# Online 3D Bin Packing with Constrained Deep Reinforcement Learning

**Hang Zhao**[1,*]     **Qijin She**[1,*]     **Chenyang Zhu**[1]     **Yin Yang**[2]     **Kai Xu**[1]
[1]National University of Defense Technology    [2]Clemson University

## Abstract

We solve a challenging yet practically useful variant of 3D Bin Packing Problem (3D-BPP). In our problem, the agent has limited information about the items to be packed into the bin, and an item must be packed immediately after its arrival without buffering or readjusting. The item's placement also subjects to the constraints of collision avoidance and physical stability. We formulate this *online 3D-BPP* as a constrained Markov decision process. To solve the problem, we propose an effective and easy-to-implement constrained deep reinforcement learning (DRL) method under the actor-critic framework. In particular, we introduce a feasibility predictor to predict the feasibility mask for the placement actions and use it to modulate the action probabilities output by the actor during training. Such supervisions and transformations to DRL facilitate the agent to learn feasible policies efficiently. Our method can also be generalized e.g., with the ability to handle lookahead or items with different orientations. We have conducted extensive evaluation showing that the learned policy significantly outperforms the state-of-the-art methods. A user study suggests that our method attains a human-level performance.

## 1  Introduction

As a classic NP-hard problem, the bin packing problem (1D-BPP) seeks for an assignment of a collection of items with various weights to bins. The optimal assignment houses all the items with the fewest bins such that the total weight of items in a bin is below the bin's capacity $c$ [1]. In its 3D version i.e., 3D-BPP, an item $i$ has a 3D "weight" corresponding to its length, $l_i$, width $w_i$, and height $h_i$. Similarly, $c$ is also in 3D including $L \geq l_i$, $W \geq w_i$, and $H \geq h_i$. It is assumed that $l_i, w_i, h_i, L, W, H \in \mathbb{Z}^+$ are positive integers. Given the set of items $\mathcal{I}$, we would like to pack all the items into as few bins as possible. Clearly, 1D-BPP is a special case of its three dimensional counter part – as long as we constrain $h_i = H$ and $w_i = W$ for all $i \in \mathcal{I}$, a 3D-BPP instance can be relaxed to a 1D-BPP. Therefore, 3D-BPP is also highly NP-hard [2].

Regardless of the difficulty of the problem itself, the bin packing problem turns out to be one of the most needed academic problems [3] (the second most needed, only after the suffix tree problem) as many real-world challenges could be much more efficiently handled if we have a good solution to it. A good example is large-scale parcel packaging in modern logistics systems (Figure. 1), where parcels are mostly in regular cuboid shapes, and we would like to collectively pack them into rectangular bins of the standard dimension. Maximizing the storage use of bins effectively reduces the cost of inventorying, wrapping, transportation,



Figure 1: 3D-BPP is widely useful in logistics, manufacture, warehousing etc.

and warehousing. While being strongly NP-hard, 1D-BPP has been extensively studied. With the state-of-the-art computing hardware, big 1D-BPP instances (with about $1,000$ items) can be exactly
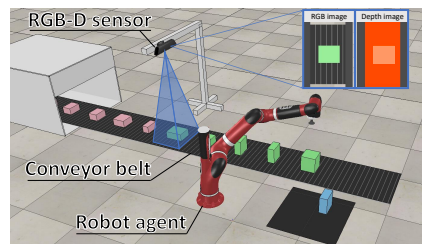
---

[*]Joint first authors

solved within tens of minutes [4] using e.g., integer linear programming (ILP) [5], and good approximations can be obtained within milliseconds. On the other hand 3D-BPP, due to the extra complexity imposed, is relatively less explored. Solving a 3D-BPP of moderate size exactly (either using ILP or branch-and-bound) is much more involved, and we still have to resort to heuristic algorithms [6, 7].

Most existing 3D-BPP literature assumes that the information of all the items is known, and the packing strategies allow backtracking i.e., one can always repack an item from the bin in order to improve the current solution [8]. In practice however, we do not know the information of all the items. For instance see Figure. 1, where a robot works beside a bin, and a conveyor forwards parcels sequentially. The robot may only have the vision of several upcoming items (similar to Tetris), and an item must be packed within a given time period after its arrival. It is costly and inefficient if the robot frequently unloads and readjusts parcels in packed bins. Such constraints further complicate 3D-BPP in its real-world applications.

As an echo to those challenges, we design a deep reinforcement learning (DRL) algorithm for 3D-BPP. To maximize the applicability, we carefully accommodate restrictions raised in its actual usage. For instance, we require each item placement being collision-free and not inducing instable stacking. An item is immediately packed upon its arrival, and no adjustment will be permitted after it is packed. To this end, we opt to formulate our problem as a constrained Markov decision process (CMDP) [9] and propose a constrained DRL approach based on the on-policy actor-critic framework [10, 11]. Concretely, we introduce a feasibility predictor, which is a CNN module predicting the feasibility mask for the item to be placed over the discretized parameterization of the bin. Meanwhile, we use the predicted feasibility to modulate the action probabilities output by the actor, and explicitly minimize the summed probabilities of all infeasible actions. These supervisions facilitate the agent to learn feasible policy efficiently, with an easier implementation than constrained policy optimization [12]. After a thorough testing and validation, we find that our algorithm outperform existing methods by a noticeable margin. It even demonstrates a human-level performance in a preliminary user study.

## 2   Related Work

**1D-BPP** is one of the most famous problems in combinatorial optimization, and related literature dates back to the sixties [13]. Many variants and generalizations of 1D-BPP arise in practical contexts such as the cutting stock problem (CSP), in which we want to cut bins to produce desired items of different weights, and minimize the total number of bins used. A comprehensive list of bibliography on 1D-BPP and CSP can be found in [14] including over 400 books, articles, dissertations, and working papers. Knowing to be strongly NP-hard, most existing literature focuses on designing good heuristic and approximation algorithms and their worse-case performance analysis [15]. For example, the well-known greedy algorithm, the next fit algorithm (NF) has a linear time complexity of $\mathbf{O}(N)$ and its *worst-case performance ratio* is 2 i.e. NF needs at most twice as many bins as the optimal solution does [16]. The first fit algorithm (FF) allows an item to be packed into previous bins that are not yet full, and its time complexity increases to $\mathbf{O}(N \log N)$. The best fit algorithm (BF) aims to reduce the residual capacity of all the non-full bins. Both FF and BF have a better worse-case performance ratio of $\frac{17}{10}$ than NF [17]. Pre-sorting all the items yields the off-line version of those greedy strategies sometimes also known as the *decreasing* version [18]. While straightforward, NF, FF, and BF form a foundation of more sophisticated approximations to 1D-BPP (e.g. see [19]) or its exact solutions [20–23]. We also refer the reader to BPPLib library [24], which includes the implementation of most known algorithms for the 1D-BPP problem.

**High-dimension BPP** such as 2D-BPP and 3D-BPP is a natural generation of the original bin packing problem. Here, an item does not only have a scalar-valued weight but a high-dimension size of width, height, and/or depth. The main difference between 1D- and 2D-/3D- packing problems is the verification of the feasibility of the packing, i.e. determining whether an accommodation of the items inside the bin exists such that items do not interpenetrate and the packing is within the bin size. The complexity and the difficulty significantly increase for high-dimension BPP instances. In theory, it is possible to generalize exact 1D solutions like MTP [20] or branch-and-bound [23] algorithms to 2D-BPP [25] and 3D-BPP [8]. However according to the timing statistic reported in [8], exactly solving 3D-BPP of a size matching an actual parcel packing pipeline, which could deal with tens of thousand parcels, remains infeasible. Resorting to approximation algorithms is a more practical choice for us. Hifi and colleagues [26] proposed a mixed linear programming algorithm for 3D-BPP by relaxing the integer constraints in the problem. Crainic and colleagues [6] refined the idea of corner points, where an upcoming item is placed [8], to the so-called *extreme points* to better explore

the un-occupied space in a bin. The local search heuristic iteratively improves an existing packing by searching within a neighbourhood function over the set of solutions, and it has been a popular strategy to design fast approximation algorithms for 2D-/3D-BPP such as the guided local search[27], greedy search [28], and tabu search [29, 30]. Compared to local search, genetic algorithms could generate better solutions thanks to the mechanism of mutation and crossover [31, 32].

**Deep reinforcement learning (DRL)** has demonstrated tremendous success in learning complex behaviour skills and solving challenging control tasks with high-dimensional raw sensory state-space [33, 34, 10]. The success can be attributed to the utilization of high-capacity deep neural networks for powerful feature representation learning and function approximation. Since the seminal work of deep Q-network (DQN) [34], a large body of literature has emerged. The existing research is largely divided into two lines: value function learning [34, 35] and policy search [36, 37]. Actor-critic methods, designed to combine the two approaches, have grown in popularity. Asynchronous Advantage Actor-Critic (A3C) [10] is a representative actor-critic method. It combines advantage updates with the actor-critic formulation and adopts asynchronously updated policy and value function networks trained in parallel with multiple agents. In A2C [38], on the other hand, the networks of multiple agents are updated synchronously. We base our actor-critic architecture on ACKTR [11] which applies trust region optimization to both the actor and the critic. We propose a series of adaption for solving online 3D BPP. To realize constrained reinforcement learning, we propose a simple approach by projecting the trajectories sampled from the actor to the constrained state-action space, instead of resorting to more involved constrained policy optimization [12].

