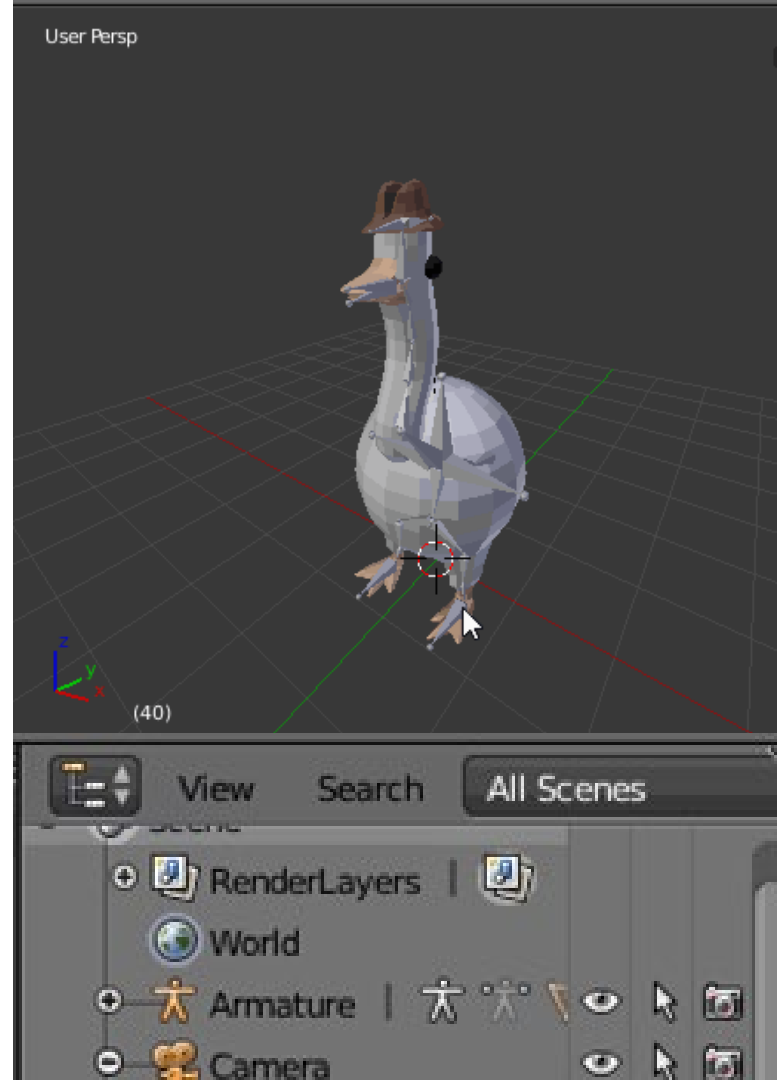


# Lecture 5

## 3D models: MLLMs application for 3D CV

Zinkovich Viktoriia

Сделайте это задание  
и получите сертификат





# Introduction: Course Plan



Day 1  
Image modality

Day 3  
Data generation  
in MLLMs



Day 5  
3D models



Day 2  
Video modality



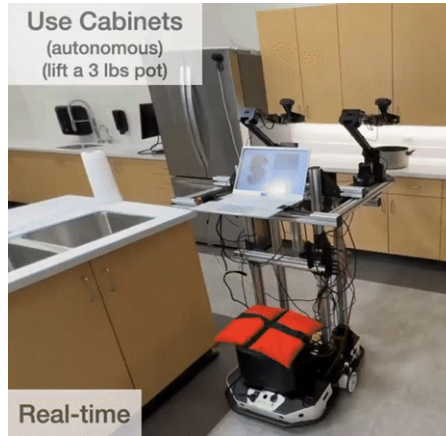
Day 4  
Action modality

# Recap: Lecture #4



many embodiments  
characterized by **how we**  
**can control it**

# Recap: Lecture #4



|| 1/2 T Y S || í || ñ e || S T E E  
L 1/2 C L I S C V S ñ E Y E || Ü H S E  
L 1/2 T H || T C || E I

evaluation can be done by  
**humans**, teleoperation,  
data-gloves

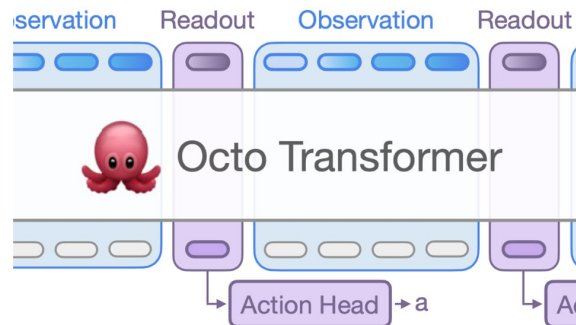
# Recap: Lecture #4



many embodiments  
characterized by **how we**  
**can control it**



evaluation can be done by  
**humans** , teleoperation,  
data-gloves

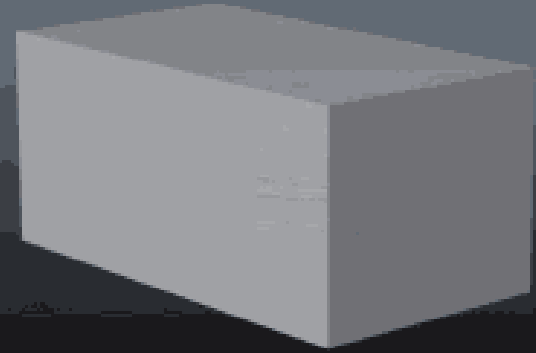


**policies** to train models for  
robotic tasks: RT-1, RT-2,  
Octo, Diffusion, OpenVLA

# ÜSLĪĴ ĆŚĒ ½

- 1 What is 3D vision ?
- 2 How we get **information** about 3D?
- 3 Tasks & **Datasets**
- 4 ° ΔĒĴ ½ ½

**modeling is easy**



# Introduction: What is 3D vision?

The Turing test (1949 originally) for **computer vision** — answer any question about the image that a human can answer



**semantic questions**  
(where / how many / text)

# Introduction: What is 3D vision?

The Turing test (1949 originally) for **computer vision** — answer any question about the image that a human can answer



**semantic questions**  
(where / how many / text)



**metric questions**  
(shape / distance)

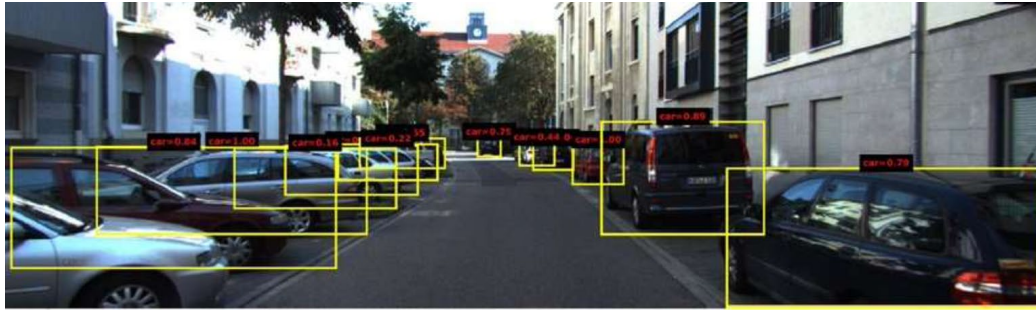
# Introduction: What is 3D vision?

The difficult task is **total 3D reconstruction**, we need to build an accurate 3D model that allows: 1) manipulation with objects; 2) image synthesis

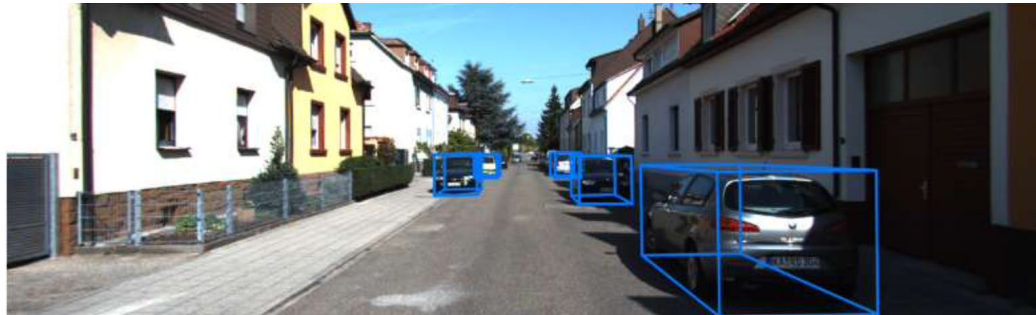


$p = \frac{1}{Z} \int \mathcal{L}(\theta) d\theta$   $\updownarrow$   $\frac{1}{2} \sum_{i=1}^n \mathcal{L}_i(\theta)$   $\rightarrow$   $\mathcal{K}$

2D



3D



3D vision =  
Metric Vision + Semantic Vision

1

---

k | Ũ Ĥ | Ę Ś Ĩ Ę | Ę | Ę | ½ Ĩ | Ę  
½ Ĩ | Ĥ Ĩ Ę Ā B Ę

Sensors, LiDARs...



# 3D Info: 1. Structure from Motion

Find key points on the images, define camera motion — “**Structure from Motion**” → sparse point cloud that can be further improved (SIFT + RANSAC)



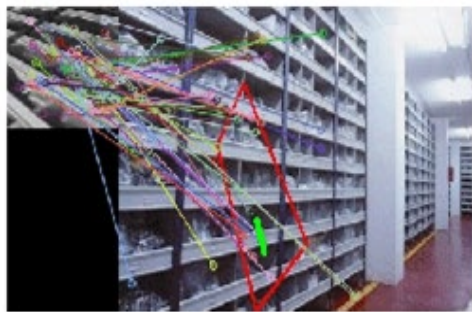
camera poses

# 3D Info: 1. Structure from Motion

Find key points on the images, define camera motion — **“Structure from Motion”** → sparse point cloud that can be further improved (SIFT + RANSAC)



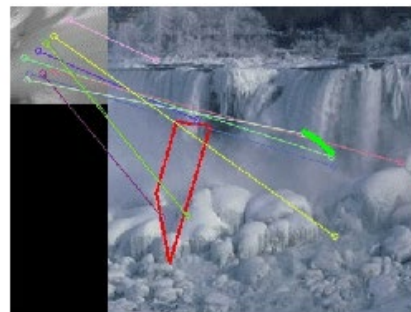
Ground truth



SIFT+RANSAC



Ground truth



SIFT+RANSAC

# 3D Info: 1. Structure from Motion

The problem with sparse reconstruction: **sparse** and **inaccurate** (noisy data) + **distortions** (the walls)



AB  $\dot{p} = \tau!$  |  $\ddot{E}$   $\dot{S} B \dot{S} \circ \ddot{I} \updownarrow \ddot{E} \frac{1}{2} || \dot{S} \hat{C} \frac{1}{2} \ddot{E}$

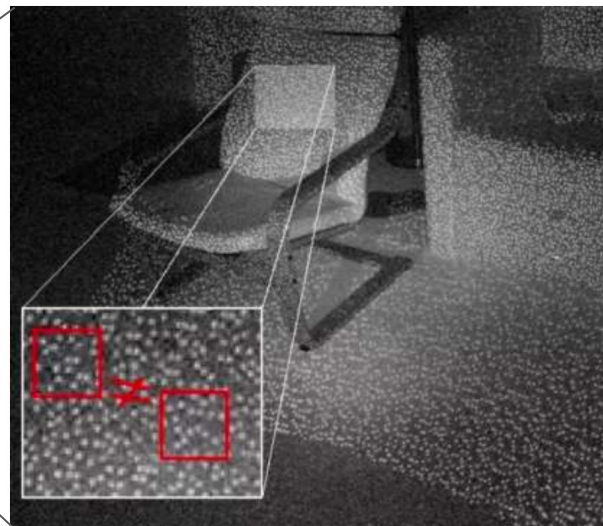
RGB image + map of distances to objects  $\rightarrow$  depth of the point pattern distortion



PrimeSense technology,  
licensed by Microsoft and  
implemented in the Kinect for  
Xbox 360 camera



Projector



"Smart" structural illumination — set of  
spots according to a tricky pattern

# 3D Info: 2. Depth Cameras

RGB image + map of distances to objects → depth of the point pattern distortion

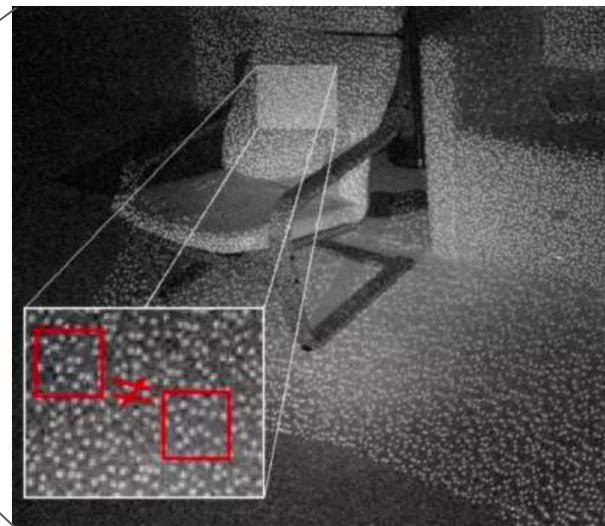


PrimeSense technology,  
licensed by **Microsoft** and  
implemented in the Kinect for  
**Xbox 360** camera



Projector

Time -of-Flight  
(Kinect-2)



"Smart" structural illumination — set of  
spots according to a tricky pattern

## 3D Info: 2. Depth Camera s

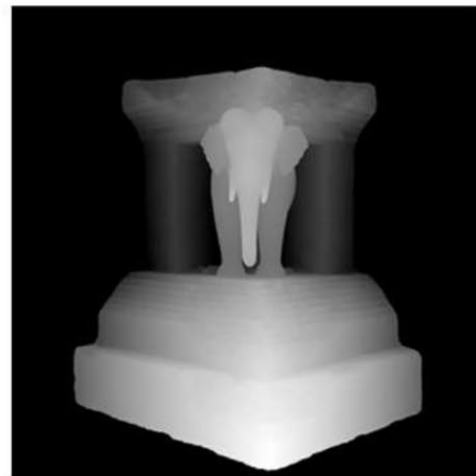
→ RGB image is combined with depth (D) to **obtain RGBD** — **CNN net!**



