

# Overcoming catastrophic forgetting in neural networks

James Kirkpatrick<sup>a,1</sup>, Razvan Pascanu<sup>a</sup>, Neil Rabinowitz<sup>a</sup>, Joel Veness<sup>a</sup>, Guillaume Desjardins<sup>a</sup>, Andrei A. Rusu<sup>a</sup>, Kieran Milan<sup>a</sup>, John Quan<sup>a</sup>, Tiago Ramalho<sup>a</sup>, Agnieszka Grabska-Barwinska<sup>a</sup>, Demis Hassabis<sup>a</sup>, Claudia Clopath<sup>b</sup>, Dharshan Kumaran<sup>a</sup>, and Raia Hadsell<sup>a</sup>

<sup>a</sup>DeepMind, London EC4 5TW, United Kingdom; and <sup>b</sup>Bioengineering Department, Imperial College London, London SW7 2AZ, United Kingdom

PNAS, Proceedings of the National Academy of Sciences  
citation: 4

2017/04/28

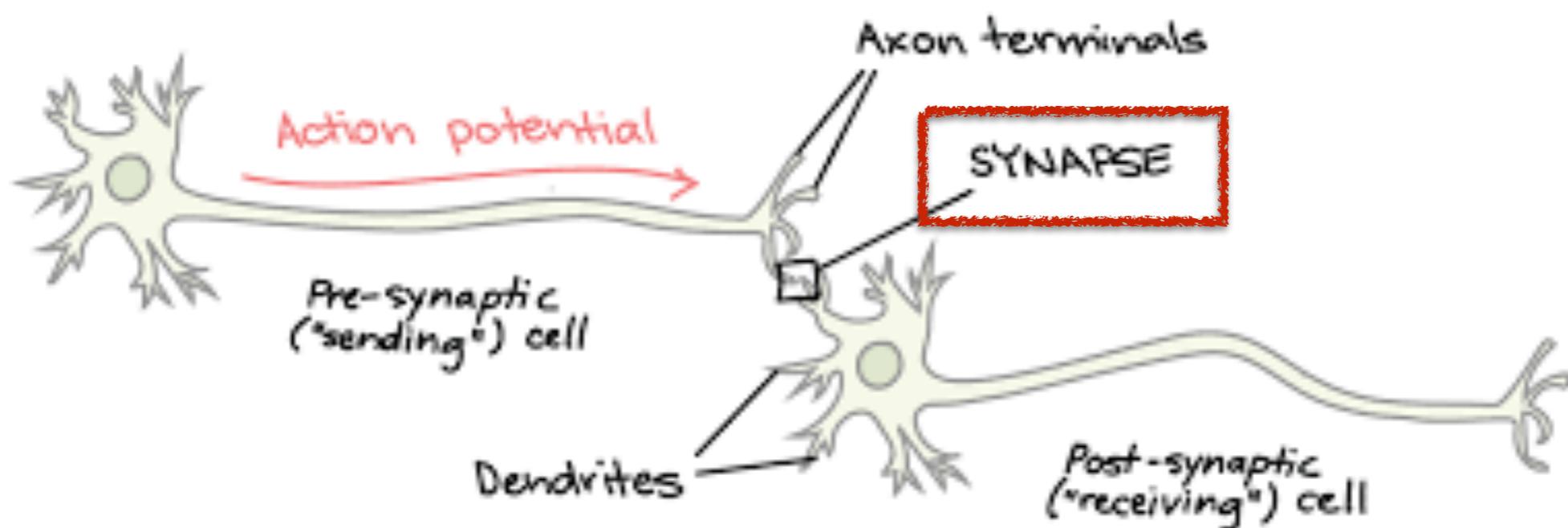
Katy@Datalab

# Background

- **Catastrophic forgetting** is forgetting key information needed to solve a previous task when training on a new task.

# The kind of consolidation that occur in the brain

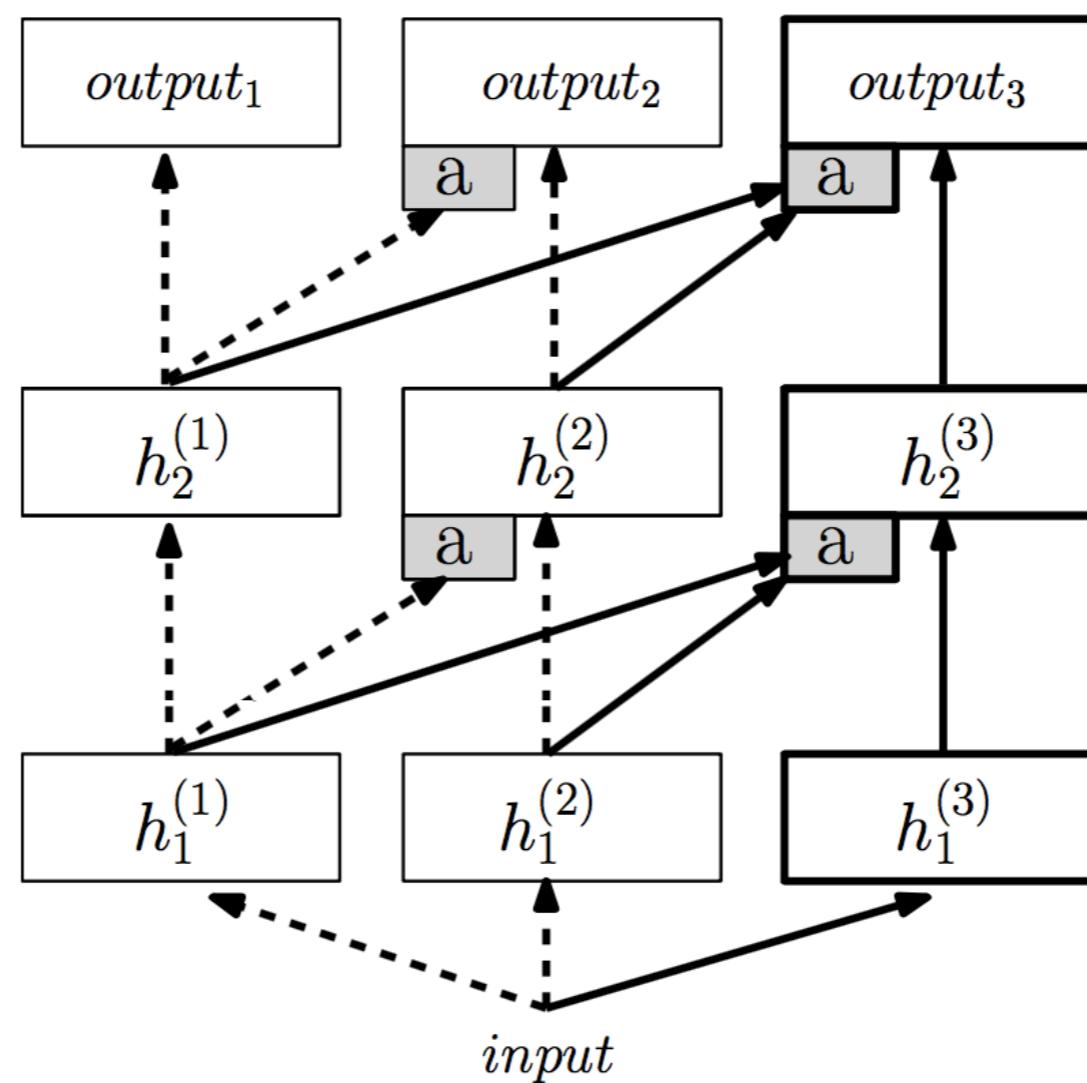
- Continual learning in the neocortex relies on **task-specific synaptic consolidation**
- **Synaptic consolidation:** connections between neurons are less likely to be overwritten if they have been important in previously learnt tasks.



