

# Towards Explainability in Monocular Depth Estimation

Vasileios Arampatzakis<sup>1,2</sup>, George Pavlidis<sup>2</sup>, Kyriakos Pantoglou<sup>2</sup>,  
Nikolaos Mitianoudis<sup>1</sup>, and Nikos Papamarkos<sup>1</sup>

<sup>1</sup>Democritus University of Thrace, Xanthi, Greece

<sup>2</sup>Athena Research Center, Xanthi, Greece



# Monocular Depth Estimation



- Inferring depth information from 2D images
- Crucial for: Robotics, Autonomous driving, Augmented reality
- An inherently ill-posed problem
- Ambiguities caused by the projection of the 3D world to 2D images
- Significance: Enhances scene understanding and 3D perception
- Deep Learning-based methods outperform traditional approaches
- CNNs, Vision Transformers capture complex patterns in images

# Typical Explainability Methods

- Unveiling the rationale behind model predictions
- Crucial for: Transparency & trust, Model improvement, Bias detection, User understanding, Regulatory compliance
- Methods: Feature visualization, Saliency maps, Attention mechanisms, LIME (Local Interpretable Model-agnostic Explanations), SHAP (Shapley Additive exPlanations), Grad-CAM (Gradient-weighted Class Activation Mapping)
- Trade-offs between simplicity and accuracy in explanation methods
- Some methods are model-specific, while others are model-agnostic
- The need to validate explanations and ensure they reflect true model behavior

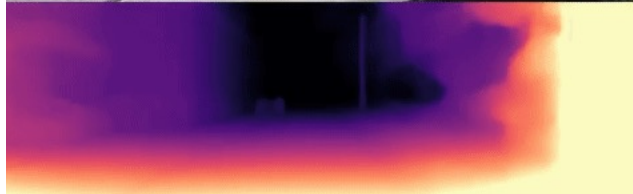
# Unique Approach: Connecting with Human Perception

- **Objective:** Enhance explainability by aligning model predictions with how humans perceive depth
- **Key idea:** *As a dataset is limited to provide only a single cue, the accuracy of the methods indirectly reflects their success in learning the selected depth cue.*

Input: RGB image  
(multiple visual depth cues)



