

Fast3R: Towards 3D Reconstruction of 1000+ Images in One Forward Pass

Jianing Yang^{♣♦} Alexander Sax[♣] Kevin J. Liang[♣] Mikael Henaff[♣] Hao Tang[♣]
 Ang Cao^{♣♦} Joyce Chai[♦] Franziska Meier[♣] Matt Feiszli[♣]

♣ Meta ♦ University of Michigan
<https://fast3r-3d.github.io/>

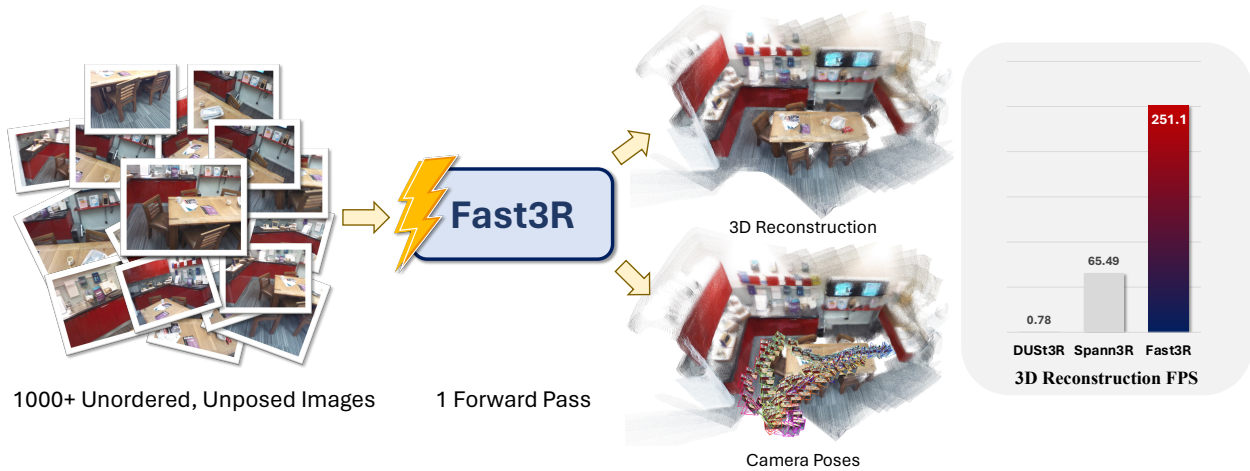


Figure 1. **Fast3R** is a method towards 3D reconstructing 1000+ unordered, unposed images in a single forward pass.

Abstract

*Multi-view 3D reconstruction remains a core challenge in computer vision, particularly in applications requiring accurate and scalable representations across diverse perspectives. Current leading methods such as DUST3R employ a fundamentally pairwise approach, processing images in pairs and necessitating costly global alignment procedures to reconstruct from multiple views. In this work, we propose Fast 3D Reconstruction (**Fast3R**), a novel multi-view generalization to DUST3R that achieves efficient and scalable 3D reconstruction by processing many views in parallel. Fast3R’s Transformer-based architecture forwards N images in a single forward pass, bypassing the need for iterative alignment. Through extensive experiments on camera pose estimation and 3D reconstruction, Fast3R demonstrates state-of-the-art performance, with significant improvements in inference speed and reduced error accumulation. These results establish Fast3R as a robust alternative for multi-view applications, offering enhanced scalability*

without compromising reconstruction accuracy.

1. Introduction

3D reconstruction from multiple views has long been a foundational task across applications in autonomous navigation, augmented reality, and robotics [31, 53]. Establishing correspondences across images, known as *multi-view matching*, is central to these applications and enables an accurate scene representation. Traditional reconstruction pipelines, such as those based on Structure-from-Motion (SfM) [44] and Multi-View Stereo (MVS) [18], fundamentally rely on image pairs to reconstruct 3D geometry. While effective in some settings, these methods require extensive engineering to manage the sequential stages of feature extraction, correspondence matching, triangulation, and global alignment, limiting scalability and speed.

This traditional “pipeline” paradigm has recently been challenged by DUST3R [61], which directly predicts 3D structure from RGB images. It achieves this with a design

that “cast[s] the pairwise reconstruction problem as a regression of pointmaps, relaxing the hard constraints of usual projective camera models” [61], yielding impressive robustness across challenging viewpoints. This represents a radical shift in 3D reconstruction, as an end-to-end learnable solution is less prone to pipeline error accumulation, while also being dramatically simpler.

On the other hand, a fundamental limitation of DUST3R is its restriction to two image inputs. While image pairs are an important use case, often one is interested in reconstructing from more than two views, as when scanning of objects [39] or scenes [4, 6, 20, 55, 67], *e.g.* for asset generation or mapping. To process more than two images, DUST3R computes $\mathcal{O}(N^2)$ pairs of pointmaps and performs a global alignment optimization procedure. This process can be computationally expensive, scaling poorly as the collection of images grows. For instance, it will lead to OOM with only 48 views on an A100 GPU.

Moreover, such a process is still fundamentally pairwise, which limits the model’s context, both affecting learning during training and ultimate accuracy during inference. In this sense, DUST3R suffers from the same pair-wise bottleneck as traditional SfM and MVS methods.

We propose **Fast3R**, a novel multi-view reconstruction framework designed to overcome these limitations. Building on DUST3R’s foundations, Fast3R leverages a Transformer-based architecture [56] that processes multiple images in parallel, allowing N images to be reconstructed in a single forward pass. By eliminating the need for sequential or pairwise processing, each frame can simultaneously attend to all other frames in the input set during reconstruction, significantly reduces error accumulation. Perhaps surprisingly, Fast3R also takes significantly less time.

Our contributions are threefold.

1. We introduce Fast3R, a Transformer-based model for multi-view pointmap estimation that obviates the need for global postprocessing; resulting in significant improvements in speed, computation overhead and scalability.
2. We show empirically that the model performance improves by scaling along the view axis. For camera pose localization and reconstruction tasks, the model improves when trained on progressively larger sets of views. Per-view accuracy further improves when more views are used during inference, and the model can generalize to significantly more views than seen during training.
3. We demonstrate state-of-the-art performance in camera pose estimation with significant inference time improvements. On CO3Dv2 [39], Fast3R gets **99.7%** accuracy within 15-degrees for pose estimation, over a **14x** error reduction compared to DUST3R *with* global alignment. Fast3R offers a scalable and accurate alternative for

real-world applications, setting a new standard for efficient multi-view 3D reconstruction.

2. Related Work

Multi-view 3D reconstruction: Almost all modern 3D reconstruction approaches are based on the traditional multi-view geometry (MVG) pipeline [21]. MVG-based methods first identify corresponding pixels between image pairs, and then use camera models and projective multiview geometry to lift these correspondences to 3D points. The process happens in sequential stages: feature extraction, finding pairwise image correspondences, triangulation to 3D and pairwise relative camera pose, and global bundle alignment. However, any pipeline approach is prone to accumulating errors, which are especially common in the hand-crafted components. Moreover, the sequential nature prevents parallelization, which limits speed and scalability.

MVG approaches have existed since the early days of computer vision, and are still in use for a reason: they can be highly accurate when they do not catastrophically fail. The latest multi-view geometry pipelines like COLMAP [44] or OrbSLAM2 [30] incorporate nearly 60 years of compounding engineering improvements, but these approaches still catastrophically fail >40% of the time on static scenes like ETH-3D [52]), which can actually be considered an easy case due to dense image coverage of the scene.

