

# 基于图卷积网络的三维手部姿态估计

彭春燕<sup>1,2</sup>, 王璇<sup>1,2</sup>, 陈杨博<sup>1,2</sup>, 何港波<sup>1,2</sup>

(1. 青海师范大学 计算机学院, 青海 西宁 810016; 2. 青海师范大学 藏语智能全国重点实验室, 青海 西宁 810016)

**摘要:** 基于单张彩色图片的三维手部姿态估计因手部存在遮挡、手部自相似性高等原因导致预测结果存在误差大、手部结构不自然等问题。针对这些问题, 首先, 提出一个基于图卷积的三维手部姿态估计方法, 使用 Keypoint R-CNN 提取图像视觉特征和手部关键点二维位置信息, 将特征信息输入到改进的自适应核图卷积模块(AK\_GraFormer)中; 其次, 引入带残差连接的 AKGNN 图核, 自适应处理图数据以增强模型的特征学习与表达; 最后, 利用提出的评估指标监控动态训练策略以获得更优的估计结果。通过在 HO3D v3 数据集与 FreiHand 数据集上实验, 结果表明在单张彩色图片手部三维姿态估计任务中, 所提方法相比其他同类方法具有明显优势, 刚性对齐后的平均每关节位置误差(PA-MPJPE)最高降低了 14.28 个百分点, 检测关节点百分比曲线下面积(AUC)最高提高了 3.33 百分点。

**关键词:** 手部三维姿态估计; 图卷积网络; 特征提取; 图核学习优化; 评估指标动态调整

**中图分类号:** TP391; TP751 **文献标志码:** A **doi:** 10.13705/j.issn.1671-6833.2026.02.013

手是我们与世界感知、交互的重要媒介, 是连接人类与世界的桥梁。计算机通常需要检测手部的姿态特征来理解人类行为, 因而获取手-物图片并从中估计手部关键点和物体的三维位置极其重要。另外, 虚拟现实技术蓬勃发展, 人机交互使其应用领域加以拓展, 如何通过获取三维手部姿态以精准操作虚拟物体成为人机交互的研究热点问题<sup>[1]</sup>。

近年来, 大量的研究聚焦于独立的三维手部姿态估计上, 而手-物三维姿态估计相较于单独的人手估计存在更多遮挡, 因此更具挑战性。手部关节多, 自由度高, 形状变换相比人体更为复杂, 手部关节相似性高, 因而对于复杂手势识别较为困难。另外, 由于手部关节灵活性高, 物体具有一定体量, 手-物图片中往往存在不同程度的自遮挡与相互遮挡, 导致手部三维姿态估计准确度低且误差较大。

目前, 手部姿态估计方法按照手部姿态的生成方式可分为以下 3 类: 生成式手部姿态估计方法、判别式手部姿态估计方法和混合式手部姿态估计方法。生成式手部姿态估计方法通过使用预先定义手部三维模型、复杂的目标函数与优化求解算法来进

行姿态估计。Oikonomidis 等<sup>[2]</sup> 将手部参数所定义的手的三维几何模型进行渲染得到合成图像, 通过约束合成图像的手部轮廓与真实图像的手部轮廓之间的误差能量函数, 结合粒子群优化算法求解模型参数。研究者们分别提出了 sphere-meshes 模型<sup>[3]</sup>、MANO 手模型<sup>[4]</sup>、SMPL-X 手模型<sup>[5]</sup>, 随后大量研究基于这些模型展开。判别式手部姿态估计方法直接从输入的手部数据中提取具有判别力的特征, 学习手部特征与其姿态空间的映射。早期的判别式手部姿态估计方法大都采用机器学习的方法, 如随机森林算法<sup>[6]</sup>。随着深度学习的发展, 出现了越来越多的基于深度学习的方法。Tompson 等<sup>[7]</sup> 是将卷积神经网络(Convolutional Neural Networks, CNN)应用到手部姿态估计的先驱者, 随后手部姿态估计领域涌现出诸多基于卷积神经网络的手部姿态估计方法<sup>[8-10]</sup>。混合式手部姿态估计方法则是将生成式和判别式方法组合到一个框架, 一般先使用判别式方法的结果作为生成式方法的初始化<sup>[11]</sup>, 更容易优化求解得到更好的结果, 但是相比前两种生成方式, 混合式模型结构相对复杂, 且计算成本较高。

收稿日期: 2025-12-10; 修订日期: 2026-01-02

基金项目: 国家自然科学基金资助项目(62441609, 62563033); 青海省重点研发与成果转化项目(2025-2J-J08)

作者简介: 彭春燕(1980—), 女, 山东菏泽人, 青海师范大学教授, 博士, 主要从事文化计算、机器学习的研究, E-mail: pcy@qhnu.edu.cn。