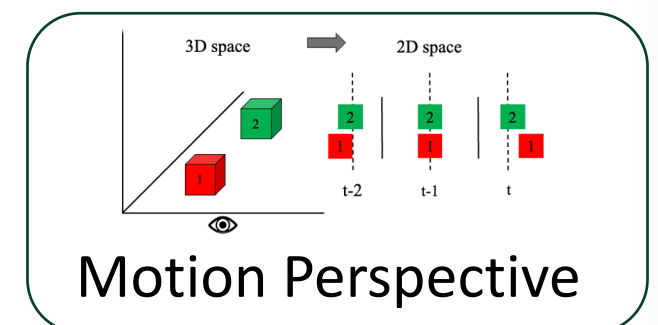
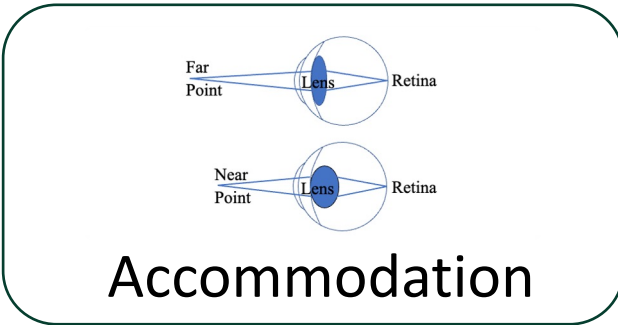
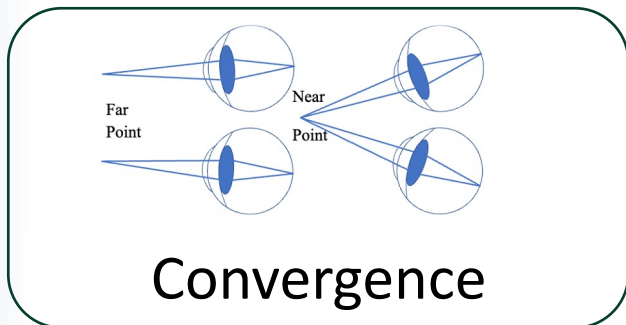
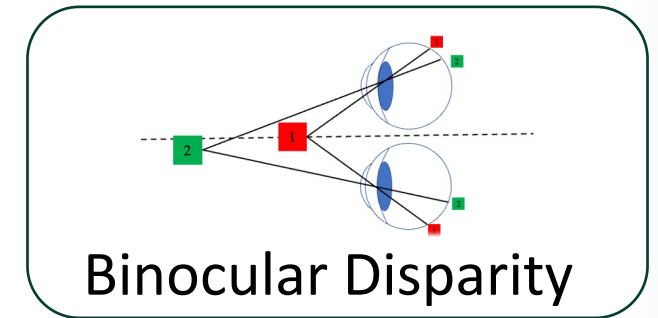
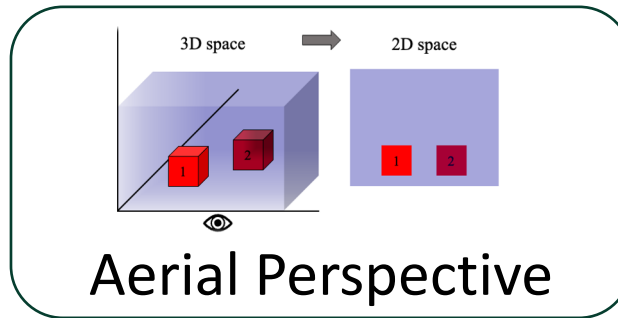
