# Leaning Heuristics for Weighted CSPs through Deep Reinforce Learning

**Dingding Chen**[1] · **Ziyu Chen**[1] · **Zhongshi He**[1] · **Junsong Gao**[1] · **Zhizhuo Su**[2] ·

**Abstract** Constraint Programming (CP) is at the core of Artificial Intelligence (AI) and has been applied to many real-world scenarios. Weighted Constraint Satisfaction Problems (WCSPs) are one of the most important CP models aiming to find a cost-minimal solution. However, due to its NP-hardness, solving a WCSP usually requires efficient heuristics to explore high-quality solutions. Unfortunately, such heuristics are hand-crafted, i.e., they rely on domain-specific knowledge and may not be generalizable across different cases. On the other hand, although Deep Learning (DL) has been proven to be a promising way to automatically learn heuristics for combinatorial optimization problems, the existing DL-based methods are unsuitable for WCSPs since they fail to exploit the problem structure of WCSPs. Besides, such methods are often based on Supervised Learning (SL), making the learned heuristics less efficient since generating optimal labels for large problem instances could be infeasible. To address the above issues, we propose a novel Deep Reinforcement Learning (DRL) framework that is able to train the model on large-scale problems, so as to the model could mine more sophisticated patterns of problems and provide high-quality solution construction heuristics for WCSPs. By exploiting the problem structure, we effectively decompose the problem by using a pseudo tree, and formulate the solution construction process as an Markov Decision Process (MDP) with multiple independent transition states. With Graph Attention Networks (GATs) parameterized deep Q-value network, we learn the optimal Q-values through a modified Bellman Equation that considers the multiple transition states, and the solution construction heuristics are extracted from the learned Q-value network. Besides constructing solutions greedily, our heuristics can also be applied to many meta-heuristics such as Beam Search (BS) and Large Neighborhood

✉ Ziyu Chen
E-mail: chenziyu@cqu.edu.cn
[1] College of Computer Science, Chongqing University, Chongqing 400044, China
[2] WMG, University of Warwick, Coventry, CV47AL, United Kingdom

Search (LNS). Finally, we demonstrate the effectiveness of our proposed DRL framework by comparing the heuristics obtained from our DRL model with those derived from the SL model. Extensive empirical evaluations show that our DRL-boosted algorithms significantly outperform their counterparts with traditional tabular-based heuristics (e.g., mini-bucket elimination, MBE) and state-of-the-art methods on benchmark problems.

**Keywords** WCSP · Incomplete WCSP Algorithm · Heuristics · DRL · GATs

## 1 Introduction

Constraint Satisfaction Problems (CSPs) [12] provide a unified framework for general problem modeling and solving in Artificial Intelligence (AI). A CSP consists of a fixed set of variables, each with an associated domain of values, and a set of hard constraints on these variables to specify the allowable assignment combinations. Weighted Constraint Satisfaction Problems (WCSPs) [38] are a generalization of CSPs where the constraints are no longer "hard". Instead, each tuple in a constraint, i.e., all variable assignments involved in that constraint, is associated with a positive real value as a weight (usually called a "cost") to specify the violation degree of constraints. The objective of solving a WCSP is to find an assignment for all variables that minimizes the aggregated constraint cost. WCSPs have been successfully applied into many real-world applications, including supply-chain management [20], judge assignment [21], bioinformatics [43, 48] and many others.

Algorithms for WCSPs can be classified into two categories, i.e., complete algorithms and incomplete algorithms. Complete algorithms guarantee to find the optimal solution and can be broadly classified into inference-based algorithms [11] and search-based algorithms [28]. Search-based algorithms systematically explore the entire solution space by branch-and-bound. Instead, inference-based algorithms employ a dynamic programming paradigm to solve a WCSP. Since solving WCSPs is NP-Hard, complete algorithms exhibit an exponentially increasing computation overhead and cannot scale up to large real-world applications. Therefore, there has been considerable research effort to develop incomplete algorithms to find high-quality solutions at an acceptable computational overhead. Incomplete algorithms generally follow three strategies, i.e., local search [52, 32, 23], belief propagation [16, 9] and sampling [33, 31]. However, existing incomplete algorithms usually rely on hand-crafted heuristics that require domain-specific knowledge to tune for different settings, such as the large neighborhood selection policy in large neighborhood search [35] and the temperature schedule in simulated annealing [7].

Deep Learning (DL) has been proven to be promising in learning algorithms for solving NP-hard problems [4]. Several works have tried to leverage DL to learn effective heuristics automatically for existing methods, such as variable selection for a branch-and-bound algorithm [5] or a stochastic local search algorithm [51], and branching heuristics for a backtracking search algorithm [26]. Unfortunately, the learned heuristics are only applicable to a particular

algorithm and usually cannot be generalized to other algorithms. Recently, Deng et al. [14] proposed a supervised pre-trained model to generate heuristics for a wide range of algorithms. However, generating optimal labeled data could be expensive, which makes the model only trained on small-scale problems and thus affects the quality of the learned heuristics.

Instead, we propose to train the model on large-scale problems with Deep Reinforcement Learning (DRL) to generate efficient solution construction heuristics for WCSPs. Specifically, our main contributions are listed as follows:

- We propose a novel Deep Reinforcement Learning (DRL) framework where the model could be trained on large-scale problems. As a result, our model could capture more sophisticated patterns of problems to provide high-quality solution construction heuristics for WCSPs. To improve computational efficiency, we propose to use a pseudo tree to decompose the problem effectively, and formulate the solution construction process of a WCSP as an Markov Decision Process (MDP) with multiple independent transition states by exploiting the problem structure. In our MDP formulation, the state is the remaining sub-problem to be solved, action space is the domain of the target variable, and reward is the change of the current cost.
- We propose to use Graph Attention Networks (GATs) [46] parameterized deep Q-value network to learn the optimal Q-values through a modified Bellman Equation that considers the multiple transition states. And, the solution construction heuristics are derived from the learned Q-value network. Besides constructing greedy solutions, we embed our heuristics into meta-heuristics including Beam Search (BS) [40] and Large Neighborhood Search (LNS) [41]. Specifically, our heuristics can be used as a quality evaluation of the partial assignments in BS or as a sub-routine of LNS to repair destroyed variables via constructing a solution greedily.
- Extensive empirical evaluations indicate that the heuristics derived from our DRL model are superior to those obtained from the SL model [14]. Furthermore, we empirically demonstrate the effectiveness of our heuristics by combining them with meta-heuristics including Beam Search (BS) and Large Neighborhood Search (LNS) on various standard benchmarks.

The rest of this paper is organized as follows. We briefly review related work in Sect. 2. The preliminaries, including WCSPs, pseudo tree, GATs, Markov Decision Process (MDP), beam search and large neighborhood search, are presented in Sect. 3. In Sect. 4, we describe our proposed method, including the MDP formulation of WCSPs, graph representation and graph embedding, the DRL training algorithm, and heuristics for beam search and large neighborhood search. Finally, we present the empirical evaluation of our proposed method in Sect. 5 and our conclusion in Sect. 6.

## 2 Related Work

In this section, we first review traditional tabular-based algorithms for WCSPs, and then introduce state-of-the-art deep learning-based methods for constraint reasoning.

2.1 Tabular-based Algorithms for WCSPs

The complete approaches for WCSPs employ either inference or systematic search. Branch-and-Bound (BnB) is the most significant landmark of search-based algorithm. Given a variable ordering, BnB sequentially explores the whole solution space in a depth-first fashion. Pruning happens when the current lower bound is higher than the known upper bound. Therefore, the bound quality is the key of efficient pruning. Currently, several works attempt to tighten the bounds via local consistency [27,10] or combining the benefit of the best-first search and depth-first search strategies [1]. On the other hand, bucket elimination [11,34] is a well-known algorithm for WCSPs that performs dynamic programming on the variable ordering (e.g., a pseudo tree). The algorithm forwards the assignment combination utilities bottom-up and then reversely propagates the optimal decisions along with the variable ordering. However, its memory consumption is exponential in the induced width [12]. Therefore, the memory-bound bucket elimination method [6] was proposed to trade the runtime for smaller memory consumption. Since solving WCSPs is NP-hard, complete algorithms require exponential computation overheads, making them unsuitable for large-scale practical applications.

Incomplete algorithms can be generally classified into local search, belief propagation-based and sampling-based algorithms. Inspired by Monte-Carlo tree search, a sampling-based algorithm with confidence bounds [33] was proposed for solving WCSPs. However, the algorithm requires an exponential memory with the number of variables, limiting its application to large-scale problems. Therefore, Nguyen et al. [31] proposed to map a WCSP to a maximum a-posteriori estimation problem and solve it by Gibbs sampling. Local search algorithms [52,32] are the most popular incomplete methods for WC-SPs, where each variable optimizes its assignment based on its local constraints and assignments of all its neighbors. Unfortunately, these algorithms tend to converge to local optima. Therefore, K-optimality (K-OPT) [47] was proposed to improve the solution of local convergence by optimizing the assignments of all variables within the K-size coalition. However, the computational effort required to find the K-optimal solution grows exponentially as K increases. Max-Sum [16] is a classic belief propagation-based algorithm that gathers global information by propagating and accumulating the maximum utility through the whole factor graph. Unfortunately, Max-Sum is only guaranteed to converge in cycle-free problems. Thus, Damped Max-Sum (DMS) [9] was proposed to improve the convergence probability of belief propagation by reducing the influence of information circulation propagation. However, incomplete algo-

rithms still either require considerable computational efforts (e.g., K-OPT) or exhibit poor performance (e.g., local search) in general problems. Therefore, by leveraging its powerful representational capability, we resort to deep neural networks to embed WCSPs to learn efficient and effective heuristics.

