

Image as an IMU: Estimating Camera Motion from a Single Motion-Blurred Image

Jerred Chen
University of Oxford
Department of Computer Science
jerred.chen@cs.ox.ac.uk

Ronald Clark
University of Oxford
Department of Computer Science
ronald.clark@cs.ox.ac.uk

Abstract

In many robotics and VR/AR applications, fast camera motions cause a high level of motion blur, causing existing camera pose estimation methods to fail. In this work, we propose a novel framework that leverages motion blur as a rich cue for motion estimation rather than treating it as an unwanted artifact. Our approach works by predicting a dense motion flow field and a monocular depth map directly from a single motion-blurred image. We then recover the instantaneous camera velocity by solving a linear least squares problem under the small motion assumption. In essence, our method produces an IMU-like measurement that robustly captures fast and aggressive camera movements. To train our model, we construct a large-scale dataset with realistic synthetic motion blur derived from ScanNet++v2 and further refine our model by training end-to-end on real data using our fully differentiable pipeline. Extensive evaluations on real-world benchmarks demonstrate that our method achieves state-of-the-art angular and translational velocity estimates, outperforming current methods like MAST3R and COLMAP.

1. Introduction

Given a sequence of images, visual odometry (VO) and Structure from Motion (SfM) methods have advanced to the point where they can accurately estimate camera poses even under moderately challenging conditions. However, these methods assume that the camera remains largely stationary during exposure, allowing each image to be treated as a snapshot of the scene that can be matched against subsequent images to compute relative camera poses. This assumption, however, ignores the continuous nature of camera motion and becomes especially problematic during fast movements where motion blur can severely degrade accuracy. Conventional approaches either discard motion-blurred frames or use inertial measurement units (IMUs) to

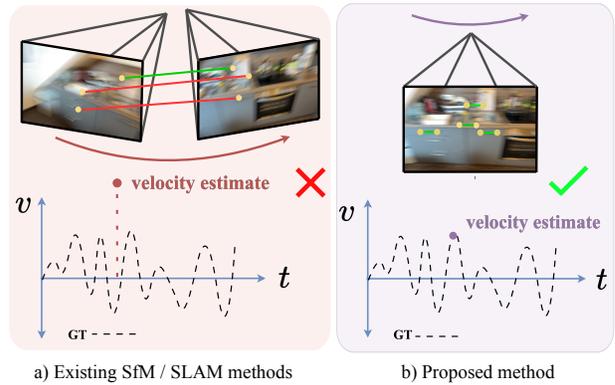


Figure 1. Existing methods rely on establishing correspondences between multiple frames to estimate inter-frame camera motion (a). This leads to failures during fast motion with lots of blur. We propose a method that can estimate intra-frame camera motion from a single image (b), making our method robust to aggressive motions.

improve robustness against blur. However, using IMUs introduces additional challenges such as sensor synchronization and drift. In this work, we therefore take a fundamentally different approach: rather than treating motion blur as an unwanted artifact, we exploit it as a rich source of information about the camera’s motion (see Figure 1).

Our key observation is that the extent and direction of motion blur provide direct cues about the camera’s motion during an exposure. Building on this idea, we propose a model that, using just the motion-blur in an image, can estimate the relative motion that the camera undergoes during the exposure. Our model works in two stages: first we predict the motion flow field and a monocular depth map, then we recover the relative pose of the camera by solving a linear least squares system. Notably, our method does not require any explicit deblurring process, and with a known exposure time, it yields an instantaneous rotational and translation velocity estimate that can be interpreted as an IMU-like measurement.

Since no existing dataset includes all the necessary data to train such a model, we construct our own dataset using a subset of ScanNet++v2 [47] and show how we obtain realistic motion blur that allows for generalization to in-the-wild images. Furthermore, the fully differentiable nature of our method enables us to refine our model using real-world motion-blurred images, even when ground-truth flows and depths are not available for training. We evaluate our model with angular and translational velocity estimates on real-world motion blurred images to test the generalization capabilities of our model. Our method runs in real-time at 30 FPS and obtains significantly more accurate velocity estimates compared to the existing state of the art methods.

In summary, our contributions are as follows:

1. A new approach for estimating camera motion that uses the blur present in a frame to provide extremely robust and accurate relative pose estimates
2. A pipeline for synthesizing the data needed to train the model (i.e. blurry images with ground-truth motion) from standard SLAM datasets
3. An evaluation of our method on real-world sequences, showing our method obtains state of the art results and is robust to aggressive camera motion

The rest of the paper is as follows. We first discuss relevant works to our method in Section 2. We then introduce preliminary information and provide a high-level intuition of our method in Section 3. In Section 4, we go into technical details about our method. Section 5 shows the dataset construction process. We evaluate our results against baseline methods in Section 6.

2. Related Works

Motion from Blur. This paper is most closely related to the line of work referred to as *motion from blur*, which aims to extract pixel flow information only from a single motion blurred image. [5] [26] are among the earliest works within this topic, where [5] introduced a linear model for computing the motion vectors and [26] independently analyzed the frequency domain to identify blurred regions. However, both methods assume that there is only linear, spatially-invariant blur in the image. [32] built on top of [26] by using the three color channels of an image. The seminal work [7] proposed a straightforward method where flow could be estimated using the alpha channel of an image. A limitation of these aforementioned methods is the inability to distinguish directional information with the motion vectors. [11] resolves this and introduces a variational approach to estimate the motion vectors and the blur segmentation in isolation. [1] proposed the first deep learned solution using a spatial transformer network. These methods only focus on estimating a dense flow field from the motion blurred image, whereas we are interested in the usefulness of the flow field for obtaining camera motion.

Our method also estimates the depth of the scene given the motion blurred image to retrieve the metric camera motion. Similarly [16], [20], and [25] demonstrate how to recover the 3D geometry even in motion-blurred scenes. [20] is the most relevant to our work and presented a visual odometry method which computes the homography during exposure time of a blurred image, which demonstrates increased robustness in severely blurred videos. However, they minimize the photometric error between a sharp keyframe and a blurred frame, where they use a off-the-shelf deblurring method to do so. We do not require any deblurring and only need a single image to estimate the motion.

