
An Investigation of Model-Free Planning

Arthur Guez^{*1} Mehdi Mirza^{*1} Karol Gregor^{*1} Rishabh Kabra^{*1} Sébastien Racanière¹ Théophile Weber¹
David Raposo¹ Adam Santoro¹ Laurent Orseau¹ Tom Eccles¹ Greg Wayne¹ David Silver¹
Timothy Lillicrap¹

Abstract

The field of reinforcement learning (RL) is facing increasingly challenging domains with combinatorial complexity. For an RL agent to address these challenges, it is essential that it can plan effectively. Prior work has typically utilized an explicit model of the environment, combined with a specific planning algorithm (such as tree search). More recently, a new family of methods have been proposed that learn how to plan, by providing the structure for planning via an inductive bias in the function approximator (such as a tree structured neural network), trained end-to-end by a model-free RL algorithm. In this paper, we go even further, and demonstrate empirically that an entirely model-free approach, without special structure beyond standard neural network components such as convolutional networks and LSTMs, can learn to exhibit many of the characteristics typically associated with a model-based planner. We measure our agent’s effectiveness at planning in terms of its ability to generalize across a combinatorial and irreversible state space, its data efficiency, and its ability to utilize additional thinking time. We find that our agent has many of the characteristics that one might expect to find in a planning algorithm. Furthermore, it exceeds the state-of-the-art in challenging combinatorial domains such as Sokoban and outperforms other model-free approaches that utilize strong inductive biases toward planning.

1. Introduction

One of the aspirations of artificial intelligence is a cognitive agent that can adaptively and dynamically form plans

^{*}Equal contribution ¹DeepMind, London, UK. Correspondence to: <{aguez, mmirza, rkabra, countzero}@google.com>.

to achieve its goal. Traditionally, this role has been filled by model-based RL approaches, which first learn an explicit model of the environment’s system dynamics or rules, and then apply a planning algorithm (such as tree search) to the learned model. Model-based approaches are potentially powerful but have been challenging to scale with learned models in complex and high-dimensional environments (Talvitie, 2014; Asadi et al., 2018), though there has been recent progress in that direction (Buesing et al., 2018; Ebert et al., 2018).

More recently, a variety of approaches have been proposed that learn to plan *implicitly*, solely by model-free training. These *model-free planning* agents utilize a special neural architecture that mirrors the structure of a particular planning algorithm. For example the neural network may be designed to represent search trees (Farquhar et al., 2017; Oh et al., 2017; Guez et al., 2018), forward simulations (Racanière et al., 2017; Silver et al., 2016), or dynamic programming (Tamar et al., 2016). The main idea is that, given the appropriate inductive bias for planning, the function approximator can learn to leverage these structures to learn its own planning algorithm. This kind of *algorithmic function approximation* may be more flexible than an explicit model-based approach, allowing the agent to customize the nature of planning to the specific environment.

In this paper we explore the hypothesis that planning may occur implicitly, even when the function approximator has *no special inductive bias toward planning*. Previous work (Pang & Werbos, 1998; Wang et al., 2018) has supported the idea that model-based behavior can be learned with general recurrent architectures, with planning computation amortized over multiple discrete steps (Schmidhuber, 1990), but comprehensive demonstrations of its effectiveness are still missing. Inspired by the successes of deep learning and the universality of neural representations, our main idea is simply to furnish a neural network with a high capacity and flexible representation, rather than mirror any particular planning structure. Given such flexibility, the network can in principle learn its own algorithm for approximate planning. Specifically, we utilize a family of neural networks based on a widely used function approximation architecture: the stacked convolutional LSTMs (ConvLSTM by

Xingjian et al. (2015)).

It is perhaps surprising that a purely model-free reinforcement learning approach can be so successful in domains that would appear to necessitate explicit planning. This raises a natural question: what is planning? Can a model-free RL agent be considered to be planning, without any explicit model of the environment, and without any explicit simulation of that model?

Indeed, in many definitions (Sutton et al., 1998), planning requires some explicit deliberation using a model, typically by considering possible future situations using a forward model to choose an appropriate sequence of actions. These definitions emphasize the nature of the mechanism (the explicit look-ahead), rather than the effect it produces (the foresight). However, what would one say about a deep network that has been trained from examples in a challenging domain to emulate such a planning process with near-perfect fidelity? Should a definition of planning rule out the resulting agent as effectively planning?

Instead of tying ourselves to a definition that depends on the inner workings of an agent, in this paper we take a behaviourist approach to measuring planning as a property of the agent’s interactions. In particular, we consider three key properties that an agent equipped with planning should exhibit.

First, an effective planning algorithm should be able to generalize with relative ease to different situations. The intuition here is that a simple function approximator will struggle to predict accurately across a combinatorial space of possibilities (for example the value of all chess positions), but a planning algorithm can perform a local search to dynamically compute predictions (for example by tree search). We measure this property using procedural environments (such as random gridworlds, Sokoban (Racanière et al., 2017), Boxworld (Zambaldi et al., 2018)) with a massively combinatorial space of possible layouts. We find that our model-free planning agent achieves state-of-the-art performance, and significantly outperforms more specialized model-free planning architectures. We also investigate extrapolation to a harder class of problems beyond those in the training set, and again find that our architecture performs effectively – especially with larger network sizes.

