
A APPENDIX

A.1 ALGORITHMS

We tested the combination of TBV with two planning algorithms: Monte Carlo Tree Search (equipped with some auxiliary mechanisms) and Best First Search. Here we present the pseudo-code for both algorithms and provide some additional explanation.

A.1.1 BESTFS

In Algorithm 3 we present the pseudocode for BestFS algorithm used in our experiments. In Data Structures 4 and 5 we present two additional data structures used by our implementation of BestFS.

Algorithm 3 BestFS planner

Require: *model* learned model used in Algorithm 1
 V value function
 C number of nodes expansions per step
 γ discount factor
Use: *F* priority queue of unexpanded nodes (fringe)
 G graph of all seen nodes

function RESET
 $F \leftarrow \emptyset$
 $G \leftarrow \emptyset$

function CHOOSE_ACTION(*state*)
 $G.ADD(state)$
 if $state \notin G$ **then**
 $state.n_visits \leftarrow 1$
 else
 $state.n_visits \leftarrow state.n_visits + 1$
 $reachable \leftarrow FIND_REACHABLE_NODES(G, state)$ \triangleright find all nodes reachable from state
 EXPAND_GRAPH($reachable$)
 $best_node \leftarrow G.FIND_BEST_NODE(reachable)$ \triangleright see Algorithm 5
 $path \leftarrow FIND_PATH_TO_NODE(best_node)$
 return $path.first_action$

function EXPAND_GRAPH($reachable$)
 for $1 \dots C$ **do**
 $n \leftarrow F.POP_BEST_NODE(reachable)$ \triangleright see Algorithm 4
 EXPAND_GRAPH_NODE(n)

function EXPAND_GRAPH_NODE($n, reachable$)
 $children \leftarrow []$
 for $a \in \mathcal{A}$ **do**
 $model.LOAD_STATE(n.state)$
 $new_state \leftarrow model.STEP(a)$
 $new_state.uncertainty \leftarrow \max_{action} model.uncertainty(new_state, action)$
 if $new_state \notin G$ **then**
 $value \leftarrow V.EVALUATE(new_state)$
 $new_state.SET_VALUE()$
 $F.ADD_NEW_NODE(new_state)$
 $G.ADD_NEW_NODE(new_state)$
 $reachable.ADD(new_state)$
 $F.REMOVE(n)$ \triangleright fringe contains only unexpanded nodes

function GET_GRAPH_OF_PLANNING()
 return G

Data structure 4 Fringe (priority queue)

Require: $S_F \triangleright$ set of nodes
1: **function** ADD_NODE(n)
2: $S_F.ADD(n)$
3: **function** REMOVE(n)
4: $S_F.REMOVE(n)$
5: **function** POP_BEST_NODE($reachable$)
6: $nodes \leftarrow S_F.INTERSECTION(reachable)$
7: **return** $\arg \max_{n \in nodes} n.uncertainty$

Data structure 5 BestFS graph

Require: $S_G \triangleright$ set of nodes
1: **function** ADD_NODE(n)
2: $S_G.ADD(n)$
3: **function** LEXICOGRAPHICAL_MAX($nodes$) \triangleright
 returns max according to keys: (solved, -n_visits, uncertainty)
4: **function** FIND_BEST_NODE($reachable, state$)
5: $nodes \leftarrow S.INTERSECTION(reachable)$
6: **return** LEXICOGRAPHICAL_MAX($nodes$)

A.1.2 MCTS WITH AVOID HISTORY COEFFICIENT DEAD ENDS AND AVOID SELF LOOPS

In Algorithm 6 we present the general structure of our MCTS planner. In Algorithms 7, 10, 8 and 9 we present our version of key functions in MCTS. In Algorithm 11 we provide pseudo code for ad auxiliary dead end mechanism used in this paper. Finally, in Data structure 12 we present the functionality of graph node used in this implementation.

Algorithm 6 MCTS planner

Require: $model$ learned model used in Algorithm 1
 V value function ensemble (assigns a vector of value estimates to state)
 C number of MCTS passes
 κ score factor for value uncertainty
 γ discount factor
 γ_{dead} dead end value
 γ_{ah} avoid history coefficient
Use: T planning graph
 $S_{visited}$ set of states visited by the agent
 S_{seen} set of states seen in one MCTS pass

function RESET
 $T \leftarrow \emptyset$
function CHOOSE_ACTION($state$)
 $root \leftarrow state$
 for $1 \dots C$ **do**
 MCTS_PASS($state$)
 return $\arg \max_{a \in A} (n.child(a).value)$
function MCTS_PASS($root$)
 $S_{seen} \leftarrow \emptyset$ \triangleright used in Algorithm 11
 $path, leaf \leftarrow TRAVERSAL(root)$
 $value \leftarrow EXPAND_LEAF(leaf, model)$
 BACKPROPAGATE($value, path$)
function GET_GRAPH_OF_PLANNING() **return** T