Camera Pose Estimation. Camera pose estimation is a fundamental computer vision task used in 3D reconstruction which can be categorized into SfM and SLAM/VO methods. SfM methods operate on an unordered set of images and often involve estimate the camera intrinsics and extrinsics. SLAM/VO methods can be viewed as a subset of SfM with the added constraint for running in real-time, with and without loop closures respectively.

For SfM methods, COLMAP [31] [30] is an incremental-SfM method that is commonly the default choice for estimating camera parameters, poses, and the scene geometry. ORB-SLAM [3] and DSO [9], on the other hand, are classic real-time methods for obtaining camera poses and mapping the environment. While highly accurate in sequences with large overlap, these traditional systems succumb to failures under adversarial conditions, including few-view reconstruction and motion blur. These failures motivate deep-learning based approaches which are more robust in these difficult scenarios. Several fully-differentiable SfM pipelines have been proposed for greatly improving 3D reconstruction and pose estimation robustness [42] [34] [2]. In particular, DUST3R [44] and MAST3R [17] have shown that directly regressing scene geometry in the form of pointmaps enables incredible performance, which has inspired a line of “-3R”-related works [48] [41] [33] [46] [43].

Similarly, SLAM and VO methods have also underwent a deep-learning revolution and have similarly adopted learned modules with improved performance [6] [45] [37] [35] [28] [24]. For example, DROID-SLAM [37] is an end-to-end solution for obtaining accurate real-time reconstructions using RAFT [36] for estimating optical flow. These approaches rely on using two or more images for estimating the camera motion, whereas our method is able to use a single motion-blurred frame to estimate camera velocity. Due to there being no existing methods that estimate camera motion from a single frame, we opt to use these standard camera pose estimation methods, specifically MAST3R, DROID-SLAM and COLMAP, and estimate the velocities with a finite-difference approximation.

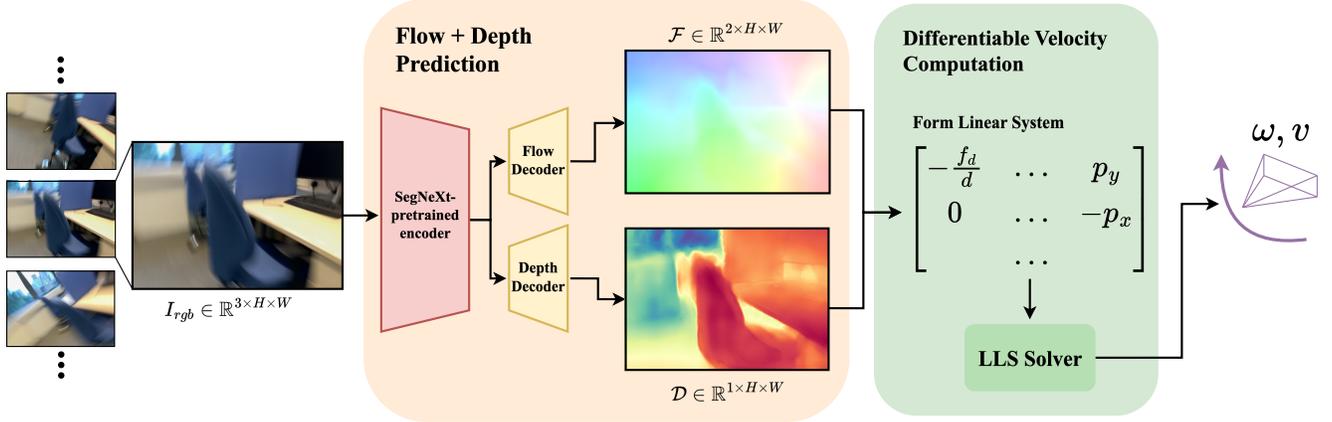


Figure 2. Overview of our method. Given a single motion blurred image, we pass it through the network to obtain a dense flow field and monocular depth prediction (Section 4.1). These are then formulated in a linear system, where the optimal velocity parameters are solved for using linear least squares (Section 4.2). Because the linear solver is fully differentiable, we can train the entire network end-to-end, supervised on the camera motion.

3. Preliminaries

Before introducing our method, we first revisit the image formation process. During the opening of the shutter, photons come into contact with the camera sensor and are collected to capture each pixel’s intensity. If the camera has moved during the exposure time, each sensor pixel receives photons from different parts of the scene and motion-blur is created. The length of the motion blur traces depends on the speed of the camera as well as the exposure time i.e. when the exposure time is lengthened, e.g. in low-illuminated settings, or if the camera moves faster relative to the scene, the blur traces are longer. More formally, let $I^B \in \mathbb{R}^{C \times H \times W}$ represent a blurry image, such that:

$$I^B = g \left(\frac{1}{\tau} \int_0^\tau I^\nu(t) dt \right) \quad (1)$$

where $I^\nu(t) \in \mathbb{R}^{C \times H \times W}$ represents the response of the sensor to incoming photons at timestep t , with the exposure time length denoted by τ , and $g(\cdot)$ representing the conversion from linear space to sRGB [27]. In practice, we approximate the creation of a blurry image in a discretized fashion:

$$I^B \approx g \left(\frac{1}{N} \sum_{i=1}^N I_i^\nu \right) \quad (2)$$

where we assume that $\{I_1^\nu, \dots, I_N^\nu\}$ are the N virtual images that existed in exposure time.

The intuition for our method is that each motion-blurred image can be seen as a composite of several “virtual images,” each capturing the camera at a different pose during the exposure. Each individual virtual image contributes a

part of the overall blur, and the blur traces therefore act as cues that represent the motion the camera has undergone during the exposure. In essence, these blur traces can be seen as providing virtual correspondences between the first and last virtual image.

We hypothesize that we can extract these “virtual” correspondences, revealing how each point in the scene moved across the image plane and capturing information similar to an optical flow field. Since we assume the scene is rigid, each virtual correspondence can be interpreted as a match between two distinct virtual images representing the camera at the beginning and end of the exposure period (see Figure 1).

Furthermore, by estimating a depth map from the blurred image, we can combine the pixel motion information extracted from the motion blur traces with the scene geometry to solve for the relative pose change from the start to the end of the exposure.