Second, a planning agent should be able to learn efficiently from relatively small amounts of data. Model-based RL is frequently motivated by the intuition that a model (for example the rules of chess) can often be learned more efficiently than direct predictions (for example the value of all chess positions). We measure this property by training our model-free planner on small data-sets, and find that our model-free planning agent still performs well and generalizes effectively to a held-out test set.

Third, an effective planning algorithm should be able to make good use of additional thinking time. Put simply, the more the algorithm thinks, the better its performance should be. This property is likely to be especially important in domains with irreversible consequences to wrong decisions (e.g. death or dead-ends). We measure this property in Sokoban by adding additional thinking time at the start of an episode, before the agent commits to a strategy, and find that our model-free planning agent solves considerably more problems.

Together, our results suggest that a model-free agent, without specific planning-inspired network structure, can learn to exhibit many of the behavioural characteristics of planning. The architecture presented in this paper serves to illustrate this point, and shows the surprising power of one simple approach. We hope our findings broaden the search for more general architectures that can tackle an even wider range of planning domains.

2. Methods

We first motivate and describe the main network architecture we use in this paper. Then we briefly explain our training setup. More details can be found in Appendix 8.

2.1. Model architectures

We desire models that can represent and learn powerful but unspecified planning procedures. Rather than encode strong inductive biases toward particular planning algorithms, we choose high-capacity neural network architectures that are capable of representing a very rich class of functions. As in many works in deep RL, we make use of convolutional neural networks (known to exploit the spatial structure inherent in visual domains) and LSTMs (known to be effective in sequential problems). Aside from these weak but common inductive biases, we keep our architecture as general and flexible as possible, and trust in standard model-free reinforcement learning algorithms to discover the capacity to plan.

2.1.1. BASIC ARCHITECTURE

The basic element of the architecture is a ConvLSTM (Xingjian et al., 2015) – a neural network similar to an LSTM but with a 3D hidden state and convolutional operations. A recurrent network f_θ stacks together ConvLSTM modules. For a stack depth of D , the state s contains all the cell states c_d and outputs h_d of each module d : $s = (c_1, \dots, c_D, h_1, \dots, h_D)$. The module weights $\theta = (\theta_1, \dots, \theta_D)$ are not shared along the stack. Given a previous state and an input tensor i , the next state is computed as $s' = f_\theta(s, i)$. The network f_θ is then repeated N times within each time-step (i.e., multiple internal ticks per

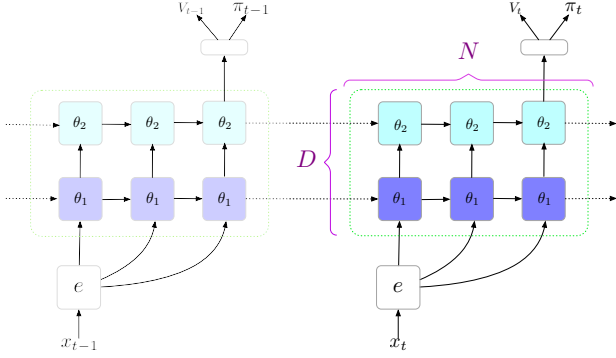


Figure 1. Illustration of the agent’s network architecture. This diagram shows DRC(2,3) for two time steps. Square boxes denote ConvLSTM modules and the rectangle box represents an MLP. Boxes with the same color share parameters.

real time-step). If s_{t-1} is the state at the end of the previous time-step, we obtain the new state given the input i_t as:

$$s_t = g_\theta(s_{t-1}, i_t) = \underbrace{f_\theta(f_\theta(\dots f_\theta(s_{t-1}, i_t), \dots, i_t), i_t)}_{N \text{ times}} \quad (1)$$

The elements of s_t all preserve the spatial dimensions of the input i_t . The final output o_t of the recurrent network for a single time-step is h_D , the hidden state of the deepest ConvLSTM module after N ticks, obtained from s_t . We describe the ConvLSTM itself and alternative choices for memory modules in Appendix 8.

The rest of the network is rather generic. An encoder network e composed of convolutional layers processes the input observation x_t into a $H \times W \times C$ tensor i_t — given as input to the recurrent module g . The encoded input i_t is also combined with o_t through a skip-connection to produce the final network output. The network output is then flattened and an action distribution π_t and a state-value v_t are computed via a fully-connected MLP. The diagram in Fig 1 illustrates the full network.

From here on, we refer to this architecture as Deep Repeated ConvLSTM (DRC) network architecture, and sometimes followed explicitly by the value of D and N (e.g., DRC(3, 2) has depth $D = 3$ and $N = 2$ repeats).

2.1.2. ADDITIONAL DETAILS

Less essential design choices in the architectures are described here. Ablation studies show that these are not crucial, but do marginally improve performance (see Appendix 11).

Encoded observation skip-connection The encoded observation i_t is provided as an input to all ConvLSTM modules

in the stack.

Top-down skip connection As described above, the flow of information in the network only goes up (and right through time). To allow for more general computation we add feedback connection from the last layer at one time step to the first layer of the next step.

Pool-and-inject To allow information to propagate faster in the spatial dimensions than the size of the convolutional kernel within the ConvLSTM stack, it is useful to provide a pooled version of the module’s last output h as an additional input on lateral connections. We use both max and mean pooling. Each pooling operation applies pooling spatially for each channel dimension, followed by a linear transform, and then tiles the result back into a 2D tensor. This operation is related to the pool-and-inject method introduced by Racanière et al. (2017) and to Squeeze-and-Excitation blocks (Hu et al., 2017).

