

# 000 GENCAPE: STRUCTURE-INDUCTIVE GENERATIVE 001 MODELING FOR CATEGORY-AGNOSTIC POSE ESTIMA- 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 GENCAPE: STRUCTURE-INDUCTIVE GENERATIVE MODELING FOR CATEGORY-AGNOSTIC POSE ESTIMA- TION

006 **Anonymous authors**

007 Paper under double-blind review

## ABSTRACT

Category-agnostic pose estimation (CAPE) aims to localize keypoints on query images from arbitrary categories, using only a few annotated support examples for guidance. Recent approaches either treat keypoints as isolated entities or rely on manually defined skeleton priors, which are costly to annotate and inherently inflexible across diverse categories. Such oversimplification limits the model’s capacity to capture instance-wise structural cues critical for accurate pixel-level localization. To overcome these limitations, we propose **GenCape**, a Generative-based framework for CAPE that infers keypoint relationships solely from image-based support inputs, without additional textual descriptions or predefined skeletons. Our framework consists of two principal components: an iterative Structure-aware Variational Autoencoder (i-SVAE) and a Compositional Graph Transfer (CGT) module. The former infers soft, instance-specific adjacency matrices from support features through variational inference, embedded layer-wise into the Graph Transformer Decoder for progressive structural priors refinement. The latter adaptively aggregates multiple latent graphs into a query-aware structure via Bayesian fusion and attention-based reweighting, enhancing resilience to visual uncertainty and support-induced bias. This structure-aware design facilitates effective message propagation among keypoints and promotes semantic alignment across object categories with diverse keypoint topologies. Experimental results on the MP-100 dataset show that our method achieves substantial gains over graph-support baselines under both 1- and 5-shot settings, while maintaining competitive performance against text-support counterparts.

## 1 INTRODUCTION

Category-Agnostic Pose Estimation (CAPE) Xu et al. (2022a); Shi et al. (2023); Hirschorn & Avidan (2024); Rusanovsky et al. (2025); Ren et al. (2024); Chen et al. (2025a) has emerged as a fundamental yet challenging task in computer vision, aiming to localize semantic keypoints on arbitrary object categories using only a handful of annotated support samples. Unlike conventional 2D pose estimation task Sun et al. (2019); Xu et al. (2022b); Yuan et al. (2021); Rao et al. (2025), which depends heavily on predefined templates or class-specific priors, CAPE requires robust generalization across semantically diverse and structurally heterogeneous object classes. This capability extends the applicability of pose estimation from closed-world scenarios to open-world scenarios, enabling scalable deployment in domains such as human motion analysis Zheng et al. (2023); Yang et al. (2023), cross-species behavior understanding Ye et al. (2024); Stoffl et al. (2024), and robotic manipulation Zheng et al. (2025); Ma et al. (2024) in dynamic environments.

In the CAPE paradigm, the objective is to estimate keypoints for objects from novel categories, conditioned on a few annotated support images. Despite recent advancements, existing CAPE approaches are hindered by two critical limitations. On the one hand, many existing approaches either treat keypoints as isolated semantic entities, neglecting the latent spatial dependencies essential for accurate pixel-level localization, or rely heavily on external priors such as manually pre-defined skeleton connections or auxiliary textual descriptions. These external priors not only incur high annotation overhead but also restrict the model’s adaptability to novel instances with large pose variations, non-rigid deformation, or diverse structural characteristics.

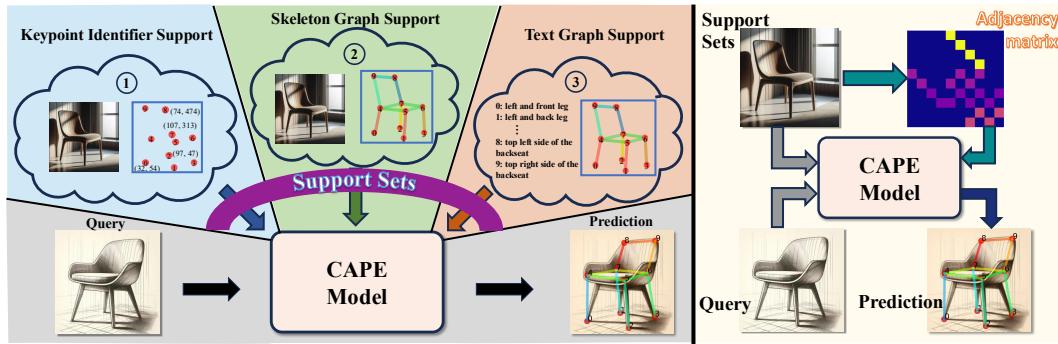


Figure 1: **Left:** Existing CAPE frameworks rely on additional structured priors within the support set, such as ① keypoint identifiers, ② fixed skeleton graphs, or ③ textual descriptions with skeleton graphs, to enhance structural reasoning. **Right:** In contrast, our framework directly infers latent keypoint relationships solely from support images, learning instance-specific *adjacency matrices* without relying on handcrafted priors.

On the other hand, the stochastic nature of support sets selection in a few-shot setting makes CAPE methods particularly vulnerable to the low quality support examples. In real-world scenarios, support images may contain severe occlusions, incomplete annotations, or structural discrepancies relative to the query, which can misguide structural inference and significantly impair prediction accuracy and generalization. Together, these limitations underscore the need for a more flexible, data-driven approach to structure modeling and robust support-query adaptation.

To address the limitations of fixed priors and structural rigidity in CAPE, we propose **GenCape**, a generative latent structure learning framework that automatically infers keypoint relationships, represented as latent adjacency matrices, exclusively from image-based support inputs, without any external priors, as illustrated in Figure 1. At its core, GenCape integrates two complementary components: an *iterative Structure-aware Variational Autoencoder* (*i-SVAE*) and a *Compositional Graph Transfer* (*CGT*) module. Specifically, the *i-SVAE* leverages variational inference to learn a distribution over instance-specific graph structures, iteratively generating and refining latent adjacency matrices that serve as flexible and data-driven structural priors. Compared to the recently related deterministic approach SDPNet Ren et al. (2024), this generative formulation captures the epistemic uncertainty in sparse and ambiguous support signals, allowing for more expressive and robust message passing within the Graph Transformer Decoder. The progressive refinement enables the model to propagate contextual cues across spatially correlated keypoints and adapt to complex object configurations. In parallel, the *CGT* module dynamically aggregates multiple latent graph hypotheses into a query-conditioned representation through a principled Bayesian fusion and an attention-based re-weighting strategy. This dynamic compositional mechanism explicitly accounts for support-query inconsistencies and mitigates the adverse impact of noisy or misleading support examples caused by occlusion, deformation, or pose variation. Together, these modules enhance structural generalization and resilience to support noise, setting a new direction for few-shot keypoint reasoning under the CAPE paradigm. Remarkably, our approach surpasses the performance of existing CAPE methods, showcasing a new state-of-the-art performance.

In summary, our contributions are as follows:

- We introduce **GenCape**, a novel generative framework for category-agnostic pose estimation, which incorporates an *iterative Structure-aware Variational Autoencoder* (*i-SVAE*) to infer latent, instance-specific skeletons solely from image-based support sets, eliminating the need for predefined anatomical priors or textual descriptions.
- We propose a *Compositional Graph Transfer* (*CGT*) mechanism that dynamically aggregates multiple structural hypotheses into a unified, query-conditioned graph through attention-guided fusion, significantly enhancing robustness under ambiguous, noisy, or structurally mismatched support scenarios.
- Our framework achieves new state-of-the-art results on the representative and challenging **MP-100** benchmark under both 1-shot and 5-shot settings, surpassing existing methods by a substantial margin of **+1.59%** mPCK averaged across evaluation thresholds, without relying on any external structural or textual annotations.

108 

## 2 RELATED WORK

109 

### 2.1 CATEGORY-AGNOSTIC POSE ESTIMATION

110 Category-agnostic pose estimation (CAPE) Xu et al. (2022a) has emerged as a compelling general-  
 111 ization of conventional pose estimation Sun et al. (2019); Yu et al. (2021); Xu et al. (2022b; 2025);  
 112 Rao et al. (2025) by localizing keypoints for arbitrary category objects with only a few annotated  
 113 support images. The pioneering POMNet Xu et al. (2022a) employed a metric-learning paradigm to  
 114 match support and query features in latent space, while CapeFormer Shi et al. (2023), a two-stage  
 115 refinement framework that first produces initial keypoint proposals and then refines their positions  
 116 in a second stage. Recent efforts Hirschorn & Avidan (2024); Liang et al. (2024) further extended  
 117 this paradigm by integrating fixed skeleton priors into graph reasoning modules to capture latent  
 118 keypoint relations. However, it remains limited by the rigidity of manually defined structures, which  
 119 constrains adaptability to novel categories with topological variation. **WeakShot** Chen et al. (2025b)  
 120 **learns category-agnostic keypoints via diffusion-based keyness prediction and correspondence trans-**  
 121 **fer.** Another line of works **Yang et al. (2024); Kim et al. (2024); Lu et al. (2024); Rusanovsky et al.**  
 122 **(2025)** also leverage textual descriptions for guidance, enhancing category-agnostic generalization  
 123 but still depending on auxiliary language priors. Most relevant to our work, SDPNet Ren et al. (2024)  
 124 adopts a discriminative approach by predicting a fixed adjacency matrix from support features. How-  
 125 ever, it lacks mechanisms to model structural uncertainty, limiting robustness under support-query  
 126 mismatch. In contrast, we propose a generative framework that infers flexible, instance-specific  
 127 graphs from support images, enabling more adaptive and resilient structure modeling.

128 

### 2.2 LATENT STRUCTURE LEARNING FOR POSE ESTIMATION

129 Latent structure learning has been widely adopted to reason inter-keypoint dependencies in pose  
 130 estimation. Early approaches Wang et al. (2020); Hassan & Hamza (2023) leverage graph convolu-  
 131 tional networks to refine keypoint predictions by modeling predefined skeletal connections, which  
 132 restrict applicability to human- or hand-specific topologies. Generative models, particularly vari-  
 133 ational frameworks like the Variational Graph Autoencoder (VGA) Kipf & Welling (2016) and  
 134 CVAM-Pose Zhao et al. (2024) have demonstrated promise in capturing structural variability and  
 135 uncertainty, but typically remain tied to specific classes or predefined topologies. More recently,  
 136 V-VIPE Levy & Shrivastava (2024) leverages a variational autoencoder framework to learn a view-  
 137 invariant latent pose representation. ProPose Han et al. (2025) reformulates 3D human pose es-  
 138 timation as a probabilistic generative task by modeling instance-level pose distributions, enabling  
 139 uncertainty-aware and sample-efficient inference. Following these advancements, we explore a gen-  
 140 erative formulation to learn instance-level latent structures, aiming to enhance generalization and  
 141 move beyond the reliance on predefined priors. While learning latent structures improves flexibil-  
 142 ity, it remains insufficient under the CAPE setting, where support sets are sampled stochastically.  
 143 This insight suggests a mechanism that aggregates multiple latent graphs into a query-conditioned  
 144 structure, dynamically emphasizing support information most aligned with the query.

145 

## 3 METHOD

146 To effectively learn optimal keypoint structural dependencies and eliminate the adverse effects of in-  
 147 appropriate support sets, we propose a generative-based framework tailored for Category-Agnostic  
 148 Pose Estimation (CAPE). This method is grounded in GraphCape Hirschorn & Avidan (2024) frame-  
 149 work without reliance on skeleton priors. We begin by presenting a concise overview of the pipeline  
 150 before introducing our generative graph learning module and query-aware fusion technique.

