

---

# Subgoal Search For Complex Reasoning Tasks

---

**Konrad Czechowski\***  
University of Warsaw  
k.czechowski@mimuw.edu.pl

**Tomasz Odrzygóźdź\***  
University of Warsaw  
tomaszo@impan.pl

**Marek Zbysiński**  
University of Warsaw  
m.zbysinski@students.mimuw.edu.pl

**Michał Zawalski**  
University of Warsaw  
m.zawalski@uw.edu.pl

**Krzysztof Olejnik**  
University of Warsaw  
k.olejnik3@student.uw.edu.pl

**Yuhuai Wu**  
University of Toronto,  
Vector Institute  
ywu@cs.toronto.edu

**Łukasz Kuciński**  
Polish Academy of Sciences  
lkucinski@impan.pl

**Piotr Miłoś**  
Polish Academy of Sciences,  
University of Oxford,  
deepsense.ai  
pmilos@impan.pl

## Abstract

Humans excel in solving complex reasoning tasks through a mental process of moving from one idea to a related one. Inspired by this, we propose Subgoal Search (kSubS) method. Its key component is a learned subgoal generator that produces a diversity of subgoals that are both achievable and closer to the solution. Using subgoals reduces the search space and induces a high-level search graph suitable for efficient planning. In this paper, we implement kSubS using a transformer-based subgoal module coupled with the classical best-first search framework. We show that a simple approach of generating  $k$ -th step ahead subgoals is surprisingly efficient on three challenging domains: two popular puzzle games, Sokoban and the Rubik’s Cube, and an inequality proving benchmark INT. kSubS achieves strong results including state-of-the-art on INT within a modest computational budget.

## 1 Introduction

Reasoning is often regarded as a defining property of advanced intelligence [39, 18]. When confronted with a complicated task, humans’ thinking process often moves from one idea to a related idea, and the progress is made through milestones, or *subgoals*, rather than through atomic actions that are necessary to transition between subgoals [15]. During this process, thinking about one subgoal can lead to a possibly diverse set of subsequent subgoals that are conceptually reachable and make a promising step towards the problem’s solution. This intuitive introspection is backed by neuroscience evidence [20], and in this work, we present an algorithm that mimics this process. Our approach couples a deep learning generative subgoal modeling with classical search algorithms to allow for successful planning with subgoals. We showcase the efficiency of our method on the following complex reasoning tasks: two popular puzzle games Sokoban and the Rubik’s Cube, and an inequality theorem proving benchmark INT [54], achieving the state-of-the-art results in INT and competitive results for the remaining two.

The deep learning revolution has brought spectacular advancements in pattern recognition techniques and models. Given the hard nature of reasoning problems, these are natural candidates to provide

---

\*equal contribution

search heuristics [4]. Indeed, such a blend can produce impressive results [43, 44, 36, 1]. These approaches seek solutions using elementary actions. Others, e.g. [29, 33, 23], utilize variational subgoals generators to deal with long-horizon visual tasks. We show that these ideas can be pushed further to provide algorithms capable of dealing with combinatorial complexity.

We present Subgoal Search (kSubS) method and give its practical implementations: MCTS-kSubS and BF-kSubS. kSubS consists of the following four components: planner, subgoal generator, a low-level policy, and a value function. The planner is used to search over the graph induced by the subgoal generator and is guided by the value function. The role of the low-level policy is to prune the search tree as well as to transition between subgoals. In this paper, we assume that the generator predicts subgoals that are  $k$  step ahead (towards the solution) from the current state, and to emphasize this we henceforth add  $k$  to the method’s abbreviation. MCTS-kSubS and BF-kSubS differ in the choice of the search engine: the former uses Monte-Carlo Tree Search (MCTS), while the latter is backed by Best-First Search (BestFS). We provide two sets of implementations for the generator, the low-level policy, and the value functions. The first one uses transformer architecture [48] for each component, while the second utilizes a convolutional network for the generator and the value function, and the classical breadth-first search for the low-level policy. This lets us showcase the versatility and effectiveness of the approach.

The subgoal generator lies at the very heart of Subgoal Search, being an implementation of reasoning with high-level ideas. To be useful in a broad spectrum of contexts, the generator should be implemented as a learnable (generative) model. As a result, it is expected to be imperfect and (sometimes) generate incorrect predictions, which may turn the search procedure invalid. Can we thus make planning with learned subgoals work? In this paper, we answer this question affirmatively: we show that the autoregressive framework of transformer-based neural network architecture [49] leads to superior results in challenging domains.

We train the transformer with the objective to predict the  $k$ -th step ahead. The main advantages of this subgoal objective are simplicity and empirical efficiency. We used expert data to generate labels for supervised training. When offline datasets are available, which is the case for the environments considered in this paper<sup>2</sup>, such an approach allows for stable and efficient optimization with high-quality gradients. Consequently, this method is often taken when dealing with complex domains (see e.g. [41, 51]) or when only an offline expert is available<sup>3</sup>. Furthermore, we found evidence of out-of-distribution generalization.

Finally, we formulate the following hypothesis aiming to shed some light on why kSubS is successful: we speculate that subgoal generation may alleviate errors in the value function estimation. Planning methods based on learning, including kSubS, typically use imperfect value function-based information to guide the search. While traditional low-level search methods are susceptible to local noise, subgoal generation allows for evaluations of the value functions at temporally distant subgoals, which improves the signal-to-noise ratio and allows a “leap over” the noise.

To sum up, our contributions are:

1. We propose Subgoal Search method with two implementations: MCTS-kSubS, BF-kSubS. We demonstrate that our approach requires a relatively little search or, equivalently, is able to handle bigger problems. We also observe evidence of out-of-distribution generalization.
2. We provide evidence that a transformer-based autoregressive model learned with a simple supervised objective to predict states  $k$ -th step ahead is an effective tool to generate valid and diverse subgoals.
3. We show in our experiments that using subgoal planning help to might mitigate the negative influence of value function errors on planning.

We provide the code of our method and experiment settings at <https://github.com/subgoal-search/subgoal-search>, and a dedicated website <https://sites.google.com/view/subgoal-search>.

---

<sup>2</sup>The dataset for INT or Sokoban can be easily generated or are publicly available. For the Rubik’s Cube, we use random data or simple heuristic (random data are often sufficient for robotic tasks and navigation.)

<sup>3</sup>For example, the INT engine can easily generate multiple proves of random statements, but *cannot* prove a given theorem.

## 2 Related work

In classical AI, reasoning is often achieved by *search* ([39]). Search rarely can be exhaustive, and a large body of algorithms and heuristics has been developed over the years, [39, Section 3.5]. It is hypothesized that progress can be achieved by combining search with learning [4]. Among notable successful examples of this approach are Alpha Zero [42], or solving Rubik’s cube using a learned heuristic function to guide the A\* algorithm (see [1]).

An eminent early example of using goal-directed reasoning is the PARADISE algorithm ([52]). In deep learning literature, [24] was perhaps the first work implementing subgoal planning. This was followed by a large body of work on planning with subgoals in the latent space for visual tasks ([23, 30, 33, 29, 21, 9, 34]) or landmark-based navigation methods ([40, 26, 14, 45, 55]).

The tasks considered in the aforementioned studies are often quite forgiving when it comes to small errors in the subgoal generation. This can be contrasted with complex reasoning domains, in which even a small variation of a state may drastically change its meaning or render it invalid. Thus, neural networks may struggle to generate semantically valid states ([4, Section 6.1]).

Assuming that this problem was solved, a generated subgoal still remains to be assessed. The exact evaluation may, in general, require exhaustive search or access to an oracle (in which case the original problem is essentially solved). Consequently, it is unlikely that a simple planner (e.g., one unrolling independent sequences of subgoals [9]) will either produce an informative outcome, could be easily improved using only local changes via gradient descent [30], or cross-entropy method (CEM) [29, 33, 34]. Existing approaches based, which are based on more powerful subgoal search methods, have their limitations, on the other hand. [13] is perhaps the closest to our method and uses MCTS to search the subgoal-induced graph. However, it uses a predefined (not learned) predicate function as a subgoal generator, limiting applicability to the problems with available high-quality heuristics. Learned subgoals are used in [31], a method of hierarchical planning. That said, the subgoal space needs to be relatively small for this method to work (or crafted manually to reduce cardinality). To the best of our knowledge, our algorithm is the first domain agnostic hierarchical planner for combinatorially complex domains.

More generally, concepts related to goals and subgoals percolated to reinforcement learning early on, leading, among others, to prominent ideas like hindsight [22], hierarchical learning [47, 7] or the Horde architecture [46]. Recently, with the advent of *deep* reinforcement learning, these ideas have been resurfacing and scaled up to deliver their initial promises. For example, [50] implements ideas of [7] and a very successful hindsight technique [2] is already considered to be a core RL method. Further, a recent paper [35] utilizes a maximum entropy objective to select achievable goals in hindsight training. Apart from the aforementioned hierarchical planning, the idea to use neural networks for subgoal generation was used to improve model-free reinforcement learning agents [11, 6] and imitation learning algorithms [32].

## 3 Method

Our method, Subgoal Search (kSubS), is designed for problems, which can be formulated as a search over a graph with a known transition model. Formally, let  $G = (\mathcal{S}, \mathcal{E})$  be a directed graph and  $\tilde{\mathcal{S}} \subset \mathcal{S}$  be the set of success states. We assume that, during the solving process, the algorithm can, for a given node  $g$ , determine the edges starting at  $g$  and check if  $g \in \tilde{\mathcal{S}}$ .

Subgoal Search consists of four components: planner, subgoal generator, low-level policy, and a value function. The planner, coupled with a value function, is used to search over the graph induced by the subgoal generator. Namely, for each selected subgoal, the generator allows for sampling the candidates for the next subgoals. Only these reachable by the low-level policy are used. The procedure continues until the solution is found or the computational budget is exhausted. Our method searches over a graph  $\tilde{G} = (\mathcal{S}, \mathcal{E}_s)$ , with the edges  $\mathcal{E}_s$  given by the subgoal generator. Provided a reasonable generator, paths to  $\tilde{\mathcal{S}}$  are shorter in  $\tilde{G}$  than in  $G$  and thus easier to find.

In this paper we provide BestFS- and MCTS- backed implementations kSubS. Algorithm 1 presents BF-kSubS; see Algorithm 4 in Appendix A.1 for MCTS-kSubS.

For INT and Rubik’s Cube, we use transformer models (see Appendix B.1) in all components other than the planner. For Sokoban, we use convolutional networks (see Appendix B.2). While transformers could also be used in Sokoban, we show that a simplified setup already achieves strong results. This showcases that Subgoal Search is general enough to work with different design choices. In Section 4.2 we describe datasets used train these neural models.

