

Emergent planning in a recurrent neural network that plays Sokoban

Adrià Garriga-Alonso¹ Mohammad Taufeque¹ Adam Gleave¹

Abstract

To predict how advanced neural networks generalize to novel situations, it is essential to understand how they reason. Guez et al. (2019, “An investigation of model-free planning”) presented a recurrent neural network (RNN) trained with model-free reinforcement learning. The network’s ability to solve Sokoban puzzles improved when given extra computation steps at episode starts. We replicate and open-source their setup, and further investigate its planning behavior. Our findings reveal that the RNN benefits from additional thinking early in training, and that adding thinking time often reduces redundant actions. The results suggest that the RNN learns to take time to think by ‘pacing’ during training, despite the per-step penalties, indicating that training incentivizes planning capabilities. The small size (1.29M parameters) and interesting behavior of this model will greatly interest the mechanistic interpretability community.

1. Introduction

Humans benefit from taking time to think for many tasks, and large language models improve from thinking step by step (Kojima et al., 2022). Thus, we should expect advanced AIs to reason, and to perform better at tasks with more thinking. Understanding reasoning in such neural networks (NNs) is key to predicting their behavior in novel situations.

Of particular relevance to AI alignment is internal reasoning which furthers a goal. Hubinger et al. (2019) term “mesa-optimizers” neural networks that have learned to explicitly pursue goals through internal reasoning, warning that these goals may not be the same as the training objective (Di Langesco et al., 2022; Shah et al., 2022).

This paper presents a 1.29M parameter recurrent neural network (RNN) that clearly benefits from extra pondering time

¹FAR AI, San Diego, California, USA. Correspondence to: Adrià Garriga-Alonso <adria@far.ai>.

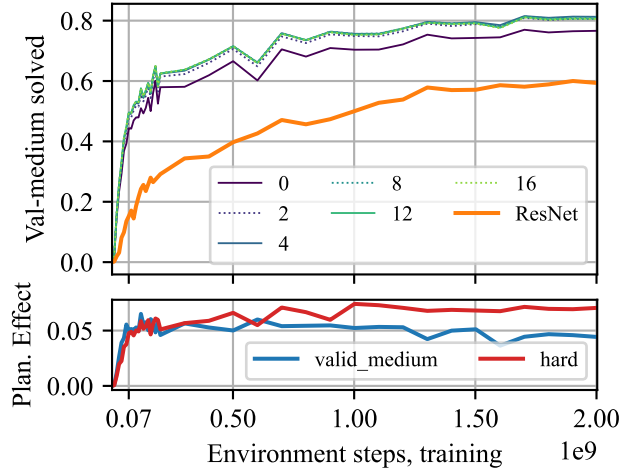


Figure 1. **Top:** Proportion of medium-difficulty validation levels solved vs. environment transitions used in training. Each curve displays the RNN model with a specific number of forced thinking steps at episode start, along with a ResNet baseline. **Bottom:** Estimate of the planning effect, the 8-steps curve minus the 0-steps curve. We can see that thinking develops in the first 70M steps and keeps increasing for the hard levels, but decreases for medium levels.

toward a concrete goal. Following Guez et al. (2019), we train a convolutional LSTM to play Sokoban, a challenging puzzle game that remains a benchmark for planning algorithms (Peters et al., 2023) and reinforcement learning (Chung et al., 2024).

We believe this NN is small, and its behavior straightforward enough for mechanistic interpretability techniques to reverse engineer it. Simultaneously, it performs reasoning in a complex environment, benefitting from extra thinking time. Thus, it is a useful first step towards understanding reasoning in sophisticated neural networks.

Our contributions are:

- We replicate and open-source an RNN that improves performance with thinking time (Guez et al., 2019).
- We present evidence that the thinking-time behavior emerges early in training and is incentivized by the

setup (Figure 1, Sections 3 and 4.1).

- We investigate the behavior changes by which pondering time improves performance (Section 3), showing that it reduces myopic behavior (Figure 4).
- We present case studies of planning behavior, suggesting the RNN naturally takes time to think (Section 4). This may explain the reduction of planning effect on medium levels (Figure 1).

2. Training the test subject

We closely follow the setup from Guez et al. (2019), which introduced the Deep Repeating ConvLSTM (DRC) recurrent architecture and trained it with reinforcement learning (RL) to play Sokoban. We make the trained networks¹ and training code² available.

Environment. Sokoban is a grid puzzle game with walls, floors, movable boxes, and target tiles. The player’s goal is to push all boxes onto target tiles while navigating walls. We use the Boxoban dataset (Guez et al., 2018), consisting of 10×10 procedurally generated Sokoban levels, each with 4 boxes and targets. The edge tiles are always walls, so the playable area is 8×8 . Boxoban separates levels into train, validation and test sets, with three difficulty levels: unfiltered, medium and hard. In this paper, we use unfiltered-train (900k levels), unfiltered-test (1k levels), and medium-difficulty validation (50k levels) sets.

The observations are 10×10 RGB images, normalized by dividing each component by 255. Each type of tile is represented by a pixel of a different color (Schrader, 2018). See Figure 5 for examples. The player can take actions (Up, Down, Left, Right). The reward is -0.1 per step, 1 for putting a box on a target, -1 for removing a box, and 10 for finishing the level. To avoid strong time correlations during learning, each episode resets at a uniformly random length between 91 and 120 time steps.

