

# DepthLM

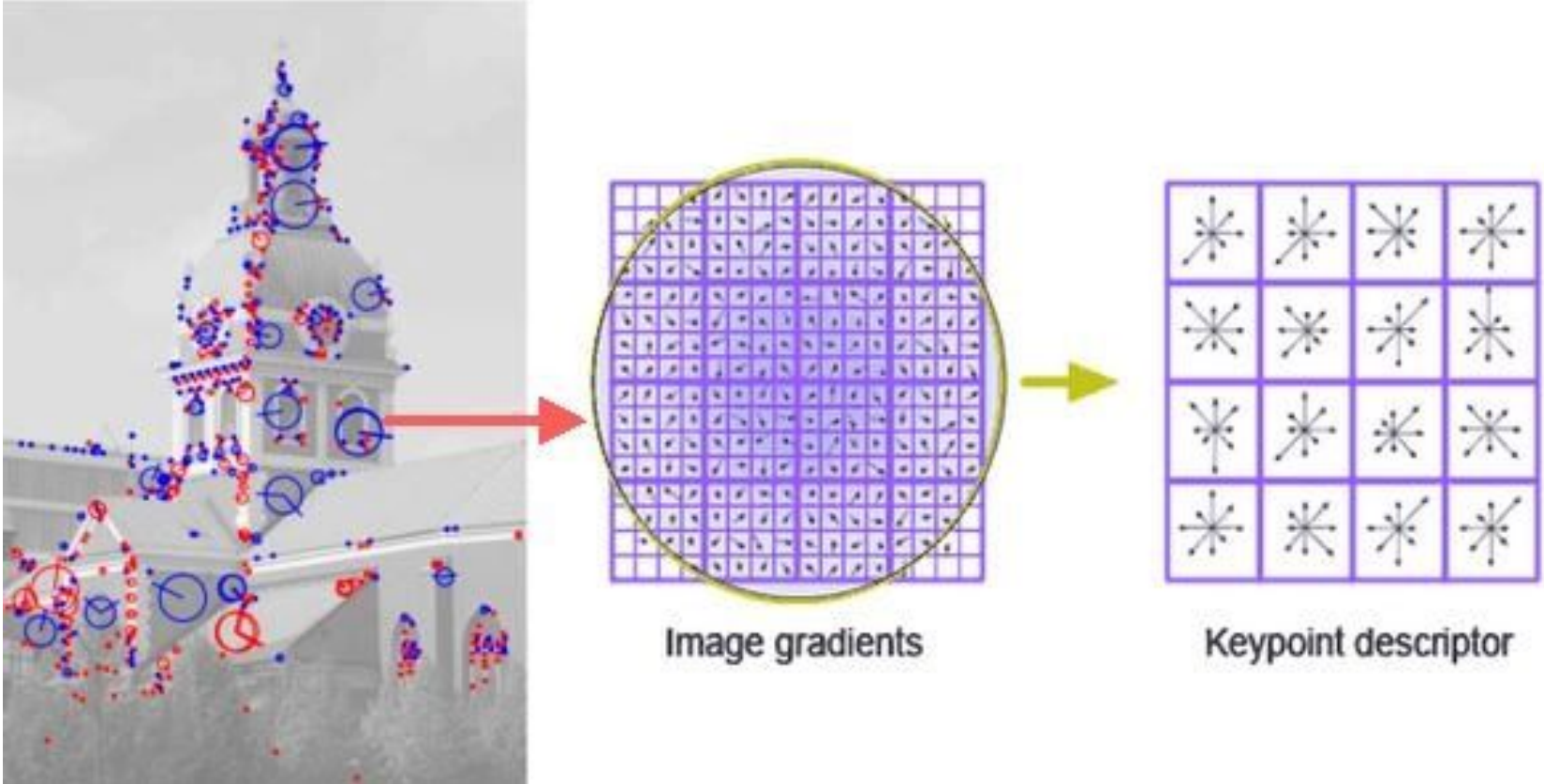
## Metric Depth From Vision Language Models

Zhipeng Cai, Ching-Feng Yeh, Hu Xu, Zhuang Liu, Gregory Meyer, Xinjie Lei,  
Changsheng Zhao, Shangwen Li, Vikas Chandra, Yangyang Shi

