

# 基于图注意力机制的自监督单目深度估计

舒 航, 丁坤岭, 王晓峰\*, 陆正霖, 冯强龙, 王海宇

重庆科技大学电子与电气工程学院, 重庆

收稿日期: 2025年1月20日; 录用日期: 2025年3月5日; 发布日期: 2025年3月19日

## 摘 要

为解决单目自监督深度估计在缺乏准确真实深度时对复杂场景几何结构的刻画不足问题, 本文在原有的基于图卷积网络(Graph Convolutional Network, GCN)的单目深度估计框架上, 引入了图注意力网络(Graph Attention Network, GAT)机制, 提出了一种GATDepth模型。该模型通过在解码器阶段采用图注意力模块, 能够自适应地为相邻节点分配不同权重, 从而更精细地保留场景中的几何拓扑关系与不连续性。DepthNet编码器利用CNN提取多层次视觉特征, 而解码器则结合转置卷积上采样和GAT模块融合节点特征。通过目标图像与重构图像之间的光度、重投影及平滑性等多重损失进行自监督训练, 模型在KITTI数据集上取得了优异的深度估计性能, 尤其在远距物体和物体边缘等关键区域表现突出。实验结果表明, 所提方法不仅在保证网络效率的同时更好地捕捉了场景关键几何信息, 而且在缺乏高质量真实深度的条件下仍能获得可靠且精细的深度预测。

## 关键词

图卷积网络, 图注意力网络, GAT模块, 自监督训练

# Self-Supervised Monocular Depth Estimation Based on Multi-Scale Graph Attention Mechanism

Hang Shu, Kunling Ding, Xiaofeng Wang\*, Zhenglin Lu, Qianglong Feng, Haiyu Wang

School of Electronics and Electrical Engineering, Chongqing University of Science & Technology, Chongqing

Received: Jan. 20<sup>th</sup>, 2025; accepted: Mar. 5<sup>th</sup>, 2025; published: Mar. 19<sup>th</sup>, 2025

## Abstract

To address the issue of insufficient depiction of complex scene geometry in monocular self-

\*通讯作者。

文章引用: 舒航, 丁坤岭, 王晓峰, 陆正霖, 冯强龙, 王海宇. 基于图注意力机制的自监督单目深度估计[J]. 人工智能与机器人研究, 2025, 14(2): 313-319. DOI: 10.12677/airr.2025.142031

supervised depth estimation due to the lack of accurate ground truth depth, this paper proposes a GATDepth model based on the existing monocular depth estimation framework using Graph Convolutional Networks (GCN). The Graph Attention Network (GAT) mechanism is introduced into the model. By adopting graph attention modules in the decoder stage, the model can adaptively assign different weights to adjacent nodes, thereby more finely preserving the geometric topology and discontinuities in the scene. The DepthNet encoder extracts multi-level visual features using CNNs, while the decoder combines transposed convolutional upsampling and GAT modules to fuse node features. The model is trained in a self-supervised manner through multiple losses such as photometric, reprojection, and smoothness losses between the target image and the reconstructed image. The model achieves excellent depth estimation performance on datasets such as KITTI, especially in key areas such as distant objects and object edges. Experimental results show that the proposed method not only better captures key geometric information of the scene while ensuring network efficiency, but also obtains reliable and fine depth predictions even in the absence of high-quality ground truth depth.

## Keywords

Graph Convolutional Networks, Graph Attention Networks, GAT Modules, Self-Supervised Training

Copyright © 2025 by author(s) and Hans Publishers Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

## 1. 引言

### 1.1. 概述

在人工智能领域,深度学习网络在深度估计和自运动预测等任务中取得了卓越的表现,并被广泛应用于自动驾驶车辆[1][2]和物体距离预测[3]。其中,单目深度估计可通过只使用单目摄像机来推断场景的几何信息,有望进一步降低硬件成本并简化部署。相比于传统的立体视觉方法,单目深度估计在获取训练数据时无需同时依赖多目相机或激光雷达,从而具有更高的灵活性。然而,如何在缺乏准确地面真值深度的情况下训练模型仍是一大挑战。