DRC(D, N) architecture. Guez et al. (2019) introduced the Deep Repeating ConvLSTM (DRC), whose core consists of D convolutional LSTM layers with 32 channels and 3×3 filters, each applied N times per time step. Our DRC(3, 3) has 1.29M parameters. Before the LSTM core, two convolutional layers (without nonlinearity) encode the observation with 4×4 filters.

The LSTM core uses 3×3 convolutional filters, and a nonstandard \tanh on the output gate (Jozefowicz et al., 2015). Unlike the original ConvLSTM (Shi et al., 2015), the input to each layer of a DRC consists of several concatenated

components:

- The encoded observation is fed into each layer.
- To allow spatial information to travel fast in the ConvLSTM layers, we apply *pool-and-inject* by max- and mean-pooling the previous step’s hidden state. We linearly combine these values channel-wise before feeding them as an input to the next step.
- To avoid convolution edge effects from disrupting the LSTM dynamics, we feed in a 12×12 channel with zeros on the inside and ones on the boundary. Unlike the other inputs, this one is not zero-padded, maintaining the output size.

ResNet architecture. This is a convolutional residual neural network, also from Guez et al. (2019). It serves as a non-recurrent baseline that cannot think for a variable number of steps but is nevertheless good at the game. The ResNet consists of 9 blocks, each with 4×4 convolutional filters. The first two blocks have 32 channels, and the others have 64. Each block consists of a convolution, then two (relu, conv) sub-blocks, which each split off and are added back to the trunk. The ResNet has 3.07M parameters.

Value and policy heads. After the image processing, an affine layer projects the flattened spatial output into 256 channels. We then apply a ReLU and two different affine layers: one for the actor (policy) and one for the critic (value function).

RL training. We train each network for 2.003 billion environment steps using IMPALA (Espeholt et al., 2018; Huang et al., 2023). For each training iteration, we collect 20 transitions on 256 actors using the network parameters from the previous iteration, and simultaneously take a gradient step. We use a discount rate of $\gamma = 0.97$ and V-trace $\lambda = 0.5$. The value and entropy loss coefficients are 0.25 and 0.01. We use the Adam optimizer with a learning rate of $4 \cdot 10^{-4}$, which linearly anneals to $4 \cdot 10^{-6}$ at the end of training. We clip the gradient norm to $2.5 \cdot 10^{-4}$. Our hyperparameters are mostly the same as Guez et al. (2019); see Appendix A.

A* solver. To obtain optimal solutions to each Sokoban puzzle, we used the A* search algorithm. The heuristic was the sum of the Manhattan distances of each box to its nearest target. Solving a single level on one CPU takes anywhere from a few seconds to 15 minutes.³

¹<https://github.com/AlignmentResearch/learned-planner>

²<https://github.com/AlignmentResearch/train-learned-planner>

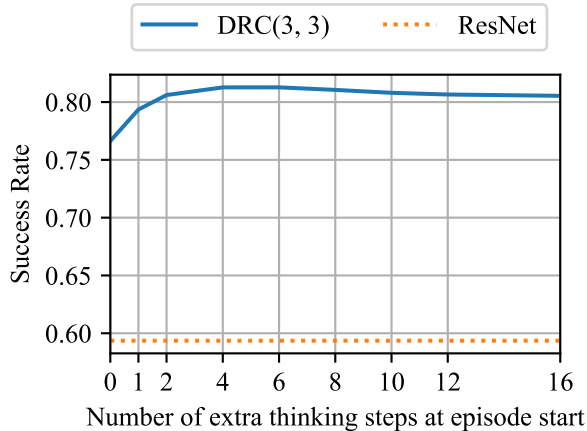


Figure 2. Proportion of solved medium-difficulty validation levels solved by DRC with x time steps to think. The dotted line shows the ResNet’s solve proportion at 2×10^9 training steps for the same levels. The proportion of solved levels increases with thinking time, but only up to 6 steps.

3. Quantitative behavior analysis

From here onward, we may refer to the DRC(3, 3) as DRC.

3.1. Are planning networks better at Sokoban?

After 2B steps from the environment, the DRC(3, 3) is able to solve 99.3% of unfiltered test levels, and the ResNet 97.9%, a score consistent with Guez et al. (2019), see Appendix B. Of the hard levels, the DRC(3, 3) solves 42.8% and the ResNet 26.2%. This gap appears larger, thus the question: are networks that plan comparatively better at harder levels? Figure 9 compares, for each checkpoint and architecture, how successful it is at various difficulties. All architectures follow the same curve. The DRC(3, 3) is better, but not comparatively better at harder levels.

3.2. Does thinking time improve performance?

Does our trained DRC(3, 3) subject replicate the thinking time effect from Guez et al. (2019)? We evaluate the DRC on the whole of medium-difficulty validation levels. To give the DRC n steps to think, we repeat the initial environment observation n times while advancing the DRC hidden state, then let the DRC policy act normally. Figure 2 shows the DRC’s performance for $n \in \{0, 1, 2, 4, 6, 8, 10, 12, 16\}$ steps, as well as the ResNet baseline. The DRC’s solving rate improves from 76.6% to 81.3% with more thinking steps, with peak performance at 6 steps. The effect (4.7%

³The A* solutions may be of independent interest, so we make them available at <https://huggingface.co/datasets/AlignmentResearch/boxoban-astar-solutions/>.

