

Efficient and Feasible Robotic Assembly Sequence Planning via Graph Representation Learning

Matan Atad^{*2}, Jianxiang Feng^{*1,2}, Ismael Rodríguez^{1,2}, Maximilian Durner^{1,2} and Rudolph Triebel^{1,2}

Abstract—Automatic Robotic Assembly Sequence Planning (RASP) can significantly improve productivity and resilience in modern manufacturing along with the growing need for greater product customization. One of the main challenges in realizing such automation resides in efficiently finding solutions from a growing number of potential sequences for increasingly complex assemblies. Besides, costly feasibility checks are always required for the robotic system. To address this, we propose a holistic graphical approach including a graph representation called Assembly Graph for product assemblies and a policy architecture, Graph Assembly Processing Network, dubbed GRACE for assembly sequence generation. With GRACE, we are able to extract meaningful information from the graph input and predict assembly sequences in a step-by-step manner. In experiments, we show that our approach can predict feasible assembly sequences across product variants of aluminum profiles based on data collected in simulation of a dual-armed robotic system. We further demonstrate that our method is capable of detecting infeasible assemblies, substantially alleviating the undesirable impacts from false predictions, and hence facilitating real-world deployment soon. Code and training data are available at <https://github.com/DLR-RM/GRACE>.

I. INTRODUCTION

Aiming for high flexibility, manufacturers around the globe are introducing automation for *Robotic Assembly Sequence Planning* (RASP) at a greater pace to respond to rapid changes in market needs for customization of novel product variants [1]. These changes cause often modifications in assembly lines, requiring time-consuming and resource-intensive re-planning, because of the NP-hard combinatorial characteristic [2] of *Assembly Sequence Planning* (ASP), where the number of possible solutions grows with the factorial of the amount of parts involved. Also, to check whether a certain assembly sequence can actually be executed on a specific robotic system is computationally expensive. For example in [3], 11 minutes were required for the assembly motion planning of an IKEA chair. This is more time than the actual execution of the plan, not to mention the case of product variants whose assembly sequence space itself must be explored, easily leading to a search of hours or days instead of minutes.

Several existing works attempt to improve the tedious ASP process by predicting feasible assembly sequences [4], [5] or inferring the underlying rules guiding their creation [6], [7]. Although those works already facilitate the assembly planning, they still lack desirable attributes such as generalization

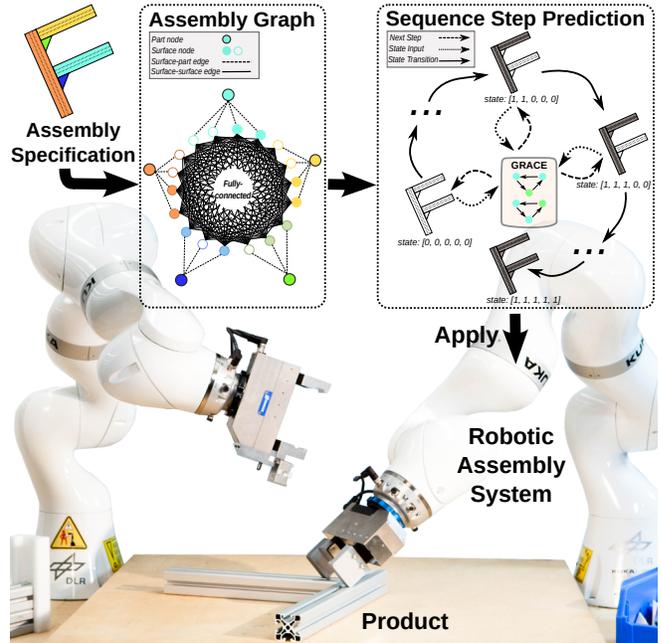


Fig. 1: **Workflow of our proposed graphical ASP method** on a dual-armed robotic system (simulated in our setting): an aluminum assembly specification is first represented as an Assembly Graph and then fed into our policy network GRACE, designed to flexibly and efficiently generate assembly sequences in a *step-by-step* manner that can be executed by the robot. Best viewed in color.

across varying product types and sizes as well as run-time efficiency. In this work, we address these problems with a graphical learning-based approach, that is able to automatically generate sequences for *unknown* assembly variants in an *efficient* way.

In a nutshell of our main idea, inspired by [8], we formulate RASP as a sequential decision-making problem with a *Markov Decision Process* (MDP), in order to break the restriction of combinatorial complexity wrt. the number of parts and thus, boost generalization performance. Hence the sequence is generated step-by-step based on the current assembly state. Meanwhile we exploit the idea of distilling previous knowledge acquired for assembling products to predict the next feasible actions with a designed policy architecture. This architecture is optimized to imitate the demonstration sequences collected in simulation which are interpreted as expert demonstrations.

Specifically, to put the aforementioned ideas into practice,

* Equal Contribution.

¹ Institute of Robotics and Mechatronics, German Aerospace Center (DLR), 82110 Wessling, Germany. <first>.<second>@dlr.de

² Department of Informatics, Technical University of Munich, 85748 Garching, Germany. <first>.<second>@tum.de

we propose to use a graphical representation to faithfully describe the spatial structure of assemblies. Our so-called Assembly Graph is adapted from and more fine-grained than the one in [7] by representing the assembly as a heterogeneous graph whose edges denote geometrical relations between the assembly part surfaces. Based on this, we further develop a policy architecture based on *Graph Neural Networks* (GNNs), called **GR**aph **A**ssembly **pr**oCessing **n**Etworks, for short GRACE, to extract useful information from the Assembly Graph and predict actions determining which parts should be assembled next. Apart from this, false predicted sequences and infeasible assemblies pose a severe problem for efficiency of learning-based assembly robots, e.g., an incorrect sequence might require the robot to perform time-consuming re-planning. Therefore, it would be beneficial to detect these beforehand, e.g., being introspective against false predictions [9], hence we further develop and analyze various schemes to enhance the performance of feasibility prediction.