Algorithm 7 <code>traversal()</code> Input: <code>root</code> 1: <code>n</code> \leftarrow <code>root</code> 2: <code>path</code> \leftarrow \emptyset 3: while <code>n</code> is not a leaf do 4: <code>a</code> \leftarrow <code>select_child(n)</code> 5: if <code>a</code> is <code>None</code> then \triangleright dead end, terminal or leaf 6: break 7: <code>path.APPEND((n, a))</code> 8: <code>n</code> \leftarrow <code>n.child(a)</code> 9: return <code>path, n</code> $\triangleright n \notin \text{path}$	Algorithm 10 <code>select_child()</code> Input: <code>n</code> 1: if <code>n</code> not expanded or terminal then 2: return <code>None</code> 3: if <code>IS_DEAD_END(n)</code> then 4: return <code>None</code> 5: else 6: $\mathcal{A}_{\text{allowed}} \leftarrow \text{ALLOWED}(n)$ 7: return $\arg \max_{a \in \mathcal{A}_{\text{allowed}}} n.\text{child}(a).\text{EVAL}()$
Algorithm 8 <code>expand_leaf()</code> Require: $\gamma_{\text{dead}}, V, \text{model}$ Input: <code>leaf</code> \triangleright MCTS tree node without children 1: if <code>leaf</code> is terminal then 2: <code>UPDATE(leaf, 0.)</code> 3: return 0. 4: else if <code>IS_DEAD_END(leaf)</code> then 5: <code>UPDATE(leaf, γ_{dead})</code> 6: return γ_{dead} 7: else 8: for <code>a</code> $\in \mathcal{A}$ do 9: <code>new_node</code> \leftarrow <code>CREATE_NODE()</code> 10: <code>model.LOAD_STATE(leaf.state)</code> 11: <code>new_state</code> \leftarrow <code>model.STEP(a)</code> 12: <code>new_node.state</code> \leftarrow <code>new_state</code> 13: if <code>new_state</code> not yet visited then 14: <code>new_state.value</code> \leftarrow <code>V(new_state)</code> 15: <code>leaf.child(a)</code> \leftarrow <code>new_node</code> 16: return <code>leaf.value</code>	Algorithm 11 Dead end detection Require: S_{seen} Input: <code>v, path</code> 1: function <code>ALLOWED(n)</code> 2: $\mathcal{A}_{\text{allowed}} \leftarrow \emptyset$ 3: for <code>a</code> $\in \mathcal{A}$ do 4: if <code>n.child(a)</code> $\notin S_{\text{seen}}$ then 5: $\mathcal{A}_{\text{allowed}}.\text{ADD}(a)$ 6: return $\mathcal{A}_{\text{allowed}}$ 7: function <code>IS_DEAD_END(n)</code> 8: $\mathcal{A}_{\text{allowed}} \leftarrow \text{ALLOWED}(n)$ 9: if $\mathcal{A}_{\text{allowed}} \neq \emptyset$ then 10: return <code>False</code> 11: else 12: return <code>True</code>
Algorithm 9 <code>backpropagate()</code> Require: γ Input: <code>v, path</code> 1: for <code>(n, a)</code> in <code>reversed(path)</code> do 2: <code>v</code> \leftarrow <code>n.reward</code> $+$ γv 3: <code>n.UPDATE(v)</code>	Data structure 12 <code>node</code> Attributes: <code>count</code> MCTS node counter <code>value</code> value ensemble estimate <code>visits</code> number of agents visits <code>child</code> list of children nodes 1: function <code>UPDATE(value)</code> 2: <code>self.value</code> \leftarrow <code>self.value</code> $+$ <code>value</code> 3: <code>self.count</code> \leftarrow <code>self.count</code> $+$ 1 4: function <code>ADD_VISIT()</code> 5: <code>self.visits</code> \leftarrow <code>self.visits</code> $+$ 1 6: function <code>EVAL()</code> 7: <code>s</code> \leftarrow <code>SCORE(self.value, self.count)</code> 8: <code>std</code> \leftarrow <code>STD(self.value)</code> \triangleright uncertainty of value ensemble 9: return <code>score</code> $+$ γ_{ah} <code>self.visits</code> $+$ κstd

A.2 TRAINING SETUP

Code for all our experiments can be accessed at <https://github.com/ComradeMisha/TrustButVerify>.