Padding The convolutional operator is translation invariant. To help it understand where the edge of the input image is, we append a feature map to the input of the convolutional operators that has ones on the boundary and zeros inside.

2.2. Reinforcement Learning

We consider domains that are formally specified as RL problems, where agents must learn via reward feedback obtained by interacting with the environment (Sutton et al., 1998). At each time-step t , the agent’s network outputs a policy $\pi_t = \pi_\theta(\cdot|h_t)$ which maps the history of observations $h_t := (x_0, \dots, x_t)$ into a probability distribution over actions, from which the action $a_t \sim \pi_t$, is sampled. In addition, it outputs $v_t = v_\theta(h_t) \in \mathbb{R}$, an estimate of the current policy value, $v^\pi(h_t) = \mathbb{E}[G_t|h_t]$, where $G_t = \sum_{t' \geq t} \gamma^{t'-t} R_{t'}$ is the return from time t , $\gamma \leq 1$ is a discount factor, and R_t the reward at time t .

Both quantities are trained in an actor-critic setup where the policy (actor) is gradually updated to improve its expected return, and the value (critic) is used as a baseline to reduce the variance of the policy update. The update to the policy parameters have the following form using the score function estimator (a la REINFORCE (Williams, 1992)): $(g_t - v_\theta(h_t)) \nabla_\theta \log \pi_\theta(a_t|h_t)$.

In practice, we use truncated returns with bootstrapping for g_t and we apply importance sampling corrections if the trajectory data is off-policy. More specifically, we used a distributed framework and the IMPALA V-trace actor-critic algorithm (Espeholt et al., 2018). While we found this training regime to help for training networks with more parameters, we also ran experiments which demonstrate that the DRC architecture can be trained effectively with A3C (Mnih et al., 2016). More details on the setup can be found in Appendix 9.2.

3. Planning Domains

The RL domains we focus on are combinatorial domains for which episodes are procedurally generated. The procedural and combinatorial aspects emphasize planning and generalization since it is not possible to simply memorize an observation to action mapping. In these domains each episode is instantiated in a pseudorandom configuration, so solving an episode typically requires some form of reasoning. Most of the environments are fully-observable and have simple 2D visual features. The domains are illustrated and explained further in Appendix 6. In addition to the planning domains listed below, we also run control experiments on a set of Atari 2600 games (Bellemare et al., 2013).

Gridworld A simple navigation domain following (Tamar et al., 2016), consisting of a grid filled with obstacles. The agent, goal, and obstacles are randomly placed for each episode.

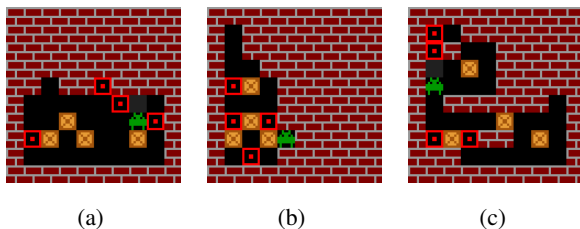


Figure 2. Examples of Sokoban levels from the (a) unfiltered, (b) medium test sets, and from the (c) hard set. Our best model is able to solve all three levels.

Sokoban A difficult puzzle domain requiring an agent to push a set of boxes onto goal locations (Botea et al., 2003; Racanière et al., 2017). Irreversible wrong moves can make the puzzle unsolvable. We describe how we generate a large number of levels (for the fixed problem size 10x10 with 4 boxes) at multiple difficulty levels in Appendix 7, and then split some into a training and test set. Briefly, problems in the first difficulty level are obtained from directly sampling a source distribution (we call that difficulty level *unfiltered*). Then the *medium* and *hard* sets are obtained by sequentially filtering that distribution based on an agent’s success on each level. We are releasing these levels as datasets in the standard Sokoban format¹. Unless otherwise specified, we ran experiments with the easier *unfiltered* set of levels.

Boxworld Introduced in (Zambaldi et al., 2018), the aim is to reach a goal target by collecting coloured keys and opening colour-matched boxes until a target is reached. The agent can see the keys (i.e., their colours) locked within boxes; thus, it must carefully plan the sequence of boxes that should be opened so that it can collect the keys that will lead to the target. Keys can only be used once, so opening

an incorrect box can lead the agent down a dead-end path from which it cannot recover.

MiniPacman (Racanière et al., 2017). The player explores a maze that contains food while being chased by ghosts. The aim of the player is to collect all the rewarding food. There are also a few power pills which allow the player to attack ghosts (for a brief duration) and earn a large reward. See Appendix 6.2 for more details.

4. Results

Paralleling our behaviourist approach to the question of planning, we look at three areas of analysis in our results. We first examine the performance of our model and other approaches across combinatorial domains that emphasize planning over memorization (Section 4.1).² We also report results aimed at understanding how elements of our architecture contribute to observed performance. Second, we examine questions of data-efficiency and generalization in Section 4.2. Third, we study evidence of iterative computation in Section 4.3.

4.1. Performance & Comparisons

In general, across all domains listed in Section 3, the DRC architecture performed very well with only modest tuning of hyper-parameters (see Appendix 9). The DRC(3,3) variant was almost always the best in terms both of data efficiency (early learning) and asymptotic performance.

