

東 南 大 學

毕业设计(论文)报告

题 目: 基于深度学习的三维重建和点云生  
成

学 号: 04216747

姓 名: 周烨凡

学 院: 信息科学与工程学院

专 业: 信息工程

指导教师: 杨绿溪

起止日期: 2020 年 1 月 -2020 年 5 月

## 东南大学毕业（设计）论文独创性声明

本人声明所呈交的毕业（设计）论文是我个人在导师指导下进行的研究工作及取得的研究成果。尽我所知，除了文中特别加以标注和致谢的地方外，论文中不包含其他人已经发表或撰写过的研究成果，也不包含为获得东南大学或其它教育机构的学位或证书而使用过的材料。与我一同工作的同志对本研究所做的任何贡献均已在论文中作了明确的说明并表示了谢意。

论文作者签名：\_\_\_\_\_ 日期：\_\_\_\_年\_\_\_\_月\_\_\_\_日

## 东南大学毕业（设计）论文使用授权声明

东南大学有权保留本人所送交毕业（设计）论文的复印件和电子文档，可以采用影印、缩印或其他复制手段保存论文。本人电子文档的内容和纸质论文的内容相一致。除在保密期内的保密论文外，允许论文被查阅和借阅，可以公布（包括刊登）论文的全部或部分内容。论文的公布（包括刊登）授权东南大学教务处办理。

论文作者签名：\_\_\_\_\_ 导师签名：\_\_\_\_\_

日期：\_\_\_\_年\_\_\_\_月\_\_\_\_日 日期：\_\_\_\_年\_\_\_\_月\_\_\_\_日

## 摘要

单视角三维物体形状重建 (Single-view 3D Reconstruction) 是三维视觉领域一直以来的核心问题之一。由于拟合非线性方程和学习模式的有效性，深度神经网络被期望在重建三维非线性形状的任务中表现出色。本课题探究了深度神经网络在三维重建任务中的有效性和内在机制，主要的工作和创新如下：

1. 针对基于自编码器结构的深度神经网络梯度过小等问题，提出采用多任务训练对编码器进行初始化，以及引入残差层这两种方法优化网络架构。在 ShapeNet 公开数据集上进行测试，取得了优于基于检索的非学习方法 Oracle Nearest Neighbor 的结果。
2. 为了衡量三维数据集的聚类程度，提出了基于 affinity propagation 和 silhouette score 的度量标准，该方法得到的量化结果与数据集降维后的可视化结果吻合。
3. 提出了影响三维重建任务中深度神经网络内在机制的主要因素是数据集的本质特征：当训练集中的三维形状集有相比于图片集更高的聚类程度的结构时，在这个数据集上训练的深度神经网络更有可能会执行识别机制而非重建机制，并定性和定量的证明了这一强相关性。

最后，本文分析了上述结论对三维重建任务中的数据收集和神经网络训练的指导意义。

关键词：三维重建，点云，深度学习，数据挖掘

## ABSTRACT

Single-view 3D reconstruction is one of the elementary tasks in the field of 3D vision. Due to the effectiveness in learning and approximating non-linear function, deep neural networks are expected to perform well on the task of reconstructing 3D nonlinear shapes. This paper investigates the effectiveness and internal mechanism of deep neural network in 3D reconstruction task. The main contribution and innovation are summarized as follows:

1. For the problem that the variance of multiple training results of deep neural network based on auto-encoder structure is too large, it proposes to leverage multi-task training to initiate the encoder and introduce residual layer to deepen the decoder in order to optimize network architecture. These methods are trained and tested on ShapeNet public dataset and they outperform the non-learning method based on retrieve, namely Oracle Nearest Neighbor.
2. Define novel way to measure clustering coefficient of 3D reconstruction dataset based on affinity propagation and silhouette score, the quantitative results of this metric is corresponding to qualitative results of visualization of dataset in low dimension.
3. It claims that the bias of internal mechanism of network is mainly affected by the intrinsic properties of dataset: when the training set of 3D shapes has a more clustered structure than images, the deep neural networks trained on this dataset become more likely to perform recognition than reconstruction, and it proves the strong correlation between the deep neural networks and dataset property both qualitatively and quantitatively.

Finally, this paper analyzes the significance of the above conclusions for data collection and neural network training.

KEY WORDS: 3D Reconstruction, Point Cloud, Deep Learning, Data mining

# 目 录

摘要 . . . . .	I
<b>ABSTRACT</b> . . . . .	II
目录 . . . . .	III
<b>第一章 绪论</b> . . . . .	1
1.1. 课题的背景和意义 . . . . .	1
1.2. 研究现状 . . . . .	1
1.2.1. 相关模型 . . . . .	1
1.2.2. 相关数据集 . . . . .	2
1.2.3. 神经网络学习机制 . . . . .	2
1.3. 本文研究内容 . . . . .	3
1.3.1. 课题关键问题以及难点 . . . . .	3
1.3.2. 主要贡献 . . . . .	3
1.4. 论文组织结构 . . . . .	3
<b>第二章 问题与定义</b> . . . . .	5
2.1. 问题描述 . . . . .	5
2.2. 识别或重建 . . . . .	6
2.3. 数据集的聚类趋势 . . . . .	6
<b>第三章 模型方法</b> . . . . .	8
3.1. 基准模型 . . . . .	8
3.1.1. 模型介绍 . . . . .	8
3.1.2. 实验结果与分析 . . . . .	10
3.2. 模型优化 . . . . .	11
3.2.1. 架构优化 . . . . .	11
3.2.2. 实验结果与分析 . . . . .	13
<b>第四章 机制探究</b> . . . . .	17
4.1. 理论分析 . . . . .	17
4.2. 实验设计与数据集生成 . . . . .	17
4.3. 实验结果 . . . . .	17
<b>参考文献</b> . . . . .	20
<b>附录 A L<sup>A</sup>T<sub>E</sub>X 实验</b> . . . . .	23
<b>附录 B MATLAB 实验</b> . . . . .	24
<b>致 谢</b> . . . . .	25

# 第一章 绪论

## 1.1 课题的背景和意义

基于单视角二维图片输入进行三维物体形状重建是三维视觉领域的核心问题之一，它能有效解决现实世界中二维数据丰富而三维数据稀少的问题，满足无人驾驶，智能建造，机器人领域的需求。点云数据具有易处理，易存储，易获得的特点，因此成为重建任务中常用的三维数据表示之一。

使用深度学习来进行单视角三维重建持续受到人们的关注。尽管很多的文章已经展示了创新的深度学习框架来提高三维重建任务中的最高水平<sup>[1-11]</sup>，但很少有文章来尝试探究这些任务的本质特性。不可否认的是，重建三维非常规形状的问题已经变成了一个新的机器学习范例，并且从事者使用的处理数据和训练网络的方法与在常规数据上的学习使用的方法不同（比如，Adam<sup>[12]</sup> 的使用频率比 SGD 要高）。因此，神经网络在三维重建学习这个新范例上怎样进行学习？是否与传统的向量分类和回归问题不同？

最近，<sup>[13]</sup> 的作者针对上述问题提出了一个让人惊讶的观点。他们尝试性的展示了当前最先进的用于三维重建任务的深度神经网络更倾向于通过首先分类输入图片到一个特定的簇，然后生成对应簇的平均三维形状，以此来进行预测。支持这一观点的主要实验证据是这些深度神经网络的三维预测结果与纯粹基于聚类和基于检索的基准模型的三维预测结果效果相近。这是一个非常有趣的观察，因为它说明了对于三维重建的任务来说，深度神经网络倾向于记住平均形状并且将其与图片输入的语义联系起来，而不是使用几何方法来生成一个形状，比如通过融合细粒度的局部结构形成一个整体形状。如果这个观点是正确的，这说明对于三维重建任务来说，最先进的深度神经网络实际执行了一个记忆任务而不是一个泛化任务<sup>[14]</sup>。

本课题旨在优化三维重建任务中深度神经网络的架构，证实其有效性。同时进一步探究深度神经网络在三维重建任务中的内在机制，分析导致其偏倚于识别或重建机制的因素，这一研究将会帮助当前应用于三维重建的深度学习方法避免陷入识别机制的误区，同时也探讨了深度学习的本质问题，即如何帮助神经网络进行更好的泛化。

## 1.2 研究现状

### 1.2.1 相关模型

<sup>[3]</sup> 首次提出了以单张图片作为输入的三维点云重建任务中深度神经网络的基本架构，即自编码器架构。编码器由卷积层和 ReLU 层构成，编码器的输入是一张图片和一个向量，向量用来模拟重建任务中的不确定性。解码器由全连接层构成，解码器输出点云的坐标，用  $N \times 3$  的矩阵表示， $N$  为一个点云中三维点的数量。<sup>[6]</sup> 提出了以三维点云作为输入的三维点云重建任务中深度自编码器的架构，创造性的提出用折叠的思想来生成点云，并重建三维形状，根据该思想实现的解码器仅使用了全连接层解码器的 7% 的参数量，却在重建效果上超过当

时的基准模型。<sup>[13]</sup> 中评估了多个效果拔尖的单视角三维重建深度神经网路模型，并提出了基于识别和检索机制的多个非深度学习方法模型，且这些模型的重建性能超过了当前的深度学习网络模型，这使得深度学习在单视角三维重建中的有效性受到了置疑。