151 

### 3.1 OVERALL PIPELINE

152 The goal of CAPE is to estimate the locations of semantic keypoints  $\hat{\mathbf{K}}_q \in \mathbb{R}^{M_c \times 2}$  for a query  
 153 image  $\mathbf{I}_q$ , given a small set of annotated exemplars from an unseen category, **where  $M_c$  denotes**  
 154 **the maximum possible number of keypoints.** In the  $N$ -shot setting, we are provided with a set of  
 155  $N$  support pairs  $\{(\mathbf{I}_i^s, \mathbf{K}_i^s)\}_{i=1}^N$ , where each support image  $\mathbf{I}_i^s$  is annotated with a set of keypoints  
 156  $\mathbf{K}_i^s \in \mathbb{R}^{M_c \times 2}$  for category  $c$  (which may vary in keypoint count  $M_c$ ). Initially, a shared backbone  
 157  $\phi(\cdot)$  extracts visual features  $F_q = \phi(\mathbf{I}_q)$  and  $F'_s = \phi(\mathbf{I}^s)$ . The support features are then aggregated

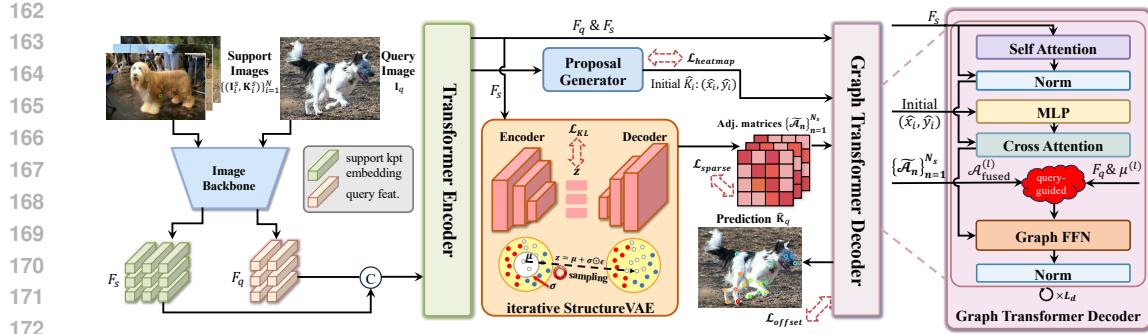


Figure 2: **Architecture Overview.** Our approach utilizes a pre-trained backbone to extract image features, which are refined by a transformer encoder through self-attention. A proposal generator is employed alongside a graph transformer decoder. Subsequently, we employ iterative StructureVAE to generate probabilistic adjacency matrices, and integrates them to a query-aware graph in the graph transformer decoder, improving localization accuracy by graph-oriented decoding.

with their corresponding keypoint targets to produce keypoint-aware embeddings  $F_s \in \mathbb{R}^{M \times D}$ , where  $M$  represents the maximum number of potential keypoints, and  $D$  is the embedding size. Then, a similarity-aware proposal generator computes correlations **between**  $F_s$  and  $F_q$ , yielding position proposals  $P \in \mathbb{R}^{M \times 2}$ . As shown in Figure 2, to model inter-keypoint dependencies and learn flexible skeleton knowledge, we introduce an *iterative Structure-aware Variational Autoencoder (i-SVAE)* that infers a latent adjacency matrix  $\mathcal{A} \in \mathbb{R}^{M \times M}$  conditioned on  $F_s$ . This probabilistic graph captures instance-specific keypoint relations and is fused with visual cues in the graph transformer decoder. We further mitigate visual uncertainty and reduce the adverse effects of improper support images through a *Compositional Graph Transfer (CGT)* strategy, which aggregates multiple latent hypotheses into a query-aware graph. This composition is injected into the graph transformer decoder to guide self- and cross-attention, progressively refining keypoint predictions.

### 3.2 ITERATIVE STRUCTURE LEARNING

In CAPE, the support and query mismatch in visibility, poses and topologies, making structural alignment essential. Our framework addresses this discrepancy jointly with *i-SVAE* that learns layer-wise, instance-specific graphs from support features, and *CGT* that adapts them to the query. Most CAPE methods make oversimplified assumptions: either modeling keypoints as independent entities or relying on manually defined priors. Such assumptions hinder the capture of topological consistency and pose variability across instances. To this end, we reformulate structure inference as a graph learning problem: keypoints are nodes, and their relationships are encoded in a latent adjacency matrix. Instead of using static, category-specific graphs, we propose an *i-SVAE* that learns and refines instance-specific keypoint graphs across decoding stages.

**Graph Formulation.** Let the graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  consist of  $M$  keypoint nodes  $v_i \in \mathcal{V}$  with initial node feature matrix  $\mathcal{X} \in \mathbb{R}^{M \times D}$ , and an initial noisy adjacency matrix  $\mathcal{A}^{(0)} \in \mathbb{R}^{M \times M}$  encoding edge set  $\mathcal{E}$ . Our goal is to learn a function  $f : \mathcal{F}_s \mapsto \mathcal{A}$  that maps support keypoint embeddings  $\mathcal{F}_s$  to a soft adjacency matrix  $\mathcal{A}$ , capturing latent inter-keypoint dependencies. We formulate this within the principled Variational Graph Autoencoder (VGAE) Kipf & Welling (2016), wherein the graph structure is treated as a latent variable:  $q_\phi(\mathbf{z} | \mathcal{F}_s)$ ,  $\mathbf{z} \in \mathbb{R}^{D_z}$ , with  $q_\phi$  denoting the approximate posterior and  $\mathbf{z}$  the latent code. The code is decoded into a soft adjacency matrix that models probabilistic keypoint connectivity.

**Iterative StructureVAE.** The i-SVAE consists of two components: a probabilistic encoder that parameterizes a latent graph distribution, and a decoder that constructs the adjacency matrix from this latent space. A detailed figure is shown in the supplementary material. At each layer  $l$ , given support node embeddings  $\mathcal{F}_s^{(l)} \in \mathbb{R}^{M \times D}$  from the previous graph transformer decoder layer, we first perform variational inference to estimate the latent structural embeddings:

$$[\boldsymbol{\mu}^{(l)}, \log(\boldsymbol{\sigma}^{(l)})] = \text{Enc}(\mathcal{F}_s^{(l)}), \quad \boldsymbol{\mu}^{(l)}, \log \boldsymbol{\sigma}^{(l)} \in \mathbb{R}^{D_z}, \quad (1)$$

$$q_\phi(\mathbf{z}^{(l)} | \mathcal{F}_s^{(l)}) = \mathcal{N}(\mathbf{z}^{(l)}; \boldsymbol{\mu}^{(l)}, \text{diag}(\boldsymbol{\sigma}^{(l)})),$$

216 where  $\text{Enc}$  denotes the graph encoder that produces the approximate posterior distribution over the  
 217 latent graph codes  $\mathbf{z}$ . We then employ the reparameterization trick to sample  $\mathbf{z}^{(l)} \in \mathbb{R}^{D_z}$ :  
 218

$$219 \quad \mathbf{z}^{(l)} = \boldsymbol{\mu}^{(l)} + \boldsymbol{\sigma}^{(l)} \odot \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I}), \quad (2)$$

220 Next, the latent code  $\mathbf{z}$  is passed through a fully connected decoder to construct the adjacency matrix  
 221  $\hat{\mathcal{A}} = \text{Dec}(\mathbf{z}) \in \mathbb{R}^{M \times M}$ , where  $\text{Dec}$  represents the graph decoder. To ensure the undirectionality and  
 222 interpretability of the adjacency matrix, we symmetrize and normalize it row-wise:  
 223

$$224 \quad \hat{\mathcal{A}}^{(l)}_{\text{sym}} = \frac{1}{2}(\hat{\mathcal{A}}^{(l)} + \hat{\mathcal{A}}^{(l)\top}), \quad \tilde{\mathcal{A}}^{(l)} = \text{norm}(\hat{\mathcal{A}}^{(l)}_{\text{sym}}), \quad (3)$$

226 Then, we take the *CGT* strategy to fuse multiple sampled adjacency matrices into a unified, query-  
 227 aware graph, which is introduced in the next subsection. This fusion strategy is performed via  
 228 Bayesian averaging, reweighted by graph-level uncertainty and relevance to the query features:  
 229

$$230 \quad \tilde{\mathcal{A}}_{\text{final}}^{(l)} = \text{CGT}(\{\tilde{\mathcal{A}}_n^{(l)}\}, \{\boldsymbol{\mu}_n^{(l)}\}, F_q). \quad (4)$$

232 To ensure meaningful latent representations and regulate structural uncertainty, we minimize the  
 233 Kullback-Leibler divergence between the approximate posterior  $q_\phi(\mathbf{z} \mid \mathbf{X})$  and a Gaussian prior  
 234  $p(\mathbf{z}) = \mathcal{N}(0, \mathbf{I})$ . Besides, we impose an  $\ell_2$  sparsity constraint on the learned adjacency to encourage  
 235 minimal and interpretable connectivity. The total *i-SVAE* loss at the  $l$ -th decoder layer is defined as:  
 236

$$237 \quad \mathcal{L}_{\text{VAE}}^{(l)} = \mathcal{L}_{KL}^{(l)} + \beta \cdot \mathcal{L}_{\text{sparse}}^{(l)} = \underbrace{D_{\text{KL}}\left[q_\phi(\mathbf{z}^{(l)} \mid F_s^{(l)}) \parallel p(\mathbf{z}^{(l)})\right]}_{\text{Prior Regularization}} + \underbrace{\frac{\beta}{M^2} \lambda \|\tilde{\mathcal{A}}_{\text{final}}^{(l)}\|_F^2}_{\text{Sparse Matrix}}, \quad (5)$$

239 where the hyper-parameter  $\beta = 0.1$ . To leverage the learned structural priors, we follow the Graph-  
 240 Cape and incorporate a graph convolutional layer conditioned on the final matrix  $\tilde{\mathcal{A}}_{\text{final}}^{(l)}$ .  
 241

$$243 \quad F_s^{(l+1)} = \sigma\left(W_{\text{adj}} F_s^{(l)} \tilde{\mathcal{A}}_{\text{final}}^{(l)} + W_{\text{self}} F_s^{(l)}\right), \quad (6)$$

245 where  $W_{\text{adj}}, W_{\text{self}} \in \mathbb{R}^{D_{\text{out}} \times D}$  are learnable weights, and  $\sigma(\cdot)$  denotes ReLU activation. The first  
 246 term aggregates features from semantically or spatially connected neighbors, while the second term  
 247 retains individual node semantics via self-transformation.

248 The final skeleton  $\tilde{\mathcal{A}}_{\text{final}}^{(l)}$  serves as the structural guidance for message passing:  
 249

$$250 \quad F_s^{(l+1)} = \text{GCN}(F_s^{(l)}, \tilde{\mathcal{A}}_{\text{final}}^{(l)}), \quad P_q^{(l+1)} = \sigma\left(\sigma^{-1}(P_q^{(l)}) + \text{MLP}(F_s^{(l+1)})\right), \quad (7)$$

252 where  $P_q^{(l)} \in \mathbb{R}^{K \times 2}$  is the predicted keypoint locations at  $l$ -th layer used for intermediate supervision,  
 253 with the output from the final layer as the final keypoint prediction. And  $\sigma$  and  $\sigma^{-1}$  are the  
 254 sigmoid and its inverse function.

255 By embedding *i-SVAE* within each decoder layer, our method performs iterative structural refinement,  
 256 progressively updating latent pose graphs in response to evolving visual semantics and local-  
 257 ization cues. This layer-wise iterative design enables the model to capture diverse structural patterns  
 258 and encode high-order keypoint dependencies, thereby strengthening relational reasoning and im-  
 259 proving generalization to novel categories. See Appendix A.1 for further details.  
 260

### 261 3.3 COMPOSITIONAL GRAPH TRANSFER

263 While *i-SVAE* enables layer-wise modeling of instance-specific pose graphs, its stochastic sampling  
 264 process introduces uncertainty across multiple latent graphs. To this end, we propose *Compositional*  
 265 *Graph Transfer (CGT)*, a query-aware graph fusion mechanism that aggregates multiple sampled  
 266 adjacency matrices into a robust and expressive structural representation. Specifically, given a set  
 267 of  $N_s$  latent graphs sampled from *i-SVAE* at the  $l$ -th decoder layer, denoted as  $\{\tilde{\mathcal{A}}_n^{(l)}\}_{n=1}^{N_s}$ , each  
 268 associated with a latent distribution parameterized by  $(\boldsymbol{\mu}^{(l)}, \boldsymbol{\sigma}^{(l)})$ , our goal is to construct a matrix  
 269  $\mathcal{A}_{\text{final}} \in \mathbb{R}^{M \times M}$  that best reflects the robust structure relationships among keypoints conditioned on  
 the query context, while simultaneously alleviating over-reliance on the support-driven guidance.