## 2.2 Deep Learning-based Methods for Constraint Reasoning

In recent years, methods that use deep learning techniques for constraint reasoning have received increasing attention, and can be generally divided into Supervised Learning (SL) based and Deep Reinforcement Learning (DRL) based methods. The classical SL-based method is NeuroSAT [39], which uses a message-passing neural network to solve SAT (Boolean satisfiability) problems in an end-to-end fashion and further decode the satisfactory assignments. In the method, a SAT problem is encoded as an unweighted bipartite graph to provide permutation and variable relabeling invariance. Since the constraint costs of WCSPs are arbitrary positive real values, the unweighted bipartite graph is unsuitable for representing WCSPs. In addition, there are some works that can be adapted to solve WCSPs, but they have their own limitations. For example, Xu et al. [50] proposed to encode a binary CSP as a matrix and train a convolutional neural network to predict the satisfiability, but the representation scheme cannot scale to arbitrary problem sizes. Instead, Galassi et al. [19] proposed to construct a feasible solution of a CSP instance by extending a partial assignment with a trained deep neural network. Still, this approach is restricted to problems of a pre-determined size.

Recently, Razeghi et al. [36] and Deng et al. [15] proposed to use Multi-layer Perceptrons (MLPs) to parameterize the high-dimensional data in bucket elimination [11] and regret matching [22], respectively. But, the online learning nature prevents these methods from generalizing to unseen instances. Furthermore, training MLPs requires a considerable time. Therefore, Deng et al. [14] proposed using SL to learn a pre-trained model to construct effective heuristics for a broad range of algorithms. However, constructing optimal labeled data in SL could be infeasible in large problem instances. Instead, there is another line of works that exploit DRL to learn heuristics for existing algorithms. For instance, Yolcu and Poczos [51] proposed to use Graph Neural Networks (GNNs) [37] to encode SAT instances and REINFORCE [49] to learn satisfying assignments inside a stochastic local search procedure. Besides, learning branching heuristics for constraint programming via DRL has been intensively investigated [5]. Unfortunately, these learned heuristics are solely devoted to a particular algorithm. Instead, our proposed DRL framework aims to learn solution construction heuristics that can be embedded in many different algorithms.
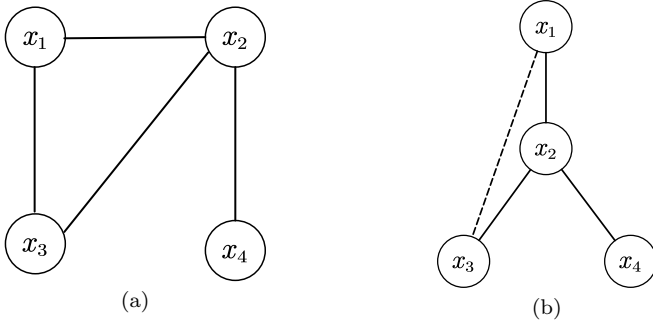
Fig. 1: An example of constraint graph and pseudo tree

## 3 Background

3.1 Weighted Constraint Satisfaction Problems

A weighted constraint satisfaction problem (WCSP) can be defined by a tuple $(X, D, F)$ such that:

- $X = \{x_1, x_2, ..., x_n\}$ is a set of variables.
- $D = \{D_1, D_2, ..., D_n\}$ is a set of finite variable domains. Each domain $D_i \in D$ consists of a set of finite allowable values for variable $x_i \in X$. Besides, we denote the maximal domain size as $d = \max_{D_i \in D} |D_i|$.
- $F = \{f_1, f_2, ..., f_m\}$ is a set of constraint functions. Each function $f_i : D_{i_1} \times D_{i_2} \times \cdots \times D_{i_k} \to \mathbb{R}_{\geq 0}$ specifies a non-negative cost to each value combination of the involved variables $x_{i_1}, x_{i_2}, ..., x_{i_k}$.

For the sake of simplicity, we assume that all constraint functions are binary (i.e., $f_{ij} : D_i \times D_j \to \mathbb{R}_{\geq 0}$). A partial assignment $\Gamma$ is a set of assignments in which each variable appears at most once. The cost of $\Gamma$ is calculated by aggregating the cost of all constraints involving only the variables that appear in the partial assignment. That is,

$$cost(\Gamma) = \sum_{f_{ij} \in F, x_i, x_j \in \Gamma} f_{ij} \left( \Gamma_{[x_i]}, \Gamma_{[x_j]} \right) \tag{1}$$

where $\Gamma_{[x_i]}$ is the assignment of $x_i$ in $\Gamma$. A partial assignment that includes all variables in $X$ is a full assignment. The goal of solving a WCSP is to find a full assignment to minimize the aggregated constraint cost. That is,

$$X^* = \underset{d_i \in D_i, d_j \in D_j}{\arg\min} \sum_{f_{ij} \in F} f_{ij} \left( d_i, d_j \right) \tag{2}$$

A WCSP can be visualized by a constraint graph where a vertex represents a variable and an edge represents a binary constraint. Fig. 1(a) gives the constraint graph of a WCSP with four variables and four binary constraints.

### 3.2 Pseudo Tree

A pseudo tree [18] arrangement of a constraint graph is a tree with the same nodes and edges as the original graph. It has the property that adjacent nodes of the original graph fall in the same branch of the tree. Therefore, the sub-problems in different branches can be solved independently. A pseudo tree can be generated by a depth-first traversal of the constraint graph, categorizing the constraints into tree edges and pseudo edges (i.e., non-tree edges). Fig. 1(b) presents a possible pseudo tree derived from Fig. 1(a), where tree edges and pseudo edges are shown as solid lines and dashed lines, respectively. For a variable $x_i$ in the pseudo tree, we denote the ancestor connected with $x_i$ through a tree edge as its parent $P(x_i)$. The ancestors connected with $x_i$ through back edges as its pseudo parents $PP(x_i)$, and the descendants connected with $x_i$ through tree edges as its children $C(x_i)$. In addition, we denote all parents of $x_i$ as $AP(x_i)$, i.e., $AP(x_i) = PP(x_i) \cup \{P(x_i)\}$. Taking $x_2$ in Fig. 1(b) as an example, we have $P(x_2) = x_1$, $PP(x_2) = \emptyset$, $C(x_2) = \{x_3, x_4\}$, and $AP(x_2) = \{x_1\}$.

### 3.3 Graph Attention Networks

Graph Attention Networks (GATs) [46] are a convolutional neural network architecture operated on graph-structured data. They are constructed by stacking several graph attention layers such that nodes can attend over the features of their neighbors. Furthermore, to indicate the importance of different nodes in the neighborhood, GATs introduce a self-attention mechanism that assigns different weights to these different nodes. In the self-attention mechanism, the calculation of the attention coefficient between each pair of neighbor nodes is:

$$e_{ij} = \mathbf{a}(\mathbf{W}h_i, \mathbf{W}h_j) \tag{3}$$

where $h_i, h_j \in \mathbb{R}^d$ are the features of nodes $i, j$, $\mathbf{W} \in \mathbb{R}^d \times \mathbb{R}^d$ is a weight matrix, and $\mathbf{a}$ is a single-layer feed-forward neural network.

Then, the attention coefficients are normalized using the softmax function to facilitate comparison between different nodes. The normalization coefficient across node $j$ is computed by:

$$\alpha_{ij} = \frac{e_{ij}}{\sum_{k \in N_i} \exp(e_{ik})} \tag{4}$$

where $N_i$ is the neighborhood of node $i$ in the graph (including $i$). These normalized attention coefficients are then computed by the linear combination of their corresponding features and served as the final output feature for each node. That is,

$$h_i' = \sigma \left( \sum_{j \in N_i} \exp(\alpha_{ij} \mathbf{W}h_j) \right) \tag{5}$$

where $\sigma$ is a nonlinear function, such as LeakyReLU or sigmoid.

In order to stabilize the learning process of self-attention, GATs adopt Multi-head attention [45] in which K independent attention mechanisms are performed. Their features are averaged as:

$$h_i' = \sigma \left( \frac{1}{K} \sum_{k=1}^{K} \sum_{j \in N_i} \exp\left(\alpha_{ij}^k \mathbf{W}^k h_j\right) \right) \tag{6}$$

where $\alpha_{ij}^k$ are normalized attention coefficients computed by the k-th attention mechanism ($\mathbf{a}_k$), and $\mathbf{W}^k$ is the corresponding input linear transformation's weight matrix.

3.4 Markov Decision Process

An Markov Decision Process (MDP) is formulated as a tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{T}, r, \gamma \rangle$, where

- $\mathcal{S}$ is a set of environment and agent states.
- $\mathcal{A}$ is a set of actions of the agent.
- $\mathcal{T}(s, a, s') = \Pr(s, a, s')$ is the transition function, specifying the probability of the transition from an initial state $s$ to a new state $s'$ under an action $a$. Here, $s, s' \in S$ and $a \in A$.
- $r(s, a)$ is the reward function, defining the immediate reward after taking an action $a$ at the state $s$.
- $\gamma \in [0, 1)$ is the discount factor, used to weight the preferences for immediate reward relative to future reward.