### 1.2.2 相关数据集

ShapeNet<sup>[15]</sup> 是一个注释丰富且规模较大的三维形状数据集，涵盖 55 个常见的类别，有大约 5 万个样本，每个样本内有一个三维模型和多张从不同视角渲染的该三维模型的图片，三维模型的数据格式为网格、体素，图片格式为 PNG，在进行单视角三维重建时会从多个视角图片中选择一个视角的图片，因此这种情况下每个样本内有一张图片与一个三维模型。该数据集的创立者在 ICCV 2017 举办了基于该数据集的单视角三维重建任务的竞赛，并将当时的基准成绩发表在<sup>[16]</sup>。<sup>[13]</sup> 提供的数据集在 ShapeNet<sup>[15]</sup> 的基础上增加了点云数据格式，每个三维形状由 9000 多个点构成。共有 52430 个样本，涵盖 55 个类。

### 1.2.3 神经网络学习机制

深度神经网络是执行记忆还是泛化一直是现代机器学习中的主要问题。与我们将要介绍的类似，众所周知的猜测是优化过程是‘内容感知’的，并且取决于数据本身的属性<sup>[14]</sup>。<sup>[14]</sup> 中还显示，训练期间的某些正则化技术可帮助深度神经网络泛化而不是记住任务。对于三维形状重建，<sup>[13]</sup> 表明神经网络倾向于记忆平均形状，而不是在几何意义上进行重建。确实，许多作品还显示了平均形状和识别信息在提高三维重建效果中的有效性<sup>[17-19]</sup>。

相比之下，也有很多作品利用三维形状的连续潜在空间中的分布信息<sup>[1, 20-25]</sup> 来提高三维重建效果，这超出了基于识别的范围。值得注意的是，<sup>[26, 27]</sup> 表明形状算术可以在三维形状的潜在空间中进行，从而排除了神经网络仅在此问题设置中执行识别的可能性（因为执行算术需要的不是均值信息离散簇的形状）。其他一些工作建议通过将每个形状分解为部分<sup>[28-30]</sup> 或通过连续过程<sup>[31, 32]</sup> 来生成三维形状，这也超出了简单的识别任务。但是，<sup>[6, 26, 33]</sup> 的作者将基于三维重建的自动编码视为点云上无监督分类的基础。尽管这些作品中的输入数据是三维形状，而不是二维图像，但是形状信息有助于分类的事实似乎确实增强了这样的概念，即三维重建更着重于识别而不是重建。尽管如此，形状信息有助于识别的事实不能成为判断三维重建网络所学知识的主要理由。有意义的未来方向是研究分别用于无监督分类和受监督三维形状重建的神经网络学习之间的差异，因为这两个方向的主要目标并不完全相同。

注意三维重建问题可以被看作是一个更普适的分布学习<sup>[34, 35]</sup> 的特殊情况。但是，与分布学习中的理论工作如拓展核方法到回归分布不同，我们的工作集中于深度学习。即便如此，使人感兴趣的是能看到分布学习的传统工作把输出分布当做一个所有训练分布样本的连续线性组合直接处理，而不是使用两步法，预测簇索引后再预测平均分布。

## 1.3 本文研究内容

### 1.3.1 课题关键问题以及难点

首先，正式定义本课题的关键问题：一、优化深度神经网络框架，使其在单视角三维重建任务中超过基于识别机制的非深度学习方法，证明其有效性。二、讨论深度神经网络在单视角三维重建任务中执行的是识别机制还是重建机制。初步猜想是影响其在两者之间偏倚的因素为训练数据集的整体特征：聚类程度。因此本课题需要考虑以下几个难点：

1. 改进当前作为基准的神经网络架构以获得更好的重建性能。我们希望能对当前作为基准线的深度神经网络进行网络结构优化，以期望其在公认的标准数据集上能接近并超越基于识别机制的非深度学习方法。因此我们考虑借鉴图像领域成熟且有效的神经网络架构优化方法和训练技巧。
2. 定义神经网络的识别与重建这两种机制的数学表达，并用具体的实验结果来描述。
3. 设计并产生具有量化特征的三维重建数据集。为了探究神经网络训练集和网络的性能之间是否有强相关性，需要能定量的操控数据集的某些整体特征，如聚类系数。可以考虑的方式就是生成自定义的数据集，同时在生成过程中通过采样改变整体特征。为此可以考虑使用计算机图形学的相关软件来合成三维模型，并探究一些能进行不同三维形状之间插值的算法。
4. 定义衡量数据集指标的度量标准。在解决问题三后，需要构建一个度量标准来衡量数据集的聚类趋势，得到量化评分。

### 1.3.2 主要贡献

在这项工作中，我们定量和定性的展示了<sup>[13]</sup> 的结论不是通用的，而且这是一个不合适的数据集收集和不合理的数据使用导致的复杂结果。我们建议使用数据挖掘中确立已久的理论测量一个数据集的“聚类趋势”，称为 affinity propagation<sup>[36]</sup> 和 silhouette score<sup>[37]</sup>，并且我们展示了当用于训练三维重建任务的数据集不具有聚类特征时，训练结束的神经网络可以学会更集中在重建而不是识别。更重要的是，我们展示了即使是现实数据集如 ShapeNet<sup>[15]</sup>，在其上进行训练的神经网络依然呈现出重建机制，且通过优化网络架构和使用训练技巧，其重建效果接近并超过了基于检索的非深度学习方法。我们关于识别与重建这两种机制的理解更丰富，将三维重建的内在机制与数据以及训练过程相关联。之更具体的是，我们展示了训练集中的三维形状之间应该比训练集中的图片间呈现出更多的差异，以此来避免产生一个基于识别的机器学习模型。这一结论很实用，因为它为三维数据集的收集以及三维重建训练提供了指导。

## 1.4 论文组织结构

本论文主要有六个章节，各章的内容安排如下：

第一章 绪论，简要介绍了单视角三维重建任务的研究背景和意义，并对现行的相关理论，模型和数据集进行简单概述与分析。介绍本论文的研究思路、关键问题以及主要贡献。

第二章 问题与定义，介绍了本论文探究的损失函数和算法的理论基础。给出了重建机制与识别机制的数学定义。

第三章 模型，详细介绍了基准模型，以及优化模型的架构，并给出了网络在 ShapeNet 大型数据集上训练的过程信息以及训练结果的分析比较。

第四章 机制分析，给出了机制探究的理论基础与实验设计，介绍了自定义数据集的制作，网络在自定义数据集上训练的过程信息以及训练结果的定性与定量分析。

第五章 总结与展望，分析并总结本课题研究成果的总体优缺点，并提出未来的研究改进方向。

## 第二章 问题与定义

### 2.1 问题描述

我们探究的问题是从单张图片重建一个三维形状。输入  $I$  是一张二维 RGB 图片。输出  $S$  由一个点云表示。我们在本课题中不考虑基于体素的体积表达。对于基于点云的表达，每个形状  $S$  是一个包含三维点的点集。一个神经网络在训练中通过减小特定损失函数定义的经验损失  $l$ ，以此来从输入图片  $I$  预测形状  $S$ 。

$$\min_f \sum_{i=0}^{n-1} l(f(I_i), S_i). \quad (2.1)$$

我们想通过优化  $f$ ，使其在训练数据集上训练后，在测试数据集上取得更低的平均损失。同时，我们还想学习  $f$  的特性来了解它执行的是一个重建任务还是识别任务。两个指标被用来测量重建结果  $\hat{S} = f(I)$  与标签点云  $S$  之间的差异，称为 Chamfer Distance 和 F-score。

**Chamfer distance.** Chamfer Distance 通过搜索另一个点集中最近点来测试一个点集到另一个点集的整体距离。

$$d_{CH}(S, \hat{S}) = \frac{1}{|S|} \sum_{\mathbf{x} \in S} \min_{\hat{\mathbf{x}} \in \hat{S}} \|\mathbf{x} - \hat{\mathbf{x}}\|_2 + \frac{1}{|\hat{S}|} \sum_{\hat{\mathbf{x}} \in \hat{S}} \min_{\mathbf{x} \in S} \|\hat{\mathbf{x}} - \mathbf{x}\|_2. \quad (2.2)$$

尽管 Chamfer Distance 是运算高效且方便的，它会被很小一部分异常点的强烈影响。所以我们要如 [26] 建议的那样采用 F-score 来预测点集。

**F-score.** 另外一个衡量形状重建的指标是 F-score，该指标是精确度和回忆度的调和平均数。在给定相应的标签  $S$ ，且在固定的距离阈值  $d$  内，重建点云  $\hat{S}$  的精确度  $Prec$  为定义为：

$$Prec(d, S, \hat{S}) = \frac{1}{|\hat{S}|} \sum_{r \in \hat{S}} [\min_{s \in S} \|r - s\| < d], \quad (2.3)$$

where  $[\cdot]$  is the Iverson bracket.