- knowledge is durably encoded by rendering a proportion of synapses less plastic and therefore stable over long timescales.

# Related Work

- Ensemble of DNNs: Rusu, Andrei A., et al.  
"Progressive neural networks." arXiv preprint arXiv: 1606.04671 (2016).



- Fernando, Chrisantha, et al. "**Pathnet**: Evolution channels gradient descent in super neural networks." arXiv preprint arXiv:1701.08734 (2017).

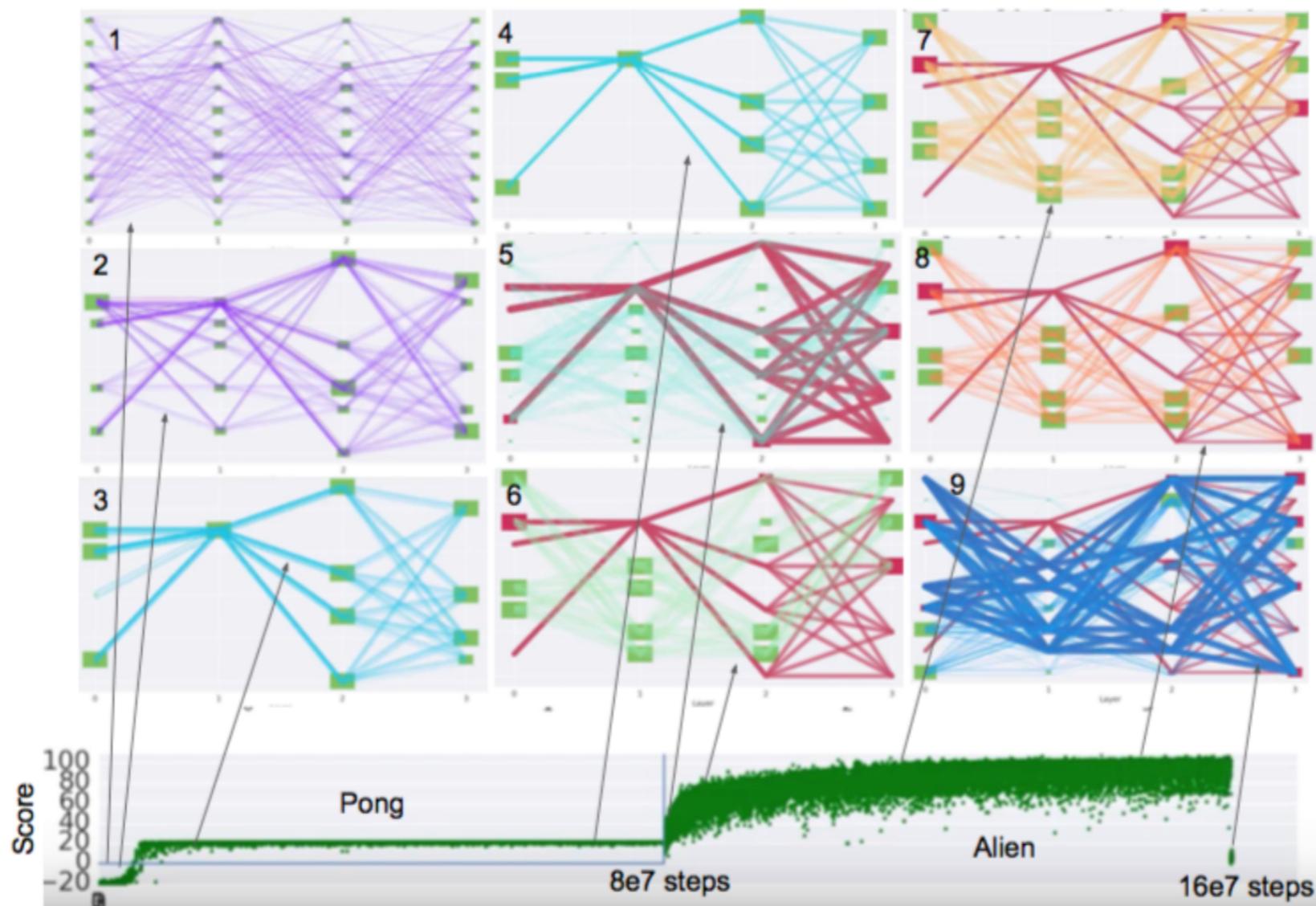


Figure 1: A population of randomly initialized pathways (purple lines in Box 1) are evolved whilst learning task A, Pong. At the end of training, the best pathway is fixed (dark red lines in Box 5) and a new population of paths are generated (light blue lines in Box 5) for task B. This population is then trained on Alien and the optimal pathway that is evolved on Alien is subsequently fixed at the end of training, shown as dark blue lines in Box 9.

