Open-World 3D Scene Understanding without Open-World 3D Data

CoRL 2022 Workshop Geometry, Physics, and Human Knowledge as Inductive Bias in Robot Learning







How to introduce prior knowledge into learning models?

- The very first inductive bias we injected to the system
- is the way we define or decompose the tasks (i.e., abstractions)



How to introduce prior knowledge into learning models? The very first inductive bias we injected to the system is the way we define or decompose the tasks (i.e., abstractions)



Abstraction: bounding box or 6DoF pose



Limitations:

- Cannot adapt based on different end tasks
- Error propagation

Limitations:

- Cannot adapt based on different end tasks
- Error propagation

Advantages:

- Decomposition -> Learn from less data
- Discover underlying structure -> Generalize beyond training data
- Abstraction as inductive bias is unavoidable.

Semantic Abstraction: Open-World 3D Scene Understanding from 2D Vision-Language Models Huy Ha Shuran Song Columbia University semantic-abstraction.cs.columbia.edu **Open Vocabulary Scene Completion** Posed RGB-D Input Visually Obscured Object Localization "Used N95s in the garbage bin" desk 🗧 bookshelf 📃 floor 📃 office chair 🔳 trash basket 📕 lego technic liebherr excavator wall with colorful cartoon murals 📒 rubiks cube 📒 harry potter and the sorcerer's stone **CoRL 2022**



"Abstraction as Inductive Bias" in the Age of Large Pre-trained Models



Recognize a dynamic set of <u>semantic</u> categories and ground these semantics with <u>spatial</u> information in the <u>3D environment</u>.







Visual Semantics Concepts X Geometry Concepts X Spatial Concepts





Visual Semantics Concepts

Geometry Concepts

Spatial Concepts

Presents different properties and can be learned in different manner and dataset

Visual Semantics Concepts

Large and non-decomposable









Geometry Concepts

Large but decomposable

Generalized Cylinder Brook & Binford 1979

Learning Shape Abstractions, Tulsiani et al, 2017

Spatial Concepts

Small and finite

Table 1. Prepositions of English

about	between	outside	
above	betwixt	over	
across	beyond	past	
after	by	through	
against	down	throughout	
along	from	to	
alongside	iπ	toward	
amid(st)	inside	under	
among(st)	into	underneath	
around	near	up	
at	nearby	upon	
atop	off	via	
behind	on	with	
below	onto	within	
beneath	opposite	without	
beside	out		
Compounds			
far from	on top of		
in back of	to the left of		
in between	to the right of		
in front of	to the side of		
in line with			
Intransitive prepe	ositions		
afterward(s)	ferward	right	
apart	here	sideways	
away	in ward	south	
back	left	there	
backward	N-ward (e.g.,	together	
downstairs	homeward)	upstairs	
downward	north	upward	
east	outward	west	

Spatial Language, Landau & Jackendoff, 1993

Visual Semantics Concepts

Requires exposure to internet-scale datasets

>>

Geometry Concepts Spatial Concepts

Tractable even with limited synthetic datasets

>>

Visual Semantics Concepts



Requires exposure to internet-scale datasets

Geometry Concepts Spatial Concepts

✓ Understands a large collection of visual-semantic concepts

Robustness by learning from internet scale data.

Visual Semantics Concepts



Requires exposure to internet-scale datasets

>>

Geometry Concepts Spatial Concepts



Finetune with 3D data !?