4. Method

An overview of our method can be seen in Figure 2. Our approach consists of two stages, the first stage takes as input a single frame and predicts a flow-field representing the motion the camera underwent during the exposure, and a monocular depth map. The second stage uses the depth and flow to solve for the instantaneous velocity of the camera using a differentiable least squares solver.

4.1. Flow and Depth Prediction

We follow GeoCalib [40] and use SegNeXt [12] as the backbone for our network. The blurry image I^B is passed through a shared SegNeXt encoder that maps to separate decoders, outputting the pixel-wise flow field $\mathcal{F} \in \mathbb{R}^{2 \times H \times W}$

and depth map $\mathcal{D} \in \mathbb{R}^{1 \times H \times W}$. Each flow vector $\mathbf{F} = [F_x, F_y]^\top \in \mathcal{F}$ is defined to be the pixel displacement from a given pixel location in the original virtual frame I_1^ν :

$$\mathbf{F} = \mathbf{p}'_2 - \mathbf{p}_1 \quad (3)$$

such that $\mathbf{p} = [p_x, p_y]^\top$, and $\mathbf{p}_1, \mathbf{p}'_2$ are the corresponding pixel coordinates from the first and last virtual frames I_1^ν and I_N^ν respectively.

We impose an L1 loss for both the flow field and the depth predictions:

$$\mathcal{L}_1 = \lambda_F \|\mathcal{F} - h_f(\hat{\mathcal{F}}_{fw}, \hat{\mathcal{F}}_{bw})\| + \lambda_D \|\mathcal{D} - \hat{\mathcal{D}}\| \quad (4)$$

where $\hat{\cdot}$ represents the ground truth labels, λ_F, λ_D are weights to balance the two losses, and h_f is the reorientation function defined as:

$$h_f(\hat{\mathcal{F}}_{fw}, \hat{\mathcal{F}}_{bw}; \mathcal{F}) = \begin{cases} \hat{\mathcal{F}}_{fw} & \text{if } \langle \hat{\mathcal{F}}_{fw}, \mathcal{F} \rangle > \langle \hat{\mathcal{F}}_{bw}, \mathcal{F} \rangle \\ \hat{\mathcal{F}}_{bw} & \text{otherwise} \end{cases} \quad (5)$$

such that $\langle \cdot, \cdot \rangle$ is the Frobenius inner product, and $\hat{\mathcal{F}}_{fw}, \hat{\mathcal{F}}_{bw}$ are the corresponding flow fields in the I_1^ν to I_N^ν direction and the I_N^ν to I_1^ν direction, respectively. Because determining the flow and depth from the true start virtual frame is ill-posed, we incorporate our reorientation function h_f so that the label $\hat{\mathcal{F}}$ is closest to the predicted direction of \mathcal{F} . This stabilizes the training and enables the model to produce globally consistent outputs.

4.2. Differentiable Velocity Computation

We can now use \mathcal{F} and \mathcal{D} to estimate the relative translation and rotation parameters $\mathbf{t} = [t_x, t_y, t_z]^\top \in \mathbb{R}^3$ and $\boldsymbol{\theta} = [\theta_x, \theta_y, \theta_z]^\top \in \mathbb{R}^3$ across exposure. To do so, we adapt the motion field equations described in Trucco and Verri [38] (see Supplementary 1), which we express as:

$$F_x = \frac{t_z p_x - t_x f}{d} - \theta_y f + \theta_z p_y + \frac{\theta_x p_x p_y}{f} - \frac{\theta_y (p_x)^2}{f} \quad (6)$$

$$F_y = \frac{t_z p_y - t_y f}{d} + \theta_x f - \theta_z p_x - \frac{\theta_y p_x p_y}{f} + \frac{\theta_x (p_y)^2}{f}. \quad (7)$$

such that $d \in \mathcal{D}$ and f is the known focal length.

These equations can be expressed in matrix form as $\mathbf{A}\mathbf{x} = \mathbf{b}$, where:

$$\mathbf{A} = \begin{bmatrix} -\frac{f}{d} & 0 & \frac{p_x}{d} & \frac{p_x p_y}{f} & -\frac{(p_x)^2 + f^2}{f} & p_y \\ 0 & -\frac{f}{d} & \frac{p_y}{d} & \frac{(p_y)^2 + f^2}{f} & -\frac{p_x p_y}{f} & -p_x \end{bmatrix}, \quad (8)$$

$$\mathbf{x} = \begin{bmatrix} t_x \\ t_y \\ t_z \\ \theta_x \\ \theta_y \\ \theta_z \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} F_x \\ F_y \end{bmatrix}. \quad (9)$$

With multiple flow vectors, the system becomes over-determined and is solved using the least squares method:

$$\mathbf{x} = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{b}. \quad (10)$$

This approach leverages the relationship between pixel displacements and camera pose, enabling closed-form estimation of the relative pose parameters \mathbf{t} and $\boldsymbol{\theta}$ from the predicted flow and depths. With a provided exposure time, we can then compute the instantaneous velocities \mathbf{v} and $\boldsymbol{\omega}$ across exposure.

A convenient property of this linear least squares formulation is that this is differentiable, which allows us to train the network fully end-to-end with clear pose supervision. Therefore, our final end-to-end loss is:

$$\mathcal{L}_2 = \lambda_R \|R - h_p(\hat{R})\|_2 + \lambda_t \|\mathbf{t} - h_p(\hat{\mathbf{t}})\|_2 + \mathcal{L}_1 \quad (11)$$

where we compute MSE on the rotations $R \in SO(3)$ and translations \mathbf{t} . λ_R and λ_t help balance the loss, and h_p , similar to Equation 5, reorients the \hat{R} or $\hat{\mathbf{t}}$ in such a way that is closest to the prediction.