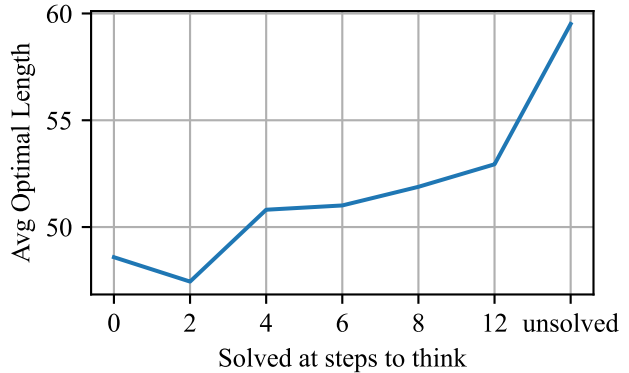


Figure 3. Average optimal solution length of levels grouped by the number of thinking steps at which the level is first solved. Each of the 50,000 medium-difficulty validation levels is assigned to a group.

difference) is about as strong as that observed by Guez et al. (2019) (4.5%).

3.3. Planning emerges early on in training

When does the DRC learn to benefit from thinking steps? Figure 1 shows the performance of the DRC on the medium-difficulty validation levels. The DRC rapidly acquires a strong planning effect by 70M environment steps, which strengthens as training progresses on the hard levels; but slowly weakens for the medium levels. This suggests that the training setup and loss reinforce planning behavior, at least for harder levels.

3.4. Do difficult levels benefit more from thinking?

Intuitively, taking extra thinking time should help more with harder levels. The A* solver gives us two proxies for the difficulty of a level: the number of nodes it had to expand to find a solution and the length of the optimal solution.

We can measure the average difficulty for levels that the DRC needs n steps of thinking to solve. Figure 3 shows the average solution length proxy, and it clearly trends up. Figure 8 (Appendix C) shows the expanded nodes proxy, and trends up but is noisy for lower n . We speculate that this is because the DRC has different heuristics than the one we use for A*, whereas the optimal solution length does not depend on heuristics.

3.5. What happens when we make the DRC think?

Figure 2 shows that, when explicitly made to think, the DRC can solve many more levels. How does that happen? Does it have any effects other than solving more levels?

Table 1. Levels at 6 thinking steps.

LEVEL CATEGORIZATION	PERCENTAGE
Solved, previously unsolved	6.87
Unsolved, previously solved	2.23
Solved, with better returns	18.98
Solved, with the same returns	50.16
Solved, with worse returns	5.26
Unsolved, with same or better returns	15.14
Unsolved, with worse returns	1.36

Thinking makes the DRC “patient.” One hypothesis is that, without the extra thinking steps, the DRC puts some boxes into targets without realizing that makes the puzzle unsolvable. This makes sense based on the training signal, which penalizes every step it takes to solve the puzzle – it need not be “irrational” for the training process to push the DRC to be “impatient” and solve part of the puzzle on-the-go if that lets it solve puzzles fast enough.

For n thinking steps, Figure 4 plots the average number of steps (*time-to-box*) until each of 4 boxes gets on a target, in the order in which they do so. The results support the hypothesis: levels solved at 6 thinking steps but not at 0 thinking steps show increased time-to-box for boxes 1–3 but decreased time for box 4. In contrast, on average over all levels, times go down for all boxes.

Thinking often improves solutions for already-solved levels but is sometimes harmful. Is an increased solution rate the only benefit of artificial thinking, or does it increase returns in other ways? Table 1 breaks down the medium-difficulty levels into types according to the behavioral effect of adding 6 steps of thinking. The breakdown shows two major benefits of thinking time: the DRC solves some levels that it did not previously solve, but also many levels that were already solved previously increase in return. The only way to improve undiscounted returns on solved levels is to take fewer steps, which matches Figure 4.

Of the levels that the DRC solved both when thinking and not thinking (81.3%), about 31.8% had a higher return, and the remaining had the same or worse return. Overall, extra thinking steps are not always good, but on balance they let the DRC solve more levels or improve return on already solved levels.

4. Case studies: how does extra thinking change DRC behavior?

Let us examine an example from each of three categories from Table 1: a solved, previously unsolved level; an unsolved, previously solved level; and a level on which the

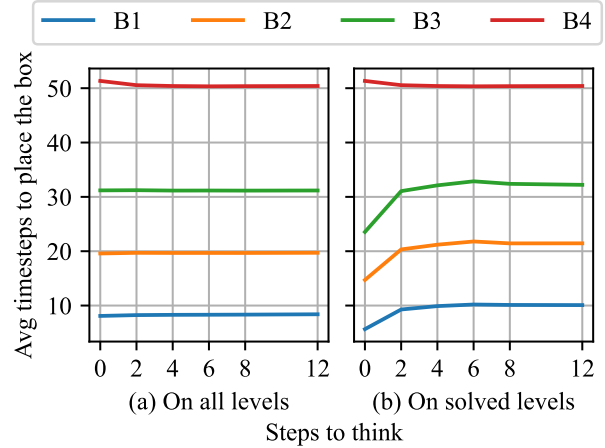


Figure 4. Average time step to place each box B_i on target for different thinking steps. (a) Averages across *all* levels where the box B_i is placed on target by DRC with n thinking steps. If box B_i is not placed (on unsolved levels) we ignore the box, but we keep the other boxes in the level. (b) Averages for all levels solved by 6 thinking steps but not solved by 0 thinking steps.