最近利用图卷积网络(graph convolutional networks, GCN)从二维姿态估计三维姿态表现出非常优秀的结果<sup>[12-14]</sup>。由于使用单个二维关键点来估计三维中的对应点是一个不确定性问题,CNN的全连接层只能回归它们的位置,无法对关键点之间的连接进行建模,而GCN是可以显式处理手关节链图结构的专用网络,能够利用学习手部结构回归节点三维位置。Doosti等<sup>[12]</sup>提出了HopeNet模型,将手部骨骼的拓扑结构看作图结构,通过学习节点关系得到手部姿态。Zhao等<sup>[13]</sup>提出的GraFormer模型,通过结合多头自注意力机制和图卷积操作,在三维姿态估计任务中取得了较好的效果。Cai等<sup>[15]</sup>通过在关节与相邻帧中的对应关节之间创建额外的边缘,从几个时间相邻的二维身体姿势中创建了一个时空图。Aboukhadra等<sup>[16]</sup>将GCNs和多头注意力层进行组合,改善了手-物体交互方面的有效性。Zhuang等<sup>[17]</sup>基于运动学的抓取稳定性,利用手和物体的初始姿态,将二维坐标与提取的特征拼接成图形。Zhang等<sup>[18]</sup>提出了一种手图和对对象图之间的交互感知方法。马胜蕾等<sup>[19]</sup>提出基于双分支多尺度注意力的手部三维姿态估计,建模了手关节之间的复杂关联关系。Yang等<sup>[20]</sup>提出一种集多头自注意力、基于空间的图卷积和基于频谱的图卷积于一体的三维手部姿态和网格估计框架。

上述方法利用图卷积对手部三维姿态估计的算法均提出了创新性方法,然而在估计的准确度以及模型的表达能力方面仍有进一步提升的空间。本文基于图卷积模型,使用Keypoint R-CNN<sup>[21]</sup>,从单张RGB图像中提取多个二维特征,如热图、边界框、特征图等,经过变换得到二维位置信息和图像视觉特征,然后将这些特征信息输入到AK\_GraFormer模块中。该模块利用GraFormer<sup>[13]</sup>和AKGNN<sup>[22]</sup>,通过在其中使用改进的res\_AKGNN层使图卷积网络进一步学习特征表示,进而提升模型的表达能力。最

后利用提出的评估指标监控策略,动态调整模型权重的更迭,避免评估指标退化,提升模型训练效果。本文主要贡献归纳如下:

(1)提出一个新的三维手部姿态估计网络,提取输入图片的二维特征,最大程度保留输入图片的手部特征信息,从数据质量角度增强后续AK\_GraFormer模块的表现。

(2)使用具有残差连接的res\_AKGNN图核自适应地学习不同节点的特征信息,增强模型的特征学习与表达能力,从而提高估计准确度。

(3)在训练过程中,设计基于评估指标监控的动态调整训练策略,根据训练情况冻结模型的敏感权重,在一定程度上能够避免过拟合问题,从而提升模型性能。

## 1 方法

基于图卷积网络,本文提出了一个新的手部姿态估计网络,如图1所示。首先,使用Keypoint R-CNN提取图片信息,得到图像视觉特征和二维位置信息,其次,将提取的信息输入到自适应核AK\_GraFormer模块中,得到手部三维关键点坐标。AK\_GraFormer模块根据节点之间的连通性从图数据中提取局部特征和全局特征,自适应地控制反向传播梯度,挖掘输入数据在图结构下更复杂、更抽象的特征表示,增强网络的特征学习和表达能力。最后,使用设计的评估指标监控的动态训练策略,在训练过程中动态冻结Keypoint R-CNN模块的参数,稳定多模块协同优化过程,增强网络的整体训练效果。

### 1.1 特征信息提取

由于手部关节存在自相似性以及严重的遮挡情况,单一的二维位置只能捕捉平面信息,这导致网络在估计三维坐标时易产生歧义,所估计的三维手部姿态难以体现真实手部结构特性。鉴于此,本文在经典方法的基础上,创新性地将从图像视觉特征和二

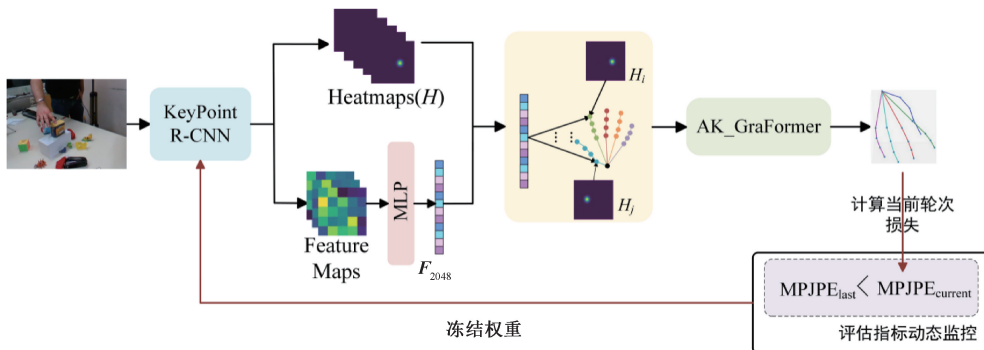


图1 本文模型整体框架

Figure 1 Overall framework of model

维位置信息同时作为后续模块的输入,旨在为网络保留更多三维空间的图像特征信息。

本文使用 KeyPoint R-CNN 提取输入图片的特征信息,如图 1 所示。通过训练 KeyPoint R-CNN 得到关键点的位置热图 (Heatmaps) 和多尺度特征图 (Feature Maps),再借助特征提取器 (MLP) 对多尺度特征图实施非线性变换,生成压缩的 2 048 维特征向量 ( $\mathbf{F}_{2048}$ )。随后,这些特征向量被附加至热图上,形成一种新的表示形式,即图像视觉特征信息。在热图中,每个通道对应一个关键点的二维位置信息,而附加的特征向量则为各关键点提供了更为丰富的上下文信息。特征信息提取模块为 AK\_GraFormer 提供了高质量输入数据,使得网络能够构建更为精准的手部节点表示。