4.3. Direction Disambiguation

While the method we have proposed so far can estimate the magnitude of the camera velocity directly, the direction has a 180° ambiguity because the blurring process erases temporal cues. For instance, both leftward or rightward motion of the camera will result in the same horizontal blur traces. To address this issue, we introduce a simple strategy to resolve the directional ambiguity. Recall τ is the exposure time length and let Δt_f be the time gap between the start of frame shutter. For frame I_i , we linearly extrapolate the flow field \mathcal{F}_i and apply the warping function Φ to I_i :

$$\mathcal{F}'_i = \frac{\Delta t_f}{\tau} \mathcal{F}_i \quad (12)$$

$$I'_i = \Phi(I_i; \mathcal{F}'_i) \quad (13)$$

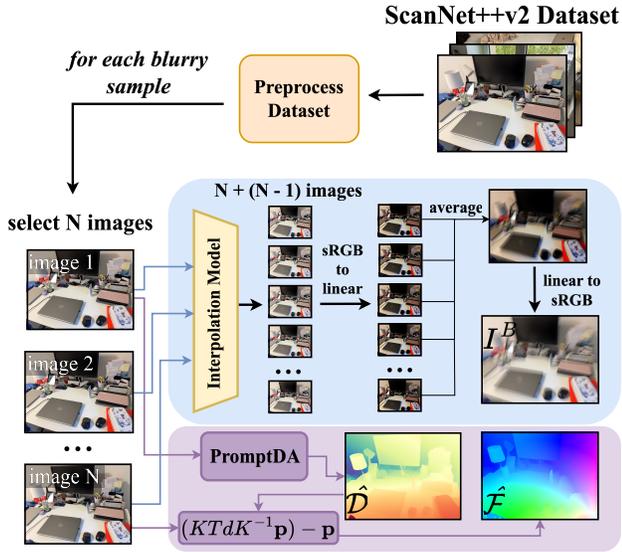


Figure 3. Overview for our synthetic dataset generation process. After preprocessing the dataset, we run selected frames through an interpolation network, which we use to synthesize our blurred image. We also take the first and last virtual frames to generate $\hat{\mathcal{D}}$, which is subsequently used for computing $\hat{\mathcal{F}}$.

For both the forward and backwards flow fields \mathcal{F}_{fw} and \mathcal{F}_{bw} , we compute the photometric error between the corresponding warped image I'_i with the next true frame in the video sequence I_{i+1} :

$$\mathcal{P}(I_1, I_2) = \frac{1}{HW} \sum_{u=0}^H \sum_{v=0}^W |I_1(u, v) - I_2(u, v)|. \quad (14)$$

We additionally do this photometric error check with warping \mathcal{F}_i to I_{i-1} and compute the sum of the photometric error between the two directions, and reorient based on the smaller photometric error.

$$e_{fw} = \mathcal{P}(I_{i+1}, I'_{i, fw}) + \mathcal{P}(I_{i-1}, I'_{i, bw}) \quad (15)$$

$$e_{bw} = \mathcal{P}(I_{i+1}, I'_{i, bw}) + \mathcal{P}(I_{i-1}, I'_{i, fw})$$

$$\omega, \mathbf{v} = \begin{cases} \omega_{fw}, \mathbf{v}_{fw} & \text{if } e_{fw} < e_{bw} \\ \omega_{bw}, \mathbf{v}_{bw} & \text{otherwise.} \end{cases} \quad (16)$$

This straightforward photometric error check enables us to reorient the direction of the camera velocities in the correct orientation.

5. Dataset Curation

5.1. Synthetic Dataset

To train our model, we require a dataset of blurred images with corresponding ground-truth flows, depths, and relative

poses across exposure time. Since no public dataset meets these requirements, we create our own by adapting a subset of ScanNet++v2 [47]. Our data pipeline is shown in Figure 3.

The ScanNet++v2 dataset consists of about 1000 iPhone videos of various indoor scenes with COLMAP [30] [31] and ARKit-captured poses. Because we are interested in obtaining relative poses that are across very short time spans (fractions of a second) and robust to fast-motion, we opt to use the ARKit poses to compute our ground truth motion. Before curating the dataset, we first preprocess the ARKit poses and filter out any ‘‘jumps’’ in the localization.

Synthesizing blur. We now detail how we select our virtual frames for creating a motion blurred image I^B . For a selected RGB image I'_i in a given video sequence, we obtain the following N images to have a collection of $N + 1$ images $\{I'_i, \dots, I'_{i+N}\}$. We then use an off-the-shelf frame interpolation model RIFE [13] to interpolate between each consecutive frame, resulting in a total of $N + (N - 1)$ virtual frames $\mathcal{I} = \{I'_i, I'_{i+\frac{1}{2}}, \dots, I'_{i+N}\}$. We follow Equation 2 and first convert \mathcal{I} in linear space $g^{-1}(\cdot)$ before averaging all images and converting back to sRGB space to create I^B . Rather than simply interpolating $N - 1$ virtual frames between two images, we intentionally interpolate between several real-world frames to 1) have more realistic motion trajectories to facilitate generalization for in-the-wild motion blur, and 2) to prevent RIFE from producing unreasonable blurred images. We curate images with $N = 3$ for small motion and $N = 10$ for large motion.

Obtaining depth and flows. We only use the first and last virtual frames I'_1 and I'_N for obtaining the ground-truth flow field. We use the low-resolution ARKit depth and I'_1 in PromptDA [18] to obtain our dense, high-resolution depth map $\hat{\mathcal{D}}$ as ground truth. For all pixels, we then perform:

$$\mathbf{p}' = K T d K^{-1} \mathbf{p} \quad (17)$$

which uses $d \in \hat{\mathcal{D}}$ to backproject all pixel points in I'_1 and subsequently project them into I'_N . Afterwards, we retrieve the flow field $\hat{\mathcal{F}}$ with the displacements of the projected points as expressed in Equation 3.

In total, we synthesize about 120k training samples and 1.2k validation samples across 150 ScanNet++v2 sequences, with corresponding blurred images, depth maps and flow fields.

