
A Novel Automated Curriculum Strategy to Solve Hard AI Planning Instances

Anonymous Author(s)
Affiliation
Address
email

1 Automated Curriculum Framework

2 Herein we provide additional details concerning the different algorithms comprising our *automated*
3 *curriculum framework*. For the sake of completeness, the overview of the workflow of our *automated*
4 *curriculum framework* is depicted in Figure 1 and the formal algorithm description can be found in
Algorithm 1. The different algorithms include pointers to their sub-routine algorithms.

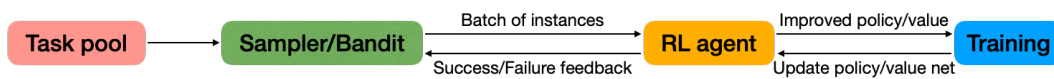


Figure 1: The workflow of our automated curriculum framework, which is formally described in Algorithm 1.

5

Algorithm 1: Automated Curriculum Learning framework overview

Input: Sokoban instance \mathcal{I} , solution length limit L , number of iterations T ;
Call $\text{TaskPool}(\mathcal{I})$ (Algorithm 2) to create a pool of sub-instances;
for $t = 1, \dots, T$ **do**
 Use uniformly sampling (baseline) or difficulty quantum momentum bandit (Algorithm 3) to
 select a batch B from the task pool;
 for $s \in B$ **do**
 for $i = 1, \dots, L$ **do**
 Use $\text{MCTS}(s)$ (Algorithm 4) to select the best move a for s ;
 $s = \text{next_state}(s, a)$;
 if s is a goal state **then**
 Generate data for training the policy/value network of the RL agent (i.e., for each
 MCTS node (board) along the branch from the input state (root) to the goal state,
 estimated distance to goal and distribution of visits to child nodes);
 Send "success" feedback to the sampler/bandit;
 Break;
 end
 end
 if Solution not found **then**
 Send "failure" feedback to the sampler/bandit;
 end
 end
 Train the policy deep neural network of the RL agent;
end

6 The master algorithm (Algorithm 1), starts by creating a task pool of sub-instances from the target
7 instance (Algorithm 2) that we need to solve, and potentially from other unsolved instances to further
8 boost performance (option MIX). Typically the task pool contains 100,000 tasks or sub-instances.
9 In each iteration, the sampler/bandit picks a batch of task sub-instances from the pool and passes
10 it to the RL agent. A batch has typically 500 tasks or sub-instances (Algorithm 3). The RL agent,
11 which is based on Monte-Carlo tree search (Algorithm 4), augmented with neural networks (CNN or
12 GNN), attempts to solve these instances. For each instance in the batch, MCTS will seek a solution
13 with a given resource budget, and for each successful solution generated, MCTS will also generate a
14 chain of new training data for the policy/value deep network (trainer) to further update its network
15 parameters. The MCTS success/failure status of each instance is sent back to the sampler/bandit to
16 adjust its weights. Each successful attempt not only generates a valid solution but also improves
17 policy/value data for the trainer to train the deep network of the agent. The trainer keeps a pool of
18 size 100000 to store the most recent training data generated by MCTS, and train the network. Each
19 training batch is uniformly randomly sampled. All experiments are done on a machine with 2x18
20 core Xeon Skylake 6154 CPUs and 5 Nvidia Tesla V100 16GB GPUs, and all training component
21 use Adam with learning rate 0.002 as the default optimizer. The number of MCTS simulation R is
22 set to 1600 and the batch size M that Exp3 samples in each iteration is set to 500.