**RL for combinatorial optimization** has a distinguished history [39, 40] and is still an active direction with especially intensive focus on TSP [41]. Early attempts strive for heuristics selection using RL [42]. Bello and colleagues [41] combined *RL pretraining* and *active search* and demonstrated that RL-based optimization outperforms supervised learning framework, when tackling NP-hard combinatorial problems. Recently, Hu et al. [43] proposed a DRL solution to 3D BPP. Different from ours, their work deals with an offline setting where the main goal is to find an optimal sequence of items inspired by the well-known Pointer Network [44].

## 3 Method

In online 3D-BPP, the agent is agnostic on $l_i$, $w_i$ or $h_i$ of all the items in $\mathcal{I}$ – only immediately incoming ones $\mathcal{I}_o \subset \mathcal{I}$ are observable. As soon as an item arrives, we pack it into the bin, and no further adjustment will be applied. As the complexity of BPP decreases drastically for bigger items, we further constrain the sizes of all items to be $l_i \leq L/2$, $w_i \leq W/2$, and $h_i \leq H/2$. We start with our problem statement under the context of DRL and the formulation based on constrained DRL. We show how we solve the problem via predicting action feasibility in the actor-critic framework.

### 3.1 Problem statement and formulation

The 3D-BPP can be formulated as a Markov decision process, which is a tuple of $(\mathcal{S}, \mathcal{A}, P, R)$. $\mathcal{S}$ is the set of environment states; $\mathcal{A}$ is the action set; $R : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ is the reward function; $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to [0, 1]$ is the transition probability function. $P(s'|s, a)$ gives the probability of transiting from $s$ to $s'$ for given action $a$. Our method is model-free since we do not learn $P(s'|s, a)$. A stationary policy $\pi : \mathcal{S} \to \mathcal{A}$ is a map from states to probability distributions over actions, with $\pi(a|s)$ denoting the probability of selecting action $a$ under state $s$. For DRL, we seek for a policy $\pi$ to maximize the accumulated discounted reward, $J(\pi) = E_{\tau \sim \pi}[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)]$. Here, $\gamma \in [0, 1]$ is the discount factor, and $\tau = (s_0, a_0, s_1, \dots)$ is a trajectory sampled based on the policy $\pi$.

The environment state of 3D-BPP is comprised of two parts: the current configuration of the bin and the coming items to be placed. For the first part, we parameterize the bin through discretizing its bottom area as a $L \times W$ regular grid along length ($X$) and width ($Y$) directions, respectively. We record at each grid cell the current height of stacked items, leading to a *height map* $\mathbf{H}_n$ (see Figure 2). Here, the subscript $n$ implies $n$ is the next item to be packed. Since all the dimensions are integers, $\mathbf{H}_n \in \mathbb{Z}^{L \times W}$ can be expressed as a 2D integer array. The dimensionality of item $n$ is given as $\mathbf{d}_n = [l_n, w_n, h_n]^\top \in \mathbb{Z}^3$. Putting together, the current environment state can be written as $s_n = \{\mathbf{H}_n, \mathbf{d}_n, \mathbf{d}_{n+1}, ..., \mathbf{d}_{n+k-1}\}$. We first consider the case where $k = |\mathcal{I}_o| = 1$, and name this special instance as BPP-1. In other words, BPP-1 only considers the immediately coming item $n$ i.e., $\mathcal{I}_o = \{n\}$. We then generalize it to BPP-$k$ with $k > 1$ afterwards.

**BPP-**1    In BPP-1, the agent places $n$'s front-left-bottom (FLB) corner (Figure 2 (left)) at a certain grid point or the loading position (LP) in the bin. For instance, if the agent chooses to put $n$ at the LP of $(x_n, y_n)$. This action is represented as $a_n = L \cdot x_n + y_n \in \mathcal{A}$, where the action set
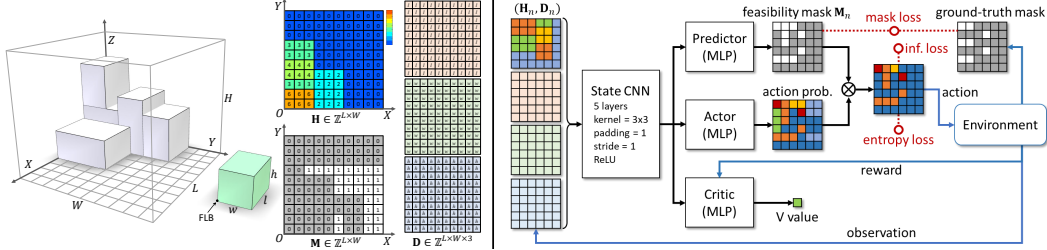
Figure 2: Left: The environment state of the agent includes the configuration of the bin (the grey boxes) and the size of the next item to be packed (green box). The bin configuration is parameterized as a height map $\mathbf{H}$ over a $L \times W$ grid. The feasibility mask $\mathbf{M}$ is a binary matrix of size $L \times W$ indicating the placement feasibility at each grid cell. The three dimensions of the next item are stored into a $L \times W \times 3$ tensor $\mathbf{D}$. Right: The network architecture (the three losses other than the standard actor and critic losses are shown in red color).

$\mathcal{A} = \{0, 1, \ldots, L \cdot W - 1\}$. After $a_n$ is executed, $\mathbf{H}_n$ is updated by adding $h_n$ for all the cells covered by $n$, i.e., $\mathbf{H}'_n(x, y) = \mathbf{H}_n(x, y) + h_n$ for $x \in [x_n, x_n + l_n], y \in [y_n, y_n + w_n]$. The probability distribution of the height map transition is a Dirac delta function: $P(\mathbf{H}|\mathbf{H}_n, a_n) = 1$ for $\mathbf{H} = \mathbf{H}'_n$ and $P(\mathbf{H}|\mathbf{H}_n, a_n) = 0$ otherwise.

During packing, the agent needs to secure enough space in the bin to host $n$. Meanwhile, it is equally important to have $n$ statically equilibrated by the underneath at the LP so that all the stacking items are physically stable. Evaluating the physical stability at a LP is involved, taking into account of $n$'s center of mass, moment of inertia, and rotational stability [45]. All of them are normally unknown as the mass distribution differs among items. To this end, we employ a conservative and simplified criterion. Specifically, a LP is considered *feasible* if it not only provides sufficient room for $n$ but also satisfies any of following conditions with $n$ placed: 1) over $60\%$ of $n$'s bottom area and all of its four bottom corners are supported by existing items; or 2) over $80\%$ of $n$'s bottom area and three out of four bottom corners are supported; or 3) over $95\%$ of $n$'s bottom area is supported. We store the feasibility of all the LPs with a *feasibility mask* $\mathbf{M}_n$, a $L \times W$ binary matrix (also see Figure 2).

Since not all actions are allowed, our problem becomes a constrained Markov decision processes (CMDP) [9]. Typically, one augments the MDP with an auxiliary cost function $C : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ mapping state-action tuples to costs, and require that the expectation of the accumulated cost should be bounded by $c_m$: $J_C(\pi) = E_{\tau \sim \pi}[\sum_{t=0}^{\infty} \gamma_C^t C(s_t, a_t)] \leq c_m$. Several methods have been proposed to solve CMDP based on e.g., algorithmic heuristics [46], primal-dual methods [47], or constrained policy optimization [12]. While these methods are proven effective, it is unclear how they could fit for 3D-BPP instances, where the constraint is rendered as a discrete mask. In this work, we propose to exploit the mask $\mathbf{M}$ to guide the DRL training to enforce the feasibility constraint without inducing excessive training complexity.

### 3.2 Network architecture

We adopt the actor-critic framework with Kronecker-Factored Trust Region (ACKTR) [11]. The actor-critic approach iteratively updates an actor and a critic module in an interleaving manner. In each iteration, the actor learns a policy network that outputs the probability of each action (i.e., placing $n$ at the each LP). The critic trains a state-value network producing the value function.

**State input** In the original ACKTR framework, both actor and critic networks take the raw state directly as input. In our implementation however, we devise a CNN or *the state CNN* to encode the raw state vector into features. To facilitate this, we "stretch" $\mathbf{d}_n$ into a rank-3 tensor $\mathbf{D}_n \in \mathbb{Z}^{L \times W \times 3}$ so that each dimension of $\mathbf{d}_n$ spans a $L \times W$ matrix with all of its elements being $l_n$, $w_n$ or $h_n$ (also see Figure 2). Consequently, state $s_n = (\mathbf{H}_n, \mathbf{D}_n)$ becomes a $L \times W \times 4$ array (Figure 2 (right)).