→ To combine depth with RGB need to **calibrate** camera poses

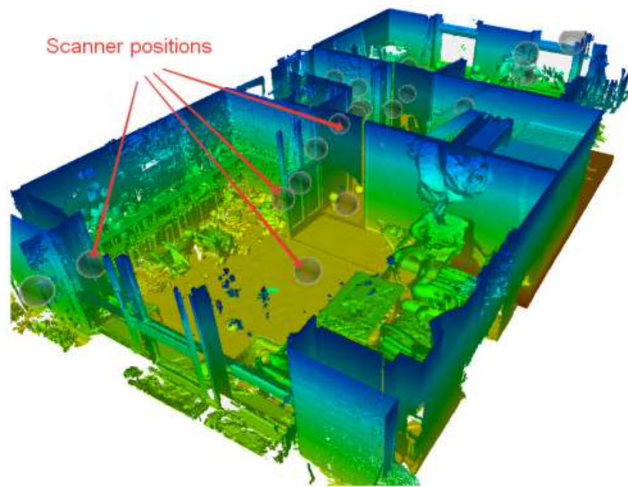


→ Sensors have **threshold** for depth estimation



# AB $\dot{p} = \dot{r} + \dot{r} \times \dot{\theta} + \dot{\theta} \times r + \dot{\theta} \times \dot{\theta} \times r$

Unlike RGB or depth cameras, **LIDARs** immediately provide a full 3D representation, not just a depth map — **time difference** between the emission and reception of the signal (phase shift)

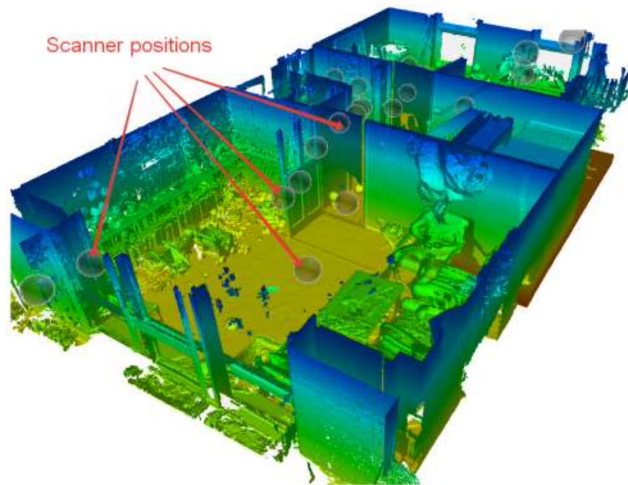


# 3D Info: 3. Laser Scanners

→  $\vec{r} = r \cdot \vec{e}_r$   
 $\vec{r} = r \cdot \begin{pmatrix} \cos\theta \cos\phi \\ \cos\theta \sin\phi \\ \sin\theta \end{pmatrix}$   
 $\vec{r} = r \cdot \begin{pmatrix} \cos\theta \cos\phi \\ \cos\theta \sin\phi \\ \sin\theta \end{pmatrix}$

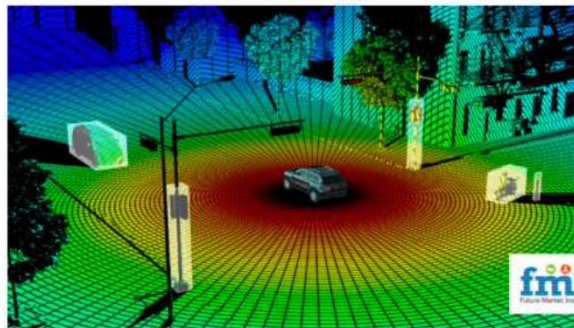
→ Sparsity of data  
(where the laser did not reach the point)

→ Solid -state LiDARs  
(autonomic vehicles)



# 3D Info: 4. Point Cloud

- $\vec{r} \approx \frac{1}{2} \left( \vec{r}_1 + \vec{r}_2 \right)$  !  
 $\vec{r}_1 = \vec{r}_2 = \vec{r}$   $\vec{r}_1 \cdot \vec{r}_2 = r^2 \cos \theta$   
 $\vec{r}_1 \cdot \vec{r}_2 = r^2 \cos \theta = r^2 \cos \theta$
- $\vec{r} = \frac{1}{2} \left( \vec{r}_1 + \vec{r}_2 \right)$   
 $\vec{r}_1 \cdot \vec{r}_2 = r^2 \cos \theta$
- $\vec{r} = \frac{1}{2} \left( \vec{r}_1 + \vec{r}_2 \right)$   
 $\vec{r}_1 \cdot \vec{r}_2 = r^2 \cos \theta$



# 2

---

## 3D tasks and datasets

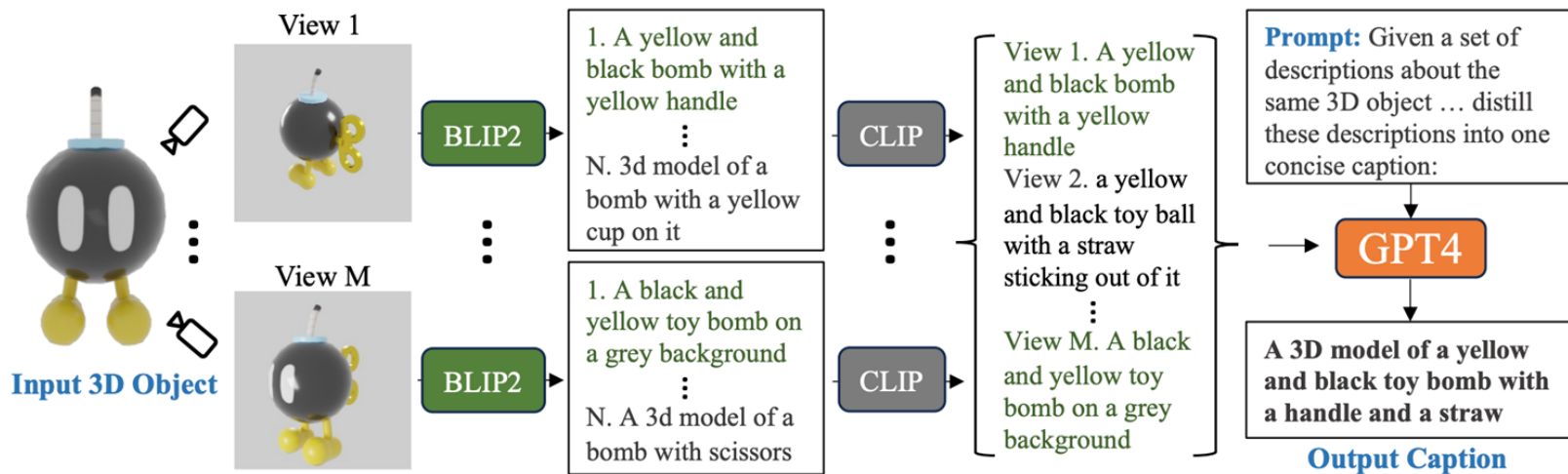
Captioning, segmentation





# Tasks & Datasets: 3D Captioning

Because all components (BLIP2, CLIP, GPT-4) are pretrained, Cap3D requires **no human annotation** and scales linearly with the number of 3D assets



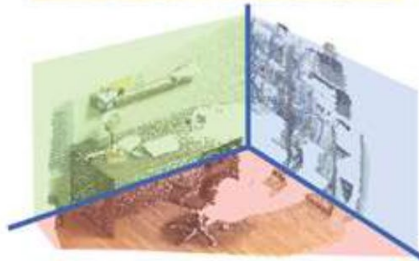
# Tasks & Datasets: SUN RGB-D (2015)

- **3D variant** of detection and segmentation tasks
- **10,000 images** from different cameras
- each scene was **shot 1 time** , without video
- **58,657** bounding boxes

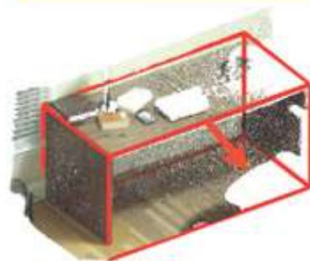
Scene Classification



Semantic Segmentation



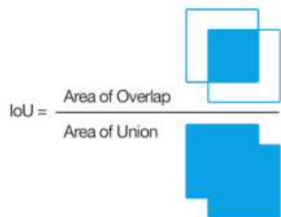
Room Layout



Detection and Pose



# Tasks & Datasets: SUN RGB-D (2015)

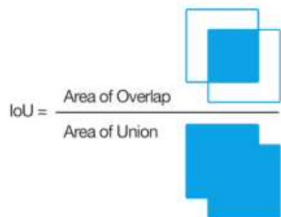
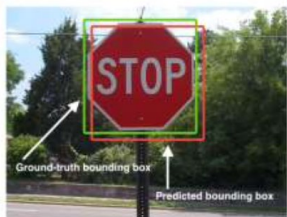


How to calculate **IoU metric in 3D** ?



Total Scene Understanding

# Tasks & Datasets: SUN RGB-D (2015)



**Oriented IoU** — matching in Z-axis + intersection of 2 projected polygons

**Average Precision (AP)** — for all classes of objects and several IoU thresholds, mAP@0.25 and mAP@0.50

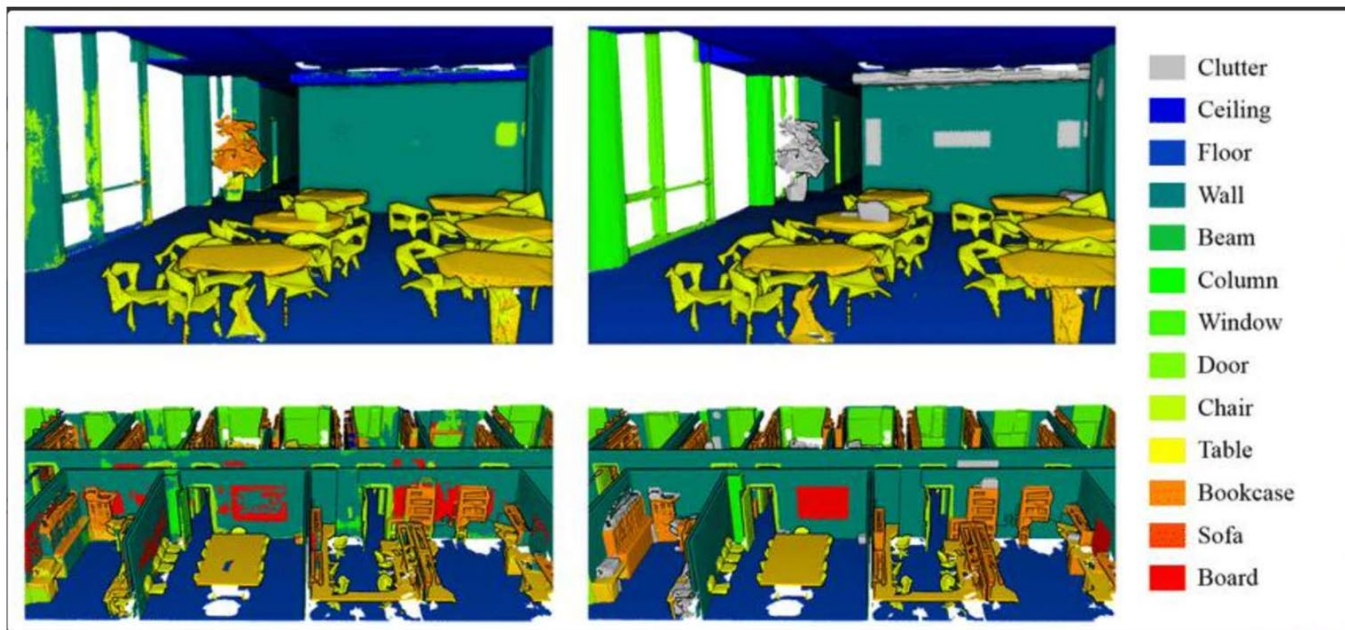


Total Scene Understanding

# Tasks & Data sets: 3D Segmentation

## 3D Segmentation:

classification of scene points – you must put a class label





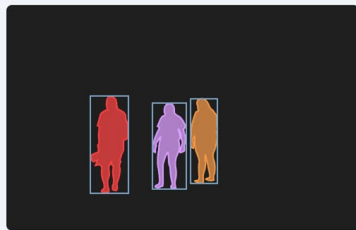
# Tasks & Datasets: Instance, Panoptic



(a) Image



(b) Semantic Segmentation



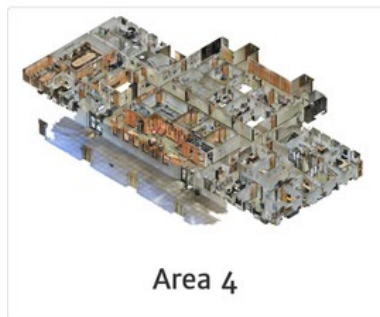
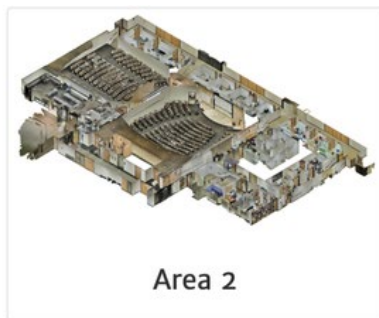
(c) Instance Segmentation