## 1.2 AK\_GraFormer 模块设计

### 1.2.1 AK\_GraFormer 模块

在图神经网络中,节点在不同任务中对信息的需求层次不同。在手部节点构成的图结构中,靠近手掌的节点在整个图结构中起着核心的连接作用,需要模型提取更丰富、精细的信息以进行三维关键点估计;而靠近指尖的节点的位置与姿态更多是从与之相连的节点的关系来间接确定,对信息的需求层次较低。因此本文提出 AK\_GraFormer 模块,在 GraFormer 模块中引入了改进的 AKGNN 层,旨在解决模型中不同节点的信息需求层次不同的问题,增强模型的特征表达与学习,模块结构如图 2 所示。其中上半部分展示了 AK\_GraFormer 模块的结构,该模块通过学习、表达特征信息,实现对手部节点的三维坐标的估计,下半部分是 AK\_GraFormer 模块中 Res\_AKGNN 层的结构,通过决策手部关节的不同信息层次需求以自适应保留更重要的特征数据。相比于 GraFormer 模块,改进后的 AK\_GraFormer 模块能够更有效地处理手部节点图结构中的信息差异。它根据不同节点的位置和在图中的重要性,自适应地调整信息的传递和处理方式,避免了不必要的信息冗余和计算资源浪费,同时使得最终的预测结果保留更多有效的手部结构特征。

### 1.2.2 Res\_AKGNN 层

Res\_AKGNN 层的结构如图所示,根据手部节点的特点,本文提出在 Res\_AKGNN 层中加入 DenseSAGEConv,以根据节点的邻居信息更新节点的特征表示;利用 AKConv 层通过自适应核学习来调整图卷积的滤波器权重,以适应不同频率的信息;使用残差连接将经过传播和变换后的特征与原始输入相加,得到传播后的特征。而后通过全局读出函数对

节点表示进行整合,利用多层感知机进行特征变换,得到保留更多信息的特征。

DenseSAGEConv 层通过有效地融合跨层级特征,从而形成一条密集的特征传递路径,进而以更有效的信息流提高模型的表征学习能力。对于每个节点,该层首先聚合其邻居节点的特征,随后通过权重矩阵和激活函数进行特征变换,最终将邻居信息整合到节点自身的嵌入向量中。

$$\mathbf{h}_i^{new} = \sigma \left( \sum_{j \in N(i)} w_j \mathbf{h}_j + \mathbf{b} \right), \quad (1)$$

式中:  $\mathbf{h}_i^{new}$  为更新后的节点  $i$  的特征,  $N(i)$  为节点的邻居节点集合,  $w_j$  为邻居节点  $j$  的权重,  $\mathbf{h}_j$  为邻居节点的特征,  $\mathbf{b}$  为偏置项,  $\sigma$  为激活函数。

进入 Res\_AKGNN 层后,节点特征在传播过程中的信息提取过程表示如下:

$$Y^p = \text{softmax}(f_{\text{MLP}}(\text{READOUT}(\mathbf{H}^{(k)}))), \quad (2)$$

式中:  $\mathbf{H}^{(k)}$  是节点特征在传播过程中生成多个中间节点  $k$  的表示矩阵, READOUT( $\cdot$ ) 为全局读出函数,  $f_{\text{MLP}}$  为多层感知机。

输入特征经过 Res\_AKGNN 层传播和变换后得到特征  $h$ , 再通过残差连接将  $h$  与原始输入相加。残差连接使得梯度能够直接从深层网络传播回浅层网络,避免了在多层网络传播过程中梯度逐渐变小甚至消失的问题。

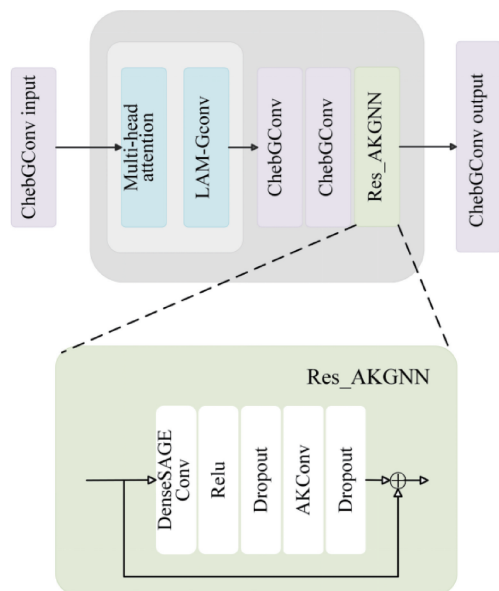


图 2 AK\_GraFormer 模块图

Figure 2 Network diagram of AK\_GraFormer

Res\_AKGNN 层使得模型能够根据节点的实际需求选择合适的信息层次进行决策。在手部姿态估计中,对于一些关键关节即靠近根节点的关节,需要更多的全局信息来确定其准确位置,而对于一些末

