

The 8-Point Algorithm as an Inductive Bias for Relative Pose Prediction by ViTs

Chris Rockwell, Justin Johnson, David F. Fouhey
University of Michigan

Abstract

We present a simple baseline for directly estimating the relative pose (rotation and translation, including scale) between two images. Deep methods have recently shown strong progress but often require complex or multi-stage architectures. We show that a handful of modifications can be applied to a Vision Transformer (ViT) to bring its computations close to the Eight-Point Algorithm. This inductive bias enables a simple method to be competitive in multiple settings, often substantially improving over the state of the art with strong performance gains in limited data regimes.

1. Introduction

Estimating the relative pose between two images is a fundamental vision problem [17], with applications including 3D understanding [22, 34, 57] and extended reality [30, 35, 38, 73]. Early work focused on robust [14] fitting of models [17, 19, 31, 42] on detected correspondences [2, 32, 54] between the images. This strategy can fail catastrophically with poor correspondence, which is especially frequent in the *wide baseline* setting, when the images have a substantial pose difference. Moreover, even when it is successful, it cannot recover the *scale* of the translation [17]. The situation is often improved in practice by obtaining more images (e.g., SfM [57] and SLAM [41]), or sensors like IMUs [15, 16, 26] and RGBD [11, 70]. Nonetheless, people routinely infer relative pose from two ordinary images with a wide baseline, and whole industries like real estate depend on this ability. Rather than use extra sensors or images, humans integrate cues like correspondence, familiar object size, and priors on scenes. This paper investigates such an ability, to estimate relative pose, including rotation and translation with scale, from two ordinary images.

Based on these observations, there has been much work applying learning to the problem. One line of attack [8, 10, 55, 60, 63] has been to follow the classic pipeline and replace classic correspondence methods [2, 32, 54] with learned ones. This approach is appealing since the learning method finds correspondence, an especially thorny

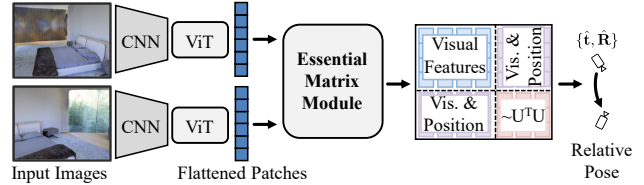


Figure 1. We propose three small modifications to a ViT via the Essential Matrix Module, enabling computations similar to the Eight-Point algorithm. The resulting mix of visual and positional features is a good inductive bias for pose estimation.

challenge in the wide-baseline setting, and the conversion of correspondences to pose is done by a provably correct method [31, 42]. However, it comes at a cost of inheriting the Essential Matrix’s intrinsic scale ambiguity, leading to translation-up-to-scale. Thus, another line of work treats relative camera pose estimation as a learning problem [5, 12, 22, 48]. These approaches have shown promise in the wide-baseline setting, but often involve multiple stages [7, 22], are not as performant as correspondence-based techniques in the settings we try [12, 48], or do not recover a translation scale [5, 7]. Moreover, since these methods learn an end-to-end mapping from images to camera pose with few inductive biases, they are often data hungry.

We propose a Vision Transformer (ViT) [9, 64] approach that estimates rotation and *translation with scale* in one forward pass by integrating the problem’s structure implicitly as an inductive bias. We reconcile the Eight-Point Algorithm [19, 31] with ViTs by showing that a ViT forward pass can be made close to [19, 31] by three minor modifications: (1) bilinear attention [24] instead of attention [64]; (2) quadratic position encodings; and (3) dual-softmax [53, 60, 63] instead of softmax. These modifications are put in one module, the *Essential Matrix Module (EMM)*, that we place atop an otherwise ordinary ViT, as shown in Fig. 1. The EMM gives an inductive bias by providing positional features that approximate a key step of [31], visual features, and features that mix the two.

We attach the Essential Matrix Module to the end of an ordinary ViT [9] described in §3. We train and evaluate this ViT on multiple relative camera pose estimation tasks and datasets as described in §4 and compare with the state of the

art for each on challenging datasets like Matterport3D [6], InteriorNet [29], and StreetLearn [39]. Our experiments on rotation+translation (§4.2) and rotation (§4.3) demonstrate: (1) that our simple approach outperforms (or occasionally matches) multiple alternate networks including concatenation [12, 48] and correlation volume methods [5, 22], techniques based on feature correspondence [8, 32], and techniques trained to optimize a full 3D reconstruction [22]; (2) that each component of the modification is important, as shown by extensive ablations; and (3) that the EMM improves data efficiency by substantially boosting performance in moderate data regimes (§4.4), suggesting that epipolar geometry is a good inductive bias.

2. Related Work

Our work introduces a learning-based approach to relative pose estimation by modifying vision transformers to perform computations similar to the Eight Point Algorithm.

Classic Work. Relative pose estimation from an image pair is a sufficiently broad problem to preclude a full account. We refer readers to [17], and focus on the closest works, which all follow a strategy of solving for pose given correspondences from local descriptors [2, 32, 54]. We revisit the 8-point algorithm [19, 31] that maps correspondences to an Essential Matrix, which was invented by Longuet-Higgins and extended to Fundamental Matrices by [13, 18]. While it often replaced by other approaches that use fewer correspondences (e.g., [27, 42]), much of the 8-point algorithm’s structure are calculations that we show can be done by Transformers. In our wide baseline setting (i.e., large pose difference), historically there are alternate descriptors and specialized techniques [36, 40, 46].

Learned Pose Estimation. Given the difficulties associated with optimization on correspondence, multiple lines of work aim to improve the pipeline with learning. For instance, many methods improve detectors and descriptors [8, 10, 63] or correspondence estimation [3, 47, 50, 55, 60, 72]. These works typically turn correspondences to pose with the Essential Matrix [31, 42], which makes it impossible to recover translation scale without additional signals [17]. In contrast, our proposed work learns a direct mapping without explicitly constructing an Essential Matrix, and therefore recovers scale using image-based cues. We note that our components are also often used in correspondence work [60, 63]; here, we use them directly for pose and show a close relationship between ViTs and [31].

Our method is closer to work that learns a mapping from images to pose. This area of research is relatively newer, and has become more complex over time (e.g., early networks concatenated data from two images [12, 28, 37], which has been supplanted by correlation volumes [5, 22]). These approaches are often data and compute hungry [7],

use multiple stages (e.g., discrete/continuous optimization in [22], two-stage networks in [7, 65]), and use little of the structure of pose estimation. In contrast, our approach brings ViT computations close to this structure, which we hypothesize helps use the data more effectively.

SLAM, SfM, and RGBD. Given the difficulties of pose estimation from two images, a wealth of other approaches have been tried that modify the problem. The most common solution is to use more images with typically high overlap, e.g. via Structure-from-Motion [57], SLAM [41, 61, 66] or localization [4, 49, 68]. In contrast, we aim to solve the two-view, wider-baseline case. Other solutions include adding sensors like an IMU [15, 16, 26] or depth data [11, 69, 70]; our approach relies only on RGB data.

Vision Transformers and Inductive Biases. Large parts of our proposed approach follow a basic recipe for Vision Transformers [9, 64]. These have emerged as a competitor to convolutional neural networks in the past few years, and we refer interested readers to [23, 43] for a more thorough summary. Our work shows that small modifications of the pipeline brings the computations close those of [31]. This is part of a broader trend of injecting geometric inductive biases to networks via layers [45] or token engineering [71].

3. Approach

Our goal is to map two overlapping images to a relative camera pose *including translation scale*, or a rotation $\mathbf{R} \in \text{SO}(3)$ and translation $\mathbf{t} \in \mathbb{R}^3$. This task requires both robustness to large view changes with limited correspondence, and handling scale ambiguity. We propose a simple approach that fuses ideas from classical multi-view geometry with large-scale learning.

At the heart of our approach is a transformer with small critical changes that mimic a computation used in the Eight Point Algorithm [19, 31]. These changes include bilinear attention [24], dual-softmaxes [53, 60, 63], and an explicit positional encoding. We first analyze the relationship between the Eight Point Algorithm and an alternate setup that is more amenable to computation by a transformer (§3.1). We then describe how we operationalize this by introducing our base transformer and our Essential Matrix Module (§3.2). We conclude by analyzing the learnability of this function with synthetic experiments (§3.3).

3.1. Transformers and the Eight Point Algorithm

The Fundamental and Essential matrices can be obtained from correspondences via the Eight-point algorithm [19, 31, 18, 13]. As input, one assumes N correspondences $[u_i, v_i] \leftrightarrow [u'_i, v'_i]$. With known intrinsics \mathbf{K} , one represents the locations of the correspondences with normalized points $\mathbf{x}_i \equiv \mathbf{K}^{-1}[u_i, v_i, 1]^\top$ and $\mathbf{x}'_i \equiv \mathbf{K}^{-1}[u'_i, v'_i, 1]^\top$ and recovers an Essential matrix (\mathbf{E}); if \mathbf{K} is unknown, one uses

standard homogeneous coordinates (i.e., $\mathbf{x}_i = [u_i, v_i, 1]^\top$) and recovers a Fundamental matrix (\mathbf{F}). Since we have the intrinsics, we will refer to the Essential matrix.

The Eight-Point Algorithm constructs a matrix $\mathbf{U} \in \mathbb{R}^{N \times 9}$ whose i th row $\mathbf{U}_{i,:}$ is the Kronecker product of the correspondences, or $\mathbf{x}_i \otimes \mathbf{x}'_i$. The matrix $\mathbf{U}^\top \mathbf{U} \in \mathbb{R}^{9 \times 9}$ captures the information needed to estimate the Essential matrix: one computes the eigenvector corresponding to the smallest eigenvalue of $\mathbf{U}^\top \mathbf{U}$, reshapes the vector, and makes the reshaped matrix rank deficient. The resulting matrix \mathbf{E} does not uniquely define the relative pose, but rather a family of solutions comprising two rotations \mathbf{R} and \mathbf{R}' and a translation direction (that can be scaled by any $\lambda \neq 0$).

Careful minor modifications of Transformer can enable the computation of the entries of $\mathbf{U}^\top \mathbf{U}$. We assume the transformer is given a set of P patches at locations $\{\mathbf{p}_j\}_{j=1}^P$ where every correspondence is at one of the patches. In addition to using these locations directly, we further define a 6D basis expansion $\phi([u, v, 1]) = [1, u, v, uv, u^2, v^2]$ that we apply to each patch to yield a matrix $\Phi \in \mathbb{R}^{P \times 6}$ such that $\Phi_{j,:} = \phi(\mathbf{p}_j)$. Finally, to represent correspondences implicitly, we define an indicator matrix $\mathbf{A} \in \{0, 1\}^{P \times P}$ such that $\mathbf{A}_{j,k} = 1$ if and only if points \mathbf{p}_k and \mathbf{p}_j are in correspondence and 0 otherwise.

Our key observation is that each unique entry of $\mathbf{U}^\top \mathbf{U} \in \mathbb{R}^{9 \times 9}$ is in the matrix $\Phi^\top \mathbf{A} \Phi \in \mathbb{R}^{6 \times 6}$. While this more compact form is not amenable to eigenvector analysis, it is all the information needed for a learned estimator. A derivation appears in the supplement, but the two critical steps are: first, to decompose the matrix as an explicit sum over correspondences $\mathbf{U}^\top \mathbf{U} = \sum_{i=1}^N \mathbf{U}_{i,:}^\top \mathbf{U}_{i,:}$ and rewrite it implicitly with \mathbf{A} ; and second, that the 36 *unique* entries in $\mathbf{U}_{i,:}^\top \mathbf{U}_{i,:}$ can be generated from $\phi(\mathbf{x}_i) \phi(\mathbf{x}'_i)^\top$.

The remaining step is estimation of \mathbf{R} and \mathbf{t} from $\mathbf{U}^\top \mathbf{U}$. MLPs are universal approximators [21], but a number of things make this easier in practice. First, often one aims to solve a subset of problems from a distribution, rather than *all* instances. Additionally, one is also using a wealth of alternate image-based cues. In addition to facilitating learning, the network can use these cues to resolve the ambiguities intrinsic to \mathbf{E} : for instance, the scale ambiguity can be resolved implicitly via recognizing familiar objects. We explore the learnability of this function in §3.3.

Together, this suggests that transformers estimating \mathbf{R} , \mathbf{t} may benefit from a few small modifications. The crux is that the computation of $\Phi^\top \mathbf{A} \Phi$, using quadratic position encodings per patch in Φ and a correspondence indicator in \mathbf{A} mimics the computation of the entries of $\mathbf{U}^\top \mathbf{U}$. Thus, a network may benefit from having $\Phi^\top \mathbf{A} \Phi$ during prediction. Moreover, \mathbf{A} also should be able to represent unmatched correspondences (i.e., $\sum_{k=1}^P \mathbf{A}_{j,k} \approx 0$). Finally, we stress that the model should also contain features beyond Φ to help learning and resolve ambiguities such as scale.

3.2. Putting things In Practice

Our approach consists of two components. The main component is an *Essential Matrix Module*, which maps from P , D -dimensional transformer tokens, one for each of the P patches in the image, to a feature that is used to predict \mathbf{R} and \mathbf{t} . This module is added to a standard ViT [9] backbone that maps images to a set of tokens. Our backbone deliberately follows a standard vision transformer recipe [9, 64]: we see backbone innovations as orthogonal to innovations in the mapping from tokens to outputs. On the other hand, our Essential Matrix Module must contain critical modifications.

Backbone and Setup. Our backbone consists of two main components that function as a learned mapping from an image to a $\mathbb{R}^{P \times D}$ matrix of features, one per patch. The first component is an encoder that uses the first blocks from a standard ResNet-18 [20], which helps the network extract good features per-patch. On top of this, we use blocks from a standard ViT [9] (ViT-Tiny) to map the the patch features to our final set of P D -dimensional tokens. Since the architectures have different feature sizes, we bridge them with a ResNet block that maps the feature dimensions. A full network description appears in the supplemental.

Standard Transformer Model. The canonical ViT maps a set of patches from one image to an output embedding used for classification. Given a patch embedding, this entails computing query, key and value matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{P \times D}$ followed by $\text{softmax}(\mathbf{Q}\mathbf{K}^\top)\mathbf{V}$. To avoid notational clutter, we drop the usual [64] scaling factor of $1/\sqrt{D}$ inside the softmax here, and in all other softmax references.