**Reward** Our reward is defined as the final capacity utilization $r = \sum 10 * l_i \cdot w_i \cdot h_i / L \cdot W \cdot H \in (0, 10]$ at the termination – no step-wise reward used. We find that this simple rewarding mechanism works quite well with our network design, saving the effort of reward engineering.
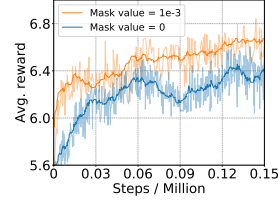
**Feasibility constraints** We design three mechanisms to enforce feasibility constraints. First, we introduce an independent MLP module, namely the *mask predictor* (Figure 2 (right)), to predict the feasibility mask $\mathbf{M}_n$ for the item $n$. The predictor takes the state CNN features of the current state as the input and is trained with the ground-truth mask as the supervision. Next, we use the

4

predicted mask to modulate the output, i.e., the probability distribution of the actions. In theory, if the LP at $(x, y)$ is infeasible for $n$, the corresponding probability $P(a_n = L \cdot x + y|s_n)$ should be set to 0. However, we find that setting $P$ to a small positive quantity like $\epsilon = 1e - 3$ works better in practice – it provides a strong penalty to an invalid action but a smoother transformation beneficial to the network training. The inset shows that softening the mask-based modulation improves the training convergence. To further discourage infeasible actions, we explicitly minimize the summed probability at all infeasible LPs: $E_{inf} = \sum P(a_n = L \cdot x + y|s_n), \forall (x, y)|\mathbf{M}_n(x, y) = \epsilon$, which is plugged into the final loss function for training.

**Loss function**    Our loss function is defined as:

$$L = \alpha \cdot L_{actor} + \beta \cdot L_{critic} + \lambda \cdot L_{mask} + \omega \cdot E_{inf} - \psi \cdot E_{entroy}. \tag{1}$$

Here, $L_{actor}$ and $L_{critic}$ are the loss functions used for training the actor and the critic, respectively. $L_{mask}$ is the MSE loss for mask prediction. To push the agent to explore more LPs, we also utilize an action entropy loss $E_{entroy} = \sum_{\mathbf{M}_n(x,y)=1} -P(a_n|s_n) \cdot \log(P(a_n|s_n))$. Note that the entropy is computed only over the set of all feasible actions whose LP satisfies $\mathbf{M}_n(x, y) = 1$. In this way, we stipulate the agent to explore only feasible actions. We use the following weights throughout our experiments: $\alpha = 1$, $\beta = \lambda = 0.5$, and $\omega = \psi = 0.01$. Figure on the right shows how softening the mask-based modulation helps improving the training convergence.
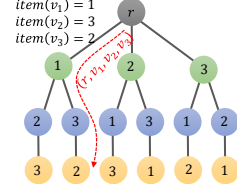
### 3.3    BPP-$k$ with $k = |\mathcal{I}_o| > 1$

In a more general case, the agent receives the information of $k > 1$ *lookahead* items (i.e., from $n$ to $n + k - 1$). Obviously, the additional items inject more information to the environment state, which should be exploited in learning the policy $\pi(a_n|\mathbf{H}_n, \mathbf{d}_n, ..., \mathbf{d}_{n+k-1})$. One possible solution is to employ sequential modeling of the state sequence $(\mathbf{d}_n, ..., \mathbf{d}_{n+k-1})$ using recurrent neural network. We found that however, such state encoding cannot well inform the agent about the lookahead items during DRL training and yields limited improvement. Alternatively, we propose a simple solution leveraging the height map $\mathbf{H}$ update and feasibility mask prediction.

The core idea is to condition the placement of the current item $n$ on that of the next $k - 1$ ones. Note that the actual placement of the $k$ items still follows the order of arrival. To make the current placement account for the future ones, we opt to "hallucinate" the placement of future items through updating the height map accordingly. Conditioned on the virtually placed future items, the decision for the current item could be globally more optimal. However, such virtual placement must satisfy the *order dependence constraint* which stipulates the earlier items should never be packed on top of the later ones. In particular, given two items $p$ and $q$, $p < q$ in $\mathcal{I}_o$, if $q$ is (virtually) placed before $p$, we require that the placement of $p$ should be spatially independent to the placement of $q$. It means $p$ can never be packed at any LPs that overlap with $q$. This constraint is enforced by setting the height values in $\mathbf{H}$ at the corresponding LPs to $H$, the maximum height value allowed: $\mathbf{H}_p(x, y) \leftarrow H$, for all $x \in [x_q, x_q + l_q]$ and $y \in [y_q, y_q + w_q]$. Combining explicit height map updating with feasibility mask prediction, the agent learns a policy with the order dependence constraint satisfied implicitly.

**Monte Carlo permutation tree search**    Having trained a policy that admits order-dependence-constrained virtual placement, we are able to search for a better $a_n$ through exploring the permutations of the sequence $(\mathbf{d}_n, ..., \mathbf{d}_{n+k-1})$. This amounts to a permutation tree search during which only the actor network test is conducted – no training is needed. The inset shows a $k$-level permutation tree in which a path from the root to a leaf forms a possible permutation of the placement of the $k$ items in $\mathcal{I}_o$. The search is conducted in a top-down manner. Suppose $v_i \neq r$ is the node being visited via the path $(r, v_1, v_2, ..., v_i)$, where $r$ is the root node, and the subscript $i$ implies the node $v_i$ is at the $i$-th level. Let $item(v_i)$ be the item associated with $v_i$. For any node $v_j$, $j < i$ along the path (i.e., $v_j$ has been visited before $v_i$ along the path), if $item(v_j) < item(v_i)$ meaning $item(v_j)$ arrives before $item(v_i)$ under this permutation, $item(v_j)$ has already been packed into the bin. Otherwise if $item(v_j) > item(v_i)$, we must additional block the position of $item(v_j)$ in $\mathbf{M}_{item(v_i)}$ when packing $item(v_i)$ to avoid place $item(v_i)$ upon $item(v_j)$. After search, we choose the action $a_n$ corresponding to the permutation with the highest accumulative critic value.

Clearly, enumerating all the permutations for $k$ items quickly becomes prohibitive with an $O(k!)$ complexity. To make the search scalable, we employ the Monte Carlo tree search (MCTS) [48].

5

With MCTS, the permutation tree is expanded in a priority-based fashion through evaluating how promising a node would lead to the optimal solution. Please refer to the supplemental material for the details on our adaption of the standard MCTS. MCTS allows a scalable lookahead for BPP-$k$ with a complexity of $O(km)$ where $m$ is the number of paths sampled from the permutation tree.

## 4 Experiments

We implement our framework on a desktop computer (`ubuntu 16.04`), which equips with an `Intel Xeon Gold 5115` CPU @ 2.40 GHz, 64G memory, and a `Nvidia Titan V` GPU with 12G memory. The DRL and all other networks are implemented with `PyTorch` [49]. The model training takes about 16 hours on our machine. The testing time of BPP-1 model (no lookahead) is less than 100 ms. Please refer to the supplemental material for more implementation details.

Table 1: In this ablation study, we quantitatively compare the space utilization and the total number of packed items with different combinations of MP, MC and FE.

| MP | MC | FE | Space uti. | # items |
|----|----|----|------------|---------|
| ✗ | ✗ | ✗ | 7.82% | 2.0 |
| ✓ | ✗ | ✓ | 27.9% | 7.5 |
| ✓ | ✓ | ✗ | 63.7% | 16.9 |
| ✗ | ✓ | ✓ | 63.0% | 16.7 |
| ✓ | ✓ | ✓ | **66.9%** | **17.5** |

Table 2: HM yields better benchmarks over 1D based parameterizations of HV and ISV.

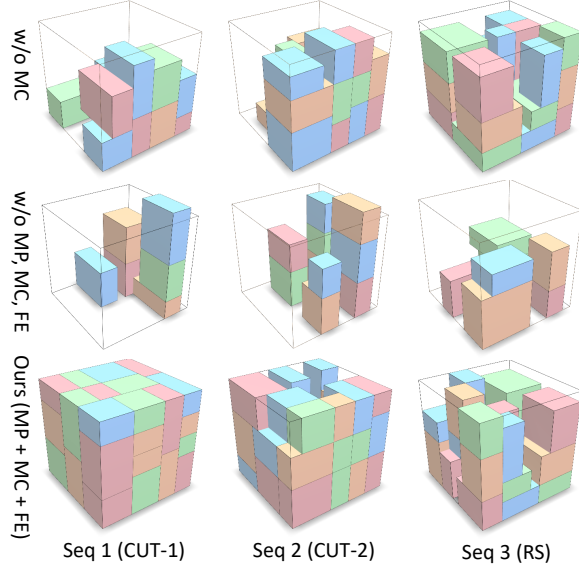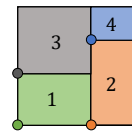| | Space uti. | # items |
|----|------------|---------|
| ISV | 54.3% | 14.1 |
| HV | 57.4% | 14.8 |
| HM (ours) | **73.4%** | **19.1** |



Figure 3: Packing results in the ablation study.

**Training and testing set** We set $L = W = H = 10$ in our experiments with 64 pre-defined item dimensions ($|\mathcal{I}| = 64$). As mentioned, we have $l_i \leq L/2$, $w_i \leq W/2$ and $h_i \leq H/2$ to avoid over-simplified scenarios. A test sequence is synthesized by generating items out of $\mathcal{I}$, and the total volume of items should be equal to or bigger than bin's volume. We first create item sequences using random sampling (RS), where each item along a sequence is picked out of $\mathcal{I}$ randomly. A disadvantage of the random sampling lies in the fact that the optimality of a sequence is unknown: we are unaware of the optimal packing configuration of a randomly generated sequence (unless performing a brute-force search). While DRL is a type of self play technique, exposing the learning network to successful sequences could significantly improve the quality of the policy.