相似的，标签点云  $S$  对重建点云  $\hat{S}$  的回忆度  $Rec$  是：

$$Rec(d, S, \hat{S}) = \frac{1}{|S|} \sum_{s \in S} [\min_{r \in \hat{S}} \|s - r\| < d]. \quad (2.4)$$

使用这两个量，如下计算 F-score：

$$F(d, S, \hat{S}) = \frac{2 \times Prec(d, S, \hat{S}) \times Rec(d, S, \hat{S})}{Prec(d, S, \hat{S}) + Rec(d, S, \hat{S})}. \quad (2.5)$$

重建的准确度被精确度量化，用以测量重建的点集离标签点集有多近。重建的整体度被回忆度量化，用来测量标签点云多少部分被重建点云覆盖。所以，一个高的 F-score 显示了重建是准确且完整的 [12]。

## 2.2 识别或重建

在这个部分，我们正式定义本文中研究的两种学习范式，被称为识别 (*recognition*) 和重建 (*reconstruction*)

**定义 1 (识别)** 一个基于识别的神经网络  $f$  用两步预测形状重建。神经网络重建方程可以写成：

$$\hat{S} = f(I) = f_1(f_2(I)), \quad (2.6)$$

方程中的  $f_2(\cdot)$  将输入图片映射到一个标量索引，并且  $f_1(f_2(I))$  将这个标量索引映射到索引  $f_2(I)$  对应的特定簇的平均形状。

**定义 2 (重建)** 一个基于重建的神经网络直接进行三维重建。即

$$\hat{S} = f(I), \quad (2.7)$$

并且重建不会明显的使用图片簇的任何信息。因为当前最流行的单视角三维重建经常使用一个编码器-解码器结构，所以一个相似的概念是由编码器获得的码字不会形成簇。

### 主要问题

在这篇论文中，我们研究神经网络执行识别机制还是重建机制。我们展示了向两者任一偏倚的趋势由数据集特性决定。

注意<sup>[3]</sup> 的主要结论是单视角三维重建中的深度神经网络主要执行识别工作。换句话说，神经网络的方程更接近于定义1而不是定义2。

## 2.3 数据集的聚类趋势

在这个部分，我们定义用来测量数据集聚类趋势的指标。具体来说，我们使用 silhouette score<sup>[37]</sup> 来测量聚类趋势。给定一个数据集  $D = \{x_i\}_{i=0}^{N-1}$  并且一个随机距离  $d(x, y)$  方程。我们考虑满足下面三个属性的任何距离方程， $d(x_1, x_2) = d(x_2, x_1)$ ,  $d(x_1, x_2) \geq 0$ , 和  $d(x, x) = 0$ 。在我们的实验中，我们使用 Chamfer Distance 作为点云的距离指标。我们使用  $\ell$ -1 距离作为图片之间的距离。我们能通过明确一个聚类特性方程  $C(\cdot)$  来确定一个数据集的聚类程度。对于每个样本  $x_i$ ，聚类方程会给出一个聚类标签  $C(x_i)$ 。我们使用  $C(x_i)$  来指示包含有  $x_i$  的簇，比如说，聚类标签  $C_i$  等价于  $C(x_i)$ 。有时候，数据集已经包含了标签聚类。更多情况下，聚类划分需要通过一个算法来获得。然后，第  $i$  个样本的轮廓系数 (silhouette score) 被定义为：

$$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}, \quad (2.8)$$

其中

$$a(i) = \frac{1}{|C_i| - 1} \sum_{j \in C_i, j \neq i} d(i, j), \quad b(i) = \min_{k \neq i} \frac{1}{|C_k|} \sum_{j \in C_k} d(i, j). \quad (2.9)$$

聚类趋势, 以轮廓系数定义, 由以下方程得出:

$$\text{Clustering-tendency} = \frac{1}{N} \sum_i s(i). \quad (2.10)$$

对于没有标签聚类的数据集, 我们需要找到聚类方程  $C(\cdot)$ 。在这篇论文中, 我们使用 affinity propagation<sup>[36]</sup> 来进行聚类。Affinity propagation 非常适合我们的设定, 因为它不需要预先确定簇的数量, 并且它使用一个预设的距离矩阵来分配数据点  $x_i$  到一个簇  $C_i$ 。我们使用 (2.2) 中的 Chamfer Distance 作为距离方程  $d(x, y)$  来处理数据集中的无序点云。

## 第三章 模型方法

在这个部分，我们介绍了本文研究的三维重建任务中深度神经网络基准模型和非深度学习基准模型，并结合在 ShapeNet 上的实验结果对其优缺点进行分析。基于这些分析，对深度神经网络基准模型进行优化，使其超过基于检索机制的非深度学习基准模型，证实了深度学习的优越性，我们使用了这两种方法：

1. 基于多任务训练的编码器初始化
2. 引入残差层增加解码器深度

我们还展示了在真实数据集 ShapeNet 上，即使是一个标准的三维重建深度神经网络也倾向于执行重建任务而不是识别任务。我们的结论基于两点观察：

1. 首先，由标准的三维重建网络获得的量化结果如 Chamfer Distance 和 F-score 与基于 Oracle-Nearest-Neighbor(Oracle-NN) 的重建方法所获得的结果相近。
2. 其次，由自编码训练获得的高维码字的 2D T-SNE 可视化，和使用标签信息进行训练相比，并没有展示出明显的形成簇的趋势。而且，没有标签的重建效果比有标签的重建效果要更好。

在这个章节，我们使用的数据集均为 ShapeNet，共计 55 个类，52430 个样本，训练集/验证集/测试集按照 70%/10%/20% 随机采样划分。图片大小为  $224 \times 224$ ，每个三维点云有 1024 个三维点构成。

### 3.1 基准模型

#### 3.1.1 模型介绍

##### 3.1.1.1 非深度学习方法

论文 [13] 中提出了以下三个基于识别机制的非深度学习模型：

- *Clustering*: 通过对样本点云集使用聚类算法，将一组点云集划分为多个集群，每个集群内部计算一个平均三维形状，接着训练一个基于输入图片预测特定集群的分类器，将集群的平均三维形状作为预测结果。
- *Retrieval*: 借鉴了现有的基于图片检索对应物体的三维形状的方法 [38]，在测试时，根据输入图片检索训练集中对应的三维点云，直接提取出来作为预测结果。
- *Oracle-Nearest-Neighbor*<sup>1</sup>: 该模型在预测时，在训练集中直接搜索与标签点云损失最小的训练点云，将其作为预测结果。因为该模型在搜索最近点云时需要标签点云，而实际测试时只有对应的二维图片作为输入，所以在实践中是不可能实现的。

<sup>1</sup>下文缩写为 Oracle-NN

因为论文<sup>[13]</sup>的实验结果表明 Oracle NN 的水平超过了所有当前最先进的深度学习模型，而且它是表征任何实际基于检索的非深度学习方法的性能极限的理论基准。所以我们将其当做本课题的非深度学习基准模型，并进行这样的实现<sup>1</sup>。我们发现在将该模型应用到 ShapeNet 时，因为测试集有 10000 个样本，训练集有 35000 个样本，如果这样计算，复杂度很高，需要在 GPU 上消耗大概一周时间，占用了大量的计算资源。于是考虑利用数据集自带的类标签，共计 55 个类，如车，床，椅子等，在每个类内数据集中运行该算法<sup>1</sup>，大约 10 小时完成。

#### 算法 1 Oracle Nearest Neighbor

输入:  $T_{i=0}^n$  测试集,  $D_{j=0}^t$  训练集

输出:  $R_{k=0}^n$  预测结果

```

1:  $i \leftarrow 0$ 
2:  $k \leftarrow 0$ 
3: while  $i < n$  do
4:    $index \leftarrow 0$ 
5:    $mindis \leftarrow inf$ 
6:    $j \leftarrow 0$ 
7:   while  $j < t$  do
8:      $dis \leftarrow Distance(T_i, D_j)$ 
9:     if  $dis < mindis$  then
10:       $index \leftarrow j$ 
11:       $mindis \leftarrow dis$ 
12:    end if
13:     $j \leftarrow j + 1$ 
14:  end while
15:   $R_{k=i} \leftarrow D_{j=index}$ 
16:   $i \leftarrow i + 1$ 
17: end while
```

#### 3.1.1.2 深度学习方法

先画流程图，可以的话画一下比较炫酷的图 根据论文 PSGN<sup>[3]</sup> 提出的基于“编码器-解码

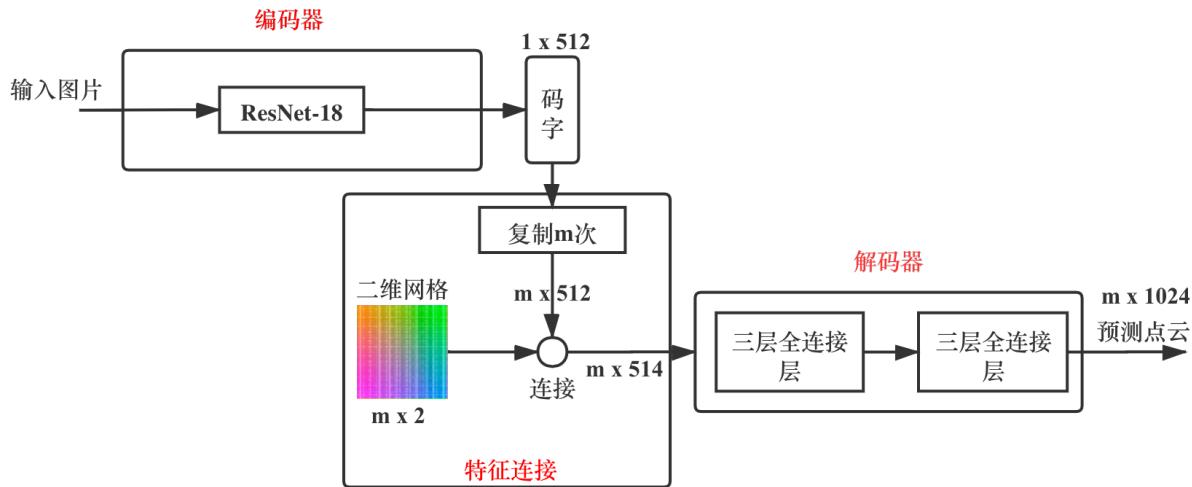


图 3-1 基于自编码器架构的深度神经网络基准模型

器”架构<sup>1</sup>的深度神经网络模型，我们进行了如图3-1的实现。网络的主要架构如下：

- 编码器：采用 ResNet-18<sup>[39]</sup> 的卷积层架构作为接受二维图像输入并提取图像特征的编码器，具体来说，移除 ResNet-18 最后一层用来分类的全连接层，将一维特征层直接输出作为码字。
- 特征连接：模型中存储有一个预先初始化好的二维正方形网格点云，长宽在 [-1,1] 之间，由  $m \times 2$  的矩阵表示。在将码字输入解码器之前，我们先将其复制  $m$  次，然后将  $m \times 512$  矩阵（码字）与  $m \times 2$  矩阵（二维网格）连接起来，连接的结果是大小为  $m \times 514$  的矩阵。
- 解码器：采用 FoldingNet 解码器<sup>[6]</sup> 作为从码字恢复到点云结构的解码器，具体来说，这个解码器由两个感知器构成，每个感知器由三层全连接层构成。