早期的单目深度估计大多采用监督学习方法[4],需要利用昂贵的 2D/3D 激光雷达(LiDAR)采集数据,或者使用计算机图形模拟引擎渲染高质量的真实深度。这些方法的局限性在于采集或生成大规模、高分辨率、高质量的地面真值数据非常困难。为克服此限制,自监督单目深度估计方法应运而生,这类方法通过比较目标图像与源图像重构结果的差异来提供监督信号,避免了对地面真值深度的显式依赖。典型地,网络通常需要在训练时同时估计深度和相机位姿,并在训练过程中最小化源帧投影到目标帧上的光度重建误差。

虽然许多近期研究表明,基于卷积神经网络(Convolutional Neural Network)的自监督深度估计能够取得令人满意的结果,但 CNN 本质上是欧几里得域(Euclidean domain)上的操作,对于图数据或需要保留复杂几何拓扑结构的任务, CNN 的局限性开始显现。为此,几何深度学习[5]将注意力聚焦于图卷积网络(Graph Convolutional Network, GCN)等方法,让网络能够处理非欧几里得域的数据,显著提升了模型对几何结构的理解能力。近期的工作已证明,将 GCN 用于深度估计可以更好地维护场景中物体间的拓扑关系并突出其几何特征。

然而,在实际应用中,并非所有邻居节点在特征传播时的重要性相同。图注意力网络(Graph Attention

Network, GAT) [6]的出现, 为图数据中的信息交互提供了更灵活、更具表现力的方式。GAT 通过可学习的注意力系数来衡量各节点对中心节点的影响程度, 在更新节点表征时赋予邻居节点不同的权重, 从而更精准地捕捉局部拓扑结构和像素间的关联性。

## 1.2. 本文主要改进点

为此, 本文在原有的 GCNDepth 框架中, 进一步引入多尺度图注意力网络(Graph Attention Networks)模块, 形成了新的 GATDepth 模型, 以期在自监督单目深度估计任务中取得更高精度。本文的主要贡献如下:

(1) 提出了基于图注意力网络(GAT)的自监督深度预测方法。在构建深度图的迭代过程中, 通过学习可区分的邻居注意力权重来传播节点信息, 有效提升深度预测精度。

(2) 在解码器网络的中引入 GAT 结构, 更全面地利用像素在不同语义层次上的空间相关性与注意力分配, 显著增强了网络对复杂场景的深度理解能力。

(3) 提出了结合光度、重投影和平滑性相关的多种损失函数, 用于提升预测深度图的质量。其中, 重投影损失可处理物体遮挡问题, 重构损失用于减少目标图像与重构图像间的差异, 而平滑性损失则可保留物体边缘并减小纹理区域对深度估计的干扰。

## 2. 背景与相关工作

### 2.1. 监督深度估计

在单目深度估计领域, 早期大多数工作都是基于监督学习, 需要使用配有地面真值的图像对来训练网络。许多研究通过端到端的方式[7]、局部预测与融合[8]以及非参数场景采样等技术, 成功地将输入图像映射到相应的深度图。然而, 地面真值的缺失或难以获取, 依然是此类方法在大规模实际应用中的主要障碍。

### 2.2. 自监督深度估计

近年来, 自监督单目深度估计方法避免了对真实深度数据的显式依赖, 而是借助目标图像与重构图像的差异来作为训练信号, 通常可分为立体训练和单目视频训练。其中, 单目视频训练场景更为通用, 通过在网络中引入相机位姿估计器[9], 将任意连续帧的相机运动纳入训练过程, 从而为深度学习模型提供几何约束。一些工作引入了边缘一致性、深度归一化层、时域信息等, 以进一步提高单目自监督深度估计的性能。