270 To achieve this, we adopt a Bayesian confidence-weighted aggregation strategy. Firstly, we define  
 271 the confidence of each sampled graph as the inverse of the total variance:  
 272

$$273 \quad w_n = \frac{1}{\sum_{i=1}^{D_z} \sigma_{n,i}^{(l)} + \epsilon}, \quad \tilde{w}_n = \frac{w_n}{\sum_{m=1}^{N_s} w_m}, \quad (8)$$

$$274$$

$$275$$

276 where  $\epsilon = 1e^{-6}$  is a small constant to ensure numerical stability. These normalized weights  $\tilde{w}_n$  reflect  
 277 the epistemic uncertainty of each latent sample and serve to guide the fusion process. The fused  
 278 adjacency is then computed via a weighted average:  $\tilde{\mathcal{A}}_{\text{fused}}^{(l)} = \sum_{n=1}^{N_s} \tilde{w}_n \cdot \tilde{\mathcal{A}}_n^{(l)}$ . To further align the  
 279 fuse structure with query-specific evidence, we incorporate query-guided gating. Let  $F_q \in \mathbb{R}^{hw \times D}$   
 280 denote the query feature, where  $[h, w]$  denotes the patch size in image backbone. We compute  
 281 attention-based gating scores  $\alpha^{(l)}$  by comparing global query descriptors with each means  $\mu^{(l)}$ :  
 282

$$283 \quad \alpha^{(l)} = \frac{\text{sim}(\text{Pool}(F_q), \mu^{(l)})}{\sum_{l=1}^{L_d} \text{sim}(\text{Pool}(F_q), \mu^{(l)})}, \quad (9)$$

$$284$$

$$285$$

286 where  $\text{sim}(\cdot, \cdot)$  denotes cosine similarity and Pool is a global average pooling operator. The final  
 287 fused graph becomes  $\tilde{\mathcal{A}}_{\text{final}}^{(l)} = \sum_{l=1}^L \alpha^{(l)} \cdot \tilde{\mathcal{A}}_{\text{fused}}^{(l)}$ , where  $L \in [1, L_d]$  is the current decoder layer. The  
 288 fusion process enhances robustness against sampling stochasticity and grounds structural reasoning  
 289 in the visual context of the query. The resulting graph  $\tilde{\mathcal{A}}_{\text{final}}^{(l)}$  is propagated into the GCN layer,  
 290 enabling structure-aware refinement of keypoint predictions. See Appendix A.1 for CGT details.  
 291

### 292 3.4 TRAINING AND INFERENCE

$$293$$

294 For the category-agnostic pose estimation task, we employ the commonly used loss Shi et al. (2023);  
 295 Hirschorn & Avidan (2024); Rusanovsky et al. (2025)  $\mathcal{L}_{\text{pred}}$ :

$$296 \quad \mathcal{L}_{\text{pred}} = \lambda_{\text{heatmap}} \cdot \mathcal{L}_{\text{heatmap}} + \mathcal{L}_{\text{offset}},$$

$$297$$

$$298 \quad \mathcal{L}_{\text{heatmap}} = \frac{1}{M_c \cdot H \cdot W} \sum_{i=1}^{M_c} \left\| \hat{\mathbf{H}}_i - \mathbf{H}_i \right\|, \quad \mathcal{L}_{\text{offset}} = \frac{1}{L_d} \sum_{l=1}^{L_d} \sum_{i=1}^{M_c} \left| \hat{\mathbf{K}}_i^l - \mathbf{K}_i \right|, \quad (10)$$

$$299$$

$$300$$

301 where  $\hat{\mathbf{H}}_i$  denotes the output similarity heatmap and  $\mathbf{H}_i$  is the ground-truth heatmap.  $\hat{\mathbf{K}}_i^l$  is the  
 302 output keypoint location from the Graph Transformer layer  $l$  and  $\mathbf{K}_i$  is the ground-truth location.  
 303 By our framework 3 and internal modules 3.2, our overall training objective is as follow:  
 304

$$305 \quad \mathcal{L} = \mathcal{L}_{\text{pred}} + \gamma \cdot \mathcal{L}_{\text{VAE}}, \quad (11)$$

$$306$$

307 where  $\gamma = 1e^{-3}$  is the hyper-parameter. During inference, our model uses the final layer output  $\hat{\mathbf{K}}_i^{L_d}$   
 308 as the predicted location. Within the i-SVAE variational inference process, the latent code  $z$  is equal  
 309 to the mean  $\mu$  of the approximate posterior, effectively collapsing the stochastic sampling. This  
 310 deterministic substitution ensures stable and consistent structural priors, aligning with the learned  
 311 feature distribution while eliminating inference-time uncertainty.  
 312

## 313 4 EXPERIMENTS

$$314$$

### 315 4.1 IMPLEMENTATION DETAILS

$$316$$

317 We train and evaluate our method on a machine with an NVIDIA A100 GPU with 40 GB of memory.  
 318 The architecture is implemented within the MMpose framework Contributors (2020). To ensure a  
 319 fair comparison, the configuration settings remain consistent with GraphCape Hirschorn & Avidan  
 320 (2024) and CapeFormer Xu et al. (2022a). During training, we use  $256 \times 256$  input images and  
 321 apply data augmentation including random scaling in the range  $([-0.15, 0.15])$  and random rotation  
 322 within  $([-15^\circ, 15^\circ])$  for fair comparisons. All models are trained for 200 epochs with a step-wise  
 323 learning rate scheduler that decreases by a factor of 10 at the 160th and 180th epochs. We use Adam  
 324 optimizer to train the model for 200 epochs with a batch size of 16. See Section A.2 for details.

324  
 325 **Table 1: Comparisons on MP-100:** PCK@0.2 performance under the 1-shot setting. GenCape  
 326 achieves the best average performance on the average of all splits, outperforming state-of-the-art  
 327 methods under all three support sets types.

| Type          | Method                                | Support              | Split 1      | Split 2      | Split 3      | Split 4      | Split 5      | Avg.         |
|---------------|---------------------------------------|----------------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Image-support | POMNet Xu et al. (2022a)              | Image                | 84.23        | 78.25        | 78.17        | 78.68        | 79.17        | 79.70        |
|               | CapeFormer Shi et al. (2023)          | Image                | 89.45        | 84.88        | 83.59        | 83.53        | 85.09        | 85.31        |
|               | ESCAPE Nguyen et al. (2024)           | Image                | 86.89        | 82.55        | 81.25        | 81.72        | 81.32        | 82.74        |
|               | MetaPoint+ Chen et al. (2024)         | Image                | 90.43        | 85.59        | 84.52        | 84.34        | 85.96        | 86.17        |
|               | CapeFormer-T Shi et al. (2023)        | Image                | 89.48        | 86.69        | 85.31        | 84.79        | 84.97        | 86.25        |
|               | SDPNet (HRNet-32) Ren et al. (2024)   | Image                | 91.54        | 86.72        | 85.49        | 85.77        | 87.26        | 87.36        |
| Text-support  | SCAPE Liang et al. (2024)             | Image                | 91.67        | 86.87        | 87.29        | 85.01        | 86.92        | 87.55        |
|               | CLAMP Zhang et al. (2023)             | Text                 | 72.37        | -            | -            | -            | -            | -            |
|               | X-Pose Yang et al. (2024)             | Image\Text           | 89.07        | 85.05        | 85.26        | 85.52        | 85.79        | 86.14        |
|               | PPM+CPT Peng et al. (2024)            | Image + Text         | 91.03        | 88.06        | 84.48        | 86.73        | 87.40        | 87.54        |
| Graph-support | CapeX-S Rusanovsky et al. (2025)      | Image + Text + Graph | 95.17        | 88.88        | 87.72        | 88.24        | 91.81        | 90.37        |
|               | GraphCape-T Hirschorn & Avidan (2024) | Image + Graph        | 91.19        | 87.81        | 85.68        | 85.87        | 85.61        | 87.23        |
|               | <b>GenCape-T (Ours)</b>               | Image (Graph)        | <b>92.05</b> | <b>88.69</b> | <b>86.89</b> | <b>85.88</b> | <b>87.02</b> | <b>88.09</b> |
|               | GraphCape-S Hirschorn & Avidan (2024) | Image + Graph        | 94.73        | 89.79        | <b>90.69</b> | 88.09        | 90.11        | 90.68        |
|               | <b>GenCape-S (Ours)</b>               | Image (Graph)        | <b>95.23</b> | <b>90.60</b> | 89.46        | <b>89.32</b> | <b>90.43</b> | <b>91.01</b> |

341  
 342 **Table 2: Performance comparisons** under 5-shot  
 343 setting with SwinV2-small as image backbone.

344 **Table 3: Performance comparison** of  
 345 CAPE methods under stricter thresholds.

| Method                | Split 1      | Split 2      | Split 3      | Split 4      | Split 5      | Avg.         |
|-----------------------|--------------|--------------|--------------|--------------|--------------|--------------|
| Fine-tune             | 71.67        | 57.84        | 66.76        | 66.53        | 60.24        | 64.61        |
| POMNet                | 84.72        | 79.61        | 78.00        | 80.38        | 80.85        | 80.71        |
| CapeFormer            | 91.94        | 88.92        | 89.40        | 88.01        | 88.25        | 89.30        |
| SDPNet                | 93.68        | 90.23        | 89.67        | 89.08        | 89.46        | 90.42        |
| PPM+CPT               | 93.64        | 92.71        | 91.76        | 92.85        | 91.94        | 92.58        |
| SCAPE                 | 95.18        | 91.25        | 91.78        | 90.74        | 91.10        | 92.01        |
| GraphCape             | 96.67        | 91.48        | <b>92.62</b> | 90.95        | 92.41        | 92.83        |
| <b>GenCape (Ours)</b> | <b>97.19</b> | <b>92.94</b> | 92.26        | <b>91.93</b> | <b>93.34</b> | <b>93.53</b> |

## 353 4.2 DATASET AND METRIC

354  
 355 We train and evaluate our method on the MP-100 Xu et al. (2022a) dataset, which is currently the  
 356 only public dataset for CAPE tasks. MP-100 contains 100 sub-categories and 8 supercategories,  
 357 with a total of 18K images and 20K annotations. The number of annotated keypoints covers a wide  
 358 range, from 8 to 68. To ensure and identify performance stability on unseen categories, following  
 359 previous methods Xu et al. (2022a), the dataset is divided into 5 splits. In each split, these categories  
 360 are split into train, validation, and test sets in a 70/10/20 ratio without any category overlap. We use  
 361 the Probability of Correct Keypoint (PCK) as the evaluation metric. We follow the standard metric  
 362 PCK@0.2 as the default metric for performance reporting. And we further evaluate model perfor-  
 363 mance under stricter threshold conditions ([0.05, 0.1, 0.15, 0.2]) for comprehensive comparisons.

## 364 4.3 BENCHMARK RESULTS

365  
 366 We conduct a comparative analysis of our approach with SwinV2 Liu et al. (2022) as backbone,  
 367 against the **graph-support method**: GraphCape Hirschorn & Avidan (2024) as our baseline; **image-**  
 368 **support methods**: POMNet Xu et al. (2022a), CapeFormer Shi et al. (2023), ESCAPE Nguyen  
 369 et al. (2024), MetaPoint Chen et al. (2024), SDPNet Ren et al. (2024), FMMP Chen et al. (2025a);  
 370 **text-support methods**: CLAMP Zhang et al. (2023), XPose Yang et al. (2024), PPM+CPT Peng  
 371 et al. (2024), CapeX Rusanovsky et al. (2025). Our evaluation is conducted on the MP-100 dataset,  
 372 considering the 1- and 5-shot settings. We denote our models as GenCape-T and GenCape-S for ab-  
 373 brevity, corresponding to employing SwinV2-tiny and small as the image backbone, respectively.

374 **1-shot results.** As reported in Table 1, GenCape-T achieves an average PCK of 88.09%, surpassing  
 375 the strong graph-based GraphCape-T baseline (87.23%) by +0.86%. Remarkably, without relying  
 376 on class-level text, our method still outperforms multimodal CAPE models such as XPose (86.14%)  
 377 and PPM+CPT (87.54%), indicating that the learned structure-aware representation serves as an  
 effective surrogate for external semantic cues. Moreover, GenCape-S attains 91.01% PCK, exceed-

378 ing CapeX-S (90.37%), which leverages both textual and skeleton support. **5-shot results.** Under  
 379 the 5-shot setting (Table 2), GenCape-S further improves, achieving an average PCK of 93.53% and  
 380 outperforming all representative CAPE methods, including PPM+CPT (92.58%), SCAPE (91.01%),  
 381 and GraphCape (92.83%). The model delivers consistent gains across all splits except Split 3, with  
 382 a maximum of 97.19% on Split 1 and a minimum of 91.93% on Split 4.

383 **More detailed comparisons.** Table 3 presents results under stricter thresholds (0.2, 0.15, 0.1, 0.05)  
 384 on Split-1 with ResNet-50 as backbone. Threshold choice strongly affects relative gaps: GenCape-  
 385 R50 outperforms FMMP by 0.92% at PCK@0.2, and the margin increases to 1.61% at PCK@0.05,  
 386 showing that coarse thresholds may mask fine-grained discriminative ability. GenCape achieves the  
 387 highest accuracy at all thresholds and improves mPCK by +1.59% over FMMP.

#### 389 4.4 ABLATION STUDY

390 In this section, we conduct all of the ablation studies of our proposed method on the MP-100 dataset  
 391 Split-1 using the SwinV2-Tiny backbone, unless otherwise specified. We now present key ablation  
 392 experiments. Additional ablations can be found in Appendix B.