(d) Panoptic Segmentation

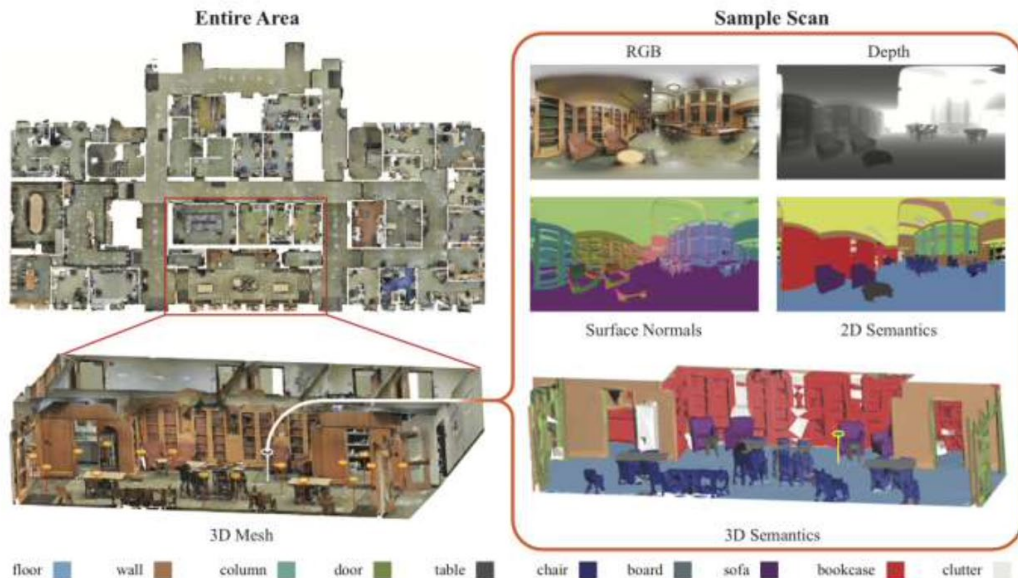
- **Instance** — cloud markup points with labels of individual instances of objects
- **Panoptic** — instance + semantic segmentations

# Tasks & Datasets: S3DIS (2016)



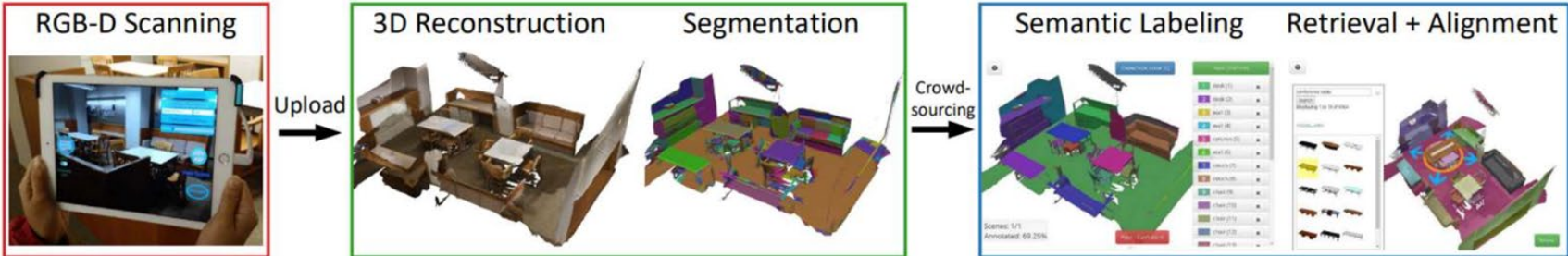
- Scanning **6 large - scale areas** on the Stanford University campus
- Total of **≈272 rooms** — offices, classrooms, corridors, conference rooms
- **215 million three - dimensional XYZ coordinates** with RGB color

# Tasks & Datasets: S3DIS (2016)



- **13 semantic classes** (wall, floor, ceiling, column, etc.)
- **Instance tags:** each individual object is assigned a unique ID
- **Markup files** are supplied for each combined cloud.txt with XYZRGB and two-label markup

# Tasks & Datasets: ScanNet (2017)

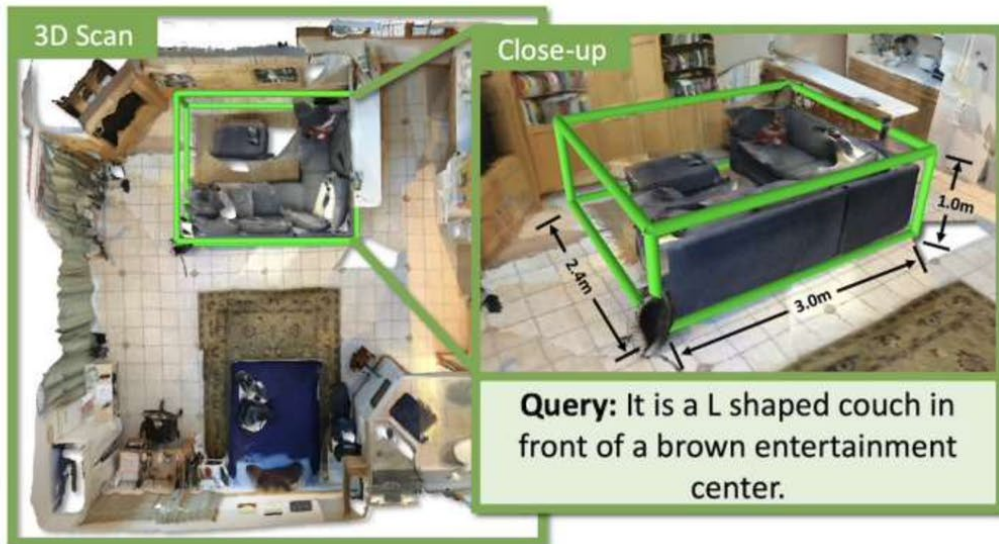
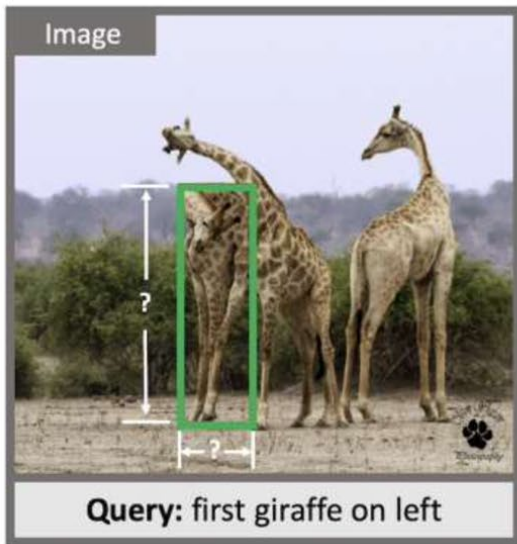


- $\hat{u} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ ,  $\hat{v} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}$ ,  $\hat{w} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$
- $\hat{u} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ ,  $\hat{v} = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ -1 \end{bmatrix}$ ,  $\hat{w} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$

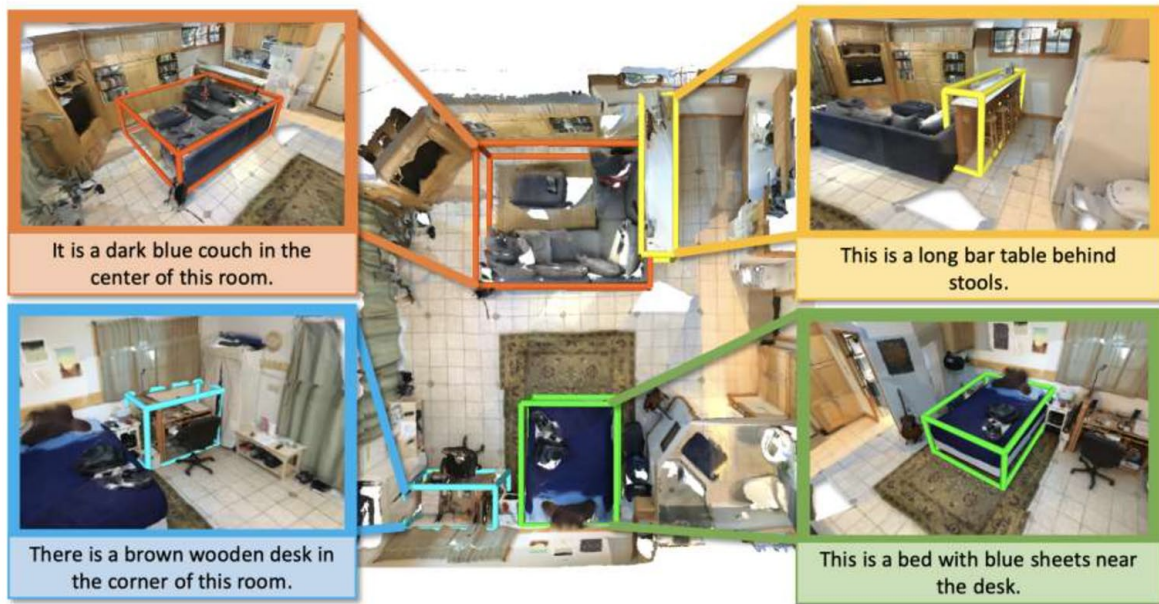


# Tasks & Datasets: 3D Visual Grounding

It's difficult to mark up rooms again, we can make **add-ons on top of this dataset** —  
localization of an object by text description



# Tasks & Datasets: ScanRefer (2020)



- 51,583 text descriptions for 11,046 objects
- 800 ScanNet scenes
- Manual markup via Amazon Turk
- Checking by **students** – they need to find the object
- Still one of the **main datasets** for 3D visual grounding

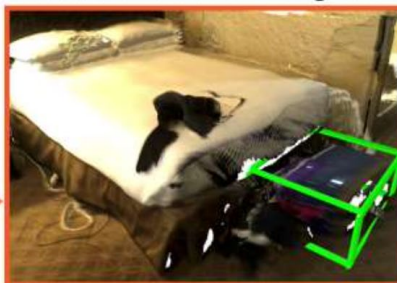
# Tasks & Datasets: 3D Question Answering

Question + 3D-Scan

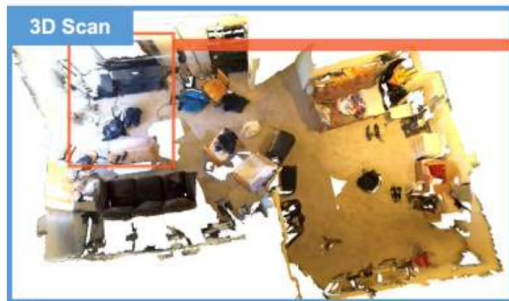


Q. Where is the medium sized blue suitcase laid?

Answer + 3D-Bounding Box



A. in front of right bed



Q. What is sitting on the floor between the tv and the wooden chair?



A. 2 black backpacks

- **Extension** of 3D visual grounding
- Properties of objects, their **relative location**, and how to **find** group of objects
- **Answer** can be **different** (not only BB): 1) text + BB, 2) text, 3) text + multiple BB

# Tasks & Datasets: ScanQA (2022)

Split	# Question	# Unique Question	# 3D Scenes
Train	25,563	20,546	562
Val	4,675	4,306	71
Test w/ objects	4,976	4,552	70
Test w/o objects	6,149	5,484	97
Total	41,363	32,337	800

Table 2. ScanQA dataset statistics.

- ScanNet + ScanRefer + additional markup
- incorrect / irrelevant questions



### Underspecified questions

- Q: What is in the corner?  
- Several objects at corners!
- Q: What color is the chair?  
- Three chairs at the scene!

### Valid questions

- Q: What is over the chair beneath the blackboard?  
- Answer: jacket
- Q: What color is the office chair next to the desk with a monitor?  
- Answer: green

# Tasks & Datasets: Irrelevant QA

## RefCOCO dataset

initially used that dataset, which have simple short prompts,  
**samples** from the dataset:



# Tasks & Datasets: Irrelevant QA

## RefCOCO dataset

Explored **6 state -of-the-art models** : HIPIE, GLEE-Pro, UniLSeg, PolyFormer, UNINEXT, VLT

Found out that all leading models are not robust to **'negative prompts'**

GLEE [CVPR 2024]



HIPIE [NeurIPS 2023]

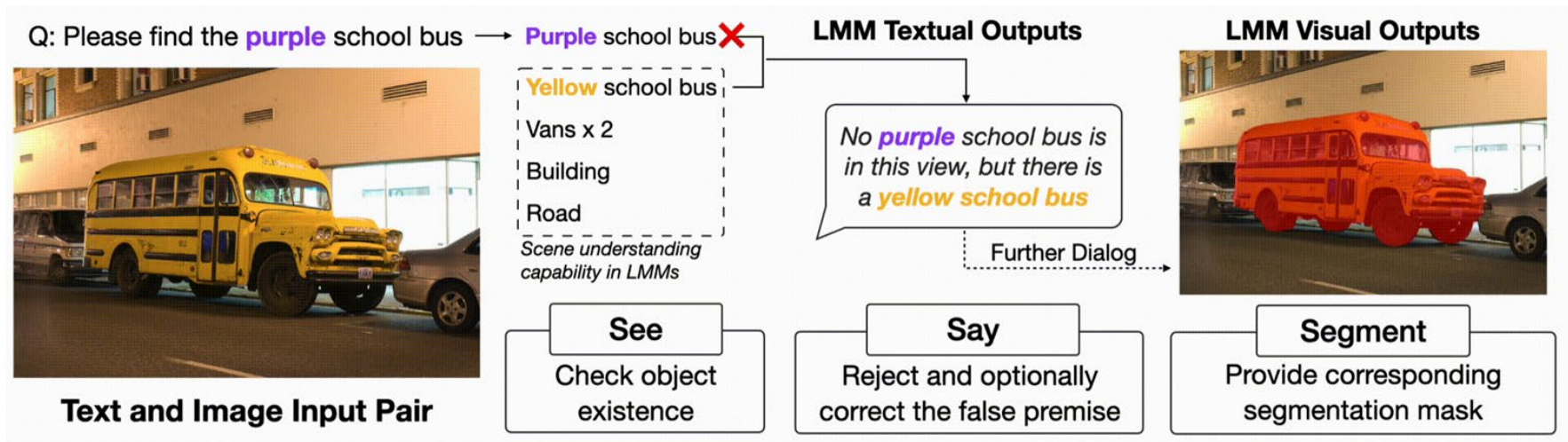


UniLSeg [CVPR 2024]



# Tasks & Datasets: Irrelevant QA

Using refCOCO for base images, authors employ an LLM to augment a false-premise referring segmentation dataset with **context -aware false premise queries**



# Tasks & Datasets: Scene Verse (2024)



## Scene Caption

### Sub-graph Context

```
{ 'scene_type': 'Bedroom',
  'object_count': {'nightstand':2, ...},
  'relation': {'nightstand', 'on', 'floor',
              {'backback', 'in front of', bed}, ...}
```



3D Sub-graph

### Summary

**Prompt:** Provide a summary for a scene from a given scene graph delimited by triple backticks, ...

**Response:** In this bedroom, there are two nightstands, ... The backpack is in front of the nightstand as well. The room appears to be functional, with the nightstands providing storage space and the telephone for communication.