### 2.3. 图神经网络与图注意力网络

图卷积网络(Graph Convolutional Networks) [5]能够在图数据上进行卷积操作, 被广泛应用于半监督节点分类和其他需要保留结构信息的任务。其核心思路是通过“邻居聚合”机制来更新中心节点表征。然而, 传统 GCN 在更新节点特征时, 默认给予邻居节点相同的权重, 无法根据具体上下文或像素差异性进行区分。

图注意力网络  $X$  则在 GCN 的基础上, 引入了可学习的注意力系数。对于节点  $i$  的邻居节点集合  $N(i)$ , GAT 首先对每个邻居节点  $j \in N(i)$  计算注意力分数  $\alpha_{ij}$ , 再将所有邻居节点特征聚合成新的中心节点特征, 公式示例为:

$$h_i = \sigma \left( \sum_{j \in N(i)} \alpha_{ij} W h_j \right) \quad (1)$$

其中  $W$  为可学习的线性变换,  $\alpha_{ij}$  由注意力机制(例如一个前馈网络 + Softmax)给出。这使得网络在更新

节点特征时，能自适应地关注最相关的邻居节点，大幅提高了对复杂几何拓扑结构的表示能力。

在自监督单目深度估计中，若能利用 GAT 对像素间的相似度与差异度进行更精准的衡量，就有望获得更高精度的深度图，同时保留场景中重要的边界、不连续性与拓扑结构信息。

### 3. 改进方法

在本节中，我们将介绍提出模型的整体框架 GATDepth。该模型由两个主要网络组成：DepthNet (带有 GAT 解码器) 和 PoseNet。DepthNet 负责生成单帧图像的深度图，而 PoseNet 负责预测连续帧之间的相机位姿。最终，通过对源帧投影到目标帧上的重构误差进行最小化来实现自监督训练。

#### 3.1. 问题定义

令  $I \in A$  表示单目 RGB 图像，对应的深度图  $D \in B$  是我们想要预测的目标。可定义一个映射函数  $\Psi_D: A \rightarrow B$ ，将输入域  $A$  的元素(RGB 图像)映射到深度域  $B$ ，公式如下：

$$D(P) = \Psi_D(I_s(P)) \quad (2)$$

其中  $P$  表示图像中的像素坐标。

为了在单目视频训练中估计相机运动，我们还定义函数  $\Psi_E: A \times A \rightarrow \mathbb{R}^3$ ，将一对连续帧  $(I_s, I_t)$  映射到旋转和平移向量  $(r^T, t^T)$ ，记为  $E_{I_s \rightarrow I_t}$ ：

$$E_{I_s \rightarrow I_t} = \Psi_E(I_s, I_t) \quad (3)$$

有了深度图  $D$  和位姿  $E_{I_s \rightarrow I_t}$  后，即可将源帧  $I_s$  投影到目标帧  $I_t$  的坐标系中，生成重构图像  $I_{rec}$ ，并通过最小化  $I_{rec}$  与  $I_t$  之间的差异实现网络的自监督训练。

##### 3.1.1. 图注意力网络

在原有 GCNDepth 中，我们使用了图卷积网络(GCN)来在解码器阶段融合图结构特征。为更好地区分每个像素(节点)在邻域聚合时的重要性，本文引入图注意力网络(GAT)。GAT 的核心是学习到注意力系数  $\alpha_{ij}$ ，并对每个邻居节点赋予不同的权重，从而在解码器阶段更精准地聚合邻居节点的信息。

同 GCNDepth 中的做法一样，我们将编码器提取到的特征图视作一个包含  $N$  个节点的图，每个节点对应特征图中的一个像素位置(或在更粗层的特征图上)。初始邻接矩阵  $A \in \mathbb{R}^{N \times N}$  可由像素间相似度或空间邻接关系确定。