Algorithm 1 Best-First Subgoal Search (BF-kSubS)	Algorithm 2 Low-level conditional policy
<p><b>Require:</b> <math>C_1</math> max number of nodes  <math>V</math> value function network  SOLVED predicate of solution</p> <p><b>function</b> SOLVE(<math>s_0</math>)  T.PUSH(<math>(V(s_0), s_0)</math>) <math>\triangleright</math> T is priority queue  paths[<math>s_0</math>] <math>\leftarrow</math> [] <math>\triangleright</math> paths is dict of lists  seen.ADD(<math>s_0</math>) <math>\triangleright</math> seen is set  <b>while</b> <math>0 &lt; \text{LEN}(T)</math> and <math>\text{LEN}(\text{seen}) &lt; C_1</math> <b>do</b>  <math>-, s \leftarrow</math> T.EXTRACT_MAX()  subgoals <math>\leftarrow</math> SUB_GENERATE(<math>s</math>)  <math>\triangleright</math> see Algorithm 3  <b>for</b> <math>s'</math> <b>in</b> subgoals <b>do</b>  <b>if</b> <math>s'</math> <b>in</b> seen <b>then</b> continue  seen.ADD(<math>s'</math>)  actions <math>\leftarrow</math> GET_PATH(<math>s, s'</math>)  <math>\triangleright</math> see Alg 2 or Alg 9  <b>if</b> actions.EMPTY() <b>then</b>  continue  T.PUSH(<math>(V(s'), s')</math>)  paths[<math>s'</math>] <math>\leftarrow</math> paths[<math>s</math>] + actions  <b>if</b> SOLVED(<math>s'</math>) <b>then</b>  <b>return</b> paths[<math>s'</math>]</p> <p><b>return</b> False</p>	<p><b>Require:</b> <math>C_2</math> steps limit  <math>\pi</math> low-level conditional policy network  <math>M</math> model of the environment</p> <p><b>function</b> GET_PATH(<math>s_0</math>, subgoal)  step <math>\leftarrow</math> 0  <math>s \leftarrow s_0</math>  action_path <math>\leftarrow</math> []  <b>while</b> step <math>&lt; C_2</math> <b>do</b>  action <math>\leftarrow</math> <math>\pi</math>.PREDICT(<math>s</math>, subgoal)  action_path.APPEND(action)  <math>s \leftarrow M</math>.NEXT_STATE(<math>s</math>, action)  <b>if</b> <math>s = \text{subgoal}</math> <b>then</b>  <b>return</b> action_path  step <math>\leftarrow</math> step + 1  <b>return</b> []</p>

**Subgoal generator** Formally, it is a mapping  $\rho: \mathcal{S} \rightarrow \mathbf{P}(\mathcal{S})$ , where  $\mathbf{P}(\mathcal{S})$  is a family of probability distributions over the environment’s state space  $\mathcal{S}$ . More precisely, let us define a trajectory as a sequence of state-action pairs  $(s_0, a_0), (s_1, a_1), \dots, (s_n, a_n)$ , with  $(s_i, a_i) \in \mathcal{S} \times \mathcal{A}$ , where  $\mathcal{A}$  stands for the action space and  $n$  is the trajectory’s length. We will say that the generator predicts  $k$ -step ahead subgoals, if at any state  $s_\ell$  it aims to predict  $s_{\min(\ell+k, n)}$ . We show, perhaps surprisingly, that this simple objective is an efficient way to improve planning, even for small values of  $k$ , i.e.  $k \in \{2, 3, 4, 5\}$ .

Operationally, the subgoal generator takes as input an element of the  $\mathbf{s} \in \mathcal{S}$  and returns *subgoals*, a set of new candidate states expected to be closer to the solution and is implemented by Algorithm 3. It works as follows: first, we generate  $C_3$  subgoal candidates with SUB\_NET\_GENERATE. Then, we prune this set to obtain a total probability greater than  $C_4$ .

For INT and Rubik’s Cube, we represent states as sequences modeled with a transformer. Following the practice routinely used in language modeling, [48], SUB\_NET\_GENERATE employs beam search to generate a set of high likelihood outputs.<sup>4</sup> With the same goal in mind, for Sokoban, we use another method described Appendix C.1. SUB\_NET\_GENERATE uses the subgoal generator network, which is trained in a supervised way on the expert data, with training examples being:  $s_\ell$  (input) and  $s_{\min(\ell+k, n)}$  (label), see Appendix D for details.

**Low-level conditional policy** Formally, it is a mapping  $\pi: \mathcal{S} \times \mathcal{S} \rightarrow \mathcal{A}^*$ . It is used to generate a sequence of actions on how to reach a subgoal starting from a given initial state. Operationally, it may return an empty sequence if the subgoal cannot be reached within  $C_2$  steps, see Algorithm 2. This is used as a pruning mechanism for the *subgoals* set in Algorithm 1.

<sup>4</sup>In language modeling, typically, only one beam search output is used. In our case, however, we utilize all of them, which turns out to be a diverse set of subgoals.

In INT and Rubik’s Cube, we use Algorithm 2, which utilizes low-level policy network  $\pi$ . Similarly to the subgoal generator, it is trained using expert data in a supervised fashion, i.e. for a pair  $(s_\ell, s_{\min(\ell+i, n)})$ , with  $i \leq k$ , its objective is to predict  $a_\ell$ .

When the branching factor is small, the low-level policy can be realized by a simple breadth-first search mechanism, see Algorithm 9 in Appendix, which we illustrate on Sokoban.

**Value function** Formally, it is a mapping  $V: \mathcal{S} \rightarrow \mathbb{R}$ , that assigns to each state a value related to its distance to the solution, and it is used to guide the search (see Algorithm 1 and Algorithm 4). For its training, we use expert data. For each state  $s_\ell$  the training target is negated distance to the goal:  $\ell - n$ , where  $n$  denotes the end step of a trajectory that  $s_\ell$  belongs to.

**Planner** This is the engine that we use to search the subgoal-induced graph.

In this paper, we use BestFS (Algorithm 1) and MCTS (Algorithm 4 in Appendix A.1). The former is a classic planning method, which maintains a priority queue of states waiting to be explored, and greedily (with respect to their value) selects elements from it (see, e.g., [39]). MCTS is a search method that iteratively and explicitly builds a search tree, using (and updating) the collected node statistics (see, e.g., [5]). In this paper, we use an AlphaZero-like [41] algorithm for single-player games.

We note that the subgoal generator can be combined with virtually any search algorithm and can benefit from an additional structure. For example, for domains providing a factored state representation, the width-based methods [25, 12] would likely be stronger search mechanisms.

## 4 Experiments

In this section, we empirically demonstrate the efficiency of MCTS-kSubS and BF-kSubS. In particular, we show that they vastly outperform their standard (“non-subgoal”) counterparts. As a testing ground, we consider three challenging domains: Sokoban, Rubik’s Cube, and INT. All of them require non-trivial reasoning. The Rubik’s Cube is a well-known 3-D combination puzzle. Sokoban is a complex video puzzle game known to be NP-hard and thus challenging for planning methods. INT [54] is a recent theorem proving benchmark.

### 4.1 Training protocol and baselines

Our experiments consist of three stages. First, we collect domain-specific expert data, see Section 4.2. Secondly, we train the subgoal generator, low-level conditional policy, and value function networks using the data and targets described in Section 3. For more details see Appendix D. Eventually, we evaluate the planning performance of MCTS-kSubS and BF-kSubS, details of which are presented below. In the second step, we use supervised learning, which makes our setup stable with respect to network initialization, see details in Appendix D.1.3.

As baselines, we use BestFS and MCTS (being a single-player implementation of AlphaZero). In INT and Rubik’s Cube, both the algorithms utilize policy networks (trained with behavioral cloning, on the same dataset, which we used to train kSubS). Note that distribution over actions induces a distribution over states; thus the policy network can be regarded as a subgoal generator for  $k = 1$ . More details about the baselines can be found in Appendix I.

---

### Algorithm 3 Subgoal generator

---

**Require:**  $C_3$  number of subgoals  
 $C_4$  target probability  
 $\rho$  subgoal generator network

**function** SUB\_GENERATE( $s$ )  
subgoals  $\leftarrow \emptyset$   
states, probs  $\leftarrow$  SUB\_NET\_GENERATE( $\rho, s; C_3$ )  
 $\triangleright$  (states, probs) is sorted wrt probs  
total\_p  $\leftarrow 0$   
**for** state, p  $\in$  (states, probs) **do**  
  **if** total\_p  $> C_4$  **then break**  
  subgoals.ADD(state)  
  total\_p  $\leftarrow$  total\_p + p  
**return** subgoals

---

## 4.2 Search Domains and Datasets

**Sokoban** is a single-player complex game in which a player controls an agent whose goal is to place boxes on target locations solely by pushing them; without crossing any obstacles or walls. Sokoban has recently been used to test the boundaries in RL [16, 28]. Sokoban is known to be hard [10], mainly due to its combinatorial complexity and the existence of irreversible states. Deciding if a given Sokoban board is solvable is an NP-hard problem [8].

We collect the expert dataset consisting of all successful trajectories occurring during the training of an MCTS agent (using an implementation of [28]). These are suboptimal, especially in the early phases of the training or for harder boards. For both expert training and kSubS evaluation, we generate Sokoban boards following the approach of [37].

**Rubik’s Cube** is a classical 3-dimensional puzzle. It is considered challenging due to the fact that the search space has more than  $4.3 \times 10^{18}$  configurations. Similarly to Sokoban, Rubik’s Cube has been recently used as a testing ground for RL methods [1].

To collect the expert dataset, we generate random paths of length 30 starting from the solved cube and take them backward. These backward solutions are highly sub-optimal (optimal solutions are proven to be shorter than 26 [38]).

**INT: Inequality Benchmark for Theorem Proving.** INT provides a generator of inequalities, which produces mathematical formulas along with their proofs, see [54, Section 3.3]. Proofs are represented as sequences of consecutively applied mathematical axioms (there are  $K = 18$  axioms in total). An action in INT is a tuple containing an axiom and its input entities. The action space in this problem can be very large, reaching up to  $10^6$  elements, which significantly complicates planning.

The INT generator constructs paths by randomly applying axioms starting with a trivial statement. Such a path taken backward constitutes the proof of its final statement (not guaranteed to be optimal). The proof length, denoted by  $L$ , is an important hyperparameter regulating the difficulty – we use 5, 10, 15.

For more details on datasets see Appendix C.

## 4.3 Main results

In this section, we present our most important finding: Subgoal Search enables for more efficient search and consequently scales up to problems with higher difficulty. Specifically, MCTS-kSubS and BF-kSubS perform significantly better than the respective methods not using subgoals, including state-of-the-art on INT.

In Figure 1, we present the performance of Subgoal Search. We measure the *success rate* as a function of the *search budget*. The success rate is measured on 1000 instances of a given problem (which results in confidence intervals within  $\pm 0.03$ ). For BF-kSubS the search budget is referred to as *graph size* and includes the number of nodes visited by Algorithm 1. For INT and Rubik’s Cube, we include both the subgoal generated by SUB\_GENERATE and the nodes visited by GET\_PATH (as they induce a significant computational cost stemming from using low-level policy  $\pi$  in Algorithm 2). For Sokoban, we use Algorithm 9 to realize GET\_PATH, as it has a negligible cost (less than 1% of the total runtime of Algorithm 1), we do not include these nodes into graph size.

For MCTS, we report *MCTS passes*, which is a common metric for MCTS, see details in Appendix A.1.

Below we discuss the results separately for each domain. We provide examples of solutions and generated subgoals in Appendix H.

**INT** The difficulty of the problems in INT increases fast with the proof length  $L$  and the number of accessible axioms. We used  $K = 18$ ; all of available axioms. We observe, that BF-kSubS scales to proofs of length  $L = 10$  and  $L = 15$ , which are significantly harder than  $L = 5$  considered in [54], see Table 2. The same holds for MCTS-kSubS, see Appendix A.2.

