

# Im2Pano3D:

### Extrapolating 360° Structure and Semantics Beyond the Field of View Shuran Song, Andy Zeng, Angel X. Chang, Manolis Savva, Silvio Savarese, and Thomas Funkhouser





### **Boundary Extension**

Intraub and Richardson, 1989







#### 20- 30 % More















"False memory 1/20th of a second later: what the early onset of boundary extension reveals about perception." Intraub and Dickinson











- This ability is critical for us in order to:
  - build a persistent understanding of the world
  - support spatial reasoning

"False memory 1/20th of a second later: what the early onset of boundary extension reveals about perception." Intraub and Dickinson







#### Can we enable machines to do the same?





#### Complete surrounding environment

Introduction

Training Data

3D Representati

n Training Objective







#### Input: RGB-D images

Introduction

Training Data

3D Representation Training Objective



#### Output1: 3D Structures



#### **Output2: Semantics**

Experiments





## Where can I move?

#### Where should I turn to find a door?

Introduction

Training Data

3D Representation Training Objective



#### Output1: 3D Structures



#### **Output2: Semantics**

Experiments



## Semantic-Structure View Extrapolation

#### Input: RGB-D images



Introduction

Training Data

3D Representation Training Objective



Experiments





# Semantic-Structure View Extrapolation

#### Input: RGB-D images



#### Nightstand-

#### Bed

### **Output:** 360° panorama with 3D structure & semantics



Introduction

**Fraining** Data

3D Representati

n Training Objective

Experiments









## Semantic-Structure View Extrapolation

#### Input: RGB-D images



#### Nightstand

#### Bed

### **Output:** 360° panorama with 3D structure & semantics



Introduction

**Fraining** Data

3D Representati

#### Behind camera

n Training Objective

Experiments









### Key idea

**Key idea:** Indoor environments are often **highly structured**. By learning over the statistics of many typical scenes, the model should be able to leverage **strong contextual cues** inside the image to predict what is beyond the FoV.

Data of indoor environments



Introduction

**Fraining** Data

3D Representati

n Training Objective



## Challenges

- How to obtain a large of amount training data?
- How to represent the 3D structure?
- How to provide meaningful supervision for training?



Introduction

Training Data

3D Representati



on Training Objective





### We need a dataset that has:

Semantic label.

- SD structure information.
- Whole room context.

Training Data

Many examples.







### **3D Houses Datasets**



### **Synthetic Houses (SUNCG):**

58,866 RGB-D panoramas Pre-train

Training Data



#### **Real-Word Houses (Matterport3D):** 5,315 RGB-D panoramas Fine-tune and test







#### 3D Scene

Introduction

Training Data

3D Representati

#### Whole Room Sky-box Panorama

n Training Objective





#### Introduction

Training Data

#### 3D Representati

#### Whole Room Sky-box Panorama

n Training Objective





Color Images

Semantics

Introduction

Training Data

3D Representation Training Objective



#### **3D Structure**





### Color Images Color = R,G,B**Standard Data Representation**

**3D** Representation



Semantics **Semantics = ClassId** 

**3D Structure** ??





#### Color Images Color = R,G,B

#### Depth?

![](_page_23_Picture_5.jpeg)

#### Hard to predict.

- •Viewpoint dependent.
- Large value variance even for nearby pixels on the same 3D plane

#### 3D Representation

![](_page_23_Picture_12.jpeg)

#### Semantics Semantic = ClassId

#### **3D Structure**

??

#### Normal?

![](_page_23_Picture_18.jpeg)

Easier to predict. Solving back depth from normal is under constrained.