In addition to RS, we also generate item sequences via *cutting stock* [50]. Specifically, items in a sequence are created by "cutting" the bin so that we understand the sequence may be perfectly packed and restored back to the bin. There are two variations of this strategy. 1) After the cutting, we sort resulting items into the sequence based on $Z$ coordinates of their FLBs, from bottom to top. If FLBs of two items have the same $Z$ coordinate, their order in the sequence is randomly determined. 2) The cut items are sorted based on their stacking dependency: an item can be added to the sequence only after all of its supporting items are there. We refer those two cutting-based strategies as CUT-1 and CUT-2, respectively. A 2D toy example is given in the right figure with FLB of each item highlighted. Under CUT-1, both $\{1, 2, 3, 4\}$ and $\{2, 1, 3, 4\}$ are valid item sequences. If we use CUT-2 on the other hand, $\{1, 3, 2, 4\}$ and $\{2, 4, 1, 3\}$ would also be valid sequences as the placement of 3 or 4 depends on 1 or 2. In our experiments, we have generated 2,000 sequences using RS, CUT-1, and CUT-2 respectively for DRL testing. The performance of the packing algorithm is quantitated with space utilization (*space uti.*) and the total number of items packed in the bin (*# items*).

**Ablation study** First, we report an ablation study. We found that the packing performance drops significantly if we do not incorporate the mask prediction (MP) during the training. The performance

is impaired if the mask constraint (MC) is not enforced, which rules this option out. The feasibility based entropy (FE) is also beneficial for both the training and final performance. The benchmark and a representative final packing can be found in Figure 3. We refer the reader to the supplementary material for more illustrative and animated results.

**Height parameterization**   Next, we show that the environment parameterization using the proposed 2D height map (HM) (i.e., $\mathbf{H}$ matrix) is necessary and effective. To this end, we compare DRL benchmark using HM with other two straightforward 1D alternatives. The first competitor is the height vector (HV), which is a $L \cdot W$-dimensional vector stacking columns in $\mathbf{H}$. The second competitor is referred to as the item sequence vector (ISV). The ISV lists all the information of items currently packed in the bin. Each packed item has 6 parameters corresponding to $X$, $Y$, and $Z$ coordinates of its FLB as well as item's dimension. The comparison is reported in left table of Figure 4, from which we can see that HM parameterization helps the DRL produce better packing actions in terms of both space utilization and the number of packed items. And we also find 2D height map (HM) is a more effective way to describe the state-action space, right 3 plots of Figure 4 compare the average reward received using different parameterizations.
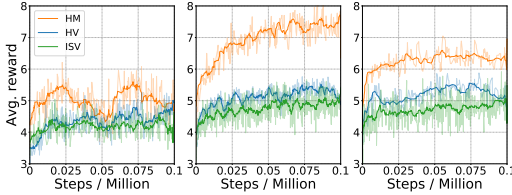


Figure 4: Comparison among different space representation. HM shows a clear advantage over vector-based parameterizations.
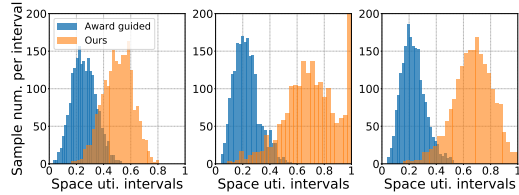
Figure 5: Comparison between our method with reward tuning. Our method obtains much better space utilization.

**Constraint vs. reward**   Normally in DRL training, one discourages low-profile moves by tuning the reward function (similar to penalty method). We show that this strategy is less effective than our constraint-based method (i.e., learn invalid move by predicting the mask). The result is reported in Figure 5, where we have the reward function yield a high penalty term to penalize unsafe placements. Constraint-based DRL seldom predicts invalid moves (99.5% placement are legit).

**Scalability of BPP-$k$**   With the capability of lookahead, it is expected that the agent better exploits the remaining space in the bin and delivers a more compact packing. On the other hand, due to the NP-hard nature, big $k$ values increase the environment space exponentially. Therefore, it is important to understand if MCTS is able to effectively navigate us in the space at the scale of $\mathbf{O}(k!)$ for a good packing strategy. The result is plotted in right of Figure 6. In this experiment, we compare our method with a brute-force permutation search, which traverses all $k!$ permutations of $k$ coming items and chooses the best packing strategy (i.e., the global optimal). We find that MCTS yields high-quality results – it does have a slightly lower space utilization rate
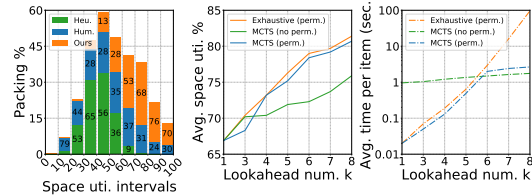


Figure 6: Left: The distribution of space utilization using boundary rule (Heu.), human intelligence (Hum.), and our BPP-1 method (Ours). Middle: The performance of our MCTS based BPP-$k$ model achieves similar performance of the brute-force search (i.e., with permutation tree). Right: Our permutation based MCTS maintain high time efficiency while keeping high performance.

($\sim 3\%$), however MCTS is far more efficient. The search time of brute-force permutation quickly surpasses 100 seconds when $k = 8$. MCTS only takes 3.6 seconds even for $k = 20$, when permutation needs hours. A slightly larger $k$ makes the brute-force search intractable on desktop computers.

**Generalization**   Our algorithm is designed specifically catering real-world applications. However, it can be generalized to handle different 3D-BPP variants such as multi-bin BPP or allowing the item to be (axis-align) oriented. For 3D-BPP with more than one bins, we can initialize multiple BPP-1 instances matching the total bin number. When an item arrives, we pack it to the bin which has the highest packing score based on BPP-1 network returns(details see supplementary material). Table 3 shows the our results on different number of bins. It is easy to understand that more bins provide more options to house an upcoming item and have a better packing benchmark than the single bin BPP. In this generalization, both space and time complexities grow linearly with the number of bins.

7

When the orientation of an item is also considered, we need to create multiple variations of feasibility mask to incorporate more freedoms induced by different item orientations. The orientation with highest score be selected. In our experiment, we find allowing orientation (i.e., BPP-RS) achieves 62.1% space utilization with average 15.1 items packed while BPP-1 only has 50.5% space utilization with average 12.2 items

Table 3: Multi-bin BPP.

| # bins | Space uti. | # items |
|--------|-----------|---------|
| 1 | 67.4% | 17.6 |
| 4 | 69.4% | 18.8 |
| 9 | 72.1% | 19.1 |
| 16 | 75.3% | 19.6 |
| 25 | 77.8% | 20.2 |

packed. This experiment demonstrates our network is capable to handle this generalization and BPP-RS yields better results.

**Our method vs. non-learning methods** To better illustrate the advantage of our method, we also design a heuristic non-learning algorithm for BPP-$k$. This heuristic replicates human's behavior by trying to place a new item side-by-side with the existing packed items and keep the packing volume as regular as possible. We refer to this heuristic strategy *boundary rule*. Comparative benchmarks are given in Table 4. Here, we also compare our method with existing non-learning method proposed in [8]. We understand that an exact 3D-BPP algorithm is also presented in [8] using branch and bounding. However, solving BPP-$k$ exactly is prohibitive for an NP problem. Therefore, we only compare our method with the approximation algorithm in [8], which is based on the *layer building principle* (LBP) [51]. Interestingly, we find that our method automatically learns the boundary rule, even without imposing such constraints explicitly. This is because leaving gaps between items decreases the space utilization of the bin. However, blindly sticking to the boundary rule could be problematic as other items (in BPP-$k$) may yield a better fit. Putting in short, boundary rule is an immediate human intuition. Yet, it does not produce good packing results consistently. This fact is further elaborated in our next experiment.

Table 4: Our method has a better performance than state-of-the-art heuristic methods. In general, we see a better benchmark when the training and testing sets match each other. Note that LBP [8] needs whole sequence information for global optimization while our method can do it online.

| Method | Space uti. | | | # items | | |
|--------|-----|-------|-------|-----|-------|-------|
| | RS | CUT-1 | CUT-2 | RS | CUT-1 | CUT-2 |
| Our method (trained on RS) | 50.5% | 60.8% | 60.9% | 12.2 | 15.7 | 16.2 |
| Our method (train on CUT-1) | 43.4% | **73.4%** | 62.4% | 10.7 | **19.1** | 16.6 |
| Our method (train on CUT-2) | 47.6% | 69.4% | **66.9%** | 11.6 | 17.9 | **17.5** |
| Boundary rule | 34.9% | 41.2% | 40.8% | 8.7 | 10.8 | 11.1 |
| LBP [8] | **54.7%** | 59.1% | 59.5% | **12.9** | 14.9 | 15.2 |

**Our method vs. human intelligence** The strongest competitor to all heuristic algorithms may be the real human. To this end, we created a simple Sokoban-like app (see the supplemental material) and asked 30 human users to pack items manually vs. AI (our method). The winner is the one with a higher space utilization rate. All the users are CS-majored undergraduate/graduate students. We do not impose any time limits to the user. The statistics are plotted in left of Figure 6. To our surprise, our method outperforms human players in general ($1,339$ AI wins vs. $406$ human wins and $98$ evens): it achieves 68.9% average space utilization over $1,851$ games, while human players only have 52.1%.

