





MACHINE LEARNING

Manual-PA: Learning 3D Part Assembly from Instruction Diagrams

Jiahao Zhang¹ Anoop Cherian² Cristian Rodriguez³ Weijian Deng¹ Stephen Gould¹

¹The Australian National University ²Mitsubishi Electric Research Labs ³The Australian Institute for Machine Learning ¹{first.last}@anu.edu.au ²cherian@merl.com ³crodriguezop@gmail.com



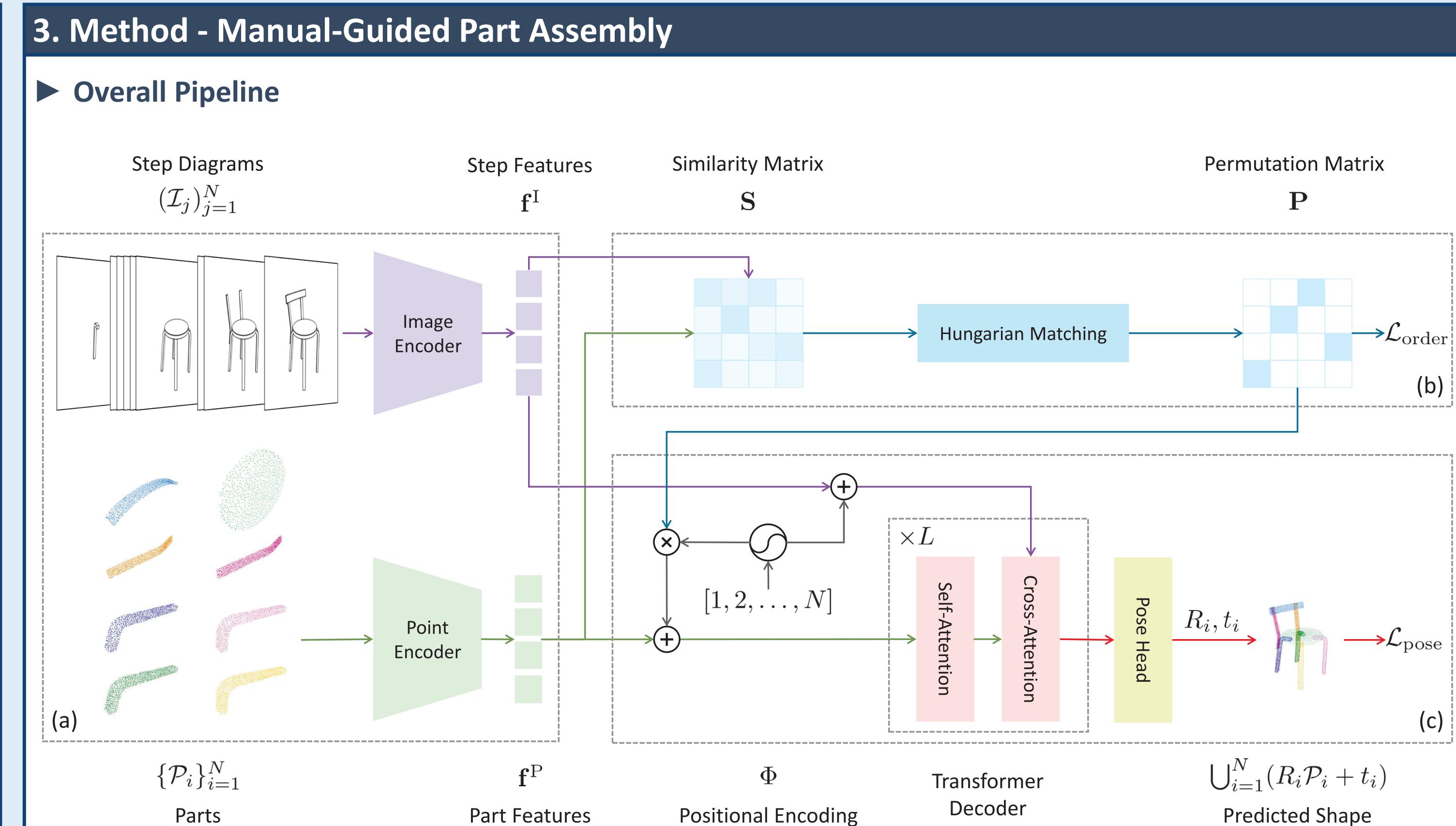


1. Introduction **Problem Statement** \rightarrow $R_3(-)+t_3$

Given (a) a diagrammatic manual book demonstrating the step-by-step assembly process and (b) a set of texture-less furniture parts, the goal is to (c) infer the order of parts for the assembly from the manual sequence and predict the 6DoF pose for each part such that the spatially transformed parts assembles the furniture described in the manual.

2. Challenges & Contributions

- Key Challenges
- 1. Cross-Modal Alignment: How to match 2D diagrams with 3D parts to infer the assembly sequence?
- 2. Guided Pose Estimation: How to leverage the inferred sequence as guidance without accumulating errors?
- Contributions
- 1. A Novel Task: We introduce and formulate the new task of manualguided 3D part assembly.
- 2. A Novel Framework (Manual-PA): We propose an end-to-end Transformer-based model that learns the assembly order via contrastive learning and uses it as soft guidance for robust pose estimation.
- 3. SOTA Results: Our method significantly outperforms prior works on the PartNet benchmark and shows strong zero-shot generalization to the real-world IKEA-Manual dataset.