DRC improved its return with extra thinking. We picked one out of ten examples to illustrate each category.

Thinking lets the DRC solve, Figure 5(a). In the no-thinking condition, the DRC first pushes box C one square to the right. It then goes back to push A to a , but it is too late: it is now impossible to push box B onto b . The network pushes C onto c and D onto d and then remains mostly still. In contrast, after thinking, the DRC pushes A to a first, which lets it solve the level.

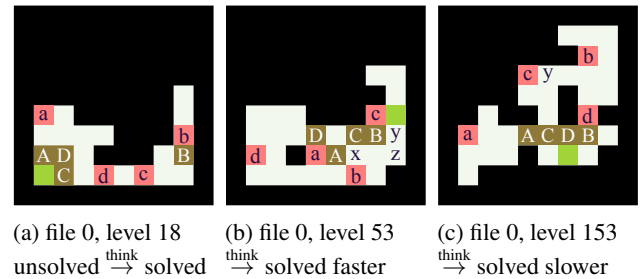


Figure 5. Case studies of three levels demonstrating different behaviors after 6 thinking steps. *Legend:* green is the player, brown (upper case) are boxes, pink (lower case) are box targets, and black are walls. Box and target letters are in the other the correct solution places them together. Videos available at [this https URL](https://www.youtube.com/watch?v=...); see Section 4 for behavior descriptions. Levels solved faster obtain higher return, the agent incurs the per-step penalty fewer times.

Thinking speeds up solving, Figure 5(b). In the no-thinking condition, the DRC spends many steps going back and forth before pushing any boxes. First it goes down to y , then up to c , then down onto z , back up to y and to z again. It then proceeds to solve the rest of the puzzle correctly: push box A onto a , prepare box B on x and box C where B originally was, push in boxes B, C and finally D . In the thinking condition, the DRC makes a beeline for A and then plays the exact same solution.

Thinking slows down solving, Figure 5(c). In the no-thinking condition, the DRC starts by pushing box C into position y , then pushes boxes A, B . On the way back down, the DRC pushes C onto c and finally D onto d . In contrast, after thinking for a bit, the DRC goes the other way and starts by pushing B onto b . The solution is the same after that: A, C , then D ; but the NN has wasted time trekking back from B to A .

4.1. Hypothesis: the DRC ‘paces’ to get thinking time

One hypothesis for why the planning effect exists is that the DRC occasionally deliberately moves around in cycles (without touching any box) until it comes up with a plan to solve the episode. To test this, we can check whether cycles in the game state occur more frequently near the start of the episode, and whether deliberately giving the DRC thinking time makes it stop going in cycles.

Figure 6 shows these results: the majority of cycles start within the first 5 steps of the episode, and forcing it to think for six steps makes about 75% of these cycles disappear. The second case study from Section 4 also displays this qualitative ‘pacing’ behavior.

Episode start is not the only time at which the DRC goes in cycles. Figure 10 shows that replacing an n -length cycle with n thinking steps leads the NN to have the *exact same behavior* for at least 65% of levels, for at least 30 steps after the cycle. For context, the median solution length for train-unfiltered is exactly 30 steps. Additionally, in 82.39% of cases, doing this prevents cycles from starting in the n steps after the thinking-steps conclude (Table 4).

Why does this behavior emerge during training? Thinking is useful for achieving higher return, so it should be reinforced. But it also has a cost, -0.1 per step, so it should be discouraged in easy levels that do not need the computation. We speculate that, as training advances and heuristics get tuned, the DRC needs to think in for fewer levels, and it is better at knowing when it needs to pace. This would doubly explain the decline in planning behavior for medium levels in Figure 1: the DRC finds levels easier, *and* also knows to pace when needed. Zero steps of planning only includes *induced* planning, and does not include planning by pacing.

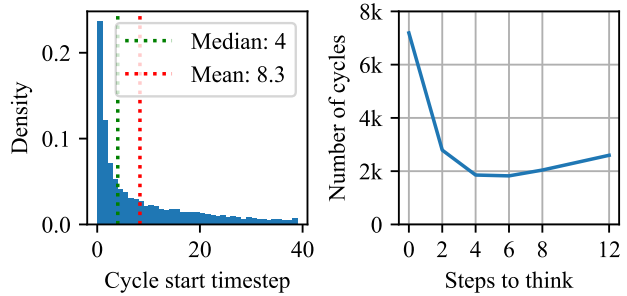


Figure 6. **Left:** Histogram of cycle start times on the medium-difficulty validation levels. **Right:** Total number of cycles the agent takes in the first 5 steps across all episodes in medium-difficulty validation levels with n initial thinking steps.