Much recent work has successfully addressed the robustness and speed by replacing increasingly large components of MVG pipelines with end-to-end learned versions that are faster and reduce the rate of catastrophic failures [48, 58, 72]. For example, [14, 19, 25, 42, 51, 68] improve feature extraction and correspondences, [27, 50, 59, 71] learn to estimate camera pose, and [52] introduce a bundle adjustment layer. [61] contains an excellent and comprehensive survey of such efforts. Overall, the trend is towards replacing increasingly large components with end-to-end solutions.

Pointmap representation: DUST3R [61] takes this evolution the furthest by proposing *pointmap* regression to replace everything in the MVG pipeline up to global pairwise alignment. Rather than first attempting to solve for camera parameters in order to triangulate corresponding pixels, DUST3R trains a model to directly predict 3D pointmaps for pairs of images in a shared coordinate system. Other MVG component tasks such as relative camera pose estimation and depth estimation can be recovered from the resulting pointmap representation. However, DUST3R’s pairwise assumption is a limitation, as it requires inference on $\mathcal{O}(N^2)$ image pairs and then a global alignment optimization, which is per-scene and does not improve with more data. Moreover, this process quickly becomes slow or crashes due to exceeded system memory, even for relatively modest numbers of images.

DUST3R has inspired several follow-ups. MAST3R [25]

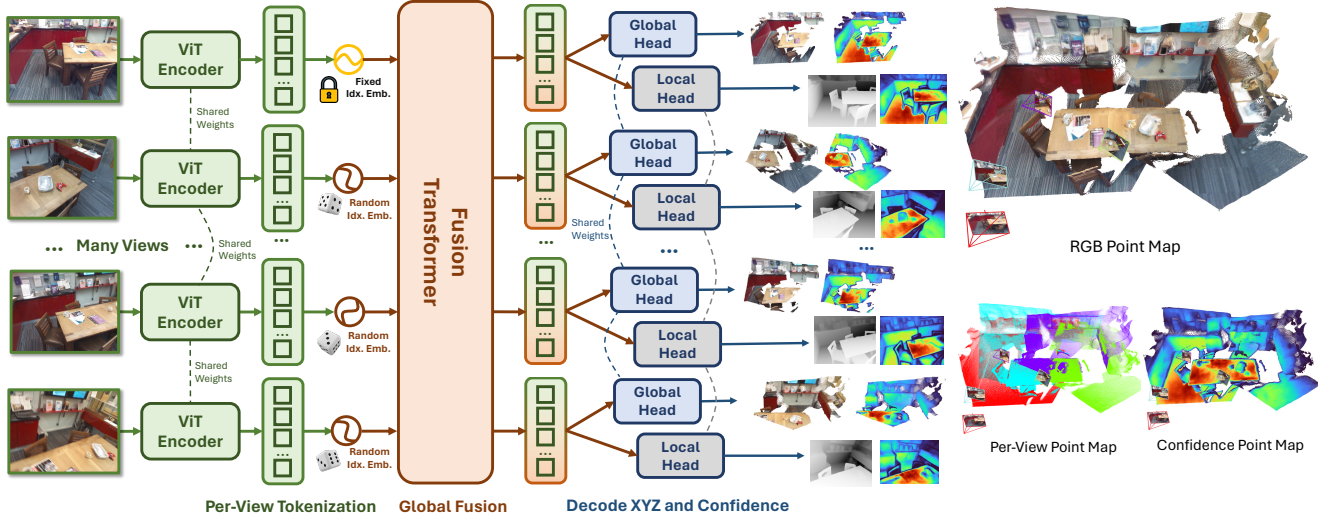


Figure 2. **Model architecture of Fast3R.** Built upon a novel Transformer-based architecture which supports bidirectional information flow, Fast3R is able to process dense input views simultaneously.

adds a local feature head to each decoder’s output, while MonST3R [69] does a data-driven exploration of dynamic scenes, but both are still fundamentally pairwise methods. MAST3R in particular does not make any changes to the global alignment methodology. Concurrently with our work, Spann3R [57] treats images as an ordered sequence (*e.g.* from a video) and incrementally reconstructs a scene using a pairwise sliding window network, along with a learned spatial memory system. This extends DUST3R to handle more images, but Spann3R’s incremental pairwise processing cannot fix reconstructions from earlier frames, which can cause errors to accumulate. Crucially, Fast3R’s transformer architecture uses all-to-all attention, allowing the model to reason simultaneously and jointly over all frames without any assumption of image order. Fast3R removes sequential dependencies, enabling parallelized inference across many devices in a single forward pass.

3. Model

Fast3R is a transformer-based model that predicts a 3D pointmap from a *set* of unordered and unposed images. The model architecture is designed to be scalable to over 1000 images during inference, though during training we use image masking to train it with far fewer. In this section, we detail our implementation of Fast3R, and discuss the design choices that enable its scalability.

3.1. Problem definition

Taking a set of (N) unordered and unposed RGB images $\mathbf{I} \in \mathbb{R}^{N \times H \times W \times 3}$ as inputs¹, Fast3R reconstructs the 3D structures of the scene by predicting the corresponding

¹We assume all images are resized to the same $H \times W$ for simplicity.

pointmap \mathbf{X} , where $\mathbf{X} \in \mathbb{R}^{N \times H \times W \times 3}$. A *pointmap* is a set of 3D locations indexed by pixels in an image \mathbf{I} , enabling the derivation of camera poses, depths, and 3D structures.

Fast3R predicts local and global pointmaps \mathbf{X}_L and \mathbf{X}_G , and corresponding confidence maps Σ_L and Σ_G (of shape $\Sigma \in \mathbb{R}^{N \times H \times W}$). Fast3R maps N RGB images to

$$\text{Fast3R} : \mathbf{I} \rightarrow (\mathbf{X}_L, \Sigma_L, \mathbf{X}_G, \Sigma_G)$$

Like MAST3R, the global pointmap \mathbf{X}_G is in the coordinate frame of the first camera and the \mathbf{X}_L is in the coordinate frame of the viewing camera, as shown in Figure 2

3.2. Training Objective

This section describes the loss, using the notation in Section 3.1 above. Fast3R’s predictions of $(\hat{\mathbf{X}}_L, \hat{\Sigma}_L, \hat{\mathbf{X}}_G, \hat{\Sigma}_G)$ are trained using generalized versions of the pointmap loss in DUST3R [61].

Our total loss is the combination of pointmap losses for the local and global pointmaps:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\mathbf{X}_G} + \mathcal{L}_{\mathbf{X}_L} \quad (1)$$

which are confidence-weighted versions of the normalized 3D pointwise regression loss.

Normalized 3D pointwise regression loss: The normalized regression loss for \mathbf{X} is a multi-view version of that in DUST3R [66] or monocular depth estimation [15, 36, 66]. It is the L_2 loss between the normalized predicted pointmaps and normalized target pointmaps, rescaled by the mean Euclidean distance to the origin:

$$\ell_{\text{regr}}(\hat{\mathbf{X}}, \mathbf{X}) = \left\| \frac{1}{\hat{z}} \hat{\mathbf{X}} - \frac{1}{z} \mathbf{X} \right\|_2, \quad z = \frac{1}{|\mathbf{X}|} \sum_{x \in \mathbf{X}} \|x\|_2 \quad (2)$$



Figure 3. **Qualitative examples of Fast3R’s output.** The text on the yellow sign says “Caution, cleaning in progress” and is legible if zoomed in.

Note that the predictions and targets are independently normalized by the mean euclidean distance to the origin.