5.2. Real-world Dataset

As our method, including the velocity solver, is fully-differentiable we can train it in an end-to-end manner. This enables us to additionally train on real-world motion blur

Method	billiards	commonroom	dining	office	avg	
	$\omega_x / \omega_y / \omega_z$					
MI	COLMAP [30] [31]	×	×	×	×	—
	↳ D+LG [39] [19]	<u>2.32 / 1.94 / 1.50</u>	<u>1.19 / 1.61 / 0.91</u>	×	<u>3.06 / 2.51 / 3.93</u>	—
	MASt3R [44] [17]	5.30 / 2.85 / 4.45	3.70 / 3.75 / 3.26	<u>2.36 / 0.84 / 1.67</u>	4.78 / 3.03 / 6.21	4.04 / <u>2.62</u> / 3.90
	DROID-SLAM [37]	5.39 / 3.33 / 5.31	3.01 / 5.89 / 3.57	2.92 / 1.20 / 1.98	6.33 / 4.90 / 5.56	4.41 / 3.83 / 4.10
SI	Ours	1.31 / 0.87 / 1.60	0.93 / 0.88 / 1.04	0.87 / 0.50 / 1.33	1.76 / 1.38 / 3.08	1.22 / 0.91 / 1.76
	Zero-Velocity baseline	5.39 / 3.43 / 5.16	3.95 / 4.50 / 2.81	4.58 / 1.53 / 3.66	5.43 / 3.19 / 6.99	4.84 / 3.16 / 4.66

Table 1. RMSE for rotational velocities across each axis, in rad/s. We evaluate against multi-image (MI) and single-image (SI) methods. The best and second-best results are **bolded** and underlined, respectively. The × represents a failure to reconstruct the poses.

Method	billiards	commonroom	dining	office	avg	
	$v_x / v_y / v_z$	$v_x / v_y / v_z$	$v_x / v_y / v_z$	$v_x / v_y / v_z$	$v_x / v_y / v_z$	
MI	COLMAP [30] [31]	×	×	×	×	—
	↳ D+LG [39] [19]	3.11 / 2.06 / 2.14	1.70 / 2.09 / 2.20	×	1.81 / 2.45 / 1.52	—
	MASt3R [44] [17]	<u>2.59 / 1.50 / 3.52</u>	<u>1.28 / 1.94 / 2.49</u>	0.82 / 0.55 / 1.22	<u>1.72 / 2.16 / 1.46</u>	<u>1.60 / 1.54 / 2.17</u>
	DROID-SLAM [37]	3.28 / 2.00 / <u>2.65</u>	2.23 / <u>1.47 / 1.25</u>	1.44 / 0.65 / <u>0.91</u>	1.98 / 2.40 / <u>1.23</u>	2.23 / 1.63 / <u>1.51</u>
SI	Ours	1.36 / 1.17 / 1.05	1.00 / 0.80 / 0.71	<u>0.95 / 0.61 / 0.81</u>	1.12 / 1.52 / 1.12	1.11 / 1.03 / 0.92
	Zero-Velocity baseline	2.80 / 1.66 / 1.41	1.98 / 1.48 / 1.26	1.30 / 1.05 / 1.20	1.94 / 2.26 / 1.07	2.01 / 1.61 / 1.24

Table 2. RMSE for translational velocities across each axis, in m/s. We evaluate both multi-image (MI) and single-image (SI) methods. The best and second-best results are **bolded** and underlined, respectively. The × represents a failure to reconstruct the poses.

to further close the reality-gap between our synthetically blurred dataset and in-the-wild images. We therefore collect 10k real-world motion-blurred images with corresponding ARKit poses, IMU measurements, and exposure times across about 30 scenes. We follow Liu et al. [21] and perform the Fast Fourier Transform to identify blurred images, where images below a specified threshold are blurry enough to be used for training or validation. We also supplement with a small amount of motionless images (with no blur) to enhance the robustness to still frames.

6. Experiments

Training. Our model is trained in three stages. We first train the flow and depth decoders without pose supervision on our synthetically blurred dataset with a batch size of 32. Pose supervision is added once the depths and flow estimates are reasonable, and train poses, flow and depths with a batch size of 8 for 300k steps. Finally, we finetune the model on the real-world collected images to close the reality gap from synthetic images for 10k steps. We train our method using the Adam optimizer [15] on a single Nvidia RTX 3090.

Evaluation. To assess the accuracy of our method, we evaluate the *rotational velocities* ω and *translational velocities* \mathbf{v} in four real-world videos (see Supplementary 2). For each video sequence, we get the ARKit-estimated camera poses, IMU measurements, and the exposure times for each image. We directly compare ω predictions to the gyroscope measurements and \mathbf{v} predictions to a finite-difference velocity approximation from the ARKit translations. We treat the ARKit translations as the ground-truth as their visual-inertial method fuses additional sensor information, which helps mitigate the poor visual input. We measure the RMSE across each axis for ω and \mathbf{v} .

Baselines. There are no existing methods that can estimate camera motion from a single blurred image, so we instead use camera pose estimation methods and estimate the relative pose (and hence velocity) between frames using a centered finite-difference approximation. We evaluate against the following SfM and SLAM baselines:

- *COLMAP* [30] [31] is an incremental-SfM method that is considered to be the de-facto method for obtaining camera poses and commonly serves as ground-truth for NeRFs [23] and Gaussian Splatting [14] methods. Because standard SfM is unreliable with small triangulation angles, we run COLMAP using all images in the se-