Our experiments adhere to the general model-based training loop logic, described in Algorithm 1. We use a distributed system with 32 workers solving distinct episodes, where data gathered across a batch of workers is collected in two common experience buffers: the replay buffer for trainable model of size 50000 and the replay buffer for value function of size 30000. Before the solving of actual episode, we collect 1000 random trajectories, that we use for initial training of the model of

environment. This initial training has critical importance for the performance of our agents. Episodes are limited to: 600 steps for ToyMontezumaRevenge and 1000 steps for the Tower of Hanoi.

To update value network for MCTS we follow approach of Milos et al. (2019) (similar to Hamrick et al. (2020)) for BestFS experiments we simply calculate future discounted returns for each transition encountered by agent.

A.3 NETWORK ARCHITECTURES FOR MODEL AND VALUE

A.3.1 SINGLE NETWORKS

In every experiments we use two neural network architectures: one to estimate the value function of a given state and the other to predict the outcome of taking a step in the environment. In ToyMR we represent the state as a vector of length 17 and in the Tower of Hanoi as a vector of length $3 \times \text{number of discs}$ (it is due to one-hot encoding of the peg number for each disc).

In both environments, for value estimation we use an MLP (multilayer perceptron) architecture with two hidden layers of 50 neurons and ReLU non-linearity. For model, in both environments, we use an MLP with four hidden layers of 250 neurons and ReLU activation functions.

The model network has three outputs: the difference between next and current observations (delta target), reward and the episode termination flag (0 or 1).

A.3.2 ENSEMBLE

For both value estimation and model we use ensemble of networks with architectures described in Section A.3.1. The usage of ensembles is however different for value network and model network. Ensemble of model network is used as described in Algorithm 2. In case of value estimation, the final result of value is an average of predictions across ensemble. The standard deviation of this predictions multiplied by κ is used by MCTS as an auxiliary score added to each vertex in the search tree.

In Algorithm 13 we present the procedure of transforming the output of ensemble of model networks to a valid prediction of environment signal. In our experiments, it turned out to be beneficial to train the model to predict change between observations rather than the ready observations. The transformation performs the following step: collects predictions of all networks in ensemble, then averages next predicted change of observations, predicted rewards and predicted end of episode flags. Averaged observation change is added to current state, then clipped to the range of possible values (some coordinates of the state vector are binary variables), rounded to integers and returned as the next predicted state. Reward and end of episode flag are predicted to be true if their average value across networks in ensemble is larger than 0.5.

To stabilize performance of value and model ensembles we used additional masks mechanism. We create some number of networks (given by ensemble size parameter), that are trained with the same data buffer, but for prediction each worker uses only a random subset of given size (given by number of masks parameter).

A.4 QUANTILE RANK VALUE ANALYSIS

Figure 8 presents performance of BestFS Agent in the Tower of Hanoi environment for different values of Quantile Rank. We observed the the best performance was for values between 0.8 and 0.95, and when the parameter was within the appropriate range of values, the performance was not very sensitive to the changes of quantile rank parameter. We observed similar behaviour in ToyMR experiments - the best performance was for quantile rank between 0.9 and 0.95 but it was hard to narrow down this interval, due to aforementioned little sensitivity.

Algorithm 13 Transforming model ensemble predictions

```
function PREDICT_STEP(state, action)
  next_observations  $\leftarrow$  []
  rewards  $\leftarrow$  []
  is_done_flags  $\leftarrow$  []
  for network  $\in$  ensemble do:
    (next_obs_delta, reward, is_done)  $\leftarrow$  network.PREDICT(state, action)
    next_obs  $\leftarrow$  state + next_obs_delta
    next_observations.APPEND(next_obs)
    rewards.APPEND(reward)
    is_done_flags.APPEND(is_done)
  predicted_reward  $\leftarrow$  AVERAGE(next_reward)
  if predicted_reward > 0.5 then
    predicted_reward  $\leftarrow$  1
  else
    predicted_reward  $\leftarrow$  0
  predicted_is_done  $\leftarrow$  AVERAGE(is_done_flags)
  if predicted_is_done > 0.5 then
    predicted_is_done  $\leftarrow$  True
  else
    predicted_is_done  $\leftarrow$  False
  predicted_next_obs  $\leftarrow$  AVERAGE(next_observations)
  predicted_next_obs  $\leftarrow$  CLIP(predicted_next_obs)
  predicted_next_obs  $\leftarrow$  ROUND(predicted_next_obs)
  return(predicted_next_obs, predicted_reward, predicted_is_done)
```