For our case of two images, there are two sets of matrices, namely $\mathbf{Q}_1, \mathbf{K}_1, \mathbf{V}_1 \in \mathbb{R}^{P \times D}$ for image 1 and $\mathbf{Q}_2, \mathbf{K}_2, \mathbf{V}_2 \in \mathbb{R}^{P \times D}$ for image 2. The simplest cross attention is to concatenate cross-attention per-image, or

$$[\text{softmax}(\mathbf{Q}_1 \mathbf{K}_2^\top) \mathbf{V}_2, \text{softmax}(\mathbf{Q}_2 \mathbf{K}_1^\top) \mathbf{V}_1]. \quad (1)$$

This approach produces good results, but a few minor modifications can substantially improve its performance.

Essential Matrix Module. We propose three changes to Eqn. 1 that help approximate the entries of $\mathbf{U}^\top \mathbf{U}$. These are shown in Fig. 2.

Bilinear Attention and Quadratic Position Encodings. We apply bilinear attention [24] to the values and quadratic positional encodings, or

$$[\mathbf{V}_2, \Phi]^\top \text{norm}(\mathbf{Q}_1 \mathbf{K}_2^\top) [\mathbf{V}_2, \Phi] \in \mathbb{R}^{(D+6) \times (D+6)} \quad (2)$$

where $\Phi \in \mathbb{R}^{P \times 6}$ contain the positional encodings $[1, u, v, uv, u^2, v^2]$ from §3.1 and norm is a normalization for the raw attention scores. Thus $\mathbf{A} = \text{norm}(\mathbf{Q}_1 \mathbf{K}_2^\top)$. To use both images, we also compute Eqn. 2 substituting in $\mathbf{Q}_2, \mathbf{K}_1$, and \mathbf{V}_1 and concatenate the results, leading to a $(2D^2 + 24D + 72)$ -dimensional feature per attention head.

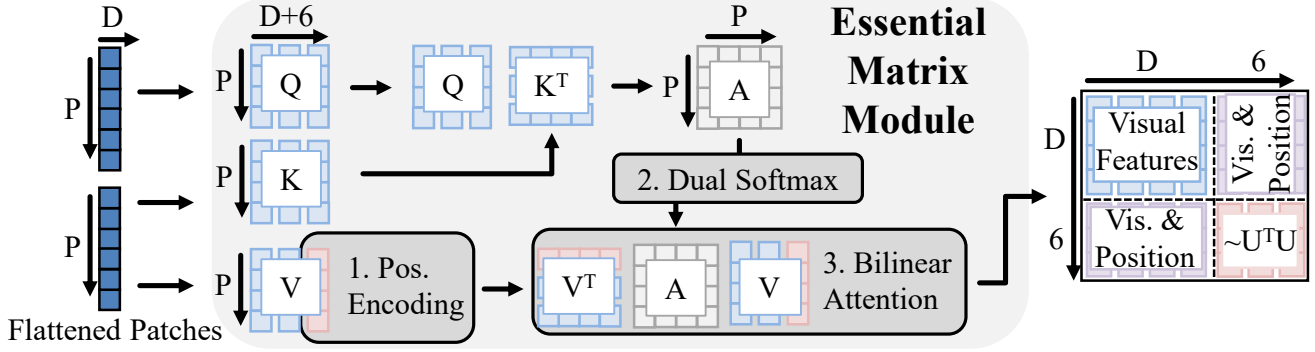


Figure 2. **Essential Matrix Module.** We make three small changes to standard ViT Cross-Attention: (1) appending positional encodings to Values, (2) applying a dual softmax on Affinities, and (3) applying bilinear attention.

Rationale. If $\mathbf{A} = \text{norm}(\mathbf{Q}_1 \mathbf{K}_2^T)$ correctly indicates correspondence, then this computation makes the bottom-right 6×6 submatrix of Eqn. 2 contain the entries of $\mathbf{U}^T \mathbf{U}$. The top-left $D \times D$ submatrix are visual features; the rest mix position and visual features. These image features are important for scale estimation, since $\mathbf{U}^T \mathbf{U}$ does not provide information about scale. They may also *implicitly* contain position encodings (e.g., due to convolutions using zero-padding as a proxy for image location). In practice, $\mathbf{U}^T \mathbf{U}$ is followed by neural layers and thus does not need to match true $\mathbf{U}^T \mathbf{U}$; though in the supplement we find non-zero rank correlation with ground truth.

Dual-softmax. The above is exact when \mathbf{A} is a correspondence indicator matrix. While attention makes this impossible to ensure exactly, we help more closely approximate it with a dual-softmax [53, 60, 63], or set $\text{norm}(\mathbf{Q}_1 \mathbf{K}_2^T)$ to

$$\text{softmax}(\mathbf{Q}_1 \mathbf{K}_2^T, 1) \odot \text{softmax}(\mathbf{Q}_1 \mathbf{K}_2^T, 2), \quad (3)$$

where $\text{softmax}(\cdot, k)$ applies softmax across the k -th axis.

Rationale. Traditional attention normalizes the matrix product $\mathbf{Q}_1 \mathbf{K}_2^T \in \mathbb{R}^{P \times P}$ by a single softmax, $\text{softmax}(\mathbf{Q}_1 \mathbf{K}_2^T, 1)$, forcing $\sum_{k=1}^P \mathbf{A}_{j,k} = 1$. This constraint means that \mathbf{A} cannot indicate correspondence for patches without matches where $\sum_{k=1}^P \mathbf{A}_{j,k} = 0$. At best, attention can be a uniform distribution; at worst, attention can latch onto a random correspondence. In all cases, all patches contribute equally to the final product in Eqn. 2. The network can mitigate this by making non-matching attention uniformly distributed and $\frac{1}{P} \sum_{k=1}^P \mathbf{V}_{k,:} = 0$, but this strategy does not work for parts of Φ , e.g., the u^2 term is non-negative and usually positive.

A dual-softmax suppresses non-matching patches while not altering bidirectional matches. If the attention to-and-from patch j is uniformly distributed, then the total attention $\sum_{k=1}^P \mathbf{A}_{j,k}$ is $\frac{1}{P}$ instead of 1 in the normal softmax case. On the other hand, if patch j and patch k both match well, then the attention approaches 1. Then, even if all

but one patches have no match, their total contribution is smaller ($\frac{P-1}{P}$) than even a single bidirectionally matching patch (1). Thus, the varying weighting helps suppress the contributions of patches without matches. While the form of \mathbf{A} intrinsically makes the computation an approximation, we stress that the consumer of $\Phi \mathbf{A} \Phi^T$ is a learned module and may be able to learn around approximation errors, especially with vision features.

Pose Regressor. Given essential matrix encodings, we regress pose using a 2 hidden layer MLP. We predict translation in real units, and predict rotation in quaternions, normalizing so scale is one. We train only using a 11 geodesic loss on pose where the geodesic loss is the magnitude of the vector between predicted and ground truth pose.

3.3. Synthetic Validation

While our approximation of $\mathbf{U}^T \mathbf{U}$ can be understood analytically, one critical component is the learned mapping from $\mathbf{U}^T \mathbf{U}$ to \mathbf{R}, \mathbf{t} that is done by the *Pose Regressor*. To better understand the learnability of the function, we show the method on synthetic examples with the entries of $\mathbf{U}^T \mathbf{U}$ but no visual features. Our scenes consist of points uniformly sampled inside a sphere with center $\sim \text{Unif}(-\frac{1}{2}, \frac{1}{2})$ and radius $\sim \text{Unif}(\frac{1}{2}, \frac{3}{2})$. We sample camera rotations and translations from distributions that we vary to analyze the learnability of the problem. For each pair of views with sufficient overlap (100 of 10K sampled 3D points projecting to the images), we compute $\mathbf{U}^T \mathbf{U}$, which is used as a feature for pose estimation by a MLP (details in supplement).

We analyze two tasks. The first is *Translation*, or estimating the generating \mathbf{t} ; due to scale-ambiguity, we assume $\|\mathbf{t}\|_2 = 1$ and $t_z > 0$. We quantify errors by the angle between the estimated and true translation. The second is *Rotation*, or estimating the rotation that generated the data, which forces the network to resolve the usual rotation ambiguity of \mathbf{E} . We quantify errors by the rotation geodesic.

We try four distributions. In 3D \mathbf{R} is sampled via uniformly distributed Euler angles, and $\mathbf{t} \sim \text{Unif}(-1, 1)$. The

Table 1. **Synthetic Validation.** We report the median angular error across two tasks (translation & rotation) and four datasets (in decreasing difficulty: 3D, 2D Large/Medium/Small (L/M/S)).

	Translation (°)				Rotation (°)			
	3D	2DL	2DM	2DS	3D	2DL	2DM	2DS
MLP	18.4	5.6	3.0	1.8	33.5	3.6	1.8	0.7
Chance	64.0	49.1	47.9	47.9	125.3	22.2	4.8	1.0

next three are *2D Small/Medium/Large*, consisting of 2D motion primarily in the xz plane with varying amounts of rotation variance: \mathbf{R} is sampled from Normally distributed Euler angles with rotation mainly in y ($y \sim N(0, r)$) and $x, z \sim N(0, \frac{r}{20})$ where $r = 1, 5, 25^\circ$ for small, medium, and large. Translation is mainly in the xz plane $\mathbf{t} \sim N(0, [\frac{1}{3}, \frac{1}{60}, \frac{1}{3}])$. To avoid epipolar degeneracies with no translation, we require $\|\mathbf{t}\| \geq \frac{1}{2}$.

We report results in Table 1 for models trained on 100K samples using only $\mathbf{U}^\top \mathbf{U}$ as features, comparing to chance for context. We compare models trained on 100K samples. Even when trained on 100K samples and estimating a general problem case, the networks learn the function. Once the data is even moderately constrained (2D Large), relative errors drop considerably. This suggests that the function is especially learnable under more constrained rotations.

3.4. Implementation Details

Full implementation details appear in the supplemental, and we will release code for reproducibility. Our encoder is a pretrained ResNet-18, which we truncate to only use the first two of four modules, producing a $24 \times 24 \times 128$ feature map; we use an additional Residual block to map to feature size of 192 for the ViT. We use the Timm [67] ViT implementation, and use ViT-Tiny with a truncated depth of 5 plus our Essential Matrix Module. Outside of the proposed changes, our Essential Matrix Module follows a standard Cross-Attention Transformer Block architecture and normalization [33]. Positional encoding locations utilize known intrinsics, $\mathbf{x}'_i \equiv \mathbf{K}^{-1}[u'_i, v'_i, 1]^\top$, a manner similar to [65]. Each head of the Essential Matrix Module produces a 64-D feature, which is 70-D after concatenating the position encodings. With 3 heads, and the bilinear attention done once per image, this results in $3 \times 70^2 \times 2 = 29\text{K}$ features. We map this large feature to a hidden size of 512 for two hidden layers in our MLP before regressing 7D pose. We implement using PyTorch [44] and use the LieTorch [62] extension for backpropagation of geodesic losses on quaternions. We use learning rate of $5e-4$ and train using Adam [25] optimizer and 1cycle scheduler [58] for 120k iterations with a batch size of 60 split over 10 GTX 1080 Tis, which takes about 1 day.

4. Experiments

We now evaluate the proposed method’s ability to estimate relative pose in comparison to the state of the art in two settings that share common metrics and evaluation settings (§4.1). Our first task (§4.2) is wide baseline rotation *and translation* estimation, or estimating a rotation in $\text{SO}(3)$ and translation in \mathbb{R}^3 (i.e., including a scale). The second task (§4.3) is wide baseline rotation, or estimating a rotation in $\text{SO}(3)$ but no translation. Finally, a crucial argument for our approach is that the modifications of the transformer architecture serve as an inductive bias for the network. We examine this empirically with experiments on substantially reduced data that test data efficiency (§4.4).

4.1. Metrics and Evaluation

For each method, we compute the rotation error (defined as the rotation geodesic to the ground-truth) and translation error (defined as the usual Euclidean distance to the ground truth), and aggregate three summary statistics: the *mean*, the *median*, and the *percent of errors within a threshold* that is task-specific (e.g. 30°) and will be described with each dataset. These capture different aspects of the problem. Specifically, due to symmetries in the data, pose estimation errors are often not unimodally distributed. Instead, often many results are highly accurate and a few are wrong by 90° or 180° . The median error captures what a typical prediction error is like and is outlier robust; the mean is the straight average and is therefore sensitive to outliers; the percent within a threshold captures a sense of how many predictions are “reasonable” for some threshold.

4.2. Wide Baseline Rotation and Translation

We begin by evaluating on our full problem, namely estimating a rotation in $\text{SO}(3)$ and translation, *including scale*, in \mathbb{R}^3 . We follow the setup of [22] to enable comparison with a variety of existing work and published baselines.

Dataset. We use data from Matterport3D [6] consisting of pairs of images with limited overlap (mean 2.3m translation, 53° rotation). This dataset is a re-rendering of a real capture, using the Habitat [56] system. The train/val/test set of the dataset consist of 32K/5K/8K image pairs, respectively. Following [22], we set the threshold for percent within a threshold to 30° for rotation and 1m for translation.

Baselines and Ablations. Our primary comparison is the Sparse Planes method of [22], a strong baseline estimating both rotation and translation (including scale). Sparse Planes does joint reconstruction and pose estimation and consists of: initial reconstruction and camera estimation, discrete optimization, and a bundle-adjustment on SIFT features [32] extracted from texture that has been made fronto-parallel. The final step adds substantial complexity, so we compare to (*SparsePlanes [22] No Bundle*) as well, which

Table 2. **Translation and Rotation Performance on Matterport.** Ours is best among methods producing translation scale. All baselines supervise depth except [22] (Camera Br) and Ours.