393 **Effects of Different Components.** To rigorously quant-  
 394 ify the contribution of each module in our framework,  
 395 we conduct a comprehensive ablation study under 1-shot  
 396 setting. We isolate and examine the effects of variational  
 397 regularization ( $\mathcal{L}_{KL}$ ), sparsity penalization ( $\mathcal{L}_{sparse}$ ), and  
 398 the CGT module. Table 4 reports the results. Starting  
 399 from the baseline that excludes all components, we ob-  
 400 serve a base PCK of 91.19%. Introducing the KL diver-  
 401 gence term  $\mathcal{L}_{KL}$  yields a modest improvement (+0.24%),  
 402 suggesting its stabilizing role by constraining posterior distributions with the Gaussian prior. Com-  
 403 bining both constraints, the performance improves substantially to 91.75%, confirming their syn-  
 404 ergistic effect in enforcing informative and interpretable structural cues. Further incorporating the  
 405 CGT mechanism yields 0.86% gains, underscoring the importance of compositional graph fusion in  
 406 consolidating uncertainty-aware structural hypotheses and enhancing query-specific robustness.

407 **Effects of Hyper-parameter.** We in-  
 408 vestigate the influence of four key hy-  
 409 perparameters: latent dimensionality  
 410  $D_z$ , sample count  $N_s$ , and loss weights  
 411  $\beta$  and  $\gamma$  (Eq.5, Eq.11). As shown  
 412 in Table 5, the model achieves opti-  
 413 mal performance at  $D_z = 32$ , with  
 414 larger dimensions leading to per-  
 415 formance degradation (e.g., 92.05% at 32  
 416 vs. lower at 64/128), indicating that  
 417 excessive latent capacity introduces re-  
 418 dundancy and weakens structural com-  
 419 pactness. Fixing  $D_z = 32$ , we vary  $N_s$   
 420 and observe the best results at  $N_s = 3$ ,  
 421 while both smaller ( $N_s = 2$ ) and larger  
 422 ( $N_s = 5$ ) values slightly reduce accu-  
 423 racy. This suggests a trade-off between uncertainty modeling and variance over-smoothing. For loss  
 424 weights, we separately tune  $\gamma$  (KL loss) and  $\beta$  (sparsity loss). The model peaks at  $\gamma = 10^{-3}$  when  $\beta$   
 425 is fixed, balancing reconstruction and regularization. Likewise,  $\beta = 0.1$  yields the best PCK, high-  
 426 lighting its role in regulating structure sparsity without over-penalizing connectivity. These results  
 427 confirm that our framework remains robust across hyperparameter variations, with optimal settings  
 428 jointly promoting expressiveness and structural regularity.

429 **Effects of Cross-Category Generalization.** To assess structural generalization beyond cat-  
 430 egory boundaries, we conduct three cross-supercategory pairs experiments: Person↔Felidae,  
 431 Felidae↔Ursidae, and Person↔AnimalFace. These settings cover diverse appearance and topology  
 432 gaps, including upright-to-quadruped transfers and full-body to face shifts. As shown in Table 6,  
 433 GenCape consistently outperforms GraphCape across all pairs, with margins up to +11.8 points (*Fe-*

Table 4: **Ablation studies** on different components (Split-1).

| $\mathcal{L}_{KL}$ | $\mathcal{L}_{sparse}$ | CGT | PCK          | $\Delta$     |
|--------------------|------------------------|-----|--------------|--------------|
|                    |                        |     | 91.19        | 0            |
| ✓                  |                        |     | 91.43        | +0.24        |
| ✓                  | ✓                      |     | 91.75        | +0.56        |
| ✓                  | ✓                      | ✓   | <b>92.05</b> | <b>+0.86</b> |

Table 5: **Ablation studies** on hyper-parameters, includ-  
 ing latent dimension  $D_z$ , number of posterior samples  $N_s$ ,  
 and weighting factors  $\beta$  and  $\gamma$  for the training objectives.

| Hyper-parameters of i-SVAE            |              |              |           |              |  |
|---------------------------------------|--------------|--------------|-----------|--------------|--|
| Parameter $D_z$ ( $N_s = 3$ )         | 32           | 64           | 96        | 128          |  |
| PCK                                   | <b>92.05</b> | 91.89        | 91.54     | 91.40        |  |
| Parameter $N_s$ ( $D_z = 32$ )        | 2            | 3            | 4         | 5            |  |
| PCK                                   | 91.60        | <b>92.05</b> | 91.47     | 91.63        |  |
| Objective weighting factor            |              |              |           |              |  |
| Parameter $\beta$ ( $\gamma = 1e-3$ ) | 1            | $1e^{-1}$    | $1e^{-2}$ | $1e^{-3}$    |  |
| PCK                                   | 91.43        | <b>92.05</b> | 91.65     | 91.74        |  |
| Parameter $\gamma$ ( $\beta = 0.1$ )  | 1            | $1e^{-1}$    | $1e^{-2}$ | $1e^{-3}$    |  |
| PCK                                   | 90.99        | 91.14        | 91.71     | <b>92.05</b> |  |

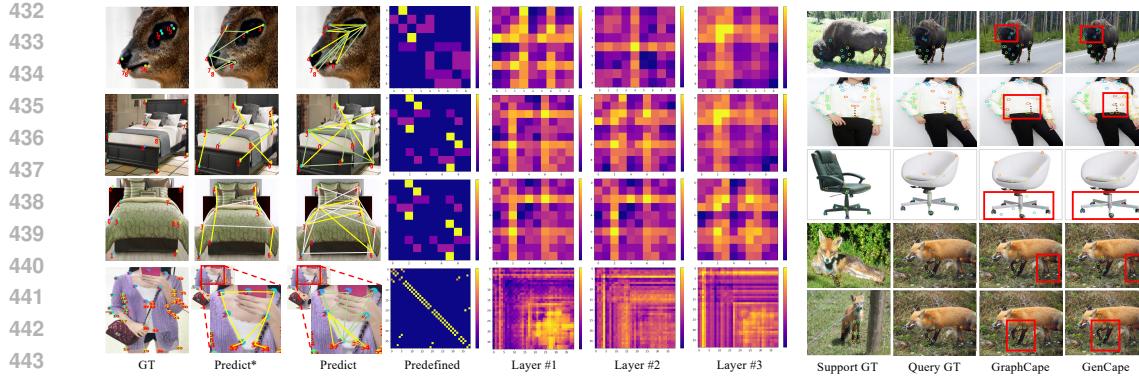


Figure 3: Comparisons on adjacency matrices inferred by i-SVAE and predefined graph. The Predict\* is the predicted locations with the prior connections, while the Predict is with learned connections.

Figure 4: Comparisons of qualitative visualization. The red boxes highlight the differences.

*lidae* → *Person*), demonstrating the robustness of compositional latent graphs under structural shifts. We further evaluate in the challenging cross-species setting, where models trained on human bodies and tested on morphologically distinct animal bodies. GraphCape suffers substantial performance drops, e.g., 31.09 on *Rabbit Body*, due to the rigidity of static priors. In contrast, GenCape maintains strong transferability, achieving +21.05 improvement and scaling effectively across species with varying skeletons on *Squirrel Body*. These results validate that generative, instance-specific graphs better capture structural uncertainty and enable robust pose reasoning across categories.

**Qualitative Analysis.** To intuitively understand how our model mines effective skeletons, we visualize the pre-normalized adjacency matrices across decoder layers with the inferred skeletons in Figure 3. Compared with fixed skeletons, whose manually chosen edges are not guaranteed to be semantically correct and can mislead message passing, the graph inferred by i-SVAE are query-conditioned and instance-specific. So pose, viewpoint and other changes are directly reflected in the adjacency matrix. The learned structures are clearly task-driven, emphasizing high-influence keypoints that most contribute to accurate localization, *e.g.*, nose and eyes

as anchors for klipspringer face, four corners for bed, and central torso for long sleeved outwear. They receive denser and stronger connections and thus provide more effective geometric constraints. Figure 4 highlights GenCape’s superior adaptability when exposed to varying support-query pairs in 1-shot setting. Compared to GraphCape, which fails under pose misalignment and occlusion (*e.g.*, bison and swivelchair), our model consistently localizes keypoints accurately by leveraging its uncertainty-aware graph. Interestingly, we further visualize the impact of varying support images for the same query instance, as shown in the fourth and fifth rows. When the fox undergoes different forms of occlusion, GraphCape suffers from inconsistent errors, while GenCape maintains stable predictions. These results suggest that our method is significantly less susceptible to the adverse effects introduced by suboptimal support sets. Additional analyzes are provided in Appendix C.

## 5 CONCLUSION

In this paper, we introduce GenCape, a generative framework for CAPE that infers keypoint relationships solely from visual inputs. GenCape integrates an iterative Structure-aware Variational Autoencoder to progressively infer instance-specific keypoint relationships, alongside a Compositional Graph Transfer module that aggregates multiple latent graph hypotheses into query-aware structural cues. Extensive experiments on MP-100 demonstrate that GenCape achieves the state-of-the-art performance. These results demonstrate the effectiveness of our framework.

486 REFERENCES  
487

488 Junjie Chen, Jiebin Yan, Yuming Fang, and Li Niu. Meta-point learning and refining for category-  
489 agnostic pose estimation. In *Conference on Computer Vision and Pattern Recognition*, pp. 23534–  
490 23543, 2024.

491 Junjie Chen, Weilong Chen, Yifan Zuo, and Yuming Fang. Recurrent feature mining and keypoint  
492 mixup padding for category-agnostic pose estimation. In *Conference on Computer Vision and*  
493 *Pattern Recognition*, pp. 22035–22044, 2025a.

494 Junjie Chen, Zeyu Luo, Zeheng Liu, Wenhui Jiang, Li Niu, and Yuming Fang. Weak-shot keypoint  
495 estimation via keyness and correspondence transfer. In *The Thirty-ninth Annual Conference on*  
496 *Neural Information Processing Systems*, 2025b.

497

498 MMPose Contributors. Openmmlab pose estimation toolbox and benchmark. <https://github.com/open-mmlab/mmpose>, 2020.

499

500 Jumin Han, Jun-Hee Kim, and Seong-Whan Lee. Propose: Probabilistic 3d human pose estimation  
501 with instance-level distribution and normalizing flow. In *Proceedings of the AAAI Conference on*  
502 *Artificial Intelligence*, volume 39, pp. 3338–3346, 2025.

503

504 Md Tanvir Hassan and A Ben Hamza. Regular splitting graph network for 3d human pose estimation.  
505 *IEEE Transactions on Image Processing*, 32:4212–4222, 2023.

506

507 Or Hirschorn and Shai Avidan. A graph-based approach for category-agnostic pose estimation. In  
508 *European Conference on Computer Vision*, pp. 469–485. Springer, 2024.

509

510 Junho Kim, Hyungjin Chung, and Byung-Hoon Kim. Capellm: Support-free category-agnostic pose  
511 estimation with multimodal large language models. *arXiv preprint arXiv:2411.06869*, 2024.

511

512 Diederick P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International*  
513 *Conference on Learning Representations*, 2015.

514

515 Thomas N Kipf and Max Welling. Variational graph auto-encoders. *NIPS Workshop on Bayesian*  
516 *Deep Learning*, 2016.

517

518 Mara Levy and Abhinav Shrivastava. V-vipe: Variational view invariant pose embedding. In *Con-*  
519 *ference on Computer Vision and Pattern Recognition*, pp. 1633–1642, 2024.

520

521 Yujia Liang, Zixuan Ye, Wenze Liu, and Hao Lu. Scape: A simple and strong category-agnostic  
522 pose estimator. In *European Conference on Computer Vision*, pp. 478–494. Springer, 2024.

523

524 Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng  
525 Zhang, Li Dong, et al. Swin transformer v2: Scaling up capacity and resolution. In *Conference*  
526 *on Computer Vision and Pattern Recognition*, pp. 12009–12019, 2022.

527

528 Changsheng Lu, Zheyuan Liu, and Piotr Koniusz. Openkd: Opening prompt diversity for zero-  
529 and few-shot keypoint detection. In *European Conference on Computer Vision*, pp. 148–165.  
530 Springer, 2024.

531

532 Xiao Ma, Sumit Patidar, Iain Haughton, and Stephen James. Hierarchical diffusion policy for  
533 kinematics-aware multi-task robotic manipulation. In *Conference on Computer Vision and Pattern*  
534 *Recognition*, pp. 18081–18090, June 2024.