##### 3.1.2. 注意力机制

对于中心节点  $i$  以及其邻居节点  $j \in N(i)$ ，我们通过线性变换将节点表征  $h_i$  和  $h_j$  分别映射到相同空间，然后计算注意力打分  $e_{ij}$ 。一个常见的做法是将  $[Wh_i \parallel Wh_j]$  送入一个前馈网络(或简单的可学习向量)并通过 LeakyReLU 获取分数：

$$e_{ij} = \text{LeakyReLU}\left(a^T [Wh_i \parallel Wh_j]\right) \quad (4)$$

对所有邻居节点的  $e_{ij}$  进行 Softmax 归一化后得到注意力系数  $\alpha_{ij}$ ：

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k \in N(i)} \exp(e_{ik})} \quad (5)$$

##### 3.1.3. 节点特征更新

有了注意力系数  $\alpha_{ij}$  后，节点  $i$  的更新公式为：

$$h_i = \sigma \left( \sum_{j \in N(i)} \alpha_{ij} W h_j \right) \quad (6)$$

其中  $\sigma$  是非线性激活函数(如 ReLU 或 LeakyReLU)。

通过注意力机制,网络在解码器阶段可以更加灵活地“选择性”关注特定邻居节点的信息,从而在深度预测时更好地保留物体边缘、纹理不连续以及复杂的几何结构。

### 3.2. 自监督 CNN-GCN-GAT 自动编码器

#### 3.2.1. DepthNet 编码器

与原先的 GCNDepth 类似, DepthNet 的编码器用于从输入图像  $I_s$  中提取多尺度视觉特征。我们采用 CNN (ResNet-50)作为编码器,连续进行卷积、批量归一化、池化等操作,得到多个不同分辨率的特征图。

#### 3.2.2. DepthNet 解码器

在解码器中,我们使用多尺度的 GCN-GAT 模块来替代传统的仅依赖转置卷积的上采样过程。具体地,每个解码层都包含以下步骤:

上采样(转置卷积): 将上一层的特征图进行上采样,得到更高分辨率的粗糙深度预测。

跨层特征拼接: 与编码器对应层的特征图进行拼接,融合高层语义与低层细节。GAT 模块: 将拼接后的特征图视为图结构数据(节点对应空间像素点),通过注意力机制对节点特征进行更新,生成更精细的深度图预测。

这样逐层上采样并叠加 GAT 操作,可在每个尺度下针对邻域像素进行更灵活、更准确的特征聚合,得到的深度图能更好地保留对象边界和几何不连续性。

#### 3.2.3. PoseNet 估计器

PoseNet 用于预测源图像  $I_s$  和目标图像  $I_t$  之间的旋转和平移向量。PoseNet 同样由编码器和解码器构成,编码器使用 ResNet-18,输入是拼接后的图像对  $[I_s, I_t]$ 。解码器则通过一系列卷积层,回归输出  $E_{I_s \rightarrow I_t} = [r^T, t^T]$ 。

### 3.3. 几何模型与损失函数

与原始 GCNDepth 类似,我们在自监督框架下定义以下主要损失函数:

#### 3.3.1. 重构损失函数 $L_{Rec}$

通过计算重构图像  $I_{Rec}$  与目标图像  $I_t$  之间的像素差异来约束深度预测准确度:

$$L_{Rec} = \sum_P |I_{Rec}(P) - I_t(P)| \quad (7)$$

#### 3.3.2. 重投影损失函数 $L_{P1}$

为了处理单目视频中可能出现的遮挡与视差变化,结合  $L_1$  范数和结构相似性指数:

$$L_{P1} = 0.15 \sum_P |I_{Rec}(P) - I_t(P)| + 0.85 \sum_P \frac{1 - SSIM(I_{Rec}, I_t)}{2} \quad (8)$$

#### 3.3.3. 平滑损失函数 $L_{Smooth}$

为保留物体边缘并减少纹理区域的影响,引入判别性损失  $L_{Dis}$  和曲率损失  $L_{Cvt}$ :

$$L_{Dis} = \sum_P e^{-\lambda \|\nabla^1 I_s(P)\|} \|\nabla^1 D(P)\| \quad (9)$$



$$L_{Cvt} = \sum_P e^{-\lambda \|\nabla^2 I_s(P)\|} \|\nabla^2 \mathbf{D}(P)\| \quad (10)$$