Method	Translation (m)			Rotation (degrees)		
	Med.↓	Avg.↓	≤1m↑	Med.↓	Avg.↓	≤30↑
[52] + [51]	3.34	4.00	8.3	50.98	57.92	29.9
Assoc.3D [48]	2.17	2.50	14.8	42.09	52.97	38.1
[22] (Camera Br)	0.90	1.40	55.5	7.65	24.57	81.9
[22] (No Bundle)	0.88	1.36	56.5	7.58	22.84	83.7
[22] (Full)	0.63	1.25	66.6	7.33	22.78	83.4
PlaneFormers [1]	0.66	1.19	66.8	5.96	22.20	83.8
Ours	0.64	1.01	67.4	8.01	19.13	85.4
SuperGlue [55]	-	-	-	3.88	24.17	77.8
LoFTR [60]	-	-	-	0.71	11.11	90.5

Table 3. **Essential Matrix Module Ablations on Matterport.** All three components of the Essential Matrix Module yield meaningful improvement across metrics.

Method	Translation (m)			Rotation (degrees)		
	Med.↓	Avg.↓	≤1m↑	Med.↓	Avg.↓	≤30↑
CNN Pose Regressor	1.53	1.83	28.6	31.31	45.05	48.8
+ViT	1.47	1.79	30.1	29.9	43.33	50.1
+Bilinear Attention	1.13	1.49	44.5	9.76	28.36	73.1
+Dual Softmax	0.70	1.06	64.8	8.62	21.23	83.3
Full	0.64	1.01	67.4	8.01	19.13	85.4

omits the final bundle adjustment, but still requires optimization. We also compare to the standalone pose estimation branch as (*Sparse Planes [22] Camera Branch*). In addition, we compare to concurrent work [1] which closely builds off of Sparse Planes.

We next report three baselines used by [22]. The first is (*Associative 3D [48] camera branch*), which is an improved version of RPNNet [12]. The second is the reconstruction-based RGBD odometry method of Raposo et al. [52] applied to [51]. Third, we compare with (*SuperGlue [55]*), using the settings from [22]. In addition, we compare to LoFTR [60]. Like SuperGlue, LoFTR supervises correspondences, and therefore requires depth supervision in addition to pose, and cannot recover translation scale.

Finally, we compare with four ablations that test the contributions of our method. All methods use the same MLP Regressor, and full descriptions of these appear in the supplement. The first is (*CNN Pose Regressor*), which predicts pose from concatenated base CNN extracted features. This gives a sense of how a simple method does. The second is (+*ViT*), which adds a ViT that is capped with standard attention (Eqn. 1) on top of the backbone. This tests the contribution of a ViT *without* the Essential Matrix Module. The third is (+*Bilinear*) which replaces standard attention (Eqn. 1) with bilinear attention, but without dual-softmax and quadratic positional encodings. Finally, we report (+*Dual Softmax*), which adds dual softmax.

Quantitative Results. We report results in Table 2. Joint prediction of rotation and translation (including scale) on wide-baseline pairs is a challenging problem. Non-trivial

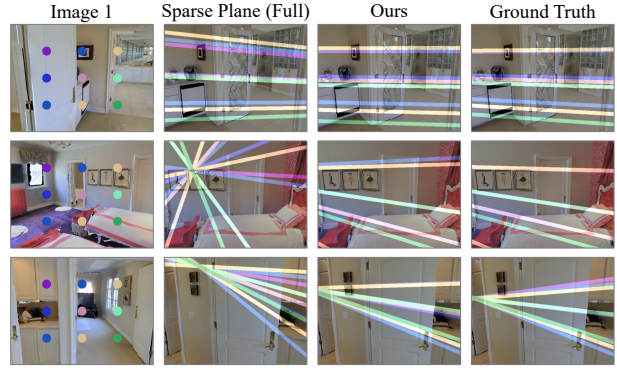


Figure 3. **Epipolar Lines on Matterport.** Our predictions better match true pose, particularly on large view changes.

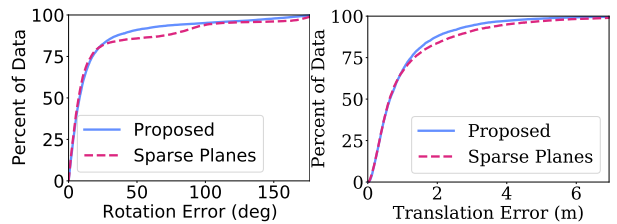


Figure 4. **Error CDFs on Matterport.** The proposed approach shows increased robustness to large view changes.

methods [48, 51, 52] have less than 40% of predictions within 30° of true rotation, and less than 20% of errors within 1m of true translation. Methods that are most competitive (LoFTR [60], SuperGlue [55], [22]) require depth supervision in addition to pose, while the best rotation results ([60], [55]) are produced by correspondence-based methods not predicting translation scale. Of methods producing translation scale, ours typically performs best, and it outperforms SuperGlue in both average rotation and percentage within 30°.

Ablations, shown in Table 3, show the reasons for success. As with [48, 51, 52], CNN and ViT models struggle at the task. Adding bilinear attention reduces errors tremendously – reducing median rotation error by two thirds – while the dual softmax reduces errors further significantly. Adding positional encodings further improves performance across all measurements.

Analysis. Qualitative results, in Figure 3, are consistent with quantitative findings. Namely, predicted pose more closely matches ground truth on difficult examples, resulting in much better mean performance than baselines. Error vs. view change is analyzed further in the Supplemental. Figure 4 displays error CDFs on Matterpor Compared to the most competitive baseline [22] (Full), the proposed method has fewer very large errors.

Table 4. **Rotation Performance on InteriorNet and StreetLearn.** We train and evaluate on only overlapping images. “*” indicates the method sometimes failed to produce pose estimation; errors were calculated only on successful image pairs. Gray text indicates failure over 50% of test pairs. The proposed method outperforms alternatives almost universally and often significantly.

Overlap	Method	InteriorNet			InteriorNet-T			StreetLearn			StreetLearn-T		
		Avg ($^{\circ}$ ↓)	Med. ($^{\circ}$ ↓)	10 (% ↑)	Avg ($^{\circ}$ ↓)	Med. ($^{\circ}$ ↓)	10 (% ↑)	Avg ($^{\circ}$ ↓)	Med. ($^{\circ}$ ↓)	10 (% ↑)	Avg ($^{\circ}$ ↓)	Med. ($^{\circ}$ ↓)	10 (% ↑)
Large	SIFT* [32]	6.09	4.00	84.86	7.78	2.95	55.52	5.84	3.16	91.18	18.86	3.13	22.37
	SuperPoint* [8]	5.40	3.53	87.10	5.46	2.79	65.97	6.23	3.61	91.18	6.38	1.79	16.45
	Reg6D [74]	5.43	3.87	87.10	10.45	6.91	67.76	3.36	2.71	97.65	12.31	6.02	69.08
	Cai <i>et al.</i> [5]	1.53	1.10	99.26	2.89	1.10	97.61	1.19	1.02	99.41	9.12	2.91	87.50
	Ours	0.48	0.40	100.00	2.90	1.83	97.91	0.62	0.52	100.00	4.08	2.43	90.13
Small	SIFT* [32]	24.18	8.57	39.73	18.16	10.01	18.52	16.22	7.35	55.81	38.78	13.81	5.68
	SuperPoint* [8]	16.72	8.43	21.58	11.61	5.82	11.73	19.29	7.60	24.58	6.80	6.85	0.95
	Reg6D [74]	17.83	9.61	51.37	21.87	11.43	44.14	7.95	4.34	87.71	15.07	7.59	63.41
	Cai <i>et al.</i> [5]	6.45	1.61	95.89	10.24	1.38	89.81	2.32	1.41	98.67	13.04	3.49	84.23
	Ours	1.81	0.94	99.32	4.48	2.38	96.30	1.46	1.09	100.00	9.19	3.25	87.70

Table 5. **Rotation Ablations InteriorNet and StreetLearn.** (Second best underlined). The ViT significantly improves over CNN only. Components of the proposed model perform in different settings, but the full model is often best and typically competitive with the best ablation, while ablations sometimes do poorly (Bilinear Att. on InteriorNet Small, ViT on StreetLearn Small).

Overlap	Method	InteriorNet-T			StreetLearn-T		
		Avg ($^{\circ}$ ↓)	Med. ($^{\circ}$ ↓)	10 (% ↑)	Avg ($^{\circ}$ ↓)	Med. ($^{\circ}$ ↓)	10 (% ↑)
Large	CNN Pose Regressor	5.29	2.6	89.85	15.25	10.00	50.00
	+ViT	<u>2.99</u>	1.64	96.72	3.52	2.56	94.74
	+Bilinear Attention	3.25	1.49	<u>97.31</u>	4.73	2.68	<u>92.76</u>
	+Dual Softmax	6.03	<u>1.63</u>	93.43	4.39	2.64	91.45
	Full	2.90	1.83	97.91	<u>4.08</u>	2.43	90.13
Small	CNN Pose Regressor	19.79	4.05	69.44	29.95	15.22	34.07
	+ViT	<u>5.43</u>	2.00	<u>94.75</u>	12.93	3.16	84.86
	+Bilinear Attention	8.54	1.79	90.43	8.70	3.41	89.59
	+Dual Softmax	10.44	<u>1.96</u>	89.51	10.74	<u>3.24</u>	87.07
	Full	4.48	2.38	96.30	<u>9.19</u>	3.25	<u>87.70</u>

4.3. Wide Baseline Rotation

We next study wide baseline rotation, where we compare with [5] and their baselines.

Datasets. We use the two datasets from [5], which were derived from panoramic photos and follow the setup of [5]. The first dataset is InteriorNet [29], which consists of 10,050 panoramic views across 112 synthetic houses. Of these, 82 houses are allocated for training and the remaining 30 houses are used for testing. The dataset has 610k image pairs (350K overlapping), with a test set of 1K pairs. StreetLearn [39], consists of panoramic outdoor images in New York City that have been scrubbed to ensure privacy (full details in Supplemental). This dataset has 1.1M train pairs (460K with overlap), and 1K test pairs from a set of 143K panoramas. We additionally evaluate on the “InteriorNet-T” and “StreetLearn-T” datasets, which select from *different* panoramas for each image in a pair, resulting in translation in addition to rotation. This translation is not, however, estimated in this setting. To facilitate comparisons we use 10° as a threshold for rotation error following [5]. We use the setup of Cai *et al.* [5] using only overlapping images, and breaking down overlap into *large* overlap (less than 45° rotation) and *small* overlap (more than 45°). Cai *et al.* also conduct experiments on non-overlapping images; we consider this beyond our scope, which is focused on the

case where correspondences may exist.

Baselines and Ablations. We compare to the state of the art (*Extreme Rotation* [5]), which computes a cross-correlation volume on paired image features, and uses a CNN to classify pose. We also report this method’s baselines: Reg6D [74], which predicts a 6D representation from concatenated image features, similar to the *Associative3D Camera Branch* from §4.2 as well as correspondence baselines SIFT [32] and SuperPoint [8]. These baselines occasionally fail. Following [5], we indicate failure on more than 50% of the test set by marking the number in gray. We report the same ablations as in §4.2.

Quantitative Results. The proposed method is typically better than all baselines across both InteriorNet and StreetLearn, for both versions and overlap settings of the dataset (Table 4). Often, the proposed method reduces error compared to competing methods by more than half (e.g., InteriorNet Mean, Median with Large Overlap; StreetLearn-T Mean with Large Overlap). Small overlap is an especially difficult setting. For instance, on InteriorNet-T, all baselines have mean error above 10° . Yet, the proposed method is within 10° more than 96% of the time. Interestingly, median error on InteriorNet-T is worse than Cai *et al.* [5]. We believe the large scale of InteriorNet is not the method’s strongest setting, and the method provides strong inductive bias for small data settings (see §4.4). Nevertheless, we consider Cai *et al.* to be a strong baseline as it is specialized to large angle changes.

Performance breakdown of the model is displayed in Table 5. Adding a ViT is quite important, likely attributable to the large scale of data available. Beyond the ViT, improvements by each step are more mixed compared to the clear improvement of each step on Matterport. For instance, adding the dual softmax without coordinate embeddings is typically not helpful compared to using only Bilinear Attention. Yet, the full model performs best (best 5 times, second best 3 times; Bilinear Attention is best 4 times, second best twice). Moreover, the full model is rarely significantly worse than any intermediate ablation. This suggests,

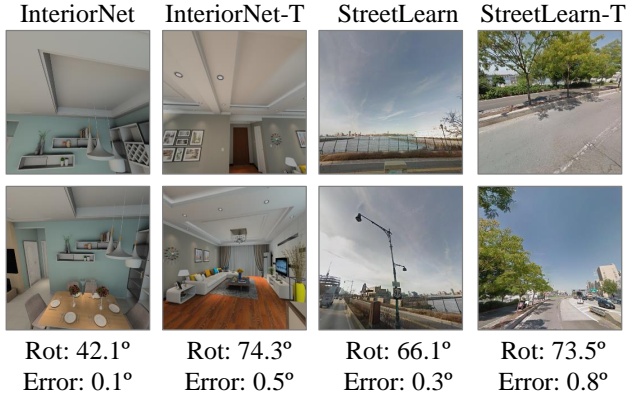


Figure 5. **Error vs. Rotation.** The proposed method produces high precision when faced with large view change.

as argued in §3, that all of the proposed components work together. We emphasize these settings have extraordinary numbers of views. §4.4 will show the substantially higher data efficiency of the essential module.

Analysis. Qualitative results validate quantitative findings in Fig. 5. While the evaluation datasets have huge rotations across indoor and outdoor settings, the proposed model is typically accurate, often even within 1% of true rotations.

4.4. Effectiveness on Smaller Datasets

One of the primary arguments for the use of the proposed network structure is that it provides a useful inductive bias by helping the network compute information that is known to constrain the set of feasible rotations and translations. In principle, since feedforward networks are universal approximators [21], networks ought to be able to learn to estimate relative pose with enough data. However, the right inductive biases ought to let them learn *faster*.

We now examine performance as a function of number of images. First, this helps empirically assess whether various networks structures provide useful inductive biases. Second, this is of practical concern since it tests data efficiency.

Datasets. We use InteriorNet-T and StreetLearn-T from §4.3, with significantly reduced 32K train image pairs. Collecting large-scale datasets such as these is challenging without a simulator or specialized company resources, so this smaller scale may be more realistic for e.g. user-collected posed images.