## Object Caption

### BLIP2 Captions

1. A bed in a hotel room. (0.85)
2. A white comforter on a bed. (0.83)
3. A bed with a striped comforter. (0.83)
- ...
- N. A picture of cat. (0.63)



Multiview Images

### Summary

**Prompt:** Summarize the captions below. The summary should be a description of the {object}. Focus on the {object}'s attributes, like color, shape, material, etc. Identify and correct the potential errors ...

**Response:** The bed is in a hotel room with a striped comforter. It has a white comforter and a blanket on it. The bed is also in a room with a bedside table.

## Object Referral

### Relationship Triplets

1. {'table', 'chair', 'left'}
2. {'bed', ('lamp', 'mini fridge'), 'between'}

### Template-based Referral

1. The table is to the left of the chair.
2. It's a bed in the middle of a lamp and the mini fridge.

### Rephrasing

**Prompt:** Rewrite the following sentence using one random sentence structure. Focus on the location and relationships about the {target\_object}, ...

#### Response:

1. The table is situated to the left of the armchair.
2. The bed occupies the space between the lamp and the mini fridge, creating a cozy atmosphere.

# Tasks & Datasets: Summary

# Tasks & Datasets: Summary

- 3D object **captioning** (3D **classification**)
- 3D **detection**
- 3D **semantic** , 3D **instance** , 3D **panoptic** segmentation
- 3D **visual grounding** as a development of the 3D detection task
- 3D **Question Answering**

# 3

---

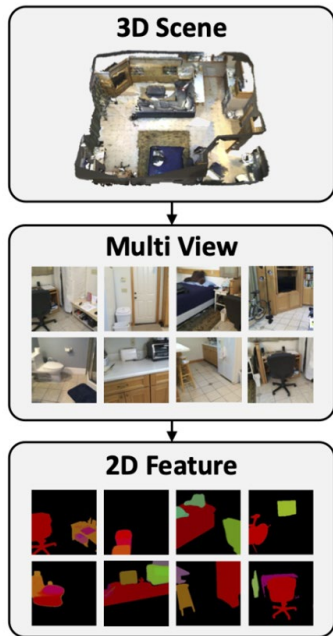
◦  $\hat{E} \hat{J}^{1/2} \hat{L} \hat{C}^n \hat{S} \hat{E} \hat{T}^n \hat{J}^{1/2} \hat{S} \hat{E}$

|| ■ ň Š Ě

3D-LLM, SpatiaIRGPT



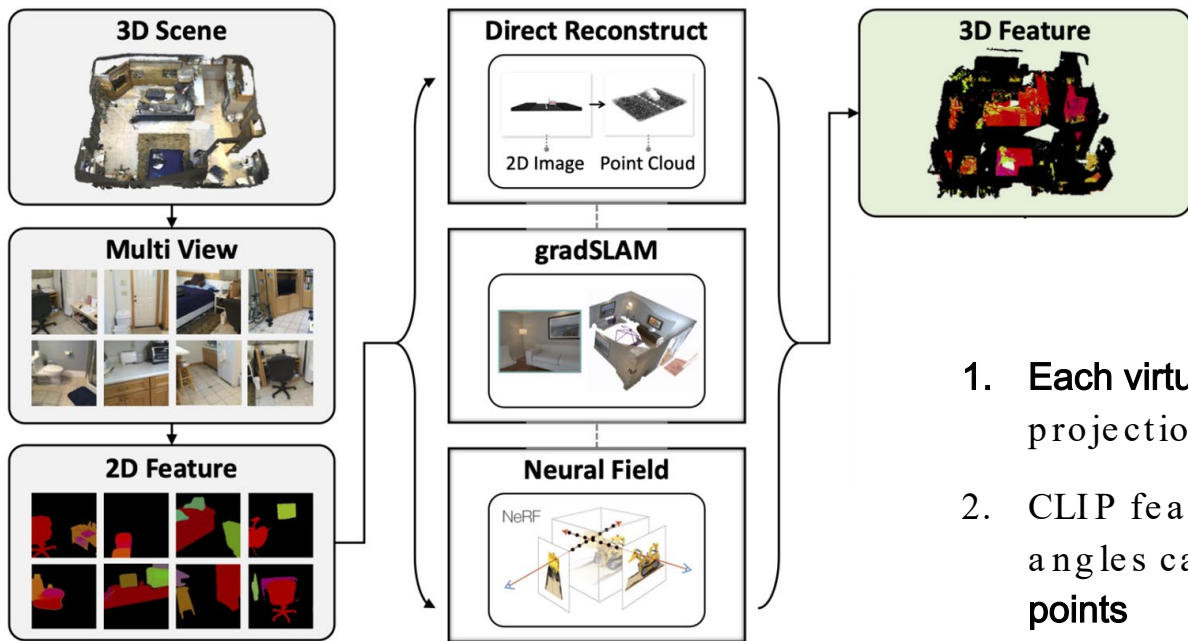
# VLLMs: 3D-LLM (first LLM for 3D)



CLIP features are **meaningful for 2D case!**

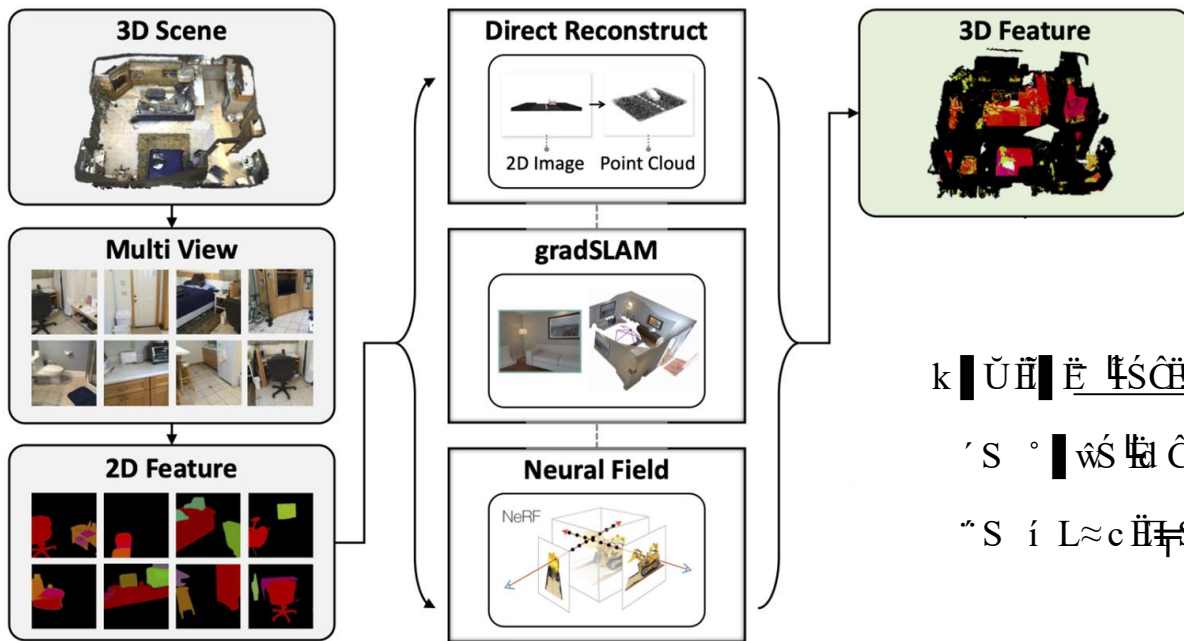
1. Take our **3D scene** (a cloud of dots or a mesh of a room).
2. Using virtual cameras, we **gather photos of the scene** from different angles
3. For each render, we run CLIP and get a feature vector

# VLLMs: 3D-LLM (first LLM for 3D)



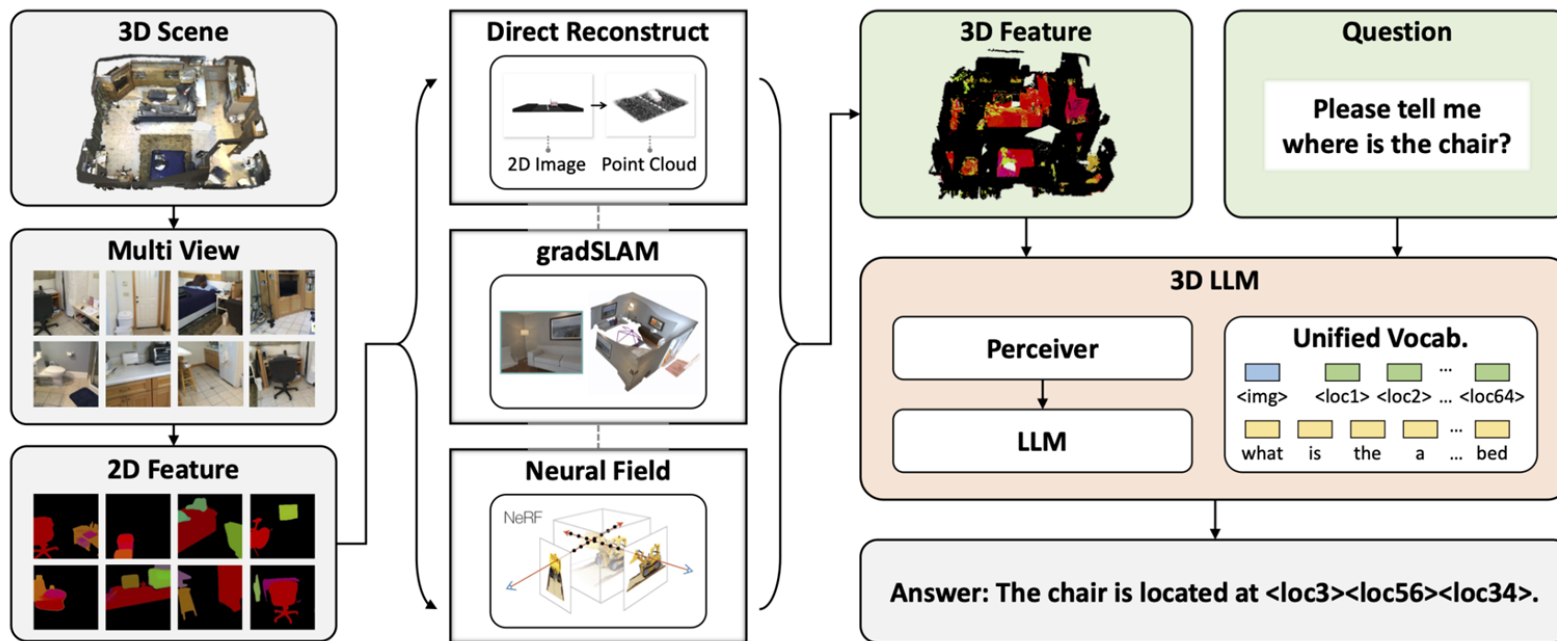
1. **Each virtual camera** has its own projection matrix and depth map
2. CLIP features from different angles can be "sewn" to **3D points**

# VLLMs: 3D-LLM (first LLM for 3D)

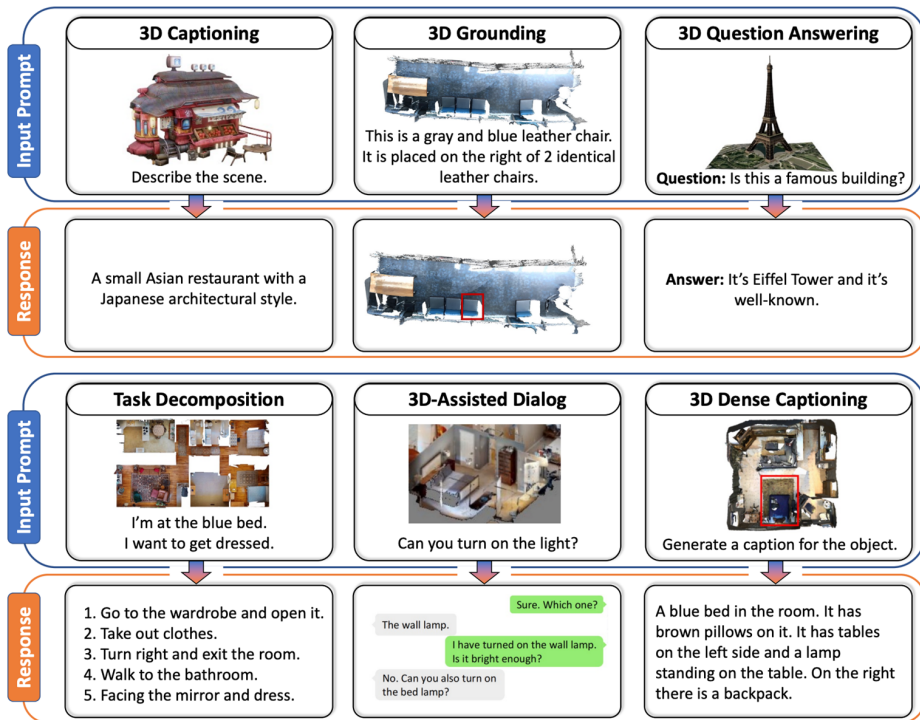


$k \parallel \tilde{U} \tilde{H} \parallel \tilde{E} \parallel \mathbb{1}_{\text{SCE}} \text{S} \frac{1}{2} \mathbb{1}_{\text{J}} \hat{\text{C}} \text{S} \mathbb{E} \parallel \frac{1}{2} \circ \tilde{I}$   
 $\text{'S} \circ \parallel \hat{w} \text{S} \parallel \hat{c} \parallel \hat{c} \parallel \hat{H} \parallel \hat{u} \parallel \hat{S} \parallel \hat{C} \parallel \hat{P} \parallel \hat{P}$   
 $\text{"S} \hat{I} \hat{L} \sim \hat{c} \parallel \hat{H} \parallel \hat{S} \parallel \hat{J} \parallel \hat{C} \parallel \hat{L} \parallel \hat{S} \parallel \hat{I} \parallel \hat{U} \parallel \hat{G} \parallel \hat{P}$

