K	2467K	-	-	
Omniglot	1773K	1884K	1995K	1995K	-	-	

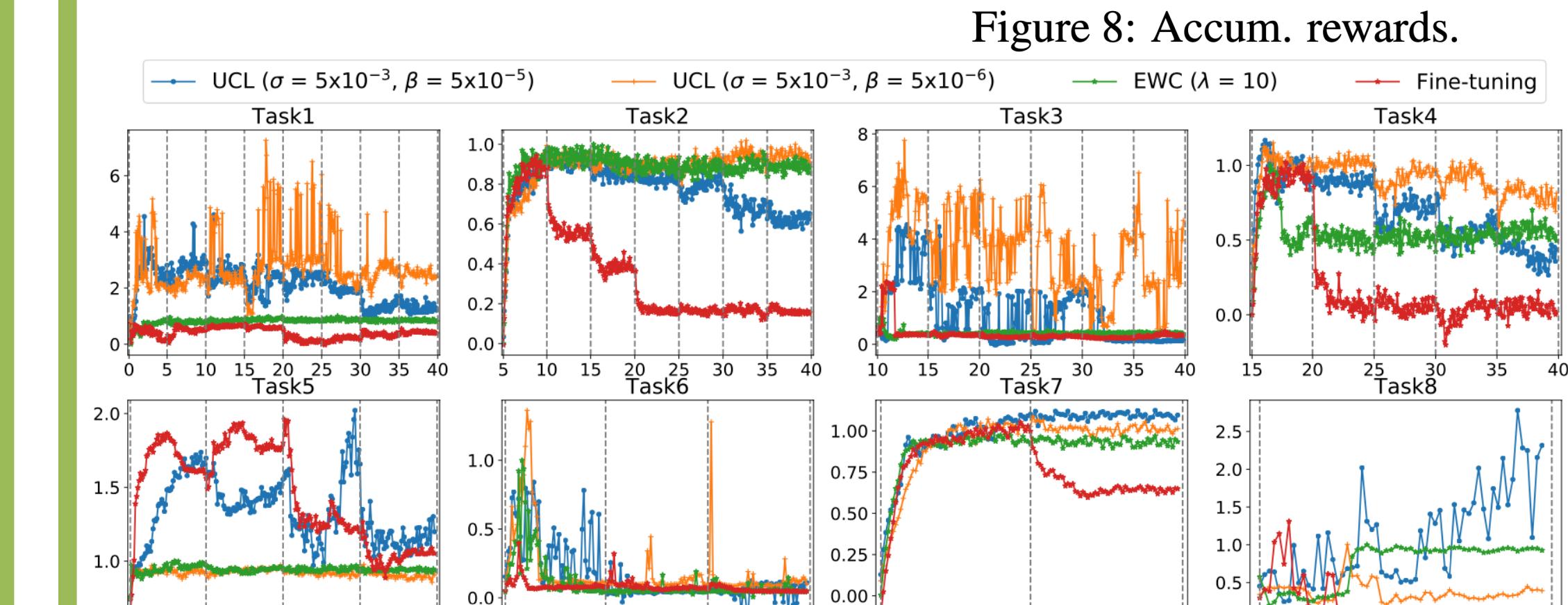
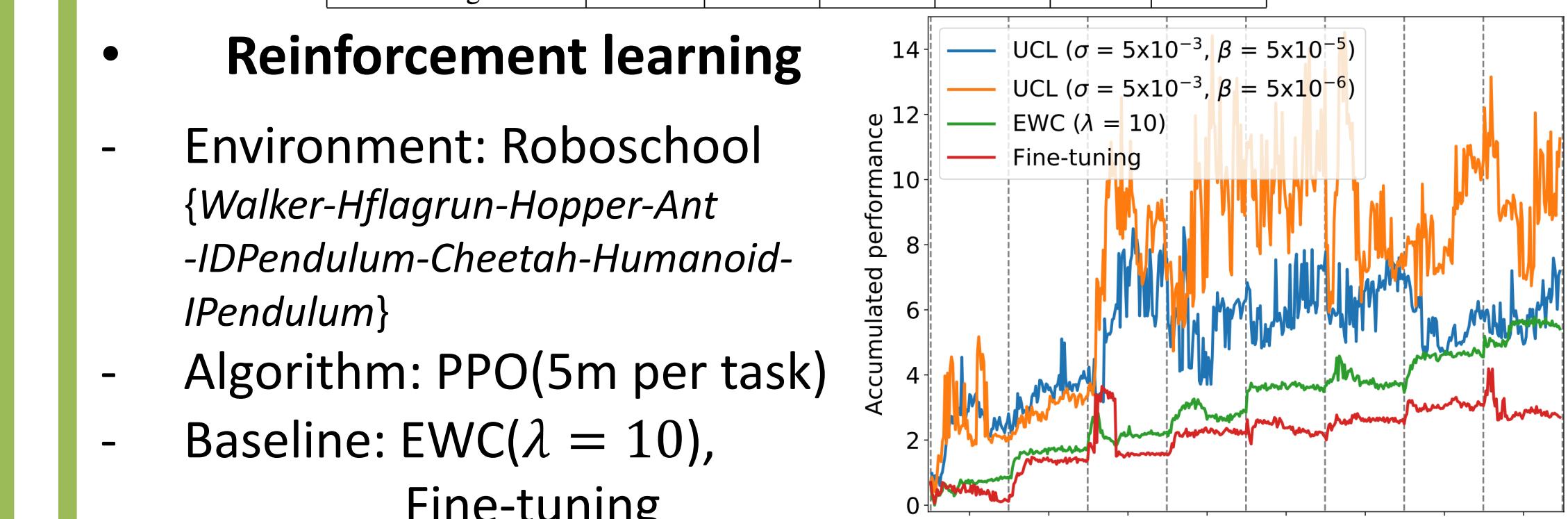


Figure 8: Accum. rewards.