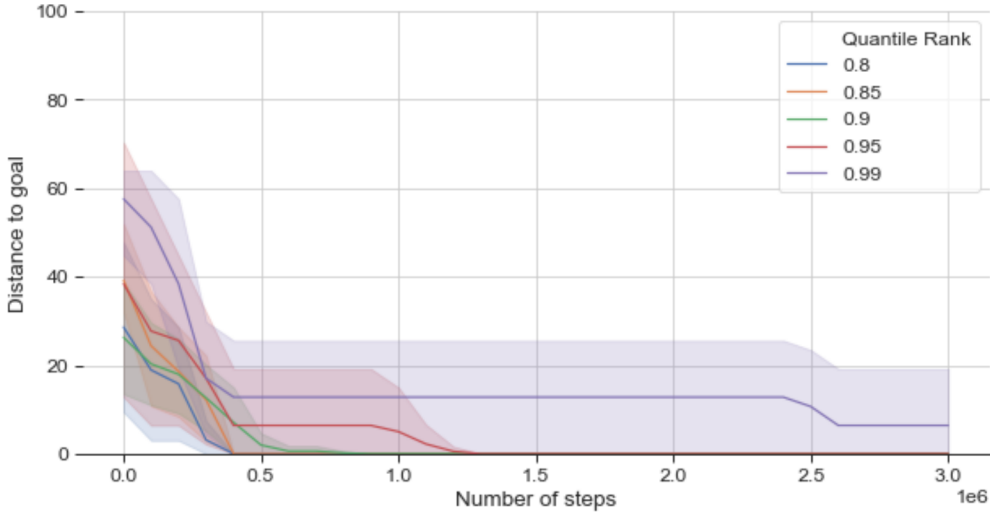


Figure 8: Performance of *TBV BestFS* on the Tower of Hanoi domain with different values of Quantile Rank

A.5 RANDOM TBV REJECTION ANALYSIS

In our preliminary experiments we have seen that without random rejection of TBV mechanism the agents often tends to get stuck in cycles with high ensemble disagreement and is unable to leave such cycle until the end of episode (since model is not updated in during the episode). In Figure 9 we present the performance in Tower of Hanoi for different values of frequency of random TBV

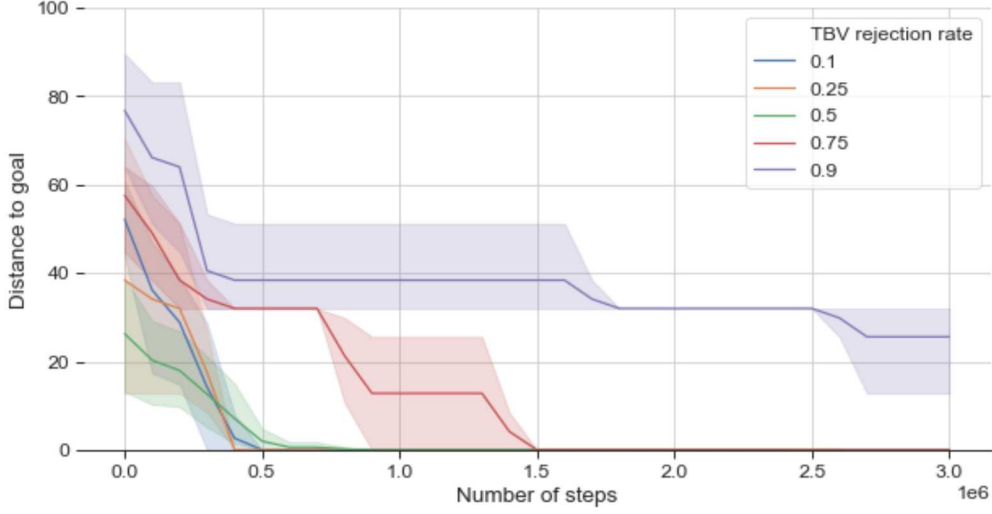


Figure 9: Performance of *TBV* BestFS on the Tower of Hanoi domain with different frequencies of random rejection of *TBV*

rejection (while the rest of parameters is same as in our best experiments). It can be seen, that if the frequency of *TBV* override (i.e. frequency with which we allow *TBV* to change planners action) is between 0.5 and 0.9 we obtained best results. Also, *TBV* is not very sensitive to this parameter as long it is in an appropriate range of values.

A.6 HYPER-PARAMETERS

In tables 1 and 2 we present hyper-parameters used in our experiments.