## 5    Conclusion

We have tackled a challenging online 3D-BPP via formulating it as a constrained Markov decision process and solving it with constrained deep RL. Such constraints encompass both collision avoidance and physical stability. Within the actor-critic framework, we achieve policy optimization subject to the complicated constraints based on a height map bin representation and action feasibility prediction. The predicted feasibility is used to modulate the actor output. In realizing BPP with multiple lookahead items, we again leverage height map input and feasibility prediction to impose the constraint of order dependence. This way, the best action can be searched over different permutations of the lookahead items. In future, we would like to investigate more relaxations of the problem. For example, one could lift the order dependence constraint by adding a buffer zone smaller than $|\mathcal{I}_o|$. Another relaxation is to allow the agent to reorient the items (with six box orientations) for a better space utilization. All these variants are interesting and practically useful problems for DRL research.

## Broader Impact

We believe the broader impact of this research is significant. BPP is considered as the one of the most needed academic problems due to its wide applicability in our daily lives [3]. Providing a good

8 of the most needed academic problems due to its wide applicability in our daily lives [3]. Providing a good

8

solution to 3D-BPP yields direct and profound impacts beyond academia. A practical application scenario of our algorithm is illustrated in Figure. 1. The conveyor belt sends a sequence of items to a robot agent. Those items are scanned by an overhead depth sensor, so that the information of their dimensions can be obtained using some basic computer vision technique. The agent places the item as soon as it arrives. We note that such setup has been already widely deployed in logistics hubs and manufacture plants. A good BPP solution benefits all the downstream operations like package wrapping, transportation, warehousing and distribution.

Aside from its real-world applications, 3D-BPP algorithm proposed in this paper is generalizable and transformative. As a member of NP-Hard family, an algorithm to 3D-BPP can be reduced to a solution to many other difficult NP-Hard problems including the famous knapsack problem, dominating set problem, nurse scheduling problem, or cut socking problem. Our constrained DRL framework and its novel solution has a good potential for solving a wide category of combinatorial problems.

# References

[1] Bernhard Korte and Jens Vygen. Bin-packing. In *Kombinatorische Optimierung*, pages 499–516. Springer, 2012.

[2] EG Co man Jr, MR Garey, and DS Johnson. Approximation algorithms for bin packing: A survey. *Approximation algorithms for NP-hard problems*, pages 46–93, 1996.

[3] Steven S Skiena. The stony brook algorithm repository. `http://www.cs.sunysb.edu/algorith/implement/nauty/implement.shtml`, 1997.

[4] Maxence Delorme, Manuel Iori, and Silvano Martello. Bin packing and cutting stock problems: Mathematical models and exact algorithms. *European Journal of Operational Research*, 255(1):1–20, 2016.

[5] Alexander Schrijver. *Theory of linear and integer programming*. John Wiley & Sons, 1998.

[6] Teodor Gabriel Crainic, Guido Perboli, and Roberto Tadei. Extreme point-based heuristics for three-dimensional bin packing. *Informs Journal on computing*, 20(3):368–384, 2008.

[7] Korhan Karabulut and Mustafa Murat İnceoğlu. A hybrid genetic algorithm for packing in 3d with deepest bottom left with fill method. In *International Conference on Advances in Information Systems*, pages 441–450. Springer, 2004.

[8] Silvano Martello, David Pisinger, and Daniele Vigo. The three-dimensional bin packing problem. *Operations research*, 48(2):256–267, 2000.

[9] Eitan Altman. *Constrained Markov decision processes*, volume 7. CRC Press, 1999.

[10] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937, 2016.

[11] Yuhuai Wu, Elman Mansimov, Roger B Grosse, Shun Liao, and Jimmy Ba. Scalable trust-region method for deep reinforcement learning using kronecker-factored approximation. In *Advances in neural information processing systems*, pages 5279–5288, 2017.

[12] Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization. In *Proceedings of the 34th International Conference on Machine Learning*, pages 22–31. JMLR. org, 2017.

[13] Leonid V Kantorovich. Mathematical methods of organizing and planning production. *Management science*, 6(4):366–422, 1960.

[14] Paul E Sweeney and Elizabeth Ridenour Paternoster. Cutting and packing problems: a categorized, application-orientated research bibliography. *Journal of the Operational Research Society*, 43(7):691–706, 1992.

[15] Edward G Coffman, Michael R Garey, and David S Johnson. Approximation algorithms for bin-packing—an updated survey. In *Algorithm design for computer system design*, pages 49–106. Springer, 1984.

[16] W Fernandez De La Vega and George S Lueker. Bin packing can be solved within $1+\varepsilon$ in linear time. *Combinatorica*, 1(4):349–355, 1981.

[17] David S. Johnson, Alan Demers, Jeffrey D. Ullman, Michael R Garey, and Ronald L. Graham. Worst-case performance bounds for simple one-dimensional packing algorithms. *SIAM Journal on computing*, 3(4):299–325, 1974.

[18] Silvano Martello. Knapsack problems: algorithms and computer implementations. *Wiley-Interscience series in discrete mathematics and optimiza tion*, 1990.

[19] Narendra Karmarkar and Richard M Karp. An efficient approximation scheme for the one-dimensional bin-packing problem. In *23rd Annual Symposium on Foundations of Computer Science (sfcs 1982)*, pages 312–320. IEEE, 1982.

[20] Silvano Martello and Paolo Toth. Lower bounds and reduction procedures for the bin packing problem. *Discrete applied mathematics*, 28(1):59–70, 1990.

[21] Armin Scholl, Robert Klein, and Christian Jürgens. Bison: A fast hybrid procedure for exactly solving the one-dimensional bin packing problem. *Computers & Operations Research*, 24(7):627–645, 1997.

[22] Martine Labbé, Gilbert Laporte, and Silvano Martello. An exact algorithm for the dual bin packing problem. *Operations Research Letters*, 17(1):9–18, 1995.

[23] Maxence Delorme, Manuel Iori, and Silvano Martello. Bin packing and cutting stock problems: mathematical models and exact algorithms. In *Decision models for smarter cities*, 2014.

[24] Maxence Delorme, Manuel Iori, and Silvano Martello. Bpplib: a library for bin packing and cutting stock problems. *Optimization Letters*, 12(2):235–250, 2018.

[25] Silvano Martello and Daniele Vigo. Exact solution of the two-dimensional finite bin packing problem. *Management science*, 44(3):388–399, 1998.

[26] Mhand Hifi, Imed Kacem, Stéphane Nègre, and Lei Wu. A linear programming approach for the three-dimensional bin-packing problem. *Electronic Notes in Discrete Mathematics*, 36:993–1000, 2010.

[27] Oluf Faroe, David Pisinger, and Martin Zachariasen. Guided local search for the three-dimensional bin-packing problem. *Informs journal on computing*, 15(3):267–283, 2003.

[28] JL de Castro Silva, NY Soma, and N Maculan. A greedy search for the three-dimensional bin packing problem: the packing static stability case. *International Transactions in Operational Research*, 10(2):141–153, 2003.

[29] Andrea Lodi, Silvano Martello, and Daniele Vigo. Approximation algorithms for the oriented two-dimensional bin packing problem. *European Journal of Operational Research*, 112(1):158–166, 1999.

[30] Teodor Gabriel Crainic, Guido Perboli, and Roberto Tadei. Ts2pack: A two-level tabu search for the three-dimensional bin packing problem. *European Journal of Operational Research*, 195(3):744–760, 2009.

[31] Xueping Li, Zhaoxia Zhao, and Kaike Zhang. A genetic algorithm for the three-dimensional bin packing problem with heterogeneous bins. In *IIE Annual Conference. Proceedings*, page 2039. Institute of Industrial and Systems Engineers (IISE), 2014.

[32] Shigeyuki Takahara and Sadaaki Miyamoto. An evolutionary approach for the multiple container loading problem. In *Fifth International Conference on Hybrid Intelligent Systems (HIS'05)*, pages 6–pp. IEEE, 2005.

[33] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.

[34] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.

[35] Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Van Hasselt, Marc Lanctot, and Nando De Freitas. Dueling network architectures for deep reinforcement learning. *arXiv preprint arXiv:1511.06581*, 2015.

[36] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. Deterministic policy gradient algorithms. 2014.

[37] Gabriel Barth-Maron, Matthew W Hoffman, David Budden, Will Dabney, Dan Horgan, Alistair Muldal, Nicolas Heess, and Timothy Lillicrap. Distributed distributional deterministic policy gradients. *arXiv preprint arXiv:1804.08617*, 2018.

[38] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

[39] Luca M Gambardella and Marco Dorigo. Ant-q: A reinforcement learning approach to the traveling salesman problem. In *Machine Learning Proceedings 1995*, pages 252–260. Elsevier, 1995.

[40] Wei Zhang and Thomas G Dietterich. Solving combinatorial optimization tasks by reinforcement learning: A general methodology applied to resource-constrained scheduling. *Journal of Artificial Intelligence Reseach*, 1:1–38, 2000.

[41] Irwan Bello, Hieu Pham, Quoc V Le, Mohammad Norouzi, and Samy Bengio. Neural combinatorial optimization with reinforcement learning. *arXiv preprint arXiv:1611.09940*, 2016.

[42] Alexander Nareyek. Choosing search heuristics by non-stationary reinforcement learning. In *Metaheuristics: Computer decision-making*, pages 523–544. Springer, 2003.