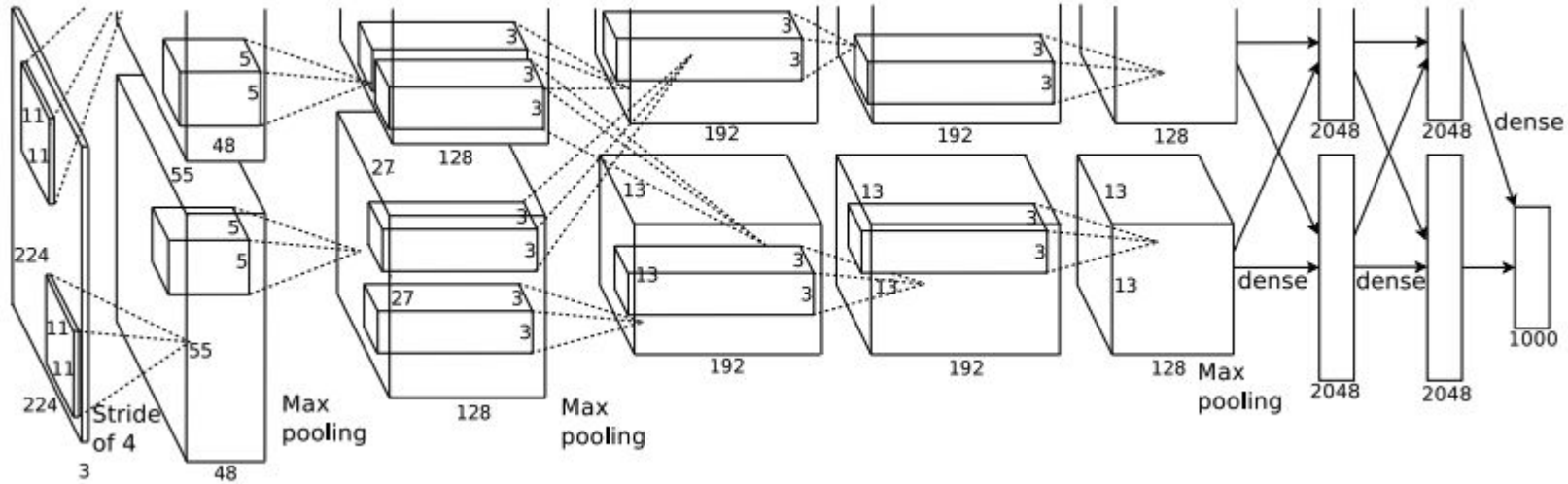


Code & Model

# History



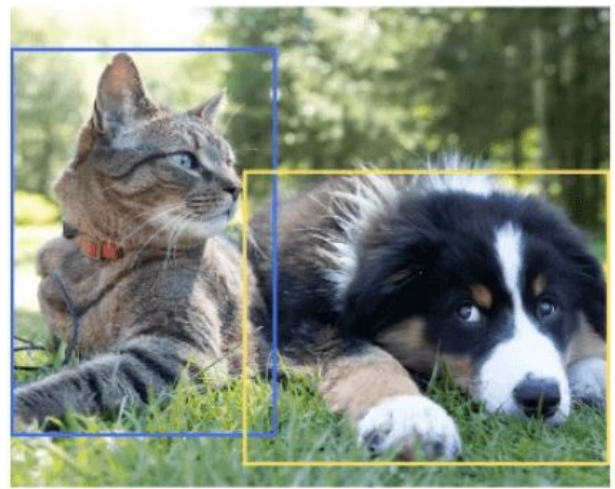
# History



classification



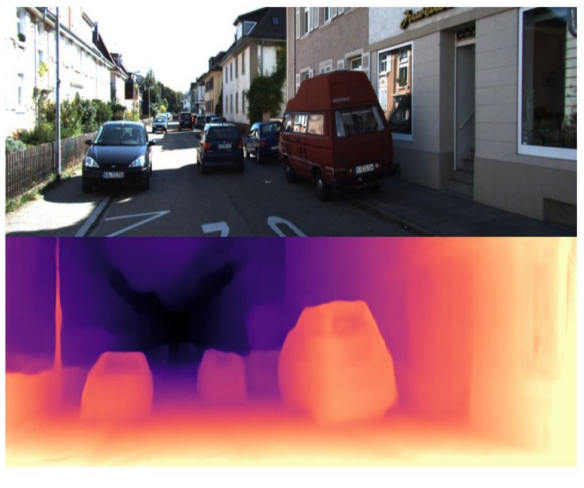
detection



segmentation

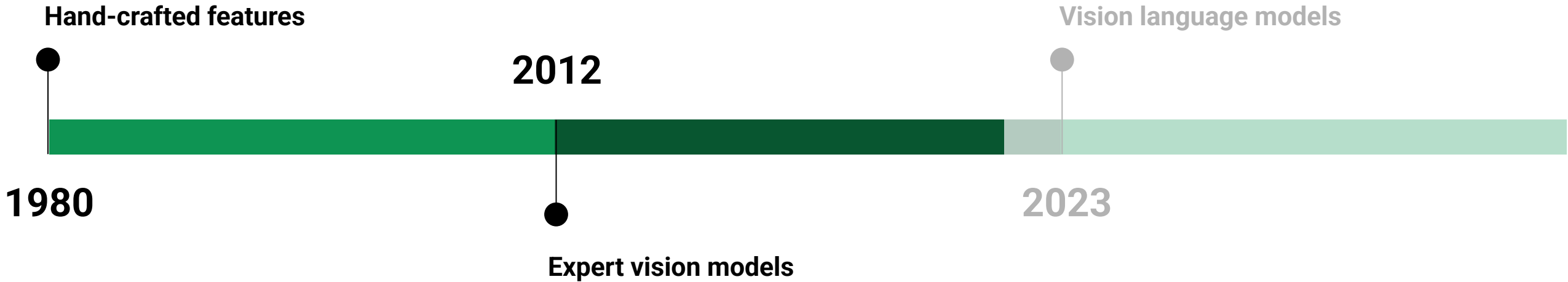


depth

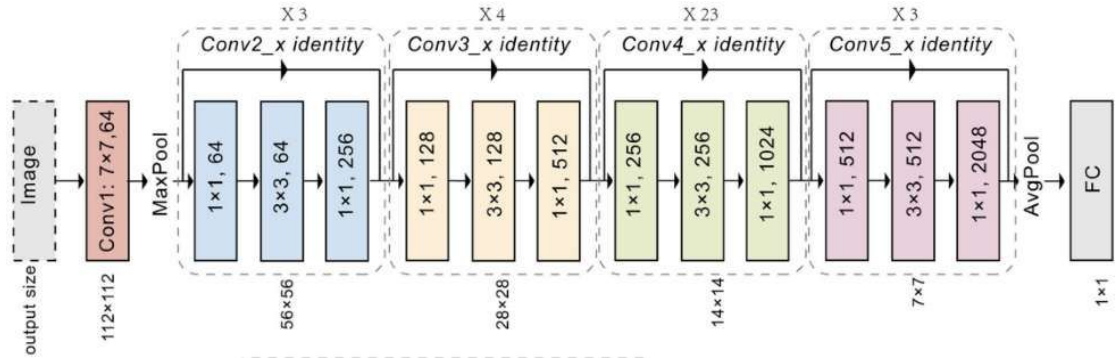


Krizhevsky et al., 2012

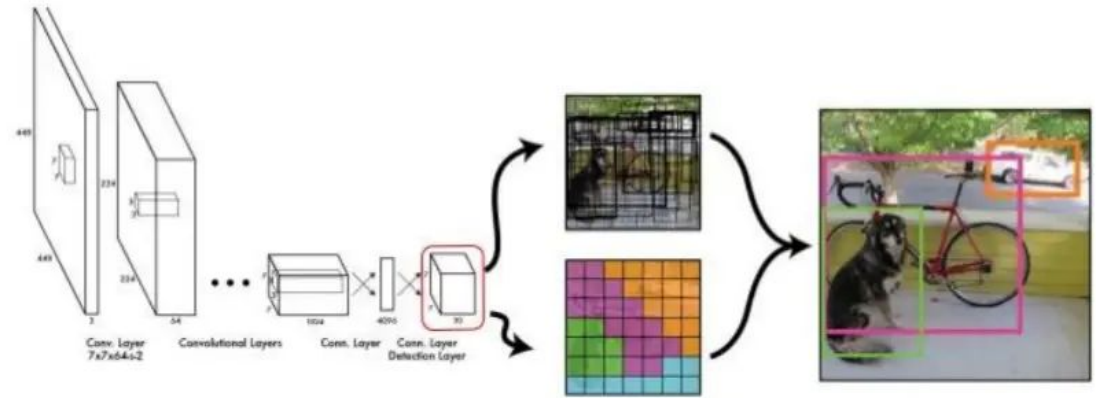
# History



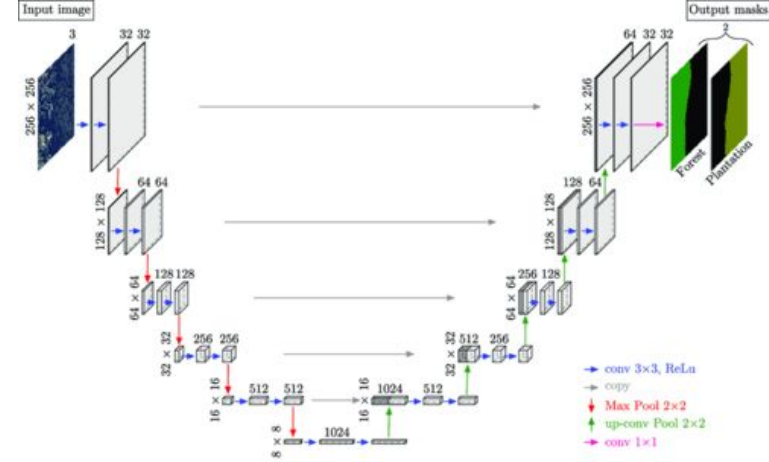
**Classification**  
(He et al. 2016)



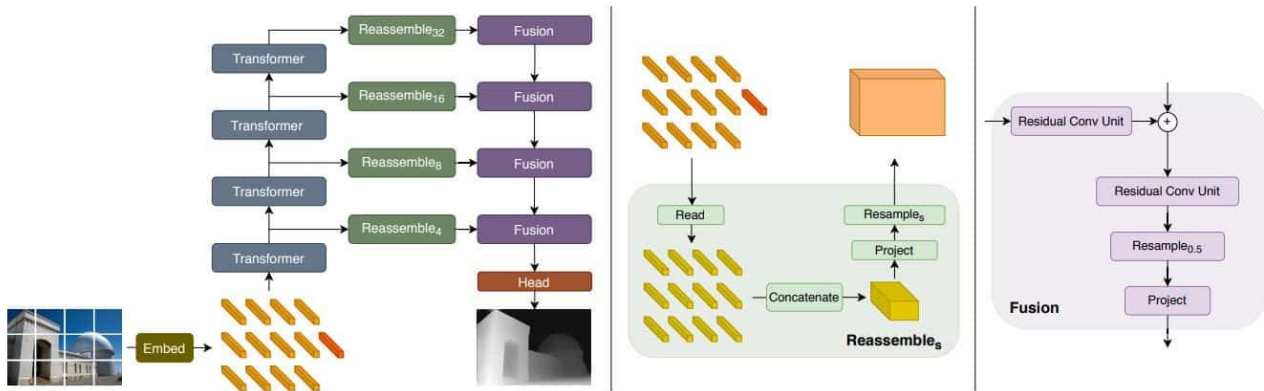
**Detection**  
(Redmon et al. 2016)



**Segmentation**  
(Ronneberger et al. 2015)



**Depth**  
(Ranftl et al. 2012)





# 3D understanding is still dominated by expert vision models

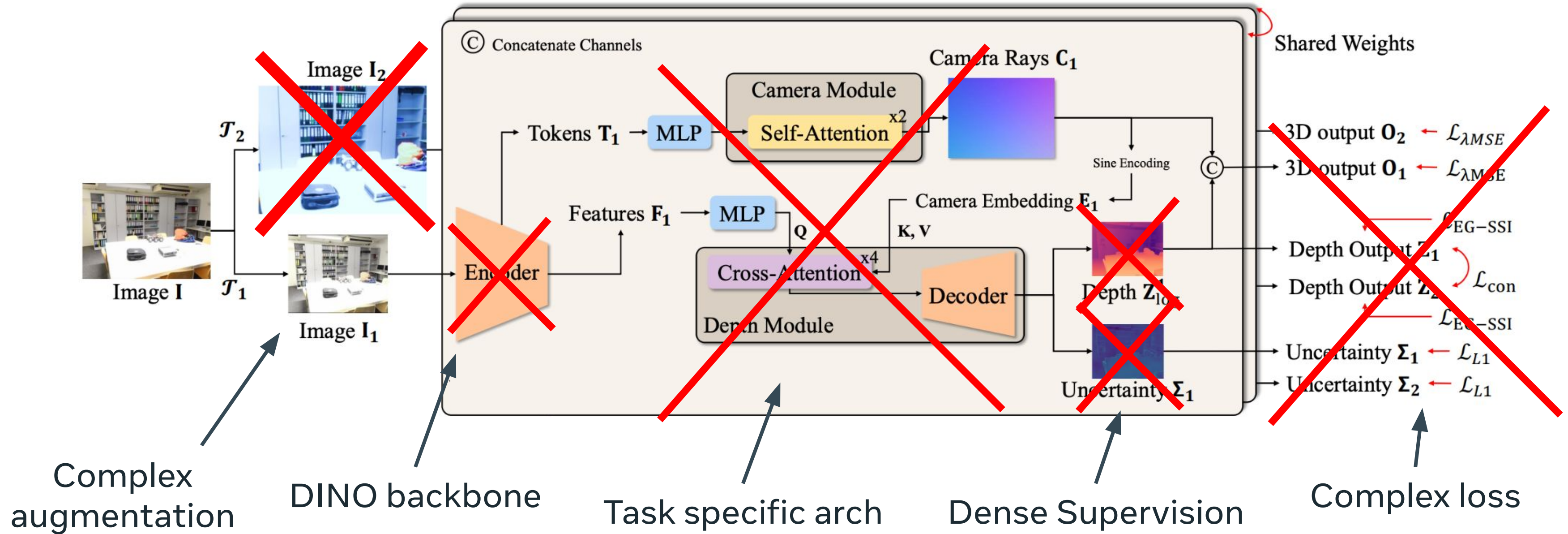
- Depth Estimation:
  - Metric3D, Unidepth, Moge, DepthAnything
- Pose Estimation:
  - VGGT, MapAnything
- Pixel Correspondence:
  - DKM, RoMa
- ...

## Key Message

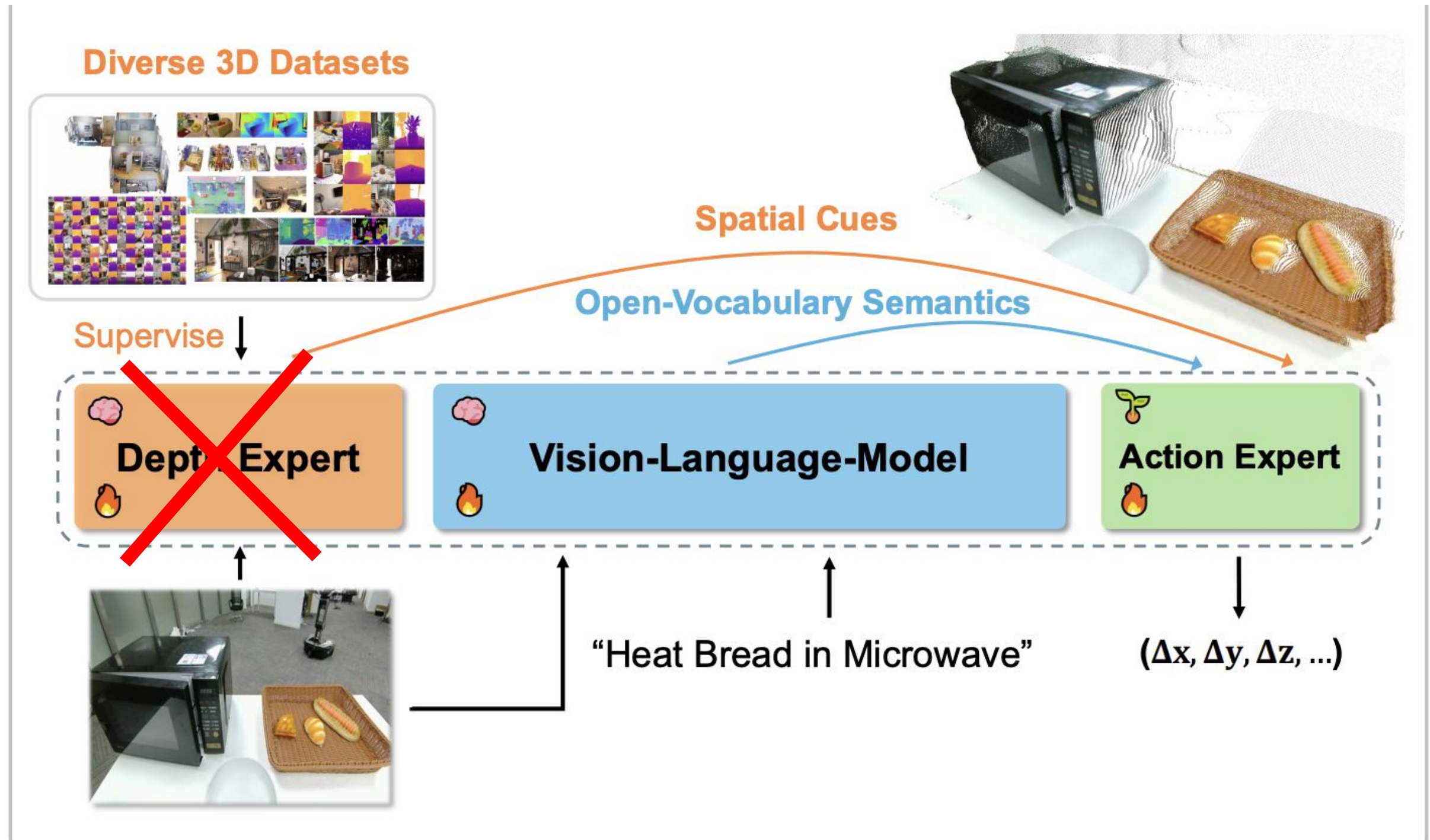
VLMs can learn 3D understanding as accurate as expert vision models.

