# Multi-View Representation is What You Need for Point-Cloud Pre-Training

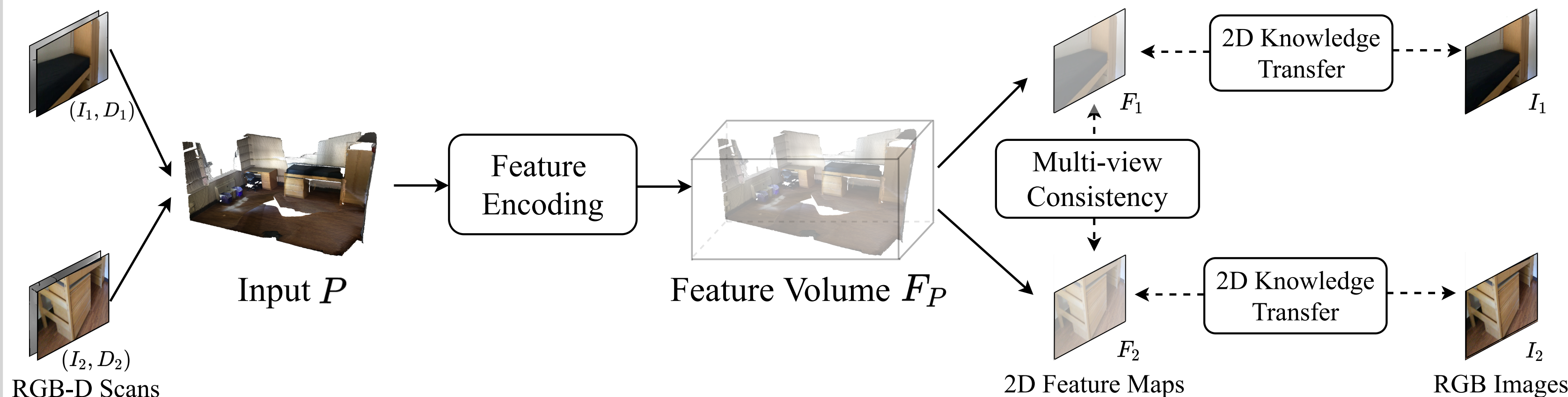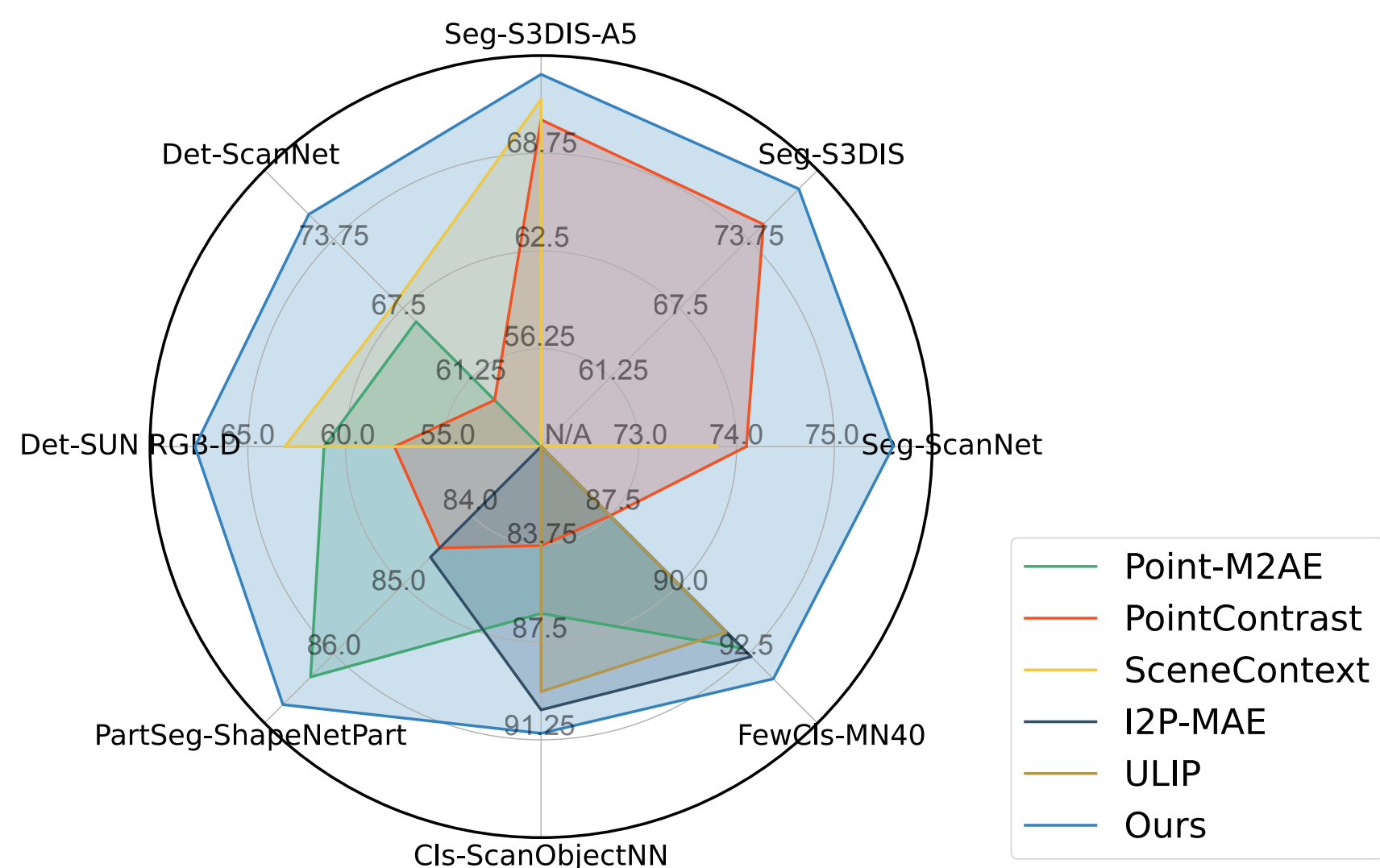Siming Yan, Chen Song, Youkang Kong, Qixing Huang