The agent interacts with the environment following a policy (i.e., $\pi : \mathcal{S} \times \mathcal{A} \to [0, 1]$) which describes the probability of taking an action from a given state. At each decision point $t$, the agent observes the current state $s_t \in \mathcal{S}$, selects an action $a_t \in \mathcal{A}$ according to the policy $\pi$, and receives an immediate reward $r(s_t, a_t)$. The goal of the agent is to learn the optimal Q-function ($Q_{\pi^*}(s, a)$), an action-value function that estimates the sum of future rewards after taking an action $a$ at the state $s$ and following the optimal policy $\pi^*$, i.e.,

$$Q_{\pi^*}(s, a) = \max_\pi Q_\pi(s, a)$$
$$= \max_\pi \mathbb{E}_{\pi, \mathcal{T}} \left[ r(s_t, a_t) + \gamma r(s_{t+1}, a_{t+1}) + \cdots \mid s_t = s, a_t = a, \pi \right] \tag{7}$$

The optimal Q-function is the fixed-point of Bellman Equation [3], i.e.,

$$Q_{\pi^*}(s, a) = \mathbb{E}_{\pi, \mathcal{T}} \left[ r(s, a) + \gamma \max_{a'} Q_{\pi^*}(s', a') \right] \tag{8}$$

According to the Bellman Equation, the Q-learning algorithm [44] can be derived. However, the Q-learning algorithm suffers from the curse of dimensionality, becoming less efficient as the dimensionality of the environment

increases. Therefore, Deep Q-network (DQN) [29] was proposed to use a multi-layered neural network (i.e., the main Q-network $Q_\theta(s,a)$) to parameterize the optimal Q-function. DQN can handle complicated decision processes with large and continuous state space by directly outputting the predicted Q-value of each state-action pair with state features as input.

To train the main Q-network efficiently, DQN introduces an experience replay memory and a separate target network. A First-In-First-Out (FIFO) experience replay memory is established to store the agent's experience to remove the correlation between consecutive transitions. As a result, the high variance in parameter updates and the instability of the training processes can be reduced. The update of the parameter of the main Q-network is based on a mini-batch of samples randomly drawn from the experience replay memory. The target Q-network ($Q_{\bar{\theta}}(s,a)$) is introduced to improve the stability of the training process that uses temporal difference error, i.e.,

$$L(\theta) = \left( Q_\theta(s,a) - r(s,a) - \gamma \max_{a'} Q_{\bar{\theta}}(s',a') \right)^2 \tag{9}$$

At each time step, the parameter of the target Q-network is updated towards the main Q-network according to $\bar{\theta} \leftarrow \rho\bar{\theta} + (1-\rho)\theta$, where $\rho$ is the interpolation factor in Polyak averaging for the target Q-network and between 0 and 1, usually close to 1.

### 3.5 Beam Search and Large Neighborhood Search

Beam Search (BS) [40] is an incomplete derivative of Branch-and-Bound (BnB) [28]. Different from the systematic search of BnB, BS selects the partial assignments that satisfy the optimal solution conditions to branch. It uses breadth-first search with an impermissible pruning rule to construct its search tree. That is, at each level of the search tree, only the predetermined number (called the beam width) of the best partial assignments is selected for further branching, while the remaining partial assignments are permanently pruned. The beam width defines the memory required to perform the search in BS. The larger the beam width, the fewer the partial assignments pruned. When the beam width is infinite, no partial assignment is pruned, and thus BS is the same as breadth-first search.

Large Neighborhood Search (LNS) [41] is a meta-heuristic that iteratively explores complex neighborhoods in a search space to find better candidate solutions. The algorithm iteratively improves an initial solution by repeatedly performing destroy and repair phases. Specifically, a subset of variables (called Large Neighborhood, LN) is selected to discard their current values in the destroy phase. During the repair phase, new assignments are found for the LN variables, provided the undestroyed variables retain their value from the previous iteration. Compared with other local search techniques, LNS is characterized by exploring a LN at each step, with the aims of making the search process away from local optima and finding better candidate solutions.
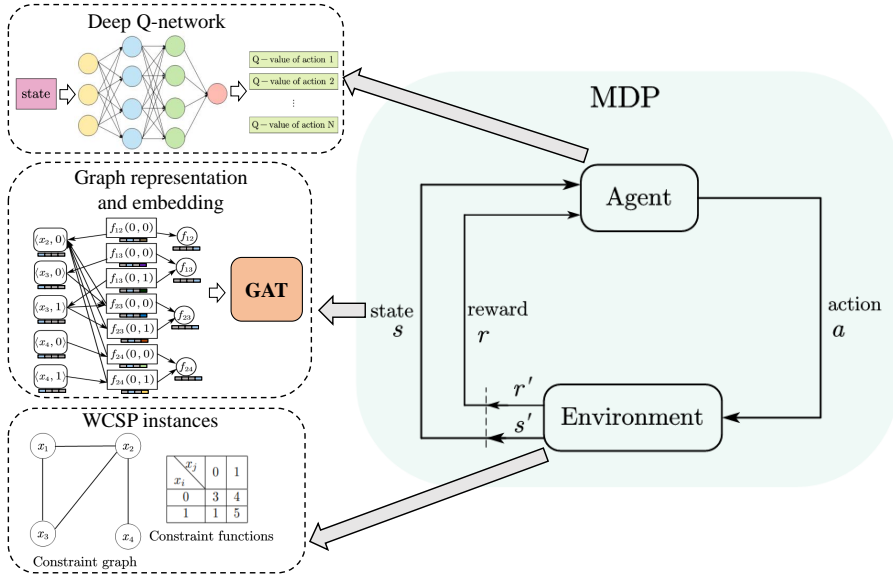
Fig. 2: The pipeline of our proposed method

## 4 The Proposed Method

The pipeline of our proposed DRL framework is shown in Fig. 2. We first re-formulate a WCSP according to MDP. Then, we give the graph representation and graph embedding of a WCSP instance, i.e., a parametric function that encodes the input problem and outputs the predicted Q-values. Next, we present a DRL training algorithm, which determines how to learn the parameter of the Q-network. Finally, we embed the learned Q-network into two meta-heuristics, including beam search and large neighborhood search.

### 4.1 The MDP Formulation

Given a variable ordering, the solution to a WCSP instance can be constructed by sequentially extending a partial assignment until a full assignment is reached. Therefore, it is natural to model the solution constructing process as an MDP where the state is the remaining sub-problem to be solved (i.e., the problem $P$ under the partial assignment $\Gamma$ and $x_i$ as the target variable), the action space is the domain of $x_i$ and the reward is the change of the current cost incurred by a new assignment of $x_i$, respectively.

A straightforward approach to order variables is assuming a global ordering over all variables, e.g., alphabetic ordering. However, such ordering ignores the problem structure and fails to decompose the problem efficiently. Taking Fig. 3(c) as an example, the problem to be solved contains $x_3$ and $x_4$. Since there is no constraint between $x_3$ and $x_4$, the problems of selecting values for $x_3$ and $x_4$
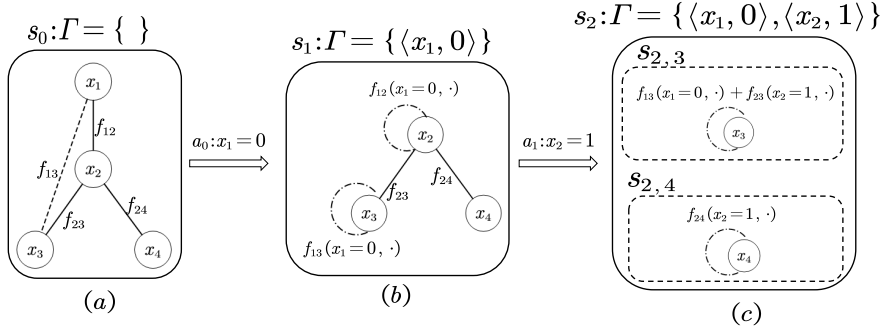
Fig. 3: The solution construction process in WCSP solving as an MDP

can be solved independently, thereby reducing the computational overheads. To cope with this issue, we propose to order variables by using a pseudo tree.

We denote the problem $P$ given the partial assignment $\Gamma$ and a target variable $x_i$ as $\langle P, \Gamma, x_i \rangle$. Since we use a pseudo tree as the ordering, after assigning $x_i$ with a value $d_i \in D_i$, the problem is reduced to multiple independent sub-problems, i.e., each of them rooted at a child of $x_i$. Therefore, the transition state in our MDP formulation consists of a set of independent sub-states, i.e., $s' = \{s_c | \forall x_c \in C(x_i)\}$ where $s_c = \langle P, \Gamma \cup \{\langle x_i, d_i \rangle\}, x_c \rangle$. As a result, the Bellman Equation in our MDP formulation needs to be modified as:

$$Q_{\pi^*}(s, a) = r(s, a) + \gamma \sum_{x_c \in C(x_i)} \max_{a_c} Q_{\pi^*}(s_c, a_c) \qquad (10)$$

where $a_c$ is the action of $x_c$.

