# Delta Schema Network
# in Model-based Reinforcement Learning

Andrey Gorodetskiy[1,2], Alexandra Shlychkova[1], and Aleksandr I. Panov[1,3]

[1] Moscow Institute of Physics and Technology (National Research University),
Moscow, Russia
[2] Bauman Moscow State Technical University, Moscow, Russia
[3] Artificial Intelligence Research Institute, Federal Research Center "Computer
Science and Control" of the Russian Academy of Sciences, Moscow, Russia
`gorodetskiyandrew@gmail.com`

**Abstract.** This work is devoted to unresolved problems of Artificial
General Intelligence - the inefficiency of transfer learning. One of the
mechanisms that are used to solve this problem in the area of reinforce-
ment learning is a model-based approach. In the paper we are expanding
the schema networks method which allows to extract the logical rela-
tionships between objects and actions from the environment data. We
present algorithms for training a Delta Schema Network (DSN), predict-
ing future states of the environment and planning actions that will lead
to positive reward. DSN shows strong performance of transfer learning
on the classic Atari game environment.

**Keywords:** Reinforcement learning · Model-based · Schema Network ·
Delta Schema Network · Transfer Learning

## 1 Introduction

For an intelligent agent acting in real-world conditions, it is necessary to general-
ize the experience gained in order not to learn from scratch after a slight change
in the environment. A human does not relearn the policy of interaction with a
familiar object, but only slightly corrects it, when object's characteristics are
changed. For this, logical relationships between objects and their characteristics
are used at different levels of generalization. For example, in the Atari game
*Breakout*, the colors of the bricks do not matter and the natural agent does not
change the policy when colors change. Artificial agent can achieve such a gen-
eralization using some universal model-based learning algorithm. In this paper,
we propose a new approach to the learning of universal models for reinforcement
learning in game environments - Delta Schema Network (DSN) - which is an
extension of the early work Schema Network [5].

A Schema Network is an object-oriented model, the main aspect of which is
a schema. In this architecture the agent receives an image from the environment,
which is parsed into a set of extracted objects. Then model learns a set of rules -

schemas, which reflect logical interconnection between objects' properties, actions and rewards. Each schema predicts some property of the objects of certain type, using information about properties of other objects and actions from past observations.

It is possible to represent this interconnection as a factor graph. A variable node in this graph is either a property of some object, potentially achievable reward or action. Factors are schemas, that have input and output nodes. The edges indicate the presence of a causal relationship between the objects and events. An agent can find a node with a positive reward in future time layers of this graph and plan actions to reach it.

Since the graph has fairly generalized properties, the trained Schema Network can be used in conjunction with some feature extractor on environments with similar interaction dynamics. This provides advantage in transfer learning. However, planning on a sufficiently large graph for a long time horizon can be a very challenging task for real-world applications.

From the point of view of creating AGI systems, the DSN algorithm can be used to automatically generate scripts for the behavior of an intelligent agent. These scripts can be used to speed up the agent's own behavior planning process [9,7], or to predict user behavior in a cognitive assistant scenario [12].

## 2   Related work

Representation of logical interconnection often helps to increase an efficiency of transfer learning. Various methods are used to represent the logical relationships of objects: in the Schema Network [5], these are specially introduced schemas with binary logic. In Logical Tensor Networks [10], it is proposed to use real logic. To further apply the obtained relationships for planning, one can use them as additional data for a neural network. For example, in [13] schemas are passed to a neural network. Authors in [2] add logical relationships to the input of a neural network using logical tensor network. Another approach is to build a dependency graph and search for a reachable state with a positive reward.

The usage of the Schema Network in reinforcement learning consists of two main stages: training a network to predict future states of environment and construct a prediction graph with the subsequent search for the best reachable reward node.

The Schema Network uses object-oriented approach described in [4]. During the training stage, due to the knowledge of the types of environment objects, a model is able to identify the logical connections between them. A similar approach was used in Interaction Network [3], for which, however, no planning algorithms were developed to obtain a positive reward. For similar problems convolutional neural networks are used as in [6]. However, this approach requires a prior knowledge about the structure of the graph, while the Schema Network allows to obtain knowledge about relations between objects automatically from the environment.

Also, during Schema Network training stage a dependency graph is constructed. Finding the reachable state of the environment in which a reward is received can be considered as a estimate of the posterior maximum and solved using the max-product belief propagation [1].