It is worthwhile to note that there are several advantages for the proposed graphical representation and the policy architecture, GRACE: (1) Invariance to number of parts: contrary to previous works such as [7] restricted by a fixed number of parts, ours is free from this limitation, as GRACE is capable of handling varying number of input graph nodes. (2) Memory efficient learning: GRACE employs shared weights across all nodes in the graph, further alleviating the burden of the aforementioned complexity. (3) Generalization: GRACE trained on assemblies of one size is able to generalize to those of smaller sizes (see results in Tab. III). (4) Multiple solutions: GRACE predicts several feasible sequences (in contrast to [8]), allowing greater flexibility and resilience during execution.

We validate the proposed method with comprehensive experiments based on a dataset of assemblies made of different numbers of aluminum parts created in simulation of a dual-armed robotic system. This setting can be mapped to various tasks in the industry [7] as it allows for construction of numerous product variations. Moreover, as shown in Fig. 3, it requires a deeper understanding of several complex relations (e.g., distances between parts, physical part characteristics). The results show that our approach is able to efficiently predict feasible assembly sequences across product variants (with few millisecond to predict the next step (V-A.3)).

To summarize, our contribution is three-fold:

- We introduce Assembly Graph, a heterogeneous graphical representation for the RASP task, which is a more fine-grained and flexible representation than our previous one in [7] by including part surfaces and parts in the same graph.
- We develop a policy architecture GRACE to process the Assembly Graph and predict feasible assembly sequences in a step-by-step manner as well as the feasibility for a given assembly specification.
- We conduct comprehensive experiments in simulation to validate the proposed approach including failure analysis and ablation studies on design choices.

II. RELATED WORK

a) Assembly Sequence Planning: A popular assembly graph representation for ASP is the AND/OR Graph [11], a formalism to encode the space of feasible assembly sequences, which can be created with the Disassembly For Assembly strategy [12]–[15]. However, these approaches are restricted on time to find a solution efficiently due to the feasibility checks. While graph search methods are impractical for larger assemblies because of the combinatorial explosion problem, heuristic intelligent search methods provide another alternative. They reject infeasible sequences and search for feasible ones close to the optimal based on manually designed termination criteria [16], [17], learned [18], [19] or hand-crafted [20] energy functions. More recently, Zhao *et al.* [4] and Watanabe *et al.* [5] applied deep *Reinforcement Learning* (RL) for ASP. Different to us, they do not have a graph representation to take into account relations between parts. Targeting at RASP, Rodriguez *et al.* [6], [7] suggested inferring assembly rules (e.g., a specific part should be assembled before another), which can be transferred from previous identified sub-assemblies to those of larger sizes to prune the search space, thus reducing planning time. Their approach only produces rules, from which the final assembly sequences need to be derived additionally. It also requires further re-training when adapting to other product variants. Enlightened by them, we refine their graph representation to a more fine-grained level and adapt their idea with a learning-based approach, aiming to mitigate these issues. Similar to us, Ma *et al.* [10] used GNN for ASP of LEGO structures. However, they differ from us in two aspects, first they do not consider assembly robots in the loop and second they model assemblies only with a coarser graph representation whose edges only consider connections among parts instead of part surfaces. To clearly show different characteristics among relevant works, we provide a concise comparison in Tab. I.

b) Graph Representation Learning in Task Planning: In this setting, graphs commonly incorporate nodes for manipulated objects [21]–[23], their target positions [8], [24] and the robot gripper [25]. Edges can represent high-level relations between objects [21], [23]. With the graph representation, Zhu *et al.* [23] and Ye *et al.* [25] generated feasible candidate paths by sampling, and trained a network that predicts a sequence of feasible actions in backward and forward search, respectively. Nguyen *et al.* [21] performed sampling to find action sequences that transform the source to target graph and then used optimization to eliminate invalid sequences subject to the environment constraints. Besides, some researchers resorted to RL methods such as [22], [24], and [26], who used GNNs for task planning. Recently, Lin *et al.* [8] utilized *Imitation Learning* (IL) to train two GNNs, one for selecting objects in the scene and another picking a suitable goal state from a set of possible goal positions for long-horizon manipulation tasks. Inspired by them, we train our GNNs for RASP task by leveraging IL for ease and efficiency in training.

TABLE I: Comparison between our proposed method and other relevant works.

	Efficient generalization across assembly sizes?	Robotic constraints involved?	Direct sequences generation?	Fine-grained graph representation?
ASPW-DRL [4]	✗	✗	✓	✗
LEGO-GRAPH [10]	✓	✗	✓	✗
KT-RASP [7]	✗	✓	✗	✗
Our proposed method	✓	✓	✓	✓

III. BACKGROUND

In this section, we briefly recap the concept of GNNs and heterogeneous graphs, which are the base for our method.

a) *Graph Neural Networks*: A GNN operates on an undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with nodes \mathcal{V} and edges \mathcal{E} , where every node $v \in \mathcal{V}$ is assigned with a feature vector $\phi(v)$. It updates node features by exchanging information between neighboring nodes. This is done with multiple Message Passing layers [27]. For each layer l , let $\mathbf{h}_i^0 = \phi(v_i)$ be the input features of node v_i and \mathcal{N}_i its set of neighboring nodes. Then we can define a three-step process to update these features:

- 1) *Gather* feature from neighboring nodes: $\{\mathbf{h}_j^{l-1}\}_{j \in \mathcal{N}_i}$.
- 2) *Aggregate* messages from the neighboring nodes:

$$\mathbf{m}_i^l = g_\omega(\{\mathbf{h}_j^{l-1}\}_{j \in \mathcal{N}_i}).$$
- 3) *Update* features of node v_i : $\mathbf{h}_i^l = f_\phi^l(\mathbf{h}_i^{l-1}, \mathbf{m}_i^l)$.

The function g_ω can be either constant (e.g. sum) or learned during training. The term f_ϕ is a *Neural Network* (NN) parameterized by ϕ . Both, f_ϕ and g_ω , are shared across all nodes in the graph, making GNNs efficient and independent of the number of nodes in the graph.