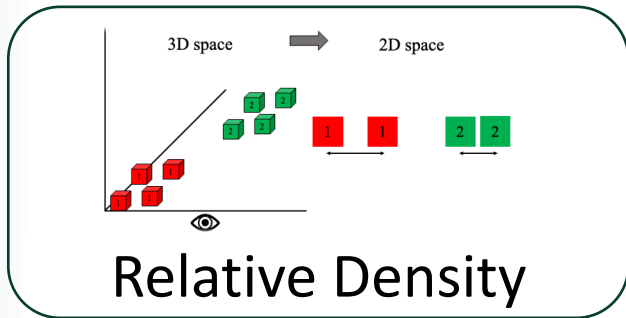
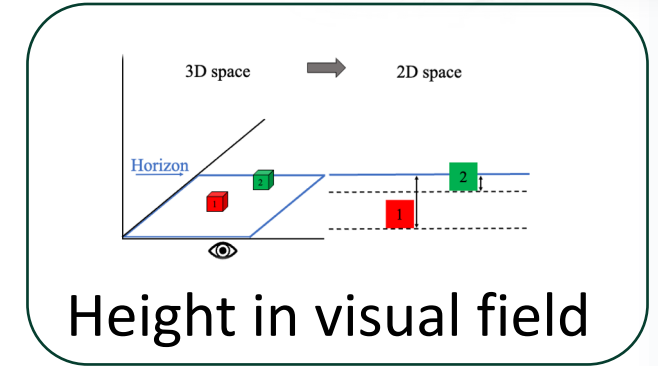
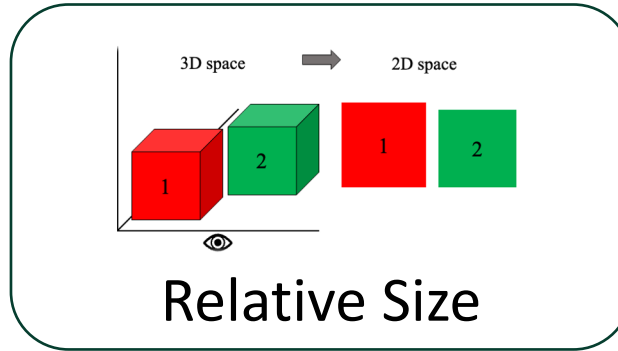
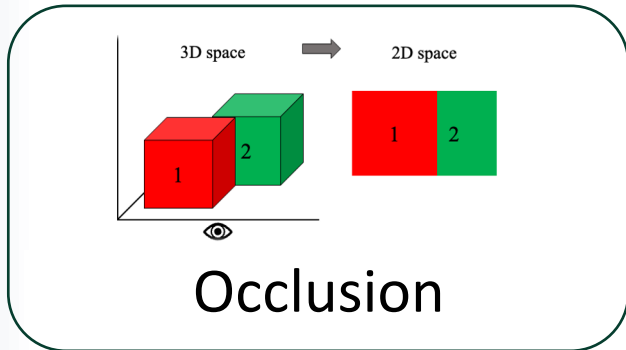
Output: depth map