## 3   Model description

### 3.1   Main concepts

Key concepts used in the Delta Schema Network (DSN) are entities, attributes and schemas. Entity is any object that can be extracted from the image. Attribute is a binary variable that indicates presence or absence of a specific property of an entity, each entity has the same $M$ number of attributes. Attribute with value of 1 or $True$ is said to be *active*. Schema is a logical AND function that predicts value of attribute or reward at time step $t$, taking as arguments arbitrary number $k$ of attributes and actions at previous time steps $\{t^* : 0 < t - t^* < d\}$.

$$Schema\colon (Attributes_{t^*} \cup Actions_{t^*})^k \to Attributes_t \cup Rewards_t$$

Schemas are represented as binary column vectors and forms parameter matrices. We define $W = (W_i : i = 1..M)$ to be a tuple of parameter matrices used for attribute prediction, one matrix per attribute type. Parameter matrix used for reward prediction we denote as R.

DSN model learns dynamics of the environment in terms of schema vectors and, from some point of view, represents both transition and reward functions of the environment. Using learned vectors, model predicts next states of the environment and plan a sequence of actions that will lead to reward.

### 3.2   State representation

In our work we considered each pixel of image as entity. However, we think this model is more suitable for reasoning on more high-level concepts. Attributes of entities have meaning of presence in this pixel object of a certain type, i.e entity's attribute vector is one-hot encoded type of this entity concatenated with void attribute that can indicate absence of any object in this pixel.

DSN relies on semantic information about observation from environment, namely which type of object each pixel belongs to. As observation at time step $t$ model gets state matrix $s_t$ of $(N, M)$ shape, where $N$ is the number of entities an $M$ is the number of attributes. This matrix is suggested to be built from image of $N$ pixels and $M - 1$ object types. We consider $s_t$ is provided by some feature extractor.

### 3.3   Prediction

Schema vectors are used to predict *changes* (deltas) in the current state $s_t$ and to predict reward $r_t$ after taking action $a_t$. If there are several schemas that predict

same attribute, their results are united using logical OR. Two types of schemas are used for attribute prediction: creating, represented by $W^+$, and destroying, represented by $W^-$. Creating schemas predict attributes that are not active in $s_t$, but should be active in $s_{t+1}$. Vice versa, destroying schemas predict destruction of attributes that are active in $s_t$, but should disappear in $s_{t+1}$.

We used $d = 2$, i.e. for computing attribute value at $t + 1$ schema can use attributes and actions only at $t$ and $t - 1$. Thus, frame stack has size 2.

To make a prediction on either attributes or rewards we construct *augmented matrix* $X_t$, which is built from the frame stack $(s_{t-1}, s_t)$ and action $a_t$ in the following way:

1. $s_{t-1}$ and $s_t$ are augmented into $s_{t-1}^*$ and $s_t^*$, correspondingly. Each row, which is an attribute vector of some entity, is horizontally concatenated with attribute vectors of these entity's $R - 1$ spatial neighbors. Referring to corresponding image, these neighbors are located in square around the central pixel, and the central pixel is represented by a row in $s$.
2. $a^*$ is built, which is broadcasted by number of rows to $s_t$ version of one-hot encoded action $a_t$.
3. horizontally concatenated $(s_{t-1}, s_t, a^*)$ result in $x_t$.

To predict next state $s_{t+1}$, we predict creating $\Delta^+$ and destroying $\Delta^-$ state changes using $W$ and then apply them to $s_t$:

$$\Delta_j^+ = \overline{\overline{X_t W_j^+}} \mathbf{1} \qquad \Delta_j^- = \overline{\overline{X_t W_j^-}} \mathbf{1},$$

where $\Delta_j$ denotes $j$th column of $\Delta$.

$$s_{t+1} = s_t - \Delta^- + \Delta^+,$$

considering elements as integers and clipping result after every operation in chain to $\{0, 1\}$.

Reward is predicted in similar way:

$$r_{t+1} = \mathbf{1}^\intercal \overline{\overline{X_t R}} \mathbf{1}$$

This matrix multiplication of augmented matrix and parameter matrix is equivalent to applying discrete convolutions to the original image, where parts of images are described by rows of augmented matrix and filters are column vectors in parameter matrix. Thus, DSN models the environment as cellular automaton - grid of entities - and reconstructs its rules as schema vectors.