稍关节即如指间关节,局部信息就足够。AKGNN的全局读出函数可以帮助模型更好地处理这种差异。因此,Res\_AKGNN层对经过前面层处理后的特征数据进行特定的特征融合、变换或者信息传递等操作,增强了整个模型对基于图结构的输入数据的特征学习和表达能力,提高模型对于三维姿态估计的准确性。

### 1.3 评估指标动态监控模块

在模型训练阶段,通常会观察到训练过程呈现一定的收敛趋势,且相关评估指标处于持续向好的状态。然而,实际训练过程中可能会出现某一训练轮次(*Epoch*)下,评估指标突然变差的情况。在本模型实际训练过程中,每一轮训练结束后即验证模型效果。随着训练*Epoch*的增加,可以明显观察到验证结果中二维损失的增速明显大于三维损失。这是因为在训练过程中,KeyPoint R-CNN的特征空间和AK\_GraFormer模块的决策边界需要协同收敛。若KeyPoint R-CNN特征持续变化,AK\_GraFormer模块需不断适应新特征分布,增加了学习难度。此外,Vasconcelos等<sup>[23]</sup>研究指出关键点检测器与后续模块联合训练时,阶段性冻结能提升最终精度。

针对这一现象,本文设计了一种评估指标监控动态训练策略,对KeyPoint R-CNN模块的模型权重采取动态冻结处理。达到一定条件后,该模块在后续的训练过程中保持固定,不再参与更新迭代。若当前*Epoch*验证集 $MPJPE_{current}$ (当前轮次的平均每关节误差)值大于上一个*Epoch*验证集 $MPJPE_{last}$ 时,则将上一个*Epoch*模型参数 $Model\ Parameters_{last}$ 赋值给当前*Epoch*模型参数 $Model\ Parameters_{frozen}$ ,并且将KeyPoint R-CNN模块的参数( $Param \in Model\ Parameters_{frozen}^{backboneUrn}$ )冻结。因此,在后续训练中不再计算该模块的梯度( $Param.requires\_grad$ )。本文通过这种方式避免异常训练轮次所带来的不良影响,保障模型训练的稳定性与有效性,进而提升手部姿态估计任务的整体性能表现。具体算法如下。

**算法 1:** 评估指标动态调节算法。

**输入:** 当前*Epoch*验证集 $MPJPE_{current}$ ;

上一个*Epoch*验证集 $MPJPE_{last}$ ;

上一个*Epoch*模型参数 $Model\ Parameters_{last}$ 。

**输出:** 冻结权重后的参数 $Model\ Parameters_{frozen}$ 。

- ① If  $MPJPE_{current} > MPJPE_{last}$  do
- ②  $Model\ Parameters_{frozen} \leftarrow Model\ Parameters_{last}$ ;
- ③ For Param in  $Model\ Parameters_{frozen}$  do
- ④ If Param  $\in Model\ Parameters_{frozen}^{backboneUrn}$  do
- ⑤ Param.requires\_grad = False

⑥ return  $Model\ Parameters_{frozen}$ .

### 1.4 损失函数

为了有效监督网络训练,减少预测坐标值与真实坐标值的差距,本文使用两个损失函数。总损失定义为

$$L = \lambda_1 L_{2D} + \lambda_2 L_{3D}, \quad (3)$$

式中:通过实验观察不同权重组合对验证效果的影响逐步调整权重值,最终权重值设置为 $\lambda_1 = 0.1$ , $\lambda_2 = 0.9$ ;  $L_{2D}$ 为二维姿态估计损失;  $L_{3D}$ 为三维姿态估计损失。

$$L_{2D} = \frac{1}{n} \sum_{i=1}^n \|J_i - \hat{J}_i\|^2. \quad (4)$$

式中:  $J_i$ 表示标注的真实的第*i*个关节的二维坐标;  $\hat{J}_i$ 表示预测的第*i*个关节二维坐标。

$$L_{3D} = \frac{1}{n} \sum_{i=1}^n \|P_i - \hat{P}_i\|^2. \quad (5)$$

式中:  $P_i$ 表示标注的真实的第*i*个关节的三维坐标,  $\hat{P}_i$ 表示预测的第*i*个关节三维坐标。

## 2 实验

### 2.1 数据集和评价方法

本文基于HO-3D\_v3数据集<sup>[24]</sup>和FreiHand数据集<sup>[25]</sup>进行实验。HO-3D\_v3数据集是一个手和物体在严重相互遮挡情况下的3D姿态注释的数据集。该数据集中的68个序列包含10个不同的人操作10个不同的物体,这些物体来自YCB物体数据集。其包含77 558张图像的注释,这些图像被分为66 034张训练图像(来自55个序列)和11 524张评估图像(来自13个序列)。FreiHand数据集于2019年发布,包括32个对象的手势。该数据集具有背景多样、视点多变、动作复杂的特点,其中部分手势是手部与物体交互的动作,其包括13 024个训练样本和3 960个评估样本。