*Cutting & Vishton (1995): Perceiving layout and knowing distances: The integration, relative potency, and contextual use of different information about depth*



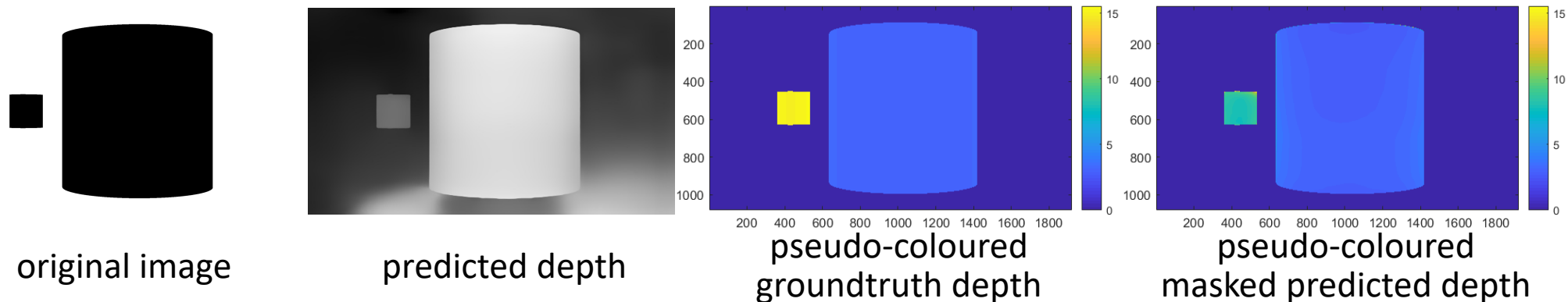
# Visual Depth Cues



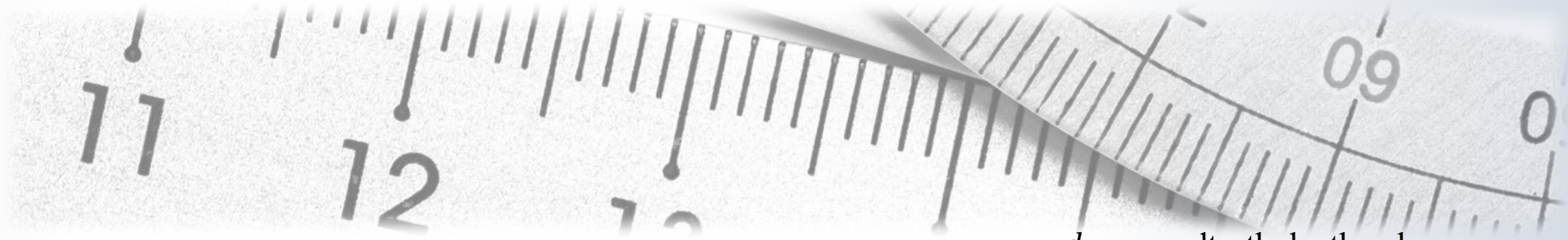
# Visual Depth Cue (VDC) Dataset



- A synthetic resource for monocular depth estimation inspired by human perception (Cutting & Vishton, Nagata)
- Exclusive depth cue representation in each image
- Relative Size ( $\approx 23800$  images ):
  - 2D images of black cylindrical objects at various distances against a white background, created through perspective projections of the corresponding virtual 3D scenes



# Metrics



$d_i$ : groundtruth depth value  
 $\hat{d}_i$ : predicted depth value  
 $N$ : number of samples

- **Absolute Relative Error:**

$$AbsRel = \frac{1}{N} \sum_{i=1}^N \frac{|d_i - \hat{d}_i|}{d_i}$$

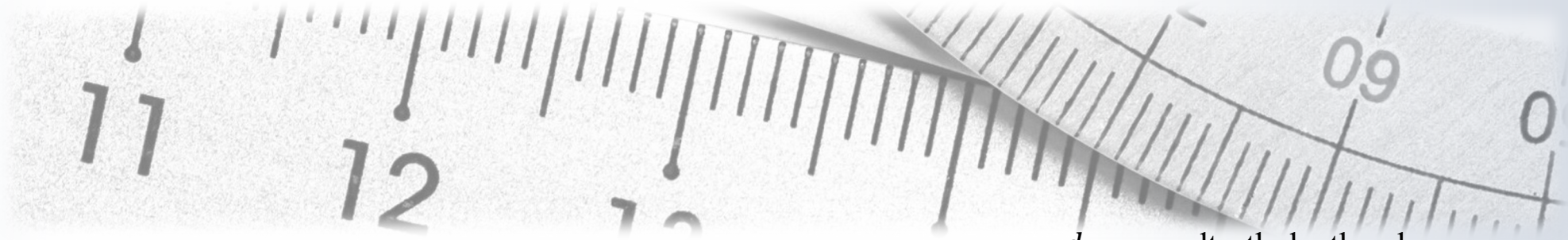
- **Squared Relative Error:**

$$SqRel = \frac{1}{N} \sum_{i=1}^N \frac{(d_i - \hat{d}_i)^2}{d_i}$$

- **Linear Root Mean Squared Error:**

$$RMSE(lin) = \sqrt{\frac{1}{N} \sum_{i=1}^N (d_i - \hat{d}_i)^2}$$

# Metrics



$d_i$ : groundtruth depth value  
 $\hat{d}_i$ : predicted depth value  
 $N$ : number of samples

- **Logarithmic Root Mean Squared Error:**

$$RMSE(\log) = \sqrt{\frac{1}{N} \sum_{i=1}^N (\log d_i - \log \hat{d}_i)^2}$$

- **Scale-invariant Mean Squared Error [Eigen]:**

$$sRMSE(\log) = \frac{1}{N} \sqrt{\sum_{i=1}^N (\log d_i - \log \hat{d}_i + a(d_i, \hat{d}_i))^2},$$

$$\text{where } a(d_i, \hat{d}_i) = \frac{1}{N} \sum_i (\log \hat{d}_i - \log d_i)$$

- **Accuracy with threshold ( $\delta_x$ ):**

$$(\%) \text{ of } d_i \text{ such that } \max\left(\frac{d_i}{\hat{d}_i}, \frac{\hat{d}_i}{d_i}\right) = \delta < thr,$$