平滑损失定义为  $L_{Smooth} = \alpha L_{Dis} + \beta L_{Cvt}$ ，其中  $\alpha$  和  $\beta$  为权重系数。

### 3.3.4. 总损失函数

总损失  $L_{Final}$  由三部分组成： $L_{P1}$  (规划损失)、 $L_{Rec}$  (重建损失)和  $L_{Smooth}$  (平滑损失)。

$$L_{Final} = L_{P1} + L_{Rec} + L_{Smooth} \quad (11)$$

## 4. 实验

本节将展示在 KITTI 数据集上的实验结果，并与基准模型 GCNDepth 进行比较，验证所提 GATDepth 的有效性。

**数据与设置：**按照 Eigen 切分，在去除静态帧后，训练集中约 39,810 张图像，验证集 4424 张，测试集 697 张。图像分辨率统一为  $1024 \times 320$  像素，深度范围限制为 80 米，并使用中值缩放对预测深度进行对齐。

**评价指标：**采用绝对误差(abs-rel)、平方相对误差(sq-rel)、均方根误差(rmse)、对数均方根误差(rmse-log)来评估模型质量。

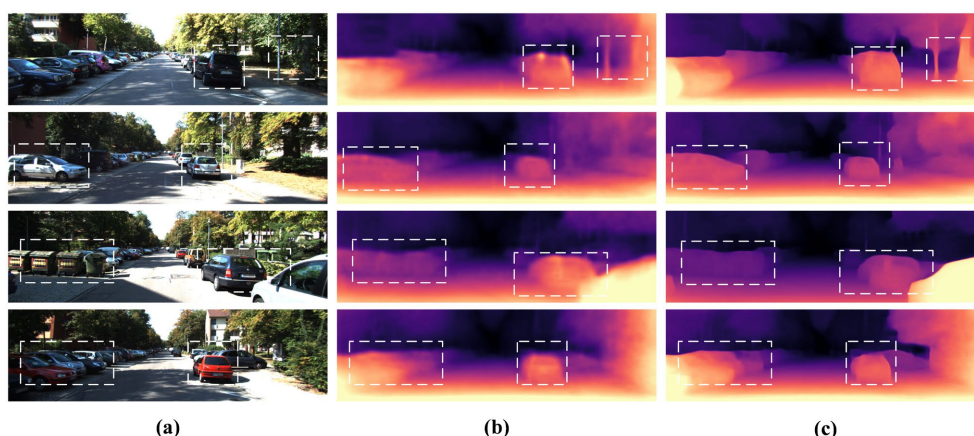
表 1 展示了与 GCNDepth 模型的对比。可以看到，GATDepth 在 abs-rel、sq-rel 等关键指标上取得了或超越了最新方法的性能，同时在 rmse 和 rmse\_log 指标上也有明显提升。与我们的先前工作 GCNDepth 相比，引入 GAT 后，模型在处理远距离目标以及场景中的遮挡与不连续区域时更具鲁棒性。

**Table 1.** Comparison of model performance between this model and GCNDepth model

**表 1.** 本文模型和 GCNDepth 模型性能对比

模型名称	abs_rel ↓	sq_rel ↓	rmse ↓	rmse_log ↓	a1 ↑	a2 ↑	a3 ↑
GCNDepth	0.115	0.882	4.701	0.190	0.879	0.961	0.982
GCN-GATDepth	<b>0.108</b>	<b>0.815</b>	<b>4.689</b>	<b>0.172</b>	<b>0.892</b>	<b>0.972</b>	<b>0.986</b>

注：↓表示数值越低越好；↑表示数值越高越好；abs\_rel 表示绝对相对误差；sq\_rel 表示平方相对误差；rmse 表示均方根误差；rmse\_log 表示对数均方根误差；a1、a2 和 a3 分别表示相对误差小于 1.25、1.25<sup>2</sup> 和 1.25<sup>3</sup> 的比例。