- **No architecture change**
- **Standard text supervision**

# Why use VLMs? —> Simplicity



# Why use VLMs? —> Simplicity



**What Is Missing?**

# Can existing VLMs understand 3D?



Question: How far is this point from the camera?



- Ground Truth: 15.45m
- Seed1.5-VL: 5.38m
- GPT-5: 8.5m

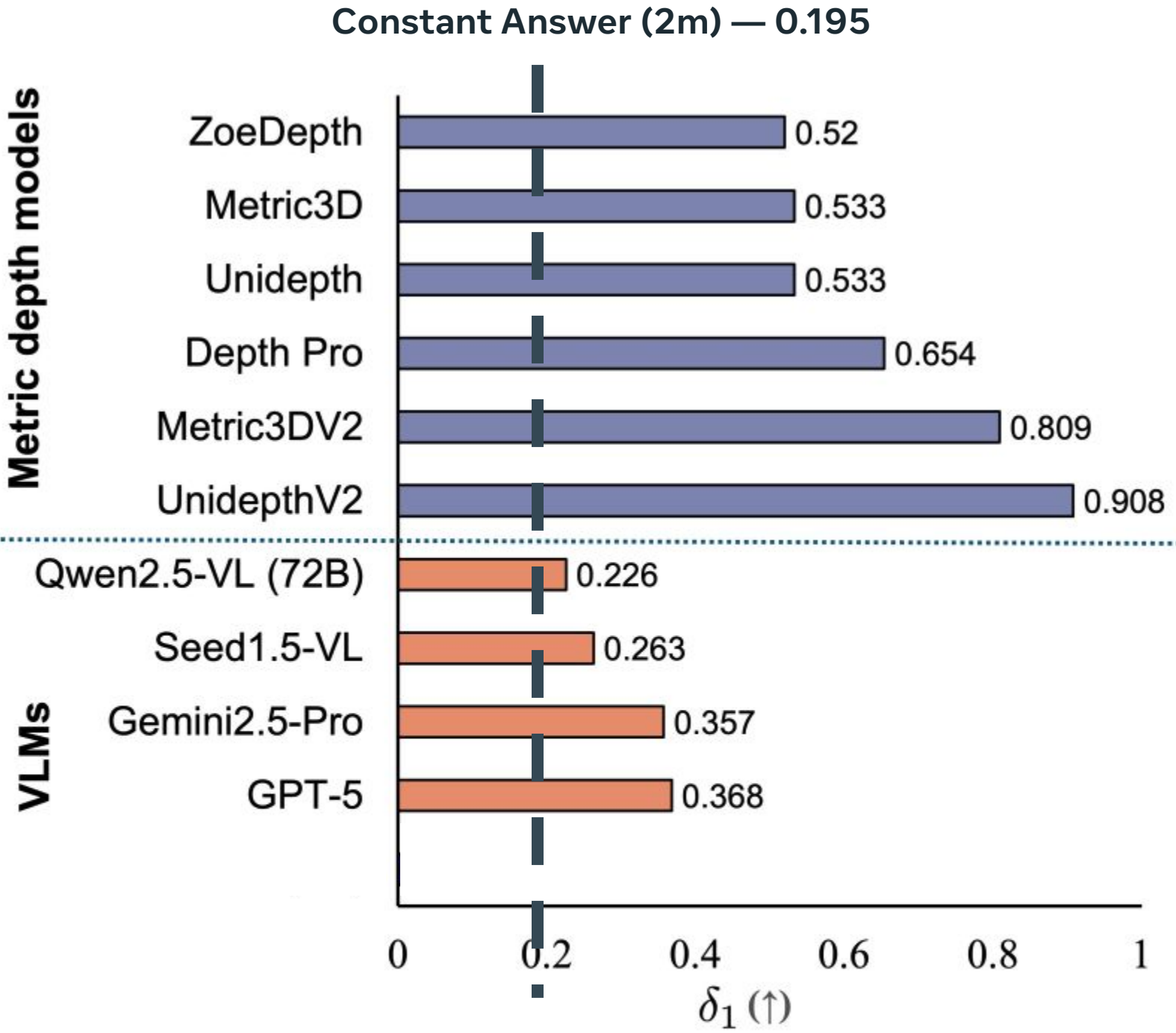
# Can existing VLMs understand 3D?



Question: How far is this point from the camera?



- Ground Truth: 15.45m
- Seed1.5-VL: 5.38m
- GPT-5: 8.5m



# DepthLM: standard VLM, expert model performance



Question: How far is this point from the camera?