$$\text{where } thr = 1.25, 1.25^2, 1.25^3$$



# Models

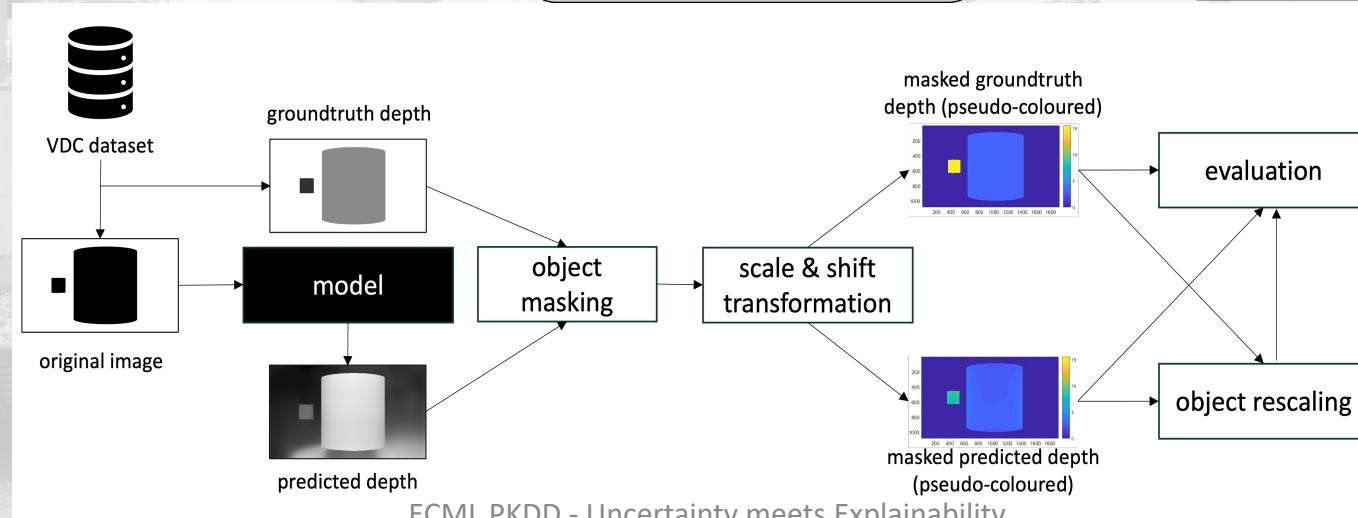
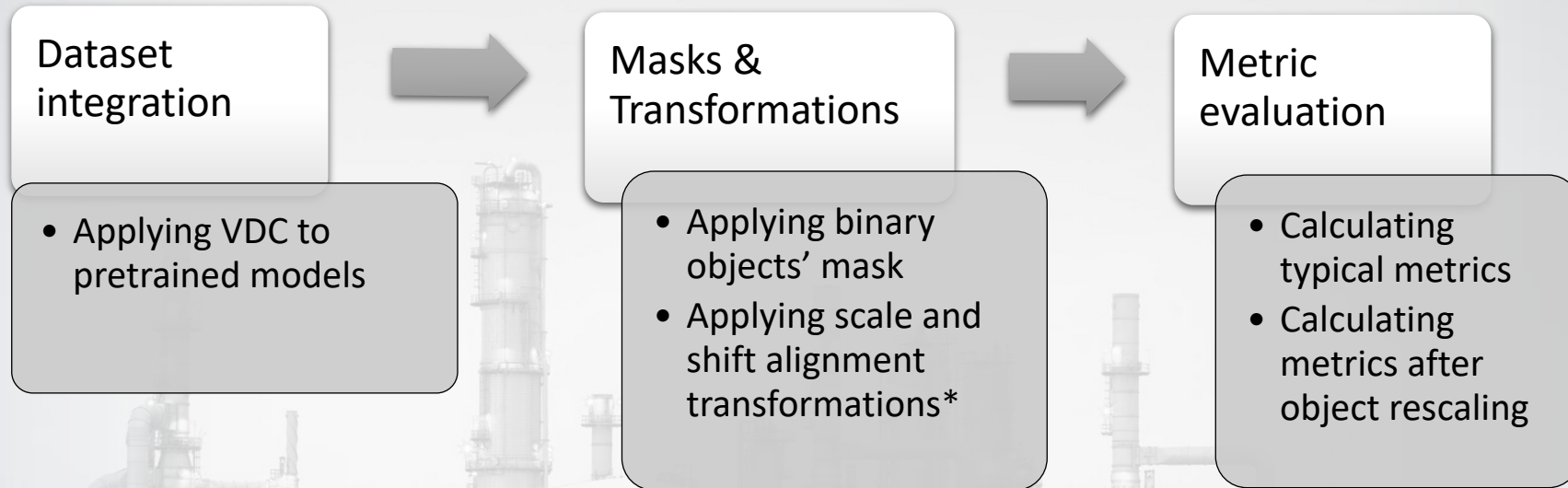