**Figure 1.** Comparison of depth prediction between the present model and GCNDepth model. (a) The input artwork; (b) The depth map predicted by the benchmark model GCNDepth; (c) The depth map predicted by the model in this paper

**图 1.** 本文模型和 GCNDepth 模型深度预测对比图。(a) 输入原图；(b) 基准模型 GCNDepth 预测的深度图；(c) 本文模型预测的深度图

定性结果: 图 1 对比了 GATDepth 与基准方法在测试集中部分样例的深度预测情况。可以看出, GAT 机制帮助模型更好地检测出物体轮廓和边缘, 使得在远处车辆、行人以及路边区域都能得到更准确的深度估计。GAT 机制在以下区域的可视化效果较为优异: 第一行中汽车前方的灌木丛、第二行中汽车前方的树木、第三行中左侧汽车前方灌木丛的阴影区域、右侧汽车前方的灌木丛阴影区域树干; 第四行中左侧汽车前方的灌木丛树干、图片右上角灌木丛轮廓。与未加入 GAT 机制相比在深度图细节信息结构性、可视化展示上表现更好。

说明本文引入 GAT 机制, 能够有效的避免细节丢失的现象, 提高模型对于全局细节特征的感知能力, 可以实现对于全局和局部多尺度上下文信息的有效利用, 有效提升模型对于深度信息估计的精度与整体的泛化性能。

## 5. 结论

本文在原有 GCNDepth 框架基础上, 进一步提出了结合图注意力网络的自监督深度估计模型——GATDepth。通过在多尺度解码器中加入 GAT 模块, 使网络能够在节点聚合时自适应地分配邻居权重, 从而更精准地捕捉场景中的几何拓扑关系并突出关键对象边界。实验结果表明, GATDepth 在 KITTI 数据集上取得了更高精度的深度预测, 同时保留了网络的高效性, 对大规模应用具有现实意义。

## 参考文献

- [1] 樊振宇. 基于深度学习的自动驾驶汽车周边车辆轨迹预测方法研究[D]: [硕士学位论文]. 镇江: 江苏大学, 2022.
- [2] 黄峻, 田永林, 戴星原, 等. 基于深度学习的自动驾驶多模态轨迹预测方法: 现状及展望[J]. 智能科学与技术学报, 2023, 5(2): 180-199.
- [3] 谭紫阳, 高忠文, 邓宇. 基于改进极限学习机和深度神经网络融合的车辆轨迹长期预测[J]. 汽车技术, 2020(11): 16-20.
- [4] 江铃燚, 郑艺峰, 陈澈, 等. 有监督深度学习的优化方法研究综述[J]. 中国图象图形学报, 2023, 28(4): 963-983.
- [5] Wu, Z., Pan, S., Chen, F., Long, G., Zhang, C. and Yu, P.S. (2021) A Comprehensive Survey on Graph Neural Networks. *IEEE Transactions on Neural Networks and Learning Systems*, **32**, 4-24. <https://doi.org/10.1109/tnnls.2020.2978386>
- [6] Rao, V.D., Zhang, Z.W. and Leskovec, J. (2021) Graph Attention on Point Clouds: A Survey. *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 11-17 October 2021, 3482-3491.
- [7] Leskovec, J., Sen, R., et al. (2020) End-to-End Graph Learning for Large-Scale Drug Discovery. *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 6-12 December 2020, 1-12.
- [8] Shao, W.Q., Bai, K., Zhang, S.Q., et al. (2021) Local Feature Fusion with Predictive Learning for 3D Object Detection. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 19-25 June 2021, 14122-14132.
- [9] Zhang, Y., Sun, Z., et al. (2019) Visual Inertial Camera Pose Estimation with Convolutional Neural Networks. *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, Montreal, 20-24 May 2019, 2155-2161.