Teaching a neural network to count: reinforcement learning with "social scaffolding"

once, which is related to the one-to-one principle [1].

demonstration.



Acknowledgement: This research is supported by the Stanford Center for the Study of Language and Information summer research program and the Rumelhart Emergent Cognition Fund.

Reward policy & Teaching strategies

Error types shown in graphs: 1: Stop early – say "done" when there exists an untouched object 2: Double touch – touching the same object twice in a trial 3: Skip – touching the (k+1)th object before touching the kth object

Note: These errors are not mutually exclusive. Also, touching empty space is not in itself an error though it can lead to delay in receipt of reward.

Teaching strategies: self-exploration trial.

Social scaffolding made learning easier, because: - Intermediate feedback makes the task more supervised. - Demonstration forces exposure to the optimal solution.

These results provide insights about how social scaffoldings support learning from a computational perspective. Further research will extend these explorations to multi-layer recurrent architectures and more complex task settings.

[1] Gelman, R., Gallistel, C. R. (1978). The child's understanding of number. Harvard U. press. [2] Sutton, R. S., & Barto, A. G. (1998). *Reinforcement learning*: An introduction. MIT press. [3] Lin, L.-J. (1993). Reinforcement learning for robots using neural networks. Technical Report, DTIC Document. [4] Van Hasselt, H., Guez, A., Silver, D. (2015). Deep Reinforcement Learning with Double Q-learning.

arXiv.

518(7540), 529-533.

Simulation source code:

There is always a final reward on successful completion. We also considered two types of "social scaffolding" and their combination: **1.** Intermediate reward: Reward the agent for touching the correct next object. The reward magnitude is one-half of the final reward. 2. Demonstration: Force the agent to execute the maximally

efficient action sequence and provide the corresponding reward. In this condition, we alternate demonstration trials and the regular

3. Intermediate reward + Demonstration: Combination of (1) and (2)

Summary

References

[5] Mnih, V., et al. (2015). Human-level control through deep reinforcement learning. *Nature*,

https://github.com/QihongL/mathCognition_PDP_RL