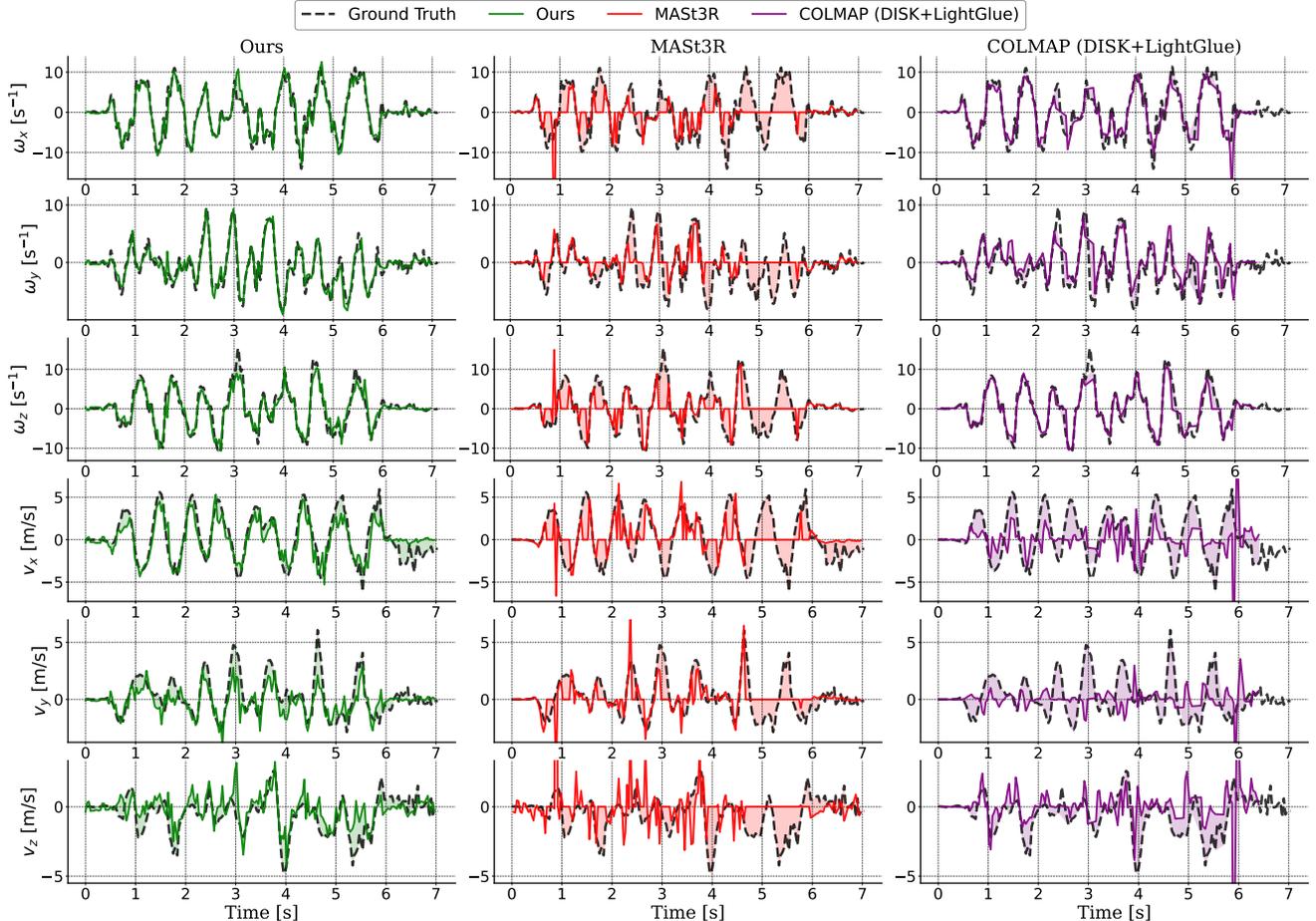


Figure 4. Visualization of the predicted velocities for the billiards sequence using our method, MAST3R and COLMAP (w/ DISK+LightGlue). The shaded area under the curve shows the error between the predicted velocity and GT velocity. Our translations and rotations are significantly better than MAST3R. While COLMAP with DISK + LightGlue feature matching does well on rotations, our method significantly outperforms it on translations.

quence. After obtaining the camera poses across all images, we scale the trajectory with the ground-truth ARKit poses. We use the `hloc` package [29] to run COLMAP with SIFT [22] and AdaLAM [4] as well as a learning baseline using DISK [39] and LightGlue [19].

- *MASt3R* [44] [17] is a recent learning-based 3D vision method that regresses the pointmaps between a pair of images. MAST3R is the current SOTA in retrieving metric-scale camera poses, particularly with extremely challenging views with little-to-no-overlap. To better represent our method, we evaluate MAST3R on pairwise images two frames apart.
- *DROID-SLAM* [37] is a learning-based SLAM method that utilizes RAFT [36] to estimate the optical flow across frames to help compute the camera poses, and obtains SOTA accuracy and robustness particularly in adversarial conditions. To obtain the camera poses for each frame, we

set DROID-SLAM to identify each image as a keyframe by setting the filter and keyframe thresholds to 0.

Note that all baselines require two or more images to estimate the velocity, whereas our method *only needs a single image*. We therefore tackle a more difficult problem, and our formulation also more closely resembles the true camera state.

Results. We report the rotational and translational RMSE in Tables 1 and 2. While the recorded video sequences suffer from severe motion blur, our method is able to handle the motion very robustly. Figure 5 highlight the overall robustness of our method, with the error CDFs demonstrating that our method suffers from noticeably less errors than the baselines. In these sequences, the baseline methods suffer from major outliers and zero-prediction failures, which can be observed in Figure 4. Velocity predictions from the

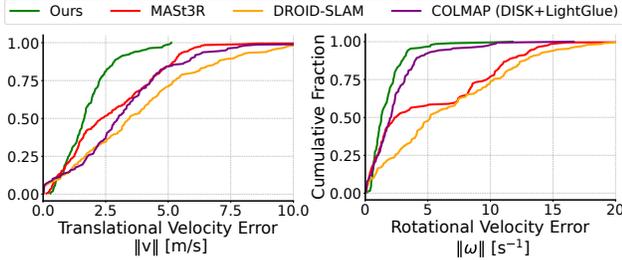


Figure 5. Error CDFs for the billiards sequence, such that the left plot shows the distribution of translational error in the sequence and the right plot the rotational error. Curves closer to the top-left corner are better.

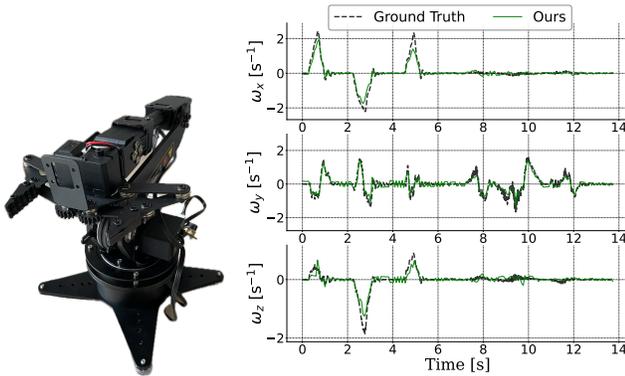


Figure 6. Real-world application example of our method using a camera attached to a fast-moving RoArm-M1 robot arm platform. (a) The robot arm used for recording. (b) The predicted and GT velocity time series for the camera attached to the end-effector.

billiards sequence are shown with the remaining sequences included in Supplementary 3. Overall, we obtain an average of 31% reduction from the second-best result in rotational velocities and a 24% reduction in translational velocities.