Gridworld: Many methods efficiently learn the Gridworld domain, especially for small grid sizes. We found that for larger grid sizes the DRC architecture learns more efficiently than a vanilla Convolutional Neural Network (CNN) architecture of similar weight and computational capacity. We also tested Value Iteration Networks (VIN) (Tamar et al., 2016), which are specially designed to deal with this kind of problem (i.e. local transitions in a fully-observable 2D state space). We found that VIN, which has many fewer parameters and a well-matched inductive bias, starts improving faster than other methods. It outperformed the CNN and even the DRC during early-stage training, but the DRC reached better final accuracy (see Table 1 and Figure 14a in the Appendix). Concurrent to our work, Lee et al. (2018) observed similar findings in various path planning settings when comparing VIN to an architecture with weaker inductive biases.

²Illustrative videos of trained agents and playable demo available at <https://sites.google.com/view/modelfreeplanning/>

³This percentage at 1e9 is lower than the 90% reported originally by I2A (Racanière et al., 2017). This can be explained by some training differences between this paper and the I2A paper: train/test dataset vs. procedurally generated levels, co-trained vs. pre-trained models. See appendix 13.5 for more details.

¹<https://github.com/deepmind/boxoban-levels>

Table 1. Performance comparison in Gridworld, size 32x32, after 10M environment steps. VIN (Tamar et al., 2016) experiments are detailed in Appendix 13.1.

Model	% solved at 1e6 steps	% solved at 1e7 steps
DRC(3, 3)	30	99
VIN	80	97
CNN	3	90

Table 2. Comparison of test performance on (unfiltered) Sokoban levels for various methods. I2A (Racanière et al., 2017) results are re-rerun within our framework. ATreeC (Farquhar et al., 2017) experiments are detailed in Appendix 13.2. MCTSnets (Guez et al., 2018) also considered the same Sokoban domain but in an expert imitation setting (achieving 84% solved levels).

Model	% solved at 2e7 steps	% solved at 1e9 steps
DRC(3, 3)	80	99
ResNet	14	96
CNN	25	92
I2A (unroll=15)	21	84 ³
1D LSTM(3,3)	5	74
ATreeC	1	57
VIN	12	56

Sokoban: In Sokoban, we demonstrate state-of-the-art results versus prior work which targeted similar box-pushing puzzle domains (ATreeC (Farquhar et al., 2017), I2A (Racanière et al., 2017)) and other generic networks (LSTM (Hochreiter & Schmidhuber, 1997), ResNet (He et al., 2016), CNNs). We also test VIN on Sokoban, adapting the original approach to our state space by adding an input encoder to the model and an attention module at the output to deal with the imperfect state-action mappings. Table 2 compares the results for different architectures at the end of training. Only 1% of test levels remain unsolved by DRC(3,3) after 1e9 steps, with the second-best architecture (a large ResNet) failing four times as often.

Boxworld: On this domain several methods obtain near-perfect final performance. Still, the DRC model learned faster than published methods, achieving $\approx 80\%$ success after 2e8 steps. In comparison, the best ResNet achieved $\approx 50\%$ by this point. The relational method of Zambaldi et al. (2018) can learn this task well but only solved $< 10\%$ of levels after 2e8 steps.

MiniPacman: Here again, we found that the DRC architecture trained faster and obtained a better score than the ResNet architectures we tried (see Figure 15a).

Atari 2600 To test the capacity of the DRC model to deal with richer sensory data, we also examined its performance on five planning-focused Atari games (Bellemare et al.,

2013). We obtained state-of-the-art scores on three of five games, and competitive scores on the other two (see Appendix 10.2 and Figure 10 for details).

4.1.1. INFLUENCE OF NETWORK ARCHITECTURE

We studied the influence of stacking and repeating the ConvLSTM modules in the DRC architecture, controlled by the parameters D (stack length) and N (number of repeats) as described in Section 2.1. These degrees of freedom allow our networks to compute its output using shared, iterative, computation with $N > 1$, as well as computations at different levels of representation and more capacity with $D > 1$. We found that the DRC(3,3) (i.e. $D = 3, N = 3$) worked robustly across all of the tested domain. We compared this to using the same number of modules stacked without repeats (DRC(9,1)) or only repeated without stacking (DRC(1,9)). In addition, we also look at the same smaller capacity versions $D = 2, N = 2$ and $D = 1, N = 1$ (which reduces to a standard ConvLSTM). Figure 3a shows the results on Sokoban for the different network configurations. In general, the versions with more capacity performed better. When trading-off stacking and repeating (with total of 9 modules), we observed that only repeating without stacking was not as effective (this has the same number of parameters as the DRC(1,1) version), and only stacking was slower to train in the early phase but obtained a similar final performance. We also confirmed that DRC(3,3) performed better than DRC(1,1) in Boxworld, MiniPacman, and Gridworld.

On harder Sokoban levels (Medium-difficulty dataset), we trained the DRC(3,3) and the larger capacity DRC(9,1) configurations and found that, even though DRC(9,1) was slower to learn at first, it ended up reaching a better score than DRC(3,3) (94% versus 91.5% after 1e9 steps). See Fig 9 in appendix. We tested the resulting DRC(9,1) agent on the hardest Sokoban setting (Hard-difficulty), and found that it solved 80% of levels in less than 48 minutes of evaluation time. In comparison, running a powerful tree search algorithm, Levin Tree Search (Orseau et al., 2018), with a DRC(1,1) as policy prior solves 94%, but in 10 hours of evaluation.