# VLLMs: 3D-LLM (first LLM for 3D)

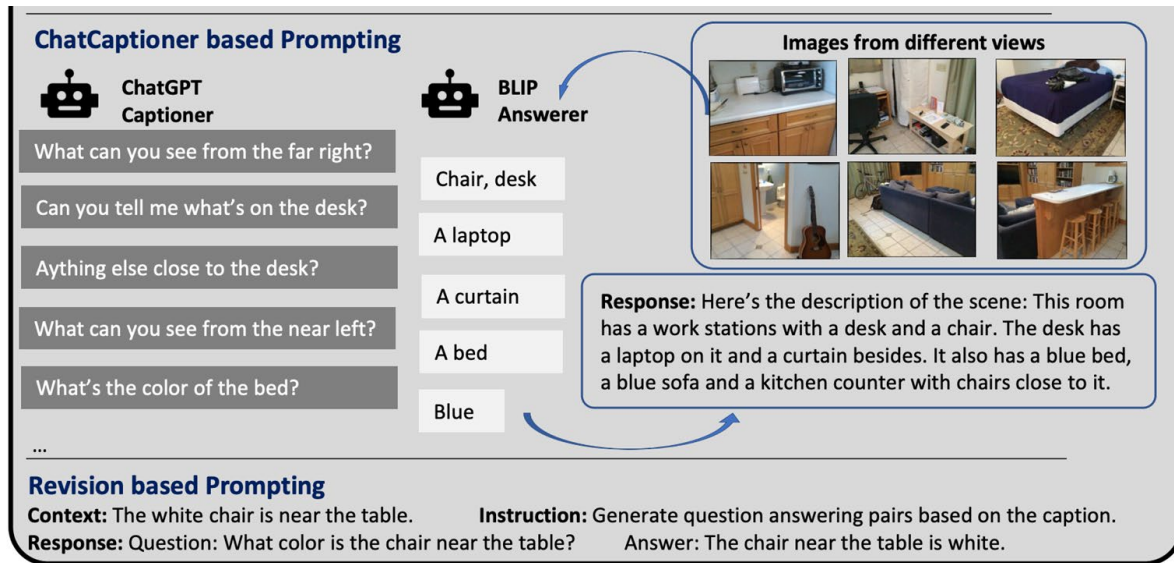


# 3D-LLM: Data Generation



- $\tilde{I} \tilde{H} \downarrow \frac{1}{2} \tilde{I} \tilde{H} \text{el} \parallel \text{SIE} \text{L} \parallel \tilde{d} \tilde{S} \tilde{n} \tilde{h} \tilde{h} \frac{1}{2} \tilde{I} \frac{1}{2} \tilde{E} \tilde{S} \tilde{I} \tilde{E} \tilde{E} \parallel \tilde{C} \tilde{E}$   
 $^{\circ} f \tilde{H} \tilde{H} \parallel \tilde{U} \tilde{H} \tilde{S} \tilde{E} \tilde{J} \frac{1}{2} \tilde{a} \tilde{S} \tilde{E} \tilde{H} \tilde{L} \tilde{S} \tilde{T} \tilde{S} \tilde{B} \parallel \tilde{C} \tilde{E} \tilde{P} \tilde{E}$
- $^{\circ} \tilde{d} \tilde{E} \tilde{J} \frac{1}{2} \tilde{H} \parallel \tilde{n} \tilde{S} \tilde{L} \tilde{E} \tilde{U} \parallel \tilde{G} \tilde{n} \tilde{S} \tilde{n} \tilde{H} \tilde{U} \tilde{d} \tilde{J} \tilde{E} \tilde{a} \tilde{n} \tilde{a} \tilde{n} \tilde{J} \frac{1}{2} \tilde{L} \tilde{E}$   
 $\text{el} \parallel \frac{1}{2} \tilde{S} \tilde{E} \tilde{I} \tilde{E} \parallel \tilde{I} \tilde{H} \tilde{S} \tilde{C} \tilde{Y} \tilde{H} \tilde{S} \parallel$
- $\tilde{d} \rightarrow \tilde{E} \tilde{3} \tilde{E}^{\circ} \uparrow \frac{1}{2} \tilde{E} \tilde{E} \parallel \tilde{I} \tilde{E} \frac{1}{2} \circ \tilde{S} \frac{1}{2} \tilde{C} \tilde{S} \tilde{n} \tilde{E} \tilde{Y} \tilde{S} \tilde{I}$

# 3D-LLM: Data Generation



- Generate questions with ChatGPT, answers with **BLIP -2**, aggregate questions and rewrite them
- Make more diverse dataset from **ScanQA**

## 3D-language data generation pipelines

# 3D-LLM: Results

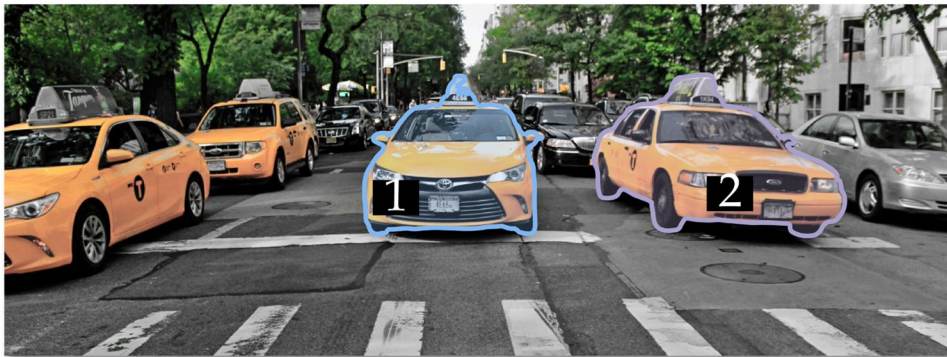
	B-1	B-2	B-3	B-4	METEOR	ROUHE-L	CIDER	EM
VoteNet+MCAN*	28.0	16.7	10.8	6.2	11.4	29.8	54.7	17.3
ScanRefer+MCAN*	26.9	16.6	11.6	7.9	11.5	30	55.4	18.6
ScanQA*	30.2	20.4	15.1	10.1	13.1	33.3	64.9	<b>21.0</b>
LLaVA(zero-shot)	7.1	2.6	0.9	0.3	10.5	12.3	5.7	0.0
flamingo-SingleImage	23.8	14.5	9.2	8.5	10.7	29.6	52	16.9
flamingo-MultiView	25.6	15.2	9.2	8.4	11.3	31.1	55	18.8
BLIP2-flant5-SingleImage	28.6	15.1	9.0	5.1	10.6	25.8	42.6	13.3
BLIP2-flant5-MultiView	29.7	16.2	9.8	5.9	11.3	26.6	45.7	13.6
3D-LLM (flamingo)	30.3	17.8	12.0	7.2	12.2	32.3	59.2	20.4
3D-LLM (BLIP2-opt)	35.9	22.5	16.0	9.4	13.8	34.0	63.8	19.3
3D-LLM (BLIP2-flant5)	<b>39.3</b>	<b>25.2</b>	<b>18.4</b>	<b>12.0</b>	<b>14.5</b>	<b>35.7</b>	<b>69.4</b>	20.5

results on **ScanQA** validation set

- Then, **2D models** in zero-shot gave poor results for 3D tasks
- **LLM gives a boost** in the 3D QA task
- 3D features are important

# VLLMs: SpatialRGPT (nvidia)

Go beyond **ScanNet -style datasets** (which are mostly indoor 3D scans)



If I am riding a motorcycle with 36" wide, do you think I can pass through the area between 1 and 2?



The distance between 1 and 2 is 3.67 feet, so yes, you can pass through them.



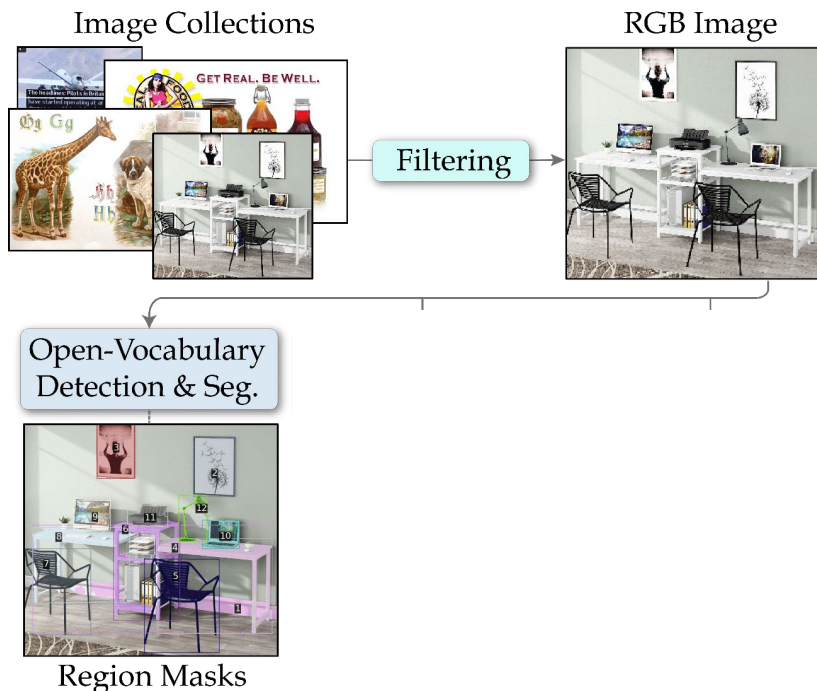
**SpatialRGPT** is a Vision-Language Model (VLM) focused on **spatial reasoning**

SpatialRGPT works **purely on 2D images** but learns to understand 3D spatial relationships



# VLLMs: Spatial IRGPT

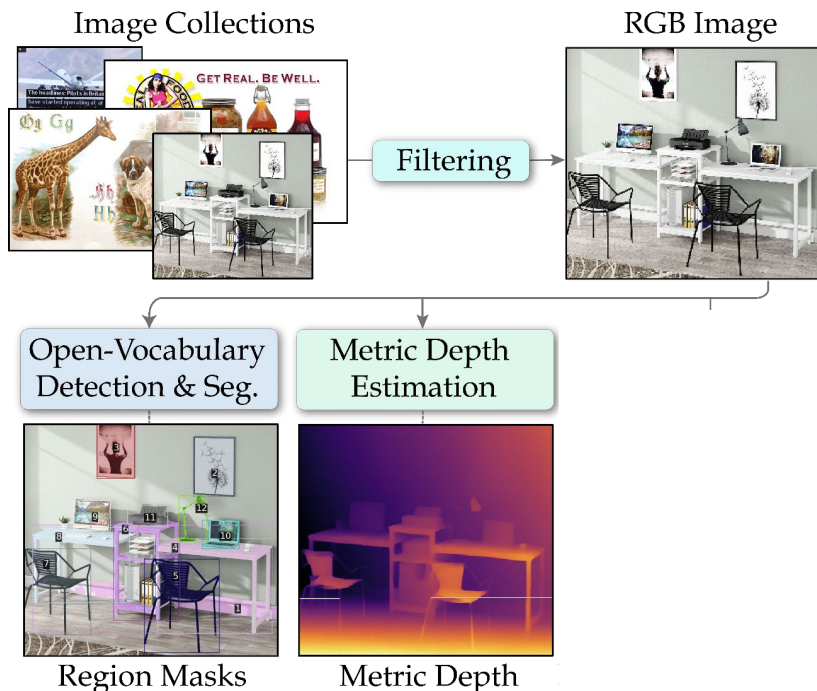
automatic data curation pipeline



- **GroundingDINO** model to detect object bounding boxes

# VLLMs: Spatial IRGPT

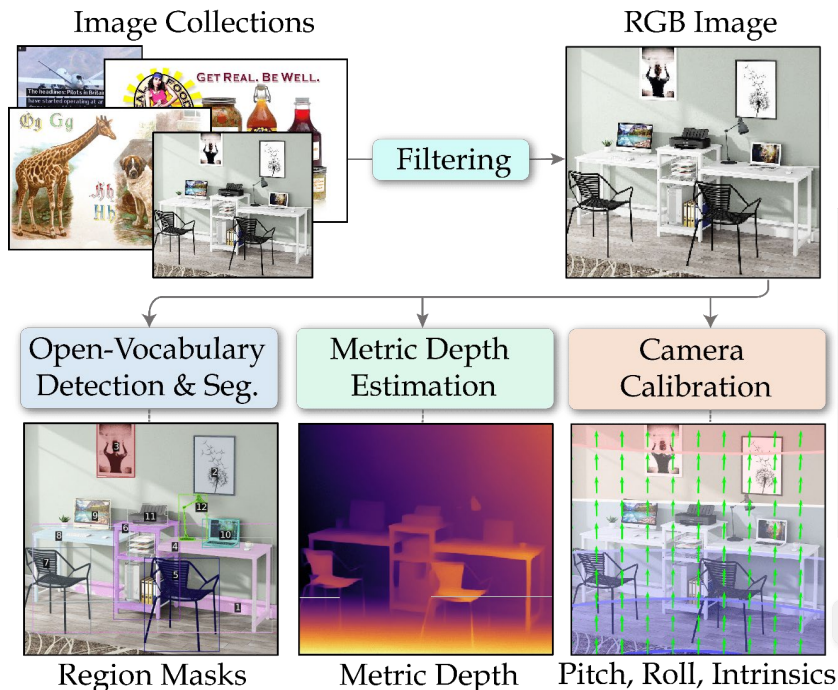
automatic data curation pipeline



- **GroundingDINO** model to detect object bounding boxes
- Estimate metric depth of every pixel using **Metric3Dv2** (focal length)