We check also that MCTS-kSubS vastly outperforms the baseline - AlphaZero algorithm, see Figure 1 (top, left). An MCTS-based agent was also evaluated in [54]. Its implementation uses graph neural

	INT	Sokoban	Rubik
$k$	3	4	4
$C_1$	400	5000	1500
$C_2$	4	4	7
$C_3$	4	4	3
$C_4$	1	0.98	1

Table 1: BF-kSubS hyperparameters.

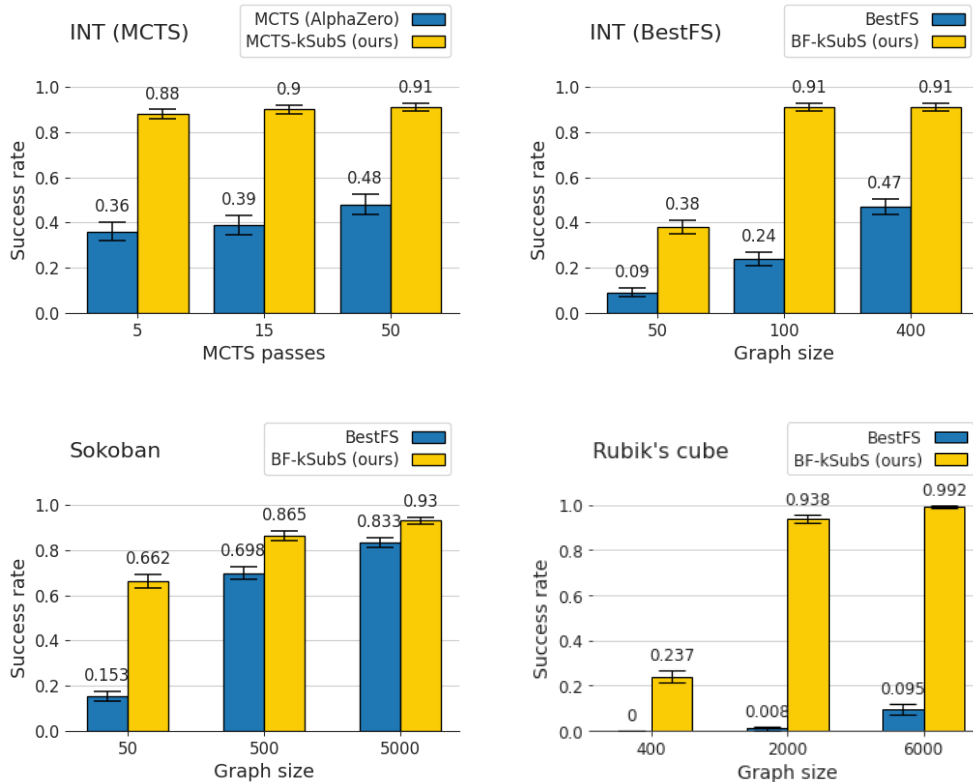


Figure 1: The performance of Subgoal Search. (top, left) comparison on INT (with the proof length 15) to AlphaZero. (top, right) BF-kSubS consistently achieves high performance even for small computational budgets. (bottom, left) similarly on Sokoban (board size 12x12 with 4 boxes) the advantage of BF-kSubS is clearly visible for small budget. (bottom, right) BestFS fails to solve Rubik’s Cube, while BF-kSubS can achieve near-perfect performance.

Proof length		5		10		15	
Method		BestFS	BF-kSubS (ours)	BestFS	BF-kSubS (ours)	BestFS	BF-kSubS (ours)
Graph size	50	0.82	<b>0.99</b>	0.47	<b>0.97</b>	0.09	<b>0.38</b>
	100	0.89	<b>0.99</b>	0.64	<b>0.99</b>	0.24	<b>0.91</b>
	200	0.92	<b>0.99</b>	0.67	<b>0.99</b>	0.35	<b>0.91</b>
	400	0.93	<b>0.99</b>	0.72	<b>0.99</b>	0.47	<b>0.91</b>

Table 2: INT success rates for various proof lengths and graphs sizes.

networks architectures and achieves 92% success rate for  $L = 5$ . Our transformed-based baseline is stronger - it solves over 99% instances on the same problem set.

**Sokoban** Using BF-kSubS allows for significantly higher success rates within the same computational budget, see Table 3. Our solution scales well to the board size as big as  $20 \times 20$ ; note that  $10 \times 10$  boards are typically used in deep RL research [16, 37]. Importantly, we observe that already for a small computational budget (graph size 1000) BF-kSubS obtains higher success rates than the expert we used to create the dataset (these are 78%, 67%, 60% for the respective board sizes).

We also tested how the quality of BF-kSubS depends on the size of the training dataset for Sokoban, the results can be found in Appendix F.

**Rubik’s Cube** BF-kSubS solves nearly 100% of cubes, BestFS solve less than 10%, see Figure 1 (bottom, right). This is perhaps the most striking example of the advantage of using a subgoal generator instead of low-level actions. We present possible explanation in Appendix K.

Board size		12 x 12		16 x 16		20 x 20	
Method		BestFS	BF-kSubS (ours)	BestFS	BF-kSubS (ours)	BestFS	BF-kSubS (ours)
Graph size	50	0.15	<b>0.66</b>	0.04	<b>0.42</b>	0.02	<b>0.46</b>
	100	0.46	<b>0.79</b>	0.23	<b>0.62</b>	0.10	<b>0.55</b>
	1000	0.75	<b>0.89</b>	0.61	<b>0.79</b>	0.47	<b>0.70</b>
	5000	0.83	<b>0.93</b>	0.69	<b>0.85</b>	0.58	<b>0.77</b>

Table 3: Sokoban success rates for various board sizes (each with 4 boxes).

**Out-of-distribution (OOD) generalization** OOD generalization is considered to be the crucial ability to make progress in hard combinatorial optimization problems [4] and automated theorem proving [54]. The INT inequality generator has been specifically designed to benchmark this phenomenon. We check that Subgoal Search trained on proofs on length 10 generalizes favourably to longer problems, see Figure 4. Following [54], we speculate that search is a computational mechanism that delivers OOD generalization.

It might be hard to compare computational budgets between various algorithms and their versions. In Appendix E we measure that BF-kSubS and MCTS-kSubS offer very practical benefits, sometimes as much as  $7\times$  faster execution.

#### 4.4 Analysis of $k$ (subgoal distance) parameter

The subgoals are trained to predict states  $k$  steps ahead of the current one. Higher  $k$  should make planning easier as the search graph is smaller. However, as  $k$  increases, the quality of the generator may drop, and thus the overall effect is uncertain. Similarly, the task of the low-level conditional policy becomes more difficult as  $k$  increases. The optimal value of  $k$  is 3 and 4 for INT and Rubik’s Cube, respectively. In these environments, increasing  $k$  further degrades performance. In Sokoban, we observe monotonic improvement up to  $k = 10$ . This is perhaps because low-level conditional policy (Algorithm 9, based on breadth-first search) never fails to fill the path from a state to the valid subgoal. The running cost of Algorithm 9 quickly becomes unacceptable (recall that for  $k = 4$ , which we used in the main experiment, it has still a negligible cost - below  $< 1\%$  of the total runtime).

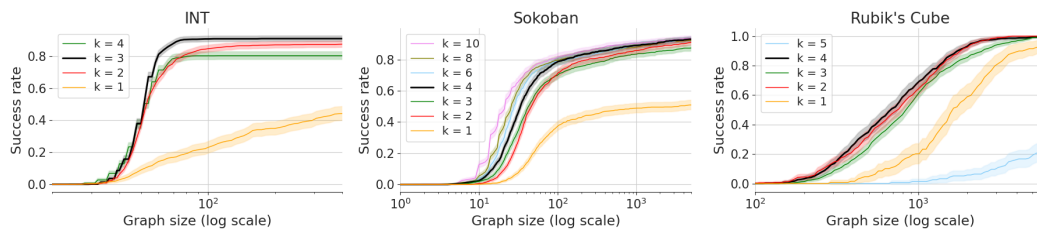


Figure 2: BF-kSubS success rates for different values of  $k$ . Black curves represent the values of  $k$  used in the main experiments (that is  $k = 4$  for Rubik’s Cube and Sokoban and  $k = 3$  for INT).

#### 4.5 Quality of subgoals

The learned subgoal generator is likely to be imperfect (especially in hard problems). We study this on  $10 \times 10$  boards of Sokoban, which are small enough to calculate the true distance  $dist$  to the solution using the Dijkstra algorithm. In Figure 3, we study  $\Delta := dist_{s_1} - dist_{s_2}$ , where  $s_1$  is a sampled state and  $s_2$  is a subgoal generated from  $s_1$ . Ideally, the histogram should concentrate on  $k = 4$  used in training. We see that in slightly more than 65% of cases subgoals lead to an improvement.

The low-level conditional policy in Algorithm 2 provides additional verification of generated states. We check that in INT and Rubik’s Cube, about 50% of generated subgoals can be reached by this policy (the rest is discarded).



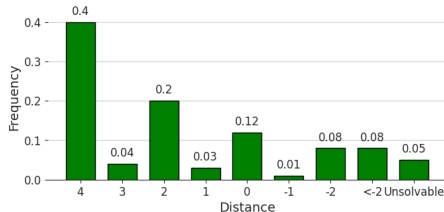


Figure 3: Histogram of  $\Delta$ . Note that 17% of subgoals increases the distance. Additional, 5% leads to unsolvable “dead states” present in Sokoban.

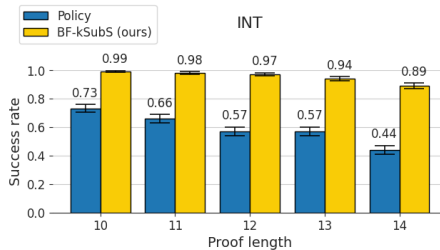


Figure 4: Out-of-distribution generalization to longer proofs. We compare with the behavioral cloning agent (Policy) studied in [54].

#### 4.6 Value errors

There might be various explanations for the success of our method. One of them is that Subgoal Search better handles errors of learned value functions. In this section, we discuss this using a synthetic grid world example and performing statistical analysis on the real value function trained to approximate the distance to the solution (as described in Section 3).

**Grid world example** Consider a grid world with the state space  $S = \{1, \dots, n\}^m$ , with  $(0, \dots, 0)$  being the initial state and  $(n, \dots, n)$  the goal state. A pair of states is connected by an edge if they are at distance 1 from each other. Let:

- Synthetic value function: the negative distance to the goal plus i.i.d Gaussian noise  $\mathcal{N}(0, \sigma^2)$ .
- Synthetic SUB\_GENERATE (instead of Algorithm 3): Let  $B_k(s)$  be the states within distance  $k$  from  $s$ . We return  $C_3 - 1$  states sampled uniformly from  $B_k(s)$  and a one "good subgoal" being a state in  $B_k(s)$  with the minimal distance to the solution.
- Node expansion in BestFS (baseline): implemented as SUB\_GENERATE above with  $k = 1$ .

In this setup, one easily sees that the probability that the good subgoal will have the highest value estimation among the generated states grows with  $k$ . Consequently, kSubS can handle higher levels of noise than the baseline BestFS, see Table 4.