[43] Haoyuan Hu, Xiaodong Zhang, Xiaowei Yan, Longfei Wang, and Yinghui Xu. Solving a new 3d bin packing problem with deep reinforcement learning method. *arXiv preprint arXiv:1708.05930*, 2017.

[44] Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. Pointer networks. In *Advances in Neural Information Processing Systems*, pages 2692–2700, 2015.

[45] Herbert Goldstein, Charles Poole, and John Safko. Classical mechanics, 2002.

[46] Eiji Uchibe and Kenji Doya. Constrained reinforcement learning from intrinsic and extrinsic rewards. In *2007 IEEE 6th International Conference on Development and Learning*, pages 163–168. IEEE, 2007.

[47] Yinlam Chow, Mohammad Ghavamzadeh, Lucas Janson, and Marco Pavone. Risk-constrained reinforcement learning with percentile risk criteria. *The Journal of Machine Learning Research*, 18(1):6070–6120, 2017.

[48] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *Nature*, 550(7676):354–359, 2017.

[49] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, pages 8024–8035, 2019.

[50] Paul C Gilmore and Ralph E Gomory. A linear programming approach to the cutting-stock problem. *Operations research*, 9(6):849–859, 1961.

[51] Fan RK Chung, Michael R Garey, and David S Johnson. On packing two-dimensional bins. *SIAM Journal on Algebraic Discrete Methods*, 3(1):66–76, 1982.

[52] Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In *European conference on machine learning*, pages 282–293. Springer, 2006.

# 6 Supplemental Material

## 6.1 Overview

In this supplemental document, we report more implementation and experiment details. Specifically, § 6.2 gives more descriptions regarding the network architecture, training selection, Monte Carlo permutation tree search, etc. § 6.3 elaborates how CUT-1, CUT-2, and RS are constructed. § 6.4 explains the heuristic baseline (i.e., the so-called boundary rule) we compared in our experiment. More experiment results are reported in § 6.5. The user study design is discussed in § 6.6.

## 6.2 Implementation Details

We report the details of our implementation in this section, and our source code is also submitted with this supplemental material.

**Network architecture and training configurations** A detailed specifications of the our network is shown in Figure 13. Our pipeline consists of three major components: an actor network, a critic network, and the feasibility mask predictor. It takes three inputs, namely *height map* $\mathbf{H}_n$, the dimensionality $\mathbf{d}_n = [l_n, w_n, h_n]^\top \in \mathbb{Z}^3$ of the next item $n$ to be packed, and the feasibility mask $\mathbf{M}_n$. Note that $\mathbf{M}_n$ is only used in the training processing, our mask predictor would generate the mask input during the testing.

The whole network is trained via a composite loss consisting of actor loss $L_{actor}$, critic loss $L_{critic}$, mask prediction loss $L_{mask}$, infeasibility loss $E_{inf}$ and action entropy loss $E_{entropy}$. These loss function are defined as:

$$
\begin{cases}
L_{actor} & = (\mathbf{R}_n - V(\mathbf{S}_n)) \log P(a_n | s_n) \\
L_{critic} & = (\mathbf{R}_n - V(\mathbf{S}_n))^2 \\
L_{mask} & = \sum_{(x,y)} (\mathbf{M}_n^{gt} - \mathbf{M}_n^{pred})^2 \\
E_{inf} & = \sum P(a_n = L \cdot x + y | s_n) \\
E_{entropy} & = \sum_{\mathbf{M}_n(x,y)=1} -P(a_n | s_n) \cdot \log\left(P(a_n | s_n)\right),
\end{cases}
\tag{2}
$$

where $\mathbf{R}_n = \sum_{i=0}^{k-1} \gamma^i r_{n+i} + \gamma^k V(\mathbf{S}_{n+k})$; $\gamma \in [0, 1]$ is the discount factor; and $r = \sum 10 \cdot l_i \cdot w_i \cdot h_i / L \cdot W \cdot H \in (0, 10]$ is our reward function. The output of critic network would help the training of actor network which outputs a possibility matrix of the next move. This probability is scaled based on $\mathbf{M}_n$ — if the move is infeasible, the possibility will be multiplied by a penalty factor of $0.001$. Afterwards, a softmax operation is adopted to output the final action distribution. Note that, the infeasibility penalty could be absent in test and our method can still work well in this situation.

**Monte Carlo permutation tree search** Our MCTS algorithm is inspired by the Upper Confidence Bounds for Trees (UCT) [52]. The main difference lies in: 1) The goal of our MCTS is finding the best packing order for the next $k$ items; 2) Max reward is used in MCTS instead of the mean reward as the evaluation strategy; 3) we ensure that every possible packing order being evaluated only once by setting a threshold during the search to limit the visit of MCTS tree nodes. Algorithm 1 outlines the entire procedure step by step, wherein $T$ is the maximum simulation time, $I$ is lookahead items, $i_0$ is current item, $s$ is height map set for each node and $a$ is action set for each node. Environment simulator $f$, which takes height map, the next item dimension, and action (position of item) as the input, and returns the updated height map. Action choosing function $p$ uses policy network from BPP-$k$ model to get the action with highest possibility. We test our method on our three benchmarks, and more results can be found in Figure 7.

**Multi-bin scoring** It is straightforward to generalize our method for multiple bins. The agent needs to determine which bin the item is to be packed. The rest part is naturally reduced to a single-bin BPP instance. To this end, we score all the bins and choose the one with highest score. This score

estimates the reward gain if the current item is packed into the bin. Our bin selection method is described in Algorithm 2. Here, $P$ is the last value estimation of bins, $value$ is the value estimation function via a value network, $i$ is the current item, $B$ is the set of bins and $H$ is the height maps of bins. The default score for each bin at beginning $s_{def} = -0.2$;

---

**Algorithm 1** Permutation MCTS

1: **function** SEARCH($s_0$)
2:     create root node $v_0$ with state $s_0$
3:     **while** $t < T$ **do**
4:         copy $I$ as R
5:         sort R according to original order
6:         $v_l, R \leftarrow$ TREEPOLICY($v_0, R$)
7:         $\Delta \leftarrow$ DEFAULTPOLICY($s(v_l), R$)
8:         BACKUP($v_l, \Delta$)
9:         $t \leftarrow t + 1$
10:     **while** $item(v)$ is not $i_0$ **do**
11:         $v \leftarrow$ BESTCHILD($v, 0$)
12:     **return** $a(v)$

13: **function** TREEPOLICY($v, R$)
14:     **while** $s(v)$ is non-terminal **do**
15:         **if** v not fully expanded then **then**
16:             **return** EXPAND($v, R$)
17:         **else**
18:             $v \leftarrow$ BESTCHILD($v, c$)
19:             $R \leftarrow R \backslash item(v)$
20:     **return** $v, R$

21: **function** DEFAULTPOLICY($s, R$)
22:     **while** $s$ is non-terminal **do**
23:         $i \leftarrow$ first item in R
24:         $a \leftarrow p(s, i)$
25:         $s \leftarrow f(s, i, a)$
26:         $R \leftarrow R \backslash i$
27:     **return** reward for state $s$

28: **function** EXPAND($v, R$)
29:     choose $i \in$ unused item from $R$
30:     add new child $v'$ to $v$
31:     with $a(v') = p(s(v), i)$
32:     with $s(v') = f(s(v), i, a(v'))$
33:     and $item(v') = i$
34:     $R \leftarrow R \backslash i$
35:     **return** $v', R$

36: **function** BESTCHILD($v, c$)
37:     **return** $\underset{v' \in children(v)}{\text{argmin}} (Q(v') + c\sqrt{\frac{N(v)}{1+N(v')}})$

38: **function** BACKUP($v, \Delta$)
39:     **while** $v$ is not null do **do**
40:         $N(v) \leftarrow N(v) + 1$
41:         $Q(v) \leftarrow max(Q(v), \Delta)$
42:         $v \leftarrow$ parent of $v$

---

Table 5: Performance comparison with and without orientation on different benchmarks.

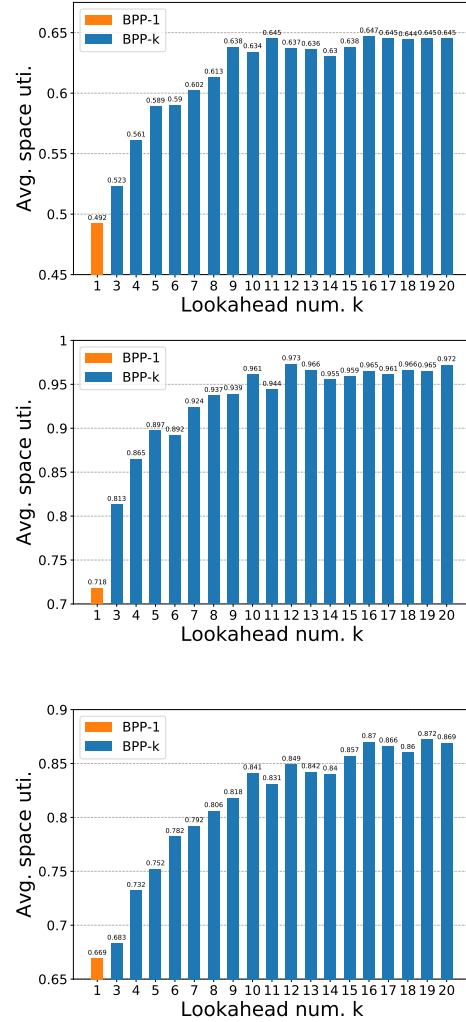|                 | RS    | CUT-1 | CUT-2 |
|-----------------|-------|-------|-------|
| w orientation   | 62.1% | 76.2% | 70.2% |
| w/o orientation | 50.5% | 73.4% | 66.9% |



Figure 7: When $k$ increases, the space utilization rate first goes up and later, enters a "plateau" zone.