- Ground Truth: 15.45m
- Seed1.5-VL: 5.38m
- GPT-5: 8.5m
- **Ours: 14.37m**

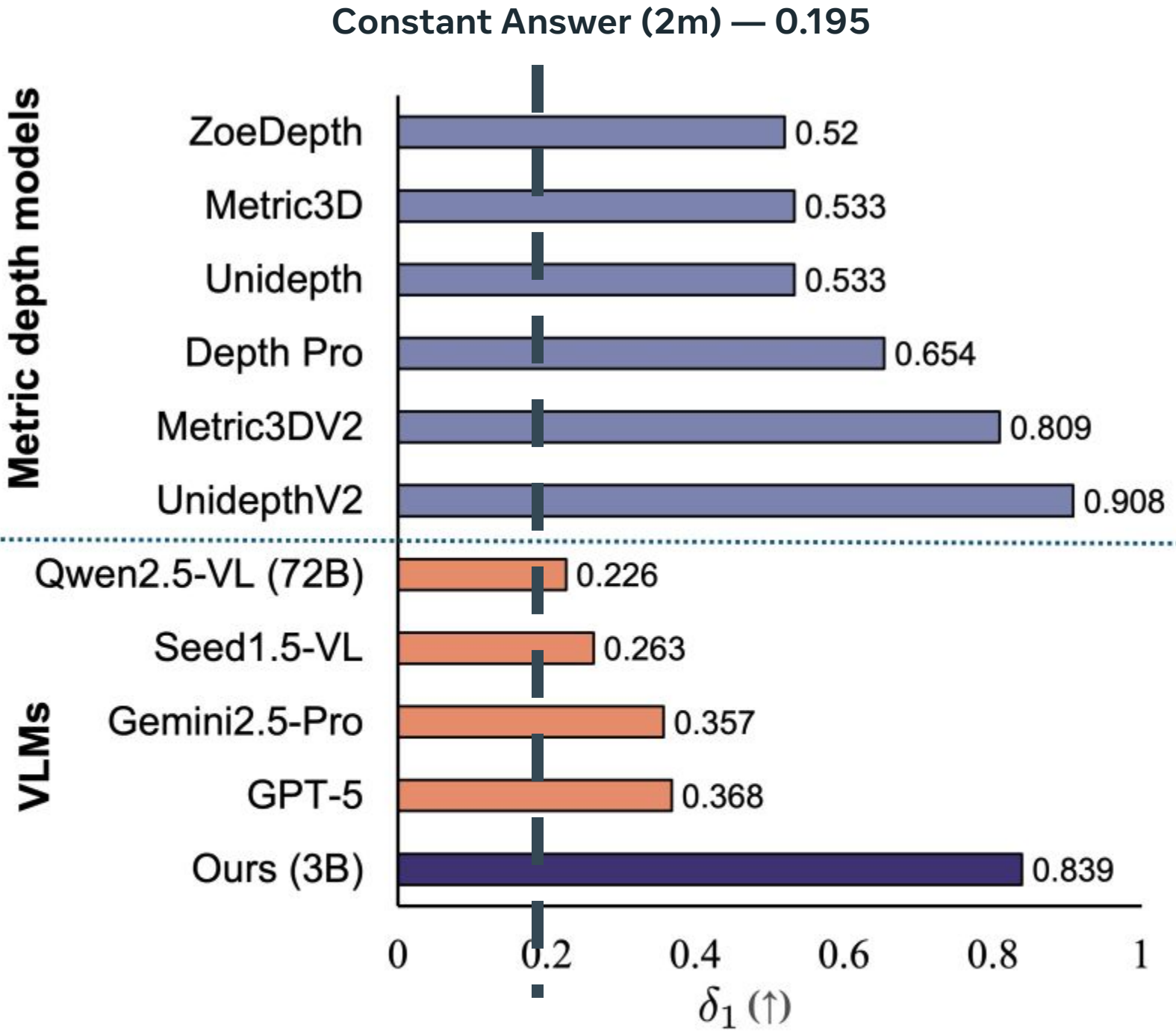


Query each pixel independently

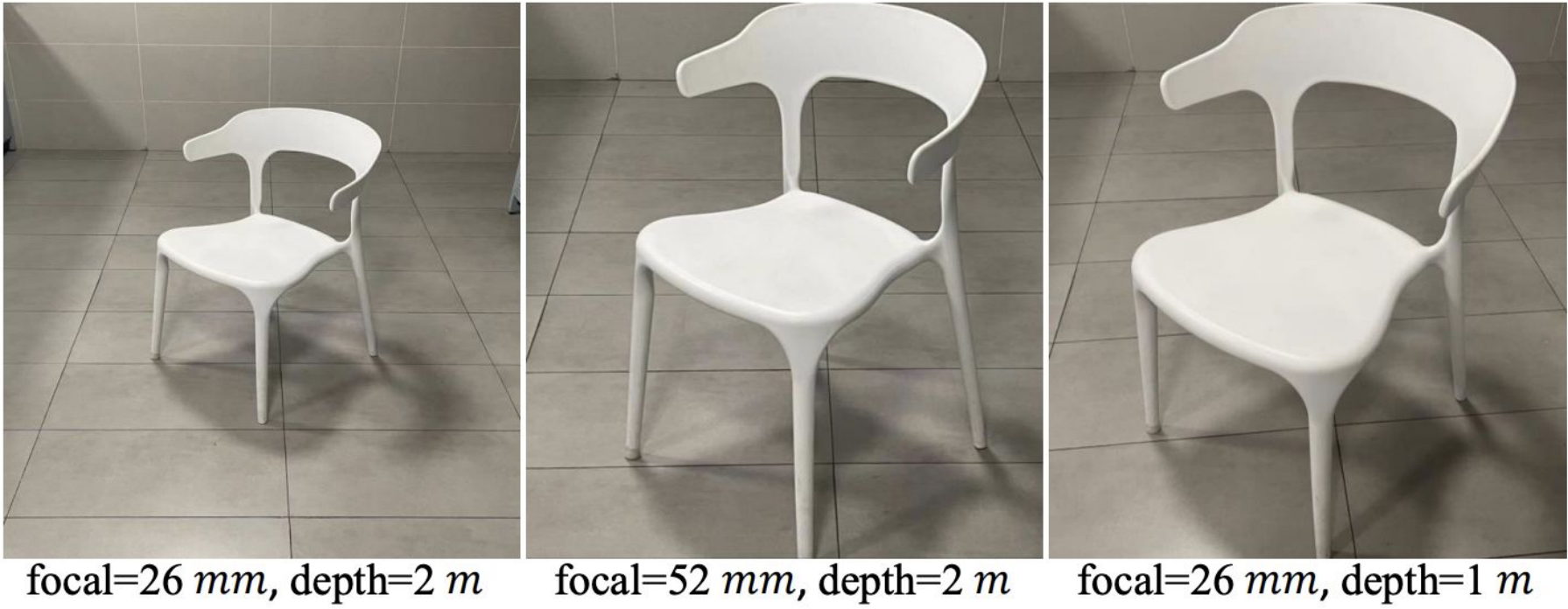
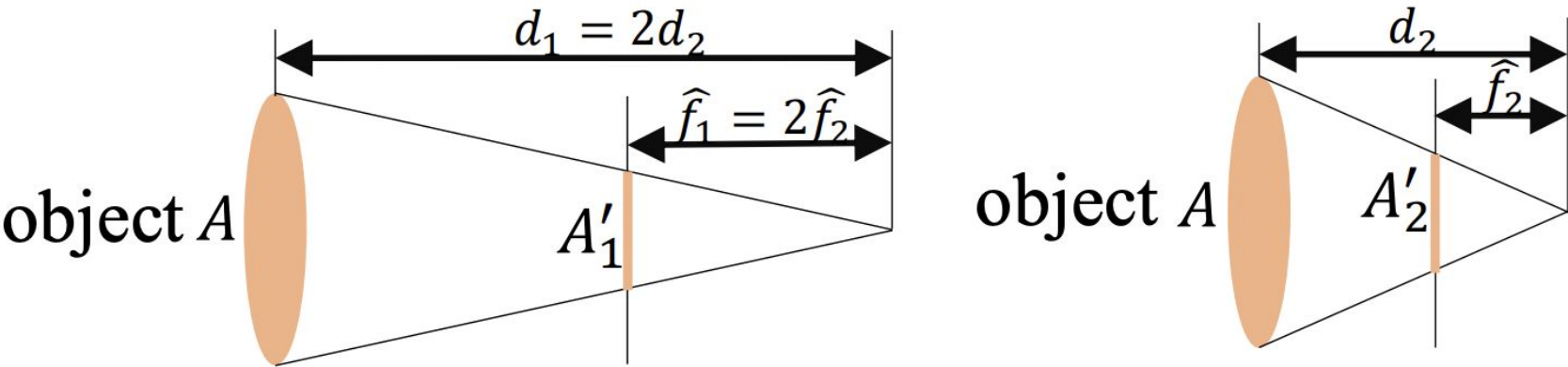


Input

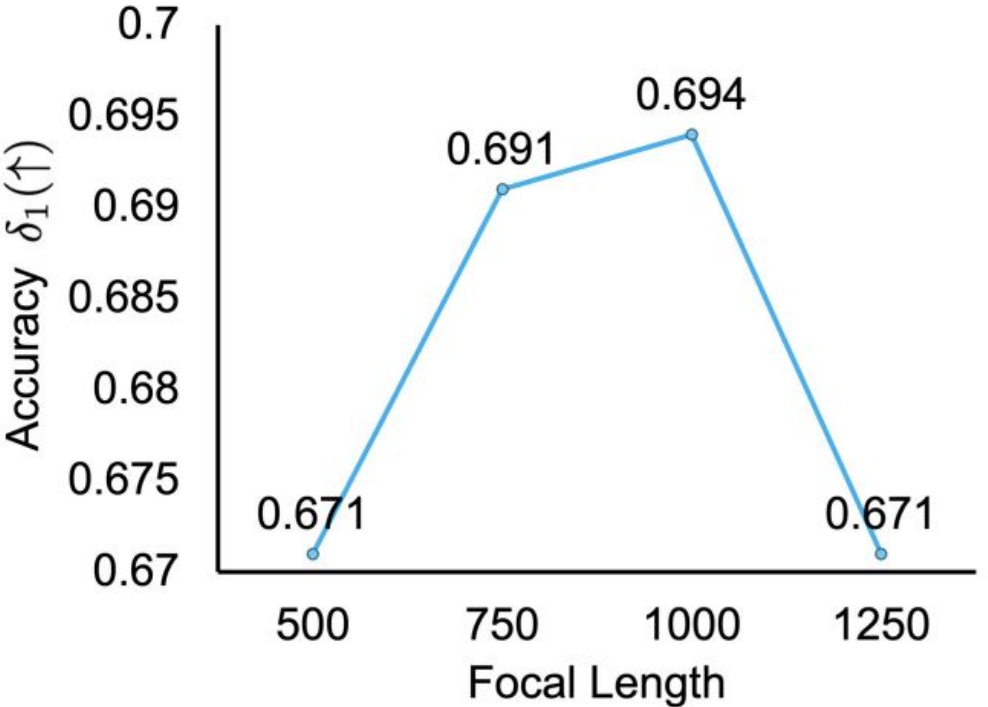
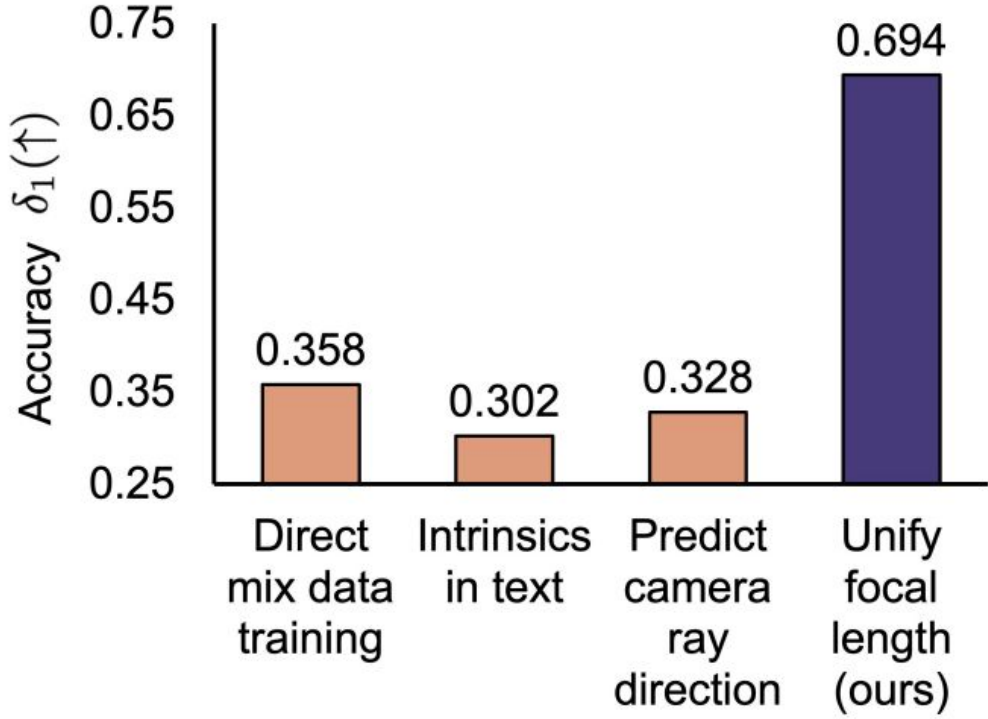
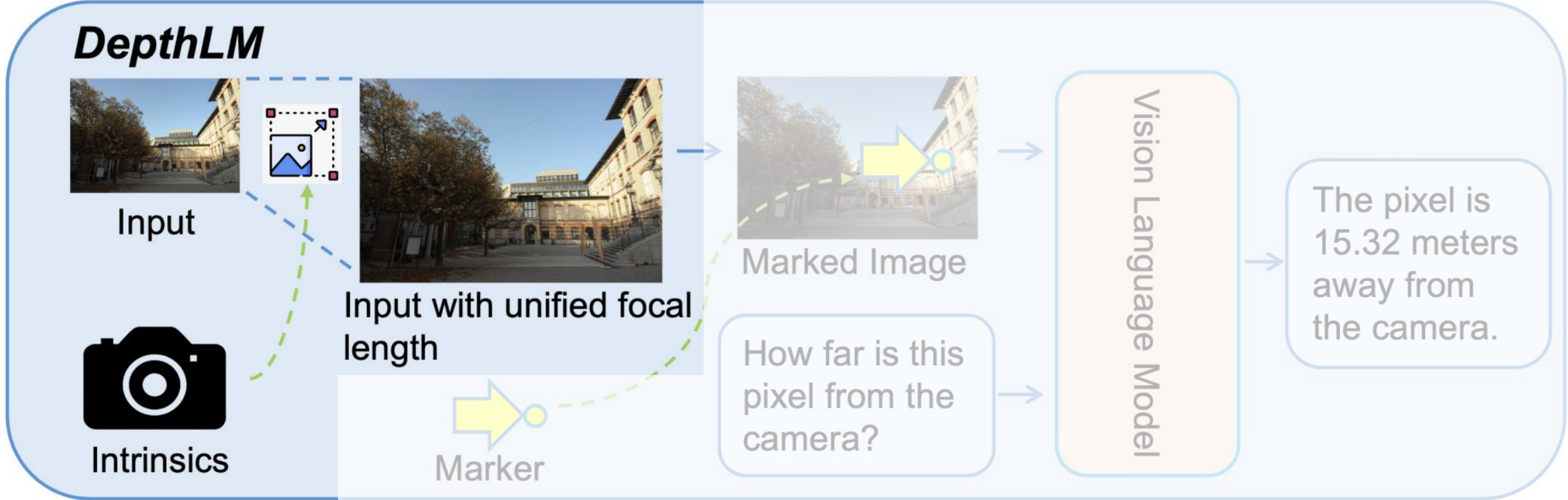
Ours