该基准模型的关键部分是提出的解码器使用两个连续的 3 层感知器将固定的 2D 网格扭曲为输入点云的形状。这个过程称为折叠操作，实质上形成了通用的二维到三维映射。为了直观地解释为什么这种折叠操作是通用的二维到三维的映射，请用矩阵  $U$  表示输入的二维网格点。 $U$  的每一行都是一个二维网格点。用  $U_i$  表示  $U$  的第  $i$  行，用  $\theta$  表示编码器输出的码字。然后，在特征连接之后，解码器的输入矩阵的第  $i$  行为  $[u_i, \theta]$ ，由于感知器并行应用于输入矩阵的每一行，因此输出矩阵的第  $i$  行可以写成  $f([u_i, \theta])$ ，其中  $f$  表示感知器进行的功能。可以将该函数视为带有码字  $\theta$  的参数化高维函数，码字  $\theta$  是指导功能结构（折叠操作）的参数。由于多层感知器擅长拟合非线性函数，因此它们可以在二维网格上执行精细的折叠操作。

### 3.1.2 实验结果与分析

在这个部分，我们展示了上述的深度学习基准模型和非深度学习基准模型在 ShapeNet 上的实验效果，并指出深度学习基准模型的问题。

对于非深度学习基准模型 Oracle-NN，我们展示了所有重建样本的损失平均值作为量化的评估结果。对于基于自编码器架构的深度学习基准模型，我们以 Chamfer Distance<sup>2</sup> 作为损失函数，使用 Adam 优化器，训练了 100 个周期，使其达到收敛，并对超参数进行调参，总共进行了 8 次训练，最终展示了由所有重建样本的损失平均值表示的最优重建效果。两个模型的最优结果展示在表3.1。从这个结果中我们可以看出，在最优平均损失方面，Oracle NN 超过了深度神经网络，这一结果与论文<sup>[13]</sup> 中展示的两者之间的差距相吻合，而且正如前文所展示的，Oracle-NN 是一个理论上的基于识别机制的基准模型，无法应用于实际任务，所以这一优越结果是被预计的。但是，在训练的过程中，我们发现了神经网络基准模型的一些问题，这意味着有可能通过一些优化或者训练技巧来继续缩小这一差距。结合图3-2展示的在不同学习率下的两次训练过程的损失变化，对主要问题的分析如下：

- 神经网络训练状态不稳定。首先可以看出，学习率仅发生了三倍的变化，就使得网络陷入了与最优效果相差较大的局部最优，且 Adam 优化器也没有能使得网络跳出这个点。在正

<sup>1</sup>下文简称为“自编码器”或“auto-encoder”

<sup>2</sup>Chamfer Distance 数值越小说明重建效果越好

常情况下，网络不应该对超参数有如此高的敏感程度。同时，这也说明重建任务也许对参数初始化有比较高的要求。

- 神经网络架构的学习能力较弱。如图3-2b所示，两个网络分别在第30周期和第60周期收敛，收敛过早说明当前的网络可能出现了训练过程中梯度过小的问题。
- 神经网络的编码器缺少显性的梯度优化。从理论角度来说，网络损失函数仅描述两个三维点云的相似程度，因此有理由怀疑该函数传播的梯度并不能给予对二维图片进行处理的编码器足够的优化指导，使其提取需要的图片特征。从实验角度来说，我们将编码器输出的码字进行了T-SNE降维可视化，如图3-8b，可以看出编码器提取的特征码字并不具有区分图像类别的特征，虽然这不能证明编码器完全没有提取有助于三维重建的图像信息。

基于上述对神经网络基准模型的问题分析，我们将在下文对网络初始化和梯度过小等问题提出优化方案。

表 3.1 基准模型的评估结果

模型	Chamfer Distance
Oracle NN	0.0719
自编码器（基准模型）	0.0806

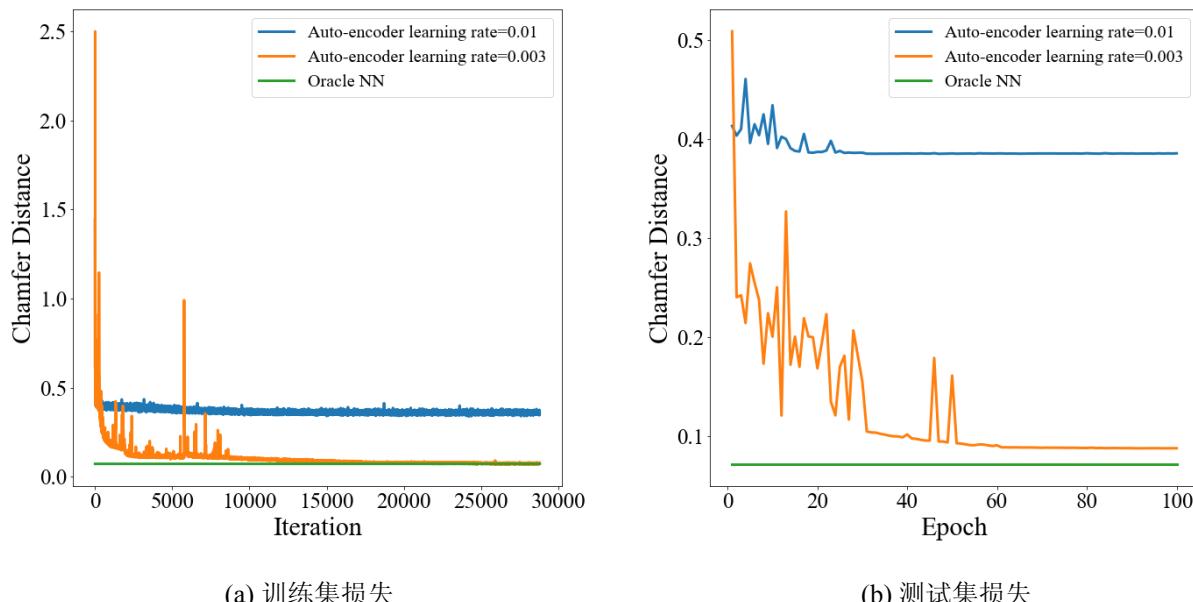


图 3-2 深度神经网络基准模型的两次训练过程。蓝色与橙色指同一网络在不同优化器学习率下的训练，绿色指非深度学习基准模型的量化结果

## 3.2 模型优化

### 3.2.1 架构优化

在这个部分，我们介绍用残差层改进译码器以及加入弱标签信息进行编码器初始化这两种方法的实现细节以及使用这两种方法的原因。

- 针对如何给网络参数提供良好的初始化，以及给编码器显性优化指导的问题，我们提出加入弱标签信息进行多任务训练以达到初始化的方法。在训练的前 5 个周期，提供类标签给网络，在预测重建点云的同时将码字从中间层抽出进行图像分类训练，如图3-3。在第 6 个周期以及之后取消这图片分类训练，只进行三维点云重建训练。这一创新方法的启发来自于图像深度学习领域常用的技巧，即使用预训练的网络参数进行网络初始化，比如在某个工程问题上需要进行图像识别，常见操作是把神经网络在 ImageNet 上训练到收敛后，再把网络在实际数据集上训练进行微调，我们没有采用这种预训练参数进行初始化的方法，一是因为 ShapeNet 的图片是由三维网格渲染而成的合成图片，ImageNet 有一部分来自现实世界图片，这两个数据域的分布不同，因此需要考虑域适应的问题。二是因为在 ShapeNet 上预先训练图像分类网络，增加三维重建的步骤和所需的计算资源，因此我们考虑使用弱标签信息来简化预训练初始化的流程。

