

# MonST3R: A Simple Approach for Estimating Geometry in the Presence of Motion

Junyi Zhang<sup>1</sup>, Charles Hermann<sup>2,+</sup>, Junhwa Hur<sup>2</sup>, Varun Jampani<sup>3</sup>, Trevor Darrell<sup>1</sup>, Forrester Cole<sup>2</sup>, Deqing Sun<sup>2,\*</sup>, Ming-Hsuan Yang<sup>2,4,\*</sup>  
<sup>1</sup>UC Berkeley, <sup>2</sup>Google DeepMind, <sup>3</sup>Stability AI, <sup>4</sup>UC Merced

## Introduction

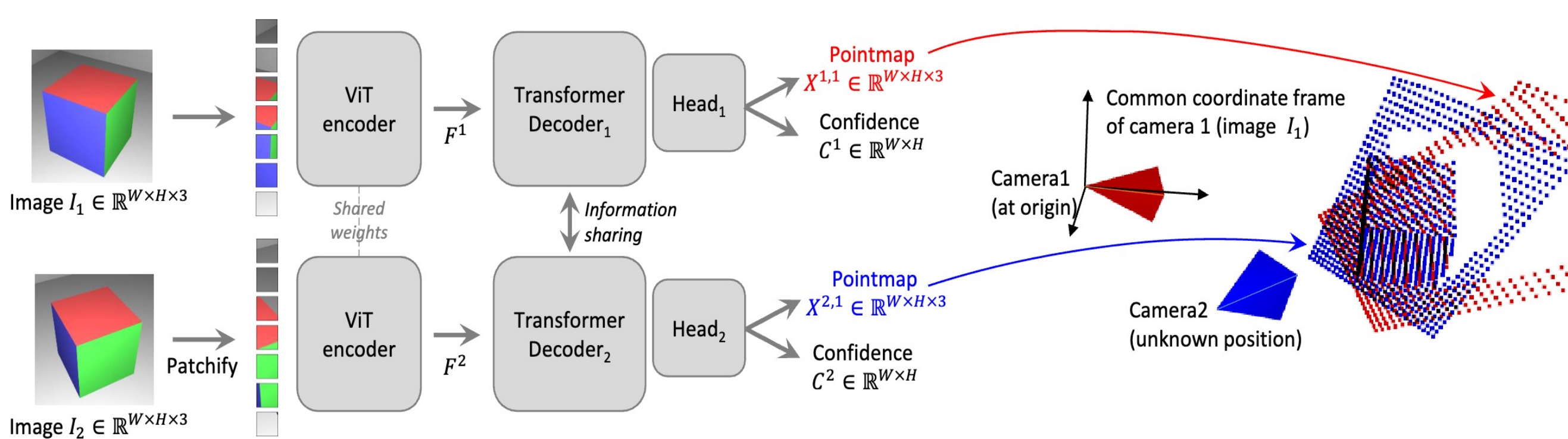
- Task: estimating global geometry** given a casually-captured **monocular video** of **dynamic scenes**, in a primarily **feed-forward** manner
- Existing methods** rely on multi-stage pipelines or global optimizations that decompose the problem into subtasks, complex and prone to errors
- How:** we take a geometry-first approach that directly estimates per-timestep geometry of dynamic scenes
- Key insight:** by simply estimating a pointmap for each timestep, we adapt DUS<sub>t</sub>3R's representation, previously used for static scenes, to dynamic scenes.
- Challenge:** despite the scarcity of training data, we show that by posing the problem as a fine-tuning task, strategically training the model on limited data can surprisingly enable it to handle dynamics

## Overview



Given a **video** of dynamic scene, MonST3R processes it to produce a time-varying **dynamic point cloud**, along with per-frame camera poses and intrinsics, in a predominantly **feed-forward** manner

## Pointmap Representation of DUS<sub>t</sub>3R



Given **two frames**, DUS<sub>t</sub>3R estimates two corresponding pointmaps (xyz coordinates for each pixel), **aligned in the camera coordinate system of the first frame**; from which, camera intrinsics, pose, and depth can be derived

No constraint on dynamic/static scenes in the representation!  
 But how does the model actually work for dynamic scenes? ->