## 4   Learning algorithm

During interaction with the environment agent stores unique transitions in replay buffer. We use learning algorithm from [5] with different target in a self-supervised manner. First, correctness of already learned schema vectors is checked on new observations. Schema vectors that produce false positive predictions are

deleted. After that, we learn new schema vectors. Targets for learning $W$ are columns of $\Delta^+$ and $\Delta^-$, which we denote as $y$. Target for learning $R$ is the reward, obtained at sample's time step.

1. We choose one random sample with false negative prediction from replay buffer and put it in the set *solved*.
2. We solve the following LP optimization problem: finding a schema vector $w$ that does not produce any false positive predictions on replay buffer, predicts positive labels for all samples in *solved* and maximizes the number of true positive predictions on replay buffer.

$$\min_{w \in \{0,1\}^D} \sum_{n:y_n=1} (1 - x_n)w$$
$$\text{s.t. } (1 - x_n)w > 1 \quad \forall_{n:y_n=0}$$
$$(1 - x_n)w = 0 \quad \forall_{n \in \text{solved}}$$

3. All samples that got predicted by obtained schema vector are added to the set *solved*.
4. We try to simplify schema vector: minimizing its $L_1$ norm with condition of absence false positive predictions on replay buffer and false negative on set *solved*.

$$\min_{w \in \{0,1\}^D} w^T \mathbf{1}$$
$$\text{s.t. } (1 - x_n)w > 1 \quad \forall_{n:y_n=0}$$
$$(1 - x_n)w = 0 \quad \forall_{n \in \text{solved}}$$

## 5   Planning algorithm

The purpose of planning is to find a sequence of actions that will lead to positive reward. The planning process consists of several stages:

1. Forward pass builds factor graph of potentially reachable nodes;
2. A set of target reward nodes are selected;
3. Sequence of actions that will activate target node are planned.

The input to the planner is the frame stack of size 2 consisting of state matrices $(s_{t-1}, s_t)$ and schema parameters $(W^+, W^-, R)$.

### 5.1   Forward pass

The future states of the environment are predicted for $T$ time steps ahead. Simultaneously, we build the factor graph $G$ in which variable nodes are attributes,

rewards or actions; factors are schemas that were activated during prediction process and edges connect schemas to their input and output nodes. Every node has assigned time step, at which it appeared in prediction. Thus, graph $G$ is said to be consisting of layers, that unite nodes within same time step.

To predict next state one need to decide which action $a$ agent takes at current state. When predicting $\Delta^+$, DSN model assumes that agent takes all possible actions, and for $\Delta^-$ it assumes agent takes "do not do anything" action. This leads to superimposing of all possible $\Delta^+$ for the next state in the single matrix $s_{t+1}$.

During forward pass we maintain graph building in the following way. After predicting $s_{t+1}$, for each predicted attribute or reward node at layer $t+1$ we instantiate on the graph concrete instances of corresponding *creating* schemas. Each attribute at $s_t$ that was not destroyed by any of the *destroying* schemas is considered to be active at $s_{t+1}$ and we mark the corresponding node at $t+1$ as having self-transition, that is like a schema with single input that activates attribute at time $t+1$ provided it was active at $t$.

### 5.2 Target nodes selection

Having predictions for the future states of the environment on $T$ ticks ahead, reward nodes of the factor graph are added to the target queue.

$q$ = sorted by time potentially reachable positive reward nodes
$= [r^+_{closest} \dots r^+_{farthest}]$

### 5.3 Finding sequence of actions

We take next reward node from queue $q$ and try to find a sequence of actions for agent to reach it. To find such sequence, one need to find a configuration of graph $G$, that satisfy following constraints:

– target reward node is active
– at each layer $t$ only *one* action node is active

In this configuration, the values of the attribute and reward nodes show their actual reachability. The action nodes $\{a_i \in G : i \in [1..T]\}$ represent the actions that must be taken to reach the target node.

Some of the learned schema vectors may depend on actions, while in the dynamics of the environment there is no such dependence. This occurs because during training events correlated, but did not have a causal relationship. For correct planning, it is necessary to find a valid configuration of the graph, constructed by predictions with such vectors. We propose the `backtrace_node` algorithm (algo. 3) to find such a configuration. It does not perform exhaustive search, but works well in our experiments.