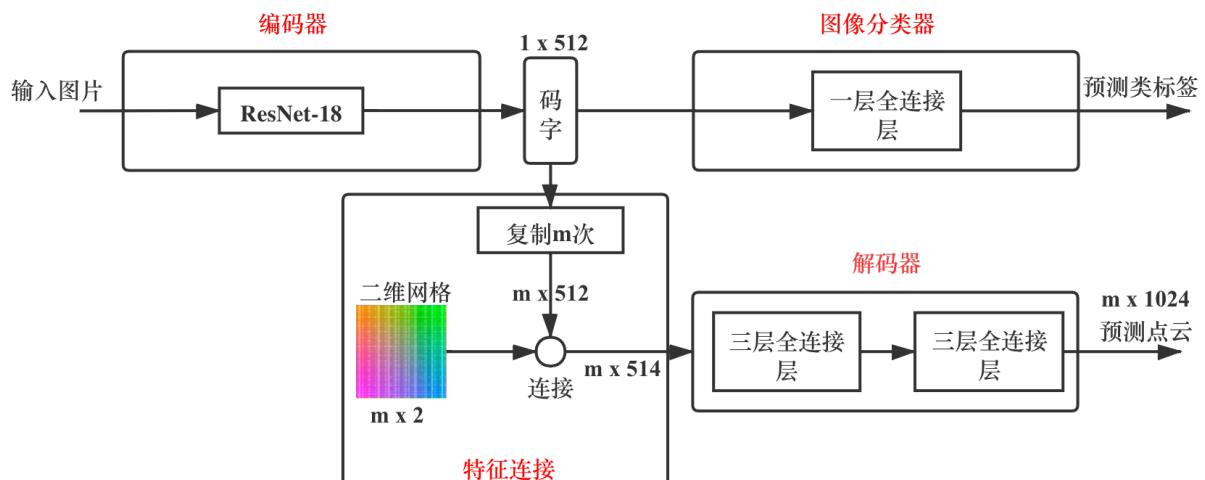


图 3-3 引入弱标签信息进行初始化

- 针对深度神经网络梯度过小和消失的问题，我们用残差层改进基准模型3-1的译码器，将三层全连接层感知器改为由三个残差块组成的感知器，如图3-4。且残差块内部全连接层之间进行批处理归一化，改进后译码器的参数量是基准模型译码器的两倍，保留译码器之前编码器和特征连接结构不变。根据图像分类领域深度学习的经验，残差层和批处理归一化都能有效的解决反向传播中梯度过小和消失的问题，且可以在此基础上增加网络深度，提高重建能力。

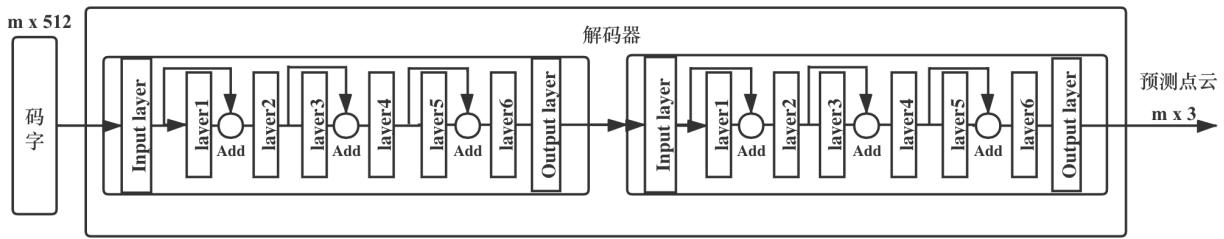


图 3-4 用残差层改进译码器

### 3.2.2 实验结果与分析

在这个部分,我们展示了将上述两种方法分别应用于优化网络,会在 Chamfer Distance2.1 和 F-score2.1 两个指标上取得接近甚至超过 Oracle-NN 的重建效果。同时, 我们进一步对深度神经网络的内在机制进行探究, 展示了一个标准的三维重建深度神经网络在公开数据集 ShapeNet 上倾向于执行重建任务而不是识别任务。

#### 3.2.2.1 重建效果

表3.2汇报了基准模型, 两种优化模型以及 Oracle NN 的最优平均重建损失 (Chamfer Distance), 数值越小说明重建效果越好, 我们可以看到加入弱标签信息进行初始化后, 基于自编码器的深度神经网络的重建水平 (0.0716) 优于 Oracle NN (0.0719), 使用残差层译码器后, 神经网络的水平 (0.0768) 也得到了显著提高。特别的, 我们将两个优化模型和 Oracle NN 的 Chamfer Distance 表示在图3-5a, 从该图可以直观的看出, 标准的基于自编码器的深度神经网络模型与 Oracle-NN 呈现出近乎相同的重建效果, 且 Oracle NN (基于识别机制的非学习模型的理论极限) 在测试时需要不应获得的标签信息。

除了 Chamfer Distance 这个指标外, 我们还使用了 F-score 这个指标, 将测试结果展现在3-5b, 这个指标下的模型比较结果与 Chamfer Distance 的比较一致。因此, 我们可以确定的得出结论: 弱标签初始化与残差层改进都成功提高了深度神经网络在 ShapeNet 整体数据集的平均重建效果, 前者上超过了基于识别机制的非学习模型的理论极限, 后者接近了这一理论极限。

表 3.2 优化模型的评估结果

模型	Chamfer Distance
自编码器 (基准模型)	0.0806
自编码器 (弱标签信息)	0.0716
自编码器 (残差层译码器)	0.0768
Oracle NN	0.0719

在平均重建效果之外, 我们还考虑了每个类的重建效果的统计特性, 如图3-6。我们可以看到基于自编码器的三维重建神经网络在不同的类的重建效果一致, 波动较小。而 Oracle NN 在一些类如 tower, 呈现出非常差的重建效果, 出现了各个类重建效果不一致的情况。这个问

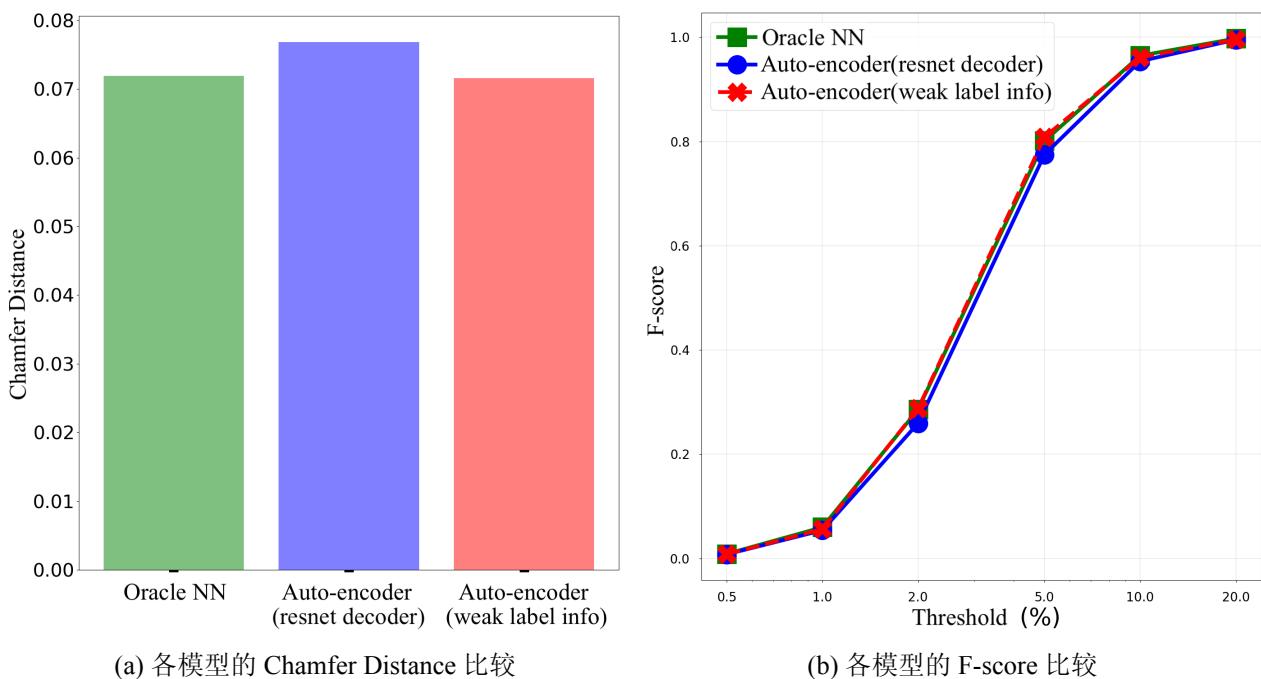


图 3-5 优化深度学习模型与非深度学习模型各指标比较

题同样也在模型预测结果可视化中呈现出来，如图 3-7第四行，我们可以看到 Oracle NN 预测出了在三维形状层面上非常不合理的点云。结合 ShapeNet 数据集的分布 TODO 附录，我们可以看到有一些类的样本量很少，Oracle NN 在样本量少的类内进行搜索时，有很大概率无法找到与预测点云相近的训练点云，因此会在这些类出现比较差的重建效果。

综上所述，我们可以看出基于识别机制的非深度学习模型的局限性，此类模型缺少对训练数据集分布的鲁棒性，当训练集无法在低维度层面解释测试集需要的形状特征时（即当训练集中没有和测试样本在损失距离上相近的点云时），该类模型的方法效果远差于深度学习模型。