- 12 Pretrained state-of-the-art models:
  - *MiDaS* (4 variations)
  - *Monodepth2* (6 variations)
  - *DenseDepth* (2 variations)

Table: Evaluation on KITTI dataset, using the Eigen split

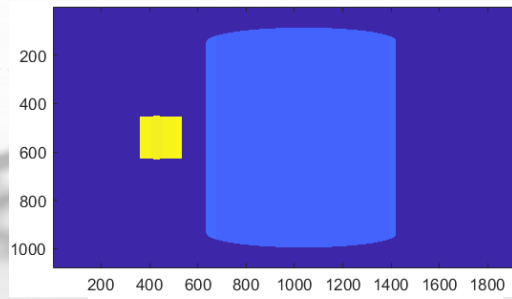
Model	Year	Citations	Version	<i>AbsRel</i>	<i>SqRel</i>	<i>RMSE</i>	<i>RMSE<sub>log</sub></i>	$\delta_1$	$\delta_2$	$\delta_3$
<i>MiDaS</i>	2020	721	<i>dpt_hybrid</i>	0.062	0.222	2.575	0.092	0.959	0.995	0.999
<i>Monodepth2</i>	2019	1708	<i>mono_640x192</i>	0.115	0.903	4.863	0.193	0.877	0.959	0.981
<i>DenseDepth</i>	2018	435	<i>kitti</i>	0.093	0.589	4.170	0.171	0.886	0.965	0.986
Eigen	2014	3782	<i>(baseline)</i>	0.190	1.511	7.156	0.270	0.692	0.899	0.967

# Experiment Pipeline Overview

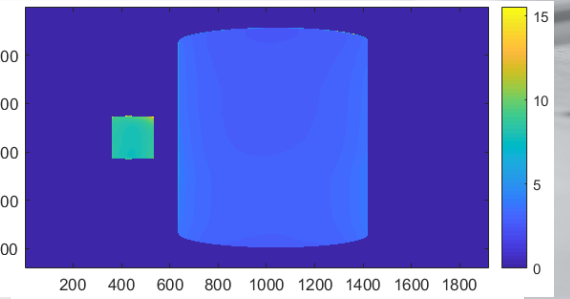


# Results

An error in the estimated depth of the far object (depicted significantly smaller) will have a negligible impact on the metrics, while the estimate should be balanced



groundtruth depth



masked predicted depth

Table: Evaluation on VDC dataset, mean values

		<i>Typical metric results</i>								<i>Rescaled object metric results</i>							
<b>Model</b>	<b>Version</b>	<i>AbsRel</i>	<i>SqRel</i>	<i>RMSE</i>	<i>RMSE<sub>log</sub></i>	<i>sRMSE</i>	$\delta_1$	$\delta_2$	$\delta_3$	<i>AbsRel</i>	<i>SqRel</i>	<i>RMSE</i>	<i>RMSE<sub>log</sub></i>	<i>sRMSE</i>	$\delta_1$	$\delta_2$	$\delta_3$
<i>MiDaS</i>	<i>midas_v21</i>	0.056	1.718	10.097	0.040	0.003	0.964	0.987	0.993	0.111	7.562	21.852	0.088	0.014	0.853	0.908	0.935
<i>Monodepth2</i>	<i>stereo1024x320</i>	0.104	57.259	21.226	0.078	0.008	0.903	0.957	0.978	0.202	228.261	42.372	0.165	0.031	0.708	0.799	0.856
<i>DenseDepth</i>	<i>kitti</i>	0.103	3.472	20.153	0.083	0.009	0.902	0.951	0.972	0.207	16.657	40.576	0.178	0.036	0.704	0.785	0.837