In our proposed method we apply a Graph Attention Network (GAT) [28], a popular variant of GNNs, that defines g_ω as attention:

$$\mathbf{m}_i^l = \sum_{j \in \mathcal{N}_i} (\alpha_{i,j} \cdot \mathbf{h}_j^{l-1}), \quad (1)$$

$$\mathbf{h}_i^l = \mathbf{W}_1 \cdot \alpha_{i,i} \mathbf{h}_i^{l-1} + \mathbf{W}_1 \cdot \mathbf{m}_i^l, \quad (2)$$

$$\alpha_{i,j} = \frac{\exp(\mathbf{a} \cdot \sigma(\mathbf{W}_2[\mathbf{h}_i^{l-1} \parallel \mathbf{h}_j^{l-1} \parallel e_{i,j}]))}{\sum_{k \in \mathcal{N}_i \cup \{i\}} \exp(\mathbf{a} \cdot \sigma(\mathbf{W}_2[\mathbf{h}_i^{l-1} \parallel \mathbf{h}_k^{l-1} \parallel e_{i,k}]))}, \quad (3)$$

where \mathbf{W}_1 , \mathbf{W}_2 , and \mathbf{a} are learned, σ is a Leaky ReLU activation function, and $[a \parallel b]$ is a concatenation operator between a and b .

b) *Heterogeneous Graph*: $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ generalizes graphs to multiple types of nodes and edges [?]. Each node $v \in \mathcal{V}$ belongs to one particular node type $\psi_n(v)$ and analogously each edge $e \in \mathcal{E}$ to an edge type $\psi_e(e)$. In [29], the authors extend *Graph Attention Network* (GAT)s to a heterogeneous graph setting. This is accomplished by obtaining for each node a different updated feature vector per group of specific neighboring source node and edge types, and aggregating the features to obtain a single result, for instance using a sum. This formulation is essential, as every type of neighboring node may have a different feature dimension.

IV. METHOD

In this section, RASP is formulated as a sequential decision-making problem with a MDP and then we present our graph representation to depict assemblies. Based on this, we elaborate the proposed network GRACE, and demonstrate the assembly sequence generation.

A. Problem Formulation

We describe the sequence prediction task for an assembly with N parts as a MDP [30] with a discrete state space \mathcal{S} and a high-level discrete action space \mathcal{A} .

Starting from state \mathbf{s}_t at time step t , executing action a_t produces a reward r_t and switches to state $\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1} | \mathbf{s}_t, a_t)$ with a transition function p . State $\mathbf{s}_t \in \{0, 1\}^N$ is a binary vector indicating which parts are already placed in their target position by 1 (i.e. assembled) otherwise by 0. Action $a_t \in \{1, \dots, N\}$ represents the next part placement among the unplaced ones. For *feasible* assemblies, there are multiple different sequences leading to the final state, in which all N parts are placed correctly. For *infeasible* assemblies, no sequence exists, due to constraints of different aspects spanning from part geometries to kinematic and dynamics regarding the robotic system. Our objective is to learn a policy network $\pi_\theta(\mathbf{s}_t) = a_t$ parameterized by θ , which is optimized to imitate the assembly demonstrations $\tau_i = \{\mathbf{s}_{i,1}, a_{i,1}^{exp}, \dots, \mathbf{s}_{i,T}, a_{i,T}^{exp}\}$ in a dataset of M sequences $\mathcal{D} = \{\tau_i\}_{i=1}^M$ and generalize across variants of different types and sizes at test time. In practice, our network predicts a set of multiple possible actions e.g. $K_t = \{a_{t,k}\}_{k=1}^{K_t}$ based on a tunable threshold to control the prediction quality.

B. Assembly Graphs

We represent the overall structure of an assembly with a heterogeneous graph. To make this representation agnostic to the rotation and mirroring of the assembly structure, we employ only relative distances instead of absolute positions for the features of edges between surfaces. More formally, given an assembly A (Fig. 2) at state \mathbf{s}_t it is modeled as a graph $\mathcal{G}_t = (\mathcal{V}, \mathcal{E})$ containing two types of nodes: part nodes \mathcal{V}^p and surface nodes \mathcal{V}^s , and two types of edges: $\mathcal{E}^{s \rightarrow s}$, connecting all surface nodes, and $\mathcal{E}^{s \rightarrow p}$, connecting each surface node to its respective part. We detail each component as follows:

1) *Part Nodes*: Responsible for encoding the current state of the assembly. A part node $v_i^p \in \mathcal{V}^p$ is associated with a feature vector $\phi(v_i^p) = [\text{assembled-flag} \in \{0, 1\}, \text{part-type} \in$

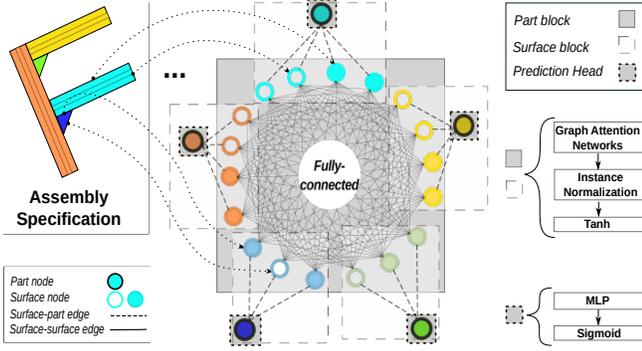


Fig. 2: **Illustration of Assembly Graph and GRACE.** Assembly Graph consists of edges connecting parts and their surfaces and edges among all part surfaces. In GRACE, the Part Block is shared for sub-graphs of part surfaces and the attached part, while Surface Block is for the sub-graph of all part surfaces. To predict scores for parts to be assembled next, we apply a prediction head on each spare part.

\mathbb{N} , $part-id \in \mathbb{R}^d$]. There are three atomic part types: *long profile*, *short profile* and *angle bracket*.

2) *Surface Nodes*: Different to the one in [7], we associate each surface node $v_i^s \in \mathcal{V}^s$ with the features $\phi(v_i^s) = [surface-type \in \mathbb{N}, surface-id \in \mathbb{R}^d]$. There are two surface types (*long* and *short*) for profiles and one (*lateral*) for brackets. Both the *part-id* and *surface-id* fields are encoded with a d -dimensional Sinusoidal Positional Encoding [31].