Parameter	Toy MR		Tower of Hanoi	
	BestFS	MCTS	BestFS	MCTS
Number of planner passes C	10	10	10	10
Discounting factor γ	0.99	0.99	0.99	0.99
ε greedy exploration for baselines	0.001	0.02	0.001	0.02
score factor for value uncertainty κ	-	3	-	3
Dead end value	-	-0.2	-	-0.2
Avoid history coefficient	-	-0.2	-	-0.2
Value ensemble size ²	20	20	20	20
Value ensemble mask size ²	10	10	10	10
Model ensemble size ²	8	8	8	8
Model ensemble mask size ²	4	4	4	4
Optimizer	RMSprop	RMSprop	RMSprop	RMSprop
Learning rate	2.5e-4	2.5e-4	2.5e-4	2.5e-4
Batch size for value training	32	32	32	32
Batch size for model training	1024	1024	1024	1024

² See Section A.3.2 for explanation of these parameters.

Table 1: Hyper-parameters values used in our experiments.

For MCTS we took hyperparameters from Milos et al. (2019). As they used setup without learned model we needed to tune model architecture and model training parameters on ToyMontezumaRe-

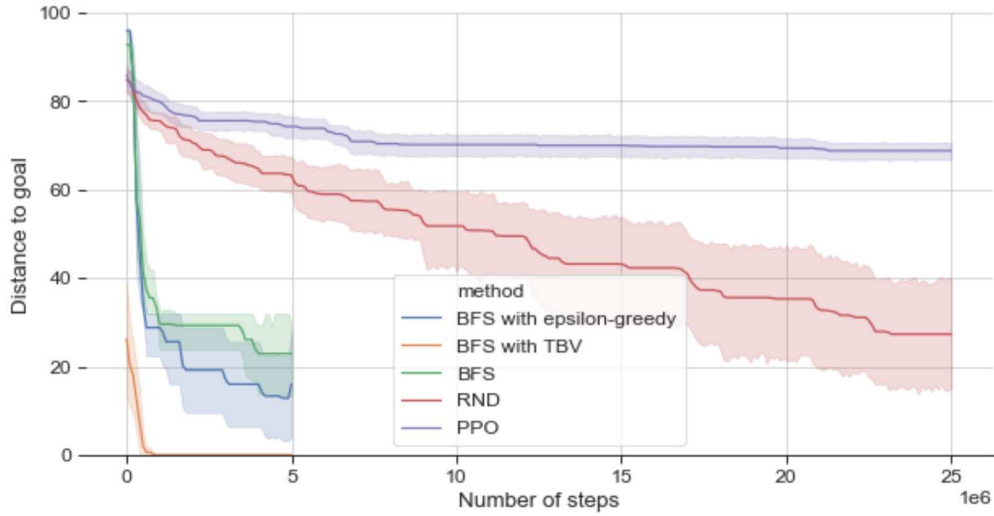


Figure 10: Performance of RND and PPO on Tower of Hanoi with longer training.

venge domain. Then, we separately tuned epsilon and Quantile Rank for each combination of domain/planner.

A.7 RND AND PPO

RND continued to further explore state space for Hanoi (figure 10), for comparison we also show performance of PPO (the same parameters as RND but without intrinsic reward).

Parameter	ToyMR	Tower of Hanoi
Rollout length	128	128
Maximal length of an episode	600	1000
Total number of rollouts per environment	6200	6200
Number of minibatches	4	4
Number of optimization epochs	16	4
Coefficient of extrinsic reward	1	1
Coefficient of intrinsic reward	100	10
Number of parallel environments	32	32
Learning rate	0.001	0.001
Optimization algorithm	Adam	Adam
λ	0.95	0.95
Entropy coefficient	0.001	0.001
Proportion of experience used for training predictor	1.0	1.0
γ_E	0.999	0.999
γ_I	0.99	0.99
Clip range	[0.9, 1.1]	[0.9, 1.1]
Policy architecture	FCN	FCN

Table 2: Default hyper-parameters for PPO and RND algorithms for ToyMR and Tower of Hanoi experiments.

Parameter	Value
Number of optimization epochs	[1, 4, 16]
Coefficient of intrinsic reward	[1, 3, 10, 30, 100, 300]
Learning rate	$[5 \cdot 10^{-2}, 10^{-2}, 5 \cdot 10^{-3}, 10^{-3}, 5 \cdot 10^{-4}, 10^{-4}]$
λ	[0.95, 0.99]
Proportion of experience used for training predictor	[0.25, 1.0]
γ_E	[0.999, 0.9999]
γ_I	[0.99, 0.999]
Policy architecture layers number	[2, 3, 4]
Policy architecture layers width	[64, 128, 256]
Random target and prediction networks last layer width	[64, 128]

Table 3: Hyper-parameters for RND algorithm checked during tuning process of ToyMontezumaRevenge.