5. Related work

Interpretability of agents and world models. Several works have attempted to find the mechanism by which a simple neural network does planning in mazes (Ivanitskiy et al., 2023; Mini et al., 2023), gridworlds (Bloom & Colognese, 2023), and graph search (Ivanitskiy et al., 2023). We believe the DRC we present is a clearer example of an agent than what these works focus on, and should be similarly possible to interpret.

Other works have found emergent world models in sequence-prediction (Li et al., 2023) and navigation (Wijmans et al., 2023) neural networks.

Goal misgeneralization and mesa-optimizers for alignment. From the alignment perspective, AIs optimizing monomaniacally for a goal have been a concern for a long time (e.g. Russell, 2019). In a machine learning paradigm (Hubinger et al., 2019), the goal of the training system is not necessarily optimized; instead, the NN may optimize for a related or different goal (Di Langosco et al., 2022; Shah et al., 2022). Or, of course, for no goal at all.

Chain-of-thought faithfulness. Large language models (LLMs) use chain of thought, but are they faithful to it or do they think about their future actions in other ways (Lanham et al., 2023; Pfau et al., 2024)? One could hope that LLMs perform all long-term reasoning in plain English, allowing unintended human consequences to be easily monitored, as in Scheurer et al. (2023).

Reasoning neural network architectures. Many papers try to enhance NN thinking by altering the training setup (Bansal et al., 2022; Graves, 2016; Chung et al., 2024).

Ethical treatment of AIs. Do AIs deserve moral consideration? Schwitzgebel & Garza (2015) argue that very human-like AIs are conceivable and clearly deserve rights.

Tomasik (2015) suggests that most AIs deserve at least a little consideration, like biological organisms of any species (Singer, 2004). Daswani & Leike (2015) argue that the way to measure pleasure and pain in a reinforcement learner is not by its absolute amount of return, but rather by the temporal difference (TD) error: the difference between its expectations and the actual return it obtained. If the internals of the NN have a potentially different objective (Hubinger et al., 2019; Di Langosco et al., 2022), then the TD error would have to come from a place *other* than the critic. This paper is an early step toward finding the learned-reward internal TD error, if it exists.

6. Conclusion

We have replicated and open-sourced a recurrent network that plays Sokoban, which benefits from additional thinking steps at test time (Guez et al., 2019).

We have shown that thinking for longer helps solve harder levels, and makes the DRC better at levels that require longer-sighted behavior. Without intervention, thinking sometimes takes the form of the DRC ‘pacing’ at the beginning or middle of the episode, in a way which can be substituted by repeating the same input; suggesting it is deliberately using more computation.

We have shown that the training setup incentivizes the planning effect at the start, and that it is disincentivized later, but only for easier levels. Finally, we offer a hypothesis about why the training process disincentivizes the planning effect: the NN finds levels easier (needs less thinking), and also learns when to do the thinking it needs (via pacing).

We believe this work will be useful to the interpretability, alignment and ethical treatment of AI communities.

7. Acknowledgements

We would like to thank ChengCheng Tan for help with editing this paper, Philip Quirke for help getting it done, and the rest of the FAR team for general support. We thank Alex F. Spies, Maximilian Li, and anonymous reviewers for improving the paper with their comments. We are also very grateful for the many questions about their RL setup that Arthur Guez and Timothy Lillicrap answered, as well as inspiration from Stephen Chung.

References

Bansal, A., Schwarzschild, A., Borgnia, E., Emam, Z., Huang, F., Goldblum, M., and Goldstein, T. End-to-end algorithm synthesis with recurrent networks: Extrapolation without overthinking. In Oh, A. H., Agarwal, A., Belgrave, D., and Cho, K. (eds.), *Advances in Neural*

Information Processing Systems, 2022. URL <https://openreview.net/forum?id=PPjSKy40XUB>.

Bloom, J. and Colognese, P. Decision transformer interpretability. , 2023.

Chung, S., Anokhin, I., and Krueger, D. Thinker: Learning to plan and act. *Advances in Neural Information Processing Systems*, 36, 2024.

Daswani, M. and Leike, J. A definition of happiness for reinforcement learning agents. *CoRR*, 2015. URL <http://arxiv.org/abs/1505.04497v1>.

Di Langosco, L. L., Koch, J., Sharkey, L. D., Pfau, J., and Krueger, D. Goal misgeneralization in deep reinforcement learning. In *International Conference on Machine Learning*, pp. 12004–12019. PMLR, 2022.

Espeholt, L., Soyer, H., Munos, R., Simonyan, K., Mnih, V., Ward, T., Doron, Y., Firoiu, V., Harley, T., Dunning, I., Legg, S., and Kavukcuoglu, K. IMPALA: Scalable distributed deep-RL with importance weighted actor-learner architectures. *CoRR*, 2018. URL <http://arxiv.org/abs/1802.01561v3>.

Graves, A. Adaptive computation time for recurrent neural networks. *CoRR*, 2016. URL <http://arxiv.org/abs/1603.08983v6>.

Guez, A., Mirza, M., Gregor, K., Kabra, R., Racaniere, S., Weber, T., Raposo, D., Santoro, A., Orseau, L., Eccles, T., Wayne, G., Silver, D., Lillicrap, T., and Valdes, V. An investigation of model-free planning: boxoban levels. <https://github.com/deepmind/boxoban-levels/>, 2018.