![](_page_23_Picture_21.jpeg)

![](_page_24_Picture_1.jpeg)

#### color images Color = R,G,B

#### semantic maps Semantic = ClassId

### Challenge 2: How to represent the 3D structure?

Introduction

Training Data

3D Representation

![](_page_24_Picture_8.jpeg)

#### **3D Structure**

on Training Objective

![](_page_24_Picture_12.jpeg)

## **Challenge2: 3D Structure Representation**

![](_page_25_Figure_1.jpeg)

#### normal (a,b,c)

3D Representation

![](_page_25_Picture_8.jpeg)

#### **3D Structure**

#### plane distance (p)

![](_page_25_Picture_13.jpeg)

![](_page_25_Picture_14.jpeg)

## **Challenge2: 3D Structure Representation**

![](_page_26_Figure_1.jpeg)

#### normal (a,b,c)

plane distance (p)

**3D** Representation

![](_page_26_Picture_8.jpeg)

![](_page_26_Picture_9.jpeg)

**3D Structure** 

 $\checkmark$  Pixels on the same planar surface share the same plane equation

✓ Representation is piecewise constant

✓ More robust

![](_page_26_Picture_15.jpeg)

![](_page_26_Picture_17.jpeg)

![](_page_26_Picture_18.jpeg)

![](_page_26_Figure_19.jpeg)

![](_page_26_Picture_20.jpeg)

![](_page_26_Picture_21.jpeg)

## **Challenge2: 3D Structure Representation**

#### Raw Depth Representation

#### **Prediction**

#### **Observation**

Training Objective **3D** Representation

#### Plane Representation

![](_page_27_Figure_7.jpeg)

![](_page_27_Picture_9.jpeg)

### Im2Pano3D Network

![](_page_28_Figure_1.jpeg)

### Challenge 3: What training objectives should we use?

Training Objective

![](_page_28_Picture_9.jpeg)

## Challenge 3: Training Objectives

![](_page_29_Figure_1.jpeg)

Introduction

Training Data

3D Representation Training Objectives

Experiments

![](_page_29_Picture_7.jpeg)

![](_page_30_Figure_0.jpeg)

Introduction

Training Data

3D Representation Training Objectives

### Challenge 3: Training Objectives

Experiments

![](_page_30_Picture_7.jpeg)

## **Challenge 3: Training Objectives**

![](_page_31_Figure_1.jpeg)

Training Data

Training Objectives

**Prediction** is Plausible

![](_page_31_Figure_8.jpeg)

![](_page_31_Picture_9.jpeg)

![](_page_31_Picture_10.jpeg)

![](_page_32_Figure_1.jpeg)

![](_page_32_Picture_6.jpeg)

## **Challenge 3: Training Objectives**

**Every Pixel is** Correct

 $L_{recon}$ 

![](_page_33_Figure_3.jpeg)

Training Objectives

Similar Scene Attribute

Prediction is Plausible

Lattribute

 $L_{adv}$ 

### $L = \lambda_1 L_{recon} + \lambda_2 L_{attribute} + \lambda_3 L_{adv}$

![](_page_33_Picture_13.jpeg)

![](_page_33_Figure_14.jpeg)

![](_page_33_Picture_15.jpeg)

### Every pixel is correct

Semantic Prediction **Pixel-wise** loU

3D Structure

**Pixel-wise** L2 distance

Training Data

### Evaluation

Similar scene attribute

Prediction is plausible

#### Probability over Groundtruth

#### Inception score (Scene classification)

Earth Mover Distance

![](_page_34_Picture_16.jpeg)

![](_page_34_Picture_17.jpeg)

### Evaluation

![](_page_35_Figure_1.jpeg)

Introduction

Training Data

3D Representati

Higher is better

Lower is better

n Training Objective

![](_page_35_Picture_9.jpeg)

![](_page_36_Picture_0.jpeg)

![](_page_36_Figure_1.jpeg)

### Evaluation

![](_page_36_Picture_16.jpeg)

### **Example Results**

Introductior

**Fraining Data** 

3D Representation Training Objective

Experiments

![](_page_37_Picture_6.jpeg)

### Observation

![](_page_38_Picture_1.jpeg)