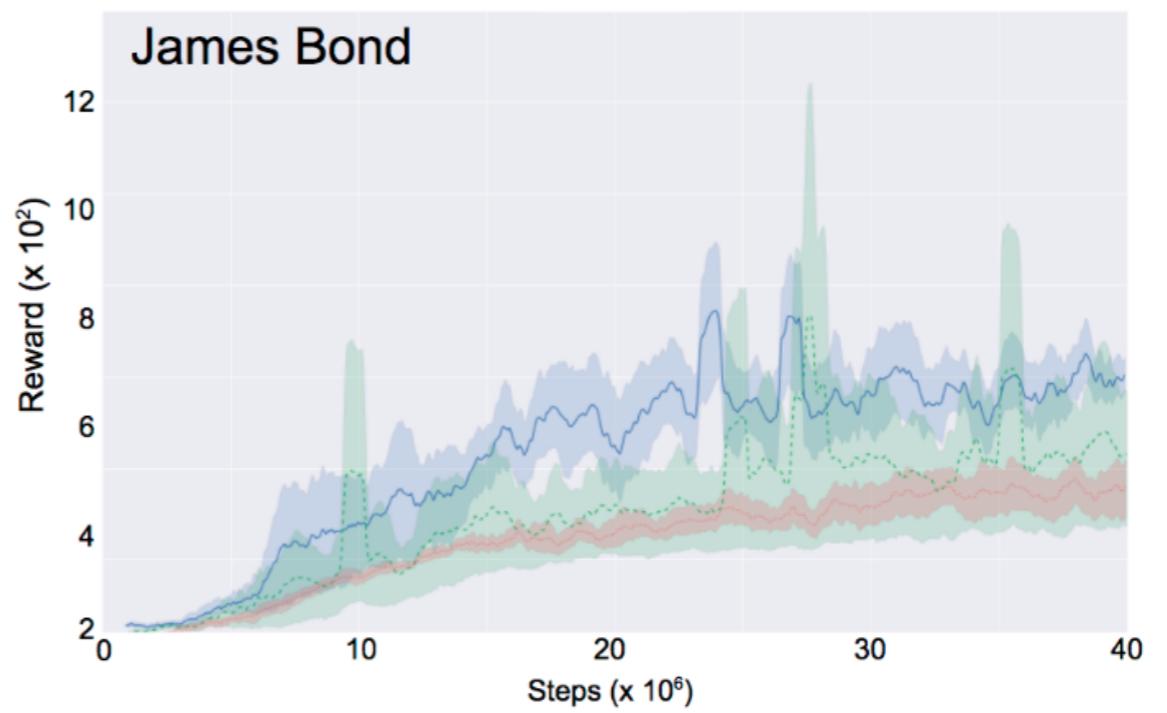
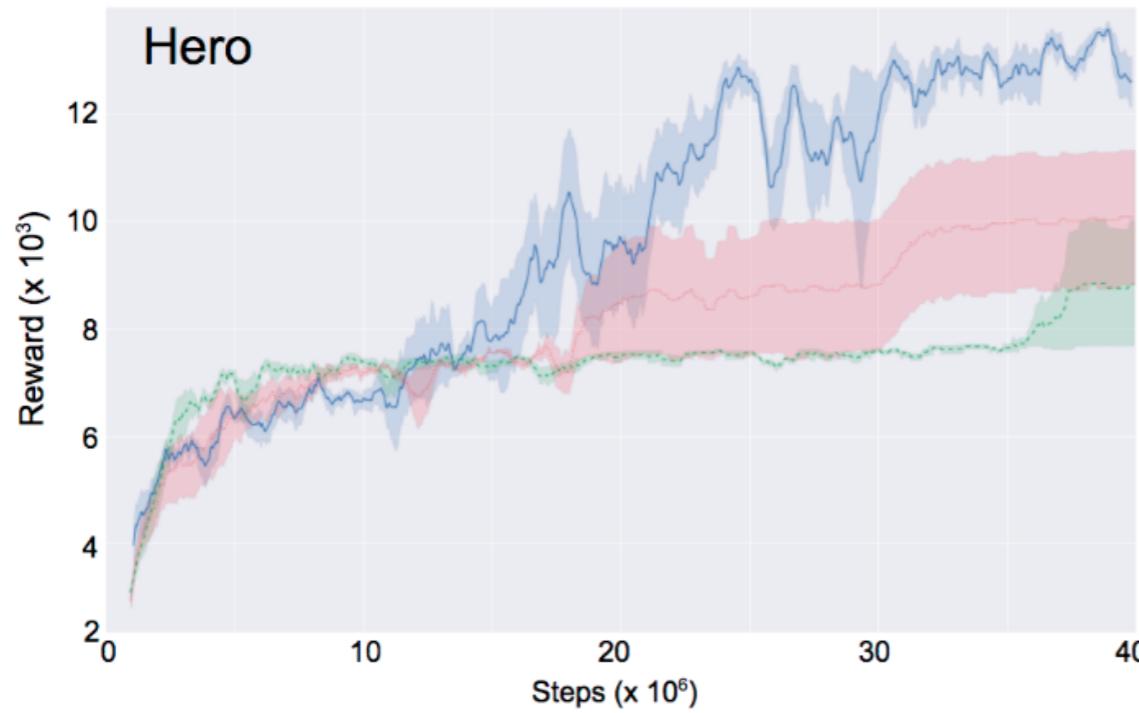
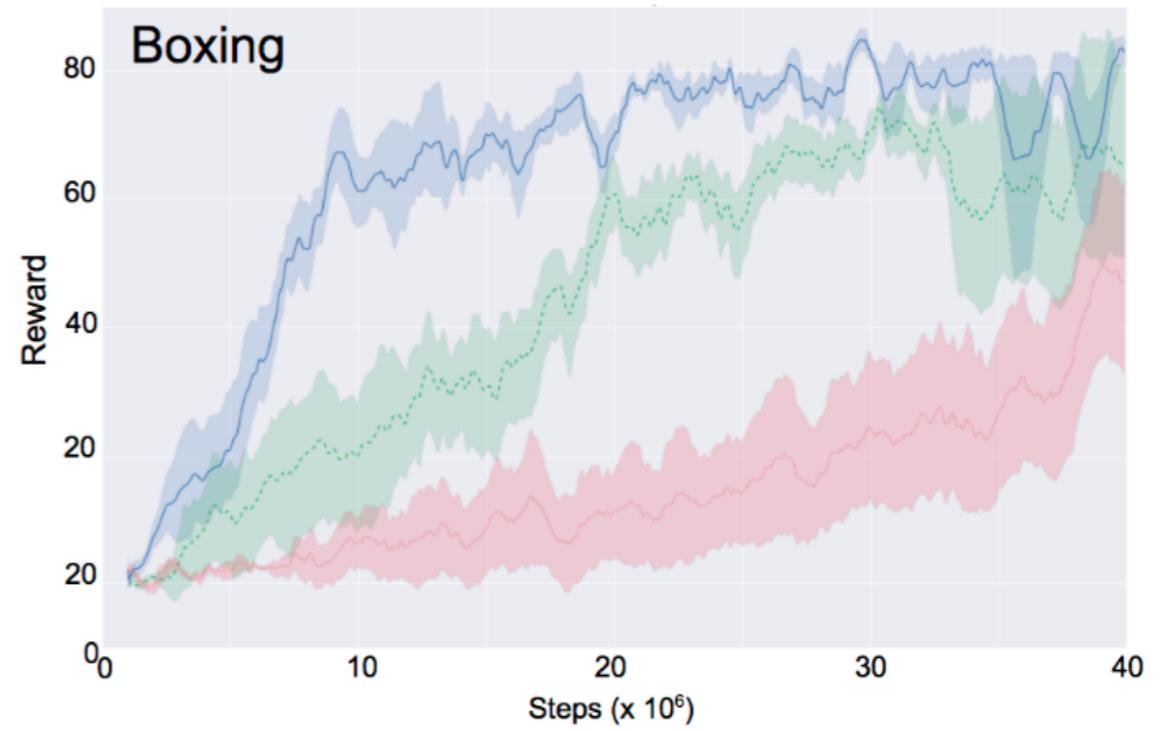