We maintain an array of joint constraints on the active action nodes for each time layer. During graph traversal, we either satisfy these constraints or replan paths to nodes committed to these constraints if there is no other path to the target. Process of node activation goes in the following order:

- try to activate the node by self-transition
- try to activate the node with an action-independent schema
- if there is no constraint on the current tick, try to activate the node with any schema
- try to select a schema that satisfies the constraint on the current tick
- replan all vertices that require current constraint
    - find a set of actions that as constraints would allow the activation of each conflicting node
    - sequentially start replanning subgraph of each conflicting node using the action acceptable by all
    - if all nodes have been replanned successfully, change the constraints at conflicting layer

During the replanning process, a new conflict situation may arise. Then new replanning process should be recursively started.

## 6    Experiments

The model was evaluated on the Atari Breakout game (see Figure 1). The goal of the game is to knock down bricks with a ball, substituting a moving platform under it. There are no random factors in the environment.

The action space consists of the following actions: do not move, move left, move right. As an observation, the agent receives an RGB image and information about a particular type of object each pixel belongs to. Rewards are distributed as follows: +1 for knocking down a brick, -1 for dropping a ball past the platform, 0 in other cases.

The number of schema vectors for each parameter matrix was limited to 500 units. The episode was limited to 5000



**Fig. 1.** Breakout

steps, agent had 3 lives after the loss of which the episode ended. Highest possible reward for episode was 36. Figure 2 shows the results of DSN evaluation.

Agent did not managed to knock down completely all bricks in part of episodes, because after destroying some part of them to the top wall, it could not longer detect future reward and hence plan actions.

A distinctive feature of the DSN is the efficient transfer of the trained model to environments with similar dynamics. We evaluated the model, trained in the previous experiment, on the same environment but with two balls. Results of transfer without additional training (see Figure 3) show similar average score.

We compared DSN model to DQN. Figure 4 shows that DQN needs significantly more time steps to reach equal performance.

## 7    Conclusion

In this paper, we proposed an original implementation of the universal logical model of environment dynamics for model-based reinforcement learning. Our approach, which we called Delta Schema Network, is a modification and extension
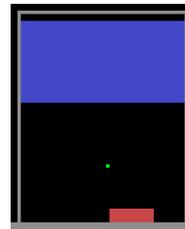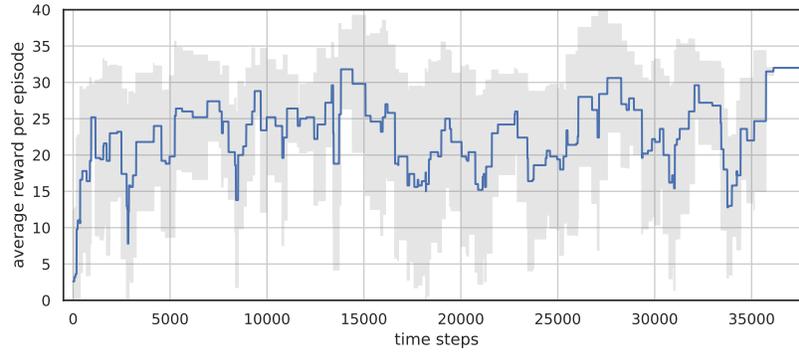
**Fig. 2.** Total reward per episode of DSN on standard Breakout. Averaged over 5 runs, shaded region represents the standard deviation.
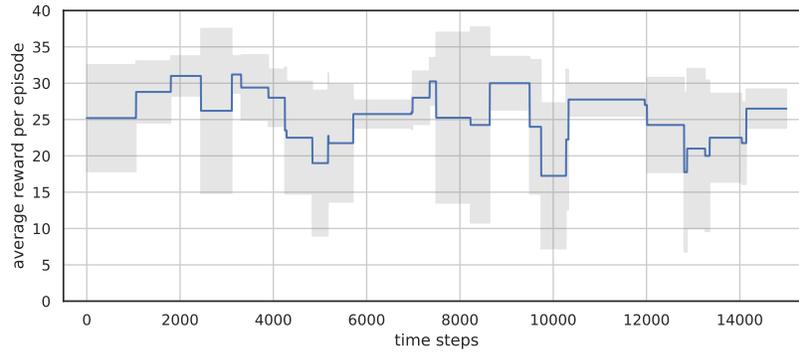