**Value monotonicity** Imagine a solution path from the starting state to the goal state. Due to the errors, the value estimates on the path may not be monotonic. This is an undesirable property, which is likely to hinder search and make finding the path harder. Now consider the subpath consisting of consecutive states spaced  $k$  actions apart, as can be constructed by Subgoal Search. For this version, the value estimates are more likely to be monotonic and easier to find. To illustrate this, we measure monotonicity on solution paths found by our algorithm for INT. The probability that value decreases when moving  $k$  steps ahead drops from 0.32 when  $k = 1$  to mere 0.02 for  $k = 4$  (see Table 7 in Appendix).

**Overoptimism** Alternatively, one can consider that erroneously positive values misguide a search method (a phenomenon known as over-optimism [19]). To illustrate this, consider  $\mathcal{S}_3(s)$ , the set of all states having the same distance to the solution as  $s$  and within distance 3 from  $s$ . Intuitively,  $\mathcal{S}_3(s)$  contains similar states with respect to the difficulty of solving. In Sokoban, the standard deviation of value function prediction for  $\mathcal{S}_3(s)$  is equal to 2.43 (averaged over different  $s$  on Sokoban boards). This is high when compared to the average increase of value for moving one step closer to the solution, which is only 1.34. Consequently, it is likely that  $\mathcal{S}_3(s)$  contains a suboptimal state, e.g., having a higher value than the best immediate neighbors of  $s$  (which by properties of the game will be closer to solution in Sokoban). Indeed, we measure that the probability of such an event is 64%. However, it drops significantly to 29% if one considers states closer by 4 steps (say given by a subgoal generator).

$\sigma$	BestFS	BF-kSubS
3	0.999	1
10	0.142	1
20	0.006	0.983

Table 4: Success rates on the grid world ( $m = 6, n = 10$ ), depending on the value function noise scale. We use the search budget of 500 nodes and  $k = 4$  for kSubS.

## 5 Limitations and future work

In this section, we list some limitations of our work and suggest further research directions.

**Reliance on expert data** In this version, we use expert data to train learnable models. As kSubS improves the performance, we speculate that training akin to AlphaZero can be used, i.e. in a planner-learner loop without any outside knowledge.

**Optimality and completeness** kSubS searches over a reduced state space, which might produce suboptimal solutions or even fail to find them. This is arguably unavoidable if we seek an efficient method for complex problems.

**Subgoals definition** We use simple  $k$ -step ahead subgoals, which is perhaps not always optimal. Our method can be coupled with other subgoal paradigms. Unsupervised detection of landmarks (see e.g. [55]) seems an attractive future research direction.

**More environments** In future work, we plan to test kSubS on more environments to understand its strengths and weaknesses better. In this work, we generate subgoals in the state space, which might be limiting for tasks with high dimensional input (e.g., visual).

**Reliance on a model of the environment** We use a perfect model of the environment, which is a common practice for some environments, e.g., INT. Extending kSubS to use learned (imperfect) models is an important future research direction.

**Determinism** Our method requires the environment to be deterministic.

**OOD generalization** A promising future direction is to investigate and leverage the out-of-distribution generalization delivered by our method and compare to (somewhat contradictory) findings of [17, 54].

**Classical planning methods** For many search problems, the state space can be represented in factored fashion (or such representation can be learned [3]). In such cases, the search can be greatly improved with width-based methods [25, 12]. It is an interesting research direction to combine kSubS with such methods.

## 6 Conclusions

We propose Subgoal Search, a search algorithm based on subgoal generator. We present two practical implementations MCTS-kSubS and BF-kSubS meant to be effective in complex domains requiring reasoning. We confirm that indeed our implementations excel in Sokoban, Rubik’s Cube, and inequality benchmark INT. Interestingly, a simple  $k$  step ahead mechanism of generating subgoals backed up by transformer-based architectures performs surprisingly well. This evidence, let us hypothesize, that our methods (and related) can be further scaled up to even harder reasoning tasks.

## Acknowledgments and Disclosure of Funding

The work of Konrad Czechowski, Tomasz Odrzygóźdź and Piotr Miłoś was supported by the Polish National Science Center grant UMO-2017/26/E/ST6/00622. We gratefully acknowledge Polish high-performance computing infrastructure PLGrid (HPC Centers: ACK Cyfronet AGH, PCSS) for providing computer facilities and support within computational grant no. PLG/2020/013619. Our experiments were managed using <https://neptune.ai>. We would like to thank the Neptune team for providing us access to the team version and technical support.

## References

- [1] Forest Agostinelli, Stephen McAleer, Alexander Shmakov, and Pierre Baldi. Solving the rubik’s cube with deep reinforcement learning and search. *Nature Machine Intelligence*, 1(8):356–363, 2019.
- [2] Marcin Andrychowicz, Dwight Crow, Alex Ray, Jonas Schneider, Rachel Fong, Peter Welinder, Bob McGrew, Josh Tobin, Pieter Abbeel, and Wojciech Zaremba. Hindsight experience replay. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural*

- Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA, pages 5048–5058, 2017.
- [3] Masataro Asai and Alex Fukunaga. Classical planning in deep latent space: Bridging the subsymbolic-symbolic boundary. In Thirty-Second AAAI Conference on Artificial Intelligence, 2018.
  - [4] Yoshua Bengio, Andrea Lodi, and Antoine Prouvost. Machine learning for combinatorial optimization: a methodological tour d’horizon. European Journal of Operational Research, 2020.
  - [5] Cameron B Browne, Edward Powley, Daniel Whitehouse, Simon M Lucas, Peter I Cowling, Philipp Rohlfshagen, Stephen Tavener, Diego Perez, Spyridon Samothrakis, and Simon Colton. A survey of monte carlo tree search methods. IEEE Transactions on Computational Intelligence and AI in games, 4(1):1–43, 2012.
  - [6] Elliot Chane-Sane, Cordelia Schmid, and Ivan Laptev. Goal-conditioned reinforcement learning with imagined subgoals. In International Conference on Machine Learning, pages 1430–1440. PMLR, 2021.
  - [7] Peter Dayan and Geoffrey E. Hinton. Feudal reinforcement learning. In Stephen Jose Hanson, Jack D. Cowan, and C. Lee Giles, editors, Advances in Neural Information Processing Systems 5, [NIPS Conference, Denver, Colorado, USA, November 30 - December 3, 1992], pages 271–278. Morgan Kaufmann, 1992.
  - [8] Dorit Dor and Uri Zwick. Sokoban and other motion planning problems. Computational Geometry, 13(4):215–228, 1999.
  - [9] Kuan Fang, Yuke Zhu, Animesh Garg, Silvio Savarese, and Li Fei-Fei. Dynamics learning with cascaded variational inference for multi-step manipulation. arXiv preprint arXiv:1910.13395, 2019.
  - [10] Alan Fern, Roni Khordon, and Prasad Tadepalli. The first learning track of the international planning competition. Machine Learning, 84(1-2):81–107, 2011.
  - [11] Carlos Florensa, David Held, Xinyang Geng, and Pieter Abbeel. Automatic goal generation for reinforcement learning agents. In International conference on machine learning, pages 1515–1528. PMLR, 2018.
  - [12] Guillem Frances, Miquel Ramrez Jvega, Nir Lipovetzky, and Hector Geffner. Purely declarative action descriptions are overrated: Classical planning with simulators. In IJCAI 2017. Twenty-Sixth International Joint Conference on Artificial Intelligence; 2017 Aug 19-25; Melbourne, Australia.[California]: IJCAI; 2017. p. 4294-301. International Joint Conferences on Artificial Intelligence Organization (IJCAI), 2017.
  - [13] Thomas Gabor, Jan Peter, Thomy Phan, Christian Meyer, and Claudia Linnhoff-Popien. Subgoal-based temporal abstraction in Monte-Carlo Tree Search. In IJCAI, pages 5562–5568, 2019.
  - [14] Wei Gao, David Hsu, Wee Sun Lee, Shengmei Shen, and Karthikk Subramanian. Intention-net: Integrating planning and deep learning for goal-directed autonomous navigation. In 1st Annual Conference on Robot Learning, CoRL 2017, Mountain View, California, USA, November 13-15, 2017, Proceedings, volume 78 of Proceedings of Machine Learning Research, pages 185–194. PMLR, 2017.
  - [15] Timothy Gowers. The importance of mathematics. Springer-Verlag, 2000.
  - [16] Arthur Guez, Mehdi Mirza, Karol Gregor, Rishabh Kabra, Sbastien Racanire, Thophane Weber, David Raposo, Adam Santoro, Laurent Orseau, Tom Eccles, et al. An investigation of model-free planning. In International Conference on Machine Learning, pages 2464–2473. PMLR, 2019.
  - [17] Jessica B. Hamrick, Abram L. Friesen, Feryal Behbahani, Arthur Guez, Fabio Viola, Sims Witherspoon, Thomas Anthony, Lars Holger Buesing, Petar Velickovic, and Theophane Weber. On the role of planning in model-based deep reinforcement learning. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net, 2021.
  - [18] Demis Hassabis, Dharshan Kumaran, Christopher Summerfield, and Matthew Botvinick. Neuroscience-inspired artificial intelligence. Neuron, 95(2):245–258, 2017.

- [19] Hado Hasselt. Double q-learning. Advances in neural information processing systems, 23:2613–2621, 2010.
- [20] Jeffrey R Hollerman, Leon Tremblay, and Wolfram Schultz. Involvement of basal ganglia and orbitofrontal cortex in goal-directed behavior. Progress in brain research, 126:193–215, 2000.
- [21] Dinesh Jayaraman, Frederik Ebert, Alexei A. Efros, and Sergey Levine. Time-agnostic prediction: Predicting predictable video frames. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019.
- [22] Leslie Pack Kaelbling. Learning to achieve goals. In Ruzena Bajcsy, editor, Proceedings of the 13th International Joint Conference on Artificial Intelligence. Chambéry, France, August 28 - September 3, 1993, pages 1094–1099. Morgan Kaufmann, 1993.
- [23] Taesup Kim, Sungjin Ahn, and Yoshua Bengio. Variational temporal abstraction. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett, editors, Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pages 11566–11575, 2019.
- [24] Thanard Kurutach, Aviv Tamar, Ge Yang, Stuart J. Russell, and Pieter Abbeel. Learning plannable representations with causal infogan. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada, pages 8747–8758, 2018.
- [25] Nir Lipovetzky and Hector Geffner. Width and serialization of classical planning problems. In ECAI 2012, pages 540–545. IOS Press, 2012.
- [26] Kara Liu, Thanard Kurutach, Christine Tung, Pieter Abbeel, and Aviv Tamar. Hallucinative topological memory for zero-shot visual planning. In Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event, volume 119 of Proceedings of Machine Learning Research, pages 6259–6270. PMLR, 2020.
- [27] Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. Multilingual denoising pre-training for neural machine translation. CoRR, abs/2001.08210, 2020.
- [28] Piotr Miłoś, Łukasz Kuciński, Konrad Czechowski, Piotr Kozakowski, and Maciek Klimek. Uncertainty-sensitive learning and planning with ensembles. arXiv preprint arXiv:1912.09996, 2019.
- [29] Suraj Nair and Chelsea Finn. Hierarchical foresight: Self-supervised learning of long-horizon tasks via visual subgoal generation. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net, 2020.
- [30] Soroush Nasiriany, Vitchyr Pong, Steven Lin, and Sergey Levine. Planning with goal-conditioned policies. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett, editors, Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pages 14814–14825, 2019.
- [31] Giambattista Parascandolo, Lars Buesing, Josh Merel, Leonard Hasenclever, John Aslanides, Jessica B Hamrick, Nicolas Heess, Alexander Neitz, and Theophane Weber. Divide-and-conquer monte carlo tree search for goal-directed planning. arXiv preprint arXiv:2004.11410, 2020.
- [32] Sujoy Paul, Jeroen Vanbaar, and Amit Roy-Chowdhury. Learning from trajectories via subgoal discovery. Advances in Neural Information Processing Systems, 32:8411–8421, 2019.
- [33] Karl Pertsch, Oleh Rybkin, Frederik Ebert, Shenghao Zhou, Dinesh Jayaraman, Chelsea Finn, and Sergey Levine. Long-horizon visual planning with goal-conditioned hierarchical predictors. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020.