Ablations and Results. Our primary comparison is with the ViT baseline. Because it is a near alternative to our proposed Essential Matrix Module, we can measure the impact of our main contributions. Results are presented in Table 6, which is a reduced version of Table 5, with results also on the 32K image train set. Across datasets, the proposed method scales significantly better to a small train set. Even in cases the ViT slightly outperformed our proposed

Table 6. **Performance with limited data.** The proposed method scales better to small data than a typical learned model (e.g. ViT), indicating better inductive biases.

		InteriorNet-T					
		Full			32K		
		Avg	Med	% < 10°	Avg	Med	% < 10°
Large	ViT	0.61	0.49	100.00	5.78	3.23	92.84
	Full	0.48	0.40	100.00	4.44	2.58	95.82
Small	ViT	1.44	1.09	100.00	11.89	4.38	78.70
	Full	1.81	0.94	100.00	8.22	4.27	89.20

		StreetLearn-T					
		Full			32K		
		Avg	Med	% < 10°	Avg	Med	% < 10°
Large	ViT	3.52	2.56	94.74	11.51	7.69	56.58
	Full	4.08	2.43	90.13	7.22	4.44	81.58
Small	ViT	12.93	3.16	84.46	29.28	14.94	36.59
	Full	9.19	3.25	87.70	13.29	5.55	71.72

full model with full set, the inductive biases of the proposed method give it substantial improvement in the small setting.

5. Discussion

In this paper we presented a simple and interpretable end-to-end approach for pose estimation. Our key technical contribution is to implicitly represent correspondences from a ViT using an essential matrix module, from which an MLP can estimate pose. Theoretical results show this formulation can approximate the matrix $U^T U$ that is analyzed in the Eight Point algorithm; empirical results show given this, the MLP can suitably estimate pose. While alternatives make additional assumptions about input or require optimization, this method requires only paired RGB images as input, and is competitive in a variety of settings and view-point changes while being computationally efficient.

Limitations and Social Impact. The model is generally robust across view change. However, other methods are better suited for the two extremes in view change. In the case of small view change, the transformer is limited in terms of precision by the number of patches. Alternative CNN-based methods such as [65] may more easily operate upon high resolution. The model is also not prepared to predict pose on images with no overlap or correspondences; classification-based work e.g. [5] is better suited for this. Using datasets such as Matterport collected in nice homes leads to models which will likely perform better in these homes and possibly not as well in less expensive homes. Using synthetic data such as InteriorNet may help combat this bias. Training and evaluating on StreetView images should be handled with special care, as these images can contain personal information. The original authors blurred faces in the dataset, and a random manual search of 500 images also revealed no personal identifying information.

Acknowledgments. Thanks to Linyi Jin, Ruojin Cai and Zach Teed for help replicating and building upon their works. Thanks to Mohamed El Banani, Karan Desai and Nilesh Kulkarni for their helpful suggestions. Thanks to Laura Fink and UM DCO for their tireless computing support. Toyota Research Institute (“TRI”) provided funds to assist the authors with their research but this article solely reflects the opinions and conclusions of its authors and not TRI or any other Toyota entity.

A. Appendix

Our supplement presents the following.

Additional experimental details for the paper’s main results. These consist of detailed network architectures (§B), descriptions of datasets and data pre-processing (§C), and additional analysis that would not fit in the main paper (§D).

Discussion of the Essential Matrix Module. We present additional discussion and exposition of the Essential Matrix Module. This consists of a more detailed derivation of the fact that the unique entries of $\mathbf{U}^T\mathbf{U}$ can be written by $\Phi^T\mathbf{A}\Phi$ as described in the main paper (§E), how accurately the network can compute $\Phi^T\mathbf{A}\Phi$ in practice (§F), as well as an explicit writing-out of one of the key steps of this (§H). It also includes a detailed description of the experiment done with synthetic data which appears in the method section (§G).

B. Detailed Network Architectures

Our architecture is outlined below and detailed in the following two-page Table 7.

Backbone. Our backbone consists of a vanilla encoder of Residual Modules and one Block, followed by a vanilla ViT. The model ends with a three-layer MLP.

Essential Matrix Module. Recall from the paper, the Essential Matrix Module is closely built off a standard Cross-Attention Block. The changes we make are carefully chosen to study their contribution to performance. Changes can be seen in the “Essential Matrix Cross-Attn” block of the architecture. First, the final dimension of query, key and values is increased by a size of 6; positional encodings fill these 6 new spaces for value matrices. Second, Softmax is computed over both the last and second to last dimension of affinities, and multiplied elementwise. Finally, Bi-Directional attention is computed, $V.T@A@V$. Note this results in output of different shape than standard Cross-Attention.

Baselines. Baselines use generally the same components as our entire architecture, with very minor changes. This helps us study our proposed contributions.

CNN Pose Regressor. CNN Pose Regressor follows our model architecture, with the exception it only uses encoder, and one pooling layer, before the MLP. The pooling layer consists of two 1x1 Convs with Batchnorm (and ReLU before the second Conv). Feature size is reduced to 96, then 43; resulting in MLP input of $24768 = 43 * 24 * 24$. We do this so the input to the MLP is comparable to our method ($29400 = 2 * 3 * (64 + 6) * (64 + 6)$), and so this method is less prone to poor overfitting. Early experiments with bigger input size to MLP hurt performance. The CNN and MLP are otherwise identical to ours.

+ViT. Starting from the CNN Pose Regressor, we simply add the ViT Layers back from our architecture. This is followed by a vanilla Cross-Attention block. This uses the same pooling layer as the CNN. The Cross Attention-Block can be distinguished from our Essential Matrix Module by changing bidirectional attention, dual softmax and positional encoding. Looking in the table, the differences can be found in Essential Matrix Cross-Attn. First, affinities are calculated using the standard $(A@V).T$. Second, Softmax is not applied over $dim = -2$. Third, values do not get positional encodings.

+Bilinear Attention. This method uses our architecture as-is, with the exception of dual softmax and positional encoding. Looking in the table, the differences can be found in Essential Matrix Cross-Attn. First, Softmax is not applied over $dim = -2$. Second, values will not get positional encodings.

+Dual Softmax. This method uses our architecture as-is, with the exception of positional encoding. Looking in the table, the difference can be found in Essential Matrix Cross-Attn – values will no longer get positional encodings.

C. Dataset Information

Matterport3D. Matterport3D is a collection of scanned indoor scenes. We use the image pairs Jin *et al.* collected using the Habitat simulator. The number of image pairs in the train, val and test set are 31932, 4707, and 7996, respectively. Images are originally 480x640, but are downsampled to 256x256 for models. Images are collected using a camera at random height of 1.5-1.6m, with a downward tilt of 11 degrees to simulate human perspective. Candidate pairs are randomly sampled cameras within each room. Next, Jin *et al.* detect planes, and select pairs such that at least 3 planes are shared between images, and at least 3 planes are unique to each image. The average rotation is 53 degrees, translation 2.3m, and overlap 21%.

InteriorNet. InteriorNet consists of 10,050 indoor panoramic views across 112 synthetic houses. 82 houses are used for training and 30 are used for testing. We sample paired images from the panoramas using the same procedure as Cai *et al.* We also follow their image selection process, which samples images over a uniform distribution of angles (yaw in [-180, 180]; pitch in [-30,30]) within panoramas. Their process samples 100 images per panorama, filters images too close to walls, and does not apply roll; arguing this does not affect performance. Images are 256x256 and have 90° FOV. For full details see the original paper.

For the InteriorNet-T dataset, pairs are selected from different panoramas, resulting in translation; for InteriorNet, pairs are selected from the same panorama, resulting in no translation. Translations in InteriorNet-T are selected to be less than 3m. The full set of extracted image pairs on Inte-

Table 7. **Model Architecture.** Detailed model architecture, broken down into sub-components. Please note, some components have more complicated structure, so we define operations at the beginning and forward pass at the bottom. For instance, in the Residual Module we define the two branches, followed by the forward pass calling each branch. We use $H = 24$, $D = 192$, $N_h = 3$ in accordance with standard ViT-Tiny. We do not define ResNet Blocks below, as we use the standard implementation available publicly.

Overview	
Operation	Output Shape
Input Image	$2 \times 3 \times 256 \times 256$
Encoder	$2 \times H \times H \times D$
ViT Layer (x5)	$2 \times H \times H \times D$
Essential Matrix Module	$2 \times N_h \times (D/N_h + 6) \times (D/N_h + 6)$
MLP	7

Encoder	
Operation	Output Shape
ResNet-18 Block 1	$2 \times 56 \times 56 \times 64$
ResNet-18 Block 2	$2 \times H \times 28 \times 128$
Residual Module	$2 \times H \times H \times D$

MLP	
Operation	Output Shape
Linear & ReLU (x2)	512
Linear & ReLU	7

Residual Module	
Operation	Output Shape
<i>Branch A</i>	
2D Conv k=3 s=1, BN, ReLU	$2 \times 28 \times 28 \times D$
2D Conv k=5 s=1, BN, ReLU	$2 \times H \times H \times D$
<i>Branch B</i>	
2D Conv k=5 s=1, BN	$2 \times H \times H \times D$
<i>Forward Pass</i>	
ReLU(Branch A + Branch B)	$2 \times H \times H \times D$

riorNet is 250K (610K for InteriorNet-T). Note these train set sizes are smaller than reported in Cai *et al.* (roughly 1M and 700k) but were supplied by the authors directly and via their repo; these sets replicate their reported paper results. Both have a test set of 1K pairs. We consider only overlapping pairs, making the train set smaller: InteriorNet: 150K, InteriorNet-T: 350K. The test set sizes are also reduced after filtering for overlap. Pairs are further broken down into large or small overlap, thresholded using rotation of 45 degrees. InteriorNet test set: 695 pairs (403 large overlap, 292 small overlap). InteriorNet-T test set: 659 pairs (335 large overlap, 324 small overlap).

StreetLearn. StreetLearn consists of 143K panoramic outdoor views in New York City and Pittsburgh; we follow the setup of Cai *et al.* which uses 56K panoramas from Man-

hattan, with 1K randomly allocated for testing. We also follow their image selection process, which samples images over a uniform distribution of angles (yaw in $[-180, 180]$; pitch in $[-45, 45]$) within panoramas. Their process samples 100 images per panorama, and does not apply roll; arguing this does not affect performance. Images are 256×256 and have 90° FOV. For full details see the original paper. For the StreetLearn-T dataset, pairs are selected from different panoramas, resulting in translation; for StreetLearn, pairs are selected from the same panorama, resulting in no translation. Translations in StreetLearn-T are selected to be less than 10m. The full set of extracted image pairs on StreetLearn is 1.1M (670K for StreetLearn-T). Both have a test set of 1K pairs. We consider only overlapping pairs, making the train set smaller: StreetLearn: 460K, StreetLearn-T: 260K. The test set sizes are also reduced

Table 8. **Model Architecture: Essential Matrix Module.** Details of essential matrix module and ViT layer.

ViT Layer		
Operation		Output Shape
<i>Attn</i>		
Q,K,V = Linear		$2 \times N_h \times (H \times H) \times (D/N_h)$
A = Softmax(Q @ K.T, dim=-1)		$2 \times N_h \times (H \times H) \times (H \times H)$
(A @ V).T		$2 \times H \times H \times D$
Linear		$2 \times H \times H \times D$
<i>MLP</i>		
Linear & GeLU		$2 \times H \times H \times (D \times 4)$
Linear		$2 \times H \times H \times D$
<i>Forward Pass</i>		
Attn(LayerNorm) + Residual		$2 \times H \times H \times D$
MLP(LayerNorm) + Residual		$2 \times H \times H \times D$

Essential Matrix Module		
Operation		Output Shape
<i>MLP</i>		
Linear & GeLU		$2 \times H \times H \times (D \times 4)$
Linear		$2 \times H \times H \times D$
<i>Forward Pass</i>		
Essential Matrix Cross-Attn(LayerNorm) + Residual		$2 \times N_h \times (D/N_h + 6) \times (D/N_h + 6)$
MLP(LayerNorm) + Residual		$2 \times N_h \times (D/N_h + 6) \times (D/N_h + 6)$

Essential Matrix Cross-Attn		
Operation		Output Shape
<i>Attn Branch</i>		
A = Softmax(Q @ K.T, dim=-1) × Softmax(Q @ K.T, dim=-2)		$N_h \times (H \times H) \times (H \times H)$
(V.T @ A @ V).T		$N_h \times (D/N_h + 6) \times (D/N_h + 6)$
Linear		$N_h \times (D/N_h + 6) \times (D/N_h + 6)$
<i>Forward Pass</i>		
$Q_1, K_1, V_1 = \text{Linear}(\text{Input}[0])$		$N_h \times (H \times H) \times (D/N_h + 6)$
$Q_2, K_2, V_2 = \text{Linear}(\text{Input}[1])$		$N_h \times (H \times H) \times (D/N_h + 6)$
$V_1[\dots, -6 :], V_2[\dots, -6 :] = \text{Pos Encoding}$		$N_h \times (H \times H) \times (D/N_h + 6)$
Concat(Attn Branch(Q_1, K_2, V_2), Attn Branch(Q_2, K_1, V_1))		$2 \times N_h \times (D/N_h + 6) \times (D/N_h + 6)$

after filtering for overlap. Pairs are further broken down into large or small overlap, thresholded using rotation of 45 degrees. StreetLearn test set: 471 pairs (170 large overlap, 301 small overlap). StreetLearn-T test set: 469 pairs (152 large overlap, 317 small overlap).

StreetLearn panoramas were captured based on Google Street View, and therefore contain real people. The authors of the dataset blurred all faces and license plates, and manually reviewed images for privacy. The dataset is distributed only upon request. If individuals request for a panorama to

be taken down or blurred, the dataset is updated to reflect the request.

D. Additional Analysis

Error against GT. In Figures 6 and 7, we report error vs. magnitude for rotation and translation, respectively; for our method and baseline Sparse Planes. Lines are fit by applying adaptive kernel density estimation. For both rotation and translation, both methods show similar general trends, increasing error as magnitude increases. In the

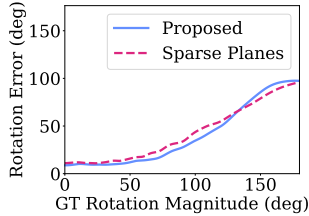


Figure 6. Mean Error as a Function of Magnitude: Rotation.