To further showcase the usefulness of our method, we test our method in estimating the velocity of the gripper on a moving robot arm. We place the camera on the robot pictured in Figure 6(a) and predict the velocity of the arm solely from the images in the recorded video. The motion-blurred images give a salient enough signal to accurately estimate the movement of the arm as seen in Figure 6(b).

In Figure 7 we show a comparison between our method and an actual IMU. To do this test we first estimate the IMU bias from a stationary calibration video, followed by estimating the initial pose of the camera to properly align with the gravity vector. We then preintegrate the IMU measurements following [10] using GTSAM [8] and compare the integrated velocities to our predictions. Figure 7 shows the drift that occurs in v_x and v_z after only 20 seconds, whereas our method is reliably drift-free for the entire video. Therefore, while our method produces inertial-like

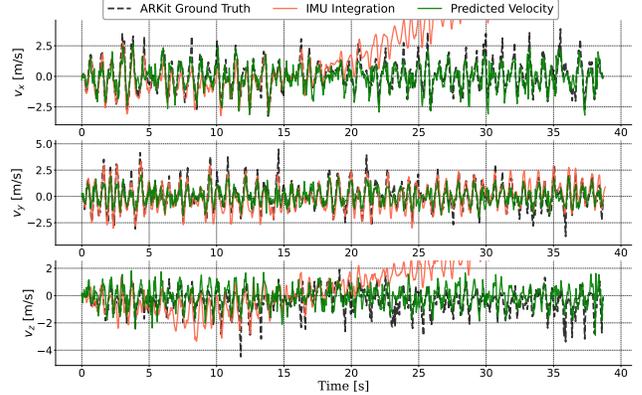


Figure 7. Comparison of our method to using IMU integration. The IMU velocity estimate is accurate for a few seconds until it drifts. Our method provides accurate and drift-free estimates throughout the sequence.

Method	Runtime (s)	Realtime factor
COLMAP (D+LG)	182.98	0.038×
MASt3R	67.90	0.103×
DROID-SLAM	41.61	0.169×
Ours	6.48	1.08×

Table 3. Runtimes of all methods on the billiards sequence showing the time taken to process the sequence and how much faster/slower the method is compared to realtime (original sequence is 7.08 s long). Our method is the only one that operates faster than realtime.

measurements, we see in this example that our method provides an *even better* motion estimate than an IMU.

We lastly measure the runtime of our method. Even when running on every frame in the sequence, we can run in real-time at 30 FPS, where the inference and direction disambiguation for each image takes on average 0.03 seconds on an RTX 3090 GPU.

7. Conclusion

In this paper we have proposed a method that can estimate the instantaneous camera velocity from individual motion-blurred images – producing inertial-like measurements without an IMU. Our evaluations indicate that not only do we obtain more accurate and robust results compared to an actual IMU and other vision-based solutions, but we also do so significantly more efficiently. We hope that this motivates the community to explore further how traditionally unwanted or overlooked cues can be used to enhance the performance of vision tasks.

References

- [1] Dawit Mureja Argaw, Junsik Kim, Francois Rameau, Jae Won Cho, and In So Kweon. Optical Flow Estimation from a Single Motion-blurred Image. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(2):891–900, 2021. 2
- [2] Eric Brachmann, Jamie Wynn, Shuai Chen, Tommaso Cavallari, Áron Monszpart, Daniyar Turmukhambetov, and Victor Adrian Prisacariu. Scene coordinate reconstruction: Posing of image collections via incremental learning of a relocalizer. In *ECCV*, 2024. 2
- [3] Carlos Campos, Richard Elvira, Juan J. Gómez Rodríguez, José M. M. Montiel, and Juan D. Tardós. Orb-slam3: An accurate open-source library for visual, visual–inertial, and multimap slam. *IEEE Transactions on Robotics*, 37(6):1874–1890, 2021. 2
- [4] Luca Cavalli, Viktor Larsson, Martin Ralf Oswald, Torsten Sattler, and Marc Pollefeys. Handcrafted outlier detection revisited. In *European Conference on Computer Vision*, 2020. 7
- [5] Wei-Ge Chen, N. Nandhakumar, and W.N. Martin. Image motion estimation from motion smear—a new computational model. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 18(4):412–425, 1996. 2
- [6] J Czarnowski, T Laidlow, R Clark, and AJ Davison. Deepfactors: Real-time probabilistic dense monocular slam. *IEEE Robotics and Automation Letters*, 5:721–728, 2020. 2
- [7] Shengyang Dai and Ying Wu. Motion from blur. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, 2008. ISSN: 1063-6919. 2
- [8] Frank Dellaert and GTSAM Contributors. borglab/gtsam, 2022. 8
- [9] J. Engel, V. Koltun, and D. Cremers. Direct sparse odometry. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018. 2
- [10] Christian Forster, Luca Carlone, Frank Dellaert, and Davide Scaramuzza. On-manifold preintegration for real-time visual–inertial odometry. *IEEE Transactions on Robotics*, 33(1):1–21, 2016. 8
- [11] Jochen Gast, Anita Sellent, and Stefan Roth. Parametric Object Motion from Blur. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1846–1854, Las Vegas, NV, 2016. IEEE. 2
- [12] Meng-Hao Guo, Cheng-Ze Lu, Qibin Hou, Zhengning Liu, Ming-Ming Cheng, and Shi-Min Hu. Segnext: Rethinking convolutional attention design for semantic segmentation. *arXiv preprint arXiv:2209.08575*, 2022. 3
- [13] Zhewei Huang, Tianyuan Zhang, Wen Heng, Boxin Shi, and Shuchang Zhou. Real-Time Intermediate Flow Estimation for Video Frame Interpolation. In *Computer Vision – ECCV 2022*, pages 624–642. Springer Nature Switzerland, Cham, 2022. Series Title: Lecture Notes in Computer Science. 5
- [14] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), 2023. 6
- [15] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. 6
- [16] Hee Seok Lee and Kuoung Mu Lee. Dense 3D Reconstruction from Severely Blurred Images Using a Single Moving Camera. In *2013 IEEE Conference on Computer Vision and Pattern Recognition*, pages 273–280, Portland, OR, USA, 2013. IEEE. 2
- [17] Vincent Leroy, Yohann Cabon, and Jerome Revaud. Grounding image matching in 3d with mast3r, 2024. 2, 6, 7
- [18] Haotong Lin, Sida Peng, Jingxiao Chen, Songyou Peng, Jiaming Sun, Minghuan Liu, Hujun Bao, Jiashi Feng, Xiaowei Zhou, and Bingyi Kang. Prompting depth anything for 4k resolution accurate metric depth estimation. 2024. 5
- [19] Philipp Lindenberger, Paul-Edouard Sarlin, and Marc Pollefeys. LightGlue: Local Feature Matching at Light Speed. In *ICCV*, 2023. 6, 7
- [20] Peidong Liu, Xingxing Zuo, Viktor Larsson, and Marc Pollefeys. MBA-VO: Motion Blur Aware Visual Odometry, 2021. arXiv:2103.13684 [cs]. 2
- [21] Renting Liu, Zhaorong Li, and Jiaya Jia. Image partial blur detection and classification. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8, 2008. 6
- [22] David G. Lowe. Distinctive image features from scale-invariant keypoints. *International Journal of Computer Vision*, 60:91–110, 2004. 7
- [23] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 6
- [24] Riku Murai, Eric Dexheimer, and Andrew J. Davison. Mast3r-slam: Real-time dense slam with 3d reconstruction priors, 2024. 2
- [25] Jiayan Qiu, Xinchao Wang, Stephen J. Maybank, and Dacheng Tao. World from blur. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8485–8496, 2019. 2
- [26] Ioannis M. Rekleitis. Optical flow recognition from the power spectrum of a single blurred image. In *International Conference in Image Processing*, Lausanne, Switzerland, 1996. IEEE Signal Processing Society. 2
- [27] Jaesung Rim, Geonung Kim, Jungeon Kim, Junyong Lee, Seungyong Lee, and Sunghyun Cho. Realistic blur synthesis for learning image deblurring. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2022. 3
- [28] Antoni Rosinol, Marcus Abate, Yun Chang, and Luca Carlone. Kimera: an open-source library for real-time metric-semantic localization and mapping. In *IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2020. 2
- [29] Paul-Edouard Sarlin, Cesar Cadena, Roland Siegwart, and Marcin Dymczyk. From coarse to fine: Robust hierarchical localization at large scale. In *CVPR*, 2019. 7
- [30] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 5, 6