- [34] Karl Pertsch, Oleh Rybkin, Jingyun Yang, Shenghao Zhou, Konstantinos G. Derpanis, Kostas Daniilidis, Joseph J. Lim, and Andrew Jaegle. Keyframing the future: Keyframe discovery for visual prediction and planning. In Alexandre M. Bayen, Ali Jadbabaie, George J. Pappas, Pablo A. Parrilo, Benjamin Recht, Claire J. Tomlin, and Melanie N. Zeilinger, editors, Proceedings of the 2nd Annual Conference on Learning for Dynamics and Control, L4DC 2020, Online Event, Berkeley, CA, USA, 11-12 June 2020, volume 120 of Proceedings of Machine Learning Research, pages 969–979. PMLR, 2020.
- [35] Silviu Pitis, Harris Chan, Stephen Zhao, Bradly C. Stadie, and Jimmy Ba. Maximum entropy gain exploration for long horizon multi-goal reinforcement learning. In Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event, volume 119 of Proceedings of Machine Learning Research, pages 7750–7761. PMLR, 2020.
- [36] Stanislas Polu and Ilya Sutskever. Generative language modeling for automated theorem proving. arXiv preprint arXiv:2009.03393, 2020.
- [37] Sébastien Racanière, Theophane Weber, David P. Reichert, Lars Buesing, Arthur Guez, Danilo Jimenez Rezende, Adrià Puigdomènech Badia, Oriol Vinyals, Nicolas Heess, Yujia Li, Razvan Pascanu, Peter Battaglia, Demis Hassabis, David Silver, and Daan Wierstra. Imagination-augmented agents for deep reinforcement learning. In NIPS, 2017.
- [38] Tomas Rokicki and M Davidson. God’s number is 26 in the quarter-turn metric, 2014.
- [39] Stuart Russell and Peter Norvig. Artificial intelligence: A modern approach. ed. 3. 2010.
- [40] Nikolay Savinov, Alexey Dosovitskiy, and Vladlen Koltun. Semi-parametric topological memory for navigation. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net, 2018.
- [41] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. nature, 529(7587):484–489, 2016.
- [42] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Thore Graepel, Timothy Lillicrap, Karen Simonyan, and Demis Hassabis. A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play. Science, 1144:1140–1144, 2018.
- [43] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, Timothy P. Lillicrap, Karen Simonyan, and Demis Hassabis. Mastering chess and shogi by self-play with a general reinforcement learning algorithm. ArXiv, abs/1712.01815, 2017.
- [44] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy Lillicrap, Fan Hui, Laurent Sifre, George Van Den Driessche, Thore Graepel, and Demis Hassabis. Mastering the game of Go without human knowledge. Nature, 2017.
- [45] Gregory J. Stein, Christopher Bradley, and Nicholas Roy. Learning over subgoals for efficient navigation of structured, unknown environments. In 2nd Annual Conference on Robot Learning, CoRL 2018, Zürich, Switzerland, 29-31 October 2018, Proceedings, volume 87 of Proceedings of Machine Learning Research, pages 213–222. PMLR, 2018.
- [46] Richard S. Sutton, Joseph Modayil, Michael Delp, Thomas Degris, Patrick M. Pilarski, Adam White, and Doina Precup. Horde: a scalable real-time architecture for learning knowledge from unsupervised sensorimotor interaction. In Liz Sonenberg, Peter Stone, Kagan Tumer, and Pinar Yolum, editors, 10th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2011), Taipei, Taiwan, May 2-6, 2011, Volume 1-3, pages 761–768. IFAAMAS, 2011.
- [47] Richard S. Sutton, Doina Precup, and Satinder P. Singh. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. Artif. Intell., 112(1-2):181–211, 1999.
- [48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman

- Garnett, editors, Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pages 5998–6008, 2017.
- [49] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pages 5998–6008, 2017.
- [50] Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Manfred Otto Heess, Max Jaderberg, David Silver, and Koray Kavukcuoglu. Feudal networks for hierarchical reinforcement learning. ArXiv, abs/1703.01161, 2017.
- [51] Oriol Vinyals, Igor Babuschkin, Wojciech M. Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H. Choi, Richard Powell, Timo Ewalds, Petko Georgiev, Junhyuk Oh, Dan Horgan, Manuel Kroiss, Ivo Danihelka, Aja Huang, Laurent Sifre, Trevor Cai, John P. Agapiou, Max Jaderberg, Alexander S. Vezhnevets, Rémi Leblond, Tobias Pohlen, Valentin Dalibard, David Budden, Yury Sulsky, James Molloy, Tom L. Paine, Caglar Gulcehre, Ziyu Wang, Tobias Pfaff, Yuhuai Wu, Roman Ring, Dani Yogatama, Dario Wünsch, Katrina McKinney, Oliver Smith, Tom Schaul, Timothy Lillicrap, Koray Kavukcuoglu, Demis Hassabis, Chris Apps, and David Silver. Grandmaster level in starcraft ii using multi-agent reinforcement learning. Nature, 575(7782):350–354, Nov 2019.
- [52] David Wilkins. Using patterns and plans in chess. Artif. Intell., 14(2):165–203, 1980.
- [53] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. Huggingface’s transformers: State-of-the-art natural language processing. CoRR, abs/1910.03771, 2019.
- [54] Yuhuai Wu, Albert Jiang, Jimmy Ba, and Roger Grosse. Int: An inequality benchmark for evaluating generalization in theorem proving. arXiv preprint arXiv:2007.02924, 2020.
- [55] Lunjun Zhang, Ge Yang, and Bradly C. Stadie. World model as a graph: Learning latent landmarks for planning. CoRR, abs/2011.12491, 2020.

## A MCTS

### A.1 MCTS-kSubS algorithm

In Algorithm 4 we present a general MCTS solver based on AlphaZero. Solver repeatedly queries the planner for a list of actions and executes them one by one. Baseline planner returns only a single action at a time, whereas MCTS-kSubS gives around  $k$  actions – to reach the desired subgoal (number of actions depends on a subgoal distance, which not always equals  $k$  in practice).

MCTS-kSubS operates on a high-level subgoal graph: nodes are subgoals proposed by the generator (see Algorithm 3) and edges – lists of actions informing how to move from one subgoal to another (computed by the low-level conditional policy in Algorithm 2). The graph structure is represented by *tree* variable. For every subgoal, it keeps up to  $C_3$  best nearby subgoals (according to generator scores) along with a mentioned list of actions and sum of rewards to obtain while moving from the parent to the child subgoal.

Most of MCTS implementation is shared between MCTS-kSubS and AlphaZero baseline, as we can treat the behavioral-cloning policy as a subgoal generator with  $k = 1$ . All the differences between MCTS-kSubS and the baseline are encapsulated in GEN\_CHILDREN function (Algorithms 5 and 6). To generate children subgoals MCTS-kSubS runs subgoal generator and low-level conditional policy, whereas the baseline uses behavioral cloning policy for that purpose.

---

**Algorithm 4** MCTS solver (common for AlphaZero baseline and MCTS-kSubS)

---

<b>Require:</b>	$L_a$ action limit $L_p$ planner calls limit $P$ planning passes $\gamma$ discount factor $c_{puct}$ exploration weight $\tau$ sampling temperature $V$ value function $env$ environment $M$ environment model	<b>function</b> SELECT(state)	$s \leftarrow state$ $path \leftarrow []$ <b>while</b> $s$ belongs to $tree$ <b>do</b> $i \leftarrow SELECT\_CHILD(s)$ $s', r, actions \leftarrow tree[s][i]$ $path.APPEND((s, i, r))$ $s \leftarrow s'$ <b>return</b> path, $s$	
<b>Use:</b>	$tree$ tree structure $N(s, i)$ visit count $W(s, i)$ total child-value $Q(s, i)$ mean child-value $\pi_e$ exploration policy	<b>function</b> EXPAND(leaf)	$children, probs \leftarrow GEN\_CHILDREN(leaf)$ $tree[leaf] \leftarrow children$ $\pi(leaf, \cdot) \leftarrow probs$ <b>for</b> $i \leftarrow 1$ to $children.LENGTH()$ <b>do</b> $s', r, actions \leftarrow tree[leaf][i]$ $W(leaf, i) \leftarrow r + \gamma * V(s')$ $N(leaf, i) \leftarrow 1$ $Q(leaf, i) \leftarrow W(leaf, i)$	
# Initialize $N, W, Q$ to zero	<b>function</b> SOLVER	<b>function</b> UPDATE(path, leaf)	$s \leftarrow env.RESET()$ $solution \leftarrow []$ $\triangleright$ List of actions <b>for</b> $1 \dots L_p$ <b>do</b> $actions \leftarrow PLANNER(s)$ <b>for</b> $a$ <b>in</b> actions <b>do</b> $s', r \leftarrow env.STEP(a)$ $solution.APPEND(a)$ $s \leftarrow s'$ <b>if</b> $solution.LENGTH() > L_a$ <b>then</b> <b>return</b> None <b>if</b> $env.SOLVED()$ <b>then</b> <b>return</b> solution <b>return</b> None	$quality \leftarrow V(leaf)$ <b>for</b> $s, i, r \leftarrow reversed(path)$ <b>do</b> $quality \leftarrow r + \gamma * quality$ $W(s, i) \leftarrow W(s, i) + quality$ $N(s, i) \leftarrow N(s, i) + 1$ $Q(s, i) \leftarrow \frac{W(s, i)}{N(s, i)}$
<b>function</b> PLANNER(state)	<b>function</b> SELECT_CHILD( $s$ )	<b>function</b> CHOOSE_ACTIONS( $s$ )	<b>for</b> $1 \dots P$ <b>do</b> $path, leaf \leftarrow SELECT(state)$ $EXPAND(leaf)$ $UPDATE(path, leaf)$ <b>return</b> CHOOSE_ACTIONS(state)	$U(s, i) \leftarrow \sqrt{\sum_{i'} N(s, i') / (1 + N(s, i))}$ $i \leftarrow argmax_i (Q(s, i) + c_{puct} \pi_e(s, i) U(s, i))$ <b>return</b> $i$ $i \sim softmax(\frac{1}{\tau} \log N(s, \cdot))$ $s', r, actions \leftarrow tree[s][i]$ <b>return</b> actions

---

---

**Algorithm 5** GEN\_CHILDREN for MCTS-kSubS **Algorithm 6** GEN\_CHILDREN for AlphaZero

---

For functions *GET\_PATH* and *SUB\_GENERATE* see Algorithms 2 and 3.

```
function GEN_CHILDREN(state)
  s ← state
  children ← []
  probs ← []
  for subgoal, prob ← SUB_GENERATE(s) do
    actions ← GET_PATH(s, subgoal)
    if actions.EMPTY() then continue
    r ← M.REWARD_SUM(s, actions)
    children.APPEND((subgoal, r, actions))
    probs.APPEND(prob)
  return children, probs
```

**Require:**  $\pi_b$  behavioral cloning policy

```
function GEN_CHILDREN(state)
  s ← state
  children ← []
  probs ← []
  for a, prob ←  $\pi_b$ .GEN_ACTIONS(s) do
    s', r ← M.NEXT_STATE_REWARD(s, a)
    children.APPEND((s', r, [a]))
    probs.APPEND(prob)
  return children, probs
```

---

Variables  $tree, N, W, Q, \pi_e$  are reused across subsequent planner invocations within a single solver run. We limit the number of planner calls  $L_p$  for better control over the computational budget for MCTS-kSubS. For MCTS pseudocode we assume a slightly modified version of *SUB\_GENERATE* function (defined originally in Algorithm 3). We presume that the function along with subgoals returns also their respective probabilities – as MCTS needs them to guide exploration.

