# Holistic 3D Human and Scene Mesh Estimation from Single View Images

Zhenzhen Wang, Serena Yeung
Stanford University
{zzwang, syyeung}@stanford.edu

## Abstract

*The 3D world limits the human body pose and the human body pose conveys information about the surrounding objects. Indeed, from a single image of a person placed in an indoor scene, we as humans are adept at resolving ambiguities of the human pose and scene layout given our prior knowledge of the physical laws and prior perception of the plausible human-object-scene configurations. However, few computer vision models fully leverage this fact. In this work, we propose a holistically trainable model that perceives the 3D scene from a single RGB image, estimates the camera pose and the room layout, and reconstructs both human body and object meshes. By imposing a set of comprehensive and well sophisticated losses on all aspects of the estimations, we show that our model outperforms existing human body and object reconstruction methods, in the meanwhile produces more plausible scene reconstruction.*

## 1. Introduction

## 2. Related Work

## 3. Model

### 3.1. Representation

### 3.2. Model Architecture

### 3.3. Loss Functions and Optimization

## 4. Experiments

### 4.1. Datasets

### 4.2. Implementation Details

### 4.3. Quantitative Results

### 4.4. Ablation Analysis

## 5. Conclusion