# VLLMs: Spatial IRGPT

automatic data curation pipeline

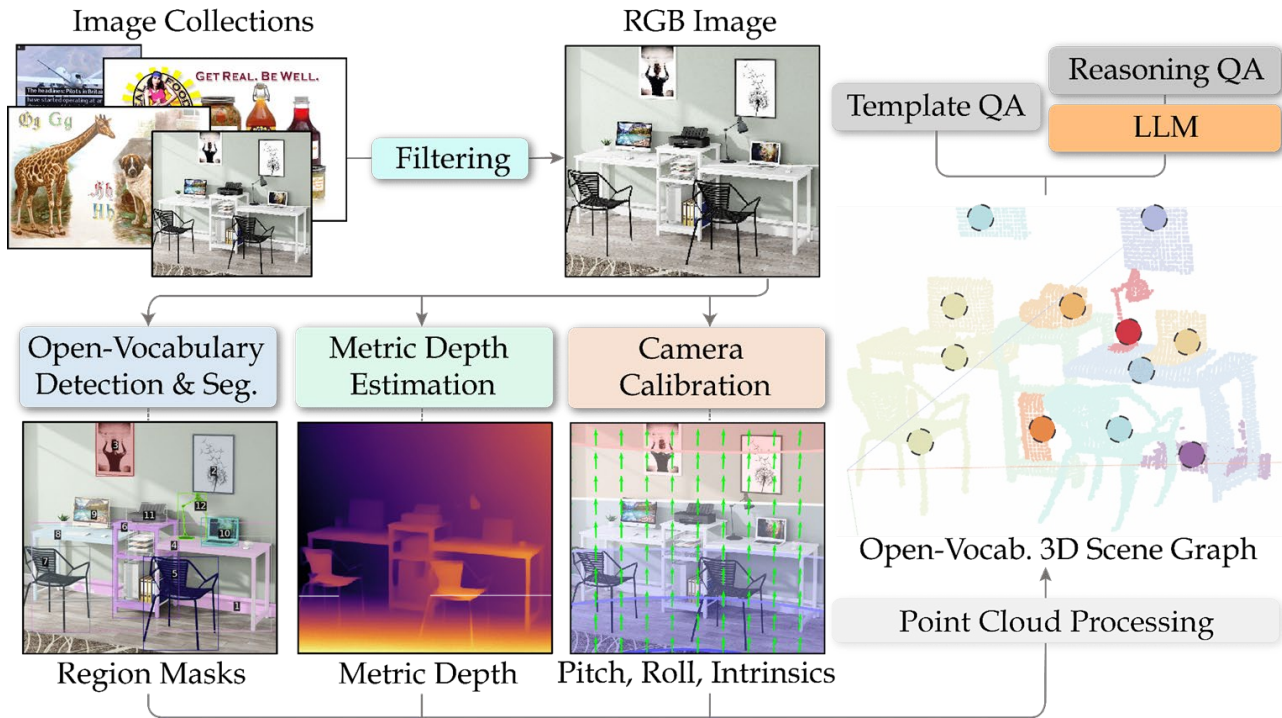


- Calibrate camera intrinsics with **WildCamera** and **PerspectiveFields**

$$K = \begin{pmatrix} f_x & s & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{pmatrix}$$

# Automatic data curation pipeline

automatic data curation pipeline



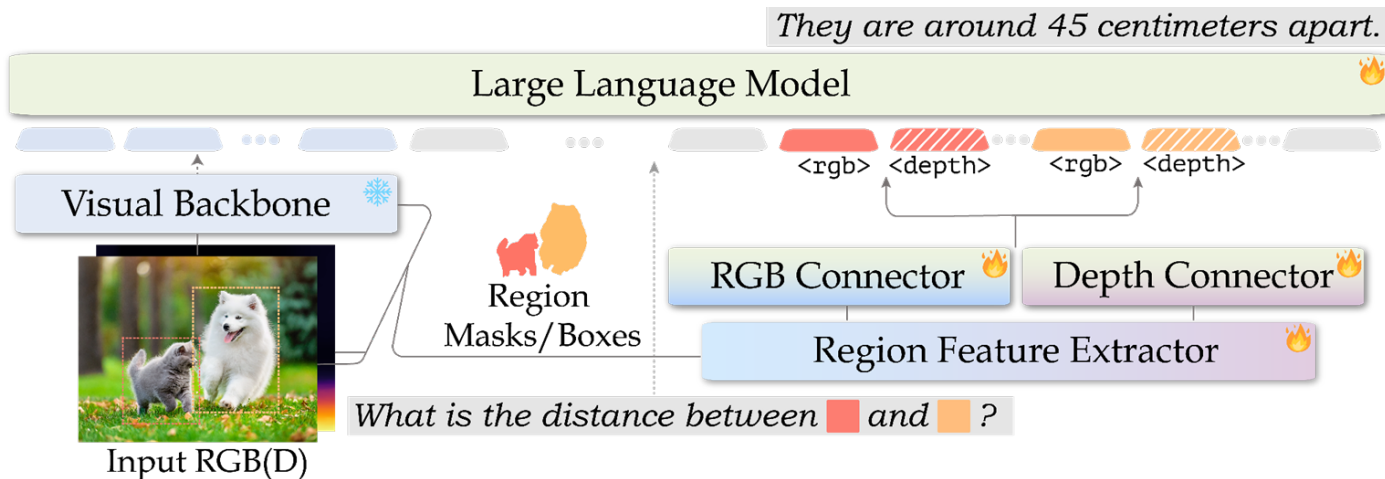
Project each object mask into **3D space** using depth + camera info

Derive a **scene graph**

Make **rich QA pairs**, ~8.7M QA concepts

# VLLMs: SpatialRGPT

**RoI -Align** for generated masks, ViT processes the raw **RGB image** into a grid of feature tokens, **depth map** passes through a small CNN



# VLLMs: SpatialRGPT



How far is 0 from 1 horizontally?

0 and 1 are 31.21 feet apart horizontally.



Could you tell me the vertical size of 0?

0 is 3.67 feet tall.



Does 0 have a larger size compared to 1?

Correct, 0 is larger in size than 1.



Does 0 have a lesser width compared to 1?

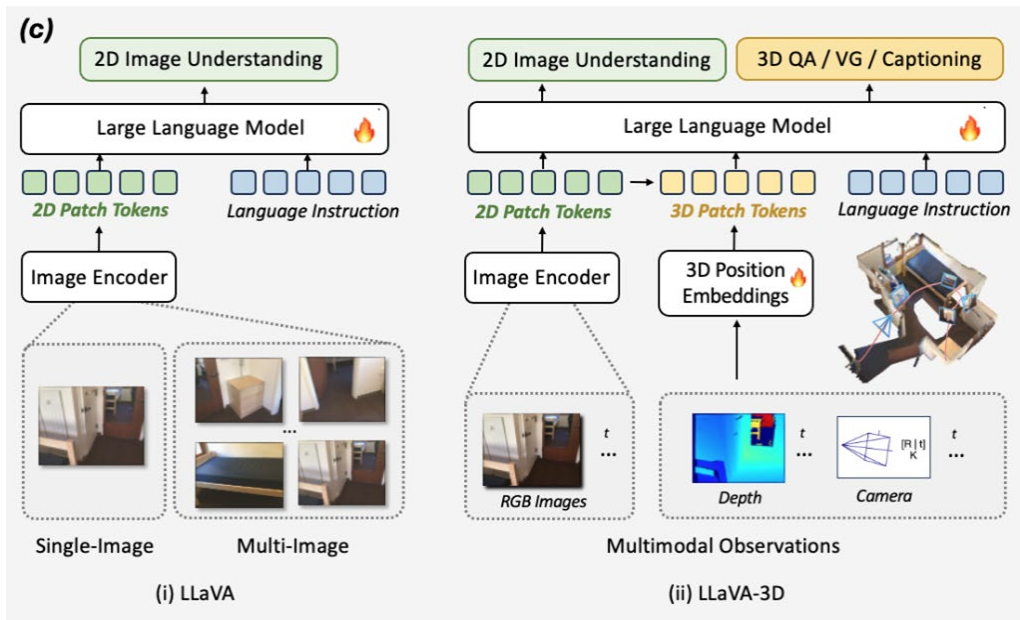
No, 0 is not thinner than 1.

First spatial dataset, for testing authors used part of their benchmark

Method	Below/ Above	Left/ Right	Big/ Small	Tall/ Short	Wide/ Thin	Behind/ Front	Avg.
GPT-4	64.1	42.8	42.8	61.6	61.6	49.0	57.8
GPT-4V	63.3	46.6	64.1	60.7	68.2	45.4	58.1
LLaVA-v1.6-34B	44.1	45.7	36.7	53.5	37.5	45.4	43.9
GPT-4V+SoM	75.0	55.2	42.4	54.4	49.0	47.2	54.3
LLaVA-v1.6-34B+SoM	44.1	40.0	33.9	47.3	41.3	46.3	42.3
Kosmos-2	28.3	15.2	4.71	26.7	12.5	12.7	17.0
RegionVILA	30.8	47.6	35.8	44.6	35.5	49.0	40.4
SpatialRGPT	99.1	99.0	79.2	89.2	83.6	87.2	89.8
SpatialRGPT-Depth	99.1	99.0	80.1	91.9	87.5	91.8	91.7

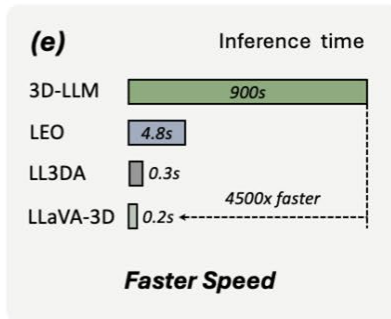
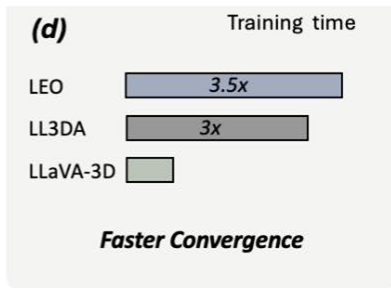
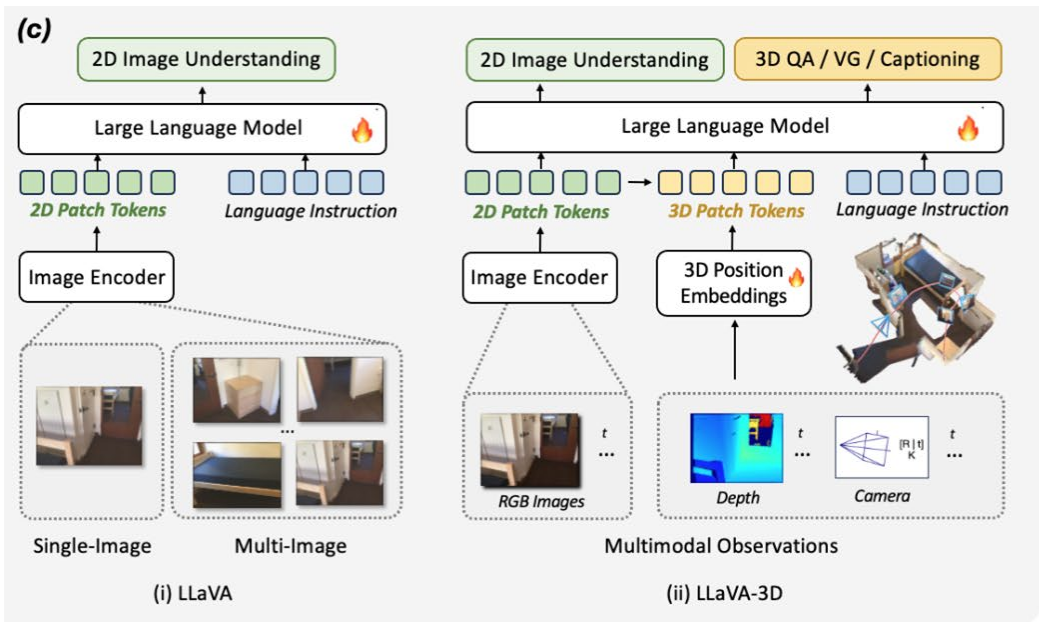
SpatialRGPT-Bench Qualitative Results: Numbers represent success rates in percentage (↑).