## A.2 Detailed results of MCTS-kSubS

We evaluate MCTS-based approaches on INT proofs of length 15. We set  $c_{puct} = \tau = 1$  and  $\gamma = 0.99$ . We tuned  $P$  on MCTS (AlphaZero) baseline and we run three variants with  $P \in \{5, 15, 50\}$ . We run MCTS-kSubS ( $k = 3$ ) with the same set of parameters and with kSubS-specific parameters fixed to  $C_2 = C_3 = 4$  (in order to match the setup for corresponding INT BF-kSubS experiments).

We limit the maximum number of actions to  $L_a = 24$  for both methods. Having the same number of planning passes  $P$ , during a single call MCTS-kSubS visits  $k$ -times more new states than the baseline (because of states visited by the low-level conditional policy). Therefore, to ensure a similar computational budget, we limit the number of planner calls to  $L_p = 8$  for MCTS-kSubS and to  $L_p = 24$  for the baseline – so the number of states visited over the course of a single solver run is similar for both methods.

Top-left part of Figure 1 illustrates results of MCTS experiments. For every number of planning passes  $P$ , MCTS-kSubS has significantly higher success rate than the corresponding baseline experiment. The highest difference is 0.52 for  $P = 5$  (0.88 for MCTS-kSubS, 0.36 for the baseline) and slightly decreases with increasing number of passes to still impressive 0.43 for  $P = 50$  (0.91 for MCTS-kSubS, 0.48 for the baseline). Comparing MCTS-kSubS for  $P = 5$  with the baseline for  $P = 50$ , shows advantage of our method still by a significant margin of 0.40, despite having 10 times smaller computational budget.

MCTS-kSubS performed better also in terms of run time. For every tested  $P$  it was at least twice as fast as the corresponding baseline experiment.

High effectiveness of MCTS-kSubS, in terms of both search success rate as well as run time, shows that our kSubS method is not specific to BestFS planner, but potentially allows to boost a wide range of other planners.

## B Architectures and hyperparameters

### B.1 Transformer

For INT and Rubik we use mBART [27] – one of the state-of-the-art sequence-to-sequence transformer architectures. To speed up training and inference we use its lightweight version. We reduced the dimensionality of the model, so the number of learned parameters decreased from the original 680M to 45M. The set of our hyperparameters matches the values proposed in [48]: we used 6 layers of encoder and 6 layers of decoder; we adjusted model’s dimension to 512 and number of attention



heads to 8; the inner-layer of position-wise fully connected networks had dimensionality 2048. The difference in our model’s size compared to 65M parameters reported in [48] results from vocabulary size. For our problems, it is enough to have 10-70 distinct tokens, whereas natural language models require a much larger vocabulary (tens of thousands of tokens).

For inference we used number of beams equal to 16 on INT and 32 on Rubik’s Cube.

## B.2 Sokoban

In Sokoban, we use three neural network architectures: for generating subgoals, for assigning value and one for baseline policy.

We took the value function network architecture from [28]. For the subgoal generator network we used the same convolutional architecture as in [28], with two exceptions. First, instead of predicting single regression target we predicted distribution over  $d \times d \times 7 + 1$  classes. Secondly, we added batch norm layers between convolutional layers to speed up training. To make the comparison between BestFS and BF-kSubS fair, we also evaluated the training of expert from [28] with additional batch norm layers, but it turned out to actually hurt the performance of the expert. The architecture for baseline policy was the same as in [28] with only one modification: it predicts one of our actions instead of a single number.

## C Data processing and datasets

### C.1 Sokoban

**Dataset.** We collected expert datasets using an RL agent (MCTS-based) from [28]. Precisely, we trained 3 agents on Sokoban boards of different sizes ( $12 \times 12$ ,  $16 \times 16$  and  $20 \times 20$ , all with four boxes). During the training process, we collected all successful trajectories, in the form of sequences of consecutive states. The number of trajectories in our datasets were: 154000 for  $12 \times 12$  boards, 45000 for  $16 \times 16$  boards and 21500 for  $20 \times 20$  boards. The difference comes from the fact that the larger boards take more time to solve, hence fewer data is collected in the same time span.

**Subgoal generation.** For a given state the generation of subgoals is depicted in Algorithm 7. We maintain a queue of modified states (MS). Iteratively we take a MS from queue, concatenate it with state and pass through subgoal generator network (`subgoal_net.SORTED_PREDICTIONS`). This produces a probability distribution over candidates for further modifications of given MS. We take the most probable candidates, apply each of them to MS, and add the new modified states to the queue. If among the best subgoal generator network predictions there is a special "valid subgoal" token (encoded with  $d \times d \times 7 + 1$ ), we put MS to subgoal candidates list (`subgoals_and_probs`). During this process, each MS is assigned the probability, which is a product of probabilities of modifications, leading to this MS. When the queue is empty, we take subgoal candidates and choose the ones with the highest probability such that the target probability ( $C_4$ ) is reached (similar to Algorithm 3). The generation of subgoals for a given state is illustrated in Figure 5.

This process is designed to be computationally efficient. The majority of subgoals differ from the input by only several pixels, which leads to short paths of point-wise modifications. Note, that we do not utilize any Sokoban-specific assumptions.

**Datapoints for training.** Preparing data points for the training of the generator is described in Algorithm 8. For each trajectory in the dataset, we choose randomly 10% of state pairs for the training (we do not use all states from a trajectory in order to reduce the correlation in the data).

**Low-level conditional policy.** In Algorithm 9 we describe the BFS-based algorithm that verifies subgoals in Sokoban.

**Performance of RL agent.** We observed that each of the three RL agents we used (for  $12 \times 12$ ,  $16 \times 16$  and  $20 \times 20$  boards), had a significantly lower success rate than our method’s counterparts (that learns from these agents). For  $12 \times 12$  boards it could solve around 78% of problems, for  $16 \times 16$  boards it dropped to 67% and for  $20 \times 20$  it was only 60%.



---

**Algorithm 8** Sokoban generate network inputs and targets

---

**Require:**  $d$  dimension of a board

**function** GENERATE\_INPUTS\_AND\_TARGETS(state, subgoal)

- inputs  $\leftarrow$  [] ▷ empty list
- targets  $\leftarrow$  [] ▷ empty list
- modified\_state  $\leftarrow$  state
- input  $\leftarrow$  CONCATENATE(state, modified\_state)
- inputs.APPEND(input)
- target\_class\_num  $\leftarrow$  0
- for**  $i \in 1 \dots d$  **do**
- for**  $j \in 1 \dots d$  **do**
- for**  $c \in 1 \dots 7$  **do**
- target\_class\_num  $\leftarrow$  target\_class\_num + 1
- if** subgoal[i, j, c] = 1 AND modified\_state[i, j, c] = 0 **then**
- targets.APPEND(target\_class\_num)
- ▷ Numpy notation, replace pixel values on position (i, j) with values
- ▷ from subgoal
- modified\_state[i, j, :]  $\leftarrow$  subgoal[i, j, :]
- input  $\leftarrow$  CONCATENATE(state, modified\_state)
- inputs.APPEND(input)
- ▷ Last target: no more changes to the modified\_state are needed (class enumerated
- ▷ with  $d \times d \times 7 + 1$ )
- targets.APPEND( $d \times d \times 7 + 1$ )
- return** inputs, targets

---

---

**Algorithm 9** BFS low-level conditional policy

---

**Require:**  $k$  limit of steps

$M$  model of the Sokoban environment

**Use:** bfs\_queue BFS queue;

stores pairs of a state and the action path to it (from the root).

# Initialize bfs\_queue to empty

**function** GET\_PATH( $s_0$ , subgoal)

- step  $\leftarrow$  0
- bfs\_queue.ADD( $(s_0, [])$ )
- while** bfs\_queue not empty **do**
- $s$ , action\_path  $\leftarrow$  bfs\_queue.POP()
- for** action  $\in$  action\_space **do**
- $s \leftarrow M.NEXT\_STATE(s, action)$
- action\_path.APPEND(action)
- if**  $s =$  subgoal **then**
- return** action\_path
- if** LEN(action\_path) <  $k$  **then**
- bfs\_queue.ADD( $(s, action\_path)$ )
- return** []

---

## C.2 INT

**State representation.** A state in INT consists of objectives to prove and ground truth assumptions, which are logic statements that are assumed to hold and can be used in proving. Each objective, as well as each ground truth assumption, is a mathematical statement. In our setup, as in the original paper [54], there is always only one objective to prove, but there may be a varying number of ground truth statements.