本文使用以下3个评估指标:平均每关节位置误差( $MPJPE$ ),通过计算预测和实际三维关节位置之间的欧几里德距离度量结果好坏; $PA-MPJPE$ 是 $MPJPE$ 的改良版,采用普氏分析对预测姿态与真实姿态进行刚性对齐(旋转、平移),消除全局位置与方向变化的影响,从而专注于评估关节间的相对姿态误差;检测关节点百分比( $PCK$ )曲线下面积( $AUC$ ),其中 $PCK$ 指正确3D关节的百分比。

### 2.2 实验环境与参数

本文实验环境:操作系统为Ubuntu 2022.04 LTS、GPU配置为RTX-4090D、CPU配置为AMD

EPYC 9754、内存大小为 60 GB。神经网络的实现框架为 PyTorch 和 PyTorch Geometric 框架。

本文的实验主要参数:输入的手部按 Mesh2D 计算训练验证检测框数据。*train\_batch* 大小设置为 32,*test\_batch* 大小设置为 1,使用 Adam 作为优化器,初始学习率设置为 0.000 1。

### 2.3 可视化分析

本文模型在 HO3D v3 数据集上的可视化结果如图 3 所示。原图中的矩形框用来表示模型检测到的手部位置,而真实值和预测值中的手部姿态通过不同的线条颜色区分不同手指。通过可视化展示,能够清晰地看到模型预测结果与真实数据之间的对比。观察这些图像时,可以发现手部关节在未被遮挡的情况下模型能够准确地预测关节位置,与真实值几乎重合。然而在发生手部和物体遮挡的情况下,预测值与真实值存在一定的偏差,但预测的关节图符合天然的手部结构的生理约束。因人工标注数据方法的成本高,HO3D v3 数据集采用机器标注方法,标注存在一定误差。因此在训练过程中,标注数据甚至对模型的学习产生了负面影响。

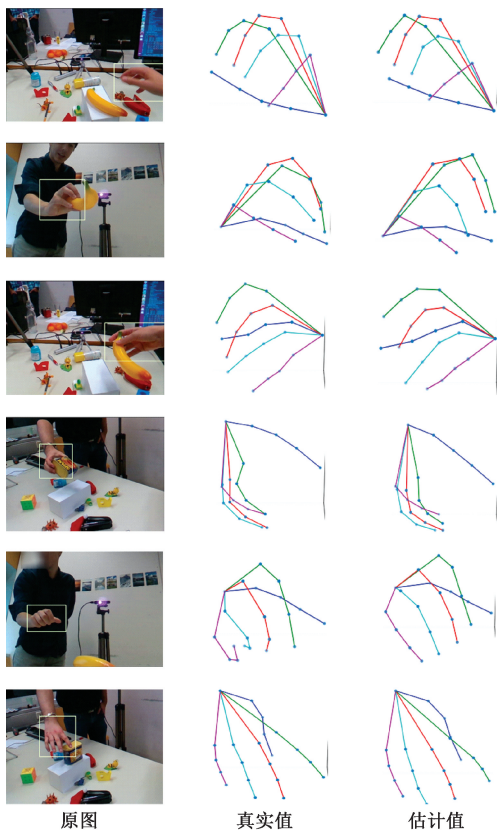


图 3 可视化结果

Figure 3 Visualization results

### 2.4 消融实验

本文在 HO-3D\_v3 数据集上进行了消融实验,以验证本文所提出模块的有效性。在本文提出的网

络中,采用分别叠加融合信息提取模块、AK\_GraFormer 模块和评估指标动态监控模块进行实验。消融实验结果如表 1 所示。其中在加入评估指标动态监控模块之前,模型训练过程中出现的验证损失值持续增长的情况如图 4 所示,可以清晰地看出 2D 验证损失随着训练轮次的增加而呈现上升趋势,模型对验证数据的拟合效果变差,其预测的准确性会随之降低,这对训练的影响是负面的。

表 1 各模块消融实验

模块名称	MPJPE	PA-MPJPE	AUC
基线模型 <sup>[13]</sup>	25.60	11.20	0.755
仅包含融合信息提取模块	24.29	9.88	0.798
融合信息提取模块+AK_GraFormer 模块	23.64	10.02	0.793
本文模型	22.36	9.45	0.812