Algorithm 2: function TaskPool(s)

Input: Sokoban instance \mathcal{I} ;

Parameters: The pool size P ;

$p = \{\}$;

$boxes =$ all initial box locations in \mathcal{I} ;

$goals =$ all initial goal locations in \mathcal{I} ;

$N = \text{size_of}(boxes)$;

for $i = 1, \dots, P$ **do**

$n = \text{UniformRand}([1, N])$;

$rand_boxes =$ A random subset of $boxes$ with size n ;

$rand_goals =$ A random subset of $goals$ with size n ;

$p = p \cup \text{SokobanInstance}(rand_boxes, rand_goals)$ (i.e., first build a empty Sokoban instance with wall and player location of \mathcal{I} unchanged, and then add $rand_boxes$ and $rand_goals$ to the board);

end

return p ;

Algorithm 3: Exp3 to sample batches of instances using difficulty quantum momentum heuristic

Input: a task pool P ;

Parameters: batch size M , exploration ratio $\gamma \in [0, 1]$, momentum $\alpha \in [0, 1]$, total iteration T ;

Initialization: $N = \text{size_of}(P)$, $\omega_i(1) = 1$, $h_i(1) = 0$ for $i = 1, \dots, N$;

for $t = 1, \dots, T$ **do**

$p_i(t) = (1 - \gamma) \frac{\omega_i(t)}{\sum_{k=1}^N \omega_k(t)} + \frac{\gamma}{N}$ for $i = 1, \dots, N$;

 Sample a non-repeated batch $B_1(t), \dots, B_M(t)$ from P according to the probability $\mathbf{p}(t)$;

 Run the RL agent on the batch $B(t)$ and train on the collected data;

$r_j(t) = (h_{B_j(t)}(t) - \mathbf{1}_{\text{succeed on } B_j(t)})^2$ for $j = 1, \dots, M$;

for $j = 1, \dots, M$ **do**

$\theta_j(t) = r_j(t) / p_{B_j(t)}(t)$;

$\omega_{B_j(t)}(t+1) = \omega_{B_j(t)}(t) \cdot \exp(\gamma \cdot \theta_j(t) / N)$;

$h_{B_j(t)}(t+1) = \alpha \cdot h_{B_j(t)}(t) + (1 - \alpha) \cdot r_j(t)$;

end

end

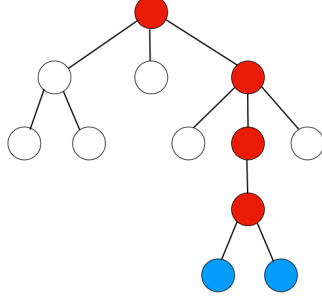


Figure 2: A whole simulation of MCTS. White and red circles correspond to the Monte Carlo tree before simulation. A simulation starts from the root node and goes down until it reaches a leaf node (the lowest red circle). Then an Expand procedure follows and adds new child nodes (blue) beneath the expanded node.

Algorithm 4: function $\text{MCTS}(s_0)$

Parameters: maximum solution length L , action set \mathcal{A} , number of MCTS simulations R , visit count $N(s, a)$;

Input: board state s_0 to seek a solution;

```

for  $l = 1, \dots, L$  do
  while  $\sum_{a \in \mathcal{A}} N(s_{l-1}, a) < R$  do
    Simulate( $s_{l-1}$ ) (Algorithm 5, see Figure 2 for the demonstration of a single simulation);
  end
  best_action =  $\text{argmax}_a N(s_{l-1}, a)$ ;
   $s_l = \text{NextState}(s_{l-1}, \text{best\_action})$ ;
  if  $s_l$  is a goal state then
    for  $i = 0, \dots, l - 1$  do
      Add  $\langle s_i, \text{Normalized}(N(s_i)), l - i \rangle$  to the trainer;
    end
    Break;
  end
end

