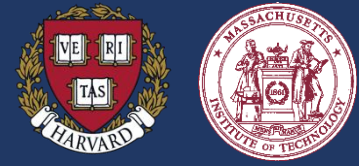


PAGE-4D: Disentangled Pose and Geometry Estimation for 4D Perception

Kaichen Zhou · Yuhan Wang · Grace Chen · Xinhai Chang · Gaspard Beaudouin · Fangneng Zhan · Paul Pu Liang † · Mengyu Wang †
Harvard · MIT · Imperial College · ENPC



INTRODUCTION

The Problem

Static-only training → Dynamic-scene failure
Camera pose errors & geometry distortions when dynamic agents (people, vehicles) are present.

MOTIVATION - 2

Empirical Observation

- Dynamic regions exhibit weaker activations compared to static ones, suggesting that static-only model tends to ignore dynamic content.
- Masking dynamic regions during the feedforward process can improve pose estimation performance.

EXPERIMENTAL RESULTS

Quantitative Results

- >30%** Pose Estimation
- >18%** Video Depth Estimation
- >17%** Monocular Depth Estimation
- >4%** View Synthesis

MOTIVATION - 1

Conflicting Task Objectives

Multi-task 4D reconstruction faces a fundamental tension:

- Pose Estimation: Needs to suppress dynamic regions
- Geometry Recon.: Needs to model dynamic regions

$$\mathbf{x}_t = \mathbf{K} \left[\mathbf{R}_{t \leftarrow r} D_r(\mathbf{x}_r) \mathbf{K}^{-1} \mathbf{x}_r + \mathbf{t}_{t \leftarrow r} \right] + \mathbf{K} \mathbf{M}_{t \leftarrow r}$$

$$\delta(\mathbf{x}_r) \equiv \tilde{\mathbf{x}}_t^T \mathbf{E} \tilde{\mathbf{x}}_r \approx \frac{1}{Z_r} \mathbf{n}(\mathbf{x}_r)^T \Delta \mathbf{X}_\perp(\mathbf{x}_r)$$

METHOD: PAGE-4D ARCHITECTURE

Overview

- Training — Training only the middle layers to improve efficiency.
- Memory-Efficient Mask - Replaces the quadratic mask with a linear-memory additive mask while preserving attention equivalence.

$$\mathbf{q}'_i = [q_i \sqrt{d'/d}, r_i \sqrt{d'/d}], \mathbf{k}'_j = [k_j, c_j], \mathbf{v}'_j = [v_j, 0]$$

Qualitative Results