Guez, A., Mirza, M., Gregor, K., Kabra, R., Racanière, S., Weber, T., Raposo, D., Santoro, A., Orseau, L., Eccles, T., Wayne, G., Silver, D., and Lillicrap, T. An investigation of model-free planning. *CoRR*, 2019. URL <http://arxiv.org/abs/1901.03559v2>.

Heek, J., Levskaya, A., Oliver, A., Ritter, M., Rondepierre, B., Steiner, A., and van Zee, M. Flax: A neural network library and ecosystem for JAX, 2023. URL <http://github.com/google/flax>.

Huang, S., Weng, J., Charakorn, R., Lin, M., Xu, Z., and Ontañón, S. Cleanba: A reproducible and efficient distributed reinforcement learning platform, 2023.

Hubinger, E., van Merwijk, C., Mikulik, V., Skalse, J., and Garrabrant, S. Risks from learned optimization in advanced machine learning systems. 2019. URL <https://intelligence.org/learned-optimization/>.

- Ivanitskiy, M., Spies, A. F., Räuker, T., Corlouer, G., Mathwin, C., Quirke, L., Rager, C., Shah, R., Valentine, D., Behn, C. D., Inoue, K., and Fung, S. W. Linearly structured world representations in maze-solving transformers. In *UniReps: the First Workshop on Unifying Representations in Neural Models*, 2023. URL <https://openreview.net/forum?id=pZakRK1QHU>.
- Jozefowicz, R., Zaremba, W., and Sutskever, I. An empirical exploration of recurrent network architectures. In Bach, F. and Blei, D. (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 2342–2350, Lille, France, 07–09 Jul 2015. PMLR. URL <https://proceedings.mlr.press/v37/jozefowicz15.html>.
- Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., and Iwasawa, Y. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35: 22199–22213, 2022.
- Lanham, T., Chen, A., Radhakrishnan, A., Steiner, B., Denison, C., Hernandez, D., Li, D., Durmus, E., Hubinger, E., Kernion, J., et al. Measuring faithfulness in chain-of-thought reasoning. *arXiv preprint arXiv:2307.13702*, 2023.
- Li, K., Hopkins, A. K., Bau, D., Viégas, F., Pfister, H., and Wattenberg, M. Emergent world representations: Exploring a sequence model trained on a synthetic task. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=DeG07_TcZvT.
- Mini, U., Grietzer, P., Sharma, M., Meek, A., MacDiarmid, M., and Turner, A. M. Understanding and controlling a maze-solving policy network. *CoRR*, 2023. URL <http://arxiv.org/abs/2310.08043v1>.
- Peters, N. S., Alexa, M., and Neurotechnology, S. F. Solving sokoban efficiently: Search tree pruning techniques and other enhancements. 2023.
- Pfau, J., Merrill, W., and Bowman, S. R. Let’s think dot by dot: Hidden computation in transformer language models. *CoRR*, 2024. URL <http://arxiv.org/abs/2404.15758v1>.
- Russell, S. *Human compatible: artificial intelligence and the problem of control*. Penguin, 2019.
- Scheurer, J., Balesni, M., and Hobbhahn, M. Large language models can strategically deceive their users when put under pressure. *CoRR*, 2023. URL <http://arxiv.org/abs/2311.07590v3>.
- Schrader, M.-P. B. gym-sokoban. <https://github.com/mpSchrader/gym-sokoban>, 2018.
- Schwitzgebel, E. and Garza, M. A defense of the rights of artificial intelligences. *Midwest Studies In Philosophy*, 39(1):98–119, 2015. doi: 10.1111/misp.12032. URL <http://www.faculty.ucr.edu/~eschwitz/SchwitzPapers/AIRights-150915.htm>.
- Shah, R., Varma, V., Kumar, R., Phuong, M., Krakovna, V., Uesato, J., and Kenton, Z. Goal misgeneralization: why correct specifications aren’t enough for correct goals. *arXiv preprint arXiv:2210.01790*, 2022.
- Shi, X., Chen, Z., Wang, H., Yeung, D.-Y., Wong, W.-K., and Woo, W.-c. Convolutional LSTM network: A machine learning approach for precipitation nowcasting. *Advances in neural information processing systems*, 28, 2015.
- Singer, P. Animal liberation. In *Ethics: Contemporary Readings*, pp. 284–292. Routledge, 2004.
- Tomasik, B. A dialogue on suffering subroutines. , 2015.
- Weng, J., Lin, M., Huang, S., Liu, B., Makoviichuk, D., Makovychuk, V., Liu, Z., Song, Y., Luo, T., Jiang, Y., Xu, Z., and Yan, S. EnvPool: A highly parallel reinforcement learning environment execution engine. In Koyejo, S., Mohamed, S., Agarwal, A., Belgrave, D., Cho, K., and Oh, A. (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 22409–22421. Curran Associates, Inc., 2022.
- Wijmans, E., Savva, M., Essa, I., Lee, S., Morcos, A. S., and Batra, D. Emergence of maps in the memories of blind navigation agents. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=lTt4KjHSsylv>.

A. Training hyperparameters

All networks were trained with the same hyperparameters, which were tuned on a combination of the ResNet and the DRC(3, 3). These are almost exactly the same as Guez et al. (2019), allowing for taking the *mean* of the per-step loss instead of the sum.