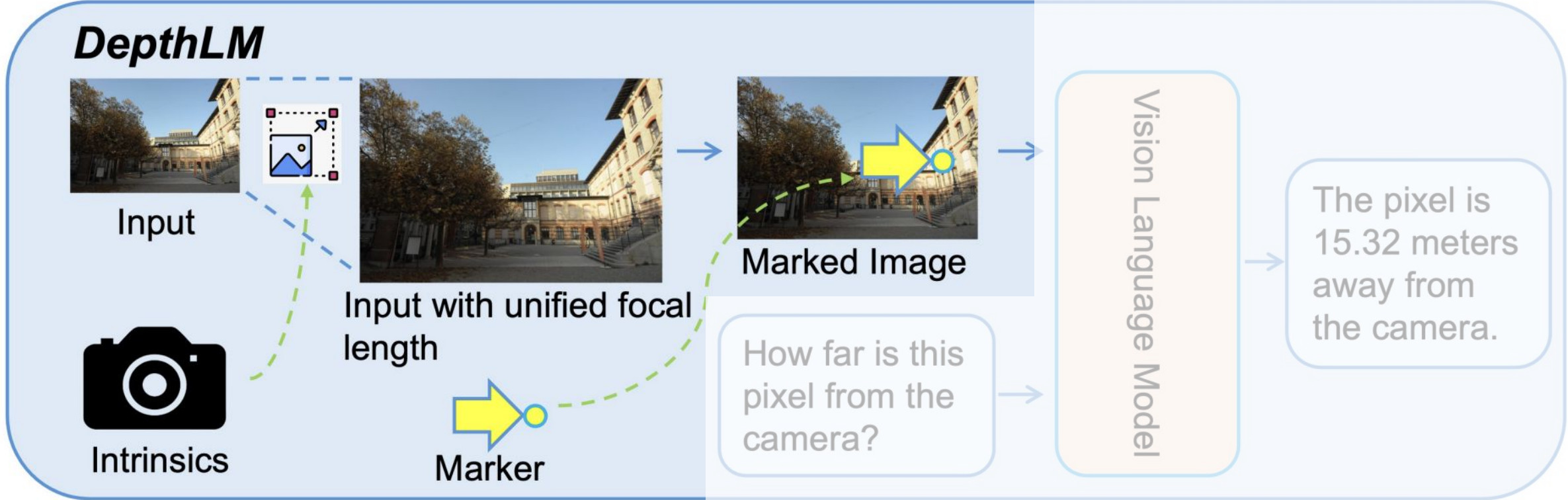
# Bottleneck 1: camera ambiguity



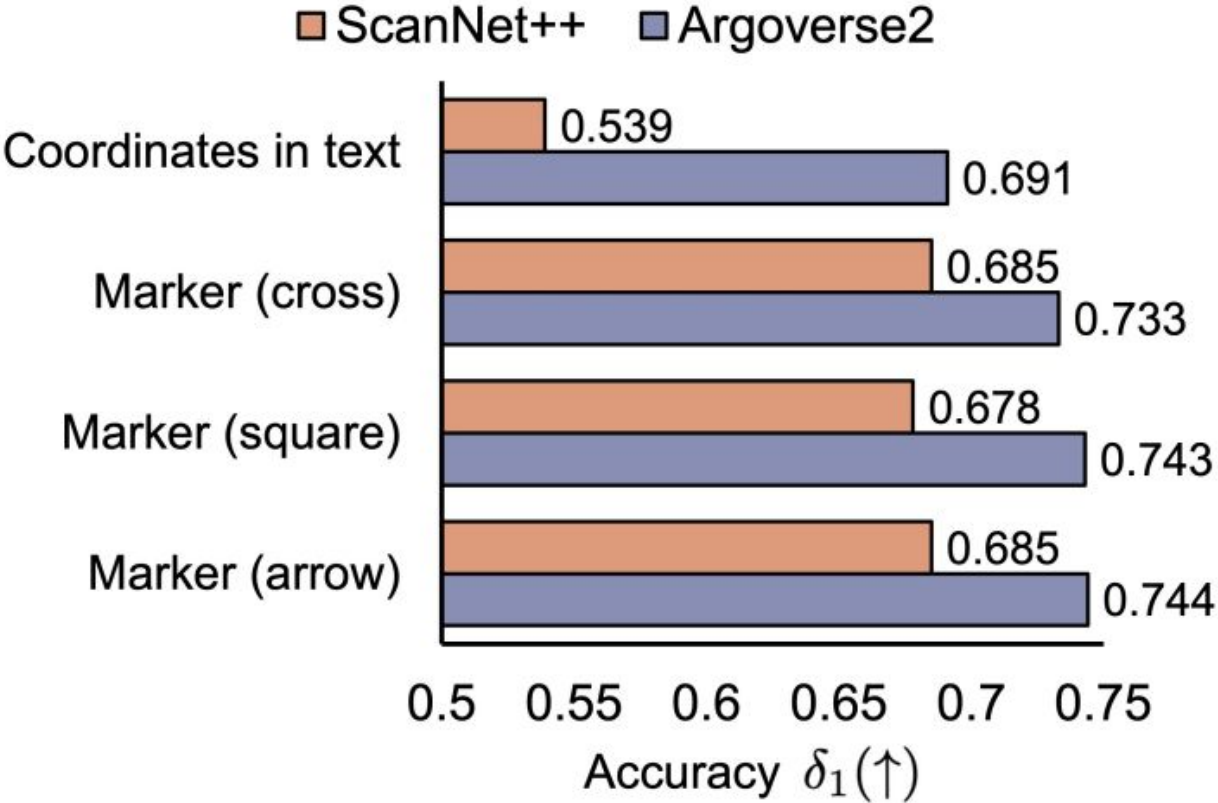
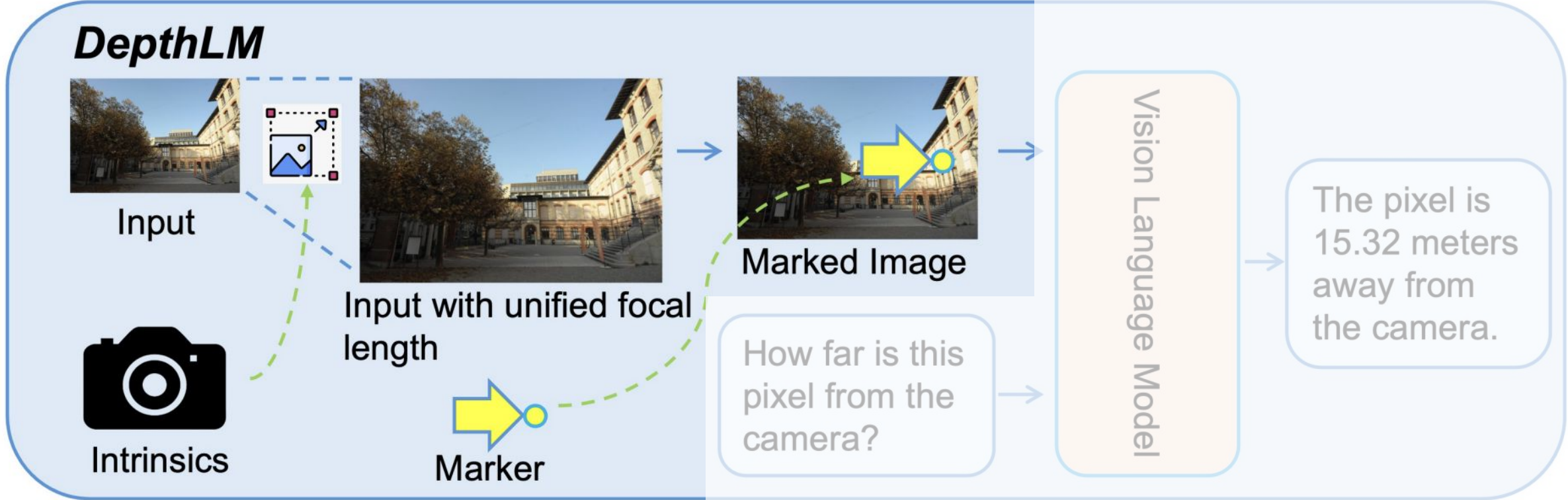
# Bottleneck 1: camera ambiguity