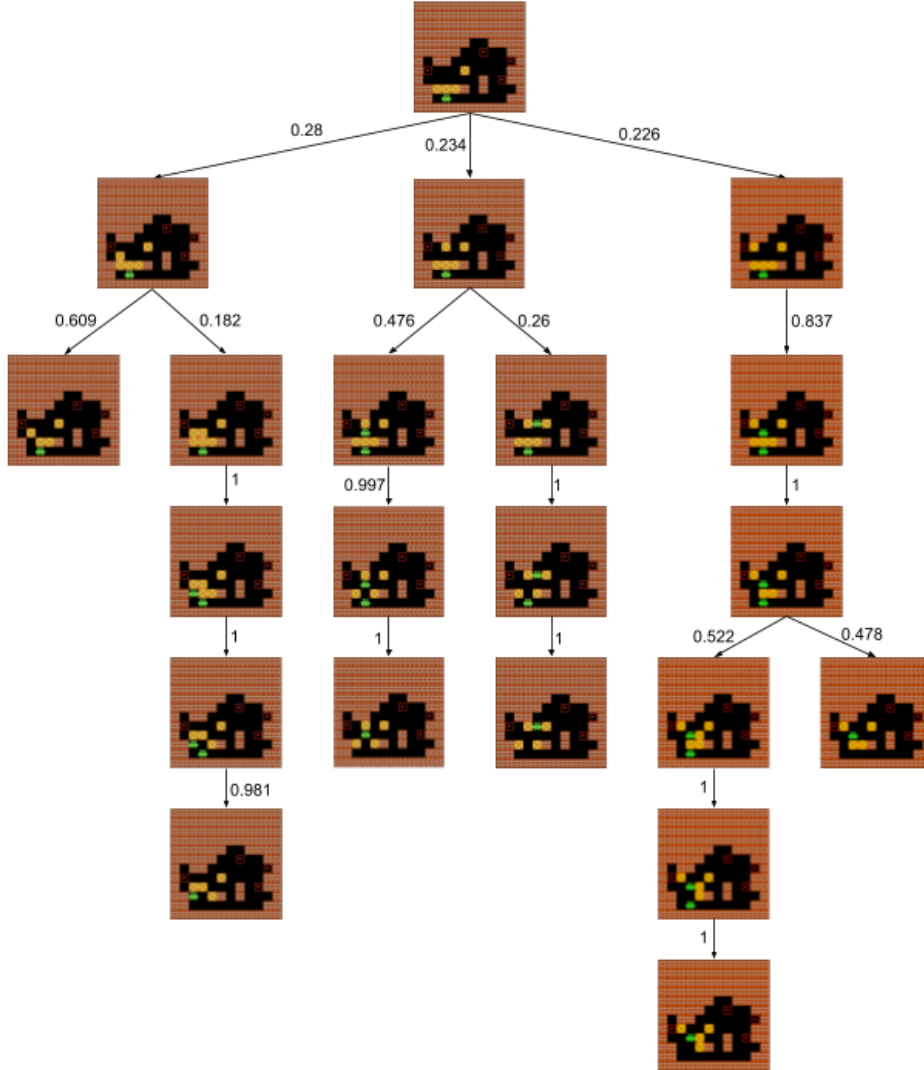


Figure 5: A detailed view of subgoal generation for Sokoban. Arrow represent probabilities of a given modification. Final subgoals are located in the leaves.

Each mathematical statement can be converted to a string by a method `logic_statement_to_seq_string` in INT library. In our code, we represent the full state as a string in the following form (all symbols are tokens, not logical operations):

$$\#[\text{objective}] \& [\text{1st ground truth}] \& [\text{2nd ground truth}] \& \dots \& [k - \text{th ground truth}] \#$$

For example, the state representation could look like this:

$$\#(((b+b)*((b+b)*(b+b)))*(((b+b)+f)*(b+b))*(b+b)) \geq 0 \& (b+f) = b \& (b+f) \geq 0 \#$$

**Action representation.** An action in INT consists of a chosen axiom (one from the set of ordered field axioms, see [54, Appendix C]) and a sequence of entities onto which the axiom will be applied. An entity is an algebraic expression that is a part of the larger statement. For example,  $(a + b)$  is an entity, which is a part of  $(a + b) \cdot c = (1 + f)$ . In our code, we represent entities by indicating the symbol of their mathematical operation, or if the entity is atomic (a single variable), by indicating the variable itself. More precisely, directly after the symbol of operation, we add a special character '~'. For example, we indicate  $(a + b)$  inside  $(a + b) \cdot c$  in the following way:  $(a + \sim b) \cdot c = (1 + f)$ . Typically, in a logic statement there may be several entities that have the same text form, but are

located in different positions, for example  $(a + b)$  appears twice in  $(a + b) \cdot (a + b) = (1 + 0)$ . Our way of encoding actions unambiguously identifies every entity. If the action has more than one input, we use more different indicators.

**Low-level conditional policy input representation.** Low-level conditional policy takes as an input two states:  $s$  and  $s'$ , where  $s$  is the initial state and  $s'$  is the subgoal. The input is constructed in the following way: first, we represent both  $s$  and  $s'$  as strings and then we find the difference (character delta) between these strings using the function `ndiff` from `diffliib`<sup>5</sup> Python library. We observed that using the character delta, instead of the concatenation of  $s$  and  $s'$ , significantly improved the performance.

### C.3 Rubik’s Cube

**State representation** The state of the Rubik’s Cube is determined by the arrangement of 54 colored labels on its faces. Therefore, to represent the observations we simply put the labels in a fixed order. An example state is as follows:

*?byywygrygobbrboorgwbowryogywbgywrrroyogyowubrwbogg\$,*

where the tokens  $b, g, o, r, w, y$  stand for *blue, green, orange, red, white, and yellow*. The consecutive blocks of 9 tokens correspond to consecutive faces of the cube. Observe, that not every permutation of colors is valid. For example, the tokens on positions 5, 14, 23, 32, 41, and 50 correspond to centers of faces, thus they are fixed. There are more such constraints, but they are irrelevant to the pipeline itself.

**Action representation** In our experiments we use quarter turns, i.e. an action corresponds to rotating a face by  $90^\circ$ , either clockwise or counterclockwise. Since the action space contains only 12 elements, we use unique tokens to represent each of them.

**Low-level conditional policy input representation.** The conditional policy takes two states  $s$  and  $s'$ , which correspond to the current state and the state to be reached. To represent such pairs, on every position we put a token corresponding to a pair of colors – one located on that position in  $s$  and the other in  $s'$ . Since there are only 6 distinct colors on the Rubik’s Cube, this requires using only 36 tokens.

## D Training details

### D.1 INT and Rubik’s Cube

#### D.1.1 Transformer training

For transformer training and inference we used HuggingFace’s Transformers library [53]. We did not use any pretrained checkpoints from HuggingFace model hub. We took mBART model class instead – and trained it from scratch in a supervised way using HuggingFace’s training pipeline. We generated (or loaded from disk) a fresh dataset for every epoch. Training batch was of size 32. For regularization, we set  $dropout = 0.1$ , but we did not use label smoothing (as opposed to [48]).

For the Adam optimizer we set  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and  $\epsilon = 10^{-8}$ . Learning rate schedule followed the formula:

$$lr = peak\_lr * \min \left( \frac{step\_num}{warmup\_steps}, \sqrt{\frac{warmup\_steps}{step\_num}} \right),$$

where  $peak\_lr = 3 \cdot 10^{-4}$  and  $warmup\_steps = 4000$ .

The schedule curve matches the one proposed in [48], but they use  $peak\_lr \approx 7 \cdot 10^{-4}$ .

#### D.1.2 Sequence generation

We used beam search with the number of beams set to 16 for INT and to 32 for Rubik’s Cube. The number of returned sequences varied from 1 to 4 depending on the network type.

<sup>5</sup><https://docs.python.org/3/library/diffliib.html>

We set softmax temperature to 1 by default. For Rubik’s Cube subgoal generator we tuned this parameter and the value of 0.5 performed best. We conducted a corresponding experiment for the baseline policy, but the temperature did not impact results in this case, as the policy outputs only a single token. For INT we did not tune the temperature, so we kept the value of 1.

### **D.1.3 Random seeds**

Due to use of the supervised learning, we observed little variance with respect to the random initialization. We tested this for the subgoal generator on proofs of length 10 and for  $k = 3$ . Namely, we trained 5 models of the subgoal generator, starting from different initializations. The success rate barely varied, as they stayed in the interval  $[0.990, 0.992]$ . In the other experiments, we used a single seed.

## **D.2 Sokoban**

For training of the convolutional networks in Sokoban we set the learning rate to  $10^{-4}$  and the number of epochs to 200.

## E Wall-time for kSubS

As indicated in Table 2, kSubS builds smaller search graphs. This has the practical advantage of making fewer neural network calls and consequently a substantially better wall-time.

The gains might be as high as 7 times due to costly sequential calls of transformer networks, see Table 5.

Proof length	5		10		15	
Method	BestFS	BF-kSubS (ours)	BestFS	BF-kSubS (ours)	BestFS	BF-kSubS (ours)
Total wall-time	4h 12m	<b>3h 44m</b>	29h 22m	<b>5h 55m</b>	69h 15m	<b>9h 22m</b>
Avg. generator calls	NA	<b>3.04</b>	NA	<b>3.89</b>	NA	<b>6.23</b>
Avg. value calls	23.34	<b>4.01</b>	112.35	<b>4.80</b>	159.46	<b>6.59</b>
Avg. policy calls	22.41	<b>8.43</b>	112.21	<b>13.09</b>	161.02	<b>20.29</b>

Table 5: Resources consumption for INT. We present evaluation on 1000 proofs and split into calls of subgoal generator network (used only in the subgoal search), value network and policy network (we report an average number of calls for a single proof).

## F Training dataset size analysis

We tested how the success rate of BF-kSubS on 12x12 Sokoban boards depends on the size of the training set. The full dataset consists of 125k trajectories. We trained subgoal generator and value network on subsets consisting of 0.5, 0.25, 0.05 and 0.01 of all trajectories. The results are presented in Table 6.

Fraction of the dataset	1	0.5	0.25	0.05	0.01
Success rate	0.93	0.86	0.84	0.48	0.14

Table 6: Sokoban success rates for different training set sizes.

## G Value errors

### G.1 INT analysis

Due to the size of state spaces, it is impossible to search over the entire space of INT formulas to find the shortest proofs. We instead analyze value estimations along proofs generated by INT engine. The monotonicity analysis in Section 4.6 was performed using 100 such proofs of length 10. The probabilities of value decrease for different step lengths  $l$  are presented in Table 7.

### G.2 Sokoban Analysis

Here, we present details related to the Sokoban analysis from Section 4.6. We sampled 500 Sokoban boards (with dimension (12, 12)). For each board, we calculated a full state-action graph and minimal distance from each state to the solution (this was needed to compute  $S(s)$  sets later on). Since Sokoban graphs can be very large we set the limit on the graph size to 200000, which left us with 119 boards. Next, for each board, we took the sequence of states from the shortest solving path from the initial state (let us call this dataset as *shortest-paths* - SP). For each pair of consecutive states in SP, we calculated the difference of the value estimation, and averaged them, which gave us a mean one-step improvement of 1.34. We calculated this metric for 5 value function networks trained with different initialization, obtaining mean one-step improvement between [1.23, 1.41].

$l$	Value decrease prob.
1	0.316
2	0.217
3	0.080
4	0.020

Table 7

To calculate the standard deviation of value function estimates for  $S(s)$  we took SP, and limit it to states  $s$  such that  $|S(s)| \geq 5$  (lets denote it as SP5). We calculated standard deviation for each

$s \in SP5$  separately. This gave us a mean deviation of 2.43. (between [2.24, 2.86] for 5 value networks trained with different initialization) The same set SP5 was used to calculate probabilities related to overoptimistic errors on Sokoban described at the end of Section 4.6.

To calculate the above statistics we used the value function trained with supervised learning to approximate the distance to the solution. We observe that similar problems arise also when using value function trained with reinforcement learning [28]. In such setup, mean variance of value function estimates for  $S(s)$  is 0.84, when one step improvement equals to 0.33. Probability that there is a state in  $S(s)$  with value higher than best immediate neighbour of  $s$  is 86% and it drops to 38%, if one considers states closer by 4 steps.

## H Example subgoals

### H.1 Example subgoals sets

In this section, we present some example outcomes of the subgoal generator for INT and Sokoban.