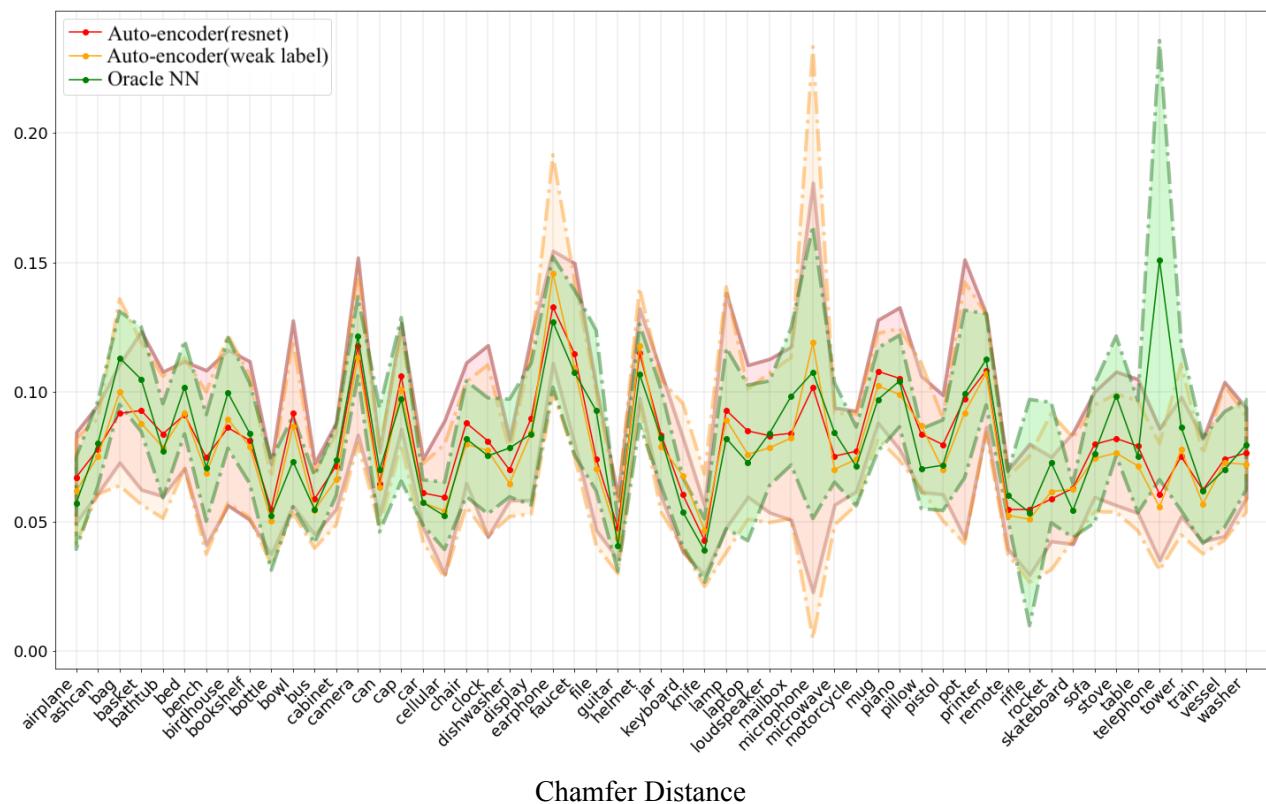


图 3-6 每个类的重建评价指标，X 轴每一列代表一个类的 Chamfer Distance 统计。实心点代表类内所有样本指标平均值，上边沿虚线代表类内所有样本指标最大值，下边沿虚线代表类内所有样本指标最小值。

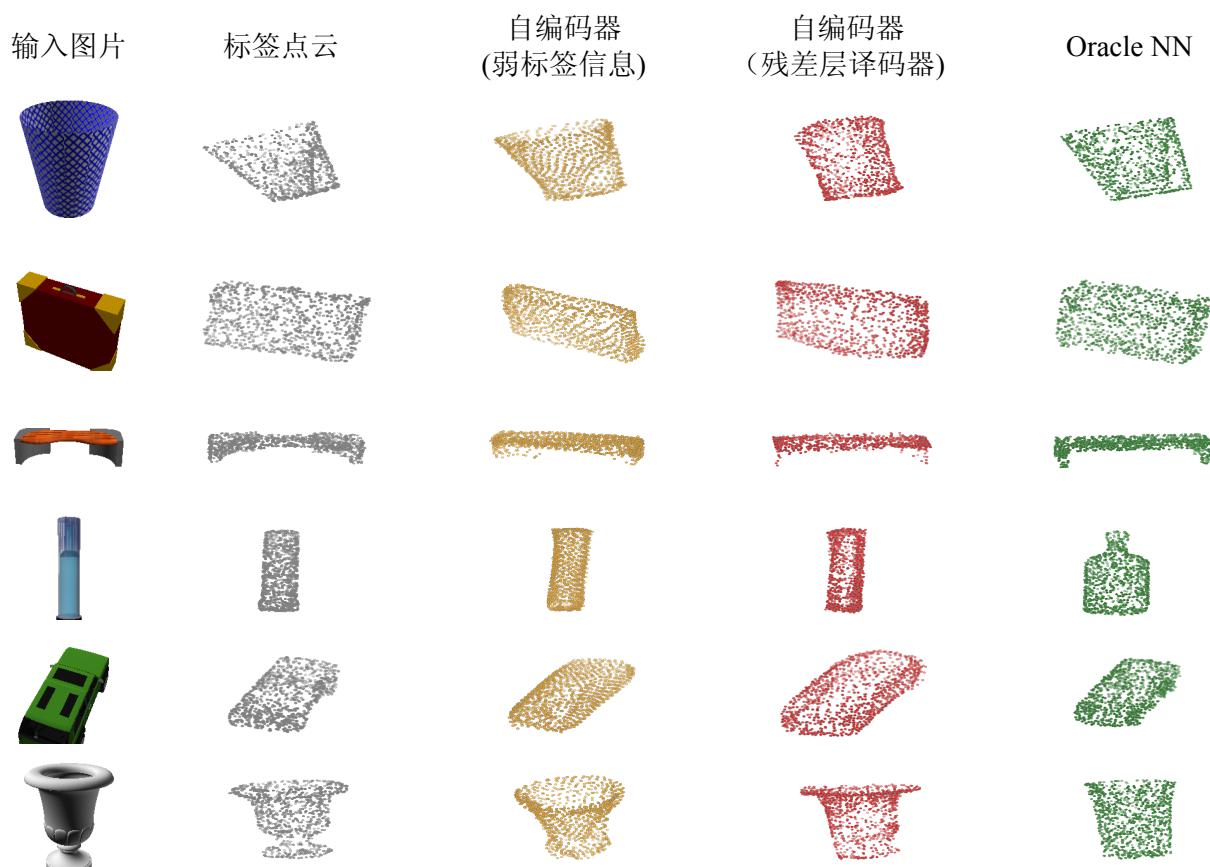


图 3-7 模型预测结果的可视化

## 3.2.2.2 机制探究

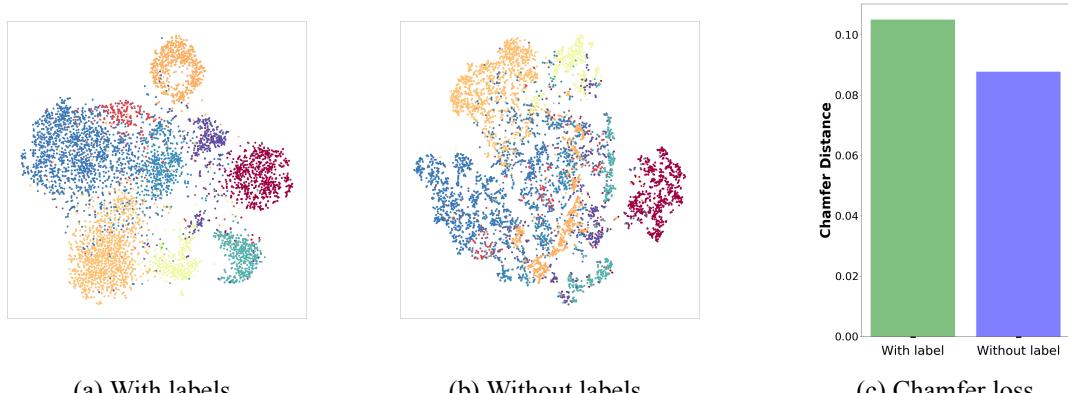


图 3-8 T-SNE and Chamfer loss of the same network trained with or without label information. The network trained without label information does not show a clear tendency to form clusters while yielding smaller Chamfer loss.

## 第四章 机制探究

### 4.1 理论分析

#### Clustering tendency relationship

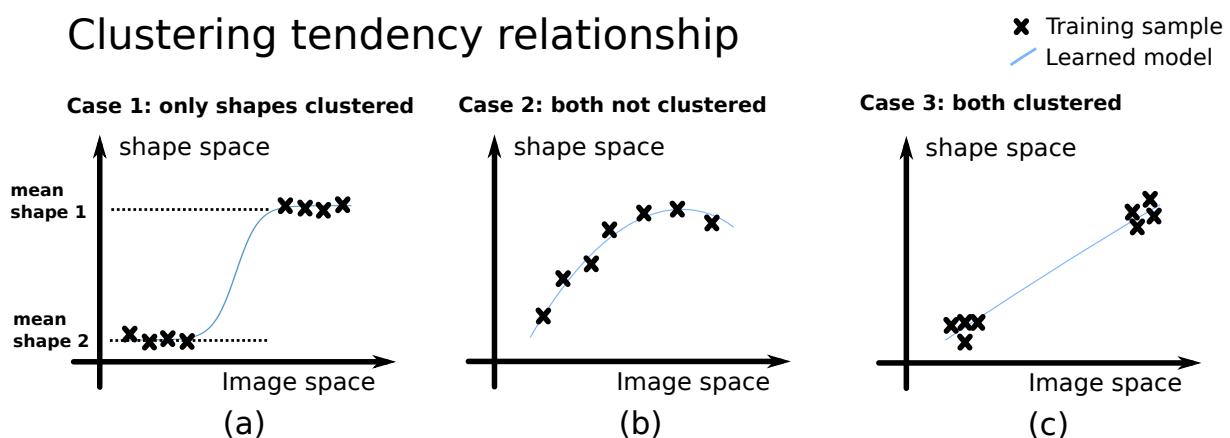


图 4-1 When the training images are less clustered than the training shapes (as shown in subfigure (a)), the learned model shows the tendency of recognition.