535

536 Khoi Duc Nguyen, Chen Li, and Gim Hee Lee. Escape: Encoding super-keypoints for category-  
537 agnostic pose estimation. In *Conference on Computer Vision and Pattern Recognition*, pp. 23491–  
23500, 2024.

538

539 Duo Peng, Zhengbo Zhang, Ping Hu, QiuHong Ke, David KY Yau, and Jun Liu. Harnessing text-  
540 to-image diffusion models for category-agnostic pose estimation. In *European Conference on*  
541 *Computer Vision*, pp. 342–360, 2024.

542

543 Jiyong Rao, Brian Nlong Zhao, and Yu Wang. Probabilistic prompt distribution learning for animal  
544 pose estimation. In *Conference on Computer Vision and Pattern Recognition*, pp. 29438–29447,  
545 2025.

540 Pengfei Ren, Yuanyuan Gao, Haifeng Sun, Qi Qi, Jingyu Wang, and Jianxin Liao. Dynamic support  
 541 information mining for category-agnostic pose estimation. In *Conference on Computer Vision*  
 542 and *Pattern Recognition*, pp. 1921–1930, 2024.

543 Matan Rusanovsky, Or Hirschorn, and Shai Avidan. Capex: Category-agnostic pose estimation from  
 544 textual point explanation. In *International Conference on Learning Representations*, 2025.

545 Min Shi, Zihao Huang, Xianzheng Ma, Xiaowei Hu, and Zhiguo Cao. Matching is not enough: A  
 546 two-stage framework for category-agnostic pose estimation. In *Conference on Computer Vision*  
 547 and *Pattern Recognition*, pp. 7308–7317, 2023.

548 Lucas Stoffl, Andy Bonnetto, Stéphane d’Ascoli, and Alexander Mathis. Elucidating the hierarchical  
 549 nature of behavior with masked autoencoders. In *European conference on computer vision*, pp.  
 550 106–125. Springer, 2024.

551 Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning  
 552 for human pose estimation. In *Conference on Computer Vision and Pattern Recognition*, pp.  
 553 5693–5703, 2019.

554 Jian Wang, Xiang Long, Yuan Gao, Errui Ding, and Shilei Wen. Graph-pcnn: Two stage human  
 555 pose estimation with graph pose refinement. In *European Conference on Computer Vision*, pp.  
 556 492–508. Springer, 2020.

557 Lumin Xu, Sheng Jin, Wang Zeng, Wentao Liu, Chen Qian, Wanli Ouyang, Ping Luo, and Xiaogang  
 558 Wang. Pose for everything: Towards category-agnostic pose estimation. In *European conference*  
 559 on computer vision, pp. 398–416. Springer, 2022a.

560 Tianyang Xu, Jiyong Rao, Xiaoning Song, Zhenhua Feng, and Xiao-Jun Wu. Learning structure-  
 561 supporting dependencies via keypoint interactive transformer for general mammal pose estima-  
 562 tion. *International Journal of Computer Vision*, pp. 1–19, 2025.

563 Yufei Xu, Jing Zhang, Qiming Zhang, and Dacheng Tao. Vitpose: Simple vision transformer base-  
 564 lines for human pose estimation. In *Advances in neural information processing systems*, vol-  
 565 ume 35, pp. 38571–38584, 2022b.

566 Jie Yang, Ailing Zeng, Ruimao Zhang, and Lei Zhang. X-pose: Detecting any keypoints. In *Euro-  
 567 pean Conference on Computer Vision*, pp. 249–268, 2024.

568 Sen Yang, Wen Heng, Gang Liu, GUOZHONG LUO, Wankou Yang, and Gang YU. Capturing  
 569 the motion of every joint: 3d human pose and shape estimation with independent tokens. In  
 570 *International Conference on Learning Representations*, 2023.

571 Shaokai Ye, Anastasiia Filippova, Jessy Lauer, Steffen Schneider, Maxime Vidal, Tian Qiu, Alexan-  
 572 der Mathis, and Mackenzie Weygandt Mathis. Superanimal pretrained pose estimation models  
 573 for behavioral analysis. *Nature Communications*, 15(1):5165, 2024.

574 Hang Yu, Yufei Xu, Jing Zhang, Wei Zhao, Ziyu Guan, and Dacheng Tao. AP-10k: A benchmark  
 575 for animal pose estimation in the wild. In *Advance Neural Information Processing System*, 2021.

576 Yuhui Yuan, Rao Fu, Lang Huang, Weihong Lin, Chao Zhang, Xilin Chen, and Jingdong Wang.  
 577 Hrformer: High-resolution vision transformer for dense predict. In *Advances in neural informa-  
 578 tion processing systems*, volume 34, pp. 7281–7293, 2021.

579 Xu Zhang, Wen Wang, Zhe Chen, Yufei Xu, Jing Zhang, and Dacheng Tao. Clamp: Prompt-based  
 580 contrastive learning for connecting language and animal pose. In *Conference on Computer Vision*  
 581 and *Pattern Recognition*, pp. 23272–23281, 2023.

582 Jianyu Zhao, Wei Quan, and Bogdan J Matuszewski. Cvam-pose: Conditional variational autoen-  
 583 coder for multi-object monocular pose estimation. In *British Machine Vision Conference*. BMVA,  
 584 2024.

585 Ce Zheng, Wenhan Wu, Chen Chen, Taojiaannan Yang, Sijie Zhu, Ju Shen, Nasser Kehtarnavaz, and  
 586 Mubarak Shah. Deep learning-based human pose estimation: A survey. *ACM computing surveys*,  
 587 56(1):1–37, 2023.

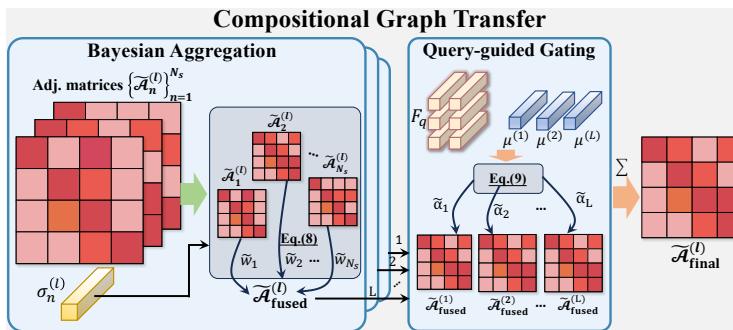
594 Ying Zheng, Lei Yao, Yuejiao Su, Yi Zhang, Yi Wang, Sicheng Zhao, Yiyi Zhang, and Lap-Pui Chau.  
595 A survey of embodied learning for object-centric robotic manipulation. *Machine Intelligence*  
596 *Research*, pp. 1–39, 2025.  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

648 APPENDIX  
649650 A ADDITIONAL IMPLEMENTATION DETAILS  
651652 A.1 DETAILS OF NETWORK STRUCTURE  
653654 **Algorithm 1** Compositional Graph Transfer (CGT)  
655

---

656 **Require:** Support keypoint embedding  $F_s$ , query feature  $F_q$   
 657 1: **for**  $l = 1, \dots, L_d$  **do**  
 658 2:   StructureVAE encoding:  
 659        $[\mu^{(l)}, \log(\sigma^{(l)})] = \text{Enc}(F_s^{(l)})$   
 660 3:   Random sampling ( $N_s$  turns)  
 661        $\mathbf{z}^{(l)}|_{n=1}^{N_s} = \mu^{(l)} + \sigma^{(l)} \odot \epsilon$   
 662 4:   StructureVAE decoding:  
 663        $\{\tilde{\mathcal{A}}_n^{(l)}\}_{n=1}^{N_s} = \text{Dec}(\mathbf{z}^{(l)}|_{n=1}^{N_s})$   
 664 5:   Computing confidence:  
 665        $w_n = \frac{1}{\sum_{i=1}^{D_z} \sigma_{n,i}^{(l)} + \epsilon}, \tilde{w}_n = \frac{w_n}{\sum_{m=1}^{N_s} w_m}$   
 666 6:   Bayesian aggregation:  $\mathcal{A}_{\text{fused}}^{(l)} = \sum_{n=1}^{N_s} \tilde{w}_n \cdot \tilde{\mathcal{A}}_n^{(l)}$   
 667 7:   Update support keypoint embedding:  
 668        $F_s^{(l+1)} = \text{GCN}(F_s^{(l)}, \mathcal{A}_{\text{fused}}^{(l)})$   
 669  
 670 8:   Compute Gating Scores:  $\alpha^{(l)} = \frac{\text{sim}(\text{Pool}(F_q), \mu^{(l)})}{\sum_{l=1}^L \text{sim}(\text{Pool}(F_q), \mu^{(l)})}, \quad L \in [1, L_d]$   
 671 9:   Fusion across Layers:  $\mathcal{A}_{\text{final}}^{(l)} \leftarrow \sum_{l=1}^L \alpha^{(l)} \cdot \mathcal{A}_{\text{fused}}^{(l)}$   
 672 10: **end for**

---



673  
 674  
 675  
 676  
 677  
 678  
 679  
 680  
 681  
 682  
 683  
 684  
 685  
 686  
 687  
 688  
 689  
 690  
 691  
 692  
 693  
 694  
 695  
 696  
 697  
 698  
 699  
 700  
 701  
 Figure 5: The pipeline of Compositional Graph Transfer. The i-SVAE in the  $l$ -th decoder layer outputs sampled adjacencies  $\{\tilde{\mathcal{A}}_n^{(l)}\}_{n=1}^{N_s}$  with posterior  $\mu^{(l)}, \sigma^{(l)}$ . Bayesian aggregation fuses samples into a single adjacency estimate  $\tilde{\mathcal{A}}_{\text{fused}}^{(l)}$ , and query-guided gating reweights layers using the query embedding  $F_q$ . The result is the final graph  $\tilde{\mathcal{A}}_{\text{final}}^{(l)}$  in that layer.