Pointmap loss: As in [61], we use a confidence-adjusted version of the loss above, using the confidence score $\hat{\Sigma}$ predicted by the model. The total loss for a pointmap is

$$\mathcal{L}_{\mathbf{X}}(\hat{\Sigma}, \hat{\mathbf{X}}, \mathbf{X}) = \frac{1}{|\mathbf{X}|} \sum \hat{\Sigma}_+ \cdot \ell_{\text{regr}}(\hat{\mathbf{X}}, \mathbf{X}) + \alpha \log(\hat{\Sigma}_+) \quad (3)$$

Since the log term requires the confidence scores to be positive, we enforce $\hat{\Sigma}_+ = 1 + \exp(\hat{\Sigma})$. Our intuition is that the confidence weighting helps the model deal with label noise. Like DUST3R, we train on real-world scans typically containing systematic errors in the underlying pointmap labels. For example, glass or thin structures are often not reconstructed properly in the ground-truth laser scans [4, 67], and errors in camera registration will cause misalignments between the images and pointmap labels [66].

3.3. Model architecture

The Fast3R meta-architecture is inspired by DUST3R, and has three components: *image encoding*, *fusion transformer*, and *pointmap decoding*, as shown in Figure 2. We emphasize that Fast3R makes no assumptions on the ordering of images in \mathbf{I} , and all output pointmaps and confidence maps $(\mathbf{X}_L, \Sigma_L, \mathbf{X}_G, \Sigma_G)$ are predicted simultaneously, not sequentially.

Image encoder: Fast3R encodes each image $I_i \in \mathbf{I}$ to a set of patch features H_i , using a feature extractor \mathcal{F} . This is done independently per image, yielding a sequence of image patch features $H_i = \{h_{i,j}\}_{j=1}^{HW/P^2}$ for each image:

$$H_i = \mathcal{F}(I_i), i \in 1, \dots, N \quad (4)$$

We follow DUST3R’s design and use CroCo ViT [63] as our encoder, though we found DINOv2 [33] works similarly.

Before passing image patch features \mathbf{H} to the fusion transformer, we add position embeddings with one-dimensional *image index positional embeddings*.

Index embeddings help the fusion transformer determine which patches come from the same image and are the mechanism for identifying I_1 , which importantly defines the global coordinate frame. This is critical for allowing the model to implicitly reason about camera pose jointly for all images from an otherwise permutationally invariant set of tokens.

Fusion transformer: Most of the computation in Fast3R happens in the fusion transformer, which has a generic architecture. We use a 12-layer transformer similar to ViT-B [12] or BERT [10], however this could be scaled up. This fusion transformer takes the concatenated encoded image patches from *all* views and performs all-to-all self-attention. This operation provides Fast3R with full context, beyond the information provided in pairs alone.

Pointmap head: Finally, Fast3R uses separate DPT-L [37] decoder heads to map these tokens to the local and global pointmaps $(\mathbf{X}_L, \mathbf{X}_G)$, and confidence maps (Σ_L, Σ_G) .

Image index positional embedding generalization: We would like Fast3R to be able to handle many views at inference, more than were used to train a model. A naïve way to embed views during testing would be to embed them in the same way as training: i.e. use the same Spherical Harmonic frequencies [49] to embed raw indices $SH(\{1, \dots, N\})$ during training, and $SH(\{1, \dots, N_{\text{test}}\})$ during inference. In LLMs this causes poor performance, and in preliminary experiments we also found that the resulting model did not work well when the number of input images exceeded that

used during training. We therefore adapt Position Interpolation [5], a solution from LLMs, where during training we randomly draw N indexes from a larger pool N' of possible samples. [5] draws samples using a regular grid since the LLM inputs form a regular ordered sequence. Our images are unordered, so we draw $N \subset \{1, \dots, N'\}$ uniformly at random. To the transformer, the strategy looks indistinguishable from masking out images, and $N' \gg N$ controls the masking ratio.² This strategy enables Fast3r to handle $N = 1000$ images during inference, even if only trained with $N = 20$ images.

3.4. Memory-Efficient Implementation

With a standard transformer architecture and a single-pass inference procedure, Fast3R is able to leverage many of the recent advances designed to improve scalability at train and inference time [2, 13, 23, 54].

For example, model size and throughput can be increased by sharding the model and/or data minibatch across multiple machines, such as through model parallelism [22, 45], data parallelism [26], and tensor parallelism [32, 46]. During training, optimizer weights, states, and gradients can also be sharded [35]. Systems-level advances have also been proposed, such as FlashAttention [7, 8], which uses GPU kernels leveraging the hardware topology to compute attention in a time and memory-efficient way. These are implemented in libraries such as FAIRScale [16], DeepSpeed [35] and Huggingface [64], and require significant engineering effort.

The Fast3R meta-architecture is explicitly designed to take advantage of these efforts. We leverage two different forms of parallelism at training and inference time, as well as FlashAttention, described in more detail in Sec. 4. More broadly, we believe that our approach will continue to benefit in the longer term as transformer-based scaling infrastructure continues to mature.

4. Experiments

Training Data We train on a mix of real-world object-centric and scene scan data: CO3D [39], ScanNet++ [67], ARKitScenes [4], and Habitat [43]. We use a subset of the datasets in DUS3R, specifically 4 of the 9 datasets, for a total of around 2000 unique scenes scans and 1300 videos of 50 object classes. Note that this is only 7% of CO3D, which is also what the baselines DUS3R [61], Spann3R [57], and MAS3R [25] use.

Baselines DUS3R [61] is the closest approach to ours, and competitive on visual odometry and reconstruction benchmarks. That paper contains extensive comparisons against other methods, and we adopt it as our main baseline. We additionally consider DUS3R’s follow-up work,

²Patches H_1 from the first image I_1 are always embedded with p_1 , since that image provides the coordinate frame for the global head.

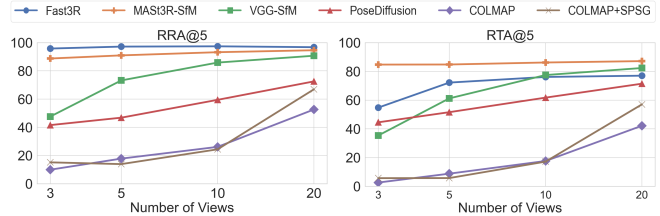


Figure 4. **Pose accuracy with more views:** Fast3R improves with the context from more views. Fast3R saturates the orientation portion of the benchmark, even using 3-5 views.

MAS3R [25], as well as a concurrent work Spann3R [57], which also seeks to replace DUS3R’s expensive global alignment stage by sequentially processing frames with a spatial memory. For camera pose estimation and 3D reconstruction, we include comparisons to task-specific methods.

Architecture Details In our experiments we use the following components for the meta-architecture:

1. The *Image Encoder* uses a CroCo ViT-B [63] architecture, initialized with DUS3R pretrained weights [61].
2. The *Fusion Transformer* is a 12-layer ViT-B model using the BERT architecture with 12 layers, 12 heads, embedding dimension size 768 and MLP ratio 4.0. We use a pool size of $N' = 1000$ image index embeddings.
3. The *Pointmap Decoder* is a dense vision Transformer following the DPT-L architecture [37].

Training Details Our models are trained on 512×512 images using AdamW [29] for 6.5K steps, with a learning rate of 0.0001 and cosine annealing schedule. Unlike DUS3R, we do not use staged training. For the most part, we use the same dataloaders used in the public baseline implementations. We train with batch size 64, with each sample consisting of a tuple of $N = 20$ views. This process takes 6.13 days on 64 A100 GPUs. We additionally make use of several strategies to enable efficient training. First, we use the FlashAttention [7, 8] to improve time and memory efficiency. Even so, a naïve implementation runs out of memory even with batch size 1 when $N > 16$, so we use DeepSpeed ZeRO stage 2 training [35], whereby optimizer states, moment estimates, and gradients are partitioned across different machines. This enables us to train with up to $N = 28$ views per data sample, with a batch size of one per GPU.