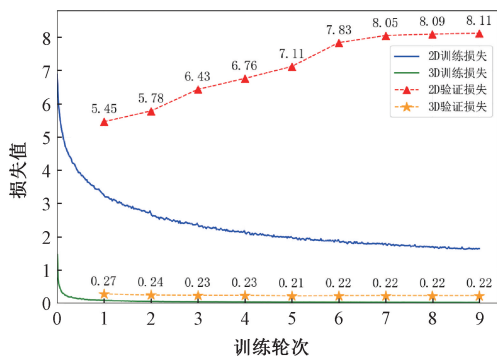


图 4 损失值变化情况

Figure 4 Changes in value of losses

从表 1 实验结果来看,基线模型误差较高,反映了其特征表达单一、关节间的拓扑结构未被充分挖掘。在加入融合信息提取模块后,模型对姿态的上下文感知能力变强,减少了因视角变化和遮挡等引起的相对姿态误差。而在此基础上加入 AK\_GraFormer 模块后,MPJPE 有所下降,但 PA-MPJPE 却略微上升。改进的 res\_AKGNN 层强化了网络的复杂姿态表达能力,从而降低了 MPJPE。但是由于两模块联合训练,导致特征漂移与参数更新异步性,最终影响模型精度。加入本文提出的全部模块后,模型误差显著降低,这表明评估指标动态监控模块的加入及时地调整了模块训练策略,确保各模块协同优化。

消融实验的可视化结果如图 5 所示,在手部姿态复杂和存在遮挡的情况下,本文方法依旧表现出良好的性能,估计出准确且符合手部结构特点的手部姿态。

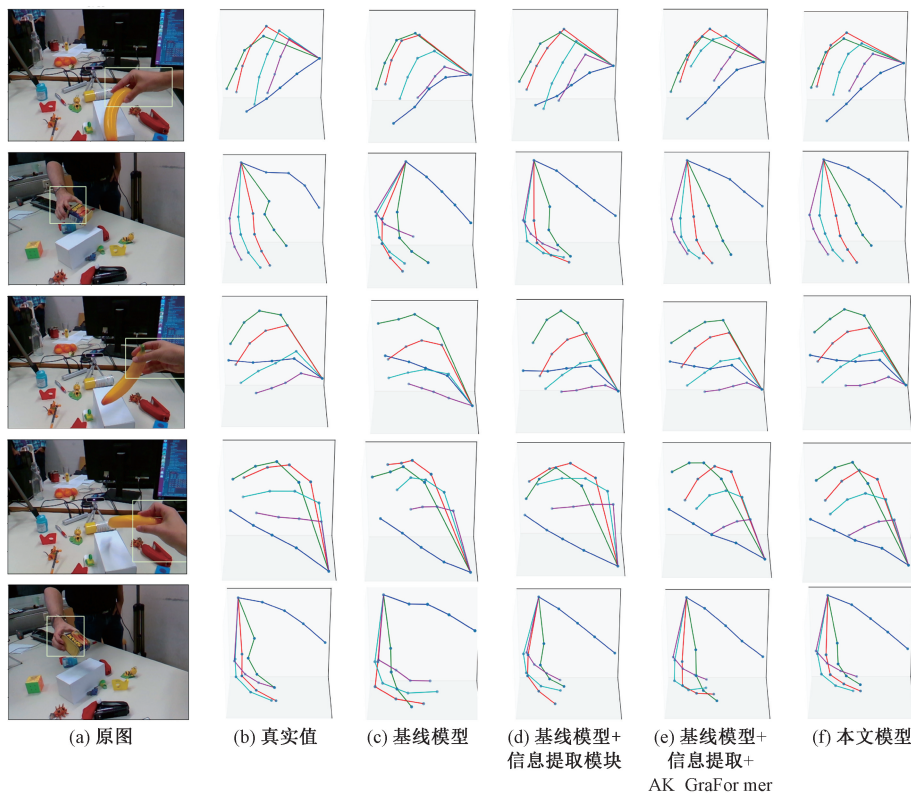


图 5 消融实验可视化结果

Figure 5 Visualization results of ablation study

2.5 定量比较

表 2 比较了本文提出的模型与其他模型在 HO3D v3 数据集上的性能。其中,为保证表格的整洁性,将杨冰等<sup>[26]</sup>提出的基于改进 FastMETRO 的三维手部姿态估计方法称为 FastMETRO。从评估的结果得知,本文模型在 *PA-MPJPE* 评估指标上取得最好结果。表 2 和表 3 是本文提出的模型与其他模型在 FreiHand 数据集上的比较结果,本文模型在所比较指标上的结果均为最优。尽管同类方法在重度遮挡场景下的手部姿态估计中展现出较低的平均关节位置误差,但却存在着手部拓扑结构失真的问题。而本文方法能够通过学习图像特征信息,优化图结构学习,提高姿态估计的准确度,同时使生成的手部姿态符合人手结构特点。

表 2 HO3D v3 数据集的定量比较

Table 2 Quantitative comparison of HO3D v3 datasets

模型名称	<i>PA-MPJPE</i>	<i>AUC</i>
S <sup>2</sup> HAND <sup>[27]</sup>	11.50	0.769
ArtiBoost <sup>[28]</sup>	10.80	0.785
PyMAF-X <sup>[29]</sup>	10.80	—
HMP <sup>[30]</sup>	10.10	—
FastMETRO <sup>[26]</sup>	10.50	—
本文模型	9.45	0.812