3) *Surface-to-Surface Edges*: We design a fully-connected graph for all surface nodes \mathcal{V}^s to capture the relation between untouched surfaces, which is more fine-grained than those in [7] with only connects between touched surfaces. These edges are assigned with a feature $\phi(e_i) \in \mathbb{R}$, indicating the *relation* between the two surfaces: $\phi(e_i) = relative\ distance$ (parallel); 1 (belong to the same part); -1 (orthogonal); 0 (same-surface loop).

4) *Surface-to-Part Edges*: These connect each surface and part node pair $(v_i^s, v_j^p) \in \mathcal{V}^s \times \mathcal{V}^p$, where surface v_i^s belongs to the part v_j^p . This type of edges is not associated with any feature vector.

C. Graph Assembly Processing Networks (GRACE)

Based on the formulation of a *step-by-step* sequential decision-making process per each part in the assembly in IV-A, we introduce **GR**aph **A**ssembly **pr**o**C**essing **n**Etworks, for short GRACE, $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$, where $a_i = \{y_i | y_i \geq \lambda\}_{i=1}^N$, to extract useful information from the Assembly Graph and predict the next action given the current state of an assembly of N parts. $\lambda \geq 0$ is a threshold used to control the quality of predicted sequences. GRACE outputs a score per part $y_i \in [0, 1], i \in \{1, \dots, N\}$, reflecting the probability of placing the i -th part next. We further articulate the main components of this network (Fig. 2), describe the algorithm for predicting the entire sequence of length N by traversing predicted steps and the way we infer the feasibility of a given assembly.

1) *Surface and Part Blocks*: The architecture is made of identical blocks, which are applied sequentially to obtain updated node features. Each block is made of a GAT [28], an Instance Normalization layer [32] and a Tanh function. We choose GAT as it allows to utilize the rich semantics of edge features for updating node features in our graph representation. Surface Blocks are applied on surface nodes \mathcal{V}^s and surface-to-surface edges \mathcal{E}^{s-to-s} for updating surface node features $\phi(v_i^s)$, while Part Blocks are applied on surface nodes \mathcal{V}^s , part nodes \mathcal{V}^p and surface-to-part edges \mathcal{E}^{s-to-p} to update part node features $\phi(v_i^p)$.

2) *Prediction Head and Loss Function*: To obtain a score per part, a fully-connected layer followed by a Sigmoid function is applied on each part node. During training, we minimize the loss between the network outputs and the ground-truth sequence steps from a dataset of assembly sequences (see IV-A) using binary cross-entropy. To note that, we apply this loss function for each part node separately. Our objective function (4) includes an additional regularization term (5), aiming at encouraging the network not to predict already placed parts:

$$L_\theta = \sum_{i=1}^M \sum_{j=1}^{N_i} (\hat{y}_{ij} \cdot \log(y_{ij}) + (1 - \hat{y}_{ij}) \log(1 - y_{ij})) + \delta L_{reg}, \quad (4)$$

$$L_{reg} = \sum_{i=1}^M \sum_{j=1}^{N_i} f_{ij} \cdot y_{ij}, \quad (5)$$

where M is the number of data examples in the dataset, N_i is the number of nodes in the i -th graph. Abusing the notations, we denote y_{ij} and \hat{y}_{ij} the output score of the model π_θ and the ground-truth step in a sequence for the j -th node in the i -th graph respectively. δ is a weighing coefficient and f_{ij} the value of the *assembled-flag* in the input features.

3) *Predicting Sequences*: As described, GRACE predicts a set of possible next steps based on the current state of an assembly. In order to generate a complete sequence (i.e. of length N), we repeatedly apply GRACE based on the current predicted state of the Assembly Graph. We devise an algorithm (Algo. 1) to traverse the assembly state tree using *Depth First Search* (DFS):

Starting with the graph in its initial state \mathcal{G}_0 – for all part nodes, *assembled-flags* are set to zero, the algorithm performs the following steps recursively: First, it checks the exit condition of the recursion – if all parts are already in place. Next, it predicts the probability for each part node y_i and picks those larger than the threshold λ , controlling the trade-off between precision and recall. Each of those nodes spawns a new branch individually. Therefore, we set the *assembled-flag* and call the recursion on the altered graph to retrieve possible sequences starting with the chosen node. Finally, we add the chosen nodes to the head of each returned sequence and return.

4) *Feasibility Prediction*: To address the issue from infeasible assemblies, we develop two schemes to infer the feasibility (defined in IV-A) of a given assembly: (1) We use the number of predicted complete sequences (output

Algorithm 1 Assembly State Tree Traversal

```
function TRAVERSE-TREE(Model  $M$ , Assembly Graph  
 $\mathcal{G}_t = (\mathcal{V}, \mathcal{E})$ , Threshold  $\lambda$ )  
   $S \leftarrow \text{list}()$   
  if ( $\forall v \in \mathcal{V} : v.\text{assembled-flag} == 1$ ) then  
    return  $S$  ▷ Exit: all parts assembled  
  end if  
   $\mathbf{y} \leftarrow M(\mathcal{G}_t)$   
  for  $i \leftarrow 1$  to  $|\mathcal{V}|$  do  
    if  $\mathbf{y}[i] < \lambda$  then  
      continue  
    end if  
     $\mathcal{G}_{t+1} \leftarrow \text{copy}(\mathcal{G}_t)$   
     $[\mathcal{V}_{t+1}]_i.\text{assembled-flag} \leftarrow 1$  ▷ assembled node  $i$   
     $S_* \leftarrow \text{TRAVERSE-TREE}(M, \mathcal{G}_{t+1}, \lambda)$   
    for  $s$  in  $S_*$  do  
       $s_* \leftarrow [i] + s$  ▷ Add current part to the sequence  
       $S.\text{append}(s_*)$   
    end for  
  end for  
  return  $S$   
end function
```

by Algo. 1) as an indicator for the feasibility of a given assembly. If no sequences were retrieved, the assembly is predicted as infeasible. (2) We aggregate the features of all part nodes from a pre-trained GRACE with a *mean-pooling* operation, creating a feature vector for the entire assembly graph. This feature vector is then used to train a binary classifier for feasibility prediction, where we analyze several classifiers i.e. Support Vector Machines (SVMs), Multi-layer Perceptrons (MLPs) and Nearest Neighbor.