4.1. Inference Efficiency

At inference time, we aim to handle 1000+ views compared to 20 during training, which requires additional optimizations. We observe the memory bottleneck at inference is in the DPT heads producing the pointmaps: with 320 views on a single A100 GPU, over 60% of VRAM is consumed by activations from the DPT heads, largely due to each needing to upsample 1024 tokens into a high-resolution 512×512 image. To address this, we implement a simple version of tensor parallelism, putting the model on GPU 0 and then

Methods	Co3Dv2 [39]					FPS
	RRA@15 \uparrow	RRA@5 \uparrow	RTA@15 \uparrow	RTA@5 \uparrow	mAA(30) \uparrow	
Colmap+SG [9, 41]	36.1	24.4	27.3	17.2	25.3	0.056
PixSfM [28]	33.7	26.1	32.9	17.6	30.1	-
RelPose [70]	57.1	-	-	-	-	0.02
PosReg [59]	53.2	-	49.1	-	45.0	0.015
PoseDiff [59]	80.5	59.5	79.8	61.7	66.5	0.015
RelPose++ [27]	(85.5)	-	-	-	-	0.02
DUST3R [60]	96.2	-	86.8	-	76.7	0.78
MASt3R [24]	94.6	93.2	91.9	86.2	81.8	0.23
Fast3R (Ours)	99.7	97.4	87.1	76.1	82.5	251.1

Table 1. **Multi-view pose regression on the CO3Dv2 [39] dataset with 10 random frames.** Parentheses denote methods that do not report results on the 10 views set; we report their best for comparison (8 views). Fast3R does not assume known camera intrinsics.

# Views	Fast3R		DUST3R	
	Time (s)	Peak GPU Mem (GiB)	Time (s)	Peak GPU Mem (GiB)
2	0.065	3.84	0.092	3.52
8	0.122	6.33	8.386	24.59
32	0.509	13.25	129.0	67.61
48	0.84	20.8	OOM	OOM
320	15.938	41.90	OOM	OOM
800	89.569	55.97	OOM	OOM
1000	137.62	63.01	OOM	OOM
1500	308.85	78.59	OOM	OOM

Table 2. **System performance metrics for different view counts on Fast3R and DUST3R on a single A100.** Time is measured in seconds (s), and memory is measured in gibibytes (GiB). Each view is 512x384 in resolution. For DUST3R, at 48 views the N^2 pairwise reconstructions eventually consume all VRAM at its global alignment stage. Note that Fast3R’s reported fastest FPS of 251.1 uses 108 views in 224x224 resolution.

copying the DPT heads to each of the $K - 1$ other GPUs. When processing a batch of $N \approx 1000$ images, we pass the entire batch through the ViT encoder and global fusion decoder, and then split the outputs across K machines for parallel DPT head inference.

Table 2 shows the inference time and memory usage as we increase the number of views. Fast3R is able to process up to 1500 views in a single pass, whereas DUST3R runs out of memory past 32. Fast3R also has a significantly faster inference time, with gains that increase with more views.

4.2. Camera Pose Estimation

We evaluate camera pose estimation on unseen trajectories from 41 object categories from CO3D [39]. Following [61], we sample 10 random views from each trajectory.

Inspired by DUST3R [61], we estimate the *focal length*, *camera rotation*, and *camera translation* from the predicted global pointmaps. We begin by initializing a set of random focal length guesses based on the image resolution, then use RANSAC-PnP to estimate the camera’s rotation and translation based on the guessed focal lengths and the global pointmap. The count of outliers from RANSAC-PnP

Method	FPS	7 scenes [47]		NRGBD [3]	
		Acc \downarrow	Comp \downarrow	Acc \downarrow	Comp \downarrow
F-Recon [65]	<0.1	7.62	2.31	20.59	6.31
DUST3R \dagger [61]	0.78	1.23	0.91	2.51	1.03
Spann3R [57]	65.4	1.48	0.85	3.15	1.10
Fast3R (Ours)	251.1	1.58	0.93	3.40	1.01

Table 3. **Quantitative reconstruction results on scene datasets:** The numbers indicate median distance to GT points on 7-Scenes [47] and NRGBD [3] datasets. These datasets contain video trajectories of 500-1500 frames, and we evaluate using the same frame skip as other baselines. For 7-Scenes this is $\text{skip}=20$, and NRGBD uses $\text{skip}=40$. DUST3R \dagger indicates using DUST3R’s final weights on 224×224 images, to fit within GPU memory. Distances are scaled $100\times$ to remove the leading 0.0. We defer the full table to the appendix.

Method	Views	Acc \downarrow	Comp \downarrow
		Med.	Med.
DUST3R [61]	all/5	1.159	0.914
DUST3R \dagger [61]	all/5	1.297	1.002
Spann3R [57]	all/5	2.268	1.295
Fast3R (Ours)	all/5	1.706	0.857

Table 4. **Quantitative results on object-centric DTU [1] dataset.** Using a $\text{skip}=5$ on trajectories of 49 frames.

is used to score each guessed focal length (lower is better), and the best-scoring focal length is selected to compute the intrinsic and extrinsic camera matrices.

During RANSAC-PnP, we only use points with the top 15% confidence scores predicted by Fast3R, ensuring efficient PnP processing and reducing outliers. If all images are known to originate from the same physical camera, we use the focal length estimated from the first view as the focal length for all cameras, as this initial estimate has been empirically found to be more reliable. Otherwise, we independently estimate the focal length for each input. It is worth noting that the camera pose estimation process is parallelized using multi-threading, ensuring minimal wall-clock time. Even for hundreds of views, the process completes in just a few seconds on standard CPUs.

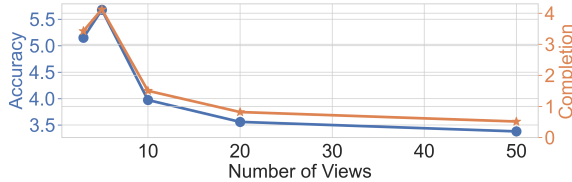


Figure 5. **DTU reconstruction quality vs. test number of views.** Accuracy and Completion (lower is better) get better as we inference with more views.

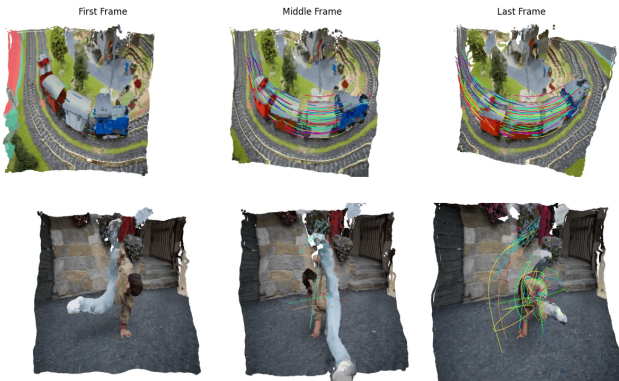


Figure 6. **Qualitative 4D reconstruction results on unseen dynamic scenes in DAVIS.** Results are obtained with one forward pass. The tracks are visualized using ground-truth track annotations from TAP-Vid-DAVIS [11].

We report Relative Rotation Accuracy (RRA) and Relative Translation Accuracy (RTA) at a threshold of 15° , mean Average Accuracy (mAA) at threshold 30° , and model frames per second (FPS) in Table 1. On Co3D, Fast3R surpasses all other methods across the RRA and mAA metrics, achieving near-perfect RRA, while remaining competitive on RTA. Importantly, it is also orders of magnitude faster: $200\times$ faster than DUST3R and $700\times$ faster than MAST3R.

