

# Temporally Consistent Online Depth Estimation in Dynamic Scenes



JOHNS HOPKINS UNIVERSITY

Zhaoshuo Li<sup>1</sup>, Wei Ye<sup>2</sup>, Dilin Wang<sup>2</sup>, Francis X. Creighton<sup>1</sup>, Russell H. Taylor<sup>1</sup>, Ganesh Venkatesh<sup>2</sup>, Mathias Unberath<sup>1</sup>

<sup>1</sup>Johns Hopkins University <sup>2</sup>Reality Labs, Meta Inc.



## Overview

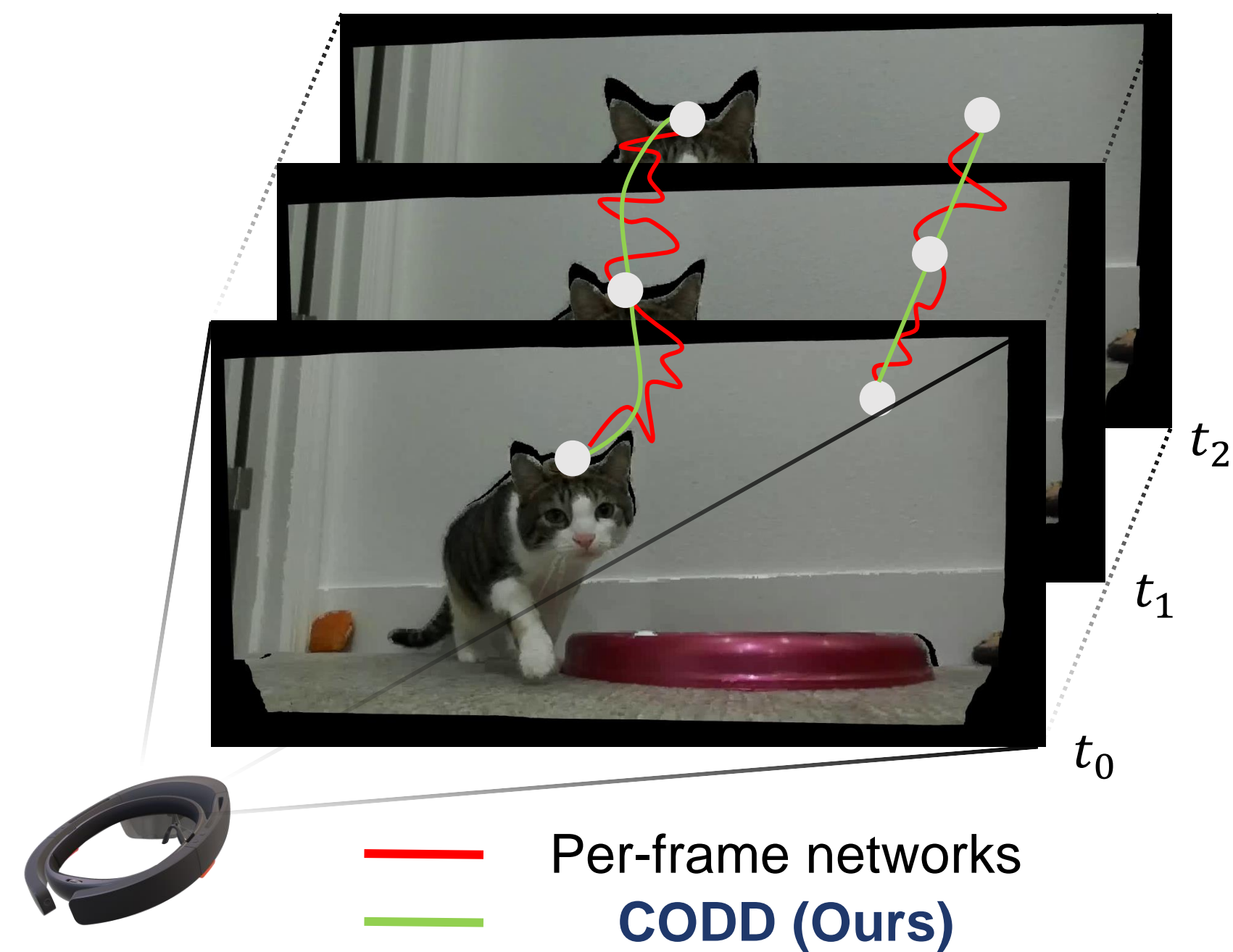
Depth estimation can be used for 3D reconstruction.

Temporally consistency has been largely overlooked. Consistency is critical for applications such as mixed reality, as jitters in depth estimates corrupt visual quality.