V. EXPERIMENT

In this section, we first describe the experimental setup such as our dataset, evaluation metrics and implementation details. To note that we use the term *size* to describe the number of parts of an assembly without prior notice. We evaluate the Sequence Prediction under two experimental protocols with 4-fold cross-validation: (1) **intra-sized**: the assemblies in training and test set share *the same* sizes; (2) **inter-sized**: the assemblies in training and test set have *different* sizes, where there are two sub-protocols: Many-to-one and One-to-Many (detailed in V-B.2) The results on Feasibility Prediction are presented before the failure analysis and ablation study.

A. Experimental Setup

1) *Dataset*: We applied our in-house simulation software MediView to randomly generate data of synthetic aluminium assemblies whose sizes range from 3 to 7 (denoted by A_i , where i is the size). The simulation software was tasked with putting together the structures by brute-forcing all part orders, while considering the restrictions of part geometries or those imposed by the capabilities of a dual-armed robotic system *KUKA LBR Med* (Fig. 1). More restrictions could

be added to this environment in a future work, e.g. taking into account grasp planning. An illustration of this process is given in Fig. 3. The resulting data consists of the following amount per size: $A_3 : 5717$, $A_4 : 2464$, $A_5 : 6036$, $A_6 : 2865$, $A_7 : 431$.

We post-processed the simulation output to obtain the *Placement Action* (required during training) – the next possible placement actions given a state of an assembly in a feasible sequence. In addition, we derive the *Feasibility* of each assembly based on the number of ground truth sequences e.g. 0 indicates an infeasible assembly.

2) *Metrics*: We use the following metrics for the sequence prediction task: (1) **Step-by-Step AUC** examines our method’s predictive performance to infer the parts that should be assembled next given the current state by comparing the ground truth binary labels with the predicted step scores. For this purpose we use the common Precision-Recall curve w.r.t. λ and finally deriving an *Area Under Curve* (AUC) score. (2) **Complete-Sequence AUC** evaluates the ability to infer the entire set of ground truth sequences, since a step-by-step evaluation only partially displays our application¹. We use Information Retrieval (IR) Precision-Recall [33], devised for set prediction evaluation, computed as $\text{IR-Precision} = |\text{RET} \cap \text{REL}| / |\text{RET}|$, $\text{IR-Recall} = |\text{RET} \cap \text{REL}| / |\text{REL}|$, where *RET* are the retrieved sequences and *REL* are the relevant sequences (i.e. ones in the ground truth set). Here, again, we plot an IR Precision-Recall curve and derive an AUC score. (3) **Precision@k (P@k)**: since in practice we only consider the highest scored predicted sequences, we also compute IR-Precision while taking into account only the top- k ones [34].

For feasibility prediction, we use common binary classification metrics *False Positive Rate FPR* and *True Positive Rate TPR*. In this setting, a positive instance is a feasible assembly and a negative an infeasible one.

3) *Implementation Details*: We use PyTorch Geometric (PyG) [35] to build the model which consists of 3 surface blocks and 1 part block with a latent-dimensionality of 94. For training, we choose a batch size of 256 and a learning rate of 0.0022 with Adam optimizer based on validation performance on 15% of the training samples during hyperparameter search. Besides, we set the regularization weight in Eq. 4 to 0.3 and length of positional encoding for node features to 16. For the training and evaluation of the feasibility classifiers a balanced dataset is used. Our model includes 51.7K trainable parameters and requires 4.06 ± 0.15 ms to infer the next feasible sequence step². More details are referred to our open-sourced code.

B. Results

1) *Sequence Prediction for Intra-sized Assemblies*: The results are shown separately per assembly size (Tab. II), including the step-by-step and complete-sequence AUC, P@k

¹Consider a method that predicts the first 97 steps correctly and fails in the 98-th step for a 100-parts-assembly.

²Measured on NVIDIA GeForce GTX 1080.

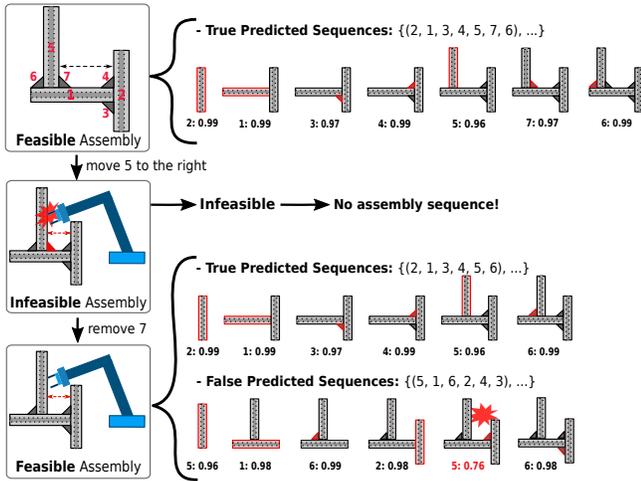


Fig. 3: **Examples of ASP for Aluminum Assemblies.** We demonstrate the complexity of our task through the predictions for three different assemblies: starting from the top, there is a feasible assembly sequence predicted based on the assembly on the left. By decreasing the distance between part 5 and 4, it becomes infeasible due to limited space for the robot arm. Further, by removing part 7, with certain sequences such as the one shown in the figure, it is feasible. But this does not work with another, i.e. the false predicted sequence, because this one would cause collision.

scores for $k \in \{1, 2, 3\}$ with a threshold of 0.5. GRACE is able to perform perfectly on step prediction for all sizes. More relevant to our goal and more challenging than step prediction, our method can reach 1.0 for small sizes (e.g. 3 and 4 parts) on the task of complete sequence prediction. However, we can observe a slight drop for larger sizes (e.g. 5, 6 and 7 parts), implying the greater complexity for large assemblies. Hence, GRACE can effectively learn an useful inductive bias from our proposed graph representation when trained with similar sizes. To note that, this performance has already reached that of the approach in [7] and GRACE is able to generalize to larger sizes which the previous one is incapable of (see Tab. I for a qualitative comparison).