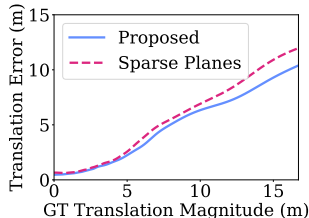


Figure 7. Mean Error as a Function of Magnitude: Translation.

Table 9. **Generalization Across Datasets**, training Cai *et al.* and Ours on the opposite of IN-T and SL-T. SuperGlue is trained on ScanNet; the only public version. In format “Large / Small” Overlap.

Method	InteriorNet-T			StreetLearn-T		
	Avg ($^{\circ}$ \downarrow)	Med. ($^{\circ}$ \downarrow)	10 (% \uparrow)	Avg ($^{\circ}$ \downarrow)	Med. ($^{\circ}$ \downarrow)	10 (% \uparrow)
SuperGlue	35.9 / 85.4	35.9 / 81.6	5.4 / 0.0	48.9 / 85.5	42.9 / 82.4	2.0 / 0.0
Cai <i>et al.</i>	31.2 / 67.4	9.3 / 84.4	52.0 / 23.8	43.0 / 71.4	21.8 / 72.5	32.4 / 9.9
Ours	18.7 / 58.6	11.1 / 66.4	45.1 / 10.5	28.6 / 40.7	14.3 / 24.0	39.5 / 27.8

case of very small rotations (i.e. $< 30^{\circ}$) or translations (i.e. $< 3m$), Sparse Planes is quite competitive with the proposed method. However, beyond this magnitude, our method tends to be more robust, outperforming until the very most extreme rotations. This is consistent with the qualitative results and CDFs plotted in the paper.

Generalization Across Datasets. A good measure of inductive bias is to evaluate methods when trained on one dataset and tested without fine-tuning on another. We evaluate our method against the most competitive baselines across datasets in Table 9. We tend to generalize better, especially in average error, though indoor \leftrightarrow outdoor does cause a large drop in performance.

Is ViT Backbone Needed? Our proposed method appends an Essential Matrix Module layer to the end of 5 ViT layers. In Table 10, we study the effect of replacing the 5 ViT layers with a comparable CNN-based backbone, a modified PWC-Net [59]. PWC-Net is a good choice as it predicts optical flow, a task closely related to correspondence estimation and therefore useful for downstream pose estimation; TartanVO [65] uses this extractor before predicting relative pose. We use PWC-Net out of the box except for modifications to feature and stride size so that it takes extracted features at the same resolution (24x24) and feature size (192) as the ViT, and produces output features at this same resolution and feature size. The CNN-based backbone performs

Table 10. **Model Backbone Ablations.** A ViT backbone is not necessary for competitive performance; a CNN replacement performs similarly if we keep the Essential Matrix Module.

Overlap	Method	InteriorNet-T			StreetLearn-T		
		Avg ($^{\circ}$ \downarrow)	Med. ($^{\circ}$ \downarrow)	10 (% \uparrow)	Avg ($^{\circ}$ \downarrow)	Med. ($^{\circ}$ \downarrow)	10 (% \uparrow)
Large	CNN	5.36	3.62	94.03	3.74	2.29	94.74
	ViT (Ours)	2.90	1.83	97.91	4.08	2.43	90.13
Small	CNN	8.31	4.58	86.42	7.30	3.03	90.54
	ViT (Ours)	4.48	2.38	96.30	9.19	3.25	87.70

similarly to the ViT backbone – better on StreetLearn-T but worse on InteriorNet-T. This indicates the Essential Matrix Module may be flexible for use with more generic models. Recall from Tables 3, 5 and 6 in the main paper, the ViT backbone is not the primary reason for our success – the Essential Matrix Module is a helpful inductive bias, particularly in the case of limited data.

Can one Construct the Essential Matrix Module from Conv Kernels? One can construct the essential matrix from $\Phi^T A \Phi$, which is similar to the attention operation performed in ViTs but not to any operation in Conv layers.

Can Essential Matrix Module Output be used for the 8-Point Algorithm? No, it is an inductive bias used for pose prediction. An analogy is convolution layers, which have the capacity to detect edges but which do not necessarily just detect edges when learned. Likewise, the EM Module has the capacity to represent information like $U^T U$, from which one could compute the Essential matrix. Our results show that this inductive bias helps learned pose estimation. We also note that some of the EM Module’s output (the upper $D \times D$ submatrix) are image features that cannot be used in the 8-Point algorithm.

How closely does EMM’s $\sim U^T U$ match true $U^T U$? We save EMM $\sim U^T U$ output for the Matterport test set and compare to true $U^T U$, computed by projecting 100k points in a 10x10x10 box into both images using ground truth pose. Note the EMM computes cross-attention in both $img1 \rightarrow img2$ and $img2 \rightarrow img1$, and has attention 3 heads for each direction. We therefore have $6 \sim U^T U$ to compare to the ground truth. We measure rank correlation between entries of $\sim U^T U$ and $U^T U$.

As we see in Figure 8 (total correlation values), there is a nontrivial rank correlation between EMM $\sim U^T U$ and the true $U^T U$. This correlation varies some on a per-head and per-direction basis, which is not surprising given directions and heads are concatenated before pose regression, and may be used for differing purposes and extents. Beyond concatenating heads, one important reason correlation is not higher is that $\sim U^T U$ is followed by linear projection and normalization layers before pose regression. Since there is no constraint on $\sim U^T U$ to match the true $U^T U$, the end-to-end system is instead encouraged to learn pose optimally, with this structure used as an inductive bias. Nevertheless, the positive correlation between $\sim U^T U$ and $U^T U$ is con-

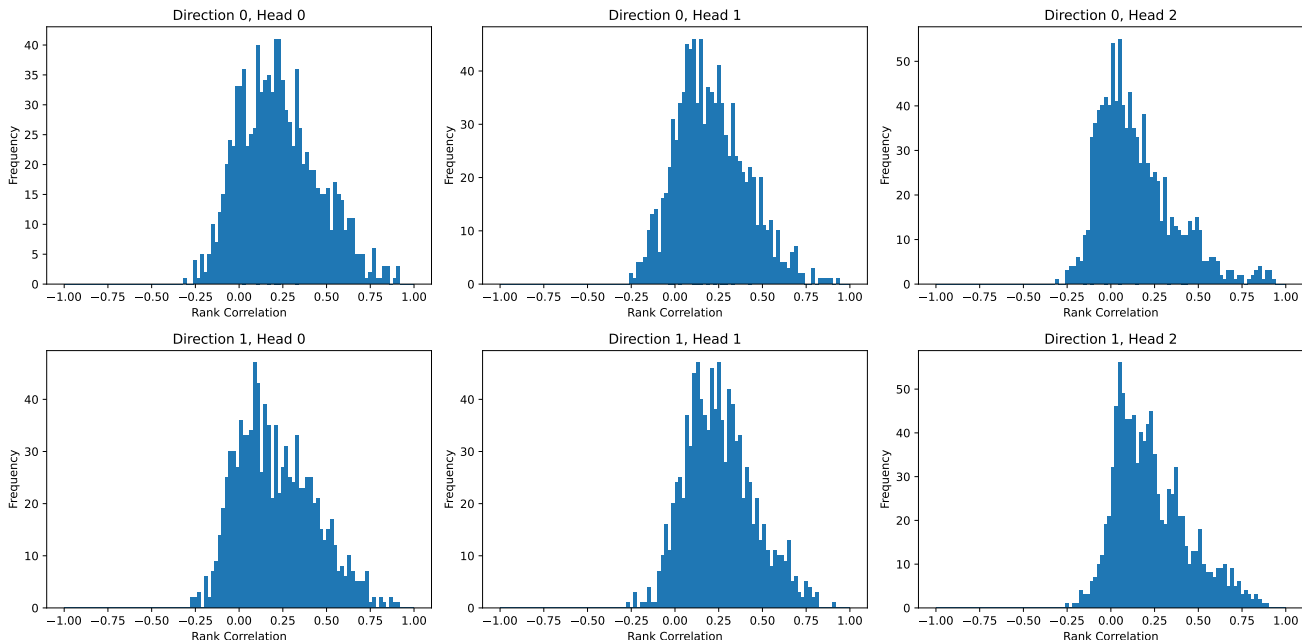


Figure 8. Distribution of rank correlations between ground-truth $\mathbf{U}^T \mathbf{U}$ and the $\mathbf{U}^T \mathbf{U}$ as computed by the transformer. We re-arrange the 6×6 bottom sub-matrix to be 9×9 and then compute the Spearman rank correlation per-image (which ranges from -1 to 1). We plot histograms of these rank correlations for each direction (i.e., image 1 or image 2) and each of the three transformer heads. Even though the matrices computed by the transformers are not used in the same way as in the 8-point algorithm, we find substantial rank correlation between the two.

sistent with our expectations the EMM is in fact computing similar positional structure to the eight point algorithm.

We visualize $\sim \mathbf{U}^T \mathbf{U}$ output of individual image pairs in Figure D. Here we see clear examples of the varying, but mostly positive and meaningful, correlation between actual $\sim \mathbf{U}^T \mathbf{U}$ and output from different heads and directions. In addition, comparison of individual outputs in the $9 \times 9 \sim \mathbf{U}^T \mathbf{U}$ tends to show similar cells with large activation across heads, regardless of actual correlation. We emphasize we should not expect $\sim \mathbf{U}^T \mathbf{U}$ to precisely match $\mathbf{U}^T \mathbf{U}$.

Why not use 8-Point Coefficients as Φ ? Using 8-point coefficients as Φ (Eqn. 13) actually leads to the same attention output entries as ours (Eqn. 11), only ours has the advantage of being more compact, as it doesn’t duplicate entries. This is detailed in Sec. H.

How is Scale Ambiguity Handled? The model is trained end-to-end using translation loss with scale, leaving the network free to learn to overcome the scale ambiguity via recognition. We believe this is due to a mix of a distribution over likely relative poses inside scenes as well as familiar objects. Essential Matrix Module output mixes learned appearance features \mathbf{V} (which can encode familiar objects), positional encodings Φ , and their interaction. Since Φ uses intrinsics (L476), features can model varying cameras.

Additional Qualitative Results: Matterport. In Figure 10, we see the proposed method generates Epipolar lines generally similar to true view changes. This is consistent with random rotation and translation errors in Figure 11, which are generally small. Not surprisingly, random results in this challenging setting are also sometimes poor, e.g. the top left Epipolar example, or the top left rotation error example. Sorting by error in Figure 12, we see the general trend of increased view change being associated with increased error, which is consistent with our Error against GT study.

Additional Qualitative Results: InteriorNet and StreetLearn. In Figure 13, we see random rotation and translation errors, which are generally quite small. Random results in these datasets are not often poor, but have some weaker examples – e.g. the bottom row of StreetLearn-T has most errors above 1° . Sorting by error in Figure 14, we again see the trend of increased view change being associated with increased error. Yet, results are often quite good even in very high rotations (e.g. InteriorNet all examples).

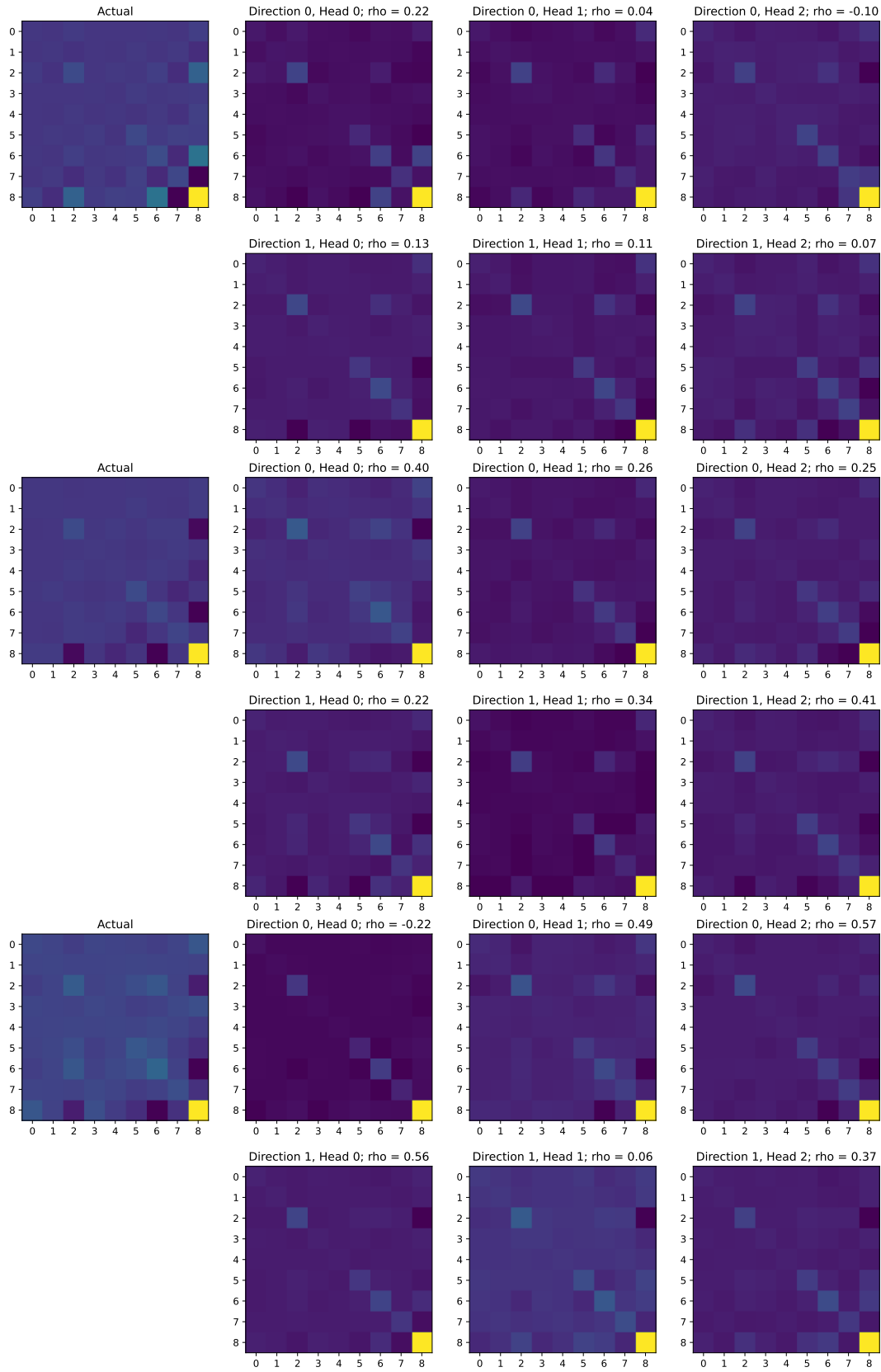


Figure 9. Selected examples of the actual and transformer $\mathbf{U}^\top \mathbf{U}$ matrices for three examples. Top left: actual $\mathbf{U}^\top \mathbf{U}$. Each row shows a different direction (i.e., image 1 or image 2) and each column shows a different transformer head (i.e., one of the three heads). The transformer matrix is actually 6×6 , but we re-arrange it to be the 9×9 $\mathbf{U}^\top \mathbf{U}$ matrix.