702  
 703  
 704  
 705  
 706  
 707  
 708  
 709  
 710  
 711  
 712  
 713  
 714  
 715  
 716  
 717  
 718  
 719  
 720  
 721  
 722  
 723  
 724  
 725  
 726  
 727  
 728  
 729  
 730  
 731  
 732  
 733  
 734  
 735  
 736  
 737  
 738  
 739  
 740  
 741  
 742  
 743  
 744  
 745  
 746  
 747  
 748  
 749  
 750  
 751  
 752  
 753  
 754  
 755  
 756  
 757  
 758  
 759  
 760  
 761  
 762  
 763  
 764  
 765  
 766  
 767  
 768  
 769  
 770  
 771  
 772  
 773  
 774  
 775  
 776  
 777  
 778  
 779  
 780  
 781  
 782  
 783  
 784  
 785  
 786  
 787  
 788  
 789  
 790  
 791  
 792  
 793  
 794  
 795  
 796  
 797  
 798  
 799  
 800  
 801  
 802  
 803  
 804  
 805  
 806  
 807  
 808  
 809  
 810  
 811  
 812  
 813  
 814  
 815  
 816  
 817  
 818  
 819  
 820  
 821  
 822  
 823  
 824  
 825  
 826  
 827  
 828  
 829  
 830  
 831  
 832  
 833  
 834  
 835  
 836  
 837  
 838  
 839  
 840  
 841  
 842  
 843  
 844  
 845  
 846  
 847  
 848  
 849  
 850  
 851  
 852  
 853  
 854  
 855  
 856  
 857  
 858  
 859  
 860  
 861  
 862  
 863  
 864  
 865  
 866  
 867  
 868  
 869  
 870  
 871  
 872  
 873  
 874  
 875  
 876  
 877  
 878  
 879  
 880  
 881  
 882  
 883  
 884  
 885  
 886  
 887  
 888  
 889  
 890  
 891  
 892  
 893  
 894  
 895  
 896  
 897  
 898  
 899  
 900  
 901  
 902  
 903  
 904  
 905  
 906  
 907  
 908  
 909  
 910  
 911  
 912  
 913  
 914  
 915  
 916  
 917  
 918  
 919  
 920  
 921  
 922  
 923  
 924  
 925  
 926  
 927  
 928  
 929  
 930  
 931  
 932  
 933  
 934  
 935  
 936  
 937  
 938  
 939  
 940  
 941  
 942  
 943  
 944  
 945  
 946  
 947  
 948  
 949  
 950  
 951  
 952  
 953  
 954  
 955  
 956  
 957  
 958  
 959  
 960  
 961  
 962  
 963  
 964  
 965  
 966  
 967  
 968  
 969  
 970  
 971  
 972  
 973  
 974  
 975  
 976  
 977  
 978  
 979  
 980  
 981  
 982  
 983  
 984  
 985  
 986  
 987  
 988  
 989  
 990  
 991  
 992  
 993  
 994  
 995  
 996  
 997  
 998  
 999  
 1000  
 1001  
 1002  
 1003  
 1004  
 1005  
 1006  
 1007  
 1008  
 1009  
 1010  
 1011  
 1012  
 1013  
 1014  
 1015  
 1016  
 1017  
 1018  
 1019  
 1020  
 1021  
 1022  
 1023  
 1024  
 1025  
 1026  
 1027  
 1028  
 1029  
 1030  
 1031  
 1032  
 1033  
 1034  
 1035  
 1036  
 1037  
 1038  
 1039  
 1040  
 1041  
 1042  
 1043  
 1044  
 1045  
 1046  
 1047  
 1048  
 1049  
 1050  
 1051  
 1052  
 1053  
 1054  
 1055  
 1056  
 1057  
 1058  
 1059  
 1060  
 1061  
 1062  
 1063  
 1064  
 1065  
 1066  
 1067  
 1068  
 1069  
 1070  
 1071  
 1072  
 1073  
 1074  
 1075  
 1076  
 1077  
 1078  
 1079  
 1080  
 1081  
 1082  
 1083  
 1084  
 1085  
 1086  
 1087  
 1088  
 1089  
 1090  
 1091  
 1092  
 1093  
 1094  
 1095  
 1096  
 1097  
 1098  
 1099  
 1100  
 1101  
 1102  
 1103  
 1104  
 1105  
 1106  
 1107  
 1108  
 1109  
 1110  
 1111  
 1112  
 1113  
 1114  
 1115  
 1116  
 1117  
 1118  
 1119  
 1120  
 1121  
 1122  
 1123  
 1124  
 1125  
 1126  
 1127  
 1128  
 1129  
 1130  
 1131  
 1132  
 1133  
 1134  
 1135  
 1136  
 1137  
 1138  
 1139  
 1140  
 1141  
 1142  
 1143  
 1144  
 1145  
 1146  
 1147  
 1148  
 1149  
 1150  
 1151  
 1152  
 1153  
 1154  
 1155  
 1156  
 1157  
 1158  
 1159  
 1160  
 1161  
 1162  
 1163  
 1164  
 1165  
 1166  
 1167  
 1168  
 1169  
 1170  
 1171  
 1172  
 1173  
 1174  
 1175  
 1176  
 1177  
 1178  
 1179  
 1180  
 1181  
 1182  
 1183  
 1184  
 1185  
 1186  
 1187  
 1188  
 1189  
 1190  
 1191  
 1192  
 1193  
 1194  
 1195  
 1196  
 1197  
 1198  
 1199  
 1200  
 1201  
 1202  
 1203  
 1204  
 1205  
 1206  
 1207  
 1208  
 1209  
 1210  
 1211  
 1212  
 1213  
 1214  
 1215  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241  
 1242  
 1243  
 1244  
 1245  
 1246  
 1247  
 1248  
 1249  
 1250  
 1251  
 1252  
 1253  
 1254  
 1255  
 1256  
 1257  
 1258  
 1259  
 1260  
 1261  
 1262  
 1263  
 1264  
 1265  
 1266  
 1267  
 1268  
 1269  
 1270  
 1271  
 1272  
 1273  
 1274  
 1275  
 1276  
 1277  
 1278  
 1279  
 1280  
 1281  
 1282  
 1283  
 1284  
 1285  
 1286  
 1287  
 1288  
 1289  
 1290  
 1291  
 1292  
 1293  
 1294  
 1295  
 1296  
 1297  
 1298  
 1299  
 1300  
 1301  
 1302  
 1303  
 1304  
 1305  
 1306  
 1307  
 1308  
 1309  
 1310  
 1311  
 1312  
 1313  
 1314  
 1315  
 1316  
 1317  
 1318  
 1319  
 1320  
 1321  
 1322  
 1323  
 1324  
 1325  
 1326  
 1327  
 1328  
 1329  
 1330  
 1331  
 1332  
 1333  
 1334  
 1335  
 1336  
 1337  
 1338  
 1339  
 1340  
 1341  
 1342  
 1343  
 1344  
 1345  
 1346  
 1347  
 1348  
 1349  
 1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403  
 1404  
 1405  
 1406  
 1407  
 1408  
 1409  
 1410  
 1411  
 1412  
 1413  
 1414  
 1415  
 1416  
 1417  
 1418  
 1419  
 1420  
 1421  
 1422  
 1423  
 1424  
 1425  
 1426  
 1427  
 1428  
 1429  
 1430  
 1431  
 1432  
 1433  
 1434  
 1435  
 1436  
 1437  
 1438  
 1439  
 1440  
 1441  
 1442  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457  
 1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503  
 1504  
 1505  
 1506  
 1507  
 1508  
 1509  
 1510  
 1511  
 1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565  
 1566  
 1567  
 1568  
 1569  
 1570  
 1571  
 1572  
 1573  
 1574  
 1575  
 1576  
 1577  
 1578  
 1579  
 1580  
 1581  
 1582  
 1583  
 1584  
 1585  
 1586  
 1587  
 1588  
 1589  
 1590  
 1591  
 1592  
 1593  
 1594  
 1595  
 1596  
 1597  
 1598  
 1599  
 1600  
 1601  
 1602  
 1603  
 1604  
 1605  
 1606  
 1607  
 1608  
 1609  
 1610  
 1611  
 1612  
 1613  
 1614  
 1615  
 1616  
 1617  
 1618  
 1619  
 1620  
 1621  
 1622  
 1623  
 1624  
 1625  
 1626  
 1627  
 1628  
 1629  
 1630  
 1631  
 1632  
 1633  
 1634  
 1635  
 1636  
 1637  
 1638  
 1639  
 1640  
 1641  
 1642  
 1643  
 1644  
 1645  
 1646  
 1647  
 1648  
 1649  
 1650  
 1651  
 1652  
 1653  
 1654  
 1655  
 1656  
 1657  
 1658  
 1659  
 1660  
 1661  
 1662  
 1663  
 1664  
 1665  
 1666  
 1667  
 1668  
 1669  
 1670  
 1671  
 1672  
 1673  
 1674  
 1675  
 1676  
 1677  
 1678  
 1679  
 1680  
 1681  
 1682  
 1683  
 1684  
 1685  
 1686  
 1687  
 1688  
 1689  
 1690  
 1691  
 1692  
 1693  
 1694  
 1695  
 1696  
 1697  
 1698  
 1699  
 1700  
 1701  
 1702  
 1703  
 1704  
 1705  
 1706  
 1707  
 1708  
 1709  
 1710  
 1711  
 1712  
 1713  
 1714  
 1715  
 1716  
 1717  
 1718  
 1719  
 1720  
 1721  
 1722  
 1723  
 1724  
 1725  
 1726  
 1727  
 1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
 1737  
 1738  
 1739  
 1740  
 1741  
 1742  
 1743  
 1744  
 1745  
 1746  
 1747  
 1748  
 1749  
 1750  
 1751  
 1752  
 1753  
 1754  
 1755  
 1756  
 1757  
 1758  
 1759  
 1760  
 1761  
 1762  
 1763  
 1764  
 1765  
 1766  
 1767  
 1768  
 1769  
 1770  
 1771  
 1772  
 1773  
 1774  
 1775  
 1776  
 1777  
 1778  
 1779  
 1780  
 1781  
 1782  
 1783  
 1784  
 1785  
 1786  
 1787  
 1788  
 1789  
 1790  
 1791  
 1792  
 1793  
 1794  
 1795  
 1796  
 1797  
 1798  
 1799  
 1800  
 1801  
 1802  
 1803  
 1804  
 1805  
 1806  
 1807  
 1808  
 1809  
 1810  
 1811  
 1812  
 1813  
 1814  
 1815  
 1816  
 1817  
 1818  
 1819  
 1820  
 1821  
 1822  
 1823  
 1824  
 1825  
 1826  
 1827  
 1828  
 1829  
 1830  
 1831  
 1832  
 1833  
 1834  
 1835  
 1836  
 1837  
 1838  
 1839  
 1840  
 1841  
 1842  
 1843  
 1844  
 1845  
 1846  
 1847  
 1848  
 1849  
 1850  
 1851  
 1852  
 1853  
 1854  
 1855  
 1856  
 1857  
 1858  
 1859  
 1860  
 1861  
 1862  
 1863  
 1864  
 1865  
 1866  
 1867  
 1868  
 1869  
 1870  
 1871  
 1872  
 1873  
 1874  
 1875  
 1876  
 1877  
 1878  
 1879  
 1880  
 1881  
 1882  
 1883  
 1884  
 1885  
 1886  
 1887  
 1888  
 1889  
 1890  
 1891  
 1892  
 1893  
 1894  
 1895  
 1896  
 1897  
 1898  
 1899  
 1900  
 1901  
 1902  
 1903  
 1904  
 1905  
 1906  
 1907  
 1908  
 1909  
 1910  
 1911  
 1912  
 1913  
 1914  
 1915  
 1916  
 1917  
 1918  
 1919  
 1920  
 1921  
 1922  
 1923  
 1924  
 1925  
 1926  
 1927  
 1928  
 1929  
 1930  
 1931  
 1932  
 1933  
 1934  
 1935  
 1936  
 1937  
 1938  
 1939  
 1940  
 1941  
 1942  
 1943  
 1944  
 1945  
 1946  
 1947  
 1948  
 1949  
 1950  
 1951  
 1952  
 1953  
 1954  
 1955  
 1956  
 1957  
 1958  
 1959  
 1960<br

702 Table 7: Architecture details of the *iterative StructureVAE*, including graph encoder, latent sampling,  
 703 and graph decoder in the Figure 2 (main manuscript).

|   | Norm & Activation             | Output Shape                       |
|---|-------------------------------|------------------------------------|
| <b>Graph Encoder (for posterior distribution)</b> |                               |                                    |
| Input node feature $x$                            | -                             | $[B, M, D]$                        |
| Reshape $x \rightarrow x'$                        | -                             | $[B, M \cdot D]$                   |
| Linear  | ReLU                          | $[B, \text{hidden\_dim}]$          |
| Linear  | -                             | $[B, 2 \cdot \text{latent\_dim}]$  |
| Split $\rightarrow (\mu, \log \sigma^2)$          | -                             | $[B, \text{latent\_dim}] \times 2$ |
| <b>Latent Sampling (Reparameterization)</b>       |                               |                                    |
| If training: $z = \mu + \epsilon \cdot \sigma$    | -                             | $[B, \text{latent\_dim}]$          |
| Else: $z = \mu$                                   | -                             | $[B, \text{latent\_dim}]$          |
| <b>Graph Decoder (to soft adjacency matrix)</b>   |                               |                                    |
| Linear  | Sigmoid                       | $[B, M \cdot M]$                   |
| Reshape   | -                             | $[B, M, M]$                        |
| Symmetrization                                    | $A = \frac{1}{2}(A + A^\top)$ | $[B, M, M]$                        |

719 priors, facilitating the downstream reasoning module to learn more expressive and structure-aware  
 720 keypoint representations.

## 722 A.2 DETAILED TRAINING CONFIGURATIONS

724 Table 8 summarizes the training and evaluation settings for the GenCape-T/S models. We adopt the  
 725 Adam Kingma & Ba (2015) optimizer with a base learning rate of  $1.0 \times 10^{-5}$ , scheduled linearly  
 726 and warmed up over 1,000 iterations with a warmup ratio of 0.001. Training is performed for  
 727 175,000 epochs with a batch size of 16 in float32 precision, ensuring stable convergence. For few-  
 728 shot evaluation, the model is assessed under both 1-shot and 5-shot settings, with each episode  
 729 comprising 15 query samples and a total of 200 episodes to ensure statistical reliability.