**Algorithm 2** Bin Selecting Algorithm

---

**Ensure:** The most fitting bin, $b_{best}$
1: Initialize the set of bin scores $S$ with $s_{def}$
2: **for all** $b \in B$ **do**
3:     **if** $P(b) \neq Null$ **then**
4:         $S(b) \leftarrow value(H(b), i) - P(b)$
5: $b_{best} \leftarrow \arg\max_{b \in B} S(b)$
6: $P(b_{best}) \leftarrow value(H(b_{best}), i)$
7: **return** $b_{best}$

---

**Orientated items** Incorporating items with different orientation is also possible. We multifold the action space and related mask based on how many different orientations are considered. Doing so induces more flexibility for the packing and potentially leads to a better result. This is observed and reported in Table 5.

## 6.3 Benchmark Construction

All 64 pre-defined items are visualized in Figure 8. Algorithm 3 outlines how the dataset is constructed given the bin size $L, W, H$ and a valid item size threshold.
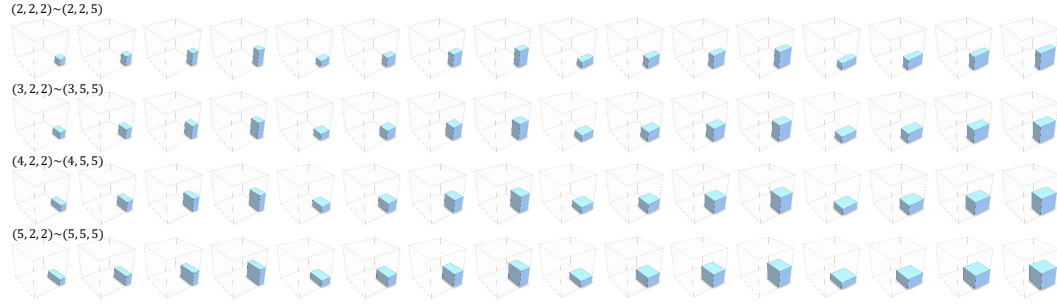


$(2, 2, 2) \sim (2, 2, 5)$

$(3, 2, 2) \sim (3, 5, 5)$

$(4, 2, 2) \sim (4, 5, 5)$

$(5, 2, 2) \sim (5, 5, 5)$

Figure 8: Pre-defined item set $\mathcal{I}$.

---

**Algorithm 3** Benchmark Construction

---

**Require:** valid item size threshold $(l_{min}, w_{min}, h_{min}) \sim (l_{max}, w_{max}, h_{max})$, bin size $(L, W, H)$
1: **function** CONSTRUCTION OF PRE-DEFINED ITEMS COLLECTION
2:     initialize invalid item list $\mathcal{L}_{invalid} = \{(L, W, H)\}$, valid item list $\mathcal{L}_{valid} = \emptyset$
3:     **while** $\mathcal{L}_{invalid} \neq \emptyset$ **do**
4:         randomly pop an $item_i$ from $\mathcal{L}_{invalid}$
5:         randomly select an axis $a_i$ of $item_i$ which $a_{min} \leq a_i \leq a_{max}, a_i \in \{x_i, y_i, z_i\}$
6:         randomly split the item into two sub items along axis $a_i$
7:         **if** $a_{min} \leq a_{sub} \leq a_{max}$ **then**
8:             add the sub item into $\mathcal{L}_{valid}$
9:         **else**
10:            add the sub item into $\mathcal{L}_{invalid}$
11:     **return** $\mathcal{L}_{valid}$
12: **function** CUT-1($\mathcal{L}_{valid}$)
13:     initialize items sequence $S = \emptyset$
14:     sort $\mathcal{L}_{valid}$ by $lz$ coordinate of each item in ascending order
15:     $s_i \leftarrow item_i$'s index in the sorted list
16:     **return** $S$
17: **function** CUT-2($\mathcal{L}_{valid}$)
18:     initialize height map $\mathcal{H}_n \in Z^{L \times W}, \mathcal{H}_n = 0^{L \times W}, S = \emptyset$
19:     **while** $\mathcal{L}_{valid} \neq \emptyset$ **do**
20:         randomly pop an $item_i$ from $\mathcal{L}_{valid}$ satisfy $lz_i = \mathcal{H}_n(item_i)$
21:         add the $item_i$ into $S$
22:         $\mathcal{H}_n(item_i) \leftarrow \mathcal{H}_n(item_i) + h_i$
23:     **return** $S$

---

The sequence in RS benchmark is generated by random sampling. Each item along the sequence is picked out of our pre-defined item set $\mathcal{I}$ randomly. However, as everything is random, we do not know the optimal packing configuration of a RS sequence ourselves (unless we run an exhaustive branch-and-bound search [8]). For a better quantitative evaluation, we also generate item sequences via cutting stock [50]. It is clear that a sequence created by cutting the bin should be packed in bin perfectly with a perfect space utilization of $100\%$. Algorithm 3 provides the detailed procedures of the data generation.

## 6.4 Heuristic Baseline Method

Online BPP is an under-investigated problem. To better demonstrate the effectiveness of our method, we design a heuristic baseline approach to evaluate the performance of our DRL method. We report details of this baseline approach in this section. This method is designed based on a simple observation, human would try to keep the volume of packed bin to be as regular as possible during the packing. Such "regularity" is used as the metric to measure a packing action.

To describe regularity of a bin, we introduce the concept of *spare cuboid*. As shown in Figure 9, a spare cuboid is an unoccupied, rectangular space in the bin, and the regularity of a bin is defined based on the *maximum* spare cuboids. Intuitively, we would like to have a bigger maximum spare cuboid. If a packed bin has many small-size spare cuboids, it implies the remaining space of this bin is not "regular". As illustrated in Figure 9, packing the blue item leads to two different remainder spaces for future packing. The



Figure 9: The maximum spare cuboid.

regularity of the bin is then defined as the maximum rectangular residual space or *maximum spare cuboid*. Since $\mathcal{I}$ is pre-defined, we know how many items can be packed into a maximum spare cuboid. Based on this, we rate each maximum spare cuboid $c$ by the number of item types can be packed in $RS_c = \|\mathcal{I}_{valid}\| + c_{volume}, \mathcal{I}_{valid} \subset \mathcal{I}$. If a maximum spare cuboid fits all the items in $\mathcal{I}$, additional reward is given as: $RS_c = \|\mathcal{I}\| + c_{volume} + 10$. The final score $BS_p$ of a bin by packing the next item at $p$ would be the sum of $RS_c$ of its maximum spare cuboid $c$. And we can find the best packing position $p_{best}$ as:
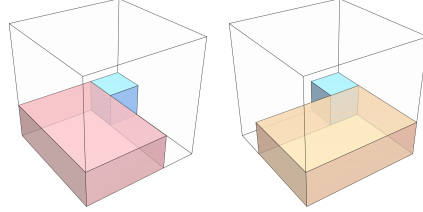
$$p_{best} = \arg\max_p \frac{1}{\|\mathcal{C}\|} \sum_{c \in \mathcal{C}} RS_c \qquad (3)$$
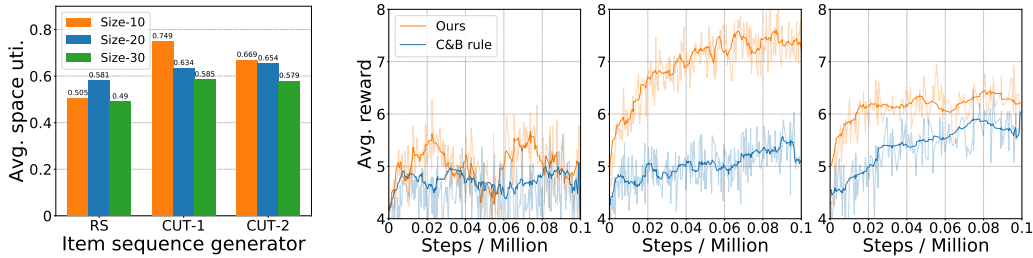


Figure 10: Left: Packing performance on different resolution. Size-10 means the resolution of the test bin is $10 \times 10 \times 10$ and etc. Second to right: Imposing the C&B rule leads to inferior performance (lower average reward).