## Introduction

- Formulate point-cloud pre-training as learning a multi-view consistent 3D feature volume.
- Leverage pre-trained 2D image-based models to supervise 3D pre-training
- Develop an auxiliary pre-task where the goal is to predict the multi-view pixel-wise correspondences from the 2D pixel embeddings.
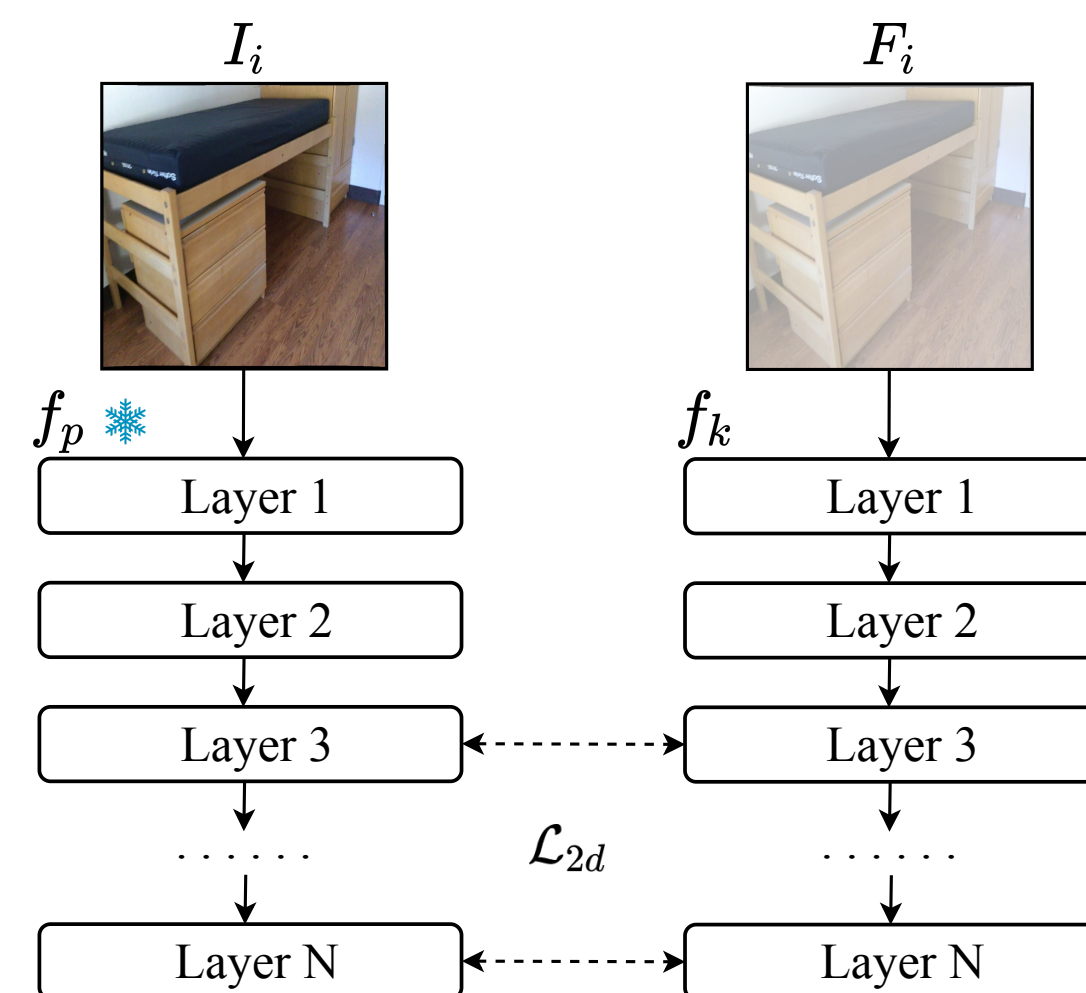- Conduct extensive experiments to demonstrate the effectiveness of our approach

## Method Overview



Input $P$ — RGB-D Scans $(I_1, D_1)$, $(I_2, D_2)$ → Feature Encoding → Feature Volume $F_P$ → Multi-view Consistency, 2D Knowledge Transfer, $F_1$, $F_2$, 2D Feature Maps, RGB Images $I_1$, $I_2$

## Experiment Results



Seg-S3DIS-A5, Seg-S3DIS, Seg-ScanNet, FewCls-MN40, Cls-ScanObjectNN, PartSeg-ShapeNetPart, Det-SUN RGB-D, Det-ScanNet

Point-M2AE, PointContrast, SceneContext, I2P-MAE, ULIP, Ours

## 2D Knowledge Transfer Module

- Frozen DINOv2 as 2D pre-trained model
- Knowledge distillation loss



$I_i$, $F_i$, $f_p$ ❄, $f_k$, Layer 1, Layer 2, Layer 3, Layer N, $\mathcal{L}_{2d}$

## Multi-view Consistency Module

- Predict feature correspondences from 2D embedding
- Cross-attention layer input: 1.concat 2-view feature maps; 2. query point position from first view



$F_1$, $F_2$, Attention-based Decoder, $x$, $x'$