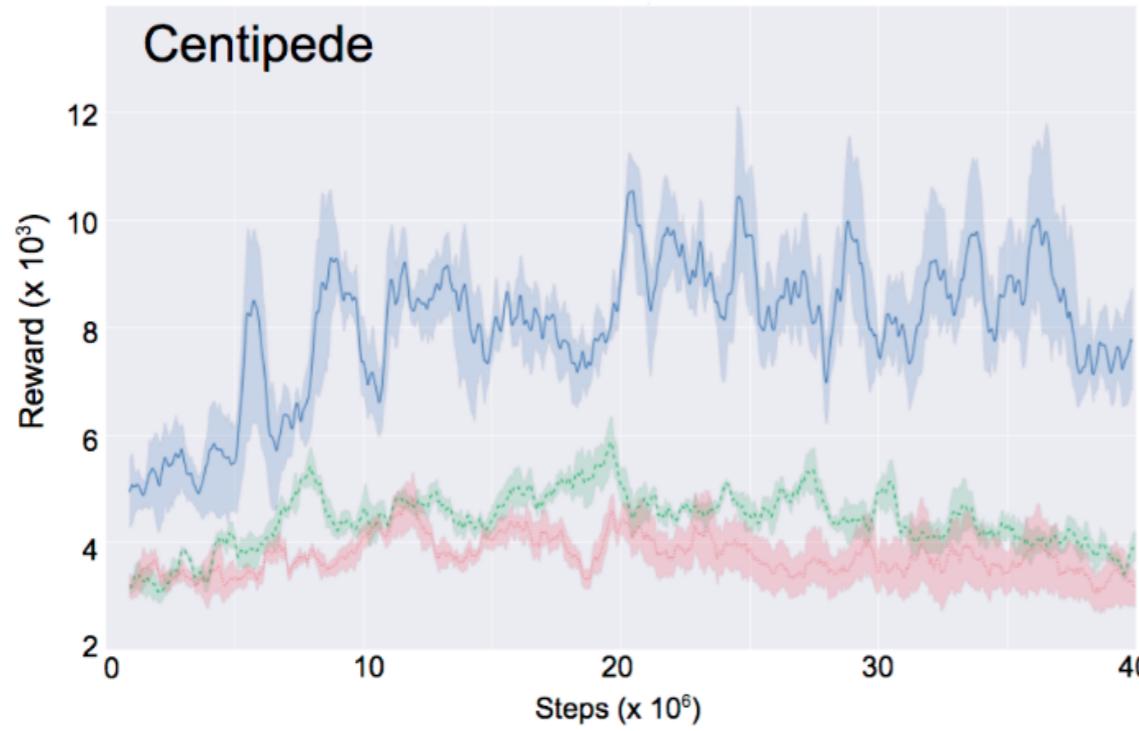
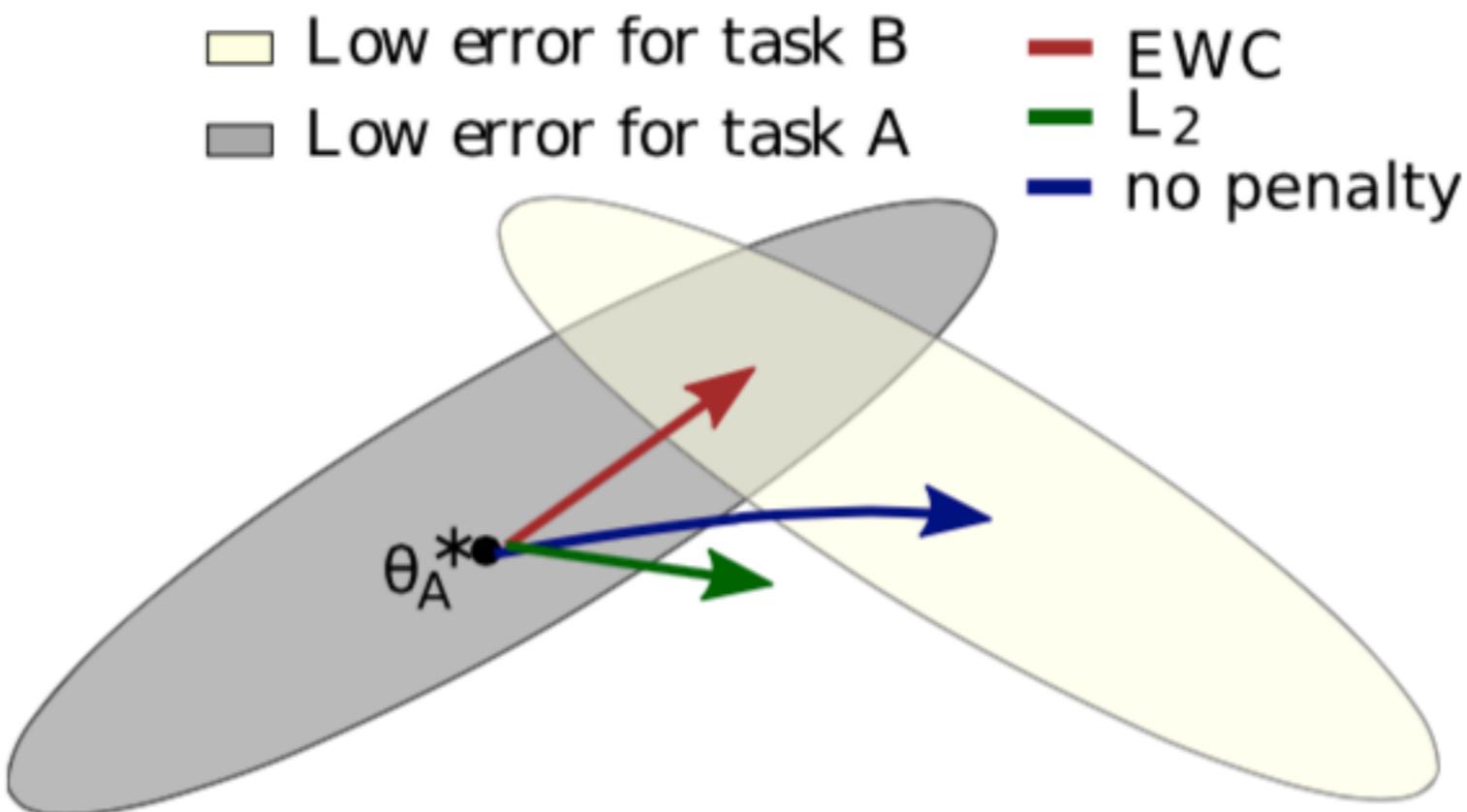