Random Results, Matterport



Figure 10. Random Results on Matterport: Epipolar Lines.

Random Results



Figure 11. Random Results on Matterport.

0 Percentile Error

Random Results by Error

99 Percentile Error

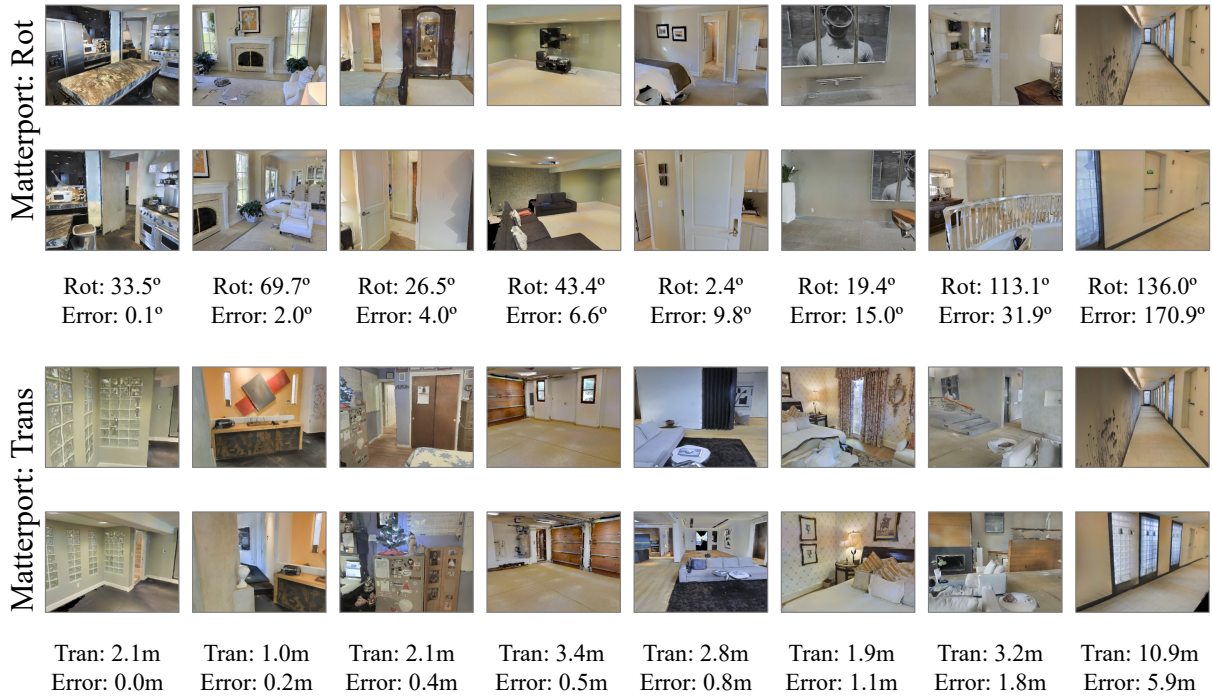


Figure 12. Results by error on Matterport.

Random Results

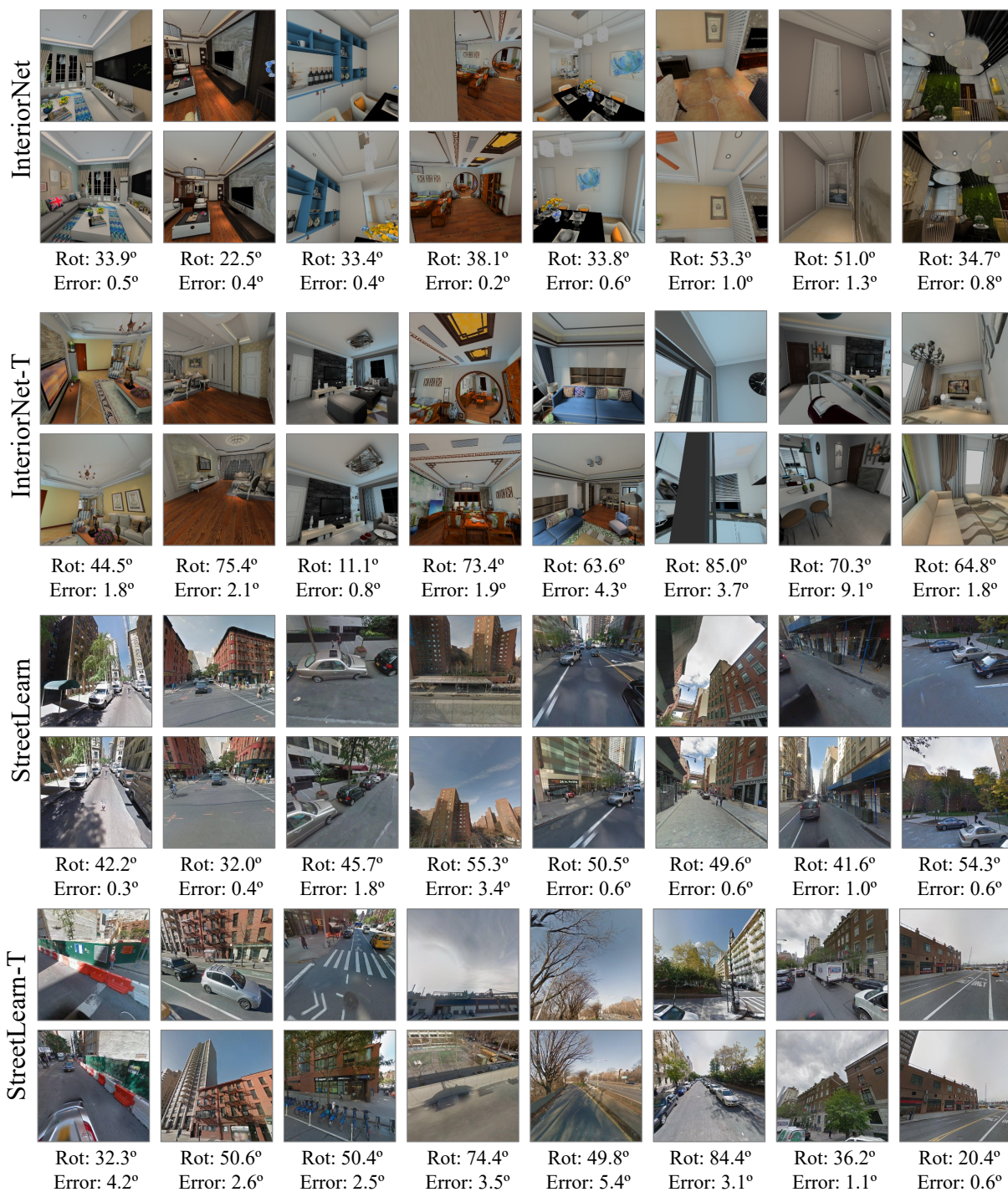


Figure 13. Random Results on InteriorNet and StreetLearn.

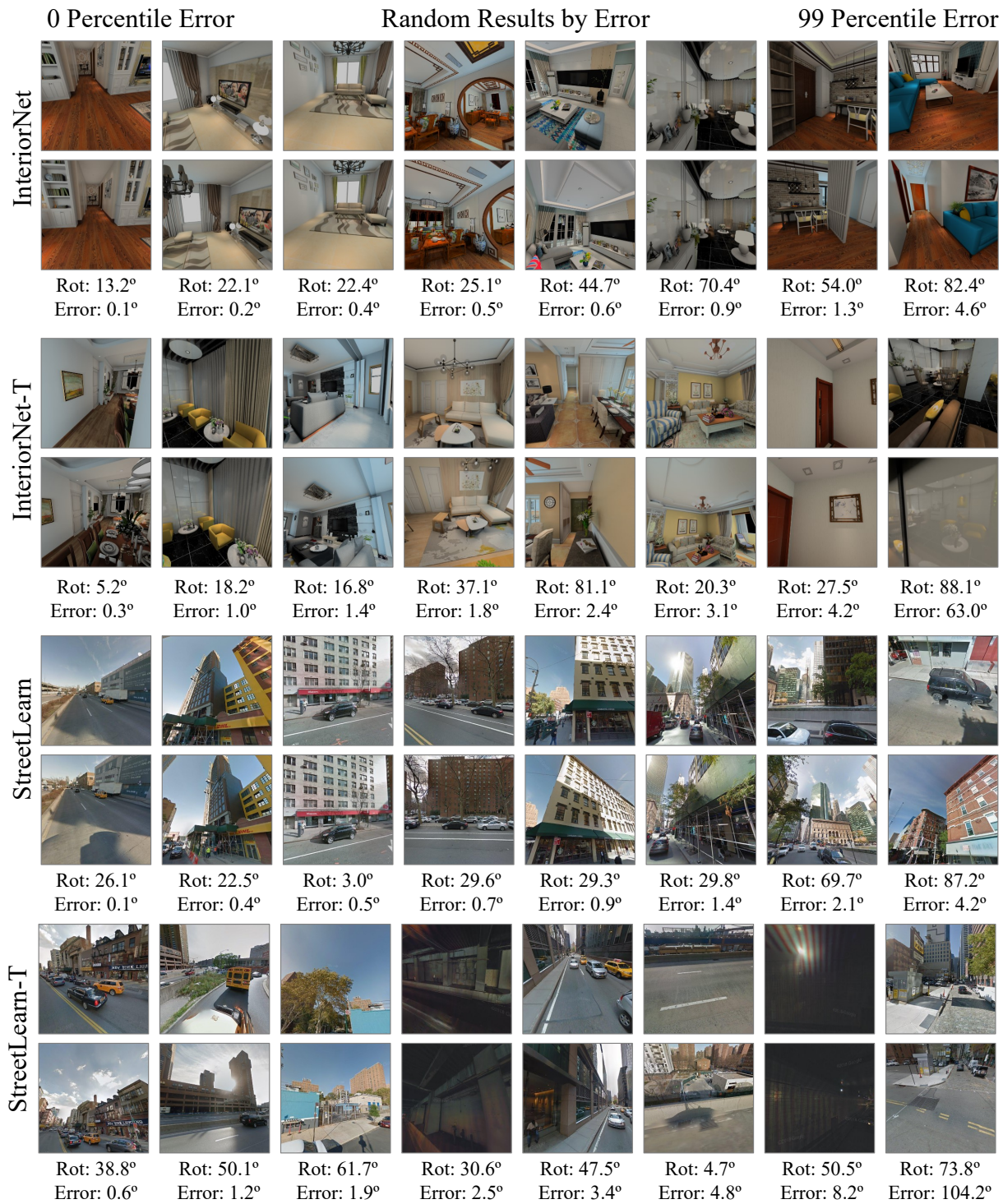


Figure 14. Results by error on InteriorNet and StreetLearn.

E. Derivation of the Unique Entries of $\mathbf{U}^\top \mathbf{U}$

Our goal is to show that the unique entries of $\mathbf{U}^\top \mathbf{U}$ that is used in the Eight Point Algorithm can be computed as $\Phi^\top \mathbf{A} \Phi$ for an attention matrix $\mathbf{A} \in \{0, 1\}^{P \times P}$ and $\Phi \in \mathbb{R}^{P \times 6}$ as defined in the main paper.

Setup. Given N correspondences, the eight-point algorithm constructs a matrix $\mathbf{U} \in \mathbb{R}^{N \times 9}$ row-wise via the Kronecker products of the homogeneous coordinates of the correspondences involved. Define $\mathbf{x}_i = [u_i, v_i, 1]$ and $\mathbf{x}'_i = [u'_i, v'_i, 1]$. Then the i th row of \mathbf{U} is

$$\mathbf{U}_{i,:} = \begin{bmatrix} u_i u'_i & u_i v'_i & u_i & v_i u'_i & v_i v'_i & v_i & u'_i & v'_i & 1 \end{bmatrix} \quad (4)$$

or more compactly,

$$\mathbf{U}_{i,:} = (\mathbf{x}_i \otimes \mathbf{x}'_i)^\top. \quad (5)$$

Note that when estimating the Essential matrix, one uses $\mathbf{x}_i \equiv \mathbf{K}^{-1}[u_i, v_i, 1]^\top$ and $\mathbf{x}'_i \equiv \mathbf{K}^{-1}[u'_i, v'_i, 1]^\top$. Since these coordinates can be rescaled by any arbitrary non-zero scalar, we assume that the last coordinate is 1.

Usual approach. Given correct correspondences, the eigenvector of $\mathbf{U}^\top \mathbf{U} \in \mathbb{R}^{9 \times 9}$ that corresponds to the smallest eigenvector is the Essential or Fundamental matrix. Usually, the matrix is not rank-deficient, and so one reshapes the eigenvector and then performs rank-reduction.

Alternate approach. We will now show that the unique entries of $\mathbf{U}^\top \mathbf{U}$ can be computed in an alternate fashion using a setup that is amenable to computation via a transformer. We'll start with the following basic substitutions and cleaning up:

$$\mathbf{U}^\top \mathbf{U} = \sum_{i=1}^N \mathbf{U}_{i,:}^\top \mathbf{U}_{i,:} = \sum_{i=1}^N (\mathbf{x}_i \otimes \mathbf{x}'_i) (\mathbf{x}_i \otimes \mathbf{x}'_i)^\top. \quad (6)$$

We'll first rewrite the interior of the sum, and then the sum itself.