### Challenges

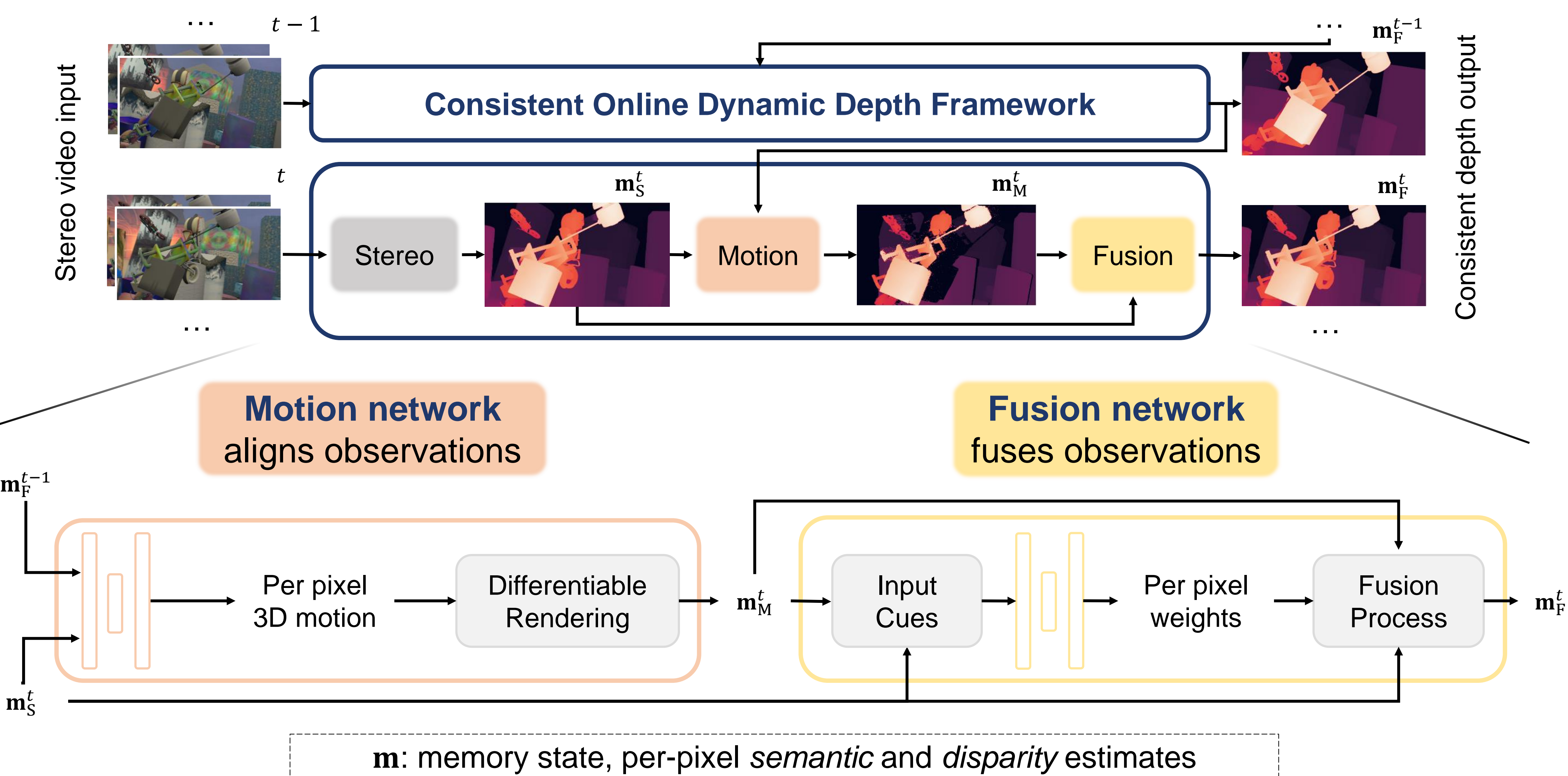
- Online – no future information is available,
- Dynamics – camera motion, object motion and deformation.

Consistent Online Dynamic Depth (CODD) framework is developed to mitigate the above challenges.



## Consistent Online Dynamic Depth Framework

CODD produces temporally consistent depth for dynamic scenes in an online setting.



## Result

### Metrics

EPE – per-frame depth accuracy.