Figure 6: PathNet shows positive transfer from River Raid to Centipede, Boxing, Hero and James Bond. The graphs show the reward over the first 40 million steps of training. Blue shows the results from the best five hyperparameter settings of PathNet out of 243. Compare these to the best five hyperparameter setting runs for independent learning (red) and fine-tuning (green) controls out of 45.

# Intuition

- They implement an algorithm that performs a similar operation in artificial neural networks by **constraining important parameters to stay close to their old values.**

# ELASTIC WEIGHT CONSOLIDATION (**EWC**)



- $\theta_A^*$  refers to the configuration of  $\theta$  that performs well at A
- slow down the learning for the weights that were important to previous task(s)

# Key Question

- How to determine which weights are most important to task A?

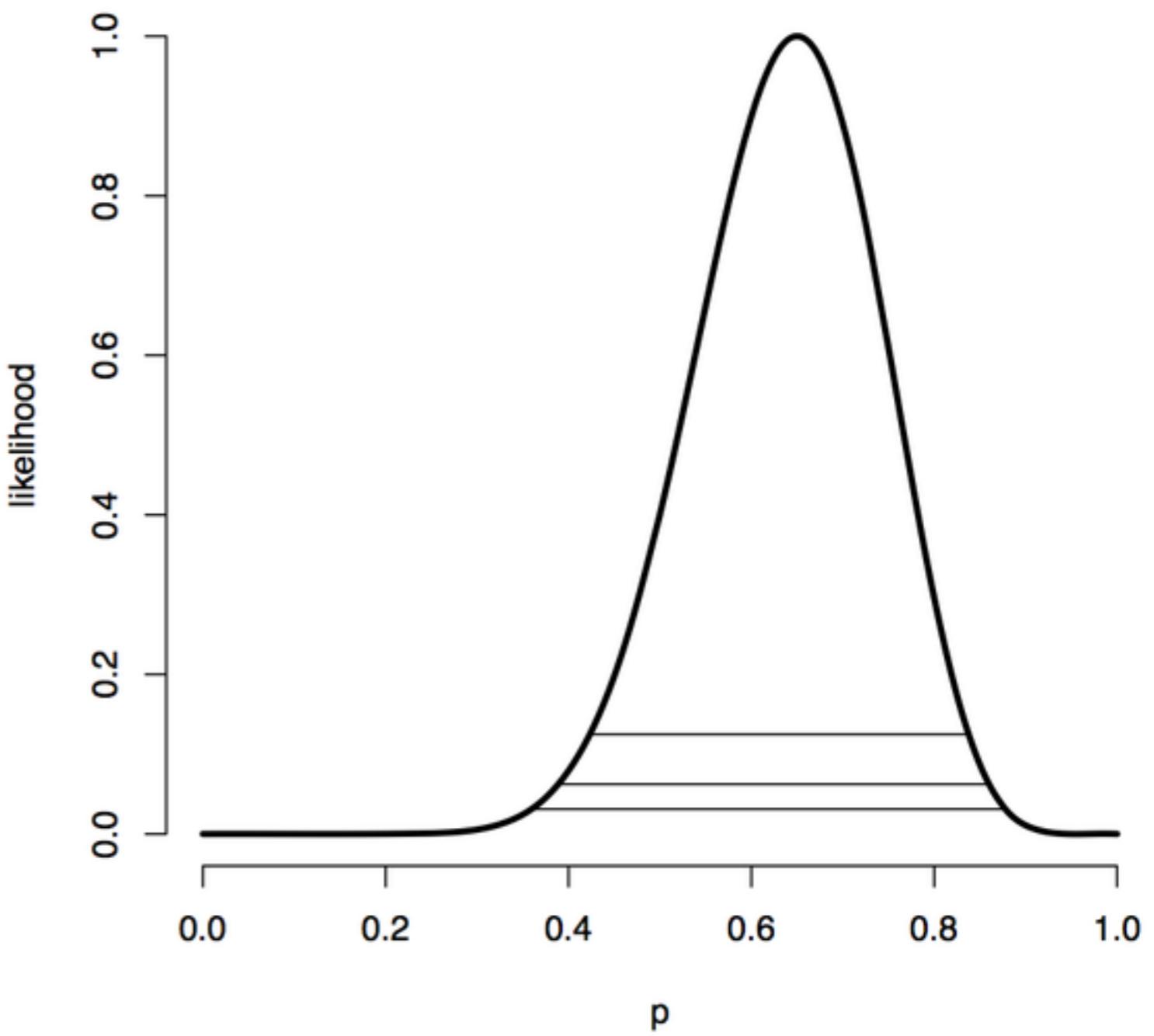
loss for B only

$$\theta^{A,B} = \operatorname{argmin}_{\theta} \mathcal{L}_B(\theta) + \frac{1}{2} \sum_i F_{i,i}^A (\theta_i - \theta_i^A)^2$$

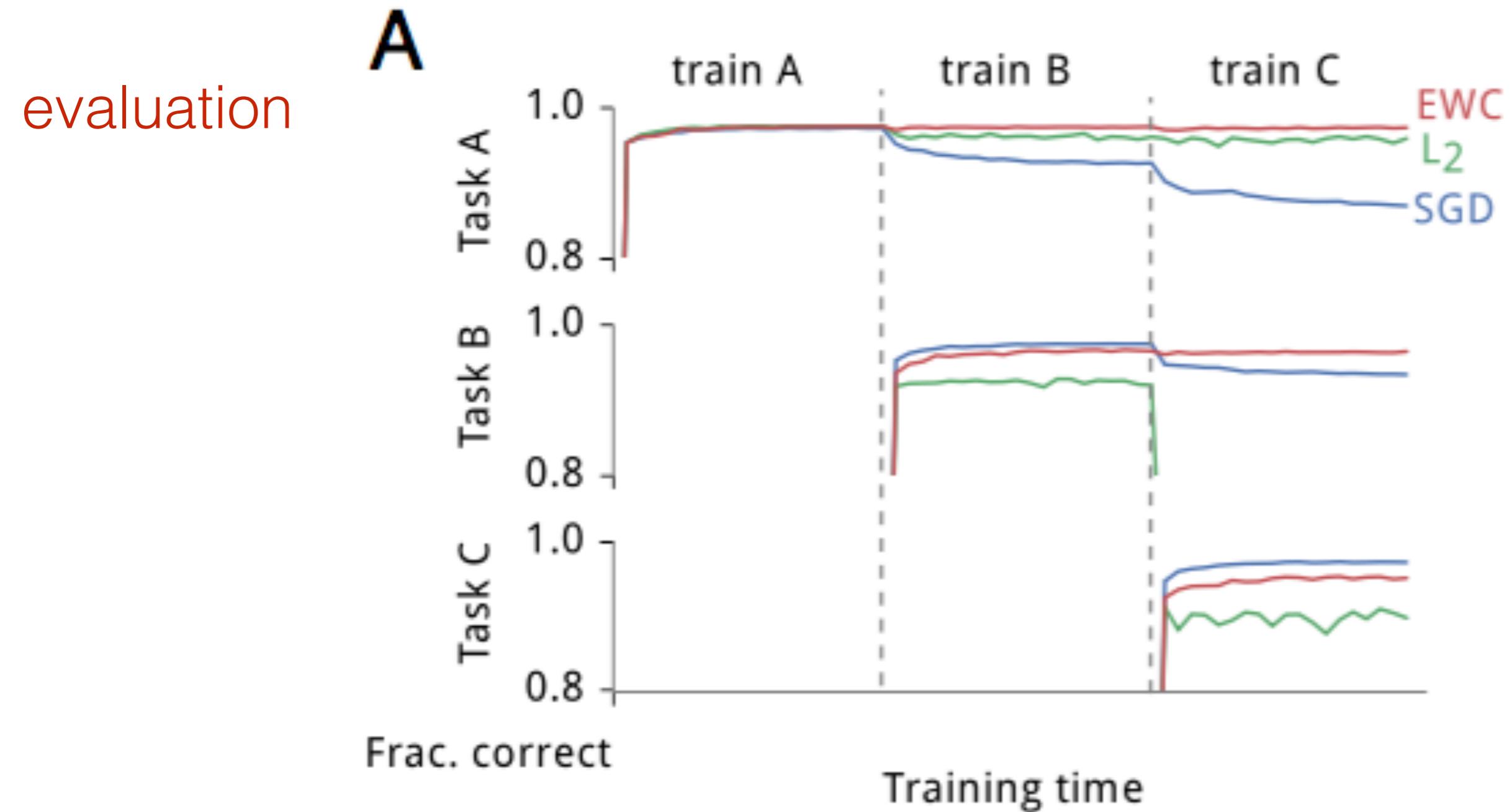
regulariser makes sure we don't catastrophically forget about A.