#### Introduction

#### Training Data

#### 3D Representati

on Training Objective

![](_page_38_Picture_7.jpeg)

![](_page_39_Picture_1.jpeg)

### Ceiling: Red indicates high probability

Introduction

Training Data

3D Representati

Training Objective

![](_page_39_Picture_8.jpeg)

![](_page_40_Picture_1.jpeg)

#### Floor: Red indicates high probability

Introduction

Training Data

3D Representati

Training Objective

![](_page_40_Picture_8.jpeg)

![](_page_41_Picture_1.jpeg)

#### Wall: Red indicates high probability

Introduction

Training Data

3D Representati

Training Objective

![](_page_41_Picture_8.jpeg)

![](_page_42_Picture_1.jpeg)

#### Bed: Red indicates high probability

Introduction

Training Data

3D Representati

Training Objective

![](_page_42_Picture_8.jpeg)

### Semantic Prediction

![](_page_43_Picture_1.jpeg)

### Semantic labels with highest probability per pixel

Introduction

**Training** Data

3D Representati

Training Objective

![](_page_43_Picture_8.jpeg)

![](_page_44_Picture_0.jpeg)

![](_page_44_Figure_2.jpeg)

### Results

![](_page_44_Picture_4.jpeg)

![](_page_44_Figure_6.jpeg)

![](_page_45_Picture_0.jpeg)

![](_page_45_Figure_2.jpeg)

### Results

![](_page_45_Picture_5.jpeg)

![](_page_45_Picture_6.jpeg)

![](_page_46_Picture_0.jpeg)

![](_page_46_Figure_2.jpeg)

### Results

![](_page_46_Picture_5.jpeg)

![](_page_46_Picture_6.jpeg)

![](_page_47_Picture_0.jpeg)

![](_page_47_Figure_2.jpeg)

### Results

![](_page_47_Picture_4.jpeg)

![](_page_47_Picture_5.jpeg)

![](_page_47_Picture_6.jpeg)

![](_page_48_Picture_0.jpeg)

![](_page_48_Figure_2.jpeg)

### Results

![](_page_48_Picture_5.jpeg)

![](_page_48_Picture_6.jpeg)

## How do we compare to people?

#### Observation

![](_page_49_Picture_2.jpeg)

![](_page_49_Picture_3.jpeg)

![](_page_49_Picture_4.jpeg)

![](_page_49_Picture_5.jpeg)

#### Completion by different users

![](_page_49_Picture_7.jpeg)

![](_page_49_Picture_8.jpeg)

![](_page_49_Picture_9.jpeg)

![](_page_49_Picture_10.jpeg)

![](_page_49_Picture_11.jpeg)

### How do we compare to people?

![](_page_50_Figure_1.jpeg)

#### Pixel distance to observation

![](_page_50_Picture_10.jpeg)

### How do we compare to people?

![](_page_51_Figure_1.jpeg)

![](_page_51_Picture_5.jpeg)

![](_page_51_Picture_6.jpeg)

![](_page_51_Picture_7.jpeg)

Baselines

![](_page_51_Picture_12.jpeg)

### Conclusion

#### New Task: Semantic-Structure View Extrapolation

![](_page_52_Picture_2.jpeg)

![](_page_52_Picture_3.jpeg)

### Contextual priors Geometric priors Meaningful supervision

#### Code & Data: im2pano3d.cs.princeton.edu

Introduction

Training Data

3D Representation Training Objectives

![](_page_52_Picture_9.jpeg)

- Two large-scale house level datasets
- 3D plane equation representation
- Multi-level loss functions

![](_page_52_Picture_14.jpeg)

# Im2Pano3D:

### Extrapolating 360° Structure and Semantics Beyond the Field of View Shuran Song, Andy Zeng, Angel X. Chang, Manolis Savva, Silvio Savarese, and Thomas Funkhouser

![](_page_53_Picture_3.jpeg)