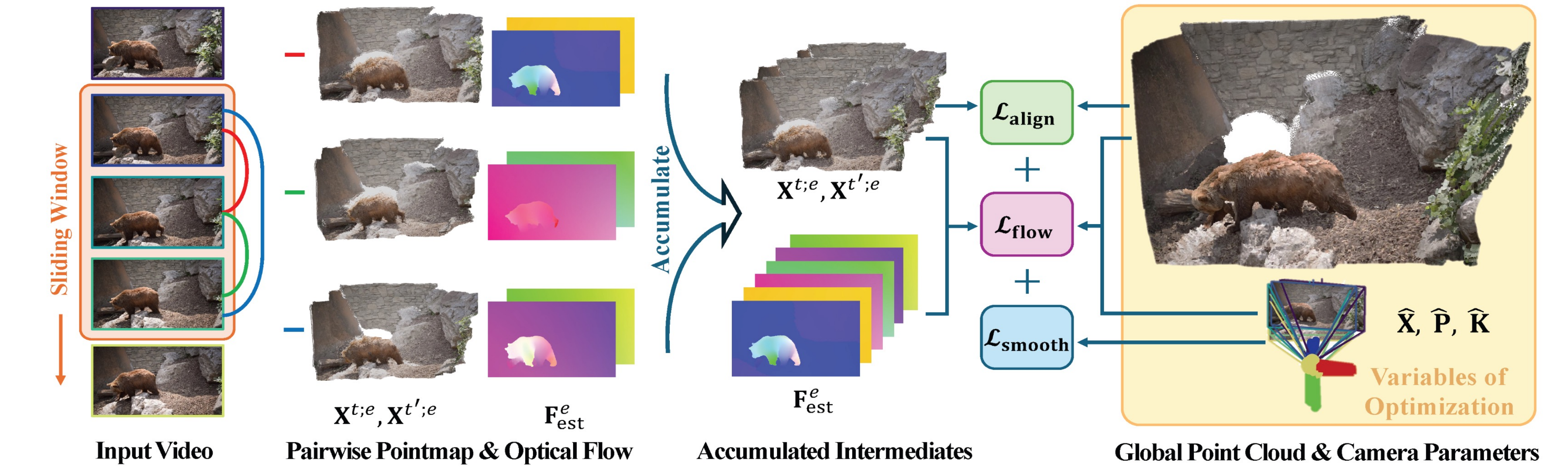
## Limitation of DUS<sub>t</sub>3R on Dynamic Scenes



As this is mainly a **data issue**, we propose a simple approach to adapt DUS<sub>t</sub>3R to dynamic scenes, by **fine-tuning** on a small set of dynamic videos, which surprisingly works well

## Dynamic Global Point Cloud

for video input, aggregate pairwise results to build global point cloud with global alignment



## Quantitative & Qualitative results

Table 1: Video depth evaluation

Alignment	Category	Method	Sintel		Bonn		KITTI	
			Abs Rel ↓ $\delta < 1.25$ ↑	Abs Rel ↓ $\delta < 1.25$ ↑	Abs Rel ↓ $\delta < 1.25$ ↑	Abs Rel ↓ $\delta < 1.25$ ↑	Abs Rel ↓ $\delta < 1.25$ ↑	Abs Rel ↓ $\delta < 1.25$ ↑
Per-sequence scale & shift	Single-frame depth	Marigold	0.532	51.5	0.091	93.1	0.149	79.6
		Depth-Anything-V2	0.367	55.4	0.106	92.1	0.140	80.4
	Video depth	NVDS	0.408	48.3	0.167	76.6	0.253	58.8
		ChronoDepth	0.687	48.6	0.100	91.1	0.167	75.9
		DepthCrafter (Sep. 2024)	<b>0.292</b>	<b>69.7</b>	<b>0.075</b>	<b>97.1</b>	<b>0.110</b>	<b>88.1</b>
		Robust-CVD	0.703	47.8	-	-	-	-
Per-sequence scale	Joint video depth & pose	CasualSAM	0.387	54.7	0.169	73.7	0.246	62.2
		<b>MonST3R</b>	<b>0.335</b>	<b>58.5</b>	<b>0.063</b>	<b>96.4</b>	<b>0.104</b>	<b>89.5</b>
	Joint depth & pose	DepthCrafter (Sep. 2024)	0.692	53.5	0.217	57.6	0.141	81.8
		<b>MonST3R</b>	<b>0.345</b>	<b>56.2</b>	<b>0.065</b>	<b>96.3</b>	<b>0.106</b>	<b>89.3</b>

Table 2: Camera pose estimation

Category	Method	Sintel			TUM-dynamics			ScanNet (static)		
		ATE ↓	RPE trans ↓	RPE rot ↓	ATE ↓	RPE trans ↓	RPE rot ↓	ATE ↓	RPE trans ↓	RPE rot ↓
Pose only	DROID-SLAM*	0.175	0.084	1.912	-	-	-	-	-	-
	DPVO*	0.115	0.072	1.975	-	-	-	-	-	-
	ParticleSfM	0.129	<b>0.031</b>	<b>0.535</b>	-	-	-	0.136	0.023	0.836
	LEAP-VO*	<b>0.089</b>	0.066	1.250	<b>0.068</b>	<b>0.008</b>	<b>1.686</b>	<b>0.070</b>	<b>0.018</b>	<b>0.535</b>
	Robust-CVD	0.360	0.154	3.443	0.153	0.026	3.528	0.227	0.064	7.374
Joint depth & pose	CasualSAM	0.141	0.035	0.615	0.071	0.010	1.712	0.158	0.034	1.618
	DUS <sub>t</sub> 3R w/ mask <sup>†</sup>	0.417	0.250	5.796	0.083	0.017	3.567	<b>0.081</b>	0.028	0.784
	<b>MonST3R</b>	<b>0.108</b>	0.042	0.732	<b>0.063</b>	<b>0.009</b>	<b>1.217</b>	<b>0.068</b>	<b>0.017</b>	<b>0.545</b>

\* requires ground truth camera intrinsics as input, <sup>†</sup> unable to estimate the depth of foreground object.