In principle, deep feedforward models should support iterative procedures within a single time-step and perhaps match the performance of our recurrent networks (Jastrzebski et al., 2017). In practice, deep ResNets did not perform as well as our best recurrent models (see Figure 3b), and are in any case incapable of caching implicit iterative planning steps over time steps. Finally, we note that recurrence by itself was also not enough: replacing the ConvLSTM modules with flat 1D LSTMs performed poorly (see Figure 3b).

Across experiments and domains, our results suggests that both the network capacity and the iterative aspect of a model drives the agent’s performance. Moreover, in these 2D puz-

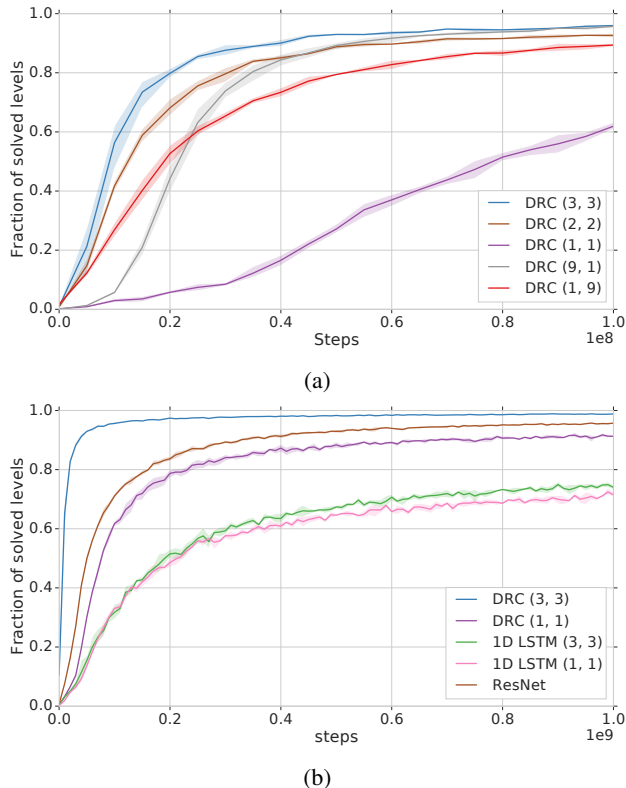


Figure 3. a) Learning curves for various configurations of DRC in Sokoban-Unfiltered. b) Comparison with other network architectures tuned for Sokoban. Results are on test-set levels.

zle domains, spatial recurrent states significantly contribute to the results.

4.2. Data Efficiency & Generalization

In combinatorial domains generalization is a central issue. Given limited exposure to configurations in an environment, how well can a model perform on unseen scenarios? In the supervised setting, large flexible networks are capable of over-fitting. Thus, one concern when using high-capacity networks is that they may over-fit to the task, for example by memorizing, rather than learning a strategy that can generalize to novel situations. Recent empirical work in SL (Supervised Learning) has shown that the generalization of large networks is not well understood (Zhang et al., 2016; Arpit et al., 2017). Generalization in RL is even less well studied, though recent work (Zhang et al., 2018a;b; Cobbe et al., 2018) has begun to explore the effect of training data diversity.

We explored two main axes in the space of generalization. We varied both the diversity of the environments as well as the size of our models. We trained the DRC architecture in various data regimes, by restricting the number of unique Sokoban levels — during the training, similar to SL, the

training algorithm iterates on those limited levels many times. We either train on a Large (900k levels), Medium-size (10k) or Small (1k) set — all subsets of the Sokoban-unfiltered training set. For each dataset size, we compared a larger version of the network, DRC(3,3), to a smaller version DRC(1,1).⁴ Results are shown in Figure 4.

In all cases, the larger DRC(3,3) network generalized better than its smaller counterpart, both in absolute terms and in terms of *generalization gap*. In particular, in the Medium-size regime, the generalization gap⁵ is 6.5% for DRC(3,3) versus 33.5% for DRC(1,1). Figures 5a-b compare these same trained models when tested on both the unfiltered and on the medium(-difficulty) test sets. We performed an analogous experiment in the Boxworld environment and observed remarkably similar results (see Figure 5c and Appendix Figure 12).

Looking across these domains and experiments there are two findings that are of particular note. First, unlike analogous SL experiments, reducing the number of training levels does not necessarily improve performance on the train set. Networks trained on 1k levels perform worse in terms of the fraction of levels solved. We believe this is due to the exploration problem in low-diversity regime: With more levels, the training agent faces a natural curriculum to help it progress toward harder levels. Another view of this is that larger networks can overfit the training levels, but only if they experience success on these levels at some point. While local minima for the loss in SL are not practically an issue with large networks, local minima in policy space can be problematic.

From a classic optimization perspective, a surprising finding is that the larger networks in our experiment (both Sokoban & Boxworld) suffer *less* from over-fitting in the low-data regime than their smaller counterparts (see Figure 5). However, this is in line with recent findings (Zhang et al., 2016) in SL that the generalization of a model is driven by the architecture and nature of the data, rather than simply as a results of the network capacity and size of the dataset. Indeed, we also trained the same networks in a purely supervised fashion through imitation learning of an expert policy.⁶ We observed a similar result when comparing the classification accuracy of the networks on the test set, with the DRC(3,3) better able to generalize — even though both networks had similar training errors on small datasets.