- [31] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016. 2, 5, 6
- [32] Yasmina Schoueri, Milena Scaccia, and Ioannis Rekleitis. Optical Flow from Motion Blurred Color Images. In *2009 Canadian Conference on Computer and Robot Vision*, pages 1–7, 2009. 2
- [33] Brandon Smart, Chuanxia Zheng, Iro Laina, and Victor Adrian Prisacariu. Splatt3r: Zero-shot gaussian splatting from uncalibrated image pairs. 2024. 2
- [34] Cameron Smith, David Charatan, Ayush Tewari, and Vincent Sitzmann. Flowmap: High-quality camera poses, intrinsics, and depth via gradient descent. In *arXiv*, 2024. 2
- [35] Edgar Sucar, Shikun Liu, Joseph Ortiz, and Andrew J. Davison. imap: Implicit mapping and positioning in real-time. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6229–6238, 2021. 2
- [36] Zachary Teed and Jia Deng. RAFT: Recurrent All-Pairs Field Transforms for Optical Flow. In *Computer Vision – ECCV 2020*, pages 402–419. Springer International Publishing, Cham, 2020. Series Title: Lecture Notes in Computer Science. 2, 7
- [37] Zachary Teed and Jia Deng. DROID-SLAM: Deep Visual SLAM for Monocular, Stereo, and RGB-D Cameras. *Advances in neural information processing systems*, 2021. 2, 6, 7
- [38] Emanuele Trucco and Alessandro Verri. *Introductory Techniques for 3-D Computer Vision*, pages 178–184. Prentice Hall PTR, USA, 1998. 4
- [39] Michał Tyszkiewicz, Pascal Fua, and Eduard Trulls. Disk: Learning local features with policy gradient. *Advances in Neural Information Processing Systems*, 33, 2020. 6, 7
- [40] Alexander Veicht, Paul-Edouard Sarlin, Philipp Lindenberger, and Marc Pollefeys. GeoCalib: Single-image Calibration with Geometric Optimization. In *ECCV*, 2024. 3
- [41] Hengyi Wang and Lourdes Agapito. 3d reconstruction with spatial memory. *arXiv preprint arXiv:2408.16061*, 2024. 2
- [42] Jianyuan Wang, Nikita Karaev, Christian Rupprecht, and David Novotny. Vggsfm: Visual geometry grounded deep structure from motion. 2023. 2
- [43] Qianqian Wang, Yifei Zhang, Aleksander Holynski, Alexei A. Efros, and Angjoo Kanazawa. Continuous 3d perception model with persistent state, 2025. 2
- [44] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vision made easy. In *CVPR*, 2024. 2, 6, 7
- [45] Wenshan Wang, Yaoyu Hu, and Sebastian Scherer. Tartanvo: A generalizable learning-based vo. 2020. 2
- [46] Jianing Yang, Alexander Sax, Kevin J. Liang, Mikael Henaff, Hao Tang, Ang Cao, Joyce Chai, Franziska Meier, and Matt Feiszli. Fast3r: Towards 3d reconstruction of 1000+ images in one forward pass. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2025. 2
- [47] Chandan Yeshwanth, Yueh-Cheng Liu, Matthias Nießner, and Angela Dai. Scannet++: A high-fidelity dataset of 3d indoor scenes. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2023. 2, 5
- [48] Junyi Zhang, Charles Herrmann, Junhwa Hur, Varun Jampani, Trevor Darrell, Forrester Cole, Deqing Sun, and Ming-Hsuan Yang. Monst3r: A simple approach for estimating geometry in the presence of motion. *arXiv preprint arxiv:2410.03825*, 2024. 2