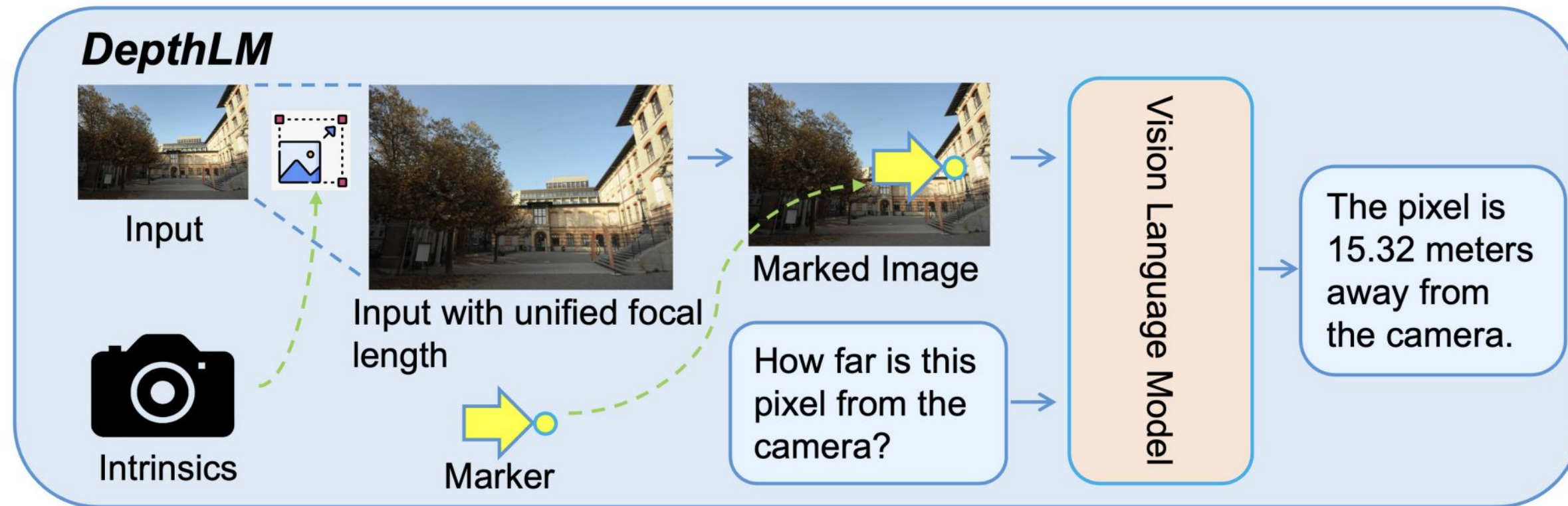
# Bottleneck 2: pixel reference



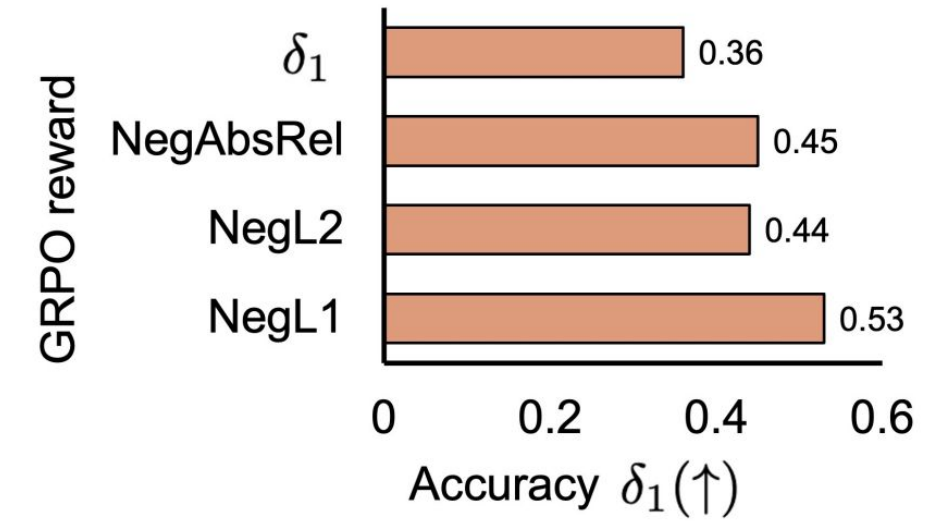
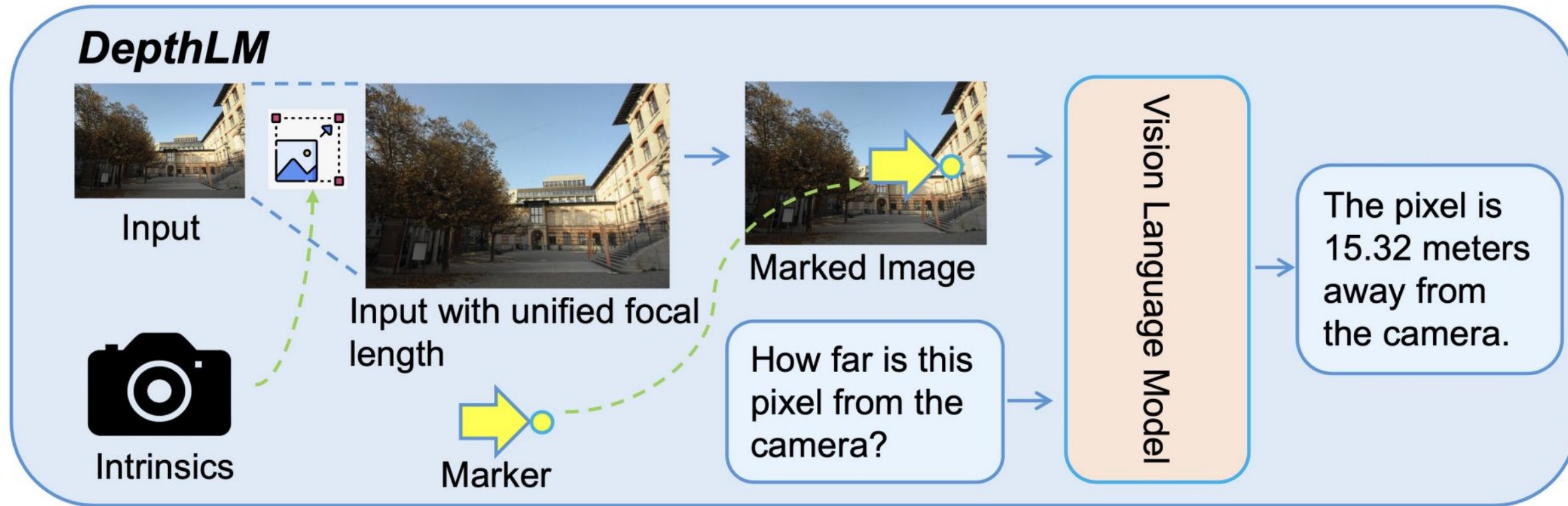
# Render Markers



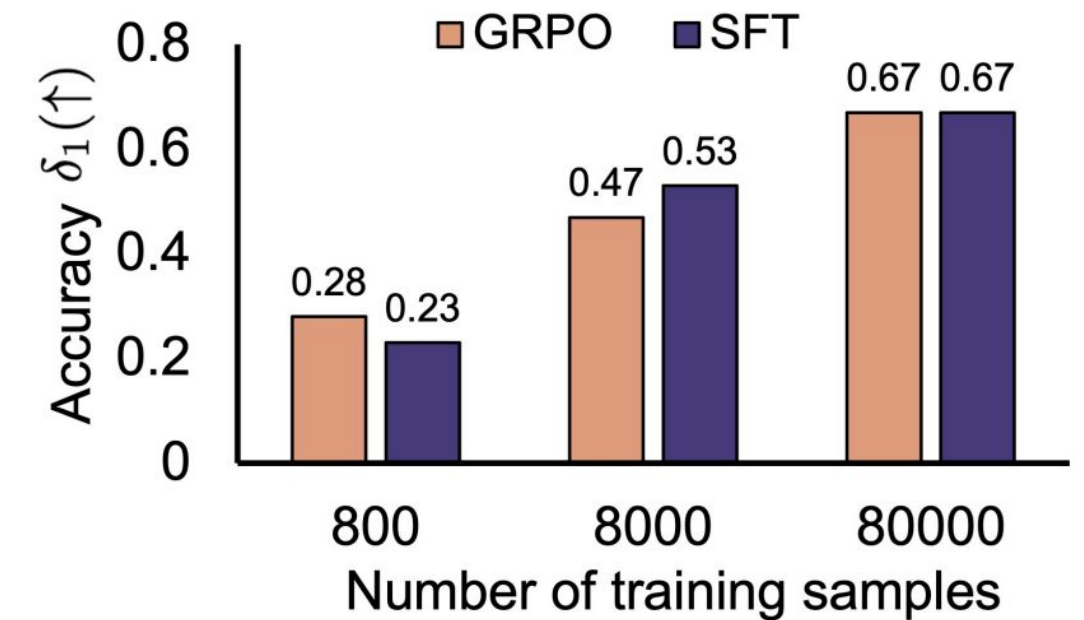
# Text Supervision is All You Need!



# SFT is More Efficient Than GRPO

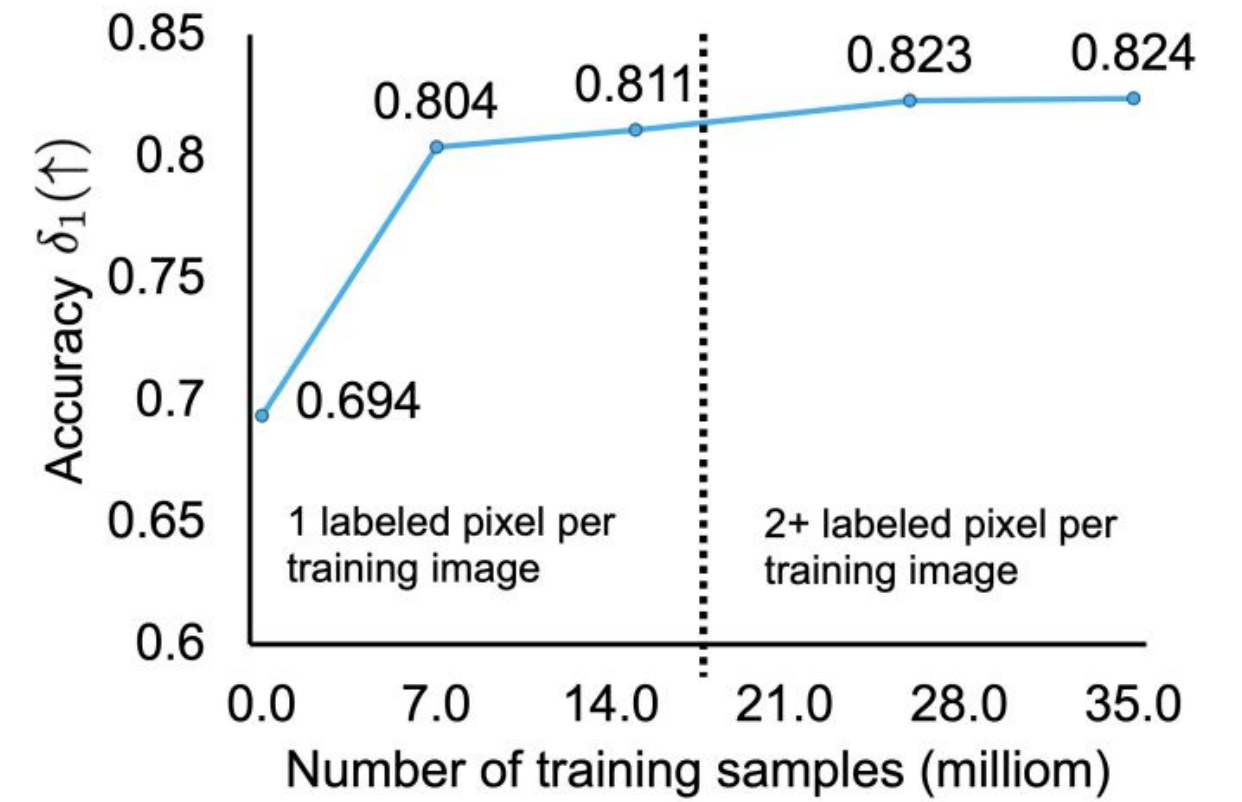
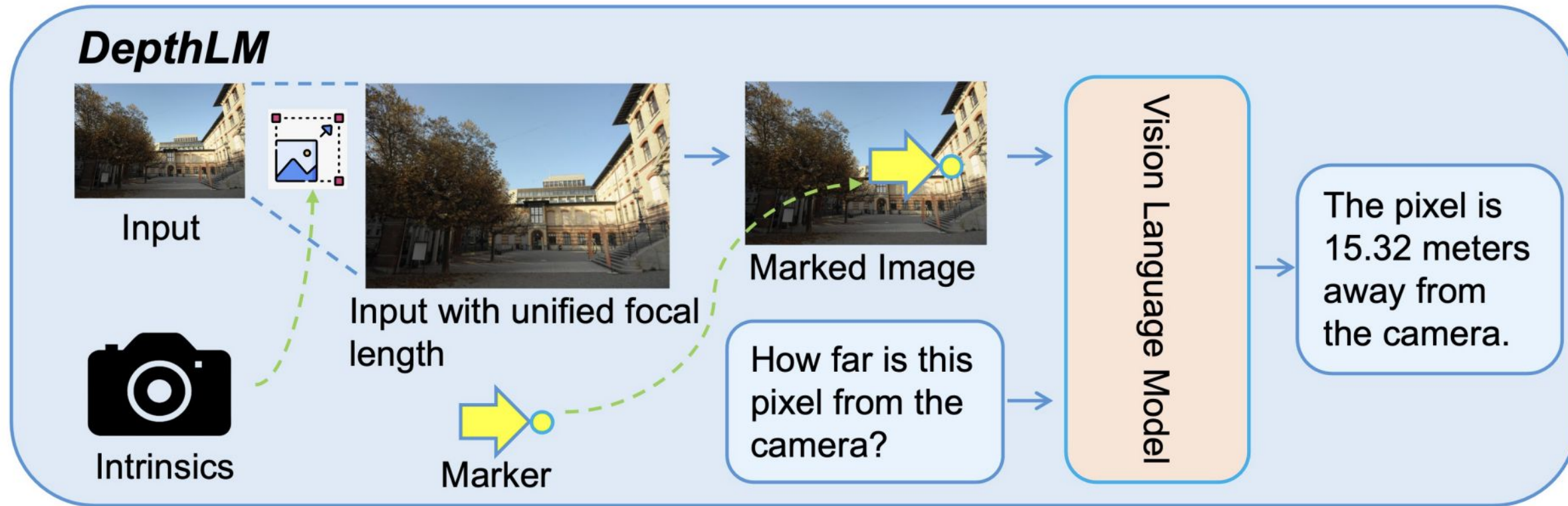