Formally, the state, action, transition, and reward in our MDP formulation are defined as follows:

− State: a state $s$ is the WCSP instance $P$ instantiated with $\Gamma$ and $x_i$ as the target variable. That is, $s = \langle P, \Gamma, x_i \rangle$.
− Action: an action $a = \langle x_i, d_i \rangle$ corresponds to assigning $x_i$ with a value $d_i \in D_i$.
− Transition: transition is deterministic, i.e, $\Pr(s, \langle x_i, d_i \rangle, s') = 1$. Since we use a pseudo tree to order variables, after taking an action $\langle x_i, d_i \rangle$ at the state $s = \langle P, \Gamma, x_i \rangle$, the transition state is $s' = \{s_c | \forall x_c \in C(x_i)\}$ where $s_c = \langle P, \Gamma \cup \{\langle x_i, d_i \rangle\}, x_c \rangle$.
− Reward: the reward function $r(s, \langle x_i, d_i \rangle)$ is defined as the change of the current cost when assigning $x_i$ with $d_i$ at $s = \langle P, \Gamma, x_i \rangle$, i.e.,

$$r(s, \langle x_i, d_i \rangle) = - \sum_{x_j \in AP(x_i)} f_{ij} \left( d_i, \Gamma_{[x_j]} \right) \qquad (11)$$

It is worth mentioning that the objective function of our MDP formulation is aligned with that of WCSPs when the discount factor $\gamma = 1$. That is because
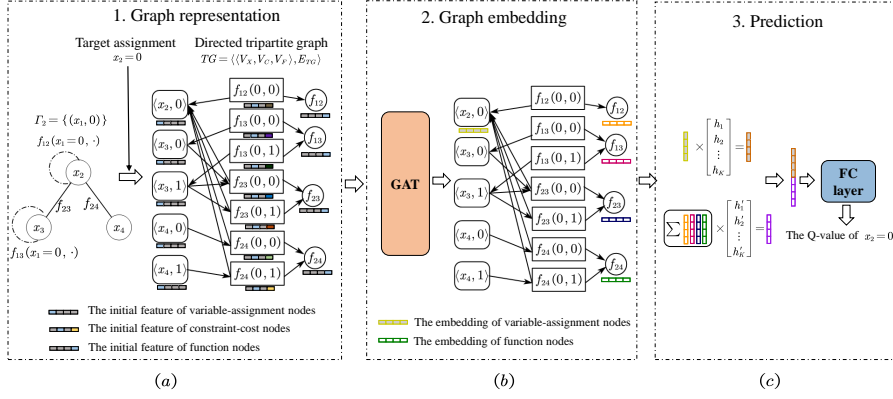
Fig. 4: The graph representation and embeddings of a partially instantiated WCSP

the negative optimal Q-value of $s = \langle P, \Gamma, x_i \rangle$ is equal to the optimal cost of the sub-problem rooted at $x_i, \forall x_i \in X$. When $x_i$ is the root variable, the negative optimal Q-value of $s$ is equal to the optimal cost of the problem.

## 4.2 Graph Representation and Graph Embedding

A critical step of applying DL to solve combinatorial optimization problems is to build an appropriate structure. Selsam et al. [39] proposed to encode a SAT as an unweighted bipartite graph, where each literal and clause corresponds to a node, and the association between the literal and clause corresponds to an edge. However, the problem structure of WCSPs is quite different from that of SAT in many ways. For instance, the constraint costs of WCSPs are arbitrary positive real values, which cannot be represented by an unweighted bipartite graph. Furthermore, a variable in a WCSP may have more than two values in its domain, while variables in a SAT are restricted to a boolean domain (i.e., Ture or False). Xu et al. [50] proposed to represent CSPs as a matrix form, where each entry indicates whether the corresponding assignment is allowed or not. Unfortunately, the representation scheme cannot scale to arbitrary problem sizes. Recently, Song et al. [42] proposed to represent the underlying constraint network of a CSP as GNNs, with variables as vertices and constraints as edges. This method also cannot explicitly handle the cost values in each constraint, making it unsuitable for WCSPs. Thus, we opt for representing a WCSP as a directed and acyclic tripartite graph [14].

The generation of the graph representation consists of the following three steps: first, the partially instantiated WCSP instances are converted into a microstructure representation [24]; then the microstructure is compiled into a tripartite graph; finally, the tripartite graph is transformed to a directed acyclic graph (cf. Fig. 4(a)). The microstructure representation of a WCSP

instance $P = (X, D, F)$ given the partial assignment $\Gamma$ is a weighted undirected graph, where each variable assignment that is compatible with $\Gamma$ corresponds to a vertex, and the constraint cost between a pair of vertices is represented by a weighted edge. Afterward, the microstructure is compiled into a tripartite graph $TG = \langle (V_X, V_C, V_F), E_{TG} \rangle$ where $V_X$ and $V_C$ correspond to variable assignment nodes and constraint costs (i.e., weighted edges) in the microstructure, respectively. $V_F$ corresponds to the constraint function in the WCSP instance. $E_{TG}$ contains two kinds of edges, the first class of edges is the undirected edges between variable assignment nodes and their neighboring constraint cost nodes. The second class of edges is the directed edges between constraint cost nodes and their constraint function nodes.

Finally, by determining the indirect direction between $V_X$ and $V_C$, the directed and acyclic tripartite graph is obtained and implemented through the following two stages. At first, the constraint graph induced by the set of unassigned variables is constructed as a pseudo tree rooted at the target variable. Then, for each constrained variable pair $x_i$ and $x_j$ with $x_i \in AP(x_j)$ and their assignments (i.e., $\langle x_i, d_i \rangle$ and $\langle x_j, d_j \rangle$), the variable-assignment node of $\langle x_i, d_i \rangle$ is set to be the precursor of the constraint-cost node of $f_{ij}(d_i, d_j)$ and the node of $f_{ij}(d_i, d_j)$ is set to be the precursor of the node of $\langle x_j, d_j \rangle$.

Since a partially instantiated WCSP instance is represented as a directed tripartite graph, we use GATs to construct a model that learns graph embeddings (cf. Fig. 4(a) and (b)). In the GAT model, we set a four-dimensional initial feature vector $h_i^{(0)}$ for each node $i$. For the feature vector $h_i^{(0)}$, its first three elements represent the type of node $i$ (i.e., variable assignment node, constraint cost node, and function node) and its last element is set to be the constraint cost of node $i$ if it is a constraint-cost node, and 0 otherwise. Subsequently, $h_i^{(0)}$ is embedded through $T$ layers of the GAT by:

$$h_i^{(t)} = \sigma \left( \frac{1}{K} \sum_{k=1}^{K} \sum_{j \in N_i} \exp \left( \alpha_{ij}^k \mathbf{W}^{k,(t)} h_j^{(t-1)} \right) \right), \forall t = 1, 2, ..., T \qquad (12)$$

where $h_i^{(t)}$ is the embeddings of node $i$ in the $t$-th time step and $\{ \mathbf{W}^{k,(t)} | \forall k = 1, 2, ..., K \}$ is the parameter of the $t$-th layer of the GAT model.

Given a target variable assignment $\langle x_i, d_i \rangle$ and a partial assignment $\Gamma$, the optimal cost of the problem $P$ instantiated with $\Gamma \cup \{ \langle x_i, d_i \rangle \}$ can be predicted according to the embeddings of the variable-assignment node $i$ (i.e., $\langle x_i, d_i \rangle$) and the cumulative embeddings of all function nodes in the tripartite graph (cf. Fig. 4(c)). That is,

$$\hat{c}_i = W_3 \left( \text{Concat} \left( W_1 h_i^{(T)}, W_2 \sum_{v_j \in V_F} h_j^{(T)} \right) \right) + b_1 \qquad (13)$$

where $W_1$, $W_2$, $W_3$ and $b_1$ are the parameters of three 1-layer Multi-layer Perceptrons (MLPs). For the sake of description, we denote the GAT model trained by our DRL framework as $Q_\theta$.

---

**Algorithm 1** Deep Q-network for training the model

---

**Require**: number of training epochs $N$, number of training iterations $K$, problem distribution $\mathcal{P}$ and experience replay memory $\mathcal{B}$ to capacity $N$

1: **for** $n = 1, ..., N$ **do**
2:      $P \equiv \langle X, D, F \rangle \sim \mathcal{P}$, $\Gamma \leftarrow \emptyset$
3:      build a pseudo tree for $P$ rooted at $x_i$
4:      **DFSDecision**$(\Gamma, x_i)$
5:      **for** $k = 1, ..., K$ **do**
6:          $B_s \leftarrow$ sample random minibatch of data from $\mathcal{B}$
7:          train the model $Q_\theta$ to minimize Eq. (14)
8:          $\bar{\theta} \leftarrow \rho\bar{\theta} + (1 - \rho)\theta$
9:      **end for**
10: **end for**
   Function **DFSDecision**$(\Gamma, x_i)$
11: $s \leftarrow \langle P, \Gamma, x_i \rangle$
12: $d_i = \begin{cases} \text{random value } d_i' \in D_i, \text{w.p. } \epsilon \\ \arg\max_{d_i' \in D_i} Q_\theta\left(s, \langle x_i, d_i' \rangle\right), \text{otherwise} \end{cases}$
13: compute $r$ by Eq. (11), $\Gamma \leftarrow \Gamma \cup \{\langle x_i, d_i \rangle\}$
14: $s' \leftarrow \{\langle P, \Gamma, x_c \rangle | \forall x_c \in C(x_i)\}$, $\mathcal{B} \leftarrow \mathcal{B} \cup \{(s, \langle x_i, d_i \rangle, r, s')\}$
15: **for** $x_c \in C(x_i)$ **do**
16:      **DFSDecision**$(\Gamma, x_c)$
17: **end for**

---

### 4.3 The DRL Training Algorithm

To learn the parameter of our DRL model, we use Deep Q-network (DQN) as illustrated in Algo. 1. In the algorithm, the experiences are generated and stored into a capacitated experience replay memory $\mathcal{B}$ (line 1-4, 11-17). Afterward, the DRL model is trained with the data from the experience memory (line 5-9). Specifically, DQN first samples a WCSP instance from the problem distribution and builds a pseudo tree to order variables (line 2-3). Then, it progressively chooses values for all variables to generate experiences (line 4, 11-17). For each variable $x_i$ given the partial assignment $\Gamma$, DQN calls Function **DFSDecision** to assign $x_i$ based on the $\epsilon$-greedy policy, compute the reward and transition state (line 11-13). The experience, including the initial state, action, reward, and transition state, is then added to the experience memory (line 14). Subsequently, DQN calls Function **DFSDecision** to select values for the children of $x_i$ (line 15-17). After all variables have been assigned, DQN uniformly samples a batch of data from the experience memory, and train the DRL model with the temporal difference error (line 6-7). That is,

$$\mathcal{L}(\theta) = \frac{1}{|B_s|} \sum_{\langle s, \langle x_i, d_i \rangle, r, s' \rangle \in B_s} \left(Q_\theta(s, \langle x_i, d_i \rangle) - y(s')\right)^2 \qquad (14)$$

where $s' = \{s_c | \forall x_c \in C(x_i)\}$ and $y(s') = r + \gamma \sum_{x_c \in C(x_i)} \max_{d_c \in D_c} Q_{\bar{\theta}}\left(s_c, \langle x_c, d_c \rangle\right)$. Lastly, the parameter of the target Q-network is updated towards to the main Q-network based on the interpolation factor $\rho$ (line 8).