Figure 4 and Figure 5 shows that Fast3R’s predictions improve with more views, indicating that the model is able to use the additional context from multiple images.

4.3. 3D Reconstruction

We evaluate Fast3R’s 3D reconstruction on scene-level benchmarks: 7-Scenes [47] and Neural RGB-D [3], and the object-level benchmark DTU [1].

We found that local pointmaps learn finer detail than the global pointmaps. Therefore we use the local pointmaps for detail and the global pointmaps for the high-level structure. Specifically, we independently align each image’s local pointmap to the global pointmap using ICP, and use aligned local pointmaps for evaluation.

Fast3R is competitive with other pointmap reconstruction methods like DUST3r and MAST3R, while being significantly faster, as shown in Table 3 and Table 4. We believe that Fast3R will continue to improve with better reconstruction

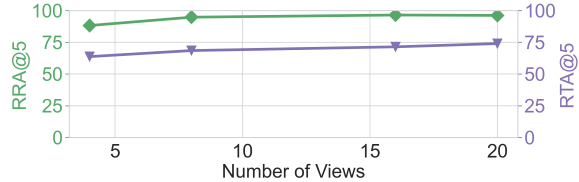


Figure 7. **Increasing # views during training: camera pose estimation on CO3D.** Estimates of both orientation (RRA@5) and translation (RTA@5) improve with more views.

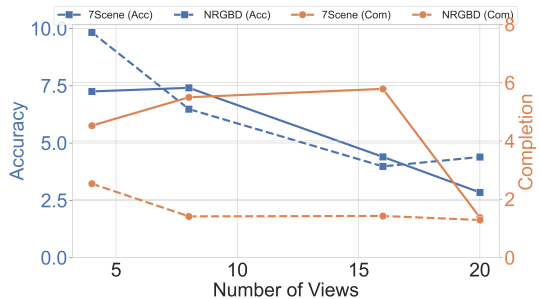


Figure 8. **Increasing # views during training: reconstruction on 7scenes and NRGD.** Accuracy and Completion (lower is better) get better as we train with more views. Normal Consistency (high is better) also gets better as we train with more views.

tion data, more compute, and better training recipes. We show supportive scaling experiments in Figure 5.1.

4.4. 4D Reconstruction: Qualitative Results

Because Fast3R can handle multiple frames naturally, one may wonder how well Fast3R can handle *dynamic* scenes. We qualitatively test Fast3R’s 4D reconstruction ability, showing examples of dynamic aligned pointmaps at multiple time steps in Figure 6. Fast3R can be trained to achieve such results by finetuning a 16 static views checkpoint on the PointOdyssey [73] and TartanAir [62] datasets, consisting of 110 dynamic and 150 static scenes, respectively. We freeze the ViT encoder, use 224×224 image resolution, and swap in a newly-initialized global DPT head. We fine-tune the model with 15 epochs with a frame length of 16, batch size per GPU of 1, and use the same learning rate schedule as Fast3R. The process takes 45 hours to finetune on 2 Nvidia Quadro RTX A6000 GPUs.

We see that our approach produces qualitatively reasonable reconstructions with minimal changes. MonST3R [69] is a concurrent work also tackling dynamic scene reconstruction that builds atop DUST3R. However, like DUST3R, it assumes a pairwise architecture and also uses a separate model to predict optical flow. We show that the same Fast3R architecture trained end-to-end with the same many-view pointmap regression (just swapping the data to dynamic scenes), can also work for 4D reconstruction. Importantly, our method remains significantly faster, opening the poten-

		Test Views				
		4	8	20	24	
Train Views	No rand. emb.	4	0.044	0.228	0.212	0.324
	8	0.033	0.033	0.145	0.297	
	20	0.038	0.055	0.045	0.265	
4 (+ rand. emb.)			0.027	0.026	0.033	0.031

Figure 9. **Effect of sampling image index PE during training.** If we train the model without sampling index embeddings, regression loss spikes (orange) when testing with more views than seen at training (top). Our embedding strategy performs comparably even with $6\times$ the number of views during training.

tial for real-time applications.

5. Ablation Studies

5.1. Scaling the number of views

Fast3R is able to use all-to-all attention during training, which lets it learn from global context. We hypothesize that the additional context provided by more views during training allows the model to learn higher-order correspondences between multiple frames, ultimately increasing model performance and increasing potential for scaling.

Figures 7 and 8 shows that training on increasingly more views consistently improves RRA and RTA for visual odometry, and reconstruction accuracy—even when the number of views used during evaluation is held constant and the model is ultimately evaluated on *fewer* views than were seen during training. We further evaluate the model’s ability to reason about additional views by increasing the number of images that Fast3R sees during inference. Figure 4 and Figure 5 indicate that as the model uses more views, the average *per-view* performance improves. This behavior holds for all evaluated metrics in both camera pose estimation and reconstruction. As shown in Figure 5, the model has a better per-view accuracy using 50 images than it does with 20, even though it was trained with 20. Many applications (*e.g.* reconstruction, odometry) require inference on many views, which is a major motivation for Fast3R removing the pairwise constraint.

5.2. Training without position interpolation

In Section 3.3, we introduced a randomized version of [5] to enable inference on more views than seen training. Without this technique, model accuracy quickly degrades for pointmap corresponding to image indexes outside the training range, as shown in Figure 9 (top). A version of Fast3R trained on $N = 4$ views still produces high qual-

Method	7-Scenes [47]		NRGBD [3]		DTU	
	Acc.↓	Comp.↓	Acc.↓	Comp.↓	Acc.↓	Comp.↓
Fast3R (ours)	2.84	1.37	4.39	1.28	3.91	1.41
w/o local head	4.81	1.64	4.85	1.32	3.88	1.41
Δ	+1.97	+0.27	+0.46	+0.04	-0.03	0.00

Table 5. **Ablation on the effect of local head on 3D reconstruction.** Red/green indicate an increase/decrease in error compared to using the local pointmap aligned to the global pointmap.

ity pointmaps for views in slot 5 to 24 (Figure 9 bottom).

5.3. Removing the local decoder head

Table 5 shows that removing the local head learns finer details before the global head does. We hypothesize that this is because the global head needs to first learn the coordinate system and then learn the fine details. As the model improves, the local head may not be necessary.

6. Conclusion

We introduce Fast3R, a transformer that predicts 3D locations for all pixels in a common frame of reference, directly in a single forward pass. By replacing the whole SfM pipeline with a generic architecture trained end-to-end, Fast3R and similar approaches should benefit from the usual scaling rules for transformers: consistent improvement with better data and increased parameters. Since Fast3R uses global attention, it avoids two potentially artificial scaling limits due to bottlenecks in existing systems. First, the bottleneck of image pair reconstruction restricts the information available to the model. Second, pairwise global optimization can only make up for this so much and does not improve with more data.

With our efficient implementation, Fast3R can operate at > 250 FPS, and process 1500 images in one forward pass, far exceeding other methods while achieving competitive results on 3D reconstruction and camera pose estimation benchmarks. We demonstrate that Fast3R can be fine-tuned to reconstruct videos by changing the data and without modifying the pointmap regression objective and architecture. In contrast with pipeline approaches bottlenecked by custom and slow operations, Fast3R inherits the benefits of future engineering improvements to efficiently serve and train large transformer-based models. For example, packages like Deepspeed-Inference [38], FlashAttention [7, 8] provide fused kernels, model parallelism, and data parallelism. These speed up inference and reduce memory requirements, allowing more images per device, and the number of images scales with the number of devices.

Limitations: A current limiting factor for scaling may be data accuracy and quantity. Synthetic data [34, 40] could be a solution as, broadly speaking, models trained for geometry estimation seem to generalize well from simulation data.