## 731 B ADDITIONAL EXPERIMENTAL RESULTS

### 732 B.1 ABLATIONS ON DIFFERENT STRATEGIES

735 To further assess the effectiveness of our proposed  
 736 framework, we investigate the impact of different  
 737 graph construction methods and different composi-  
 738 tional graph transfer strategies. Table 9 compares  
 739 different methods for constructing adjacency matrix.  
 740 Random initialization results in a significantly per-  
 741 formance drop, confirming the importance of struc-  
 742 tural priors. Learnable graphs consistently outper-  
 743 form static counterparts, and removing the symme-  
 744 try constraint leads to only a slight decrease, sug-  
 745 gesting that bidirectionality is more critical than di-  
 746 rectional specificity. We replace our layer-wise i-  
 747 SVAE updates with the first-layer adjacency matrix  
 748 predicted only once from the encoder output. We ob-  
 749 serve: iter (92.05) vs. non-iter (91.48). This +0.57  
 750 improvement demonstrates that iterative refinement  
 751 is indeed beneficial. A fixed adjacency estimated  
 752 once from static support features cannot adapt to the  
 753 evolving decoder representations. Furthermore, we  
 754 conducted an experiment that directly used the self-  
 755 attention weights from Figure 2 as the adjacency matrix. This variant achieves only 89.33 PCK@0.2  
 (2.72 drop vs. 92.05), indicating that attention-induced “connectivity” fails to provide meaningful  
 structure. This is expected: self-attention mainly captures appearance-driven correlations, lacks  
 uncertainty modeling, especially when support features are ambiguous. Table 10 investigates com-

756 Table 8: Training configuration used for  
 757 GenCape-T/S.

| Training recipe:                          |         |
|---|---------|
| optimizer                                 | Adam    |
| <b>Learning hyper-parameters:</b>         |         |
| base learning rate                        | 1.0E-05 |
| learning rate schedule                    | linear  |
| batch size                                | 16      |
| training steps                            | 175,000 |
| lr warmup iters                           | 1,000   |
| warmup ratio                              | 0.001   |
| warmup schedule                           | linear  |
| data type                                 | float32 |
| norm epsilon                              | 1.0E-06 |
| <b>Few-shot testing hyper-parameters:</b> |         |
| shots                                     | 1 / 5   |
| num_query                                 | 15      |
| num_episodes                              | 200     |

756 positional graph transfer strategies. The combination of query weighting at the layer level and  
 757 Bayesian fusion at the sampling level achieves the highest 92.05% PCK. In contrast, using the same  
 758 strategy at both levels (e.g., query weighting for both) yields suboptimal results. We attribute this  
 759 improvement to the complementary strengths of the two mechanisms: query weighting enables dy-  
 760 namic alignment of structural importance across layers based on the semantic context of the query,  
 761 while Bayesian fusion effectively mitigates uncertainty introduced by latent graph sampling. Con-  
 762 versely, a mismatch between Bayesian fusion across layers and query weighting across samples  
 763 performs worse (91.55%). These findings highlight the importance of hybrid fusion strategies that  
 764 jointly consider semantic relevance and structural reliability. Layer-wise representations encode  
 765 semantically distinct structural abstractions that benefit from query-adaptive weighting rather than  
 766 confidence-based averaging. Applying Bayesian fusion at this level can obscure semantically salient  
 767 but uncertain layers, effectively flattening meaningful hierarchical distinctions.

768  
 769 Table 9: Comparisons of different adjacency  
 770 matrix construction strategies under 1-shot  
 771 setting with SwinV2-Tiny backbone.

| Type                     | Symmetric | PCK          |
|--------------------------|-----------|--------------|
| Static Graph             | ✓         | 91.19        |
| Learnable Graph          | ✓         | <b>92.05</b> |
| Learnable Graph          | ✗         | 91.71        |
| Random Initialized Graph | ✓         | 84.39        |
| Non-iter                 | ✓         | 91.48        |
| Self-attention           | ✓         | 89.33        |

768  
 769 Table 10: Comparisons on compositional  
 770 graph transfer strategies under 1-shot setting  
 771 with SwinV2-Tiny backbone.

| Layer-wise      | Sampling-wise   | PCK          |
|-----------------|-----------------|--------------|
| Bayesian Fusion | Bayesian Fusion | 91.36        |
| Query Weighting | Query Weighting | 91.74        |
| Bayesian Fusion | Query Weighting | 91.55        |
| Query Weighting | Bayesian Fusion | <b>92.05</b> |

## 781 B.2 ADDITIONAL CROSS-SUPERCATEGORY RESULTS

782  
 783 Table 11: **Cross-supercategory results.** PCK@0.2 performance under the 1-shot setting on Split-1.  
 784 Following the standard super-category partitioning protocol, our method achieves the best perfor-  
 785 mance across all splits, demonstrating its strong generalization.

| Method         | HumanBody    | HumanFace    | Vehicle      | Furniture    |
|----------------|--------------|--------------|--------------|--------------|
| ProtoNet       | 37.61        | 57.80        | 28.35        | 42.64        |
| MAML           | 51.93        | 25.72        | 17.68        | 20.09        |
| Fine-tune      | 52.11        | 25.53        | 17.46        | 20.76        |
| POMNet         | 73.82        | 79.63        | 34.92        | 47.27        |
| CapeFormer     | 83.44        | 80.96        | 45.40        | 52.49        |
| GraphCape      | 88.38        | 83.28        | 44.06        | 45.56        |
| <b>GenCape</b> | <b>89.69</b> | <b>93.76</b> | <b>47.74</b> | <b>66.63</b> |

796 We follow prior works Liang et al. (2024); Hirschorn & Avidan (2024); Rusanovsky et al. (2025)  
 797 and perform a cross-supercategory evaluation to rigorously assess the generalization ability of our  
 798 model across semantically diverse object classes. Concretely, in Table 11 we treat one of the four  
 799 MP-100 supercategories—**HumanBody**, **HumanFace**, **Vehicle**, **Furniture**—as the test domain and  
 800 train on the remaining three, creating four disjoint train–test splits. GenCape consistently achieves  
 801 the highest accuracy across all cross-supercategory splits. These improvements highlight the strong  
 802 generalization of our generative structural modeling, which captures keypoint dependencies that  
 803 remain robust under large variations.

804 We further evaluate a more challenging setting: training on one category and testing on a differ-  
 805 ent, structurally mismatched category. Across all Train→Test pairs, Table 12 shows that Gen-  
 806 Cape consistently surpasses GraphCape, indicating stronger resilience to cross-category discrep-  
 807 ancies. Notably, GenCape achieves substantial gains in HumanBody→HumanFace (+11.36), Vehi-  
 808 cle→HumanBody (+8.40), and Chair→HumanBody (+4.23), showing its ability to adapt to entirely  
 809 different topologies. Overall, these results highlight that GenCape learns transferable, generative  
 keypoint relations that generalize reliably across heterogeneous object categories.

810  
811 Table 12: **Cross domain transfer evaluation.** PCK@0.2 performance under the 1-shot setting on  
812 **Split-1.** Training on one super-category and testing on the other.

| Train     | Test      | GraphCape    | GenCape      |
|-----------|-----------|--------------|--------------|
| HumanBody | HumanFace | 33.87        | <b>45.23</b> |
| HumanFace | HumanBody | 55.90        | <b>56.43</b> |
| Furniture | HumanBody | <b>73.79</b> | 73.09        |
| HumanBody | Furniture | 31.11        | <b>50.47</b> |
| Vehicle   | HumanBody | 50.13        | <b>58.53</b> |
| HumanBody | Vehicle   | 28.52        | <b>32.64</b> |
| Chair     | HumanBody | 49.10        | <b>53.33</b> |

821 **B.3 ADDITIONAL METRICS RESULTS**

822  
823 Table 13: **Additional Metrics.** AUC, EPE and NME performance under 1-shot setting.

| Method           | AUC % ( $\uparrow$ ) | EPE ( $\downarrow$ ) | NME ( $\downarrow$ ) | PCK % ( $\uparrow$ ) |
|------------------|----------------------|----------------------|----------------------|----------------------|
| GraphCape-T      | 89.10                | 41.04                | 0.08                 | 91.19                |
| <b>GenCape-T</b> | <b>89.50</b>         | 39.65                | 0.08                 | 92.05                |
| GraphCape-S      | 91.16                | 30.05                | 0.06                 | 94.73                |
| <b>GenCape-S</b> | <b>91.37</b>         | 29.62                | 0.06                 | 95.23                |

831  
832 We further evaluate our model using three standard keypoint localization metrics, as summarized  
833 in Table 13. AUC measures the area under the PCK curve and reflects overall accuracy across a  
834 range of distance thresholds. EPE computes the Euclidean distance between predicted and ground-  
835 truth keypoints. NME reports the mean localization error normalized by object scale. These metrics  
836 provide a comprehensive assessment of both absolute and scale-invariant localization performance.

837  
838 **B.4 COMPARISONS ON COMPUTATIONAL COMPLEXITY**

839  
840 Table 14: **Comparison of computational complexity and accuracy across methods.**

| Method           | GFLOPs       | Params | FPS   | PCK   |
|------------------|--------------|--------|-------|-------|
| POMNet           | 38.01        | 48.21M | 6.80  | 46.05 |
| One-Stage        | 22.65        | 26.86M | 36.90 | –     |
| CapeFormer       | 23.68        | 31.14M | 26.09 | 89.45 |
| GraphCape-T      | 15.48        | 43.68M | 15.36 | 91.19 |
| <b>GenCape-T</b> | <b>15.66</b> | 44.47M | 14.89 | 92.05 |
| GraphCape-S      | 27.75        | 65.06M | 10.45 | 94.73 |
| <b>GenCape-S</b> | <b>27.93</b> | 65.85M | 10.44 | 95.23 |

851  
852 To further clarify computational efficiency, we provide a comparison of GFLOPs, parameter counts,  
853 and inference speed. GraphCape and GenCape are tested on A100, while the results of POMNet,  
854 One-Stage, and CapeFormer are taken from CapeFormer Shi et al. (2023) where all measurements  
855 were obtained on a RTX 3090. As shown in Table 14, despite this hardware discrepancy, the  
856 comparison still reveals a clear trend: GenCape introduces negligible computational overhead relative to  
857 GraphCape (e.g., +0.18 GFLOPs and +0.8M params for the Tiny variant), while consistently deliv-  
858 ering higher accuracy. This confirms that our generative structural modeling improves performance  
859 without sacrificing efficiency.

860  
861 **C ADDITIONAL VISUALIZATION RESULTS**

862  
863 In this section, we present more qualitative results. As shown in the first row of Figure 6, the support  
864 image depicts an inverted top-view of a dog, where GraphCape Hirschorn & Avidan (2024) exhibits

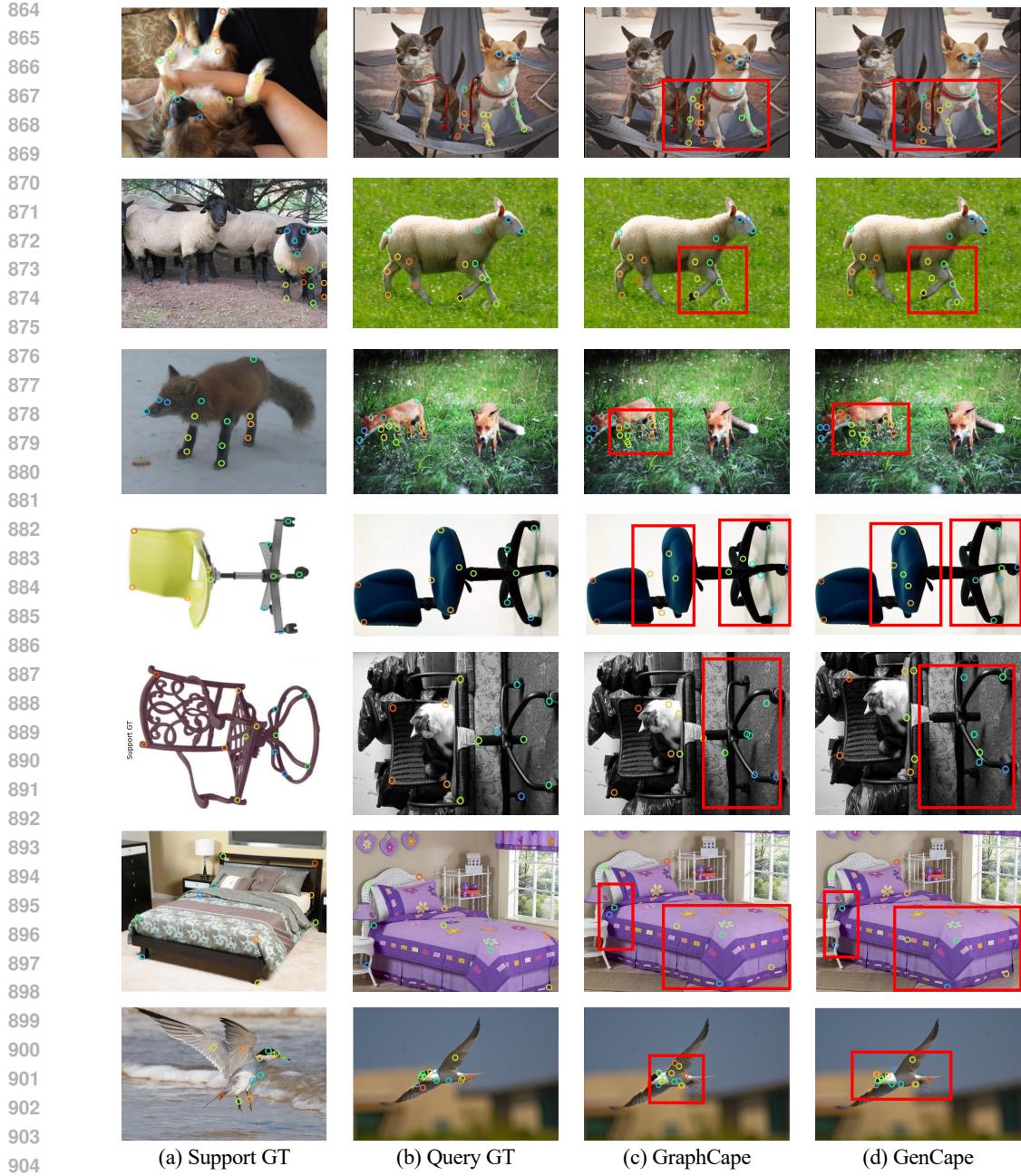


Figure 6: **Comparative qualitative results.** We compare more keypoint predictions with GraphCape Hirschorn & Avidan (2024) under the 1-shot setting. The red boxes highlight the regions with significant differences.

noticeable prediction drift on the left hind leg and right foreleg, while our method achieves more accurate localization. In the second and third rows, GraphCape overly relies on support instances with dissimilar poses, resulting in incorrect keypoint predictions. The fourth row presents a challenging case involving a swivel chair, where structural reasoning becomes critical for precise keypoint inference. Despite preserving the overall skeleton shape, GraphCape relies solely on manually defined connections and fails to localize the seat correctly, producing an upward shift as highlighted by the red box. A similar issue is observed in the fifth row, where the predicted seat location is misaligned to the edge of the cat. Figure 7 provides additional category-agnostic pose estimation examples to further illustrate the effectiveness of our approach.