# VLLMs: LLaVA-3D



- Preserves the ability to understand **2D images**
- Expand the model with the part connected with the **3D understanding**
- Give the images from **different angles**

# VLLMs: LLaVA-3D



$\hat{e} \hat{J} \hat{L} \hat{\uparrow} \hat{E} \hat{M} \hat{E} \hat{I} \hat{S} \hat{C} \hat{E}$   
 $\hat{I} \hat{\uparrow} \hat{1} \hat{2} \hat{=} \hat{E} \hat{A} \hat{B} \hat{3} \hat{U} \hat{U} \hat{e} \hat{I} \hat{E}$   
 $\hat{S} \hat{E} \hat{E} \hat{E} \hat{S} \hat{1} \hat{J} \hat{C} \hat{S} \hat{E} \hat{I} \hat{E}$   
 $\hat{S} \hat{w} \hat{S} \hat{=} \hat{E} \hat{=} \hat{S} \hat{E} \hat{A} \hat{B} \hat{E}$   
 $\hat{C} \hat{S} \hat{L} \hat{=} \hat{E} \hat{I} \hat{C} \hat{L} \hat{I} \hat{4}$   
 $\hat{=}$

# VLLMs: LLaVA-3D

Generated their **own dataset** (combination of 2D and 3D) – 1M instructions, **~sota**

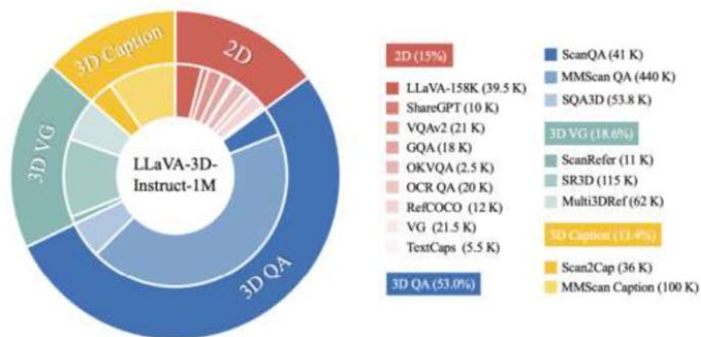


Figure 3. **LLaVA-3D-Instruct-1M**. The hybrid 2D and 3D Dataset Collection. Left: Distribution of data across categories, with the outer circle representing all categories and the inner circle illustrating data subset distribution. Right: Detailed dataset quantities.

	ScanQA (val)					SQA3D (test)
	C	B-4	M	R	EM@1	EM@1
<b>Task-specific models</b>						
Scan2Cap [11]	-	-	-	-	-	41.0 <sup>†</sup>
ScanRefer+MCAN [49]	55.4	7.9	11.5	30.0	18.6	-
ClipBERT [25]	-	-	-	-	-	43.3
ScanQA [3]	64.9	10.1	13.1	33.3	21.1	47.2
3D-VisTA [53]	69.6	10.4	13.9	35.7	22.4	48.5
<b>Task-specific fine-tuned 3D LMMs</b>						
3D-LLM (FlanT5) [16]	69.4	12.0	14.5	35.7	20.5	-
LL3DA [35]	76.8	13.5	15.9	37.3	-	-
Chat-3D v2 [17]	87.6	14.0	-	-	-	54.7
LEO [18]	<b>101.4</b>	13.2	20.0	49.2	24.5 (47.6)	50.0 (52.4)
Scene-LLM [15]	80	12.0	16.6	40.0	<b>27.2</b>	54.2
<b>Zero-shot 2D LMMs</b>						
VideoChat2 [30]	49.2	9.6	9.5	28.2	19.2	37.3
LLaVA-NeXT-Video [26]	46.2	9.8	9.1	27.8	18.7	34.2
GPT-4V	59.6	-	13.5	33.4	-	-
Gemini	68.3	-	11.3	35.4	-	-
Claude	57.7	-	10.0	29.3	-	-
<b>LLaVA-3D</b>	<b>91.7</b>	<b>14.5</b>	<b>20.7</b>	<b>50.1</b>	27.0 (45.0)	<b>55.6</b> (57.6)

4

---

; B3-|| ■ ň Š Ě

3D-LLM, SpatiaIRGPT



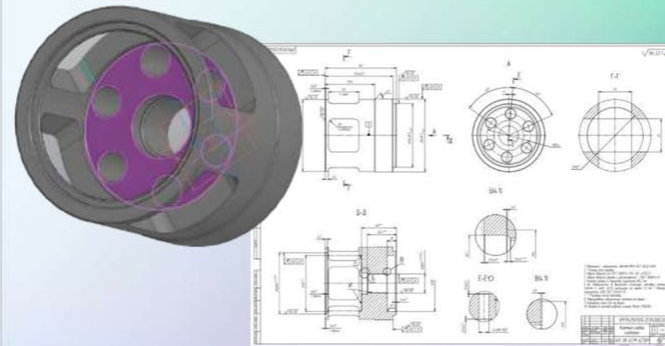
# CAD-models: Reverse Engineering

Build a full-fledged **constructive model (CAD)** of the selected part in the engineering software, in order to produce an exact copy of the part

Реальный объект и его 3D скан



Цифровая модель и чертёж



# CAD-models: Reverse Engineering

**Mesh approximation** (piecewise linear approximation of a cylinder) →

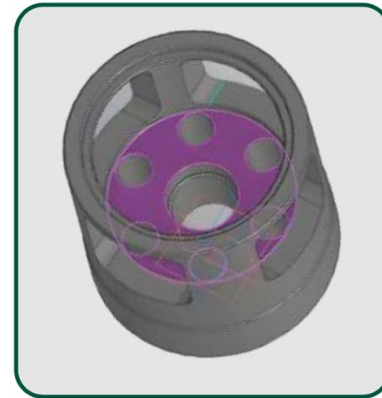
It can be done using **splines** — boundary-representation (exact mathematical model)



photo



Mesh

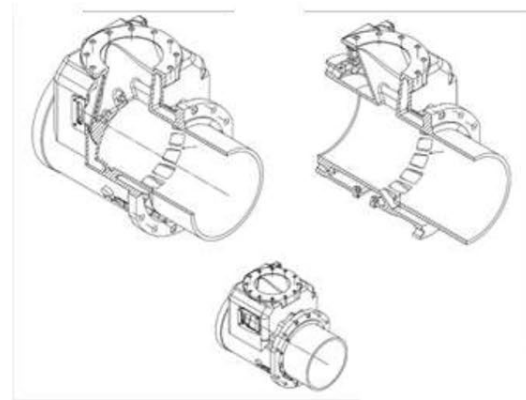


B-rep

# CAD-models: Reverse Engineering

**Mesh approximation** (piecewise linear approximation of a cylinder) →

It can be done using **splines** — boundary-representation (exact mathematical model)



real device with multiple details

; B3|| ■ ņ Š ĽĚĚ Š ū Š Ć Ě Š Ĭ Ľ Ľ Š Š Ć Ľ Ľ Ě



real device with multiple details

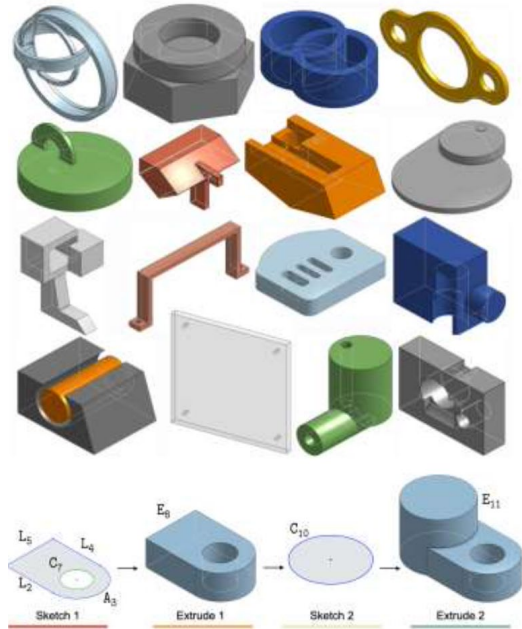
1. device is sequentially **disassembled** into parts, **record** the process
2. **calibration** marks + **3D** scanner
3. **manual measurements** of individual elements
4. scan (**mesh**) of the model is loaded into the software
5. **engineer manually designs** the part so that the result matches the scan
6. **manager validates** the result

# CAD-models: Reverse Engineering

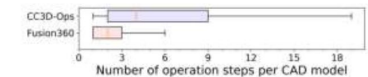
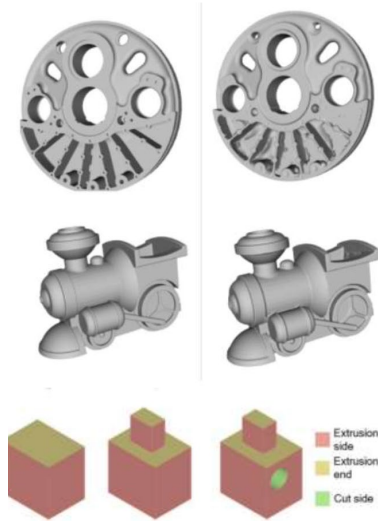
The part can be significantly **damaged**



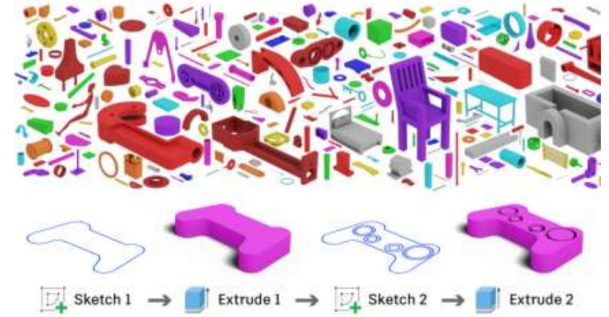
# CAD-models: Benchmarks



BŠŠ° ; B



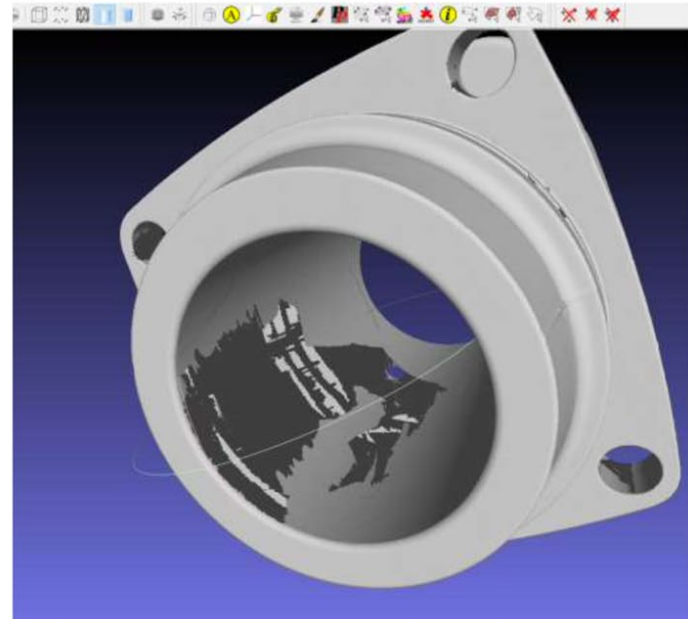
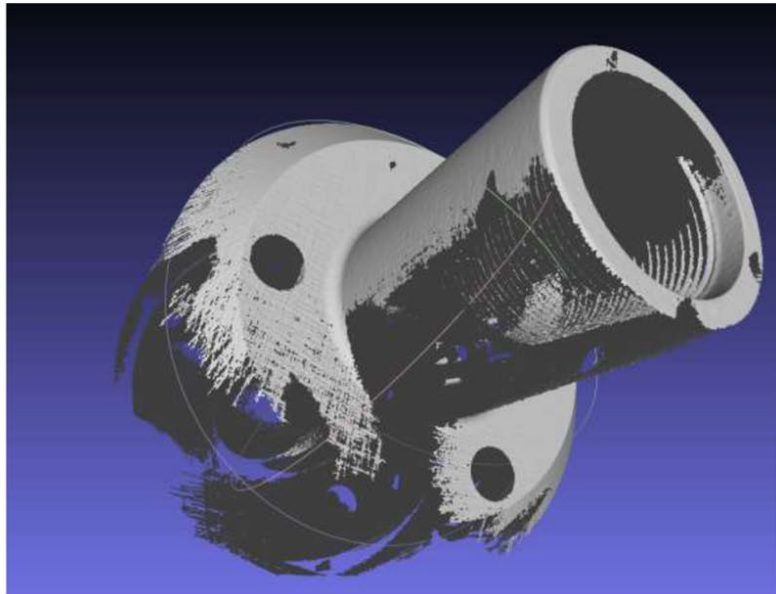
cvi2



AutodeskAILab

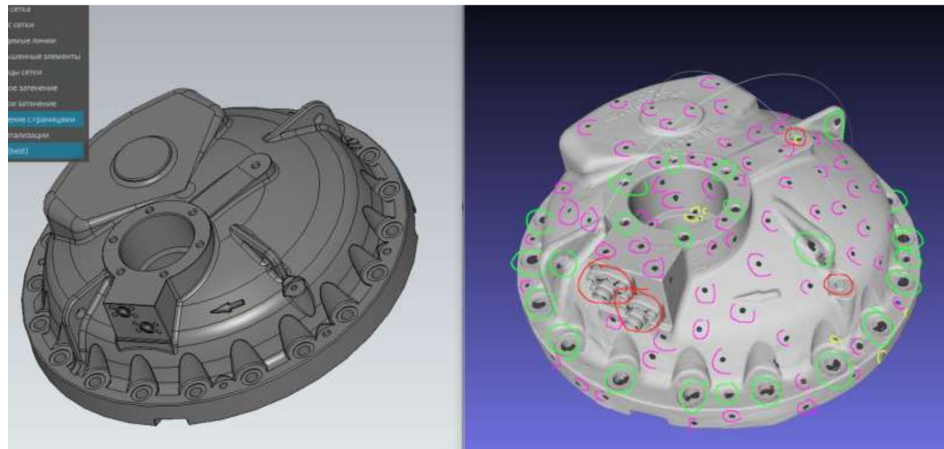
# CAD-models: Benchmarks

There are **no surface pieces in the scan** due to the limited scanning angles in practice



# CAD-models: Benchmarks

1. **random holes** in the mesh  
(need to be repaired)
2. round holes in the mesh due  
to **marks** (need to be  
repaired)
3. **round holes in the part** itself  
that look almost like round  
holes in the mesh (can not be  
repaired)
4. part of the mesh belongs to  
**another part**