**Fig. 3.** DSN performance after zero-shot transfer. Averaged over 5 runs, shaded region represents the standard deviation.
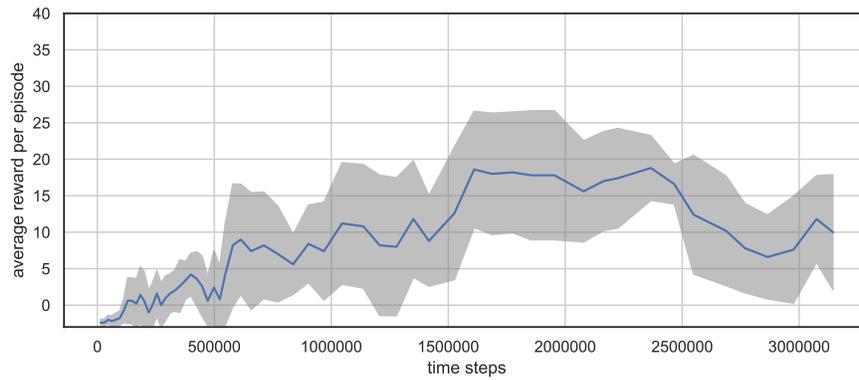


**Fig. 4.** DQN training process on standard Breakout. Single run, rolling mean with window size = 5. Shaded region represents standard deviation of window samples.

of Schema Network for RL. We described in detail the algorithmic implementation of the proposed method and conducted basic experimental studies on the Atari Breakout environment.

Future works include the use of logical model-based approaches for real-world robotic tasks, such as controlling a robotic manipulator [8,14] or a car at an road intersection [11]. Code of the DSN model can be obtained in the repository: github.com/cog-isa/schema-rl.

# References

1. Attias, H.: Planning by probabilistic inference. In: AISTATS. Citeseer (2003)
2. Badreddine, S., Spranger, M.: Injecting prior knowledge for transfer learning into reinforcement learning algorithms using logic tensor networks. arXiv preprint arXiv:1906.06576 (2019)
3. Battaglia, P., Pascanu, R., Lai, M., Rezende, D.J., et al.: Interaction networks for learning about objects, relations and physics. In: Advances in neural information processing systems. pp. 4502–4510 (2016)
4. Diuk, C., Cohen, A., Littman, M.L.: An object-oriented representation for efficient reinforcement learning. In: Proceedings of the 25th international conference on Machine learning. pp. 240–247 (2008)
5. Kansky, K., Silver, T., Mély, D.A., Eldawy, M., Lázaro-Gredilla, M., Lou, X., Dorfman, N., Sidor, S., Phoenix, S., George, D.: Schema networks: Zero-shot transfer with a generative causal model of intuitive physics. In: Proceedings of the 34th International Conference on Machine Learning-Volume 70. pp. 1809–1818. JMLR. org (2017)
6. Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. arXiv preprint arXiv:1609.02907 (2016)
7. Kiselev, G., Panov, A.: Hierarchical Psychologically Inspired Planning for Human-Robot Interaction Tasks. In: Ronzhin, A., Rigoll, G., Meshcheryakov, R. (eds.) Interactive Collaborative Robotics. ICR 2019. Lecture Notes in Computer Science. vol. 11659, pp. 150–160. Springer (2019). https://doi.org/10.1007/978-3-030-26118-4_15
8. Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G., Graves, A., Riedmiller, M.A., Fidjeland, A.K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., Hassabis, D.: Human-level control through deep reinforcement learning. Nature **518**, 529–533 (2015)
9. Panov, A.I., Yakovlev, K.S.: Psychologically Inspired Planning Method for Smart Relocation Task. In: Procedia Computer Science. vol. 88, pp. 115–124. Elsevier (2016). https://doi.org/10.1016/j.procs.2016.07.414
10. Serafini, L., Garcez, A.d.: Logic tensor networks: Deep learning and logical reasoning from data and knowledge. arXiv preprint arXiv:1606.04422 (2016)
11. Shikunov, M., Panov, A.I.: Hierarchical Reinforcement Learning Approach for the Road Intersection Task. In: Samsonovich, A.V. (ed.) Biologically Inspired Cognitive Architectures 2019. BICA 2019. Advances in Intelligent Systems and Computing. vol. 948, pp. 495–506. Springer (2020). https://doi.org/10.1007/978-3-030-25719-4_64