Extrapolation Another facet of generality in the strategy

⁴DRC(3,3) has around 300K more parameters, and it requires around 3 times more computation

⁵We compute the generalization gap by subtracting the performance (ratio of levels solved) on the training set from performance on the test set.

⁶Data was sampled on-policy from the expert policy executed on levels from the training datasets.

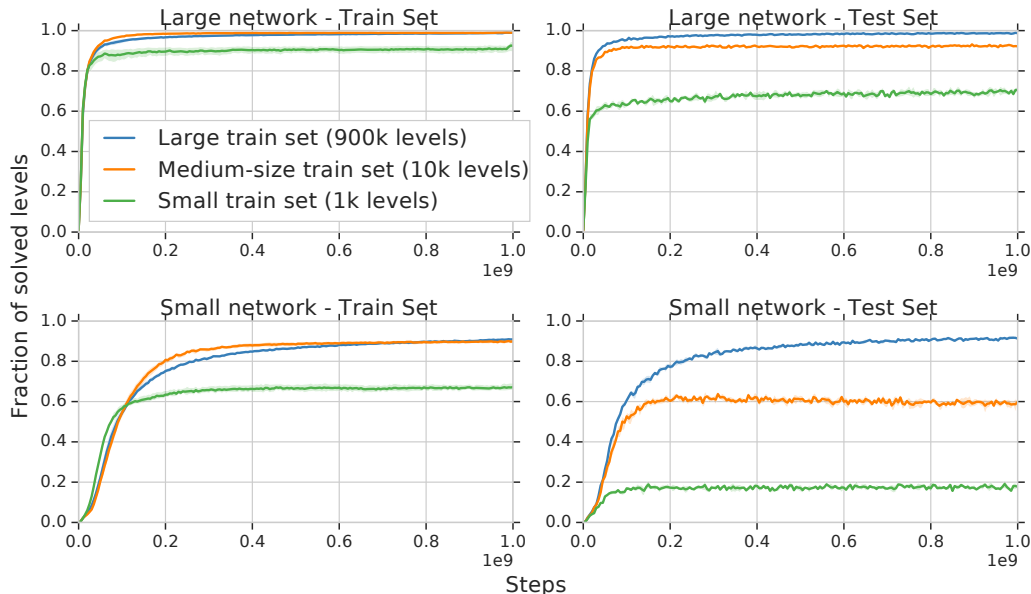


Figure 4. Comparison of DRC(3,3) (Top, Large network) and DRC(1,1) (Bottom, Small network) when trained with RL on various train set sizes (subsets of the Sokoban-unfiltered training set). Left column shows the performance on levels from the corresponding train set, right column shows the performance on the test set (the same set across these experiments).

found by the DRC network is how it performs outside the training distribution. In Sokoban, we tested the DRC(3,3) and DRC(1,1) networks on levels with a larger number of boxes than those seen in the training set. Figure 13a shows that DRC was able to extrapolate with little loss in performance to up to 7 boxes (for a fixed grid size). The performance degradation for DRC(3,3) on 7 boxes was 3.5% and 18.5% for DRC(1,1). In comparison, the results from [Racanière et al. \(2017\)](#) report a loss of 34% when extrapolating to 7 boxes in the same setup.

4.3. Iterative Computation

One desirable property for planning mechanisms is that their performance scales with additional computation without seeing new data. Although RNNs (and more recently ResNets) can in principle learn a function that can be iterated to obtain a result ([Graves, 2016](#); [Jastrzebski et al., 2017](#); [Greff et al., 2016](#)), it is not clear whether the networks trained in our RL domains learn to amortize computation over time in this way. To test this, we took trained networks in Sokoban (unfiltered) and tested post hoc their ability to improve their results with additional steps. We introduced ‘no-op’ actions at the start of each episode – up to 10 extra computation steps where the agent’s action is fixed to have no effect on the environment. The idea behind forced no-ops is to give the network more computation on the same inputs, intuitively akin to increasing its search time. Under these testing conditions, we observed clear performance improvements

on medium difficulty levels of about 5% for DRC networks (see Figure 6). We did not find such improvements for the simpler fully-connected LSTM architecture. This suggests that the DRC networks have learned a scalable strategy for the task which is computed and refined through a series of identical steps, thereby exhibiting one of the essential properties of a planning algorithm.

5. Discussion

We aspire to endow agents with the capacity to plan effectively in combinatorial domains where simple memorization of strategies is not feasible. An overarching question is regarding the nature of planning itself. Can the computations necessary for planning be learned solely using model-free RL, and can this be achieved by a general-purpose neural network with weak inductive biases? Or is it necessary to have dedicated planning machinery — either explicitly encoding existing planning algorithms, or implicitly mirroring their structure? In this paper, we studied a variety of different neural architectures trained using model-free RL in procedural planning tasks with combinatorial and irreversible state spaces. Our results suggest that general-purpose, high-capacity neural networks based on recurrent convolutional structure, are particularly efficient at learning to plan. This approach yielded state-of-the-art results on several domains – outperforming all of the specialized planning architectures that we tested. Our generalization and scaling analyses, together with the procedural nature of the studied