Loss. The value and entropy coefficients are 0.25 and 0.01 respectively. It is very important to *not* normalize the advantages for the policy gradient step.

Gradient clipping and epsilon The original IMPALA implementation, as well as Huang et al. (2023), *sum* the per-step losses. We instead average them for more predictability across batch sizes, so we had to scale down some parameters by a factor of 1/640: Adam ϵ , gradient norm for clipping, and L2 regularization).

Weight initialization. We initialize the network with the Flax (Heek et al., 2023) default: normal weights truncated at 2 standard deviations and scaled to have standard deviation $\sqrt{1/\text{fan_in}}$. Biases are initialized to 0. The forget gate of LSTMs has 1 added to it (Jozefowicz et al., 2015). We initialize the value and policy head weights with orthogonal vectors of norm 1. Surprisingly, this makes the variance of these unnormalized residual networks decently close to 1.

Adam optimizer. As our batch size is medium-sized, we pick $\beta_1 = 0.9$, $\beta_2 = 0.99$. The denominator epsilon is $\epsilon = 1.5625 \cdot 10^{-7}$. Learning rate anneals from $4 \cdot 10^{-4}$ at the beginning to $4 \cdot 10^{-6}$ at 2,002,944,000 steps.

L2 regularization. In the training loss, we regularize the policy logits with L2 regularization with coefficient 1.5625×10^{-6} . We regularize the actor and critic heads' weights with L2 at coefficient 1.5625×10^{-8} . We believe this has essentially no effect, but we left it in to more closely match Guez et al. (2019).

Number of training steps. In the body of the paper we state the networks train for 2.003B steps. The exact number is 2,002,944,000 steps. This is a holdover from earlier on in development, when we tested with $20\text{M} \approx 20\,029\,440 = 5120 \cdot 3912$ steps, which is divisible by 256 environments \times 20 steps collected. This number of steps is slightly larger than it could have been: 2B is also divisible by 5120.

Software. We base our IMPALA implementation on Cleanba (Huang et al., 2023). We implemented Sokoban in C++ using Envpool (Weng et al., 2022) for faster training, based on gym-sokoban (Schrader, 2018).

B. Learning curve comparison

It is difficult to fully replicate the results by Guez et al. (2019). (Chung et al., 2024) propose an improved method for RL in planning-heavy domains. They employ the IMPALA DRC(3, 3) as a baseline, but in Figure 5 one can see separate curves for what was reported by Guez et al. (2019) and what a decently-tuned DRC baseline looks like.

We did not innovate in RL, so were able to spend more time on the replication. We compare our replication to (Guez et al., 2019) in Figure 7, which shows that the learning curves for DRC(3, 3) and ResNet are compatible, but not the one for DRC(1, 1). Our implementation also appears much less stable, with large error bars and large oscillations over time. We leave addressing that to future work.

The success rate in Figure 7 is computed over 1024 random levels, unlike the main body of the paper. Table 3 reports test and validation performance for the DRC and ResNet seeds which we picked for the paper body.

The parameter counts (Table 2) are very different from what Guez et al. (2019) report.

C. Additional figures

Table 2. Parameter counts for each architecture.

Architecture	Parameter count	
DRC(3, 3)	1,285,125	(1.29M)
DRC(1, 1)	987,525	(0.99M)
ResNet:	3,068,421	(3.07M)

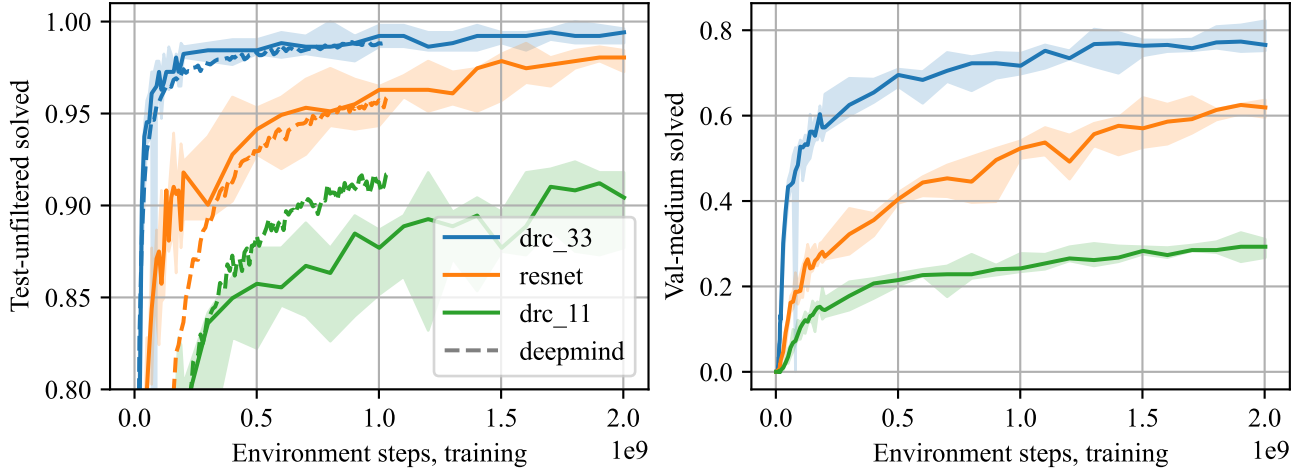


Figure 7. Success rate for Test-unfiltered and Validation-medium levels vs. environment steps of training. Each architecture has 5 random seeds, the solid line is the pointwise median and the shaded area spans from the minimum to the maximum. The dotted lines are data for the performance of architectures extracted from the (Guez et al., 2019) PDF file. The values are slightly different from what Figure 1 and Section 3 report because they are calculated on a random sample of 1024 levels (24 levels are repeated for test-unfiltered).