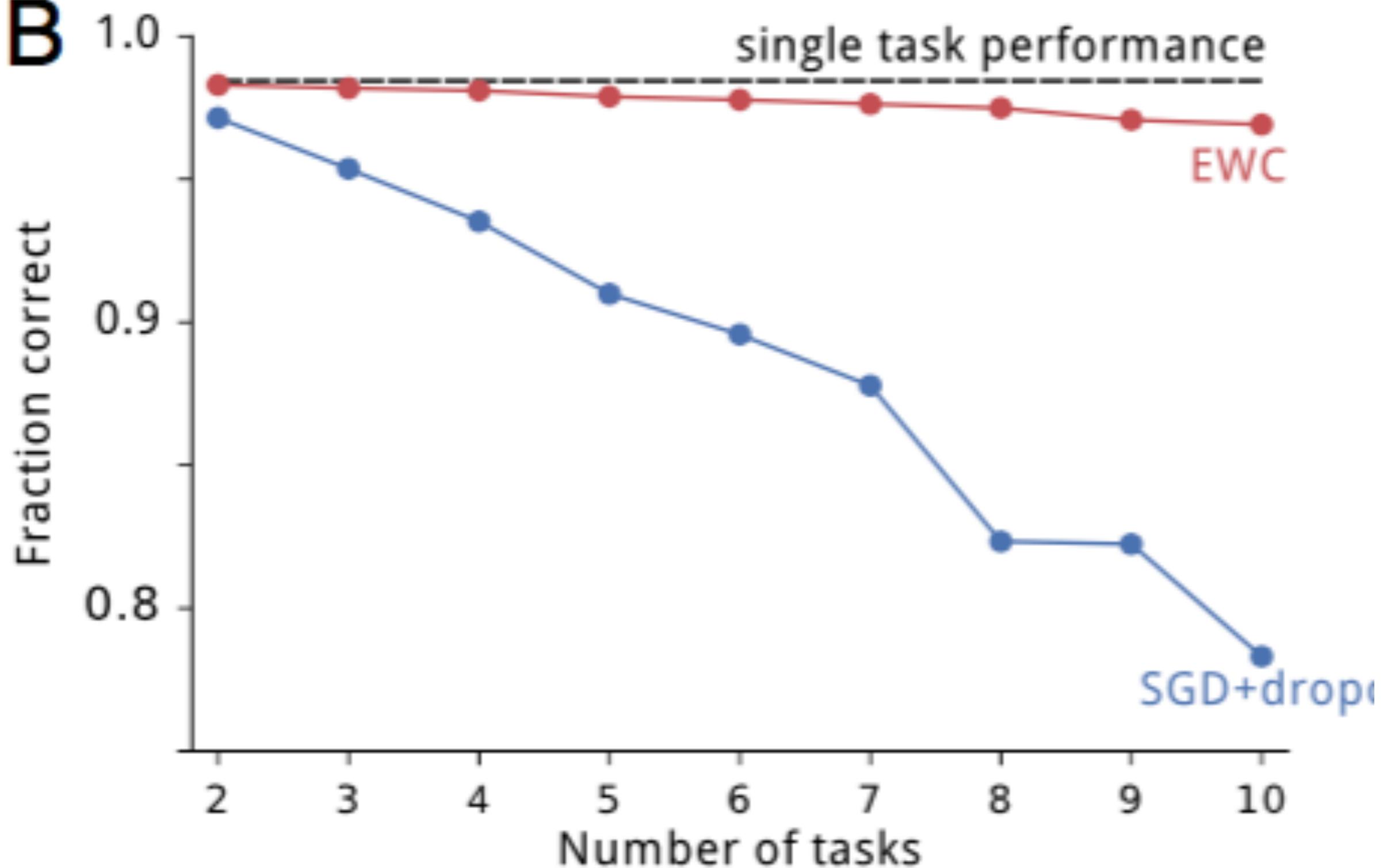
$$\textcolor{blue}{F}_{\theta} = \mathbb{E}_{\theta}[\dot{\ell}_{\theta}\dot{\ell}_{\theta}] = -\mathbb{E}_{\theta}[\nabla^2 \log p_{\theta}(X)].$$

- The bigger the fish information(negative second-order derivatives), the more pointy of the distribution of the likelihood
- likelihood is more sensitive to this parameter