```

23 **2 Network Architecture**

24 We use convolution neural network (CNN) as the baseline and compare its performance with graph
 25 network (GN) to show how different architecture setting affects the final result. The input to the CNN
 26 network is a $7 \times H \times W$ image stack consisting of 7 features planes with height H and width W .
 27 Each feature plan corresponds to walls, empty squares, empty goal squares, boxes, boxes on goal
 28 square, player-reachable squares, player-reachable squares on goal square, respectively. We use the
 29 standard ResNet-18 to extract the feature of the input Sokoban instance. For graph network, we build
 30 the graph by assigning a node to each cell of the board input and adding edges to each pair of adjacent
 31 cells. Horizontal and vertical edges have two different labels to further enhance spatial information
 32 to GN. Each input board cell lies in seven different categories as the same in the CNN architecture.
 33 We learn an embedding from these seven categories to a feature vector of length 128 as the starting of
 34 the GraphNet (Algorithm 7). We set the number of iterations D of graph network to 10. The output
 35 feature of graph network is further sent to two different multiple perceptions (MLP) to predict action
 36 probability (\mathbf{p}) and remaining distance v of the input state s .

37 The output of CNN and GN consists of two predictions: action probability \mathbf{p} and estimated remaining
 38 step v . We use *Softmax* activation for the probability \mathbf{p} . Since we set a maximum solution length L
 39 throughout the experiment and $v \in [0, L]$, we normalize the step prediction to $[0, 1]$ and use *Tanh*
 40 activation for the value v .

Algorithm 5: function Simulate(s) (See also see Figure 2).

Parameters: visit count $N(s, a)$, mean action value $Q(s, a)$;
while s not a goal state **do**
 if s is a leaf node **then**
 Expand(s) (Algorithm 3);
 end
 else
 best_action = $\operatorname{argmax}_a Q(s, a) + \text{cput} \cdot \frac{\sqrt{1 + \sum_b N(s, b)}}{1 + N(s, a)} \cdot \mathbf{p}_a$;
 $s = \text{NextState}(s, \text{best_action})$;
 end
end

Algorithm 6: function Expand(s)

Parameters: visit count $N(s, a)$, mean action value $Q(s, a)$;
 $\mathbf{p}, v = f_\theta(s)$;
for $a \in \mathcal{A}$ **do**
 $N(s, a) = 0$;
 $Q(s, a) = \text{CuriosityReward}(\text{NextState}(s, a))$;
end
while v not root **do**
 $r = \text{Parent}(v)$;
 $a = \text{PreviousAction}(v)$;
 $Q(r, a) = (Q(r, a) \cdot N(r, a) + v) / (N(r, a) + 1)$;
 $N(r, a) = N(r, a) + 1$;
 $v = r$;
end

41 3 Curiosity Reward

42 We use random network distillation (RND) as our curiosity reward generator. Specifically, we build a
43 graph network f_δ with randomized parameters and fix the parameters throughout the whole procedure.
44 We then try to learn another graph network f_τ with different randomized initialization and try to make
45 the prediction of f_τ as similar as the one of f_δ . For each state s that is requested for a curiosity reward,
46 we set $\text{CuriosityReward}(s) = l_2$ distance between $f_\tau(s)$ and $f_\delta(s)$. After each reward prediction, the
47 input state s is sent to a training pool of size 100000. At the end of each iteration, we train $f_{-\tau}$ for
48 100 epochs of batch size 64 to make its output closer to that of f_δ using squared error loss between
49 the outputs of the two networks.

Algorithm 7: Graph Neural Network extracting feature from a Sokoban board

Input: graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$, input features $\{x_v, \forall v \in \mathcal{V}\}$, depth D , neighborhood function $\mathcal{N} : v \rightarrow 2^{\mathcal{V}}$;

Output: A global feature of the graph $h_v^0 = x_v \forall v \in \mathcal{V}$;

for $d = 1, \dots, D$ **do**

$g_v^d = \text{Aggregate}_k(h_v^{d-1}, \forall u \in \mathcal{N}(v))$;

$h_v^d = \text{Normalized}(\text{ReLU}(W^d \cdot \{h_v^{d-1}, g_v^d\}))$;

end

return $\text{Average}(h_v^D)$ for $v \in \mathcal{V}$;