Many more artifacts

# CAD-models: Benchmarks

## Этап 1: Формообразующие операции

1. Выдавливание (Extrude)
2. Вращение (Revolve)
3. Элемент по траектории (Sweep)
4. Элемент по сечениям (Loft)
5. Булевы операции (Boolean)
6. Обечайка (Shell)
7. “Операция без истории”
8. Отдельные детали по ГОСТу

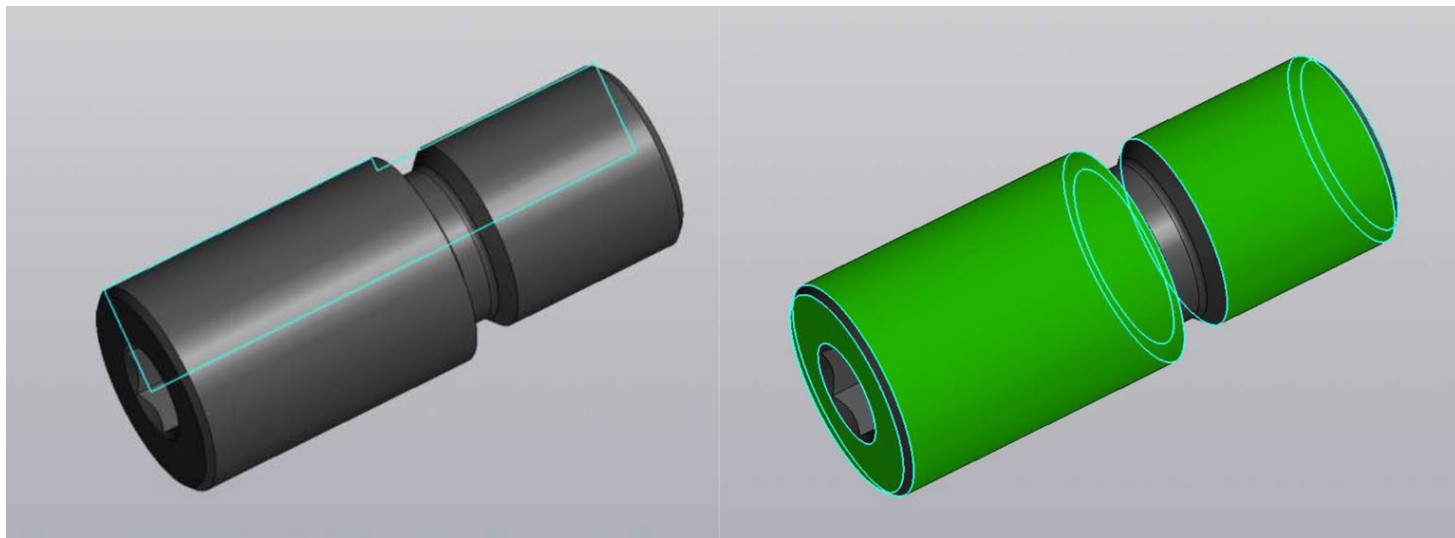
- очень часто
- довольно часто
- редко
- один раз

## Этап 2: Дорабатывающие операции

1. Отверстие (Hole)
2. Фаски (Chamfer)
3. Скругления (Fillet)
4. Массивы (Array, Pattern)
5. Надписи / гравировки (Text)
6. Резьба
7. Прочие детали оформления: “отверстие ГОСТ 14034-74” (Model 30.7), условные обозначения и т. п.

# CAD-models: Benchmarks

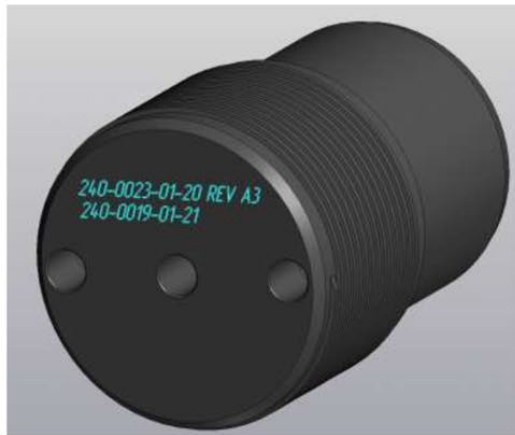
**Вращение** как формообразующая операция



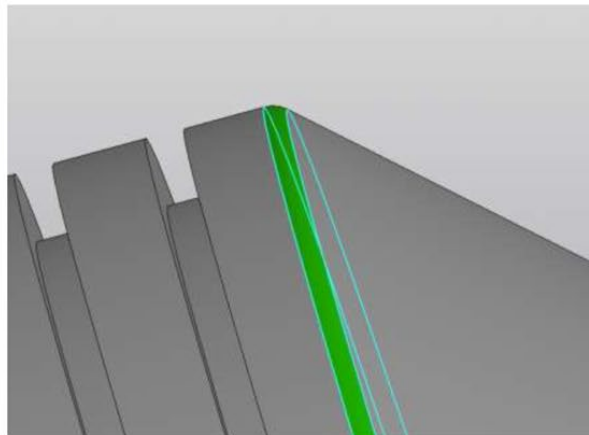
Винт сердечника

# CAD-models: Benchmarks

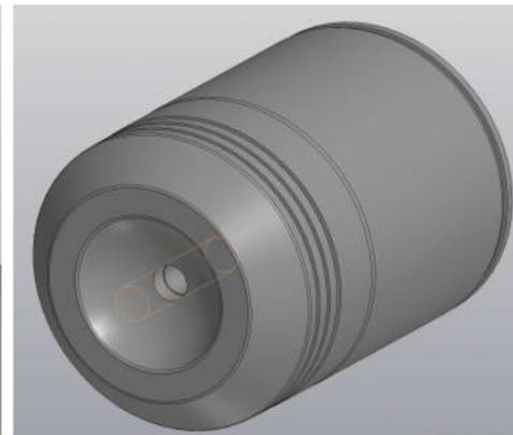
Примеры дорабатывающих операций



Надпись



Скругление

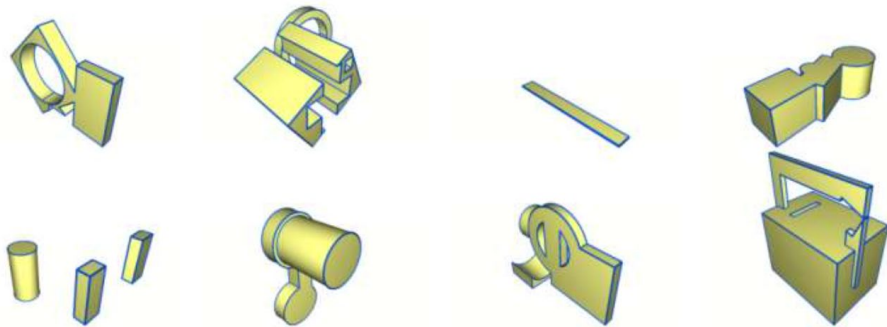


Отверстие



# CAD-models: Datasets

Ĉ½řň || ΔŮň



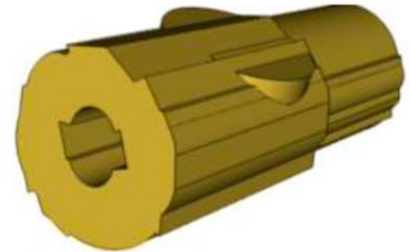
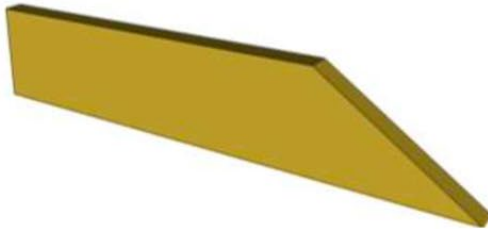
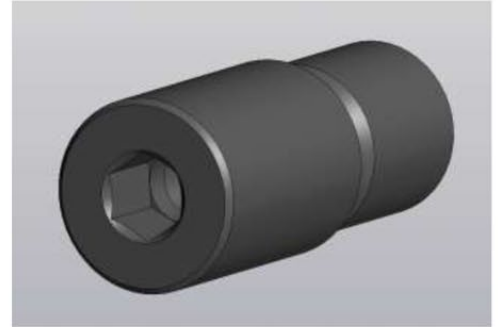
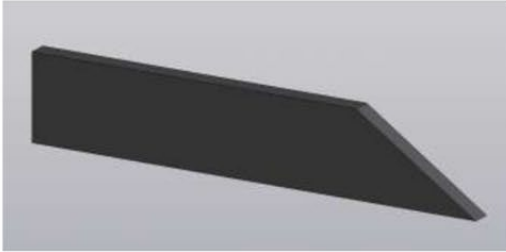
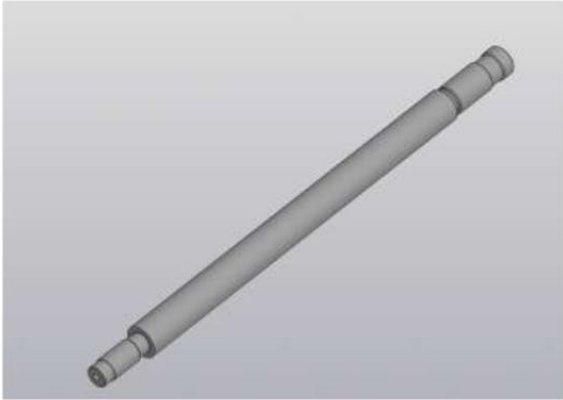
- Randomized generation of examples based on operations (profile, extrude)
- The examples are not very meaningful

LLM-based

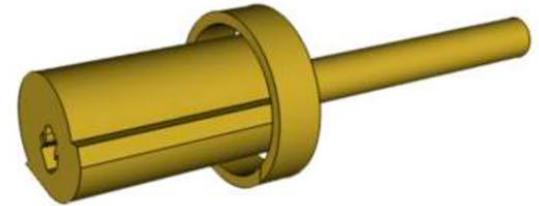
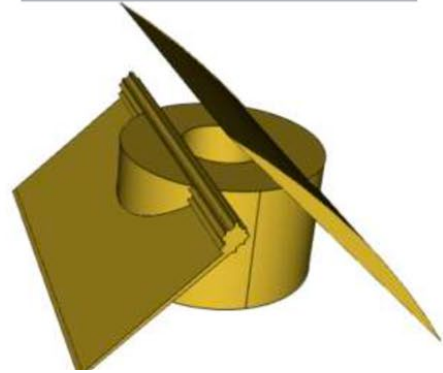
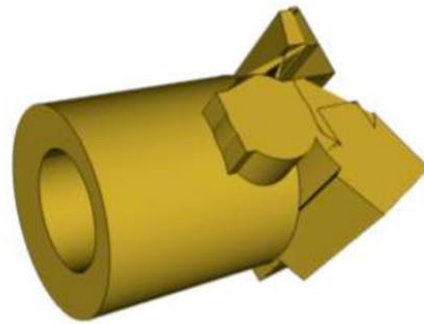
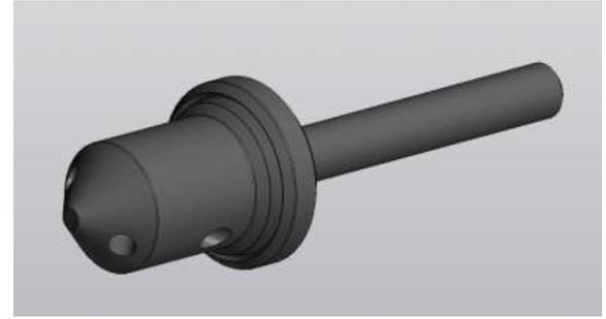
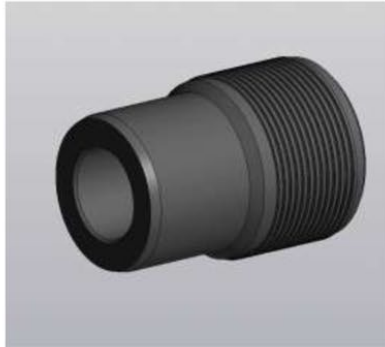


- Identify 10-20 common types of parts
- Use LLM to get the parameterization of the parts based on the description of the part class

# CAD-models: MLLMs



; B3|| ■ ň Š Ľ ě ě ě ě ě ě ě ě



# Conclusions

1 3D vision = semantic + metric visions

2  $c \mid \hat{C} \hat{E} \frac{1}{2} \uparrow \hat{S} \hat{C} \hat{A} \hat{B} \hat{H} \frac{1}{2} \frac{1}{2} \hat{U} \hat{S} \hat{L} \frac{1}{2} \hat{H} \hat{E} \hat{S} \hat{E}$   
 $\hat{n} \hat{S} \circ \hat{I} \uparrow \hat{H} \frac{1}{2} \parallel \hat{S} \hat{C} \hat{A} \hat{E} \hat{I} \hat{U} \hat{B} \approx \hat{E} \hat{I} \hat{E} \mid \hat{A} \hat{I} \hat{H} \mid \hat{J} \hat{n} \hat{E} \hat{E}$

3  $\approx \hat{S} \hat{L} \hat{S} \hat{I} \hat{E} \hat{n} \hat{u} \frac{1}{2} \hat{L} \hat{S} \parallel \hat{S} \hat{I} \mid ! \hat{E} \hat{A} \hat{B} \hat{E} \hat{U} \hat{E} \hat{I}$   
 $\hat{n} \frac{1}{2} \hat{I} \frac{1}{2} \hat{E} \hat{S} \hat{I} \hat{E} \hat{C} \hat{S} \circ \frac{1}{2} \hat{C} \hat{I} \hat{I} \mid \hat{H} \hat{E} \hat{A} \hat{B} \hat{E} \hat{S} \frac{1}{2} \hat{J} \hat{C} \hat{S} \hat{E}$   
 $\hat{S} \hat{I} \mid \hat{n} \hat{A} \hat{H} \hat{H} \mid \hat{H} \mid \hat{S} \hat{I} \hat{E} \hat{E}$



Когда прочитал весь материал на курсе, который смог, и студенты пошли сдавать бизнес-проекты