2) Sequence Prediction for Inter-sized Assemblies:

To comprehensively evaluate the *generalization* ability of GRACE across different sizes, a distinct limitation of previous works [7], [36], we further design two more *challenging* sub-protocols under the inter-sized protocol.

- **Many-to-one:** GRACE is trained on assemblies of mixed sizes but i , i.e. $A_{\forall j \neq i}$, and tested on A_i .
- **One-to-many:** GRACE is trained on a single-sized dataset A_i and tested on all the other, i.e. $A_{\forall j \neq i}$.

1. *Many-to-one:* This setting is similar to the intra-sized one except that we excluded assemblies of the size evaluated at test time from the training set. When comparing the results in this setting to the intra-sized ones (Tab. II), we observe a slight performance decrease in AUC on step and sequence prediction. However, note that $P@1$ and $P@2$ can still reach ~ 1.0 for small sizes (3 and 4 parts) and ~ 0.9 for large

sizes, indicating that GRACE is capable on generalizing to assembly variants with different sizes that have not been seen before.

2. *One-to-many:* This setting is an inverse version of the previous one, which is more challenging, since the amount and diversity of the training set are much lower than before³. The results (lower triangular block in Tab. III) provide a clear pattern that GRACE is able to obtain comparably better results for assemblies with less parts. For instance, trained with only A_5 , GRACE preforms well for A_3 and A_4 which is reasonable as the constraints guiding smaller assembly structures are *contained* in larger ones. This shows the *generalization* capability and *sample efficient* learning ability (trained on single size and worked on smaller sizes) of our method. Nevertheless, the performance drops for larger assemblies (see the upper triangular block in Tab. III). We hypothesize that an increasing amount of items introduces new constraints that are not covered by the training data. Thus, there might be a critical number of items containing all possible constraints that if included in the training data can lead to an overall generalizing model.

3) *Feasibility Prediction:* In this setting, we examine the ability of our approach to detect infeasible assemblies. For this experiment we consider GRACE trained on multiple sizes and test on the A_5 set. As mentioned in IV-C, we compare the implicit approach via the number of predicted sequences (Algo. 1) and alternative schemes exploiting the graph representation of the pre-trained GRACE. Therefore, we explore several binary classifiers i.e. SVMs, a Multi-layer Perceptron (MLP) and Nearest Neighbor. As seen in Fig. 4, GRACE (*#sequences*) is able to detect infeasible assemblies (AUC of 0.97). However, training our method exclusively with feasible assemblies (*#sequences, feasible only*) results in a poor detection performance. We hypothesize that by missing infeasible structures during training the method learns to always assemble an item leading to overconfidence. Exploiting an additional scheme by adding one of the classifiers (except SVM with RBF kernel) maintains or even slightly improves the performance.

C. Failure Analysis

To better understand the limitations of our method, we conduct an analysis of falsely predicted assembly sequences by our baseline model for A_5 and A_6 . Each of these false predicted sequences includes a *false step*, i.e. action from which the sequence deviates from the corresponding ground truth sequences. Fig. 5 depicts the histogram of false steps binned by their predicted probability. One can observe a large amount false steps performed in the beginning of the sequence (steps 1 and 2) with a high confidence. On the other hand, wrong step predictions at the end of an assembly (steps 4 and 5) exhibit lower confidence scores. We hypothesize that this bias is a result of an inherit imbalance in our training setting. Our dataset samples could be thought of as nodes in

³We do not perform this experiment on A_7 , as there are relatively small amount of assemblies with 7 parts in the dataset.

TABLE II: Sequence Prediction Results for intra-sized and inter-sized (many-to-one) assemblies.

Metrics	Intra-sized					Inter-sized				
	A ₃	A ₄	A ₅	A ₆	A ₇	A ₃	A ₄	A ₅	A ₆	A ₇
Step-by-Step AUC (↑)	1.00±0.00	1.00±0.00	1.00±0.00	1.00±0.00	1.00±0.00	0.99±0.02	1.00±0.00	0.98±0.10	0.98±0.00	1.00±0.00
Complete Sequence AUC (↑)	1.00±0.00	1.00±0.00	0.96±0.02	0.93±0.03	0.97±0.02	0.97±0.03	1.00±0.10	0.87±0.03	0.90±0.04	0.95±0.07
P@1 (↑)	1.00±0.00	1.00±0.00	0.95±0.04	0.96±0.06	0.99±0.01	0.99±0.01	0.98±0.03	0.90±0.09	0.88±0.07	0.96±0.05
P@2 (↑)	1.00±0.00	1.00±0.00	0.94±0.04	0.95±0.07	0.99±0.01	0.99±0.01	0.97±0.03	0.87±0.12	0.87±0.07	0.96±0.04
P@3 (↑)	–	1.00±0.00	0.99±0.01	0.95±0.08	0.99±0.02	–	0.95±0.06	0.93±0.10	0.85±0.08	0.96±0.04

TABLE III: Sequence Prediction Results for inter-sized assemblies in one-to-many setting.

Training Set	Step-by-Step AUC (↑) on assemblies of various sizes					Complete Sequence AUC (↑) on assemblies of various sizes				
	A ₃	A ₄	A ₅	A ₆	A ₇	A ₃	A ₄	A ₅	A ₆	A ₇
A ₄	0.92±0.11	–	0.48±0.12	0.41±0.11	0.43±0.12	0.93±0.09	–	0.28±0.12	0.25±0.15	0.25±0.15
A ₅	0.93±0.06	0.89±0.07	–	0.78±0.14	0.59±0.16	0.83±0.07	0.70±0.09	–	0.36±0.12	0.24±0.07
A ₆	0.90±0.10	0.89±0.11	0.93±0.04	–	0.59±0.16	0.73±0.16	0.68±0.24	0.71±0.13	–	0.24±0.07