表 3 FreiHand 数据集的定量比较

Table 3 Quantitative comparison of FreiHand datasets

模型名称	<i>PA-MPJPE</i>	<i>AUC</i>
I2UV-HandNet <sup>[31]</sup>	6.70	0.856
Mesh Graphormer <sup>[32]</sup>	5.90	0.874
HandMIM <sup>[33]</sup>	6.29	0.871
HaMeR <sup>[34]</sup>	6.00	—
本文模块	5.78	0.878

3 结论

本文针对手部姿态估计提出基于图卷积的三维手部姿态估计算法。为了减少中间的二维位置对最终的三维位置的影响,同时为姿态估计网络提供更多辅助信息,本文通过提取输入图片的特征信息,得到图像视觉特征和二维位置信息;将提取的信息输入到 AK\_GraFormer 模块中,引入改进的 res\_AKGN 层,自适应学习输入数据更深层的特征表示,提高模块的表达能力,使得估计结果更精准;设计评估指标监控的动态训练策略,以应对训练过程中评估指标退化现象,避免过拟合,有效提升了模型性能。实验表明,本文方法取得了较好结果,相较于其他方法误差更小、手部结构更自然,为手部姿态估计提供了更有效的解决方案。

## 参考文献:

- [1] Sridhar S, Feit A M, Theobalt C, et al. Investigating the dexterity of multi-finger input for mid-air text entry [C]// Proceedings of the 33rd Annual ACM Conference on Human Factors in Computing Systems. New York: ACM, 2015: 3643–3652.
- [2] Oikonomidis I, Kyriazis N, Argyros A A. Tracking the articulated motion of two strongly interacting hands [C]// Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. Providence: IEEE, 2012: 1862–1869.
- [3] Tkach A, Pauly M, Tagliasacchi A. Sphere-meshes for real-time hand modeling and tracking [J]. ACM Transactions on Graphics. New York: ACM, 2016, 35 (6): 1–11.
- [4] ROMERO J, TZIONAS D, BLACK M J. Embodied Hands: Modeling and Capturing Hands and Bodies Together [J]. ACM Transactions on Graphics. New York: ACM, 2017, 36(6): 1–17.
- [5] Pavlakos G, Choutas V, Ghorbani N, et al. Expressive body capture: 3d hands, face, and body from a single image [C]// Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. Long Beach: IEEE, 2019: 10975–10985.
- [6] KESKIN C, KİRAÇ F, KARA Y E, et al. Hand Pose Estimation and Hand Shape Classification Using Multilayered Randomized Decision Forests [C]// Proceedings of the European Conference on Computer Vision. Berlin: Springer, 2012: 852–863.
- [7] TOMPSON J, STEIN M, LECUN Y, et al. Real-time Continuous Pose Recovery of Human Hands Using Convolutional Networks [J]. ACM Trans Graph, 2014, 33(5): 1–10.
- [8] Pan X, Li S, Wang H, et al. LGCANet: lightweight hand pose estimation network based on HRNet [J]. The Journal of Supercomputing, 2024(80): 1–23.
- [9] Hoang D C, Tan P X, Pham D L, et al. Efficient Multimodal Fusion For Hand Pose Estimation With Hourglass Network [J]. IEEE Access, 2024(12): 113810–113825.
- [10] Zhan Z, Luo G. Multiscale feature fusion network for monocular complex hand pose estimation [J]. Electronics Letters, 2023, 59(24): 1–4.
- [11] Panteleris P, Oikonomidis I, Argyros A. Using a single rgb frame for real time 3d hand pose estimation in the wild [C]// Proceedings of the 2018 IEEE winter conference on applications of computer vision. Lake Tahoe: IEEE, 2018: 436–445.
- [12] Doosti B, Naha S, Mirbagheri M, et al. Hope-net: A graph-based model for hand-object pose estimation [C]// Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. Seattle: IEEE, 2020: 6608–6617.
- [13] Zhao W, Wang W, Tian Y. Graformer: Graph-oriented transformer for 3d pose estimation [C]// Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. New Orleans: IEEE, 2022: 20438–20447.
- [14] 李志新, 商樊淇, 郇战, 等. 基于混合特征图卷积神经网络的人体行为识别方法 [J]. 郑州大学学报 (工学版), 2024, 45(04): 46–52.
- [15] Cai Y, Ge L, Liu J, et al. Exploiting spatial-temporal relationships for 3d pose estimation via graph convolutional networks [C]// Proceedings of the IEEE/CVF international conference on computer vision. Seoul: IEEE, 2019: 2272–2281.
- [16] Aboukhadra A T, Robertini N, Malik J, et al. Shape-GraFormer: GraFormer-Based Network for Hand-Object Reconstruction from a Single Depth Map [J]. IEEE Access, 2024.
- [17] Zhuang N, Mu Y. Joint hand-object pose estimation with differentially-learned physical contact point analysis [C]// Proceedings of the 2021 international conference on multimedia retrieval. New York: ACM, 2021: 420–428.
- [18] Zhang M, Li A, Liu H, et al. Coarse-to-fine hand-object pose estimation with interaction-aware graph convolutional network [J]. Sensors, 2021, 21(23): 8092.
- [19] 马胜蕾, 李敬华, 孔德慧, 等. 基于双分支多尺度注意力的手三维姿态估计 [J]. 计算机学报, 2023, 46(07): 1383–1395.
- [20] Yang W, Xie L, Qian W, et al. Coarse-to-fine cascaded 3D hand reconstruction based on SSGC and MHSA [J]. The Visual Computer, 2025, 41(1): 11–24.
- [21] He K, Gkioxari G, Dollár P, et al. Mask r-cnn [C]// Proceedings of the IEEE international conference on computer vision. Venice: IEEE, 2017: 2961–2969.
- [22] Ju M, Hou S, Fan Y, et al. Adaptive kernel graph neural network [C]// Proceedings of the AAAI Conference on Artificial Intelligence. Online: AAAI, 2022, 36(6): 7051–7058.
- [23] Vasconcelos C, Birodkar V, Dumoulin V. Proper reuse of image classification features improves object detection [C]// Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. New Orleans: IEEE, 2022: 13628–13637.
- [24] Hampali S, Sarkar S D, Lepetit V. Ho-3D\_v3: Improving the accuracy of hand-object annotations of the ho-3D dataset [J]. arxiv preprint arxiv:2107.00887, 2021.
- [25] Zimmermann C, Ceylan D, Yang J, et al. Freihand: A dataset for markerless capture of hand pose and shape-