---

**Algorithm 2** Beam search with our DRL model

---

**Require**: a WCSP instance $P$, the DRL model $Q_\theta$, the paramters $k_{bw}$ and $k_{ext}$
1: $PT \leftarrow$ build a pseudo tree for $P$
2: $B \leftarrow \{\emptyset | \forall i = \{1, ..., k_{bw}\}\}$
3: **for** $l = 1, ..., l^*_{PT}$ **do**
4:     **for** $x_i \in X, \text{level}_{PT}(x_i) = l$ **do**
5:         $B' \leftarrow \emptyset, E' \leftarrow \emptyset$
6:         **for** $j = 1, ..., k_{bw}$ **do**
7:             $\Gamma \leftarrow B_{[j]}, E \leftarrow \{q(\Gamma, \langle x_i, d_i \rangle) \text{ computed by Eq. (15)} | \forall d_i \in D_i\}$
8:             $Ind \leftarrow$ the index of the $k_{ext}$ smallest elements in $E$
9:             $D'_i \leftarrow \{d_i | \forall d_i \in D_i, i \in Ind\}$
10:            $B' \leftarrow B' \cup \{\langle \Gamma \cup \{\langle x_i, d_i \rangle\} | \forall d_i \in D'_i\}$
11:            $E' \leftarrow E' \cup \{cost(\Gamma) + q(\Gamma, \langle x_i, d_i \rangle) | \forall d_i \in D'_i\}$
12:         **end for**
13:         $Ind' \leftarrow$ the index of the $k_{bw}$ smallest elements in $E'$
14:         $B \leftarrow \left\{ B'_{[i]} | \forall i \in Ind' \right\}$
15:     **end for**
16: **end for**
17: **return** $\arg\min_{\Gamma \in B} cost(\Gamma)$

---

### 4.4 Heuristics for Beam Search and Large Neighborhood Search

Besides constructing solutions greedily from scratch, in this section, we show that our DRL model can also be embedded into various meta-heuristics. In particular, we consider two well-known meta-heuristics, i.e., Beam Search (BS) and Large Neighborhood Search (LNS) and show how to apply the predicted Q-values to these heuristics.

Since the Q-values in our model approximate the optimal cost of the sub-problem rooted at a variable, one can use the predicted Q-values as a proxy to evaluate the quality of each assignment. Thus, our DRL model can be embedded into BS to guide the expansion of partial assignments and determine which extended partial assignments are eligible for further branching. Algo. 2 gives the pseudo-code of BS with our DRL model, where the beam $B$ is introduced to track the most promising partial assignments and initialized as $k_{wb}$ (called the beam width) empty sets at the beginning of the algorithm (line 2). BS first orders variables by using a pseudo tree $PT$ (line 1) where $l^*_{PT}$ is the maximal level of variables in $PT$ (line 3) and $\text{level}_{PT}(x_i)$ is the level of $x_i$ in $PT$ (line 4). Then, it explores the search tree according to the level of each variable in $PT$ (lines 3-16). For $x_i$, a variable at the $l$-th level of $PT$, BS works by extending each partial assignment $\Gamma \in B$ at most $k_{ext}$ possible ways according to their quality (line 6-9). The quality of each value $d_i \in D_i$ under $\Gamma$ is evaluated by:

$$q(\Gamma, \langle x_i, d_i \rangle) = -Q_\theta(\langle P, \Gamma, x_i \rangle, \langle x_i, d_i \rangle) \tag{15}$$

The extended partial assignments and their evaluation values are then stored in a new beam $B'$ and $E'$, respectively (line 10-11). After all partial

Table 1: The intermediate results of $x_2$ when executing BS

| $\Gamma$ | $q(\Gamma, \langle x_2, 0\rangle)$ | $q(\Gamma, \langle x_2, 1\rangle)$ | $d_2$ | $cost(\Gamma)$ | $B'$ | $E'$ |
|---|---|---|---|---|---|---|
| $\{\langle x_1, 0\rangle\}$ | 9 | 11 | 0 | 0 | $\{\{\langle x_1, 0\rangle, \langle x_2, 0\rangle\},$ | $\{9, 10\}$ |
| $\{\langle x_1, 1\rangle\}$ | 10 | 11 | 0 | 0 | $\{\langle x_1, 1\rangle, \langle x_2, 0\rangle\}\}$ | |

assignments in $B$ have been processed, the algorithm updates $B$ with the top $k_{bw}$ optimal partial assignments in $B'$ (line 13-14).

Next, we will take the variable $x_2$ in Fig. 1(b) as an example to show how our DRL model is embeded into BS with $k_{wb} = 2$ and $k_{ext} = 1$. Assuming that $B = \{\{\langle x_1, 0\rangle\}, \{\langle x_1, 1\rangle\}\}$ and $D_2 = \{0, 1\}$, $x_2$ expands each partial assignment $\Gamma \in B$ with its corresponding best value $d_2 \in D_2$ and the intermediate results can be found in Table 1. For $\Gamma = B_{[1]} = \{\langle x_1, 0\rangle\}$, it computes the quality evaluation for each value in $D_2$ given $\Gamma$ by Eq. (15). Thus, we have $q(\Gamma, \langle x_2, 0\rangle) = 9$ and $q(\Gamma, \langle x_2, 1\rangle) = 11$ as shown in the first row of Table 1. Since $q(\Gamma, \langle x_2, 0\rangle) < q(\Gamma, \langle x_2, 1\rangle)$ and $k_{ext} = 1$, $x_2$ assigns itself with 0 and extends $\Gamma$ with its assignment, i.e., the extended assignment is $\Gamma' = \Gamma \cup \{\langle x_2, 0\rangle\}$. Then, it adds $\Gamma'$ and the corresponding evaluation value (i.e., $cost(\Gamma) + q(\Gamma, \langle x_2, 0\rangle) = 9$) to $B'$ and $E'$, respectively. Thus, we have $B' = \{\langle x_1, 0\rangle, \langle x_2, 1\rangle\}$ and $E' = \{9\}$.

For $\Gamma = B_{[2]} = \{\langle x_1, 1\rangle\}$, $x_2$ performs a similar procedure as above. As shown in the second row of Table 1, it selects 0 for itself (since $q(\Gamma, \langle x_2, 0\rangle) < q(\Gamma, \langle x_2, 1\rangle)$) and expands $\Gamma$ to $\Gamma' = \Gamma \cup \{\langle x_2, 0\rangle\}$. Afterward, $\Gamma'$ and the corresponding evaluation value are added to $B'$ and $E'$, respectively. Now, $B' = \{\{\langle x_1, 0\rangle, \langle x_2, 1\rangle\}, \{\langle x_1, 1\rangle, \langle x_2, 0\rangle\}\}$ and $E' = \{9, 10\}$. Since $|B'| = k_{wb} = 2$, no partial assignment in $B'$ will be abandoned and thus the beam for $x_3$ is $B = B' = \{\{\langle x_1, 0\rangle, \langle x_2, 1\rangle\}, \{\langle x_1, 1\rangle, \langle x_2, 0\rangle\}\}$.