Table 3. Success rate and return of DRC and ResNet on the unfiltered test set at various training environment steps.

TRAINING ENV STEPS	TEST UNFILTERED				VALID MEDIUM			
	RESNET		DRC(3, 3)		RESNET		DRC(3, 3)	
	SUCCESS	RETURN	SUCCESS	RETURN	SUCCESS	RETURN	SUCCESS	RETURN
100M	87.8	8.13	95.4	9.58	18.6	-6.59	47.9	-0.98
500M	93.1	9.24	97.9	10.21	39.7	-2.64	66.6	2.62
1B	95.4	9.75	99.2	10.47	50.0	-0.64	70.4	3.40
2B	97.9	10.29	99.3	10.52	59.4	1.16	76.6	4.52

Table 4. The percentage of cycles that disappear when the network is given N artificial thinking steps before it is about to begin an N -length cycle.

Total cycles	13702
No cycles at the end of thinking steps	86.74%
No cycles in the next N timesteps	82.39%

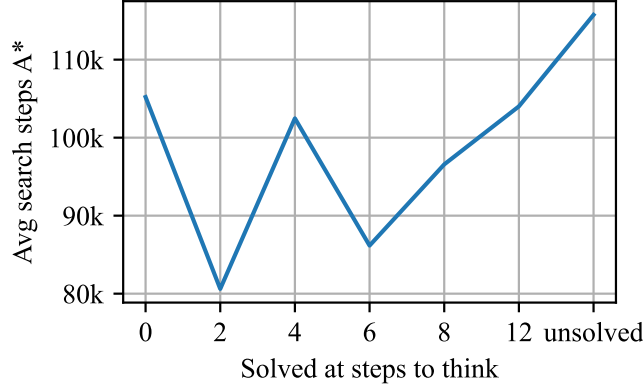


Figure 8. Number of thinking steps required to solve the level vs. number of nodes A^* needs to expand to solve it. The trend is somewhat increasing but much less clear, indicating different heuristics used by the NN and A^* .

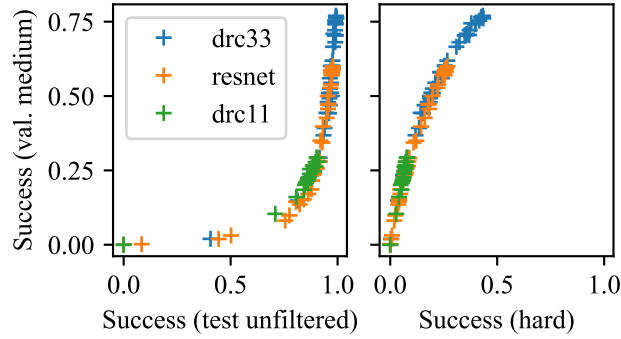


Figure 9. Success rate on datasets of various difficulty, for various checkpoints of each architecture. This deviates very little from a curve, which shows that ResNets and DRCs which are equally good at the easier sets are also equally good at the harder sets. Perhaps DRC(1, 1) is a slight exception, but it also performs much worse than the others overall (see Figure 7).

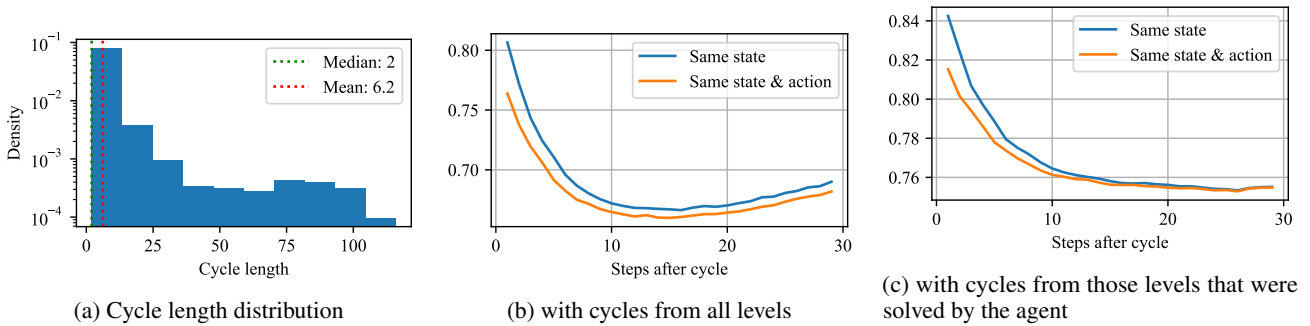


Figure 10. We replace N -length cycles with N thinking steps and checking for the same state after some timesteps. (a) A histogram of cycle lengths in the medium-validation set. (b, c) After replacing an n -length cycle with n steps of thinking, are the current state and action the same as the reference agent without extra thinking?