>>

Visual Semantics Concepts

	Combined S	caling for Open-Vocabulary Ima	ge Classification		
Hieu Pham*			HYREU@ GOOGLE.COM		
Zihang Dai" Golnaz Ghiasi" Kenji Kawaguchi" Hansian Liu Adams Wei Yu		Robust fine-tuning of zero-shot models			
5	Jiahui Yu Yi-Ting Chen		Mitchell Wortsman*"	Gabriel Ilhanco ^{*†}	Jong Wook Kim [§] Mike Li [†]
r 202	Minh-Thang Luong Yonghui Wu Mingying Tan		Simon Komblith [*]	Rebecca Roelofs*	Raphael Gontijo-Lopes ^a
Quoe V. Le *: Equal contributions. 60 Corresponding authors: {HYP 61 Editor: To be assigned.	The Evoluti	on of Out-of-Distribution Robus Throughout Fine-Tuning	tness	ong ^{*‡} Ludwig Schmidt [†] curacy across a range of data in a specific dataset). Although larget distribution, they often aducing a simple and effective of the zero-shot and fine tunad	
		Anders Andreass	Anders Andreassen, Yasaman Bahri, Behnam Neyshabar, Reberen Raciofs Google Research {ajandreassen, yasamanb, neyshabur, rofis}@google.com		i large accuracy improvements ibotion. On ImageNot and five irm shift by 4 to 6 percentage b. WISE-FT achieves similarly in shifts, and accuracy gains of ansfer learning datasets. These is or inference.
		m 2(Abstract		
		Although machin on out-of-distribu widely observed t models. Models th baseline exhibit " models, and under oerformatice. We	e learning models typically experience a drop in perform ation data, accuracies on in-versus out-of-distribution dato o follow a single linear trend when evaluated across a test bat are more accurate on the out-of-distribution data relative te effective robustness" and are exceedingly rate. Mentifying rstanding their properties, is key to improving out-of-distrib conduct a thereach empirical investigation of effective robustness	tance to are ed of o this such ution duess	

Naïvely finetuning VLMs on a limited and task-specific dataset results in reduced model's robustness and generality.

Geometry Concepts Spatial Concepts



Finetune with 3D data !?



Visual Semantics Concepts

Geometry Concepts Spatial Concepts

Goal: Reason about 3D concepts in a semantic-agnostic manner.

Visual Semantics Concepts



2D localization of the object

Geometry Concepts Spatial Concepts

Goal: Reason about 3D concepts in a semantic-agnostic manner.

Learn different concepts for

"behind the Harry Potter book"

"behind the trashcan"

. . .

Learn one semantic-agnostic concept of ``behind that object"



Visual Semantics Concepts



Pixel-wise Relevancy Map: represents the probability of the object's locations.

Geometry Concepts Spatial Concepts

Visual Semantics Concepts



Frozen 2D VLM (Semantic-aware)

Visual Semantics Concepts



Frozen 2D VLM (Semantic-aware) **Geometry Reasoning** (Shape Completion)



Spatial Reasoning (Localization)



Relevancy Map From 2D LVM

Semantic-Agnostic Does not observe color or semantic label



Visual Semantics Concepts



Frozen 2D VLM (Semantic-aware)

From 2D LVM

Tools for Attention Extraction (GradCam)



Grad-CAM for "Cat"





Grad-CAM for "Dog"



Grad-CAM. Selvaraju et al.

Class Activation Map Australian terrier











Generic Attention-model Explainability. Chefer et al.

Noisy & Low-res relevancy maps





Issues: 1. Noisy & highlights irrelevant regions 2. Misses small objects



Multi-scale Relevancy Extractor



endo Switch Ninte



1. Apply many image augmentations, and then average the resulting relevancy maps to reduce noise

Multi-scale Relevancy Extractor



endo Switch Ninte





3) Extracting relevancy at a multiple scale.



Semantic-Abstracted 3D Module



Relevancy Map

interpret it as a rough indicator to the 3D network about which object is referred and need to be complete in the scene.

Relevancy value does not need to be perfect, completion network is trained to correct small errors

3D Shape Completion

Semantic-Abstracted 3D Module



Together with the text description provided in the beginning, we can get the completed shape for a particular object referenced in text.



Open-Vocabulary Scene Completion



Note: If the input object is not visible or cannot be find the 2D VLM, the relevancy map produce low relevancy score for the whole image.

Provide additional spatial reference to help localize the object



Provide additional spatial reference to help localize the object

THE P







Step 1: Use the same shape completion module to find and complete the reference (fireplace) and target object (CoRL ticket)





Step 2: Localize target object based on spatial relations (e.g., on top of)

Groundtruth from Simulation (AI2-THOR)







Dense 3D label for both the shape completion task and the spatial localization task.