---

**Algorithm 3** Large neighborhood search with our DRL model

---

**Require**: a WCSP instance $P$, the DRL model $Q_\theta$, the proportion of destroyed variables $p$, number of iterations $T$
1: $sol \leftarrow$ a random solution
2: **for** t=1,...,T **do**
  **Destroy phase**
3:　　$X_{des} \leftarrow$ uniformly select $|X|p$ variables from $X$
  **Repair phase**
4:　　$\Gamma \leftarrow \{\langle x_i, sol_{[x_i]}\rangle | \forall x_i \notin X_{des}\}$
5:　　**for all** connected sub-problem $X_{cr} \subseteq X_{des}$ **do**
6:　　　　$PT \leftarrow$ build a pseudo tree for $X_{cr}$
7:　　　　**for** $l = 1, ..., l^*_{PT}$ **do**
8:　　　　　　**for** $x_i \in X, \text{level}_{PT}(x_i) = l$ **do**
9:　　　　　　　　compute $d^*_i$ by Eq. (16), $\Gamma \leftarrow \Gamma \cup \{\langle x_i, d^*_i\rangle\}$
10:　　　　　**end for**
11:　　　**end for**
12:　　**end for**
13:　　update $sol$ with $\Gamma$
14: **end for**

---

Our model can also be embedded into LNS as a repair policy which solves each connected sub-problem by assigning variables according to the predicted Q-values. The pseudo-code of LNS with our DRL model can be found in Algo. 3. The algorithm iteratively improves an initial solution by repeatedly executing destroy and repair phases (line 1-14). During the destroy phase, LNS selects the Large Neighborhood (LN) variables whose current assignment will be discarded (line 3). Then, in the repair phase, LNS finds a new assignment for LN variables given the assignment of undestroyed variables from the previous iteration (line 4-12). The assignment for a LN variable $x_i$ is computed by:

$$d_i^* = \arg \max_{d_i \in D_i} Q_\theta(\langle P, \Gamma, x_i \rangle, \langle x_i, d_i \rangle) \tag{16}$$

## 5 Experimental Evaluations

In this section, we conduct extensive empirical studies. We first present the details of the experiments and training stage. Then, we compare the heuristics obtained from our DRL model against those derived from the SL model to illustrate the superiority of our DRL framework. Finally, we compare our DRL-based algorithms with their counterparts with traditional tabular-based heuristics (e.g., mini-bucket elimination, MBE) and state-of-the-art incomplete WCSP algorithms on various standard benchmarks, with the aims of showing the capability of our DRL model to boost WCSP algorithms.

### 5.1 Benchmarks

In this experiment, the algorithms are benchmarked on four types of problems, including random WCSPs, weighted graph coloring problems, scale-free networks and random meeting scheduling problems.

- *Random WCSPs* are a general form for WCSPs, where a set of variables are randomly constrained with one another. In the experiment, the performance of the algorithms is evaluated by varying the number of variables and graph density (see the detailed configurations in Subsection 5.3). The constraint costs are uniformly selected from [0, 100].
- *Weighted graph coloring problems* are the problems in which every vertex should be colored, and two adjacent vertexes should have different colors. In the experiments, we consider weighted graph coloring problems with 3 available colors for each variable, the graph density of 0.05, and varying the number of variables from 100 to 300. The cost of each violation is chosen uniformly from 0 to 100.
- *Scale-free networks* [2] are networks whose degree distribution follows a power law. The experiment uses Barabási-Albert (BA) model to generate the constraint graph topology. In the beginning, we set the number of initial variables to 10 and connect them randomly. In an iteration of the

BA process, we add a new variable and connect it with 3 variables (for spare problems) or 10 variables (for dense problems) with a probability proportional to the current number of links. We set the variable number to 120, and the other parameters are the same as that of random WCSPs.

– *Random meeting scheduling problems* [53] are problems of persons scheduling a set of meetings. We randomly select a travel time for each pair of meetings. When the difference between the time-slots of two meetings with overlapping persons is less than the travel time, the persons in both meetings are considered overbooked, and the cost is defined as the number of the overbooked persons. In the experiments, we set the number of participants to 90, the number of meetings to 20, and the available time-slots to 20. Each person randomly participates in two meetings, and travel times are uniformly selected from 6 to 10.

5.2 Baselines

In the expriment, we focus on three meta-heuristics for combinational optimization problems, including Greedy Search (GS), Beam Search (BS) and Large Neighborhood Search (LNS). To demonstrate the superiority of our DRL framework, we benchmark the heuristics derived from our DRL model, those obtained from the SL model [14] and traditional tabular-based heuristics (e.g., mini-bucket elimination, MBE [13]). We set the memory budget of MBE to $d^6$, and $k_{wb} = 4$ and $k_{ext} = 2$ for BS. Besides, we consider three types of baselines: local search, belief propagation and large neighborhood search. We use a stochastic algorithm (SA) [52] with $p = 0.8$ and GDBA [32] with $\langle M, NM, T \rangle$ as two representative local search methods, Damped Max-Sum (DMS) [9] with the damping factor to 0.9 as a representative belief propagation method, and T-LNS [23] with the destroy probability to 0.5 as a representative large neighborhood search method.

For a fair comparison, the configuration of the GAT model we use is the same as that in [14]. Specifically, our model consists of 4 GAT layers (i.e. $T = 4$), where the first 3 layers have 8 output channels and 8 attention heads, and the latter layer has 16 outputs channel and 4 attention heads. Besides, we use ELU [8] as the activation function for each GAT layer. The training set and validation set in the expriment are derived from a WCSP random distribution with $|X| \in [40, 60]$, $d \in [3, 15]$, $p_1 \in [0.1, 0.4]$ where $p_1$ is the graph density of the WCSP instance. The GAT model is implemented with the PyTorch geometry framework [17] and trained with the Adam optimizer [25] with a learning rate of 0.0001 and a weight decay rate of 0.00001. For the hyper-parameters of the deep Q-network algorithm, we set the discount factor $\gamma$ to 0.99, the batch size to 64, and the number of training steps to 15000. Finally, we average the experimental results over 50 independently generated problems and evaluate the solution quality of all algorithms after running for the same wall-clock time. All experiments are carried out on an Intel i7-7820X workstation with GeForce RTX 3090 GPUs.
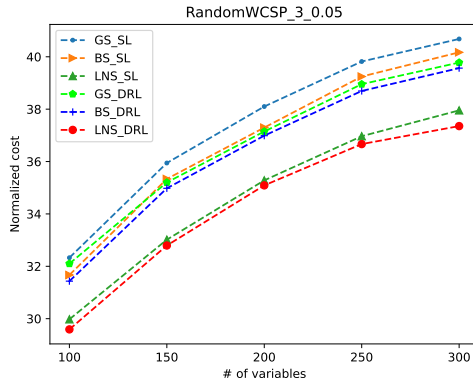
Fig. 5: Solution quality comparisons of the heuristics derived from our DRL model and the SL model on random WCSPs

## 5.3 Performance Comparisons

Fig. 5 shows the comparison of the DRL-based and SL-based algorithms on random WCSPs with the graph density of 0.05, the domain size of 3, and the number of variables varying from 100 to 300. It can be seen from the figure that, the DRL-based algorithms have great advantages over the SL-based algorithms on all the problems. Specifically, GS_DRL, BS_DRL and LNS_DRL outperform GS_SL, BS_SL, LNS_SL by about 0.71∼5.51%, 0.94∼3.11% and 0.81∼3.63%, respectively. This phenomenon shows that the heuristics derived from our DRL model is better than those obtained from the SL model, and the gap become wider when solving large-scale problems. That is because training on large-scale problems helps the model to mine more sophisticated patterns of WCSPs, and the scale of training problems for DRL is much larger than that for SL. Specifically, SL uses optimal labeled data to train the model. However, obtaining optimal labeled data is quite challenging, making SL only trainable on small-scale problems. In the experiment, the training data for SL comes from a random WCSP distribution with $|X| \in [15, 30]$, $d \in [3, 15]$ and $p_1 \in [0.1, 0.4]$. The DRL training process can be seen as an alternation between finding "high-quality labeled data" and performing SL on the collected data. Thus, DRL can train models on large-scale problems. The training problems for DRL comes from a random WCSP distribution with $|X| \in [40, 60]$, $d \in [3, 15]$ and $p_1 \in [0.1, 0.4]$.

Fig. 6 shows solution quality results for all the baselines on random WCSPs and Weighted Graph Coloring (WGC) problems. We set the graph density to 0.05, the domain size to 10, and the number of variables varied from 100 to 300 for random WCSPs. It can be seen from Fig. 6, the performance of the DRL-based algorithms on small-scale problems (e.g., $|X|$=100) is similar to their counterparts with MBE. The advantage of our DRL-based algorithms comes out as the number of variables increases. Specifically, our DRL-based
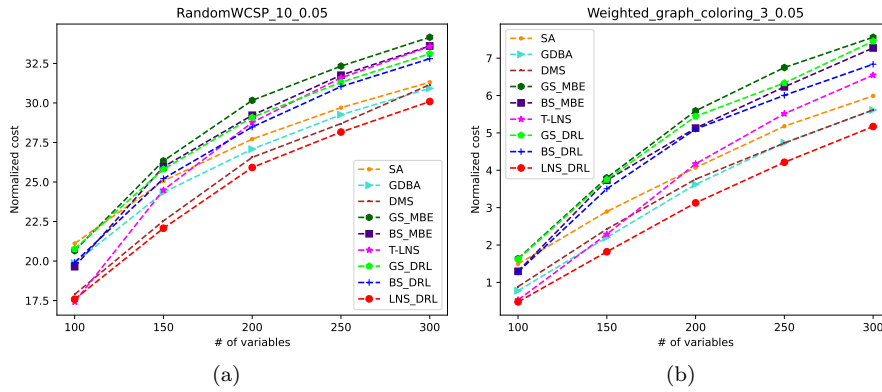
Fig. 6: Solution quality comparisons: (a) random WCSPs; (b) weighted graph coloring problems

algorithms improve their counterparts with MBE by about 2.54∼10.85% on random WCSPs with 300 variables and 1.17∼25.4% on WGC problems with 300 variables. That is because the induced width of the built pseudo tree increases with the number of variables. Concretely, the average induced width is 43 for the problems with 100 variables and 223 for the problems with 300 variables. Therefore, in the face of complex problems, MBE needs to approximate more dimensions in the variable elimination process and cannot provide tight enough bounds. T-LNS attempts to optimize by optimally solving tree-structure relaxations of sub-problems induced by destroyed variables in each round. It also performs poorly when solving large-scale problems (e.g., problems with more than 200 variables) due to ignoring many constraints. In contrast, our LNS_DRL solves the induced problem without relaxation, which is a significant improvement over state-of-the-art methods such as DMS. Precisely, LNS_DRL outperforms DMS by about 1.69∼3.42% on random WCSPs and 7.93∼25.03% on WGC problems.