# Results

- *MiDaS* demonstrate superior performance (rescaled  $\delta_1 \approx 0.85$ ):  
Potential for partial learning of the relative size cue
- *Densedepth* (*nyu*) also exhibit enhanced accuracy:  
Pivotal role of training datasets: *Densedepth* (*kitti*) displays weaker results

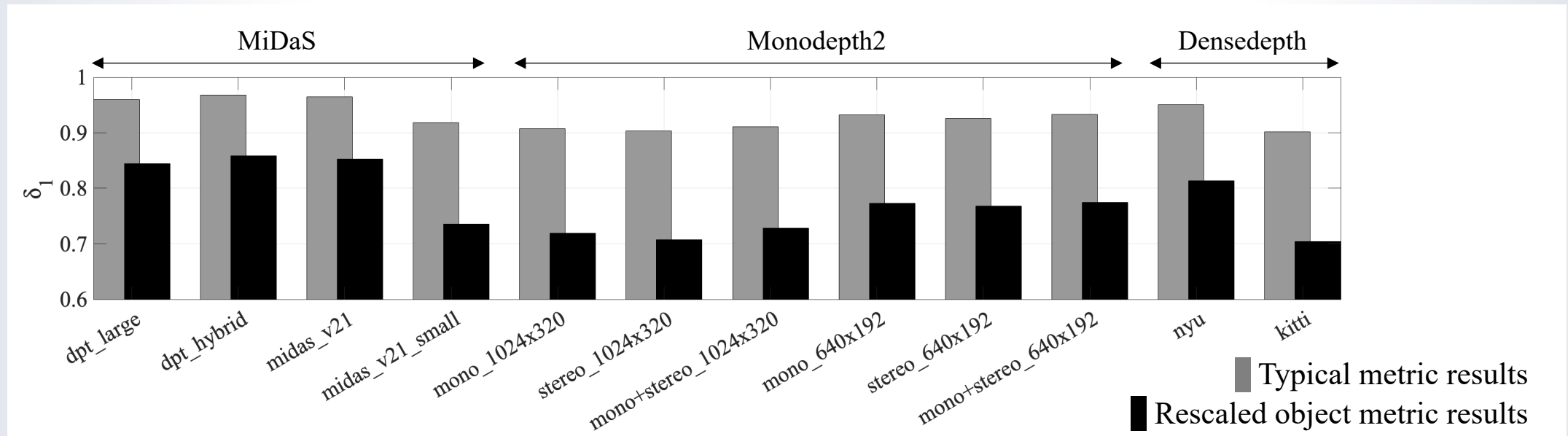


Figure: Accuracy with threshold  $\delta_1$  of the pertained models tested on our dataset



# Conclusion

- Preliminary study: New explainability concept for monocular depth estimation
- Creating a novel dataset:
  - Introducing the Visual Depth Cue Dataset (VDC)
- Testing pretrained methods on a single visual depth cue:
  - Exploring relative size
- Assessing indirect success:
  - Metrics unveil monocular depth estimation performance
- Balancing metrics:
  - The role of rescaled object assessments
- Future directions:
  - Expanding VDC: Incorporating other visual depth cues → benchmark dataset
  - Evaluating state-of-the-art method efficiency
  - Bridging Deep Learning with Human Perception
  - New depth estimation models aligned with human perception