(a) Effect of GRPO rewards (8K training samples).



(b) Accuracy vs training data size.

# 1 Labeled Pixel Per Image Works!



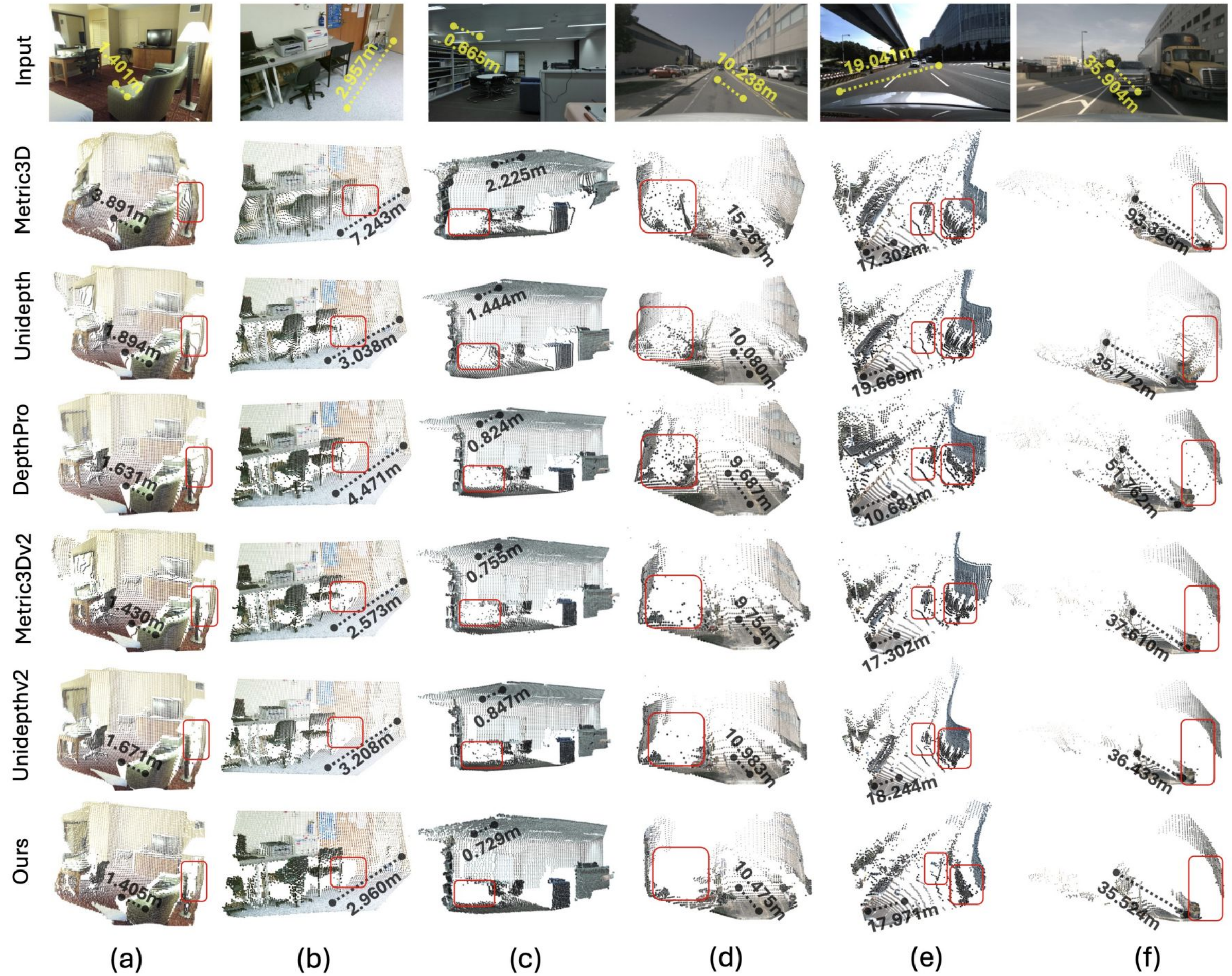
# Result

| $\delta_1(\uparrow)$ of different methods      | <i>Out</i>   |              |              | <i>Out+In</i> | <i>In</i>    |              |              |              | Average      |
|--|--------------|--------------|--------------|---------------|--------------|--------------|--------------|--------------|--------------|
|  | Argoverse2   | DDAD         | NuScenes     | ETH3D         | ScanNet++    | sunRGBD      | iBims1       | NYUv2        |              |
| <i>Naive Prediction with Constant Answers</i>  |              |              |              |               |              |              |              |              |              |
| ALWAYS OUTPUT 2.0M                             | 0.006        | 0.010        | 0.010        | 0.106         | 0.305        | 0.384        | 0.280        | 0.383        | 0.186        |
| <i>VLMs</i>                                    |              |              |              |               |              |              |              |              |              |
| QWEN2.5-VL (3B)                                | 0.133        | 0.083        | 0.090        | 0.087         | 0.120        | 0.134        | 0.080        | 0.128        | 0.106        |
| QWEN2.5-VL (7B)                                | 0.077        | 0.120        | 0.070        | 0.126         | 0.135        | 0.089        | 0.160        | 0.168        | 0.118        |
| QWEN2.5-VL (72B)                               | 0.119        | 0.140        | 0.186        | 0.220         | 0.272        | 0.276        | 0.212        | 0.324        | 0.219        |
| MOLMO (7B-D)                                   | 0.200        | 0.132        | 0.200        | 0.126         | 0.244        | 0.299        | 0.200        | 0.225        | 0.203        |
| PIXTRAL (12B)                                  | 0.157        | 0.132        | 0.118        | 0.141         | 0.318        | 0.308        | 0.270        | 0.145        | 0.199        |
| GEMINI-2.5-PRO                                 | 0.280        | 0.252        | 0.365        | 0.328         | 0.380        | 0.270        | 0.466        | 0.394        | 0.342        |
| GPT-o3   | 0.208        | 0.283        | 0.309        | 0.305         | 0.375        | 0.426        | 0.375        | 0.470        | 0.344        |
| GPT-5  | 0.218        | 0.302        | 0.382        | 0.313         | 0.428        | 0.471        | 0.307        | 0.540        | 0.370        |
| <i>Spatial VLMs</i>                            |              |              |              |               |              |              |              |              |              |
| SPACELLAVA (13B)                               | 0.100        | 0.067        | 0.083        | 0.090         | 0.269        | 0.233        | 0.208        | 0.178        | 0.154        |
| SPATIALRGPT (8B)                               | 0.055        | 0.046        | 0.100        | 0.220         | 0.346        | 0.369        | 0.240        | 0.265        | 0.205        |
| <i>VLMs Trained on Metric Depth Estimation</i> |              |              |              |               |              |              |              |              |              |
| SEED1.5-VL (OFFICIAL SETUP)                    | 0.009        | 0.012        | 0.013        | 0.219         | 0.495        | 0.321        | 0.459        | 0.412        | 0.243        |
| SEED1.5-VL (OUR PROMPT)                        | 0.040        | 0.074        | 0.028        | 0.309         | 0.593        | 0.689        | 0.627        | 0.841        | 0.400        |
| OURS (3B)                                      | 0.808        | 0.724        | <b>0.870</b> | <b>0.745</b>  | 0.838        | 0.850        | 0.890        | 0.868        | 0.824        |
| OURS (7B)                                      | <b>0.833</b> | <b>0.747</b> | 0.865        | 0.718         | <b>0.850</b> | <b>0.859</b> | <b>0.920</b> | <b>0.915</b> | <b>0.838</b> |
| OURS - PIXTRAL (12B)                           | 0.734        | 0.670        | 0.819        | 0.653         | 0.834        | 0.786        | 0.870        | 0.799        | 0.771        |