## 6.5 More Results

Figure 14 shows more packing results on three different benchmarks. An animated packing can be found in the supplemental video. We also investigate how the resolution of the bin would affect the our performance. In this experiment, we increase the spatial discretization from $10 \times 10 \times 10$ to $20 \times 20 \times 20$ and $30 \times 30 \times 30$. As shown in Figure 10, the performance slightly decreases. This is because we do not increase the item types in $\mathcal{I}$ ($|\mathcal{I}|$ remains 64). Increased discretization widens the distribution of possible action space and dilutes the weight of the optimal action. However, our

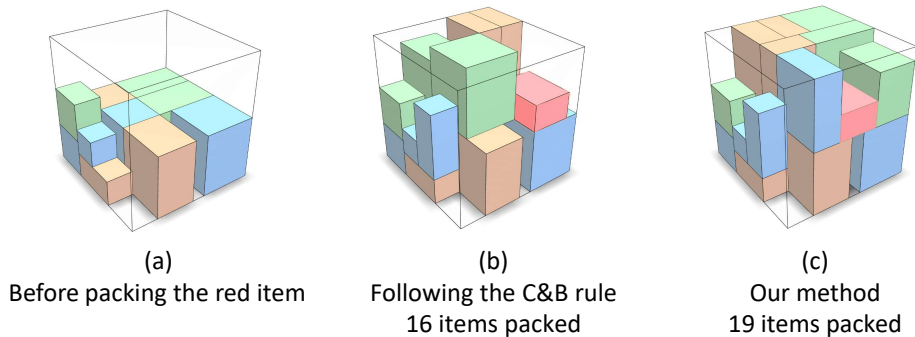|  (a) | (b) | (c) |
|:---:|:---:|:---:|
| Before packing the red item | Following the C&B rule<br>16 items packed | Our method<br>19 items packed |

Figure 11: Visual comparison between learned vs. heuristic strategy. Different packing strategy for the red item would lead to different performance.

method remains efficient even when the problem complexity is $\sim 27\times$ bigger. This experiment demonstrates a good scalability of our method in a high-resolution environment.

**Learned vs. heuristic strategy**    In real-world bin packing, human tends to place a box to touch the sides or corners of the bin or the other already placed boxes. We refer to this intuitive strategy as *corner & boundary rule* (C&B rule). An interesting discovery from our experiment is that our method can automatically learn the C&B rule. In addition, imposing such constraints explicitly leads to inferior performance. We found that the performance (average reward) drops about $20\%$ when adding such constraints, as shown in right of Figure 10. By checking into the placement sequences, we found that our agent can decide when to follow the rule to obtain a globally more optimal packing; see supplemental material for details. This can also be verified by the experiment in Table 6.

Table 6: Evaluating the effect of boundary rule.

|  | Space uti. | # packed items |
|---|:---:|:---:|
| w/o corner&boundary rule | 66.9% | 17.5 |
| w corner&boundary rule | 60.9% | 16.2 |

Table 7: Space utilization of unseen items.

|  | $64 \rightarrow 64$ | $64 \rightarrow 125$ |
|---|:---:|:---:|
| CUT-2 | 66.9% | 41.4% |

To illustrate why our agent can decide when to follow the C&B rule to obtain a globally more optimal packing, we give a visual example here. As shown in Figure 11 (b), if the agent exactly follow the C&B rule when packing the red item, it will leave gaps around the item. However, our method (Figure 11 (c)) can make a decision of packing the item in the middle upon the yellow and blue ones. Our method is trained to consider the whether there is enough room for next moves but not only takes the current situation into consideration. This move reserves enough space around the red item for the following item and this decision makes our method packing 3 more items when dealing with a same sequence.

**Generalizability with unseen items**    We also test our method with unseen items. In this experiment, we increase the number of pre-defined item (i.e., $|\mathcal{I}|$) from $64$ to $125$ in the test. All the items are generated with CUT-2, but their dimensions in the test may not be seen in the training. The result is presented in Table 7. It shows that our method does demonstrate some generalizability and provides a reasonable benchmark. However, we do not want to count on that: the training should cover as many different items as possible in order to secure a satisfactory result in general. After all, MDP is only for the fixed environment.

## 6.6   User Study

Figure 12 is the interface of our user study app, which consists of two parts: visualization and action space. The test sequences is randomly picked from CUT-2 test set. To make the competition fair, the ground-truth is also provided to human users as suggestions. Users can drag our UI to change the angle of view thus having a full observation of the packed items. To help users make better decision,
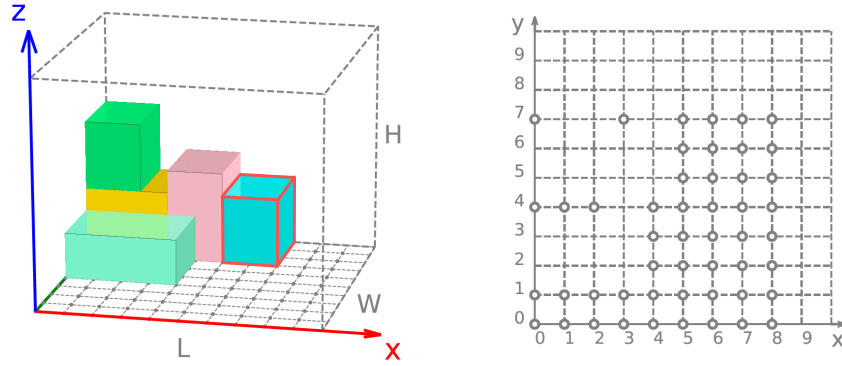
Figure 12: Left: The 3D visualization of packed items in our user study app. The current item is highlighted with red frame. Right: The action space of the bin — the empty circles indicate suggested placements for the current item.
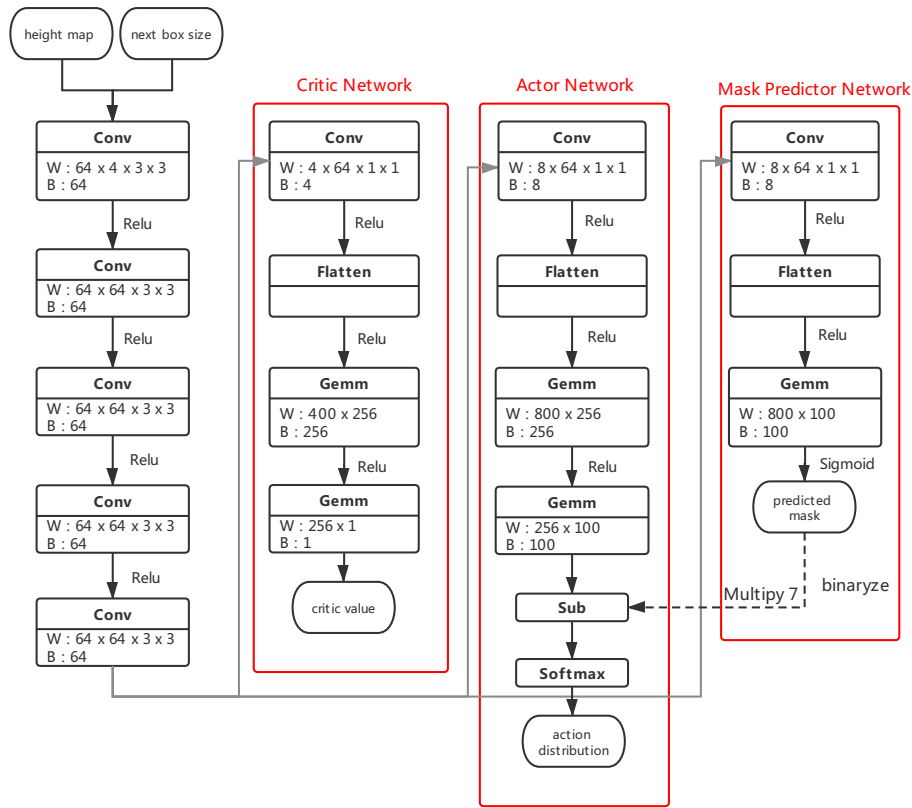


Figure 13: Detailed network architecture.

our app allow them choose any suggestion circle in action space and virtually place item before they make the final decision. No time limits is given. When there is no suitable place for the current item, the test will reset and the selected sequence will be saved.
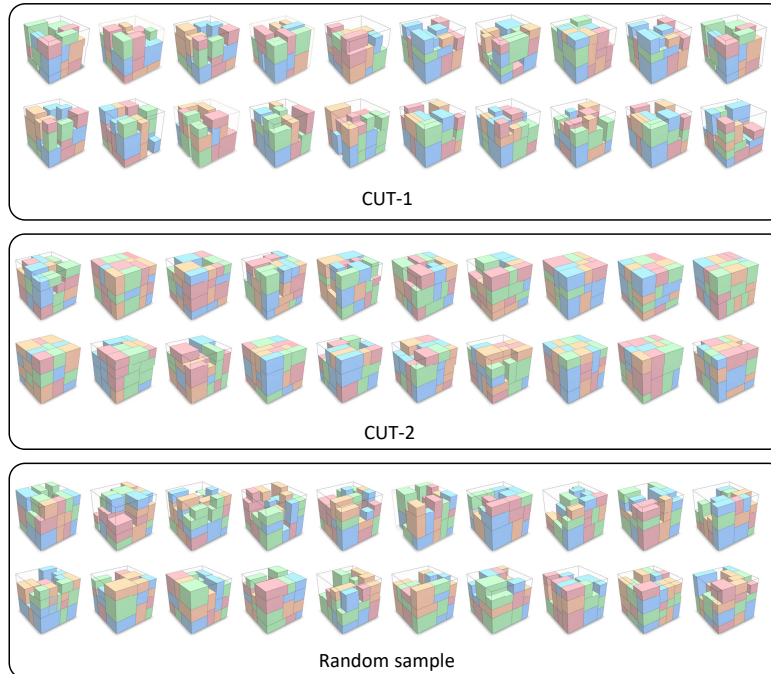
Figure 14: Packing results of our BPP-1 model.