#### H.1.1 INT

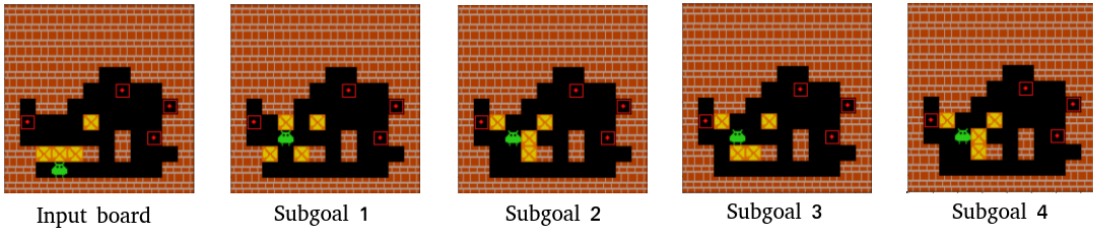
In (1) and (2) we provide two examples of applying the subgoal generator (trained on proofs of length 5) to the given states in INT. The number of subgoals varies since not all of the outputs generated by the network could be reached by the conditional low-level policy.

$$\begin{aligned}
 &\text{Input state: } (((b \cdot b) + ((b + (b + f)) \cdot b)) + (f + f)) \geq (((b + (b + b)) \cdot b) + 0) + c \\
 &\text{Ground truth 1: } (b + f) = b \\
 &\text{Ground truth 2: } (f + f) \geq c \\
 &\text{Subgoal 1: } ((b \cdot b) + ((b + (b + f)) \cdot b)) = (((b + (b + b)) \cdot b) + 0) \\
 &\text{Subgoal 2: } (((b + (b + f)) \cdot b) + (b \cdot b)) = (((b + (b + b)) \cdot b) + 0) \\
 &\text{Subgoal 3: } ((b \cdot b) + ((b + (f + b)) \cdot b)) = (((b + (b + b)) \cdot b) + 0)
 \end{aligned} \tag{1}$$

$$\begin{aligned}
 &\text{Input state: } (((((b \cdot (\frac{1}{b})) + a)^2) + (c + (\frac{1}{b})))) = ((((((\frac{1}{b}) \cdot b) + a) \cdot (a + 1)) + c) + (\frac{1}{b})) \\
 &\text{Subgoal 1: } (((((b \cdot (\frac{1}{b})) + a)^2) + (c + (\frac{1}{b})))) = (((((b \cdot (\frac{1}{b})) + a) \cdot (a + 1)) + c) + (\frac{1}{b})) \\
 &\text{Subgoal 2: } (((((b \cdot (\frac{1}{b})) + a) \cdot ((b \cdot (\frac{1}{b})) + a)) + (c + (\frac{1}{b})))) = (((((b \cdot (\frac{1}{b})) + a) \cdot (a + 1)) + c) + (\frac{1}{b}))
 \end{aligned} \tag{2}$$

#### H.1.2 Sokoban

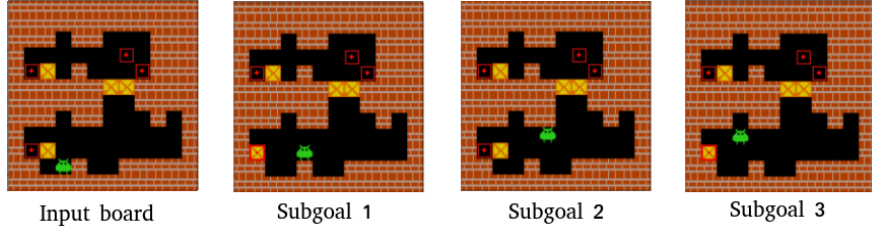
Here we present two examples of the outcomes of the subgoal generator trained for  $12 \times 12$  boards:



## H.2 Example solutions with subgoals

In this section, we present several examples of solutions obtained with our method. For simplicity, we only show the subgoal states on which the successful trajectories were constructed. In our setup, the last subgoal is always a solution.





### H.2.1 INT

An example solution of INT problem of length 5:

$$\text{Problem: } (((0 \cdot ((a+0) + (-a \cdot 1))) \cdot (\frac{1}{(0^2)}) + ((a+0) + (0^2))) \geq (((0^2) + (1 + (a+0))) + b) + (-((a+0) + f)))$$

$$\text{1st subgoal: } (((0 \cdot ((a+0) + (-a \cdot 1))) \cdot (\frac{1}{(0^2)}) + ((a+0) + (0^2))) = ((0^2) + (1 + (a+0))))$$

$$\text{2nd subgoal: } (((0 \cdot ((0+a) + (-a \cdot 1))) \cdot (\frac{1}{(0^2)}) + ((a+0) + (0^2))) = (1 + ((a+0) + (0^2))))$$

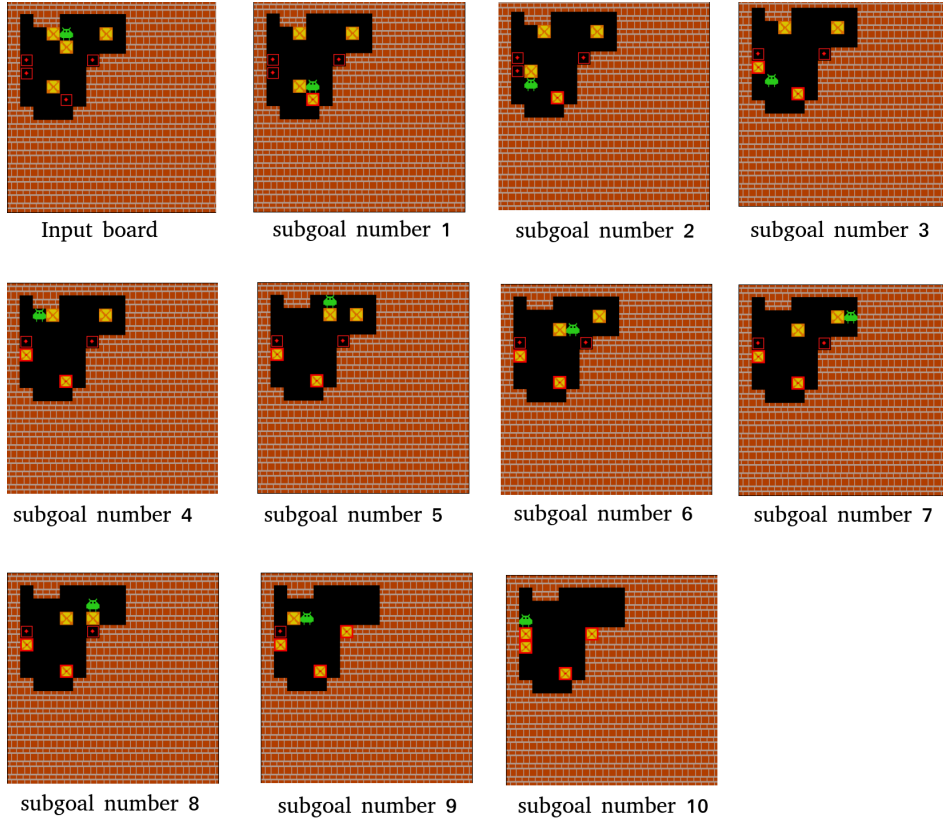
$$\text{3rd subgoal: } (0+a) = (a \cdot 1)$$

$$\text{4th subgoal: } a = a$$

(3)

### H.2.2 Sokoban

An example solution of Sokoban board:



## I Baselines

Our first baseline is the low-level policy trained with behavioral cloning from the expert data trained to solve the problem (contrary to the low-level conditional policy  $P$ , which aims to achieve subgoals). Such policy was used [54]. We verified that our behavioral cloning policy reproduces the results from [54] for proofs of lengths 5.

**MCTS.** As a baseline for MCT-kSubS we used an AlphaZero-based MCTS planner described in Appendix A.1.

**BestFS.** The baseline for BF-kSubS is a low-level planning. We substitute SUB\_GENERATE with a function returning adjacent states indicated by the most probable actions of behavioral cloning policy, see Algorithm 10.

---

**Algorithm 10** Low-level generator

---

<b>Require:</b>	$C_3$	number of states to produce
	$NB$	number of beams in sampling
	$\pi_b$	behavioral cloning policy
	$M$	model of the environment

```

function SUB_GENERATE(s)
  actions ←
  BEAM_SEARCH( $\pi_b, \mathbf{s}; C_3, NB$ )
  subgoals  $\leftarrow \square$ 
  for action  $\in$  actions do
     $s \leftarrow M.NEXT\_STATE(\mathbf{s}, \text{action})$ 
    subgoals.APPEND( $\mathbf{s}$ )
  return subgoals

```

---

## J Simple planners

An appropriate search mechanism is an important design element of our method. To show this, we evaluate an alternative, simpler procedure used by e.g. [9] for subgoal-based planning. It works by sampling independent sequences of subgoals and selects the best one. This method solved none of 1000 Rubik’s Cube instances despite using the same subgoal-generator as BF-kSubS (which has a success rate of 0.999 with a comparable computational budget).

## K Investigation of baseline-BestFS on Rubik’s Cube

To obtain the training datasets on the Rubik’s Cube environment, we generated random paths starting from a solved cube and stored them in the reversed order. These backward solutions are highly sub-optimal: for example, the states obtained by 10 random moves are usually in a distance of about 6 - 7 steps from the final state and the gap gets much larger for a higher number of moves. This means that on collected trajectories only some of the actions indeed lead to the solution and the majority of them only introduce noise.

We observed that the baseline is much weaker than BF-kSubS even for  $k = 1$ , despite the effort on tuning it. We extended the training of behavioral cloning policy and did a grid-search over parameters to further improve its success rate. We managed to reach no more than 10% success rate for the best version. To provide a meaningful evaluation of the baseline, we also trained it on very short trajectories consisting of 10 random moves. Such a curriculum allowed the behavioral cloning policy to learn, however still suffered from randomness in data (for states located far from the solution).

Interestingly, BF-kSubS for  $k = 1$  turned out to perform much better than the baseline. Full understanding of this phenomenon requires additional research, however we hypothesize that in our setup learning to predict states is an easier task than predicting an action. A potential reasons are that: states prediction provides a denser learning signal and Transformers perform better when dealing with sequences (then with predicting a single token).

Note that both kSubS and baseline policy learn from fully off-policy data. The problem of solving the Rubik’s Cube is challenging, thus learning from noisy off-policy trajectories can be simply too hard for standard algorithms.

## L Technical details

### L.1 Infrastructure used

We had 2 types of computational nodes at our disposal, depending on whether a job required GPU or not. GPU tasks used a single Nvidia V100 32GB card (mostly for transformer training) and Nvidia

RTX 2080Ti 11GB (for evaluation) with 4 CPU cores and 16GB RAM. The typical configuration of CPU job was the Intel Xeon E5-2697 2.60GHz processor (28 cores) with 128GB memory.

Each transformer for INT was trained on a single GPU for 3 days (irrespective of proof length and the network type). INT evaluation experiments used a single GPU for a time period varying from several hours to 3 days – with baselines being the bottleneck.

Transformers for Rubik’s Cube required more training – every network was trained for 6 days on a single GPU. Because of relatively short sequences representing the cube’s state, we were able to run evaluations without GPU. We used 20 CPU cores with 20GB memory, while still fitting in a 3-day run time.

We trained and evaluated Sokoban on CPU only mostly because of rather small neural network sizes. Training time varied from 2 to 3 days, whereas evaluation took only an hour.