Fast3R can successfully use simulated data to train for 4D reconstruction, showing generalization results on DAVIS. Similarly, DepthAnythingV2 [66] showed the potential of this approach to scale for monocular depth estimation.

The architecture of Fast3R allows for parallel processing of many views, and its positional embedding design enables “train short, test long” in terms of context length of views. However, we observed that for scenes with very large reconstruction areas, when the number of views becomes extreme (e.g., more than 200), the point map of some views (particularly those with a low confidence score) begins to exhibit drifting behavior. One current way to address this issue is to drop frames with low confidence scores. In dense reconstruction, this approach typically does not hurt reconstruction quality too much. However, to fundamentally address this problem, we hypothesize that future work could explore the following avenues: (1) incorporating more data containing large scenes to improve generalization to such cases; (2) designing better positional embeddings inspired by state-of-the-art long-context language models [74], which can handle very long context lengths and exploit the temporal structure of ordered image sequences (e.g., video).

References

- [1] Henrik Aanæs, Rasmus Ramsbøl Jensen, George Vogiatzis, Engin Tola, and Anders Bjarholm Dahl. Large-scale data for multiple-view stereopsis. *International Journal of Computer Vision*, 120:153–168, 2016. 6, 7
- [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmerschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report, 2024. 5
- [3] Dejan Azinović, Ricardo Martin-Brualla, Dan B Goldman, Matthias Nießner, and Justus Thies. Neural rgb-d surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6290–6301, 2022. 6, 7, 8
- [4] Gilad Baruch, Zhuoyuan Chen, Afshin Dehghan, Tal Dimry, Yuri Feigin, Peter Fu, Thomas Gebauer, Brandon Joffe, Daniel Kurz, Arik Schwartz, and Elad Shulman. Arkitscenes: A diverse real-world dataset for 3d indoor scene understanding using mobile rgb-d data. *arXiv preprint arXiv:2111.08897*, 2021. 2, 4, 5
- [5] Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuan-dong Tian. Extending context window of large language models via positional interpolation. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024. 5, 8
- [6] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2017. 2
- [7] Tri Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning. In *International Conference on Learning Representations (ICLR)*, 2024. 5, 8
- [8] Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 5, 8
- [9] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 337–33712, 2017. 6
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, 2019. Association for Computational Linguistics. 4
- [11] Carl Doersch, Ankush Gupta, Larisa Markeeva, Adrià Recasens, Lucas Smaira, Yusuf Aytar, João Carreira, Andrew Zisserman, and Yi Yang. Tap-vid: A benchmark for tracking any point in a video, 2022. 7
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021. 4, 13
- [13] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models, 2024. 5
- [14] Mihai Dusmanu, Ignacio Rocco, Tomás Pajdla, Marc Pollefeys, Josef Sivic, Akihiko Torii, and Torsten Sattler. D2-net: A trainable cnn for joint description and detection of local features. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8084–8093, 2019. 2
- [15] Ainaz Eftekhari, Alexander Sax, Jitendra Malik, and Amir Zamir. Omnidata: A scalable pipeline for making multi-task mid-level vision datasets from 3d scans. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10786–10796, 2021. 3
- [16] FairScale authors. Fairscale: A general purpose modular pytorch library for high performance and large scale training. <https://github.com/facebookresearch/fairscale>, 2021. 5
- [17] Zhiwen Fan, Wenyan Cong, Kairun Wen, Kevin Wang, Jian Zhang, Xinghao Ding, Danfei Xu, Boris Ivanovic, Marco Pavone, Georgios Pavlakos, Zhangyang Wang, and Yue Wang. Instantsplat: Unbounded sparse-view pose-free gaussian splatting in 40 seconds, 2024. 13
- [18] Silvano Galliani, Katrin Lasinger, and Konrad Schindler. Massively parallel multiview stereopsis by surface normal diffusion. In *Proceedings of the IEEE international conference on computer vision*, pages 873–881, 2015. 1

- [19] Pierre Gleize, Weiyao Wang, and Matt Feiszli. Silk: Simple learned keypoints. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 10932–10942, 2023. 2
- [20] Kristen Grauman, Andrew Westbury, Lorenzo Torresani, Kris Kitani, Jitendra Malik, Triantafyllos Afouras, Kumar Ashutosh, Vijay Baiyya, Siddhant Bansal, Bikram Boote, et al. Ego-exo4d: Understanding skilled human activity from first- and third-person perspectives. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2024. 2
- [21] Richard Hartley and Andrew Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, New York, NY, USA, 2 edition, 2003. 2
- [22] Yanping Huang, Youlong Cheng, Ankur Bapna, Orhan Firat, Dehao Chen, Mia Chen, HyoukJoong Lee, Jiquan Ngiam, Quoc V Le, Yonghui Wu, and zhifeng Chen. Gpipe: Efficient training of giant neural networks using pipeline parallelism. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2019. 5
- [23] Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L elio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th eophile Gervet, Thibaut Lavril, Thomas Wang, Timoth e Lacroix, and William El Sayed. Mixtral of experts, 2024. 5
- [24] Vincent Leroy, Yohann Cabon, and J er ome Revaud. Grounding image matching in 3d with mast3r. *arXiv preprint arXiv:2406.09756*, 2024. 6
- [25] Vincent Leroy, Yohann Cabon, and Jerome Revaud. Grounding image matching in 3d with mast3r, 2024. 2, 5
- [26] Shen Li, Yanli Zhao, Rohan Varma, Omkar Salpekar, Pieter Noordhuis, Teng Li, Adam Paszke, Jeff Smith, Brian Vaughan, Pritam Damania, and Soumith Chintala. Pytorch distributed: Experiences on accelerating data parallel training. *CoRR*, abs/2006.15704, 2020. 5
- [27] Amy Lin, Jason Y Zhang, Deva Ramanan, and Shubham Tulsiani. Relpose++: Recovering 6d poses from sparse-view observations. *arXiv preprint arXiv:2305.04926*, 2023. 2, 6
- [28] Philipp Lindenberger, Paul-Edouard Sarlin, Viktor Larsson, and Marc Pollefeys. Pixel-perfect structure-from-motion with featuremetric refinement. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5967–5977, 2021. 6
- [29] Ilya Loshchilov and Frank Hutter. Fixing weight decay regularization in adam. *CoRR*, abs/1711.05101, 2017. 5
- [30] Ra ul Mur-Artal and Juan D. Tard os. ORB-SLAM2: an open-source SLAM system for monocular, stereo and RGB-D cameras. *IEEE Transactions on Robotics*, 33(5):1255–1262, 2017. 2
- [31] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. Orb-slam: a versatile and accurate monocular slam system. *IEEE transactions on robotics*, 31(5):1147–1163, 2015. 1
- [32] Deepak Narayanan, Mohammad Shoeybi, Jared Casper, Patrick LeGresley, Mostofa Patwary, Vijay Korthikanti, Dmitri Vainbrand, Prethvi Kashinkunti, Julie Bernauer, Bryan Catanzaro, Amar Phanishayee, and Matei Zaharia. Efficient large-scale language model training on GPU clusters. *CoRR*, abs/2104.04473, 2021. 5
- [33] Maxime Oquab, Timoth e Darcet, Theo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Russell Howes, Po-Yao Huang, Hu Xu, Vasu Sharma, Shang-Wen Li, Wojciech Galuba, Mike Rabbat, Mido Assran, Nicolas Ballas, Gabriel Synnaeve, Ishan Misra, Herve Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2023. 4
- [34] Alexander Raistrick, Lahav Lipson, Zeyu Ma, Lingjie Mei, Mingzhe Wang, Yiming Zuo, Karhan Kayan, Hongyu Wen, Beining Han, Yihan Wang, Alejandro Newell, Hei Law, Ankit Goyal, Kaiyu Yang, and Jia Deng. Infinite photorealistic worlds using procedural generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12630–12641, 2023. 8
- [35] Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimization towards training A trillion parameter models. *CoRR*, abs/1910.02054, 2019. 5
- [36] Ren e Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. In *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2020. 3
- [37] Ren e Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. *arXiv preprint arXiv:2103.13413*, 2021. 4, 5
- [38] Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3505–3506, 2020. 8
- [39] Jeremy Reizenstein, Roman Shapovalov, Philipp Henzler, Luca Sbordone, Patrick Labatut, and David Novotny. Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10901–10911, 2021. 2, 5, 6
- [40] Mike Roberts, Jason Ramapuram, Anurag Ranjan, Atulit Kumar, Miguel Angel Bautista, Nathan Paczan, Russ Webb, and Joshua M. Susskind. Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding. In *International Conference on Computer Vision (ICCV) 2021*, 2021. 8
- [41] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. SuperGlue: Learning feature matching with graph neural networks. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4937–4946, 2019. 6