$$d_{GT} - d_{pred}$$

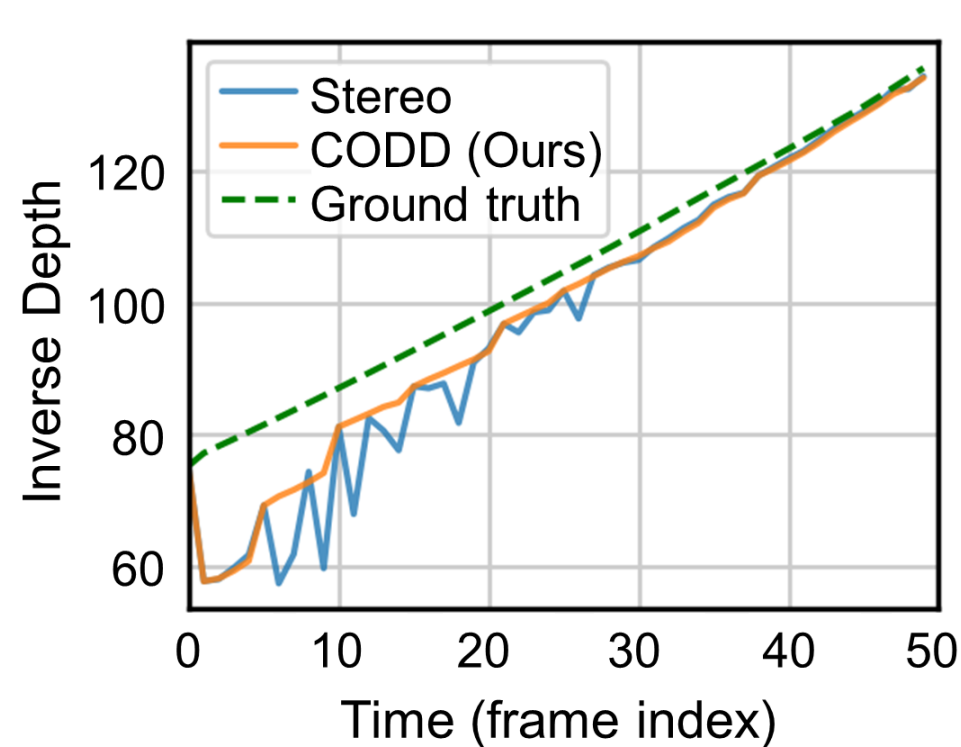
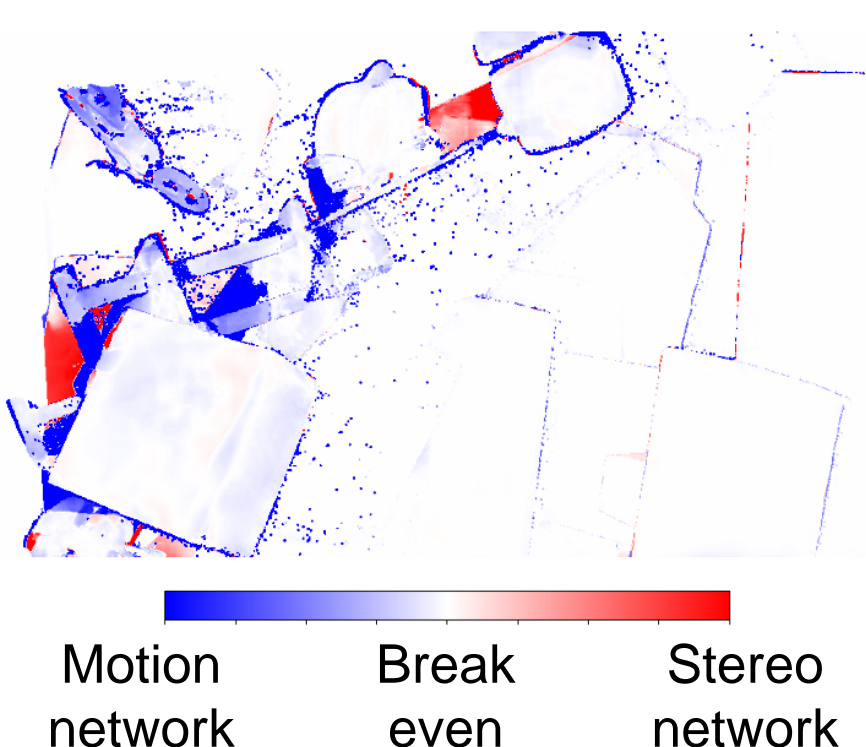
TEPE – temporal depth consistency, capturing consistency of fast-moving objects.

$$|(d_{GT,t} - d_{GT,t-1}) - (d_{pred,t} - d_{pred,t-1})|$$

TEPE<sub>r</sub> – relative temporal depth consistency, capturing consistency of static objects.

$$\frac{|(d_{GT,t} - d_{GT,t-1}) - (d_{pred,t} - d_{pred,t-1})|}{|d_{GT,t} - d_{GT,t-1}|}$$

### Qualitative Results



### Experiments

HITNet [1] is used as the per-frame stereo network. Our fusion network is compared against classical Kalman Filter algorithm [2].

CODD performs better by up to 31% for temporal depth consistency and performs on par for per-frame depth accuracy.



[1] Tankovich et al. Hitnet: Hierarchical iterative tile refinement network for real-time stereo matching. CVPR 2021.  
[2] Kalman. A new approach to linear filtering and prediction problems. JFE 1960.