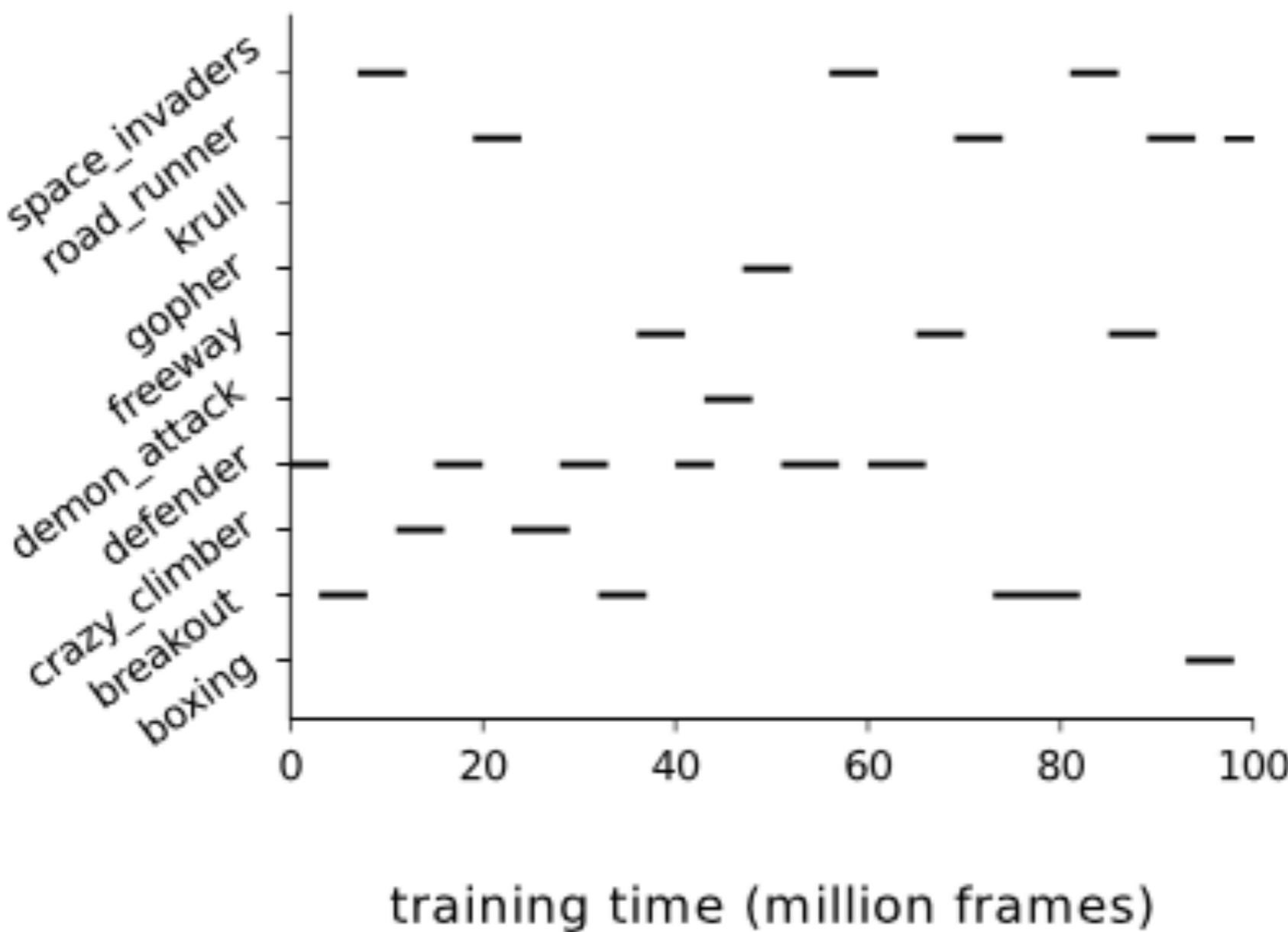
# Experiment 1 over MNIST



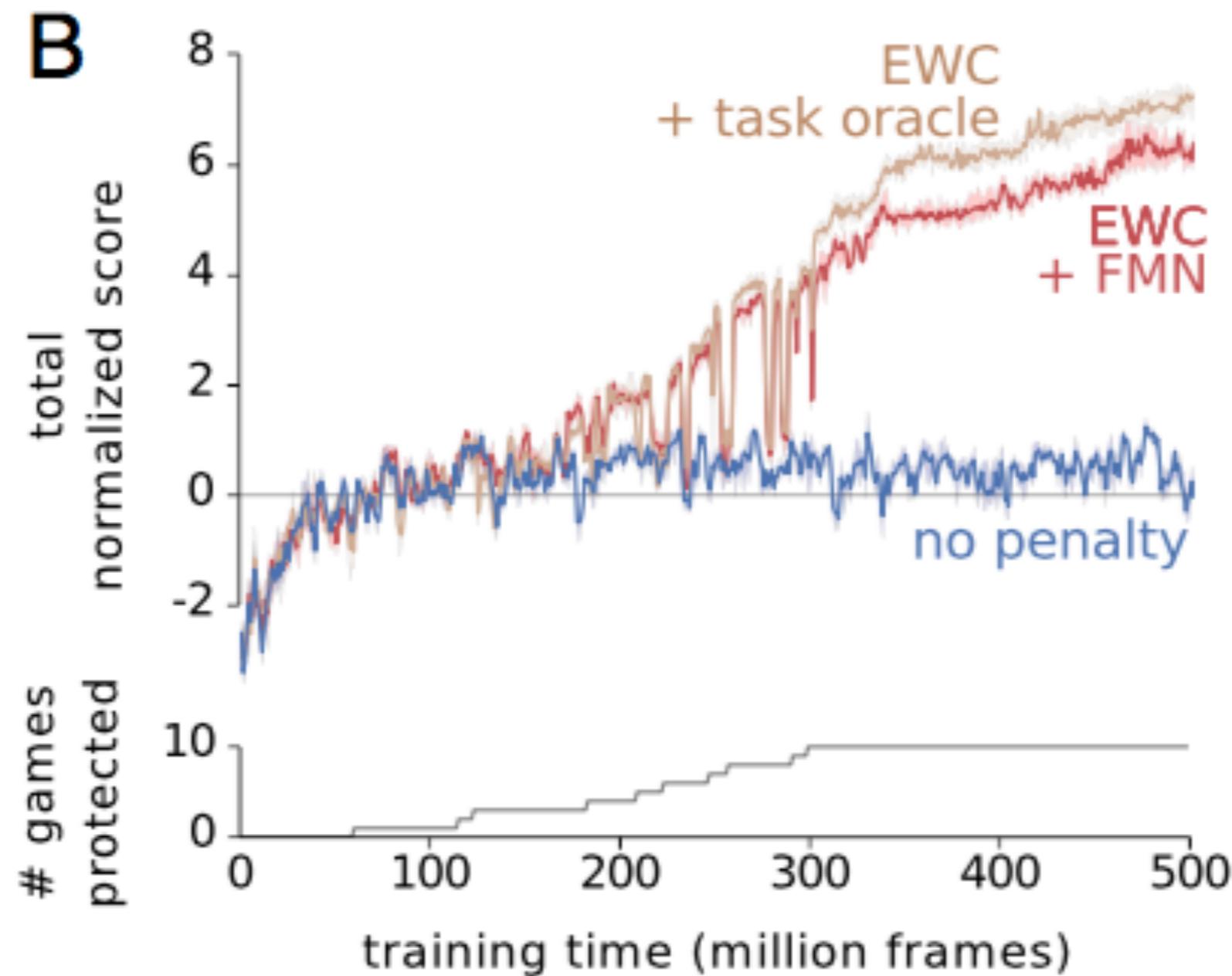
**B**

# Experiment 2 on Atari Training Schedule

**A**



10: human level,  
1: random agent



# Conclusion

- The ability to learn tasks in succession without forgetting is a core component of biological and artificial intelligence.
- Weight uncertainty should inform learning rates.
- EWC allows the network to effectively squeeze in more functionality into a network with fixed capacity,