Fig. 7 shows the comparison results of random WCSPs with the number of variables of 70, where the domain size is set to 20, the graph density is set as 0.1 for the sparse configuration, and 0.6 for the dense configuration. It can be seen from the figure that, the DRL-based algorithms have obvious advantages over their counterparts with MBE on these problems. Precisely, GS_DRL and BS_DRL excel GS_MBE and BS_MBE by about 2.32% and 2.03% on the sparse problems, respectively, and 1.39% and 1.58% on the dense problems, respectively. LNS_DRL outperforms T-LNS by about 2.64% on the sparse problems and 4.6% on the dense problems. These phenomena show that our heuristics can effectively improve the performance of the original algorithms. In addition, it can be seen that the local search algorithms (e.g., SA and GDBA) quickly converge to poor local optima, while DMS can find a better solution in 25 seconds for the sparse problems and 200 seconds for the dense problems. In contrast, our LNS_DRL improves more steadily, outperforming
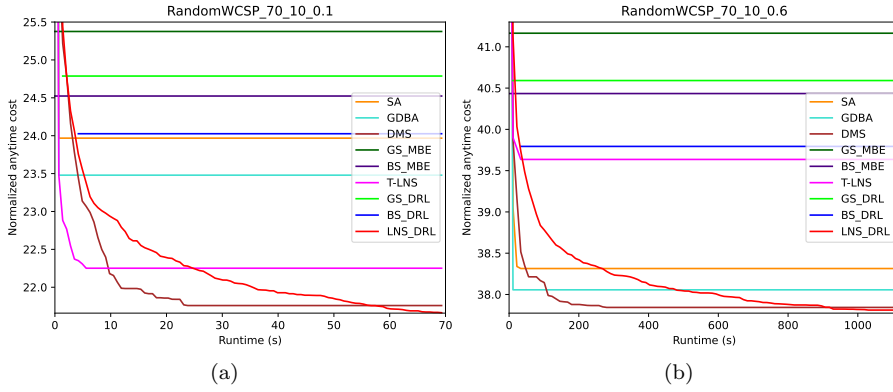
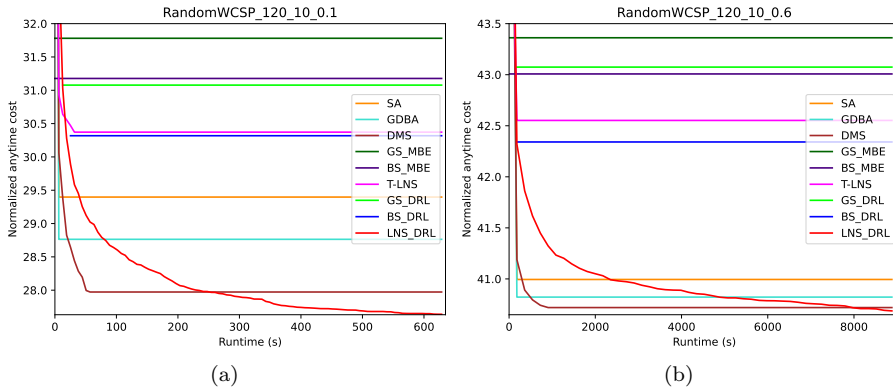Fig. 7: Solution qualities on random WCSPs (70 variables)



Fig. 8: Solution qualities on random WCSPs (120 variables)

all the baselines after 60 seconds on the sparse problems and 900 seconds on the dense problems. Specifically, LNS_DRL outperforms state-of-the-art algorithms (e.g., DMS) by about 0.43% on the sparse problems and 0.11% on the dense problems.

Fig. 8 presents the comparisons on the random WCSPs with 120 variables, where the domain size is set to 10, the graph density is set as 0.1 for the sparse configuration, and 0.6 for the dense configuration. Similarly, the DRL-based algorithms have great superiorities over their counterparts with MBE. Concretely, GS_DRL and BS_DRL surpass GS_MBE and BS_MBE by about 2.21% and 2.76% on the sparse problems, respectively, and 0.66% and 1.55% on the dense problems, respectively. LNS_DRL outperforms T-LNS by about 9.01% on the sparse problems and 4.38% on the dense problems. Different from the experimental results shown in Fig. 7, we can see from Fig. 8 that our BS_DRL outperforms T-LNS by about 1.8% on the sparse problems, and the advantage on the dense problems extends to 5.12%. That is because the
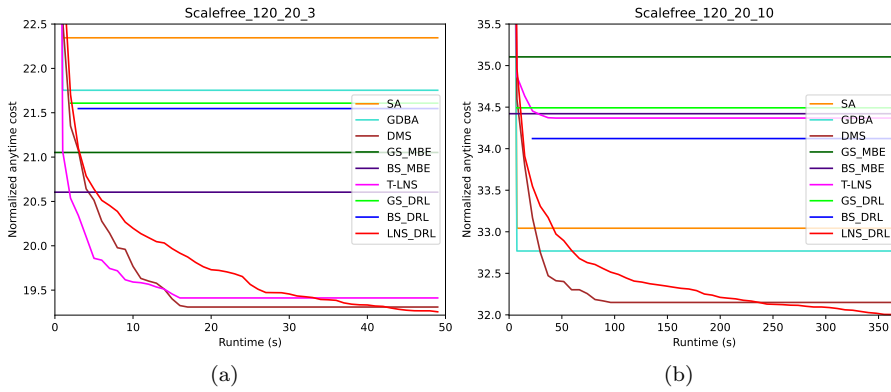
Fig. 9: Solution qualities on scale-free problems

proportion of constraints ignored by T-LNS when solving the problems with 120 variables is much larger than when solving the problems with 70 variables.

Fig. 9 shows the comparisons of normalized anytime costs on the sparse and dense scale-free networks, respectively. In the experiment, the average induced width is 32 for the spare problems and 77 for the dense problems. Unlike the experiment results on random WCSPs, GS_DRL and BS_DRL perform worse than their counterparts with MBE on the spare scale-free networks. That is because the induced width of spare problems is relatively small compared to that of random WCSPs, so MBE can provide more efficient bounds for greedy search and beam search. Also, it can be seen from Fig. 9(a) that T-LNS exhibits great superiorities over local search algorithms (i.e., SA and GDBA) and is slightly lower than DMS. However, due to the effectiveness of our heuristics, LNS_DRL still outperforms state-of-arts methods (e.g., DMS) by about 2.48% on the spare problems. On the dense problems, the advantage of our heuristics comes out in greedy search and beam search. Specifically, GS_DRL and BS_DRL outperform GS_MBE and BS_MBE by about 1.75% and 0.87%, respectively. Although T-LNS exhibits its huge advantage on the spare problems, its performance on the dense problems is even 0.72% worse than BS_DRL and 6.88% worse than LNS_DRL. These results demonstrate the necessity of our heuristics for LNS to solve structured problems, especially when facing large-scale problems.

Fig. 10 presents the comparisons of normalized anytime costs on the random meeting scheduling problems. The average induced width for this experiment is 14, which is much smaller than that of random WCSPs. Thus, it can be seen that the performance of GS_DRL and BS_DRL is also inferior to their counterparts with MBE, similar to the experimental results on the sparse scale-free networks. Furthermore, T-LNS achieves better performance than all the algorithms expected LNS_DRL. That is because the variables in the problem are under-constrained, and T-LNS only needs to drop a few edges to obtain a tree-structured problem. However, our LNS_DRL still exhibits great superi-
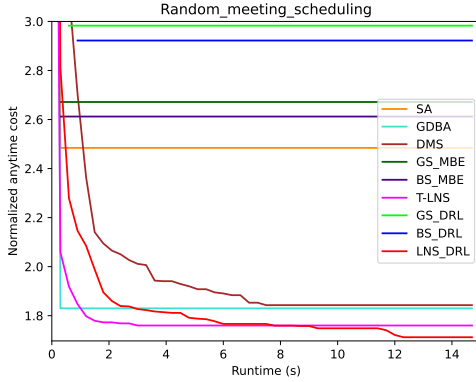
Fig. 10: Solution qualities on random meeting scheduling problems

orities over other competitors. Specifically, LNS_DRL is superior to DMS by about 7.87%, GDBA by abour 7.05% and T-LNS by about 3.53%. The results demonstrate that our heuristics can improve the performance of LNS when solving real-world problems.

## 6 Conclusion

In this paper, we propose a DRL framework that enables the model to be trained on large-scale problems. As a result, the learned model could capture more sophisticated patterns of the problems to generate effective solution construction heuristics for WCSPs. In the framework, we propose to effectively decompose the problem by using a pseudo tree, and formulate the solution construction process as an MDP with multiple independent transition states. Through a modified Bellman Equation, we use GATs parameterized deep Q-value network to learn the optimal Q-values, and the solution construction heuristics are extracted from the learned Q-value network. In addition to constructing greedy solutions, we embed our heuristics into BS by evaluating the quality of partial assignments, and LNS by finding new assignments for destroyed variables via constructing a solution greedily. The extensive empirical evaluations confirm the effectiveness of our proposed DRL framework.