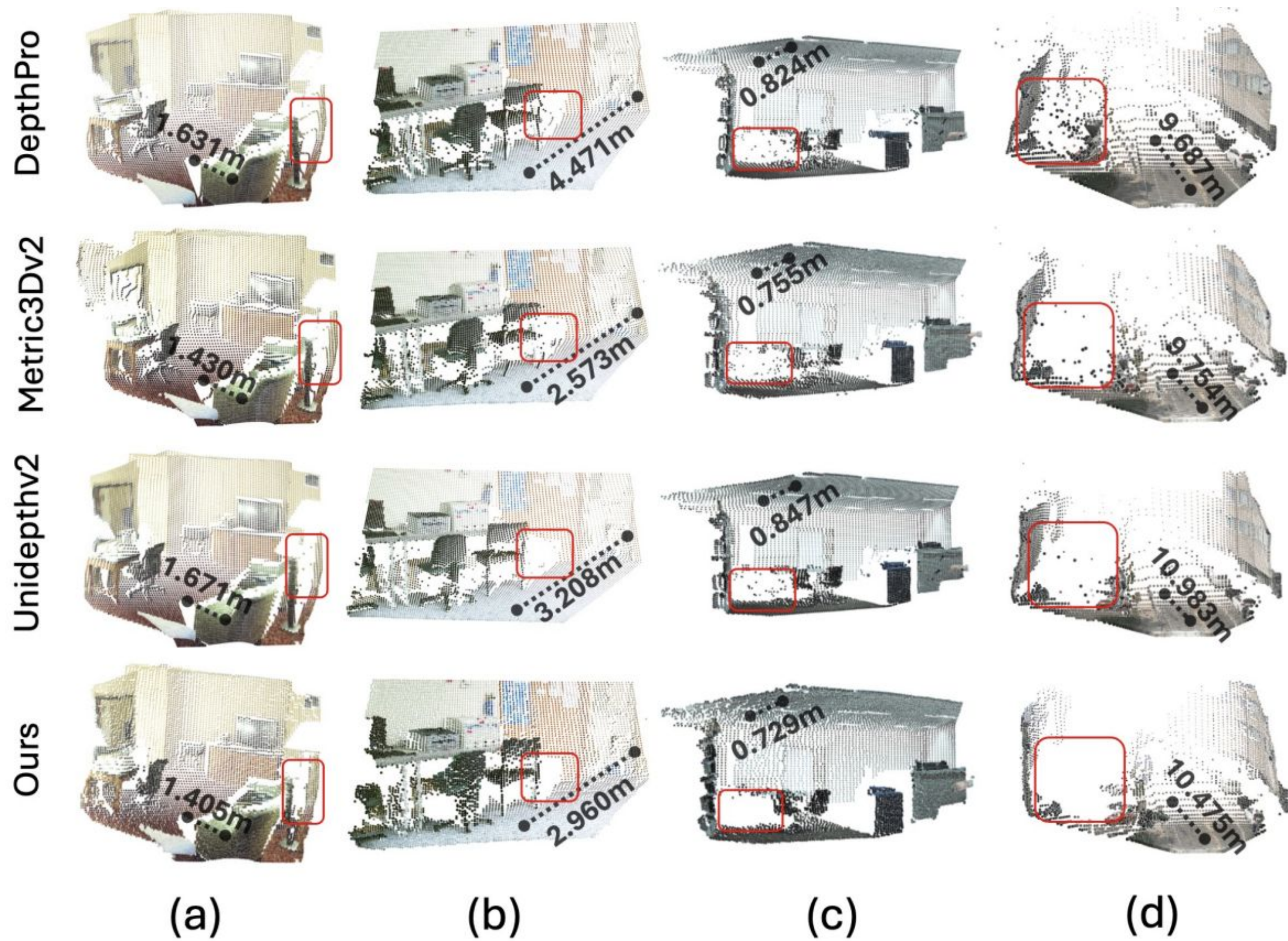
# Result

| $\delta_1(\uparrow)$ of different methods | <i>Out</i> |          | <i>Out+In</i> | <i>In</i> |        | vs Ours ( $\uparrow$ ) |
|---|------------|----------|---------------|-----------|--------|------------------------|
|   | DDAD       | Nuscenes | ETH3D         | sunRGBD   | ibims1 |                        |
| ZOEDPTH                                   | 0.272      | 0.283    | 0.350         | 0.867     | 0.580  | -42.8%                 |
| DEPTHANYTHING                             | -          | 0.354    | 0.093         | 0.850     | 0.714  | -40.3%                 |
| DEPTHANYTHINGV2                           | -          | 0.171    | 0.363         | 0.724     | -      | -48.5%                 |
| METRIC3D                                  | -          | 0.723    | 0.456         | 0.154     | 0.797  | -36.6%                 |
| UNIDPTH                                   | 0.858      | 0.846    | 0.185         | 0.943     | 0.157  | -27.3%                 |
| DEPTH PRO                                 | 0.299      | 0.566    | 0.397         | 0.831     | 0.823  | -29.1%                 |
| METRIC3DV2                                | -          | 0.841    | 0.900         | 0.812     | 0.684  | -3.8%                  |
| UNIDPTHV2                                 | 0.882      | 0.870    | 0.852         | 0.964     | 0.945  | +9.2%                  |
| OURS (7B)                                 | 0.747      | 0.865    | 0.718         | 0.859     | 0.920  | -                      |




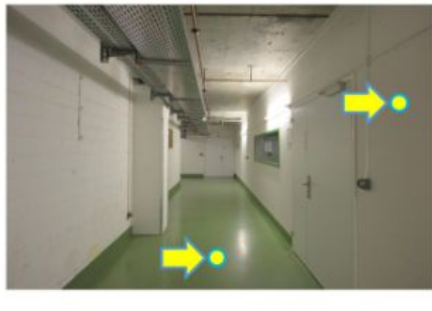

# Visualizations



# Visualizations



# Scale Up to Diverse Tasks!

|   |   |   |  |  |
|---|---|---|--|--|
|  <p><b>Principal axis distance:</b> How far is this point from the camera in the forward backward direction?</p> <ul style="list-style-type: none"> <li>• Ground Truth: 4.40m</li> <li>• GPT-5: 2.7m</li> <li>• <b>Ours: 4.30m</b></li> </ul> <p>Single Image Single Point</p> |  <p><b>Speed:</b> How many meters per second should we move in order to reach this point in exactly 4.0 seconds?</p> <ul style="list-style-type: none"> <li>• Ground Truth: 1.76m/s</li> <li>• GPT-5: 250m/s</li> <li>• <b>Ours: The point is around 7.23 meters away. Hence, the speed should be around <math>7.23 / 4.0 = 1.81\text{m/s}</math></b></li> </ul> <p>(Single Image Single Point) Reasoning</p> |  <p><b>Time:</b> How many seconds do we need to reach this point if we move towards it with the speed of 6.0m/s?</p> <ul style="list-style-type: none"> <li>• Ground Truth: 7.75s</li> <li>• GPT-5: 2.5s</li> <li>• <b>Ours: The point is around 48.28 meters away. Hence, we need around <math>48.28 / 6.0 = 8.05\text{s}</math></b></li> </ul> |  <p><b>Two point distance:</b> How far are these 2 points from each other?</p> <ul style="list-style-type: none"> <li>• Ground Truth: 2.75m</li> <li>• GPT-5: 9.87m</li> <li>• <b>Ours: 2.58m</b></li> </ul> <p>Multi-Point</p> |  <p><b>Metric Scale Camera Pose:</b> How many meters has the camera moved between these 2 images?</p> <ul style="list-style-type: none"> <li>• Ground Truth: 5.94m</li> <li>• GPT-5: 0m</li> <li>• <b>Ours: 5.62m</b></li> </ul> <p>Multi-Image</p> |
|---|---|---|--|--|

| $\delta_1(\uparrow)$ on different tasks | <i>Single image single point</i> |                         | <i>Reasoning</i> |              | <i>Multi-point</i> | <i>Multi-image</i> | Average      |
|---|----------------------------------|-------------------------|------------------|--------------|--------------------|--------------------|--------------|
|   | Distance                         | Principal axis distance | Speed            | Time         | Two point distance | Pose               |              |
| ALWAYS OUTPUT 2.0                       | 0.186                            | 0.172                   | 0.087            | 0.094        | 0.119              | 0.189              | 0.141        |
| QWEN2.5-VL (7B)                         | 0.118                            | 0.085                   | 0.136            | 0.087        | 0.066              | 0.048              | 0.09         |
| SPACELLAVA (13B)                        | 0.154                            | 0.163                   | 0.116            | 0.122        | 0.157              | 0.047              | 0.127        |
| SPATIALRGPT (8B)                        | 0.205                            | 0.132                   | 0.167            | 0.122        | 0.143              | 0.195              | 0.161        |
| SEED1.5-VL (our prompt)                 | 0.400                            | 0.174                   | 0.223            | 0.119        | 0.101              | 0.000              | 0.170        |
| Gemini-2.5-Pro                          | 0.342                            | 0.209                   | 0.213            | 0.209        | 0.140              | 0.025              | 0.189        |
| GPT-5                                   | 0.370                            | 0.241                   | 0.199            | 0.181        | 0.150              | 0.120              | 0.210        |
| <b>OURS (7B)</b>                        | <b>0.828</b>                     | <b>0.831</b>            | <b>0.817</b>     | <b>0.816</b> | <b>0.657</b>       | <b>0.876</b>       | <b>0.804</b> |

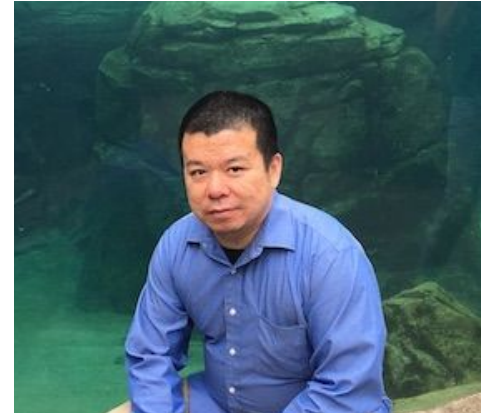
# DepthLM Team



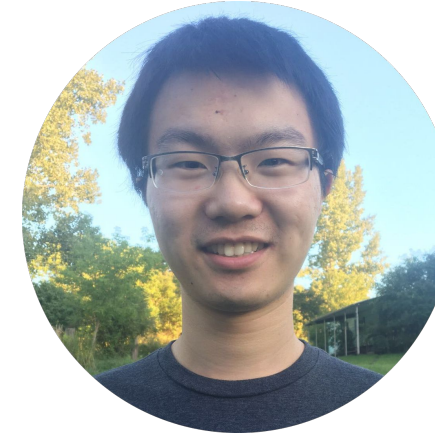
Zhipeng Cai



Ching-Feng Yeh



Hu Xu



Zhuang Liu



Gregory Meyer



Xinjie Lei



Changsheng Zhao



Shangwen Li



Vikas Chandra



Yangyang Shi



Code & Model