- [42] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. SuperGlue: Learning feature matching with graph neural networks. In *CVPR*, 2020. 2
- [43] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv Batra. Habitat: A platform for embodied AI research. *CoRR*, abs/1904.01201, 2019. 5
- [44] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 2
- [45] Noam Shazeer, Youlong Cheng, Niki Parmar, Dustin Tran, Ashish Vaswani, Penporn Koanantakool, Peter Hawkins, HyoukJoong Lee, Mingsheng Hong, Cliff Young, Ryan Sepassi, and Blake A. Hechtman. Mesh-tensorflow: Deep learning for supercomputers. *CoRR*, abs/1811.02084, 2018. 5
- [46] Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-lm: Training multi-billion parameter language models using model parallelism. *CoRR*, abs/1909.08053, 2019. 5
- [47] Jamie Shotton, Ben Glocker, Christopher Zach, Shahram Izadi, Antonio Criminisi, and Andrew Fitzgibbon. Scene coordinate regression forests for camera relocalization in rgb-d images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2930–2937, 2013. 6, 7, 8
- [48] Cameron Smith, David Charatan, Ayush Kumar Tewari, and Vincent Sitzmann. Flowmap: High-quality camera poses, intrinsics, and depth via gradient descent. *ArXiv*, abs/2404.15259, 2024. 2
- [49] Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. In *Advances in Neural Information Processing Systems*, pages 7537–7547, 2020. 4
- [50] Hao Tang, Weiyao Wang, and Matt Feiszli. Aden: Adaptive density representations for sparse-view camera pose estimation. *arXiv preprint arXiv:2408.09042*, 2024. 2
- [51] Zachary Teed and Jia Deng. Droid-slam: Deep visual slam for monocular, stereo, and rgb-d cameras. In *Neural Information Processing Systems*, 2021. 2
- [52] 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
- [53] Sebastian Thrun. Probabilistic robotics. *Communications of the ACM*, 45(3):52–57, 2002. 1
- [54] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models, 2023. 5
- [55] Vadim Tschernezki, Ahmad Darkhalil, Zhifan Zhu, David Fouhey, Iro Larina, Diane Larlus, Dima Damen, and Andrea Vedaldi. EPIC Fields: Marrying 3D Geometry and Video Understanding. In *Proceedings of the Neural Information Processing Systems (NeurIPS)*, 2023. 2
- [56] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017. 2
- [57] Hengyi Wang and Lourdes Agapito. 3d reconstruction with spatial memory. *arXiv preprint arXiv:2408.16061*, 2024. 3, 5, 6
- [58] Jianyuan Wang, Nikita Karaev, Christian Rupprecht, and David Novotny. Vggsfm: Visual geometry grounded deep structure from motion. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 21686–21697, 2023. 2
- [59] Jianyuan Wang, C. Rupprecht, and David Novotný. Posediffusion: Solving pose estimation via diffusion-aided bundle adjustment. *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9739–9749, 2023. 2, 6
- [60] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jérôme Revaud. Dust3r: Geometric 3d vision made easy. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20697–20709, 2023. 6
- [61] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vision made easy. In *CVPR*, 2024. 1, 2, 3, 4, 5, 6
- [62] Wenshan Wang, DeLong Zhu, Xiangwei Wang, Yaoyu Hu, Yuheng Qiu, Chen Wang, Yafei Hu, Ashish Kapoor, and Sebastian Scherer. Tartanair: A dataset to push the limits of visual slam. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020. 7
- [63] Philippe Weinzaepfel, Vincent Leroy, Thomas Lucas, Romain Bréquier, Yohann Cabon, Vaibhav Arora, Leonid Antsfeld, Boris Chidlovskii, Gabriela Csurka, and Jérôme Revaud. Croco: Self-supervised pre-training for 3d vision tasks by cross-view completion. *Advances in Neural Information Processing Systems*, 35:3502–3516, 2022. 4, 5
- [64] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. Huggingface’s transformers: State-of-the-art natural language processing. *CoRR*, abs/1910.03771, 2019. 5
- [65] Guangkai Xu, Wei Yin, Hao Chen, Chunhua Shen, Kai Cheng, and Feng Zhao. Frozenrecon: Pose-free 3d scene reconstruction with frozen depth models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9310–9320, 2023. 6
- [66] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *CVPR*, 2024. 3, 4, 9
- [67] Chandan Yeshwanth, Yueh-Cheng Liu, Matthias Nießner, and Angela Dai. Scannet++: A high-fidelity dataset of 3d indoor scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12–22, 2023. 2, 4, 5
- [68] Kwang Moo Yi, Eduard Trulls, Vincent Lepetit, and Pascal V. Fua. Lift: Learned invariant feature transform. In *European Conference on Computer Vision*, 2016. 2

- [69] 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. [3](#), [7](#)
- [70] Jason Y. Zhang, Deva Ramanan, and Shubham Tulsiani. Rel-pose: Predicting probabilistic relative rotation for single objects in the wild. *ArXiv*, abs/2208.05963, 2022. [6](#)
- [71] Jason Y. Zhang, Amy Lin, Moneish Kumar, Tzu-Hsuan Yang, Deva Ramanan, and Shubham Tulsiani. Cameras as rays: Pose estimation via ray diffusion. *ArXiv*, abs/2402.14817, 2024. [2](#)
- [72] Wang Zhao, Shaohui Liu, Hengkai Guo, Wenping Wang, and Y. Liu. Particlesfm: Exploiting dense point trajectories for localizing moving cameras in the wild. In *European Conference on Computer Vision*, 2022. [2](#)
- [73] Yang Zheng, Adam W Harley, Bokui Shen, Gordon Wetstein, and Leonidas J Guibas. Pointodyssey: A large-scale synthetic dataset for long-term point tracking. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19855–19865, 2023. [7](#)
- [74] Dawei Zhu, Nan Yang, Liang Wang, Yifan Song, Wenhao Wu, Furu Wei, and Sujian Li. Pose: Efficient context window extension of llms via positional skip-wise training. In *The Twelfth International Conference on Learning Representations*. [9](#)

Fast3R: Towards 3D Reconstruction of 1000+ Images in One Forward Pass

Supplementary Material

A. Model Scaling Effect

We investigate the effect of scaling model size by trying three model sizes for the Fusion Transformer: ViT-base, ViT-large, and ViT-huge, according to the settings in the original ViT paper [12]. The results are shown in Figure 11. This experiment demonstrates that larger model size continually benefits 3D tasks including camera pose estimations and 3D reconstruction. Note that the Fusion Transformer size used in the main text for all experiments is a ViT-base.