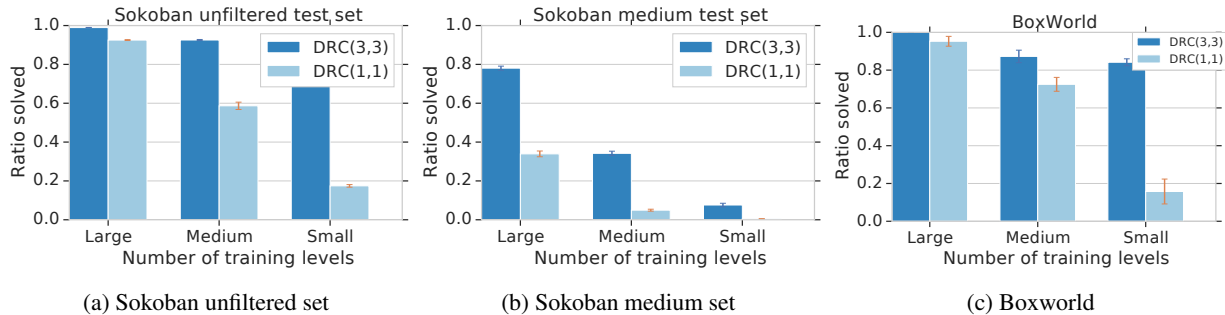


Figure 5. Generalization results from a trained model on different training set size (Large, Medium and Small subsets of the unfiltered training dataset) in Sokoban when evaluated on (a) the unfiltered test set and (b) the medium-difficulty test set. (c) Similar generalization results for trained models in Boxworld. (These figures show a summary of results in Figure 4 and Appendix Fig. 12.)

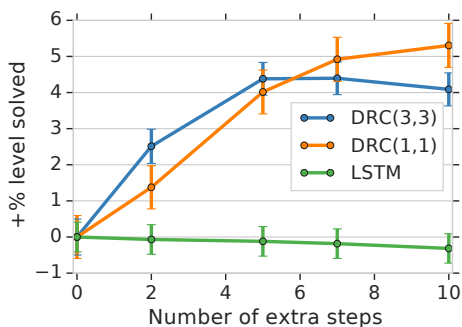


Figure 6. Forcing extra computation steps after training improves the performance of DRC on Sokoban medium-difficulty set (5 networks, each tested on the same 5000 levels). Steps are performed by overriding the policy with no-op actions at the start of an episode. The green line is the LSTM(1,1) model.

domains, suggests that these networks learn an algorithm for approximate planning that is tailored to the domain. The algorithmic function approximator appears to compute its plan dynamically, amortised over many steps, and hence additional thinking time can boost its performance.

There are, of course, many approaches to improving the efficacy of model-free algorithms. For example, DARLA, ICM, and UNREAL improve performance and transfer in RL by shaping representations using an unsupervised loss (Higgins et al., 2017; Pathak et al., 2017; Jaderberg et al., 2016). Our work hints that one of the most important approaches may be to study which inductive biases allow networks to learn effective planning-like behaviours. In principle these approaches are straightforward to combine.

Recent work in the context of supervised learning is pushing us to rethink how large neural network models generalize (Zhang et al., 2016; Arpit et al., 2017). Our results further demonstrate the mismatch between traditional views on generalisation and model size. The surprising efficacy of our planning agent, when trained on a small number of scenar-

ios across a combinatorial state space, suggests that any new theory must take into account the algorithmic function approximation capabilities of the model rather than simplistic measures of its complexity. Ultimately, we desire even more generality and scalability from our agents, and it remains to be seen whether model-free planning will be effective in reinforcement learning environments of real-world complexity.

References

- Arpit, D., Jastrzbski, S., Ballas, N., Krueger, D., Bengio, E., Kanwal, M. S., Maharaj, T., Fischer, A., Courville, A., Bengio, Y., et al. A closer look at memorization in deep networks. *arXiv preprint arXiv:1706.05394*, 2017.
- Asadi, K., Misra, D., and Littman, M. L. Lipschitz continuity in model-based reinforcement learning. *arXiv preprint arXiv:1804.07193*, 2018.
- Bellemare, M. G., Naddaf, Y., Veness, J., and Bowling, M. The arcade learning environment: An evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279, 2013.
- Botea, A., Mller, M., and Schaeffer, J. Using abstraction for planning in sokoban. In *Computers and Games*, volume 2883, pp. 360, 2003.
- Buesing, L., Weber, T., Racaniere, S., Eslami, S., Rezende, D., Reichert, D. P., Viola, F., Besse, F., Gregor, K., Hassabis, D., et al. Learning and querying fast generative models for reinforcement learning. *arXiv preprint arXiv:1802.03006*, 2018.
- Cobbe, K., Klimov, O., Hesse, C., Kim, T., and Schulman, J. Quantifying generalization in reinforcement learning. *arXiv preprint arXiv:1812.02341*, 2018.
- Ebert, F., Finn, C., Dasari, S., Xie, A., Lee, A., and Levine, S. Visual foresight: Model-based deep reinforcement