### 4.2 实验设计与数据集生成



图 4-2 Base dataset # 1: interpolation between a sphere and a cube.

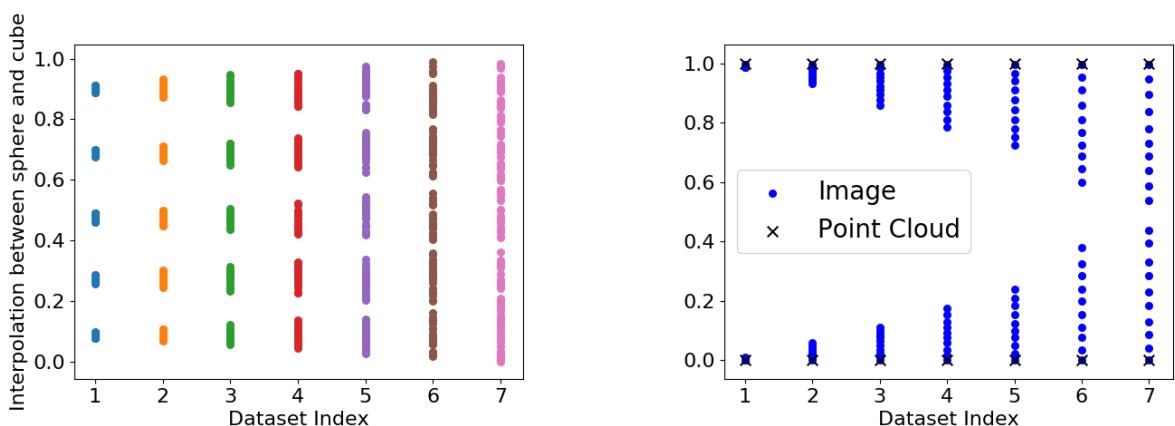


图 4-3 Coverage of two subsampled data sets 1.1.n and 1.2.n.

### 4.3 实验结果

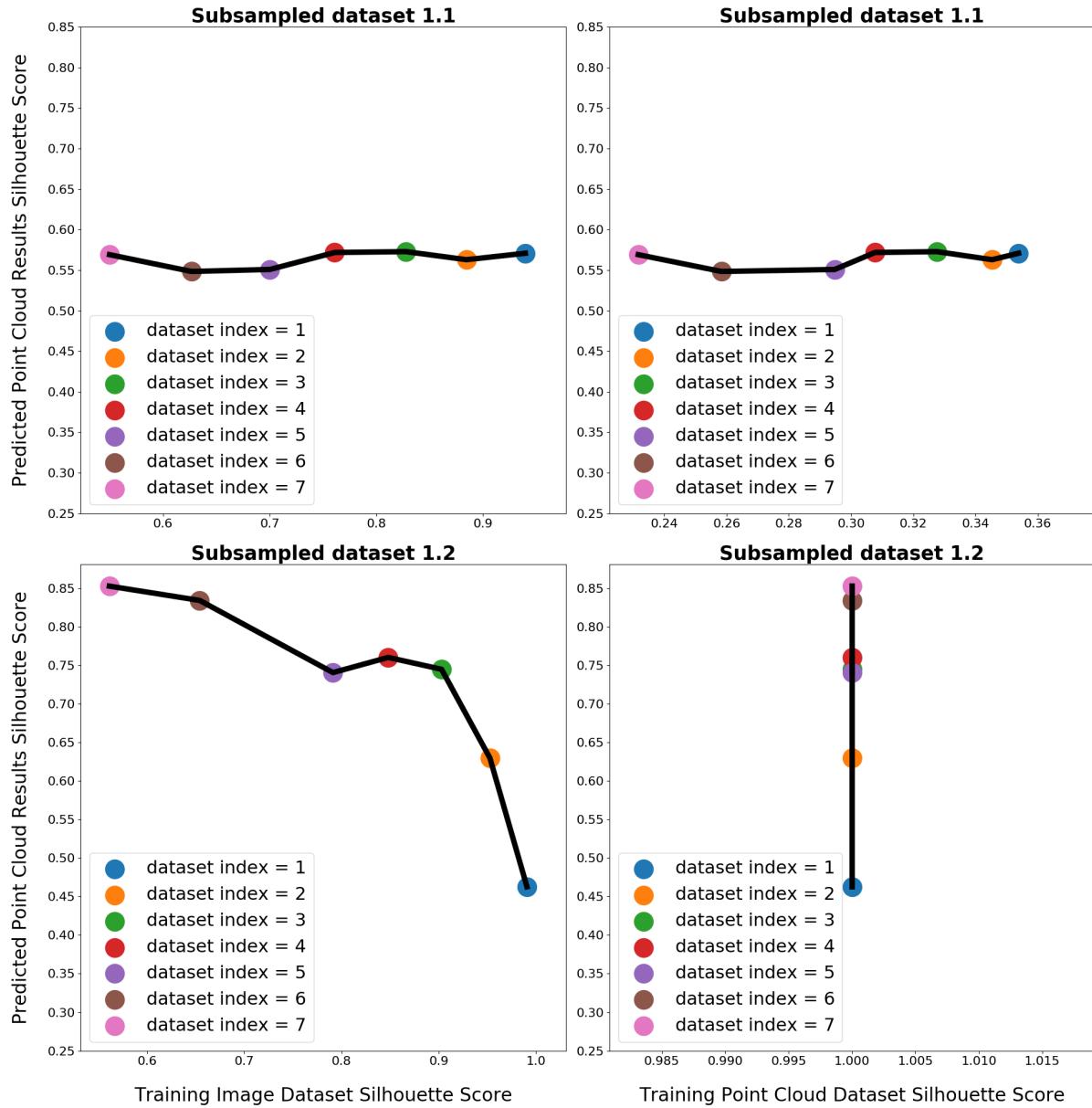


图 4-4 Training-vs-prediction silhouette scores. Only Case 1 (the points with large dataset indices shown in the bottom-left figure) shows high tendency towards recognition.