B. Data Scaling Effect

We study the effect of scaling the data using 4 different scales of data, 12.5%, 25%, 50%, and 100%, to train the model. The results are shown in Figure 12. The training settings for all models are kept the same except for how much data they have access to. The results demonstrate that Fast3R continually benefits from more data, suggesting Fast3R could achieve better results in the future given more data.

C. Gaussian Splatting

We qualitatively demonstrate the potential of using Fast3R’s output for downstream novel view synthesis tasks. A visualization of the Gaussian Splatting generated by adopting the pipeline of InstantSplat [17] is shown in Figure 13.

D. Bundle Adjustment (via Gaussian splatting)

While not necessary, using bundle-adjustment at inference time can also improve Fast3R’s performance. We show an example of bundle adjustment using Gaussian Splatting (GS-BA).

Specifically, we use InstantSplat [17] to optimize a set of gaussians per-scene, using initializations from a pointcloud, and updates the locations and poses in order to minimize reprojection error. We show an example of the Gaussian reconstruction in Figure 13 shows an example reconstruction on CO3D.

We can compare against ground-truth trajectories from COLMAP. We found that GS-BA significantly reduces both the pose and translation error. Table 6 quantifies this, showing over a 2.5x reduction in translation error and 4x reduction in rotational error on the “Family” scene from Tanks and Temples, which we found to be representative. We show a visualization of the original reconstruction and the poses pre- and post- bundle adjustment. There are only 8 scenes in the evaluation set in InstantSplat.

Method	RPE Rotation (\downarrow)	RPE Translation (\downarrow)
Fast3R	27.9	7.64
Fast3R w/ GS-BA	11.0	1.80

Table 6. **Pose estimation can further improve with Bundle Adjustment.** We show an example on the “Family” scene from Tanks and Temples, using InstantSplat [17].

E. More Visualizations

We show more visualizations of Fast3R’s performance on indoor scenes in Figure 10. Fast3R learns the regularity of indoor rooms (square-like shapes) and demonstrates “loop closure” capabilities.

In Figure 15, we visualize the point cloud produced using local vs. global point maps.



Figure 10. **Visualizations of results on NRGBD scenes.** Fast3R learns the regularity of indoor rooms (square-like shapes) and demonstrates “loop closure” capabilities.

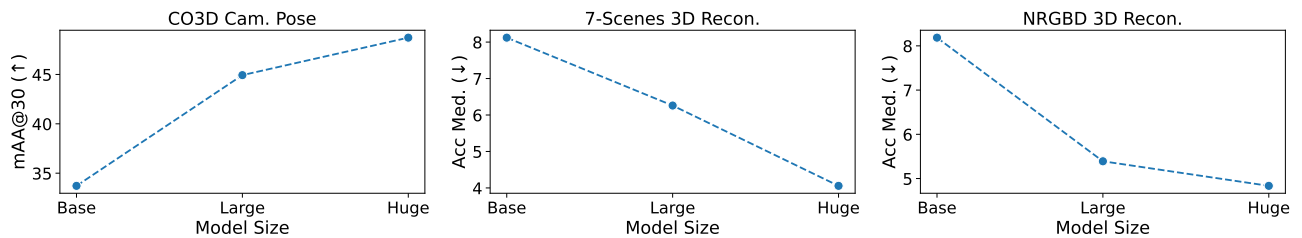


Figure 11. **Model scaling effect.** Increasing the size of the Fusion Transformer leads to better camera pose estimation (↑) and 3D reconstruction (↓). All models are trained for 60k steps (equivalent to 60 epochs; the main paper uses 100 epochs).

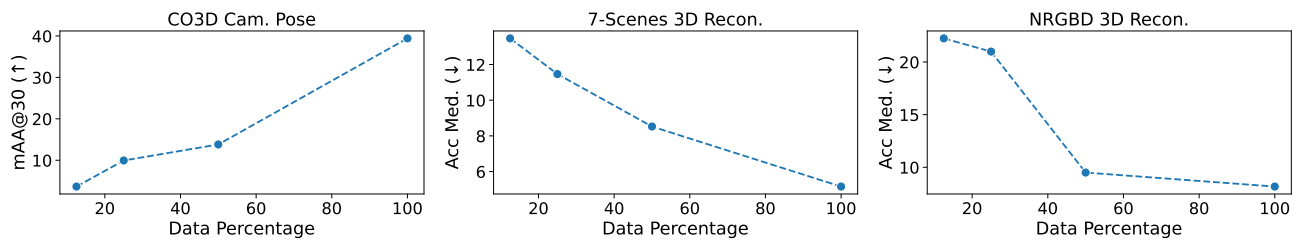


Figure 12. **Data scaling effect.** More training data leads to better camera pose estimation (↑) and 3D reconstruction (↓). All models are trained for 60k steps (equivalent to 60 epochs; the main paper uses 100 epochs).



Figure 13. **Visualization of Gaussians from unseen poses.** The frames are ordered temporally along the direction of the arrows. The middle frames show poses very different from those used for reconstruction, as is evidenced by the large areas with no Gaussians. The scene is fit from 7 images from CO3D.

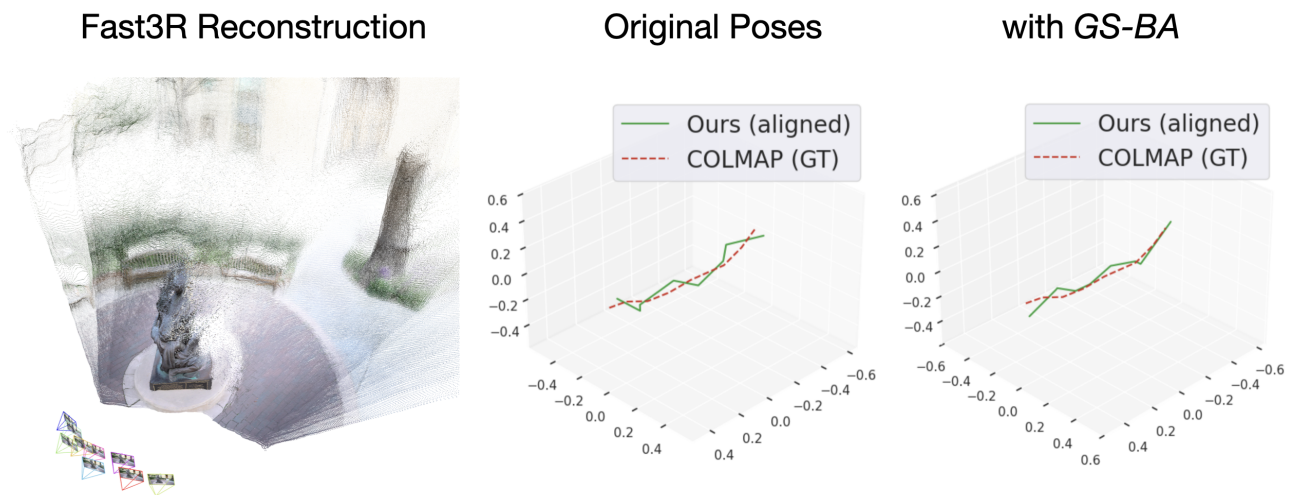


Figure 14. **Bundle adjustment further improves pose.** Left: reconstruction from Fast3R. Middle: Original poses pre-GS-BA. Right: Poses after GS-BA.



Reconstruction using Global Point Map



Reconstruction using Local (aligned to Global) Point Map

Figure 15. **Effect of using local vs. global point map.** Global point maps provide good anchors for locations of points while local point maps use those anchors (by aligning using ICP on the anchor points to the global point map) to provide more accurate point locations.