AI2-THOR: An Interactive 3D Environment for Visual AI, Kolve et al.



Open-World Evaluation

We hope to inherent the generality from the 2D VLMs



Open-world generalization



Novel Domain (Real)



(+) OOD Nouns

Novel semantic classes



(+) OOD 3D Geometry



Qualitative Results on Matterport3D dataset

ARKitScenes - A Diverse Real-World Dataset for 3D Indoor Scene Understanding Using Mobile RGB-D Data, Baruch et al

Input RGB-D

Scene Completion



OOD visual (e.g., painted wall)



"used N95s in the garbage bin"

Semantic Labels

office chair

floor

desk

Small

Longta

trash basket

book shelf

rubiks cube

harry potter book

lego technic excavator

wall with colorful cartoon murals

"COVID-19 rapid test behind Harry Potter book"







COVID-19 rapid test







Input RGB-D





"vaccination card in the red folder"

Scene Completion



"vaccination card in the **blue folder**"

Semantic Labels

floor

wall

ceiling

light

book shelf

house plant

globe

black aluminum file cabinet









Low relevancy value

vaccination card





Qualitative Results on Apple ARKitScenes dataset **Different environments, Different camera ... Direct test without fine-tuning**

ARKitScenes - A Diverse Real-World Dataset for 3D Indoor Scene Understanding Using Mobile RGB-D Data, Baruch et al

Input RGB-D









"Hair dryer with its wires tangled behind chair legs" "sunscreen bottle in pink make up bag"

Scene Completion

Semantic Labels

lamp

wall

table

light switch

mirror

carpet

pink make up bag

upholstered chair in faux leather

woven chair

Can handle long descriptions

















Results on NYUv2 for <u>all 894 classes</u>



Application: CLIP on Wheels



The modularized design (i.e., localization & exploration) allow us to directly apply VLMs to object navigation without additional training on navigation.

Language Driven Zero-Shot Object Navigation



Evaluating CoWs on Pasture Benchmark



~20 different variant of CoWs

Study capabilities that closed-vocabulary object navigation agents do not possess



1) Find uncommon objects

llama wicker basket



4) Hidden objects



"mug under the bed"



Evaluating CoWs on Pasture Benchmark



3) Finding object based on attributes

"...small, green apple..." (appearance) "...apple on a coffee table near a laptop..." (spatial)



4) Zero-shot to different domains















1) Find uncommon objects

llama wicker basket



4) Hidden objects



"mug under the bed"



Evaluating CoWs on Pasture Benchmark



3) Finding object based on attributes

"...small, green apple..." (appearance) "...apple on a coffee table near a laptop..." (spatial)



Require a deeper understanding of the descriptions than treating them as a bag of words

Evaluating CoWs on Pasture Benchmark





CoWs on PASTURE: Baselines and Benchmarks for Language-Driven Zero-Shot Object Navigation. Samir Gadre et al



Zero-shot CoW (no training on navigation task) is better than fine-tuned SOTA model!



Advantage:

- Learn from less data
- Better generalization

Semantic Abstraction

Advantage:

- Learn from less data
- Better generalization

Address limitation:

- Adapt to different tasks with different reusable modules
 - e.g., shape completion and localization
- Learnable downstream model to reduce error propagation

Semantic Abstraction

In the age large pre-trained models ...



- Large models learned from Internet knowledge
- Fine-tuning them for specific application can be impractical (e.g., compute) or challenging (e.g., hurting robustness)
- The right abstraction allow us to extract the relevant knowledge from them in a zero-shot manner (w.o. fine-tuning)

Extend New skills ...





Manipulation

Do As I Can, Not As I Say: Grounding Language in **Robotic Affordances**

I would:

- 1. find a sponge
- 2. pick up the sponge
- come to you 4. put down the sponge
- 5. done

Combine Different Expertises ...



Socratic Models: Composing Zero-Shot Multimodal Reasoning with Language

What's Next? Your Work!

plant inside WALL-E



Thank You! Question?