- learning for vision-based robotic control. *arXiv preprint arXiv:1812.00568*, 2018.
- Espeholt, L., Soyer, H., Munos, R., Simonyan, K., Mnih, V., Ward, T., Doron, Y., Firoiu, V., Harley, T., Dunning, I., et al. Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures. *arXiv preprint arXiv:1802.01561*, 2018.
- Farquhar, G., Rocktäschel, T., Igl, M., and Whiteson, S. Treeqn and atrec: Differentiable tree planning for deep reinforcement learning. In *ICLR*, 2017.
- Graves, A. Adaptive computation time for recurrent neural networks. *arXiv preprint arXiv:1603.08983*, 2016.
- Greff, K., Srivastava, R. K., and Schmidhuber, J. Highway and residual networks learn unrolled iterative estimation. *arXiv preprint arXiv:1612.07771*, 2016.
- Guez, A., Weber, T., Antonoglou, I., Simonyan, K., Vinyals, O., Wierstra, D., Munos, R., and Silver, D. Learning to search with mctsnets. *arXiv preprint arXiv:1802.04697*, 2018.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Higgins, I., Pal, A., Rusu, A., Matthey, L., Burgess, C., Pritzel, A., Botvinick, M., Blundell, C., and Lerchner, A. Darla: Improving zero-shot transfer in reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 1480–1490. JMLR. org, 2017.
- Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- Horgan, D., Quan, J., Budden, D., Barth-Maron, G., Hessel, M., Van Hasselt, H., and Silver, D. Distributed prioritized experience replay. *arXiv preprint arXiv:1803.00933*, 2018.
- Hu, J., Shen, L., and Sun, G. Squeeze-and-excitation networks. *arXiv preprint arXiv:1709.01507*, 7, 2017.
- Jaderberg, M., Mnih, V., Czarnecki, W. M., Schaul, T., Leibo, J. Z., Silver, D., and Kavukcuoglu, K. Reinforcement learning with unsupervised auxiliary tasks. *arXiv preprint arXiv:1611.05397*, 2016.
- Jastrzebski, S., Arpit, D., Ballas, N., Verma, V., Che, T., and Bengio, Y. Residual connections encourage iterative inference. *arXiv preprint arXiv:1710.04773*, 2017.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Lee, L., Parisotto, E., Chaplot, D. S., Xing, E., and Salakhutdinov, R. Gated path planning networks. *arXiv preprint arXiv:1806.06408*, 2018.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540): 529, 2015.
- Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., and Kavukcuoglu, K. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pp. 1928–1937, 2016.
- Oh, J., Singh, S., and Lee, H. Value prediction network. In *Advances in Neural Information Processing Systems*, pp. 6118–6128, 2017.
- Orseau, L., Lelis, L., Lattimore, T., and Weber, T. Single-agent policy tree search with guarantees. In *Advances in Neural Information Processing Systems*, 2018.
- Pang, X. and Werbos, P. Neural network design for j function approximation in dynamic programming. *arXiv preprint adap-org/9806001*, 1998.
- Pathak, D., Agrawal, P., Efros, A. A., and Darrell, T. Curiosity-driven exploration by self-supervised prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 16–17, 2017.
- Racanière, S., Weber, T., Reichert, D., Buesing, L., Guez, A., Rezende, D. J., Badia, A. P., Vinyals, O., Heess, N., Li, Y., et al. Imagination-augmented agents for deep reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 5690–5701, 2017.
- Schmidhuber, J. Making the world differentiable: On using self-supervised fully recurrent neural networks for dynamic reinforcement learning and planning in non-stationary environments. 1990.
- Silver, D., van Hasselt, H., Hessel, M., Schaul, T., Guez, A., Harley, T., Dulac-Arnold, G., Reichert, D., Rabinowitz, N., Barreto, A., et al. The predictron: End-to-end learning and planning. *arXiv preprint arXiv:1612.08810*, 2016.
- Sutton, R. S., Barto, A. G., et al. *Reinforcement learning: An introduction*. 1998.
- Talvitie, E. Model regularization for stable sample rollouts. In *Proceedings of the Thirtieth Conference on Uncertainty in Artificial Intelligence*, pp. 780–789. AUAI Press, 2014.

- Tamar, A., Wu, Y., Thomas, G., Levine, S., and Abbeel, P. Value iteration networks. In *Advances in Neural Information Processing Systems*, pp. 2154–2162, 2016.
- Wang, J. X., Kurth-Nelson, Z., Kumaran, D., Tirumala, D., Soyer, H., Leibo, J. Z., Hassabis, D., and Botvinick, M. Prefrontal cortex as a meta-reinforcement learning system. *Nature neuroscience*, 21(6):860, 2018.
- Williams, R. J. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- Xingjian, S., Chen, Z., Wang, H., Yeung, D.-Y., Wong, W.-K., and Woo, W.-c. Convolutional lstm network: A machine learning approach for precipitation nowcasting. In *Advances in neural information processing systems*, pp. 802–810, 2015.
- Zambaldi, V., Raposo, D., Santoro, A., Bapst, V., Li, Y., Babuschkin, I., Tuyls, K., Reichert, D., Lillicrap, T., Lockhart, E., et al. Relational deep reinforcement learning. *arXiv preprint arXiv:1806.01830*, 2018.
- Zhang, A., Ballas, N., and Pineau, J. A dissection of overfitting and generalization in continuous reinforcement learning. *arXiv preprint arXiv:1806.07937*, 2018a.
- Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. Understanding deep learning requires rethinking generalization. *arXiv preprint arXiv:1611.03530*, 2016.
- Zhang, C., Vinyals, O., Munos, R., and Bengio, S. A study on overfitting in deep reinforcement learning. *arXiv preprint arXiv:1804.06893*, 2018b.