Rewriting $\mathbf{U}_{i,:}^\top \mathbf{U}_{i,:}$ with a basis expansion. We'll tackle the interior of the sum first. While $\mathbf{U}_{i,:}^\top \mathbf{U}_{i,:} = (\mathbf{x}_i \otimes \mathbf{x}'_i) (\mathbf{x}_i \otimes \mathbf{x}'_i)^\top \in \mathbb{R}^{9 \times 9}$ and thus has 81 entries, there are only 36 unique entries. The smaller number of entries can be seen mechanically via direct expansion (see §H to see this). It can also be reasoned out by distributing transposes and using the mixed product property to rewrite it as

$$(\mathbf{x}_i \otimes \mathbf{x}'_i) (\mathbf{x}_i \otimes \mathbf{x}'_i)^\top = (\mathbf{x}_i \mathbf{x}_i^\top) \otimes (\mathbf{x}'_i \mathbf{x}'_i{}^\top). \quad (7)$$

Note that while $\mathbf{x}_i \mathbf{x}_i^\top$ has 9 entries, it only has 6 *unique* entries $(1, u, v, uv, u^2, v^2)$. Likewise, $\mathbf{x}'_i \mathbf{x}'_i{}^\top$ has 6 unique entries $(1, u', v', u'v', u'^2, v'^2)$. Therefore, their Kronecker product $(\mathbf{x}_i \mathbf{x}_i^\top) \otimes (\mathbf{x}'_i \mathbf{x}'_i{}^\top)$ has only 36 unique entries.

We can create a 6×6 matrix containing the unique entries of $\mathbf{U}_{i,:}^\top \mathbf{U}_{i,:}$ by applying a basis expansion to the coordinates. Let us define $\phi([u, v, 1]) = [1, u, v, uv, u^2, v^2]$.

Then the unique entries of $\mathbf{U}_{i,:}^\top \mathbf{U}_{i,:} \in \mathbb{R}^{9 \times 9}$ can be written as $\phi(\mathbf{x}_i) \phi(\mathbf{x}'_i)^\top \in \mathbb{R}^{6 \times 6}$. This means that the unique entries of $\mathbf{U}^\top \mathbf{U}$ are given by

$$\sum_{i=1}^N \phi(\mathbf{x}_i) \phi(\mathbf{x}'_i)^\top. \quad (8)$$

As an additional benefit, this factorization separates the terms involving each image into separate components.

Making the sum implicit. We next rewrite the sum implicitly by assuming each correspondence lies on a fixed grid. Given a grid of P patches in each image, we assume that \mathbf{p}_j is the j th patch's location. Then, rather than have N correspondences, we can define the correspondences implicitly via an indicator matrix $\mathbf{A} \in \{0, 1\}^{P \times P}$ such that $\mathbf{A}_{j,k} = 1$ if and only if points \mathbf{p}_k and \mathbf{p}_j are in correspondence and 0 otherwise. If each correspondence is on each patch, then we can rewrite

$$\sum_{i=1}^N \phi(\mathbf{x}_i) \phi(\mathbf{x}'_i)^\top = \sum_{j=1}^P \sum_{k=1}^P \phi(\mathbf{p}_k) \mathbf{A}_{j,k} \phi(\mathbf{p}'_j)^\top. \quad (9)$$

This can be further simplified by gathering the basis expanded coordinates of the grid in a matrix $\Phi \in \mathbb{R}^{P \times 6}$ such that $\Phi_{j,:} = \phi(\mathbf{p}_j)^\top$. Then $\phi(\mathbf{p}_k) \mathbf{A}_{j,k} \phi(\mathbf{p}'_j)^\top = \Phi_{k,:}^\top \mathbf{A}_{j,k} \Phi_{j,:}$, and so Equation 9 can be rewritten as

$$\sum_{j=1}^P \sum_{k=1}^P \Phi_{k,:}^\top \mathbf{A}_{j,k} \Phi_{j,:} = \Phi^\top \mathbf{A} \Phi, \quad (10)$$

and therefore the unique entries of $\mathbf{U}^\top \mathbf{U}$ can be compactly written as $\Phi^\top \mathbf{A} \Phi$.

F. Discussion of Limitations

The $\Phi^\top \mathbf{A} \Phi$ expression is exact when: (1) every correspondence can be represented as one of the P patches; and (2) the attention matrix \mathbf{A} produced by the ViT represents correspondence and is binary. Without an explicit binarization of attention and infinitely small patches, the Essential Matrix Module can, at best, compute an approximation. We now discuss how close this approximation can get.

The closeness of these approximation depends in part on the network architecture and field of view. Throughout, we use patches that are arrayed in a 24×24 grid.

Representing Each Correspondence as a Patch. We replace each correspondence with its equivalent patch, effectively quantizing the correspondence locations. With patches that are the size of pixels, this has close to no impact on the accuracy of estimating pose; if one represents the image with a handful of patches, this clearly ought to have a large impact on the accuracy. We now analyze the impact empirically.

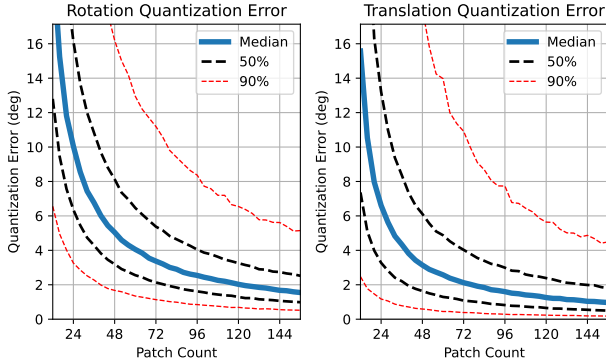


Figure 15. Correspondence quantization error as a function of the number of patches for rotation and translation. We use a patch count of 24×24 , which has a moderate quantization error. However, since transformers can contain sub-patch information implicitly, this quantization may be considerably lower.

Defining Quantization Error. For a given quantization level q (i.e., number of patches that uniformly divide the image along each axis), we generated 10,000 instances of synthetic correspondence by: sampling a relative camera pose with uniform Euler angles, and translation $\sim \text{Unif}(-1, 1)$, as well as a set of 3D points $\sim \text{Unif}(-1, 1)$. We render these points to the images using the Matterport3D intrinsics producing a set of correspondences $(x_i, y_i) \leftrightarrow (x'_i, y'_i)$. We then compute the relative pose two ways: first, we do this with the original correspondences, yielding \mathbf{R}_o and \mathbf{t}_o ; second, we do it with the correspondences uniformly quantized to q levels, which yields estimates \mathbf{R}_q and \mathbf{t}_q . We then define the quantization error as the rotation geodesic between \mathbf{R}_o and \mathbf{R}_q as well as the angle between \mathbf{t}_o and \mathbf{t}_q .

We then plot the median quantization error per quantization level, along with 50% and 90% intervals in Figure 15. The quantization error rapidly decreases as patch size increases. We use a patch count of 24×24 in this work, which corresponds to a moderate quantization error ($d(\mathbf{R}_o, \mathbf{R}_q) \approx 10^\circ$, $d(\mathbf{t}_o, \mathbf{t}_q) \approx 7^\circ$). Transformer tokens can represent sub-patch information, and once the patch count reaches 96×96 , the errors become quite small ($d(\mathbf{R}_o, \mathbf{R}_q) \approx 2.5^\circ$, $d(\mathbf{t}_o, \mathbf{t}_q) \approx 1.6^\circ$).

Producing \mathbf{A} . We now discuss how closely a transformer can its attention $\mathbf{A} = \text{norm}(\mathbf{QK}^\top)$ match the binarized matrix that our setup uses. We divide this into two cases: patches that have correspondence and patches without correspondence. We refer to the total contribution as the total size of the weights for a patch j , or $\sum_{k=1}^P \mathbf{A}_{j,k}$.

Patches with correspondence. If patch j has a correspondence with patch k , then we would like $\mathbf{A}_{j,k} = 1$ and $\mathbf{A}_{j,k'} = 0$ for all $k' \neq k$. Standard attention cannot exactly reach this, but can get arbitrarily close by making its dot product $(\mathbf{QK}^\top)_{j,k}$ as high as possible. Thus the total

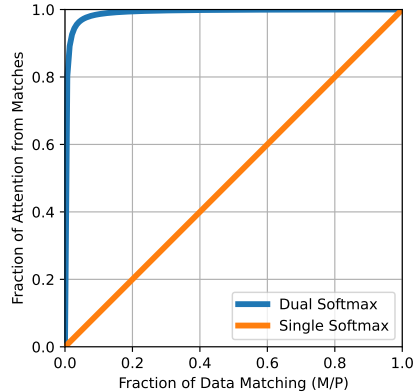


Figure 16. Fraction of attention contributed by matches as a fraction of prevalence of matches. Dual softmax enables matches to rapidly dominate the attention matrix’s entries.

contribution of a matching patch $j \sum_{k=1}^P \mathbf{A}_{j,k} \approx 1$.

Patches without correspondence. If patch j has no correspondence, then we would like $\mathbf{A}_{j,k'} = 0$ for all k' . This is impossible under standard attention. We can get this to be as close to 0 as possible, by the following: make $(\mathbf{QK}^\top)_{j,:}$ and $(\mathbf{QK}^\top)_{:,j}$ all equal, which in turn makes the resulting softmax distributions uniform. In turn this makes $\text{softmax}(\mathbf{QK}^\top)_{j,k'} = \text{softmax}(\mathbf{QK}^\top)_{k',j} = 1/P$. Under dual-softmax, then $\mathbf{A}_{j,k'} = 1/P^2$ for all k' . Thus, the total contribution of a patch j is $\sum_{k'=1}^P 1/P^2 = 1/P$.

Together, this means that with dual softmax, the vast majority of the attention matrix’s energy comes from matches. We do a simple experiment, assuming that $(\mathbf{QK}^\top)_{j,k} = 100$ if patches j and k match and 1 otherwise. We can quantify the fraction of the attention that is from matches by examining the fraction of the resulting matrix $\mathbf{A} = \text{norm}(\mathbf{QK}^\top)$ that corresponds to matches. We plot this as a function of the prevalence of matches in Fig. 16, comparing regular and dual softmax. With dual softmax, a handful of matches dominate the attention, whereas for regular softmax, attention increases linearly.

G. Synthetic Experimental Details

Datasets. For each dataset, we generate an instance consisting of a scene and relative camera pose. These can be used to derive features that are suitable for training.

Scenes: Each scene consists of a points drawn uniformly inside a sphere with center \mathbf{c} with its coordinates independently and identically sampled from $\text{Unif}(-\frac{1}{2}, \frac{1}{2})$ and radius $r \sim \text{Unif}(\frac{1}{2}, \frac{3}{2})$. This sampling is done with rejection sampling. *Relative Camera Pose:* We generate Euler angles $(\theta_x, \theta_y, \theta_z)$ for the three axes, and a translation vector \mathbf{t} . In all cases, we reject samples with $\|\mathbf{t}\| \leq \frac{1}{2}$.

1. **3D (General 3D Motion):** $\theta_x, \theta_y, \theta_z \sim \text{Unif}(0, 360^\circ)$;

$\mathbf{t}_x, \mathbf{t}_y, \mathbf{t}_z \sim \text{Unif}(-1, 1)$.

2. **2D Large** (Large-Rotation Motion In XZ plane): $\theta_y \sim \text{Normal}(0, 25^\circ)$, $\theta_x, \theta_z \sim \text{Normal}(0, 1.25^\circ)$; $\mathbf{t}_x, \mathbf{t}_z \sim \text{Normal}(0, \frac{1}{3})$, $\mathbf{t}_y \sim \text{Normal}(0, \frac{1}{60})$.
3. **2D Medium** (Medium-Rotation Motion In XZ plane): $\theta_y \sim \text{Normal}(0, 5^\circ)$, $\theta_x, \theta_z \sim \text{Normal}(0, 0.25^\circ)$; $\mathbf{t}_x, \mathbf{t}_z \sim \text{Normal}(0, \frac{1}{3})$, $\mathbf{t}_y \sim \text{Normal}(0, \frac{1}{60})$.
4. **2D Small** (Small-Rotation Motion In XZ plane): $\theta_y \sim \text{Normal}(0, 1^\circ)$, $\theta_x, \theta_z \sim \text{Normal}(0, 0.05^\circ)$; $\mathbf{t}_x, \mathbf{t}_z \sim \text{Normal}(0, \frac{1}{3})$, $\mathbf{t}_y \sim \text{Normal}(0, \frac{1}{60})$.

Given a scene and relative camera pose, we project the points onto a virtual camera with height and width 800 units, a focal length of 800 units and principal point of 400 units. We record the point if it is in front of the camera, and on the virtual sensor. We reject the image pair and scene if fewer than 100 points out of 10K random points are valid for both cameras.

Input Feature. Given a 3D point, we denote its projection into image 1 as \mathbf{x} and its projection into image 2 as \mathbf{x}' . Given the set of valid points, we compute the explicit form of $\mathbf{U}^\top \mathbf{U}$ with two modifications. First, for numerical stability we divide each coordinate by the width of the image and subtract $1/2$, which centers the data. Second, to make the feature independent of the number of correspondences, we normalize by the number of points to obtain $\frac{1}{N} \mathbf{U}^\top \mathbf{U}$ rather than $\mathbf{U}^\top \mathbf{U}$. These are identical from the perspective of eigenvectors, but normalizing makes the feature independent of the number of points.

Method. For each task, we train a multilayer perceptron consisting of 3 hidden layers with 4096 units each. Each hidden layer is capped with a leaky ReLU. We predict a normalized vector (3D for translation direction, 4D for rotation).