Figure 8 shows more visualizations between predefined skeletons and the latent adjacency matrices inferred across different decoder layers. In the first example (giraffe), the learned adjacency matrices

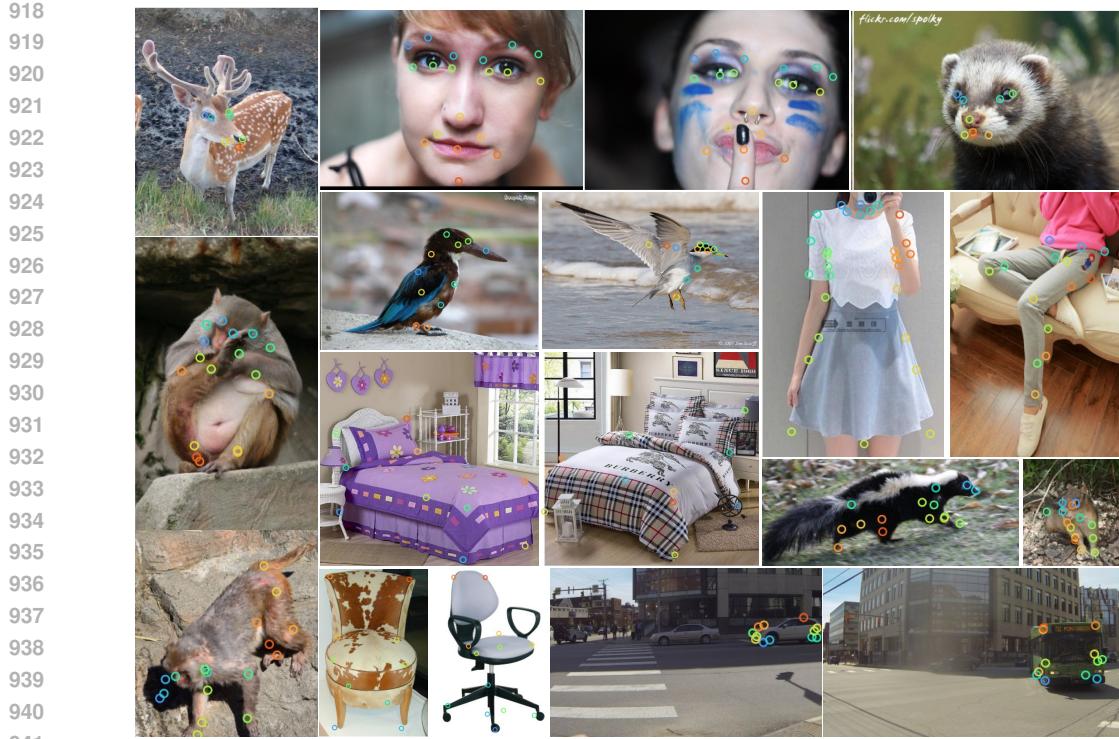


Figure 7: Qualitative visualization. We visualize more keypoint predictions across different splits.

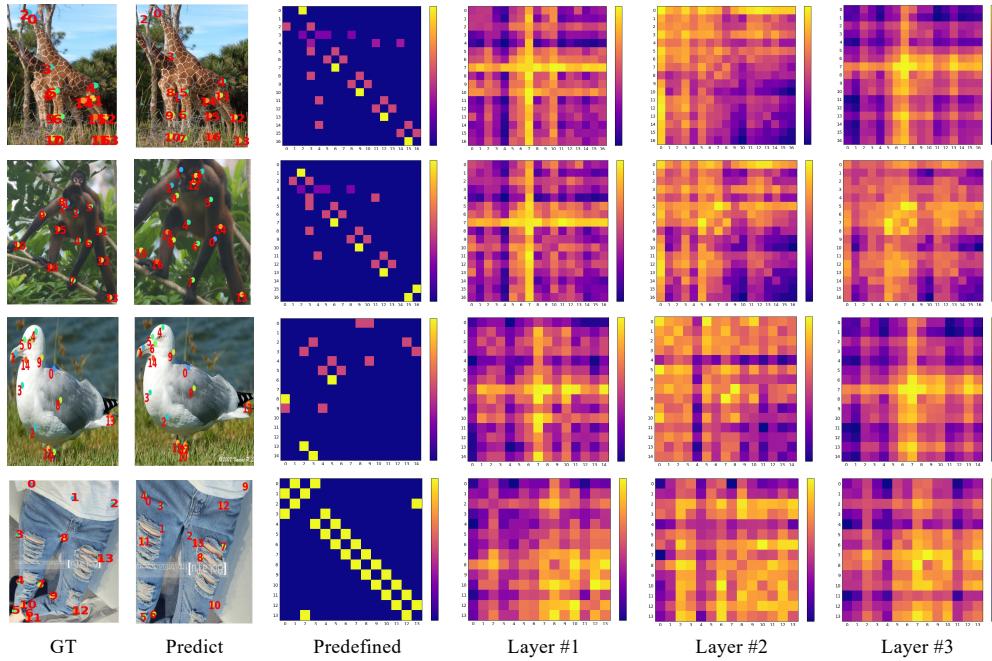
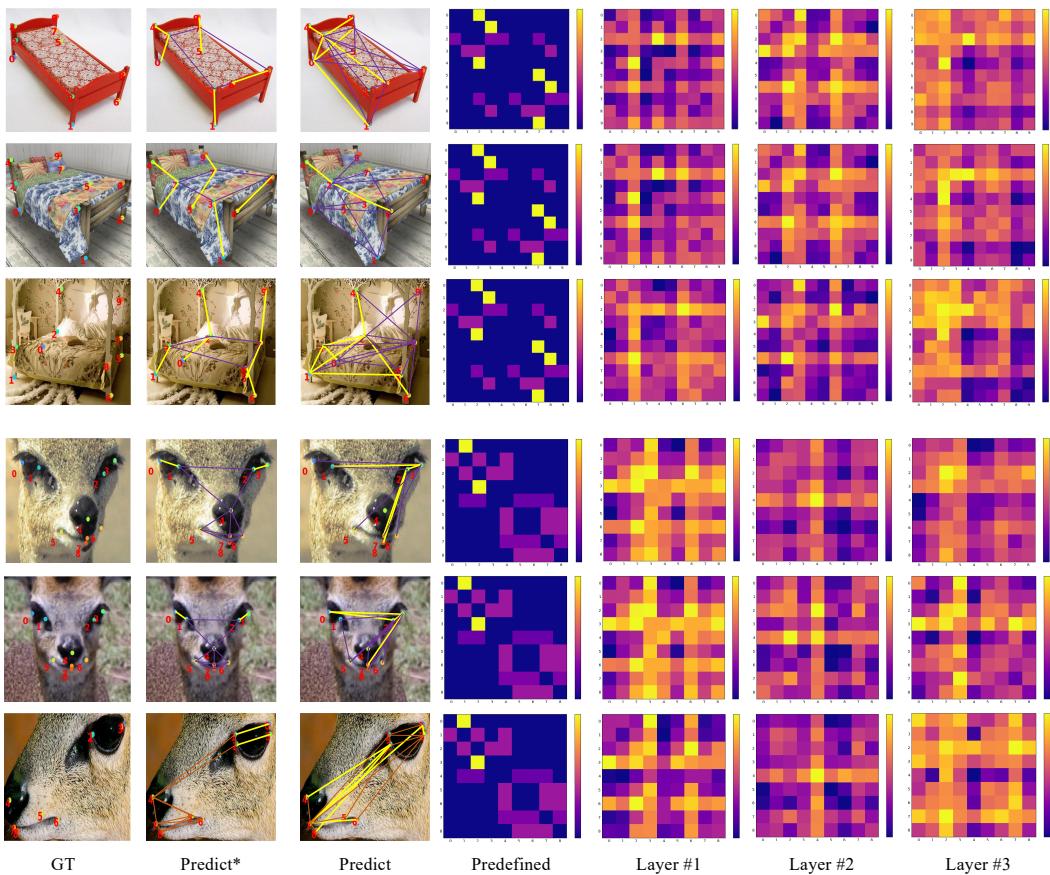


Figure 8: Adjacency matrix visualization. We visualize more latent graph structures across different splits.

progressively capture long-range dependencies critical for reasoning along the elongated neck and legs, outperforming the sparse predefined skeleton. The second case (monkey) presents occlusion challenges from tree branches, where the model adaptively strengthens cross-limb connections in deeper layers to improve structural consistency, though minor prediction drift remains. For the bird, although predefined structures already offer reasonable symmetry, latent graphs further refine bilateral dependencies, enhancing accuracy. The trouser case represents a particularly challenging case due to its dense, repetitive keypoints and weak structural cues. These visualizations also enhance

972 the interpretability of our GenCape. However, from Figure 4 and 9, we observe that localization  
 973 errors are primarily caused by visual feature ambiguity. Predictions on the swivelchair category show  
 974 large structural deviations, suggesting that the weak and homogeneous textures of this class hinder  
 975 the transfer of support-driven structural priors. We further evaluate robustness by randomly masking  
 976 the query image. When 25% of the image is masked, GenCape-S shows only a mild drop (93.05),  
 977 but performance collapses at 50% masking (76.81), showing that heavy occlusion severely disrupts  
 978 the visual evidence required for localization. This sharp degradation reinforces the need for gen-  
 979 erative, uncertainty-aware structural modeling to cope with missing keypoints and support–query  
 980 mismatch.



1009 **Figure 9: Additional adjacency matrix and skeleton visualization.** We visualized samples from two cate-  
 1010 gories (bed and klipspringer face) in Split-1 for a more comprehensive illustration. The Predict\* column is the  
 1011 predicted locations with the prior connections, while the Predict column is with learned connections.

1012 **More qualitative analysis of the adjacency matrix of the same category.** Figure 9 presents ad-  
 1013 dditional adjacency-matrix and skeleton visualizations for two same categories in Figure 3. Each  
 1014 row corresponds to the layer-wise adjacency matrices progression produced and finally the skeleton.  
 1015 Brighter colors in the matrices indicate stronger relational dependencies between corresponding  
 1016 keypoints. We observe that adjacency patterns remain similar across samples of the same category. For  
 1017 instance, in first 3 rows of bed category, the adjacency matrices consistently shows a core keypoint  
 1018 showing strong influence on the others. From left to right: initially the core keypoint fires broadly,  
 1019 then adjacent points start to dominate local neighborhoods, and the final layer produces localized  
 1020 high-response clusters and suppressed irrelevant edges, indicating a confident and discriminative  
 1021 dependency graph. For the klipspringer face category in the last three rows, the progression follows  
 1022 a different pattern: early layers emphasize local smoothness among neighboring keypoints, then the  
 1023 model increasingly attends to the central nose region as an anchor, and the final layer converges to  
 1024 a distinct pattern. These differences across categories reflect that the layer-wise graphs represent a  
 1025 coarse-to-fine refinement of functional dependencies. The evidences that the learned graphs encode  
 structural dependencies useful for CAPE, rather than merely aiding optimization convergence.