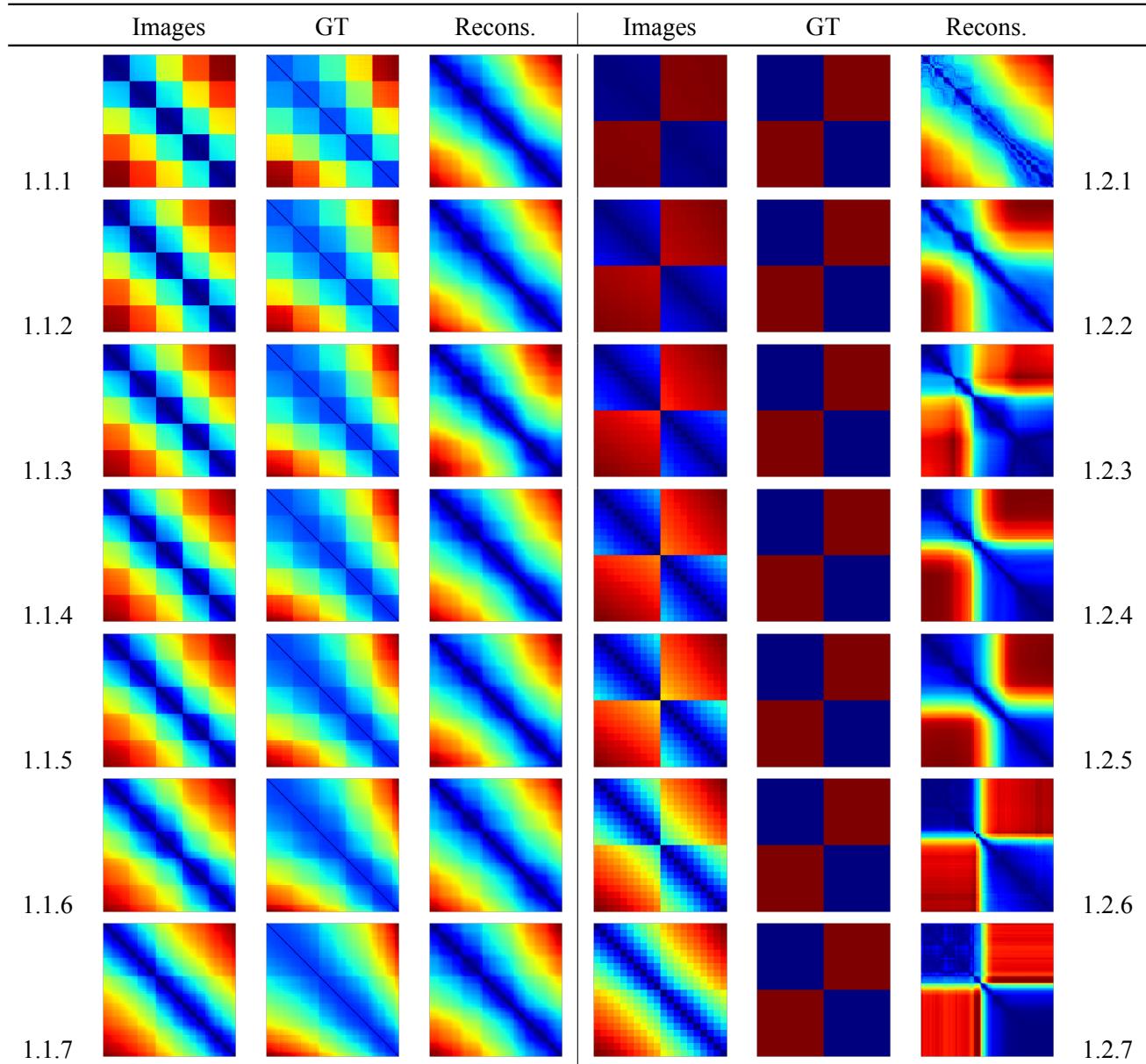


图 4-5 Distance matrices. Blue represents small distance and red represents large distance. (**Left**) predicted shapes of subsampled dataset 1.1 are not clustered. (**Right**) predicted shapes of subsampled dataset 1.2 are clustered if the training images are not clustered (i.e., if the dataset belongs to Case 1). Top to bottom: dataset index from 1 to 7.

## 参考文献

- [1] Li C L, Zaheer M, Zhang Y, et al. Point cloud gan[J]. arXiv preprint arXiv:1810.05795, 2018..
- [2] Park J J, Florence P, Straub J, et al. Deepsdf: Learning continuous signed distance functions for shape representation. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019. 165–174.
- [3] Fan H, Su H, Guibas L J. A point set generation network for 3d object reconstruction from a single image. Proceedings of the IEEE conference on computer vision and pattern recognition, 2017. 605–613.
- [4] Tatarchenko M, Dosovitskiy A, Brox T. Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. Proceedings of the IEEE International Conference on Computer Vision, 2017. 2088–2096.
- [5] Groueix T, Fisher M, Kim V G, et al. A papier-mâché approach to learning 3d surface generation. Proceedings of the IEEE conference on computer vision and pattern recognition, 2018. 216–224.
- [6] Yang Y, Feng C, Shen Y, et al. Foldingnet: Point cloud auto-encoder via deep grid deformation. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018. 206–215.
- [7] Wang N, Zhang Y, Li Z, et al. Pixel2mesh: Generating 3d mesh models from single rgb images. Proceedings of the European Conference on Computer Vision (ECCV), 2018. 52–67.
- [8] Sun X, Wu J, Zhang X, et al. Pix3d: Dataset and methods for single-image 3d shape modeling. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018. 2974–2983.
- [9] Tulsiani S, Zhou T, Efros A A, et al. Multi-view supervision for single-view reconstruction via differentiable ray consistency. Proceedings of the IEEE conference on computer vision and pattern recognition, 2017. 2626–2634.
- [10] Wu J, Wang Y, Xue T, et al. Marrnet: 3d shape reconstruction via 2.5 d sketches. Advances in neural information processing systems, 2017. 540–550.
- [11] Yan X, Yang J, Yumer E, et al. Perspective transformer nets: Learning single-view 3d object reconstruction without 3d supervision. Advances in neural information processing systems, 2016. 1696–1704.
- [12] Kingma D P, Ba J. Adam: A method for stochastic optimization[J]. arXiv preprint arXiv:1412.6980, 2014..
- [13] Tatarchenko M, Richter S R, Ranftl R, et al. What do single-view 3d reconstruction networks learn? Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019. 3405–3414.
- [14] Arpit D, Jastrz̄bski S, Ballas N, et al. A closer look at memorization in deep networks. Proceedings of the 34th International Conference on Machine Learning-Volume 70. JMLR. org, 2017. 233–242.

- [15] Chang A X, Funkhouser T, Guibas L, et al. Shapenet: An information-rich 3d model repository[J]. arXiv preprint arXiv:1512.03012, 2015..
- [16] Yi L, Shao L, Savva M, et al. Large-scale 3d shape reconstruction and segmentation from shapenet core55[J]. ArXiv, 2017, abs/1710.06104.
- [17] Kanazawa A, Tulsiani S, Efros A A, et al. Learning category-specific mesh reconstruction from image collections. Proceedings of the European Conference on Computer Vision (ECCV), 2018. 371–386.
- [18] Pontes J K, Kong C, Sridharan S, et al. Image2mesh: A learning framework for single image 3d reconstruction. Asian Conference on Computer Vision. Springer, 2018. 365–381.
- [19] Kurenkov A, Ji J, Garg A, et al. Deformnet: Free-form deformation network for 3d shape reconstruction from a single image. 2018 IEEE Winter Conference on Applications of Computer Vision (WACV). IEEE, 2018. 858–866.
- [20] Li K, Garg R, Cai M, et al. Single-view object shape reconstruction using deep shape prior and silhouette[J]. arXiv preprint arXiv:1811.11921, 2018..
- [21] Gwak J, Choy C B, Chandraker M, et al. Weakly supervised 3d reconstruction with adversarial constraint. 2017 International Conference on 3D Vision (3DV). IEEE, 2017. 263–272.
- [22] Sinha A, Unmesh A, Huang Q, et al. Surfnet: Generating 3d shape surfaces using deep residual networks. Proceedings of the IEEE conference on computer vision and pattern recognition, 2017. 6040–6049.
- [23] Wu J, Zhang C, Zhang X, et al. Learning shape priors for single-view 3d completion and reconstruction. Proceedings of the European Conference on Computer Vision (ECCV), 2018. 646–662.
- [24] Yang B, Wen H, Wang S, et al. 3d object reconstruction from a single depth view with adversarial learning. Proceedings of the IEEE International Conference on Computer Vision Workshops, 2017. 679–688.
- [25] Yang G, Huang X, Hao Z, et al. Pointflow: 3d point cloud generation with continuous normalizing flows. Proceedings of the IEEE International Conference on Computer Vision, 2019. 4541–4550.
- [26] Wu J, Zhang C, Xue T, et al. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. Advances in neural information processing systems, 2016. 82–90.
- [27] Girdhar R, Fouhey D F, Rodriguez M, et al. Learning a predictable and generative vector representation for objects. European Conference on Computer Vision. Springer, 2016. 484–499.
- [28] Tulsiani S, Su H, Guibas L J, et al. Learning shape abstractions by assembling volumetric primitives. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017. 2635–2643.
- [29] Niu C, Li J, Xu K. Im2struct: Recovering 3d shape structure from a single rgb image. Proceedings of the IEEE conference on computer vision and pattern recognition, 2018. 4521–4529.

- [30] Sung M, Su H, Kim V G, et al. Complementme: weakly-supervised component suggestions for 3d modeling[J]. ACM Transactions on Graphics (TOG), 2017, 36(6):1–12.
- [31] Sun Y, Wang Y, Liu Z, et al. Pointgrow: Autoregressively learned point cloud generation with self-attention[J]. arXiv preprint arXiv:1810.05591, 2018..
- [32] Choy C B, Xu D, Gwak J, et al. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. European conference on computer vision. Springer, 2016. 628–644.
- [33] Achlioptas P, Diamanti O, Mitliagkas I, et al. Learning representations and generative models for 3d point clouds[J]. arXiv preprint arXiv:1707.02392, 2017..
- [34] Oliva J, Póczos B, Schneider J. Distribution to distribution regression. International Conference on Machine Learning, 2013. 1049–1057.
- [35] Póczos B, Rinaldo A, Singh A, et al. Distribution-free distribution regression[J]. 2013..
- [36] Wang K, Zhang J, Li D, et al. Adaptive affinity propagation clustering[J]. arXiv preprint arXiv:0805.1096, 2008..
- [37] Van Craenendonck T, Blockeel H. Using internal validity measures to compare clustering algorithms[J]. Benelux 2015 Poster presentations (online), 2015. 1–8.
- [38] Li Y, Su H, Qi C R, et al. Joint embeddings of shapes and images via cnn image purification[J]. ACM Trans. Graph., 2015, 34(6).
- [39] He K, Zhang X, Ren S, et al. Deep residual learning for image recognition[J]. 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016..

## 附录 A L<sup>A</sup>T<sub>E</sub>X 实验

技术实验结果在这里写

## 附录 B MATLAB 实验

技术实验结果在这里写

## 致 谢

这次的毕业论文设计总结是在我的指导老师 xxx 老师亲切关怀和悉心指导下完成的。从毕业设计选题到设计完成，x 老师给予了我耐心指导与细心关怀，有了莫老师耐心指导与细心关怀我才不会在设计的过程中迷失方向，失去前进动力。x 老师有严肃的科学态度，严谨的治学精神和精益求精的工作作风，这些都是我所需要学习的，感谢 x 老师给予了我这样一个学习机会，谢谢！

感谢与我并肩作战的舍友与同学们，感谢关心我支持我的朋友们，感谢学校领导、老师们，感谢你们给予我的帮助与关怀；感谢肇庆学院，特别感谢计算机科学与软件学院四年来为我提供的良好学习环境，谢谢！