- from single rgb images [C] // Proceedings of the IEEE/CVF International Conference on Computer Vision. Seoul: IEEE, 2019: 813–822.
- [26] 杨冰, 徐楚阳, 姚金良, 等. 基于单目 RGB 图像的三维手部姿态估计方法 [J]. 浙江大学学报(工学版), 2025, 59(01): 18–26.
- [27] Chen Y, Tu Z, Kang D, et al. Model-based 3d hand reconstruction via self-supervised learning [C] // Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. Online: IEEE, 2021: 10451–10460.
- [28] Yang L, Li K, Zhan X, et al. Artiboost: Boosting articulated 3d hand-object pose estimation via online exploration and synthesis [C] // Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. New Orleans: IEEE, 2022: 2750–2760.
- [29] Zhang H, Tian Y, Zhang Y, et al. Pymaf-x: Towards well-aligned full-body model regression from monocular images [J]. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2023, 45(10): 12287–12303.
- [30] Duran E, Kocabas M, Choutas V, et al. HMP: Hand motion priors for pose and shape estimation from video [C] // Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. Waikoloa: IEEE, 2024: 6353–6363.
- [31] Chen P, Chen Y, Yang D, et al. I2UV-HandNet: Image-to-uv prediction network for accurate and high-fidelity 3d hand mesh modeling [C] // Proceedings of the IEEE/CVF international conference on computer vision. Montreal: IEEE, 2021: 12929–12938.
- [32] Lin K, Wang L, Liu Z. Mesh graphormer [C] // Proceedings of the IEEE/CVF international conference on computer vision. Montreal: IEEE, 2021: 12939–12948.
- [33] Liu Z, Lin G, Wang C, et al. HandMIM: Pose-Aware Self-Supervised Learning for 3D Hand Mesh Estimation [J]. arxiv preprint arxiv:2307.16061, 2023.
- [34] Pavlakos G, Shan D, Radosavovic I, et al. Reconstructing hands in 3d with transformers [C] // Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. Seattle: IEEE, 2024: 9826–9836.

### 3D Hand Pose Estimation Based on Graph Convolution Network

PENG Chunyan<sup>1,2</sup>, WANG Xuan<sup>1,2</sup>, CHEN Yangbo<sup>1,2</sup>, HE Gangbo<sup>1,2</sup>

(1. College of Computer, Qinghai Normal University, Xining 810016, China; 2. The State Key Laboratory of Tibetan Intelligence, Qinghai Normal University, Xining 810016, China)

**Abstract:** In the task of 3D hand pose estimation from a single color image, challenges such as occlusion and high self-similarity of hand parts are faced, which lead to large prediction errors and unnatural hand structures. To address these issues, a graph convolution-based 3D hand pose estimation method is firstly proposed. Visual features and 2D keypoint positions are extracted from the input image using Keypoint R-CNN. These features are then fed into an improved Adaptive Kernel Graph Convolution module (AK\_GraFormer). Subsequently, a residual-connected AKGNN graph kernel is introduced to adaptively process graph-structured data, thereby enhancing the model's feature learning and representation. Finally, a dynamic training strategy is employed, which is monitored by a proposed evaluation metric, to optimize estimation performance. Experimental results on the HO3D v3 and FreiHand datasets demonstrate that the proposed method outperforms existing approaches in monocular 3D hand pose estimation. Specifically, the Procrustes-Aligned Mean Per Joint Position Error (PA-MPJPE) is reduced by up to 17.83 percentage points, and the Area Under the Curve (AUC) of the Percentage of Correct Keypoints (PCK) metric is improved by up to 5.59 percentage points compared to state-of-the-art methods.

**Keywords:** 3D hand pose estimation; graph convolution networks; feature extraction; optimisation of graph kernel learning; dynamic adjustment of assessment indicators