H. Verifying that $\phi(\mathbf{x})\phi(\mathbf{x}')^\top$ contains all the terms needed for $\mathbf{U}^\top \mathbf{U}$

To enable visually verifying Equation 8, we'll show a visual expansion. To avoid notational clutter, we will drop the i th index and deal with $\mathbf{x} = [u, v, 1]$ and $\mathbf{x}' = [u', v', 1]$. Our goal is to show that $(\mathbf{x} \otimes \mathbf{x}')^\top (\mathbf{x} \otimes \mathbf{x}')$ has the same entries as $\phi(\mathbf{x})\phi(\mathbf{x}')^\top$. We'll first expand out $\phi(\mathbf{x})\phi(\mathbf{x}')^\top$. When factored out,

$$\phi(\mathbf{x})\phi(\mathbf{x}')^\top = \begin{bmatrix} 1 & u' & v' & u'v' & u'^2 & v'^2 \\ u & uu' & uv' & uu'v' & uu'^2 & uv'^2 \\ v & vu' & vv' & vu'v' & vu'^2 & vv'^2 \\ uv & uvu' & uvv' & uvu'v' & uvu'^2 & uvv'^2 \\ u^2 & u^2u' & u^2v' & u^2u'v' & u^2u'^2 & u^2v'^2 \\ v^2 & v^2u' & v^2v' & v^2u'v' & v^2u'^2 & v^2v'^2 \end{bmatrix}. \quad (11)$$

We color the terms according to which row they appear in $\phi(\mathbf{x})\phi(\mathbf{x}')^\top$. The matrix created for each correspondence is $(\mathbf{x} \otimes \mathbf{x}')^\top (\mathbf{x} \otimes \mathbf{x}')$. We'll first define $(\mathbf{x} \otimes \mathbf{x}')$:

$$(\mathbf{x} \otimes \mathbf{x}') = [uu' \quad uv' \quad u \quad vu' \quad vv' \quad v \quad u' \quad v' \quad 1]. \quad (12)$$

We can then compute the outer product $(\mathbf{x} \otimes \mathbf{x}')^\top (\mathbf{x} \otimes \mathbf{x}')$. This is highly redundant – note that the i th row and i th column are identical. More specifically,

$$(\mathbf{x} \otimes \mathbf{x}')^\top (\mathbf{x} \otimes \mathbf{x}') = \begin{bmatrix} u^2u'^2 & u^2u'v' & u^2u' & uvu'^2 & uvu'v' & uvu' & uu'^2 & uu'v' & uu' \\ u^2u'v' & u^2v'^2 & u^2v' & uvu'v' & uvv'^2 & uvv' & uu'v' & uv'^2 & uv' \\ u^2u' & u^2v' & u^2 & uvu' & uvv' & uv & uu' & uv' & u \\ uvu'^2 & uvu'v' & uvu' & v^2u'^2 & v^2u'v' & v^2u' & vu'^2 & vu'v' & vu' \\ uvu'v' & uvv'^2 & uvv' & v^2u'v' & v^2v'^2 & v^2v' & vu'v' & vv'^2 & vv' \\ uvu' & uvv' & uv & v^2u' & v^2v' & v^2 & vu' & vv' & v \\ uu'^2 & uu'v' & uu' & vu'^2 & vu'v' & vu' & u'^2 & u'v' & u' \\ uu'v' & uv'^2 & uv' & vu'v' & vv'^2 & vv' & u'v' & v'^2 & v' \\ uu' & uv' & u & vu' & vv' & v & u' & v' & 1 \end{bmatrix}. \quad (13)$$

Note that the rows of $\phi(\mathbf{x})\phi(\mathbf{x}')^\top$ appear in 3×3 blocks inside $(\mathbf{x} \otimes \mathbf{x}')^\top (\mathbf{x} \otimes \mathbf{x}')$. In particular, inside each 3×3 block, the columns of the $\phi(\mathbf{x})\phi(\mathbf{x}')^\top$ appear in the following order:

$$\begin{bmatrix} 5 & 4 & 2 \\ 4 & 6 & 3 \\ 2 & 3 & 1 \end{bmatrix}. \quad (14)$$

References

- [1] Samir Agarwala, Linyi Jin, Chris Rockwell, and David F. Fouhey. Planeformers: From sparse view planes to 3d reconstruction. In *ECCV*, 2022. 6
- [2] Herbert Bay, Tinne Tuytelaars, and Luc Van Gool. Surf: Speeded up robust features. In *ECCV*, 2006. 1, 2
- [3] Eric Brachmann, Alexander Krull, Sebastian Nowozin, Jamie Shotton, Frank Michel, Stefan Gumhold, and Carsten Rother. Dsac-differentiable ransac for camera localization. In *CVPR*, 2017. 2
- [4] Samarth Brahmhatt, Jinwei Gu, Kihwan Kim, James Hays, and Jan Kautz. Geometry-aware learning of maps for camera localization. In *CVPR*, 2018. 2
- [5] Ruojin Cai, Bharath Hariharan, Noah Snavely, and Hadar Averbuch-Elor. Extreme rotation estimation using dense correlation volumes. In *CVPR*, 2021. 1, 2, 7, 8
- [6] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. In *3DV*, 2017. 2, 5
- [7] Kefan Chen, Noah Snavely, and Ameesh Makadia. Wide-baseline relative camera pose estimation with directional learning. In *CVPR*, 2021. 1, 2
- [8] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. In *CVPRW*, 2018. 1, 2, 7
- [9] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021. 1, 2, 3
- [10] Mihai Dusmanu, Ignacio Rocco, Tomas Pajdla, Marc Pollefeys, Josef Sivic, Akihiko Torii, and Torsten Sattler. D2-Net: A Trainable CNN for Joint Detection and Description of Local Features. In *CVPR*, 2019. 1, 2
- [11] Mohamed El Banani, Luya Gao, and Justin Johnson. Unsuperviseddr: Unsupervised point cloud registration via differentiable rendering. In *CVPR*, 2021. 1, 2
- [12] Sovann En, Alexis Lechervy, and Frédéric Jurie. RpNet: An end-to-end network for relative camera pose estimation. In *ECCVW*, 2018. 1, 2, 6
- [13] Olivier D Faugeras. What can be seen in three dimensions with an uncalibrated stereo rig? In *ECCV*, pages 563–578. Springer, 1992. 2
- [14] Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 1981. 1
- [15] Banglei Guan, Qifeng Yu, and Friedrich Fraundorfer. Minimal solutions for the rotational alignment of imu-camera systems using homography constraints. *Computer vision and image understanding*, 170:79–91, 2018. 1, 2
- [16] Christopher Ham, Simon Lucey, and Surya Singh. Hand waving away scale. In *ECCV*, 2014. 1, 2
- [17] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003. 1, 2
- [18] Richard I Hartley. Estimation of relative camera positions for uncalibrated cameras. In *ECCV*. Springer, 1992. 2
- [19] Richard I Hartley. In defense of the eight-point algorithm. *TPAMPI*, 19(6):580–593, 1997. 1, 2
- [20] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 3
- [21] Kurt Hornik. Approximation capabilities of multilayer feed-forward networks. *Neural networks*, 4(2):251–257, 1991. 3, 8
- [22] Linyi Jin, Shengyi Qian, Andrew Owens, and David F Fouhey. Planar surface reconstruction from sparse views. In *ICCV*, 2021. 1, 2, 5, 6
- [23] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *ACM Computing Surveys (CSUR)*, 2021. 2
- [24] Jin-Hwa Kim, Jaehyun Jun, and Byoung-Tak Zhang. Bilinear attention networks. *NeurIPS*, 31, 2018. 1, 2, 3
- [25] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015. 5
- [26] Laurent Kneip, Margarita Chli, and Roland Siegwart. Robust real-time visual odometry with a single camera and an imu. In *BMVC*, 2011. 1, 2
- [27] Viktor Larsson, Magnus Oskarsson, Kalle Astrom, Alge Wallis, Zuzana Kukelova, and Tomas Pajdla. Beyond grobner bases: Basis selection for minimal solvers. In *CVPR*, 2018. 2
- [28] Zakaria Laskar, Iaroslav Melekhov, Surya Kalia, and Juho Kannala. Camera relocalization by computing pairwise relative poses using convolutional neural network. In *ICCVW*, 2017. 2
- [29] Wenbin Li, Sajad Saeedi, John McCormac, Ronald Clark, Dimos Tzoumanikas, Qing Ye, Yuzhong Huang, Rui Tang, and Stefan Leutenegger. Interiornet: Mega-scale multi-sensor photo-realistic indoor scenes dataset. In *BVMC*, 2018. 2, 7
- [30] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *ICCV*, 2021. 1
- [31] H Christopher Longuet-Higgins. A computer algorithm for reconstructing a scene from two projections. *Nature*, 293(5828):133–135, 1981. 1, 2
- [32] David G Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*. 1, 2, 5, 7
- [33] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *NeurIPS*, 32, 2019. 5
- [34] Xuan Luo, Jia-Bin Huang, Richard Szeliski, Kevin Matzen, and Johannes Kopf. Consistent video depth estimation. In *SIGGRAPH*, 2020. 1
- [35] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *CVPR*, 2021. 1

- [36] Jiri Matas, Ondrej Chum, Martin Urban, and Tomás Pa-jdla. Robust wide-baseline stereo from maximally stable extremal regions. *Image and vision computing*, 22(10):761–767, 2004. 2
- [37] Iaroslav Melekhov, Juha Ylioinas, Juho Kannala, and Esa Rahtu. Relative camera pose estimation using convolutional neural networks. In *ACIVS*, 2017. 2
- [38] 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*. Springer, 2020. 1
- [39] Piotr Mirowski, Andras Banki-Horvath, Keith Anderson, Denis Teplyashin, Karl Moritz Hermann, Mateusz Malinowski, Matthew Koichi Grimes, Karen Simonyan, Koray Kavukcuoglu, Andrew Zisserman, et al. The streetlearn environment and dataset. *arXiv preprint arXiv:1903.01292*, 2019. 2, 7
- [40] Dmytro Mishkin, Jiri Matas, Michal Perdoch, and Karel Lenc. Wxbs: Wide baseline stereo generalizations. *BMVC*, 2015. 2
- [41] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. Orb-slam: a versatile and accurate monocular slam system. *T-RO*. 1, 2
- [42] David Nistér. An efficient solution to the five-point relative pose problem. *TPAMI*, 26(6):756–770, 2004. 1, 2
- [43] Namuk Park and Songkuk Kim. How do vision transformers work? In *ICLR*, 2021. 2
- [44] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *NeurIPS*, 32, 2019. 5
- [45] Omid Poursaeed, Guandao Yang, Aditya Prakash, Qiuren Fang, Hanqing Jiang, Bharath Hariharan, and Serge Belongie. Deep fundamental matrix estimation without correspondences. In *ECCVW*, 2018. 2
- [46] Phil Pritchett and Andrew Zisserman. Wide baseline stereo matching. In *ICCV*, 1998. 2
- [47] Thomas Probst, Danda Pani Paudel, Ajad Chhatkuli, and Luc Van Gool. Unsupervised learning of consensus maximization for 3d vision problems. In *CVPR*, 2019. 2
- [48] Shengyi Qian, Linyi Jin, and David F Fouhey. Associative3d: Volumetric reconstruction from sparse views. In *ECCV*, 2020. 1, 2, 6
- [49] Noha Radwan, Abhinav Valada, and Wolfram Burgard. Vlocnet++: Deep multitask learning for semantic visual localization and odometry. *RA-L*, 2018. 2
- [50] René Ranftl and Vladlen Koltun. Deep fundamental matrix estimation. In *ECCV*, 2018. 2
- [51] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *TPAMI*, 2020. 6
- [52] Carolina Raposo, Miguel Lourenço, Michel Antunes, and Joao Pedro Barreto. Plane-based odometry using an rgb-d camera. In *BMVC*. 6
- [53] I. Rocco, R. Arandjelović, and J. Sivic. End-to-end weakly-supervised semantic alignment. In *CVPR*, 2018. 1, 2, 4
- [54] Ethan Rublee, Vincent Rabaud, Kurt Konolige, and Gary Bradski. Orb: An efficient alternative to sift or surf. In *ICCV*, 2011. 1, 2
- [55] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superglue: Learning feature matching with graph neural networks. In *CVPR*, 2020. 1, 2, 6
- [56] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A platform for embodied ai research. In *ICCV*, 2019. 5
- [57] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 1, 2
- [58] Leslie N Smith and Nicholay Topin. Super-convergence: Very fast training of neural networks using large learning rates. In *Artificial intelligence and machine learning for multi-domain operations applications*, 2019. 5
- [59] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *CVPR*, 2018. 13
- [60] Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. Loftr: Detector-free local feature matching with transformers. In *CVPR*, 2021. 1, 2, 4, 6
- [61] Zachary Teed and Jia Deng. Droid-slam: Deep visual slam for monocular, stereo, and rgb-d cameras. *NeurIPS*, 2021. 2
- [62] Zachary Teed and Jia Deng. Tangent space backpropagation for 3d transformation groups. In *CVPR*, 2021. 5
- [63] Michał Tyszkiewicz, Pascal Fua, and Eduard Trulls. Disk: Learning local features with policy gradient. *NeurIPS*, 33, 2020. 1, 2, 4
- [64] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 2017. 1, 2, 3
- [65] Wenshan Wang, Yaoyu Hu, and Sebastian Scherer. Tartanvo: A generalizable learning-based vo. In *CoRL*, 2021. 2, 5, 8, 13
- [66] 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 *IROS*, 2020. 2
- [67] Ross Wightman. Pytorch image models. <https://github.com/rwightman/pytorch-image-models>, 2019. 5
- [68] Fei Xue, Xin Wang, Zike Yan, Qiuyuan Wang, Junqiu Wang, and Hongbin Zha. Local supports global: Deep camera relocalization with sequence enhancement. In *ICCV*, 2019. 2
- [69] Zhenpei Yang, Jeffrey Z Pan, Linjie Luo, Xiaowei Zhou, Kristen Grauman, and Qixing Huang. Extreme relative pose estimation for rgb-d scans via scene completion. In *CVPR*, 2019. 2
- [70] Zhenpei Yang, Siming Yan, and Qixing Huang. Extreme relative pose network under hybrid representations. In *CVPR*, 2020. 1, 2

- [71] Wang Yifan, Carl Doersch, Relja Arandjelović, João Carreira, and Andrew Zisserman. Input-level inductive biases for 3d reconstruction. *arXiv preprint arXiv:2112.03243*, 2021. [2](#)
- [72] Qunjie Zhou, Torsten Sattler, Marc Pollefeys, and Laura Leal-Taixe. To learn or not to learn: Visual localization from essential matrices. In *ICRA*, 2020. [2](#)
- [73] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. In *SIGGRAPH*, 2018. [1](#)
- [74] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In *CVPR*, 2019. [7](#)