12. Smirnov, I., Panov, A.I., Skrynnik, A., Isakov, V., Chistova, E.: Personal Cognitive Assistant: Concept and Key Principals. Informatika i ee Primeneniya **13**(3), 105–113 (2019). https://doi.org/10.14357/19922264190315
13. Toyer, S., Trevizan, F., Thiébaux, S., Xie, L.: Action schema networks: Generalised policies with deep learning. In: Thirty-Second AAAI Conference on Artificial Intelligence (2018)
14. Younes, A., Panov, A.I.: Toward Faster Reinforcement Learning for Robotics : Using Gaussian Processes. In: Osipov, G.S., Panov, A.I., Yakovlev, K.S. (eds.) RAAI Summer School 2019. Lecture Notes in Computer Science, vol. 11866, pp. 160–174. Springer (2019). https://doi.org/10.1007/978-3-030-33274-7_11

# 8 Appendix

Algorithms 1 and 2 are used in algorithm 3. Node in the graph is considered to have next attributes:

- node.is_reachable - the actual reachability of the node, subject to currently selected actions, or `None` if the reachability is not known.
- node.schemas - map from actions to node's schemas requiring these actions
- node.transition - self-transition node, if any

---

**Algorithm 1:** backtrace_node_by_schemas

---

**Input** : node - target node
            schemas - set of available schemas
**Output:** actual node reachability, planned actions

---

1 **for** schema in schemas **do**
2     backtrace_schema(schema)
3     **if** schema.is_reachable **then**
4        node.is_reachable ← True
5        break
6 **end**

---

---

**Algorithm 2:** backtrace_schema

---

**Input** : schema - target schema
            preconditions - input nodes of schema
**Output:** actual schema reachability

---

1 schema.is_reachable ← True
2 **for** precondition in preconditions **do**
3     **if** precondition.is_reachable *is None* **then**
4        backtrace_node(precondition)
5     **if** *not* precondition.is_reachable **then**
6        schema.is_reachable ← False
7        break
8 **end**

---

---

**Algorithm 3:** backtrace_node(node, desired_constraint=None)

---

**Input :**
- node - target node to backtrace
- desired_constraint=None - desired constraint at node.t − 1 to satisfy, if any

**Output:** actual node reachability, planned actions in joint_constraints

**1** node.is_reachable ← *False*
**2** **if** desired_constraint *is not None* **then**
**3**   │  backtrace_node_by_schemas(node, schemas[desired_constraint])
**4**   │  **return**
  // try to activate the node by self-transition
**5** **if** node.transition.is_reachable *is None* **then**
**6**   │  backtrace_node(node.transition)
**7** node.is_reachable ← node.transition.is_reachable
**8** **if** node.is_reachable **then return**
  // try to activate the node with an action-independent schema
**9** backtrace_node_by_schemas(node, schemas[All action independent])
**10** **if** node.is_reachable **then return**
**11** **if** no current constraint **then**
**12**   │  backtrace_node_by_schemas(node, schemas[All action dependent])
**13**   │  **if** node.is_reachable **then**
**14**   │    │  set new constraint for current layer in joint_constraints
**15**   │  **return**
  // try to select a schema that satisfies the current constraint
**16** backtrace_node_by_schemas(node, schemas[current_constraint])
**17** **if** node.is_reachable **then**
**18**   │  add current node as committed to current constraint
**19**   │  **return**
  // replan all vertices that require current constraint
**20** negotiated_actions ← actions acceptable by all conflicting nodes
**21** is_success = False
**22** **for** action in negotiated_actions **do**
**23**   │  backtrace_node(node, desired_constraint=action)
**24**   │  **if** node.is_reachable **then**
**25**   │    │  is_success = True
**26**   │    │  **for** curr_node in committed_nodes **do**
**27**   │    │    │  backtrace_node(curr_node, desired_constraint=action)
**28**   │    │    │  **if** *not* curr_node.is_reachable **then**
**29**   │    │    │    │  curr_node.is_reachable = True
**30**   │    │    │    │  is_success = False
**31**   │    │    │    │  break
**32**   │    │  **end**
**33**   │  **if** is_success **then**
**34**   │    │  change constraints for current layer to new ones
**35**   │    │  break
**36** **end**

---