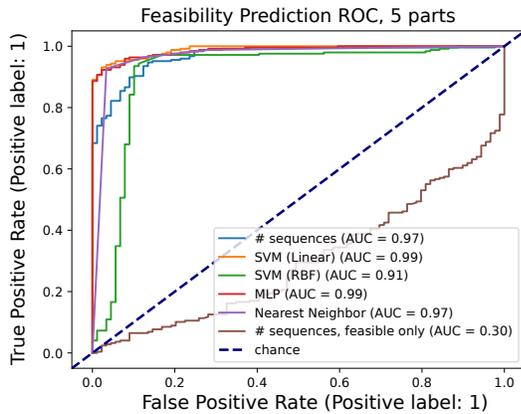


Fig. 4: **Results of Feasibility Prediction.** Comparison of different proposed schemes for feasibility classification and an ablation study in which infeasible ones are unavailable based on assemblies with 5 parts.

a state tree, where earlier steps share state nodes closer to the root and later ones have independent nodes towards the leaves. As each of these nodes is represented only once, there are fewer samples in the dataset attributed to earlier steps. This problem could be solved by balancing the training set based on the sequence step.

D. Ablation Study

In Assembly Graph (IV-B), both the part and surface node embeddings contain a $16d$ sinusoidal positional encoding [31]. We conducted an ablation study to investigate the impacts from the values and permutation order thereof based on A_5 . (1) **Values:** Initializing the positional encoding with random values dramatically decrease the performance of our method (Tab. IV). We hypothesize that these positions introduce *geometrical bias*, which is helpful in our task. (2) **Permutation Order:** We number assembly parts and surfaces in a constant order. Parts are counted beginning from the one closest to the environment origin. Surfaces, on the other hand, are always numbered clockwise, starting from the respective part top. Permuting both part and surface

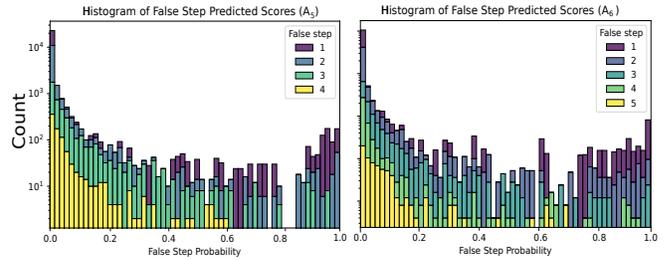


Fig. 5: **Failure Analysis: false predicted sequences.** Histogram of predicted probability (for A_5 and A_6) in false steps reveals a drift in which GRACE is overconfident in mistakes preformed early. This is an evidence for an inherit bias in our training setting.

TABLE IV: Ablation study into the contribution of positional encodings to our method.

Positional Encoding	A ₅ AUC (↑)
Baseline, sinusoidal encoding [31]	0.94
Random values	0.37
No encodings	0.07
Part permutations (test time)	0.60
Surface permutations (test time)	0.27
Part permutations (training and test time)	0.97
Surface permutations (training and test time)	0.10

orders *only* at test time causes severe performance degradation, indicating constant numbering during training harms the model’s ability to generalize (Tab. IV). Interestingly, allowing permutations for only part order can boost the performance while this is not the case for surface permutations. This demonstrates the importance of these features for the network to extract information from the parts’ geometrical structure.

VI. CONCLUSION

In this work, we addressed the RASP problem with a learning-based framework. Concretely, we propose a graph representation, called Assembly Graphs for the aluminum profile assemblies, which is flexible to represent different

2d structures and meanwhile agnostic to rotation and mirroring. Based on this, a novel policy network – GRACE is introduced to extract meaningful information for assembly sequence prediction. Extensive experiments in simulation verify the capability of transferring knowledge between different assembly tasks, on which previous methods fall short. Further, our method can generalize knowledge gained on larger assemblies and then apply it to smaller ones. Last but not least, it is worth to mention, though only validated in simulation, our method should address the challenges during the real-world deployment like not finding a valid motion or a feasible grasping point if these cases are enclosed in the training data and learned to reject by GRACE. Meanwhile encouraged by the superior results on objects with simple geometries, our holistic graphical method lays a solid basis for handling complex 3d objects like curve blocks in the future.

ACKNOWLEDGMENTS

The paper received partial funding by the DLR internal project Factory of the Future Extended (FoF-X). Jianxiang Feng is supported by the Munich School for Data Science (MUDS). Rudolph Triebel is a member of MUDS.