In the future, we plan to improve our DRL framework in the following aspects: at first, we will exploit curriculum learning [30] to make the training processes more stable; then, we will devote to extending our DRL model to handle the problems with hard and higher-arity constraints; finally, we will explore to embed our model into complete WCSP algorithms, such as extracting the domain ordering heuristics from the model for branch-and-bound algorithms.

# References

1. Allouche, D., Givry, S.d., Katsirelos, G., Schiex, T., Zytnicki, M.: Anytime hybrid best-first search with tree decomposition for weighted CSP. In: CP, pp. 12–29. Springer (2015)
2. Barabási, A.L., Albert, R.: Emergence of scaling in random networks. Science **286**(5439), 509–512 (1999)
3. Bellman, R.: A markovian decision process. Journal of mathematics and mechanics pp. 679–684 (1957)
4. Bengio, Y., Lodi, A., Prouvost, A.: Machine learning for combinatorial optimization: a methodological tour d'horizon. European Journal of Operational Research **290**(2), 405–421 (2021)
5. Chalumeau, F., Coulon, I., Cappart, Q., Rousseau, L.M.: Seapearl: A constraint programming solver guided by reinforcement learning. In: CPAIOR, pp. 392–409. Springer (2021)
6. Chen, Z., Zhang, W., Deng, Y., Chen, D., Li, Q.: RMB-DPOP: Refining MB-DPOP by reducing redundant inference. In: AAMAS, pp. 249–257 (2020)
7. Cicirello, V.A.: On the design of an adaptive simulated annealing algorithm. In: CP Workshop on Autonomous Search (2007)
8. Clevert, D.A., Unterthiner, T., Hochreiter, S.: Fast and accurate deep network learning by exponential linear units (ELUS). In: ICLR (2016)
9. Cohen, L., Galiki, R., Zivan, R.: Governing convergence of Max-sum on DCOPs through damping and splitting. Artificial Intelligence **279**, 103212 (2020)
10. De Givry, S., Heras, F., Zytnicki, M., Larrosa, J.: Existential arc consistency: Getting closer to full arc consistency in weighted CSPs. In: IJCAI, vol. 5, pp. 84–89 (2005)
11. Dechter, R.: Bucket elimination: A unifying framework for reasoning. Artificial Intelligence **113**(1-2), 41–85 (1999)
12. Dechter, R., Cohen, D., et al.: Constraint processing. Morgan Kaufmann (2003)
13. Dechter, R., Rish, I.: Mini-buckets: A general scheme for bounded inference. Journal of the ACM (JACM) **50**(2), 107–153 (2003)
14. Deng, Y., Kong, S., An, B.: Pretrained cost model for distributed constraint optimization problems. In: AAAI (2022)
15. Deng, Y., Yu, R., Wang, X., An, B.: Neural regret-matching for distributed constraint optimization problems. In: IJCAI (2021)
16. Farinelli, A., Rogers, A., Petcu, A., Jennings, N.R.: Decentralised coordination of low-power embedded devices using the max-sum algorithm. In: AAMAS, pp. 639–646 (2008)
17. Fey, M., Lenssen, J.E.: Fast graph representation learning with PyTorch geometric. In: ICLR Workshop on Representation Learning on Graphs and Manifolds (2019)
18. Freuder, E.C., Quinn, M.J.: Taking advantage of stable sets of variables in constraint satisfaction problems. In: IJCAI, vol. 85, pp. 1076–1078 (1985)
19. Galassi, A., Lombardi, M., Mello, P., Milano, M.: Model agnostic solution of CSPs via deep learning: A preliminary study. In: CPAIOR, pp. 254–262. Springer (2018)
20. Gaudreault, J., Frayret, J.M., Pesant, G.: Distributed search for supply chain coordination. Computers in Industry **60**(6), 441–451 (2009)
21. Givry, S.d., Lee, J.H., Leung, K.L., Shum, Y.W.: Solving a judge assignment problem using conjunctions of global cost functions. In: CP, pp. 797–812. Springer (2014)
22. Hart, S., Mas-Colell, A.: A simple adaptive procedure leading to correlated equilibrium. Econometrica **68**(5), 1127–1150 (2000)
23. Hoang, K.D., Fioretto, F., Yeoh, W., Pontelli, E., Zivan, R.: A large neighboring search schema for multi-agent optimization. In: CP, pp. 688–706. Springer (2018)
24. Jégou, P.: Decomposition of domains based on the micro-structure of finite constraint-satisfaction problems. In: AAAI, vol. 93, pp. 731–736 (1993)
25. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: ICLR (2015)
26. Kurin, V., Godil, S., Whiteson, S., Catanzaro, B.: Can Q-learning with graph networks learn a generalizable branching heuristic for a SAT solver? In: NeurIPS, vol. 33, pp. 9608–9621 (2020)
27. Larrosa, J., Schiex, T.: In the quest of the best form of local consistency for weighted CSP. In: IJCAI, vol. 3, pp. 239–244 (2003)

28. Lawler, E.L., Wood, D.E.: Branch-and-bound methods: A survey. Operations Research **14**(4), 699–719 (1966)
29. Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G., Graves, A., Riedmiller, M., Fidjeland, A.K., Ostrovski, G., et al.: Human-level control through deep reinforcement learning. Nature **518**(7540), 529–533 (2015)
30. Narvekar, S., Peng, B., Leonetti, M., Sinapov, J., Taylor, M.E., Stone, P.: Curriculum learning for reinforcement learning domains: A framework and survey. Journal of Machine Learning Research **21**, 181:1–181:50 (2020)
31. Nguyen, D.T., Yeoh, W., Lau, H.C., Zivan, R.: Distributed Gibbs: A linear-space sampling-based DCOP algorithm. Journal of Artificial Intelligence Research **64**, 705–748 (2019)
32. Okamoto, S., Zivan, R., Nahon, A., et al.: Distributed breakout: Beyond satisfaction. In: IJCAI, pp. 447–453 (2016)
33. Ottens, B., Dimitrakakis, C., Faltings, B.: DUCT: An upper confidence bound approach to distributed constraint optimization problems. ACM Transactions on Intelligent Systems and Technology **8**(5), 1–27 (2012)
34. Petcu, A., Faltings, B.: DPOP: A scalable method for multiagent constraint optimization. In: IJCAI, pp. 266–271 (2005)
35. Pisinger, D., Ropke, S.: Handbook of Metaheuristics. Springer (2010)
36. Razeghi, Y., Kask, K., Lu, Y., Baldi, P., Agarwal, S., Dechter, R.: Deep bucket elimination. In: IJCAI, pp. 4235–4242 (2021)
37. Scarselli, F., Gori, M., Tsoi, A.C., Hagenbuchner, M., Monfardini, G.: The graph neural network model. IEEE Transactions on Neural Networks **20**(1), 61–80 (2009)
38. Schiex, T., Fargier, H., Verfaillie, G., et al.: Valued constraint satisfaction problems: Hard and easy problems. In: IJCAI, vol. 95, pp. 631–639 (1995)
39. Selsam, D., Lamm, M., Benedikt, B., Liang, P., de Moura, L., Dill, D.L., et al.: Learning a SAT solver from single-bit supervision. In: ICLR (2019)
40. Shapiro, S.C.: Encyclopedia of artificial intelligence, second edition. Wiley-Interscience (1992)
41. Shaw, P.: Using constraint programming and local search methods to solve vehicle routing problems. In: CP, pp. 417–431. Springer (1998)
42. Song, W., Cao, Z., Zhang, J., Xu, C., Lim, A.: Learning variable ordering heuristics for solving constraint satisfaction problems. Engineering Applications of Artificial Intelligence **109**, 104603 (2022)
43. Strokach, A., Becerra, D., Corbi-Verge, C., Perez-Riba, A., Kim, P.M.: Fast and flexible protein design using deep graph neural networks. Cell systems **11**(4), 402–411 (2020)
44. Sutton, R.S., Barto, A.G.: Reinforcement learning: An introduction. MIT press (2018)
45. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. In: NeurIPS, pp. 5998–6008 (2017)
46. Veličković, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., Bengio, Y.: Graph attention networks. In: ICLR (2018)
47. Vinyals, M., Shieh, E., Cerquides, J., Rodriguez-Aguilar, J.A., Yin, Z., Tambe, M., Bowring, E.: Quality guarantees for region optimal DCOP algorithms. In: AAMAS, pp. 133–140 (2011)
48. Vucinic, J., Simoncini, D., Ruffini, M., Barbe, S., Schiex, T.: Positive multistate protein design. Bioinformatics **36**(1), 122–130 (2020)
49. Williams, R.J.: Simple statistical gradient-following algorithms for connectionist reinforcement learning. Machine Learning **8**(3), 229–256 (1992)
50. Xu, H., Koenig, S., Kumar, T.S.: Towards effective deep learning for constraint satisfaction problems. In: CP, pp. 588–597. Springer (2018)
51. Yolcu, E., Póczos, B.: Learning local search heuristics for boolean satisfiability. In: NeurIPS, pp. 7990–8001 (2019)
52. Zhang, W., Wang, G., Xing, Z., Wittenburg, L.: Distributed stochastic search and distributed breakout: properties, comparison and applications to constraint optimization problems in sensor networks. Artificial Intelligence **161**(1-2), 55–87 (2005)
53. Zivan, R., Parash, T., Cohen, L., Peled, H., Okamoto, S.: Balancing exploration and exploitation in incomplete min/max-sum inference for distributed constraint optimization. Autonomous Agents and Multi-Agent Systems **31**(5), 1165–1207 (2017)