- ► Manual-Guided Part Permutation Learning
- 1. Similarity

$$\sigma = \operatorname{argmax}(\mathbf{P})$$

2. Hungarian Matching

$$egin{aligned} & \min egin{aligned} & \sum_{i=1}^N \sum_{j=1}^N \mathbf{C}_{ij} \mathbf{P}_{ij} \end{aligned} \ & \sup \mathbf{C}_{ij} \mathbf{P}_{ij} = 1, \quad orall i \in \{1,\ldots,N\} \end{aligned} \ & \sum_{i=1}^N P_{ij} = 1, \quad orall j \in \{1,\ldots,N\} \end{aligned} \ & P_{ij} \in \{0,1\}, \quad orall i, j \in \{1,\ldots,N\} \end{aligned}$$

3. Order

$$\mathbf{S}_{ij} = \mathrm{sim}(\mathbf{f}_i^{\mathrm{P}}, \mathbf{g}_j^{\mathrm{I}})$$

- Manual-Guided Part Pose Estimation
- . Sinusoidal Encoding

$$\phi(x)_{2i}=\sinigg(rac{x}{ au_p^{2i/d}}igg),\;\phi(x)_{2i-1}=\cosigg(rac{x}{ au_p^{2i/d}}igg)$$

2. Permutation

$$\hat{\mathbf{p}}^{\mathrm{P}}=\hat{\mathbf{P}}^{T}\Phi$$

Losses

- 1. Loss for Permutation Learning
 - (a) Contrastive Loss

$$\mathcal{L}_{ ext{order}} = -rac{1}{B} \sum_{i=1}^{B} \log rac{\exp(ext{sim}(\mathbf{f}_{~\sigma(i)}^{ ext{P}}, \mathbf{g}_{i}^{ ext{I}})/ au)}{\sum_{j=1}^{B} \exp(ext{sim}(\mathbf{f}_{~\sigma(i)}^{ ext{P}}, \mathbf{g}_{j}^{ ext{I}})/ au)}$$

- 2. Losses for Pose Estimation (a) Translation Loss (L2)
 - $\mathcal{L}_C = rac{1}{N} \sum_{i=1}^{N} ext{CD}(\hat{R}_{[i]} \mathcal{P}_i, R_i \mathcal{P}_i)$
- (b) Rotation Loss (Chamfer)

$$\mathcal{L}_T = rac{1}{N} \sum_{i=1}^N \lVert \hat{t}_{[i]} - t_i
Vert_2$$

(c) Rotation Loss (L2)

$$\mathcal{L}_E = rac{1}{N} \sum_{i=1}^N \lVert \hat{R}_{[i]} \mathcal{P}_i - R_i \mathcal{P}_i
Vert_2$$

(d) Shape Loss (Chamfer)

$$\mathcal{L}_S = ext{CD}igg(igl_{i=1}^N (\hat{R}_{[i]}\mathcal{P}_i + \hat{t}_{[i]}), \ igcup_{i=1}^N (R_i\mathcal{P}_i + t_i)igg)$$

(e) Final Loss

$$\mathcal{L}_{ ext{pose}} = \lambda_T \mathcal{L}_T + \lambda_C \mathcal{L}_C + \lambda_E \mathcal{L}_E + \lambda_S \mathcal{L}_S$$

. Results							
Method	Condition	SCD↓		PA†		SR↑	
		Chair	Table	Chair	Table	Chair	Table
Fully-Supervised on Partl	Vet [26]						
DGL _{NIPS'20} [52] IET _{RA-L'22} [55] Score-PA _{BMVC'23} [8] CCS _{AAAI'24} [56]	- - -	9.1 5.4 7.4 7.0	5.0 3.5 4.5	39.00 62.80 42.11 53.59	49.51 61.67 51.55	- 8.320 -	- 11.23 -
RGL _{WACV'22} [27] SPAFormer _{ArXiv'24} [51] Joint-PA _{CVPR'24} [24] Image-PA _{ECCV'20} [22] Image-PA [†] _{ECCV'20}	Sequence Sequence Joint Image Diagram	5.1 8.7 6.7 6.0 6.7 5.9	2.8 4.8 3.8 7.0 3.7 3.9	49.06 55.88 72.80 45.40 62.67	54.16 64.38 67.40 71.60 70.10	- 16.40 - - 19.97	- 33.50 - - 32.83
Manual-PA w/o Order Manual-PA (Ours)	Manual Manual	3.0 1.7	3.6 1.8	79.07 89.06	74.03 87.41	34.13 58.03	37.71 56.66
Zero-Shot on IKEA-Manu	al [47]						
3DHPA _{CVPR'24} Image-PA [†] _{ECCV'20}	- Diagram	34.3 17.3	37.8 14.7	1.914 19.07	4.027 36.74	0.000	0.000 <u>10.53</u>
Manual-PA w/o Order Manual-PA (Ours)	Manual Manual	12.8 11.4	8.9 4.8	38.36 42.51	42.01 49.72	1.754 3.509	10.53 15.79
Input Last Step Diagram	Input Point Clouds	3DHPA	Image-PA		anual-PA 'o Order	Manual-PA (Ours)	Ground Truth
(a)							
(b)							