REFERENCES

- [1] W. C. Shih, “Global supply chains in a post-pandemic world,” *Harvard Business review*, 98(5), 82-89, 2020.
- [2] M. F. F. Rashid, W. Hutabarat, and A. Tiwari, “A review on assembly sequence planning and assembly line balancing optimisation using soft computing approaches,” *International Journal of Advanced Manufacturing Technology*, vol. 59, pp. 335–349, 2011.
- [3] F. Suárez-Ruiz, X. Zhou, and Q.-C. Pham, “Can robots assemble an ikea chair?” *Science Robotics*, vol. 3, no. 17, p. eaat6385, 2018.
- [4] M. Zhao, X. Guo, X. Zhang, Y. Fang, and Y. Ou, “Aspw-drl: assembly sequence planning for workpieces via a deep reinforcement learning approach,” *Assembly Automation*, 2019.
- [5] K. Watanabe and S. Inada, “Search algorithm of the assembly sequence of products by using past learning results,” *International Journal of Production Economics*, vol. 226, p. 107615, 2020.
- [6] I. Rodríguez, K. Nottensteiner, D. Leidner, M. Kaßecker, F. Stulp, and A. Albu-Schäffer, “Iteratively refined feasibility checks in robotic assembly sequence planning,” *IEEE Robotics and Automation Letters (RAL)*, vol. 4, no. 2, pp. 1416–1423, 2019.
- [7] I. Rodríguez, K. Nottensteiner, D. Leidner, M. Durner, F. Stulp, and A. Albu-Schäffer, “Pattern recognition for knowledge transfer in robotic assembly sequence planning,” *IEEE RAL*, vol. 5, no. 2, pp. 3666–3673, 2020.
- [8] Y. Lin, A. S. Wang, E. Undersander, and A. Rai, “Efficient and interpretable robot manipulation with graph neural networks,” *IEEE RAL*, vol. 7, no. 2, pp. 2740–2747, 2022.
- [9] J. Feng, M. Durner, Z.-C. Márton, F. Bálint-Benczédi, and R. Triebel, “Introspective robot perception using smoothed predictions from bayesian neural networks,” in *Robotics Research: The 19th International Symposium ISRR*. Springer, 2022, pp. 660–675.
- [10] L. Ma, J. Gong, H. Xu, H. Chen, H. Zhao, W. Huang, and G. Zhou, “Planning assembly sequence with graph transformer,” *arXiv preprint arXiv:2210.05236*, 2022.
- [11] L. Homem de Mello and A. Sanderson, “And/or graph representation of assembly plans,” *IEEE Transactions on Robotics and Automation*, vol. 6, no. 2, pp. 188–199, 1990.
- [12] T. De Fazio and D. Whitney, “Simplified generation of all mechanical assembly sequences,” *IEEE Journal on Robotics and Automation*, vol. 3, no. 6, pp. 640–658, 1987.
- [13] U. Thomas, M. Barrenscheen, and F. Wahl, “Efficient assembly sequence planning using stereographical projections of c-space obstacles,” in *Proceedings of the IEEE International Symposium on Assembly and Task Planning, 2003.*, 2003, pp. 96–102.
- [14] K. Nottensteiner, T. Bodenmueller, M. Kassecker, M. A. Roa, A. Stemmer, T. Stouraitis, D. Seidel, and U. Thomas, “A complete automated chain for flexible assembly using recognition, planning and sensor-based execution,” in *Proceedings of ISR 2016: 47st International Symposium on Robotics*, 2016, pp. 1–8.
- [15] U. Thomas, T. Stouraitis, and M. A. Roa, “Flexible assembly through integrated assembly sequence planning and grasp planning,” in *2015 IEEE International Conference on Automation Science and Engineering (CASE)*, 2015, pp. 586–592.
- [16] B. Li, Y. Wu, H. Sun, Z. Cheng, and J. Liu, “Unity 3d-based simulation data driven robotic assembly sequence planning using genetic algorithm,” in *2022 14th International Conference on Computer and Automation Engineering (ICCAE)*, 2022, pp. 1–7.
- [17] Iwankowicz and R.R., “An efficient evolutionary method of assembly sequence planning for shipbuilding industry,” *Assembly Automation*, vol. 36, no. 1, pp. 60–71, 2016.
- [18] W.-C. Chen, P.-H. Tai, W.-J. Deng, and L.-F. Hsieh, “A three-stage integrated approach for assembly sequence planning using neural networks,” *Expert Systems with Applications*, vol. 34, no. 3, pp. 1777–1786, 2008.
- [19] C. Sinanoğlu and H. R. Börklü, “An assembly sequence-planning system for mechanical parts using neural network,” *Assembly Automation*, 2005.
- [20] A. Rashid and M.F.F., “A hybrid ant-wolf algorithm to optimize assembly sequence planning problem,” *Assembly Automation*, vol. 37, no. 2, pp. 238–248, 2017.
- [21] S. Nguyen, O. S. Oguz, V. N. Hartmann, and M. Toussaint, “Self-supervised learning of scene-graph representations for robotic sequential manipulation planning,” in *Conference on Robot Learning (CoRL)*, 2020, pp. 2104–2119.
- [22] V. Bapst, A. Sanchez-Gonzalez, C. Doersch, K. Stachenfeld, P. Kohli, P. Battaglia, and J. Hamrick, “Structured agents for physical construction,” in *ICML*. PMLR, 2019, pp. 464–474.
- [23] Y. Zhu, J. Tremblay, S. Birchfield, and Y. Zhu, “Hierarchical planning for long-horizon manipulation with geometric and symbolic scene graphs,” in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 6541–6548.
- [24] N. Funk, G. Chalvatzaki, B. Belousov, and J. Peters, “Learn2assemble with structured representations and search for robotic architectural construction,” in *CoRL*. PMLR, 2022, pp. 1401–1411.
- [25] Y. Ye, D. Gandhi, A. Gupta, and S. Tulsiani, “Object-centric forward modeling for model predictive control,” in *CoRL*. PMLR, 2020, pp. 100–109.
- [26] R. Li, A. Jabri, T. Darrell, and P. Agrawal, “Towards practical multi-object manipulation using relational reinforcement learning,” in *2020 IEEE ICRA*. IEEE, 2020, pp. 4051–4058.
- [27] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, “Neural message passing for quantum chemistry,” in *ICML*. PMLR, 2017, pp. 1263–1272.
- [28] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, “Graph attention networks,” *Stat*, vol. 1050, no. 20, pp. 10–48 550, 2017.
- [29] X. Wang, H. Ji, C. Shi, B. Wang, Y. Ye, P. Cui, and P. S. Yu, “Heterogeneous graph attention network,” in *The world wide web conference*, 2019, pp. 2022–2032.
- [30] R. Bellman, “A markovian decision process,” *Journal of mathematics and mechanics*, pp. 679–684, 1957.
- [31] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [32] D. Ulyanov, A. Vedaldi, and V. Lempitsky, “Instance normalization: The missing ingredient for fast stylization,” *arXiv preprint arXiv:1607.08022*, 2016.
- [33] W. B. Croft, D. Metzler, and T. Strohman, *Search engines: Information retrieval in practice*. Addison-Wesley Reading, 2010, vol. 520.
- [34] J. L. Herlocker, J. A. Konstan, L. G. Terveen, and J. T. Riedl, “Evaluating collaborative filtering recommender systems,” *ACM Transactions on Information Systems (TOIS)*, vol. 22, no. 1, pp. 5–53, 2004.
- [35] M. Fey and J. E. Lenssen, “Fast Graph Representation Learning with PyTorch Geometric,” 5 2019. [Online]. Available: https://github.com/pyg-team/pytorch_geometric
- [36] A. M. Wells, N. T. Dantam, A. Shrivastava, and L. E. Kavraki, “Learning feasibility for task and motion planning in tabletop environments,” *IEEE RAL*, vol. 4, no. 2, pp. 1255–1262, 2019.