

Tangent Images for Mitigating Spherical Distortion

Marc Eder, Mykhailo Shvets, John Lim, and Jan-Michael Frahm
 University of North Carolina at Chapel Hill
 Chapel Hill, NC
 {meder, mshvets, jlim13, jmf}@cs.unc.edu

Abstract

In this work, we propose “tangent images,” a spherical image representation that facilitates transferable and scalable 360° computer vision. Inspired by techniques in cartography and computer graphics, we render a spherical image to a set of distortion-mitigated, locally-planar image grids tangent to a subdivided icosahedron. By varying the resolution of these grids independently of the subdivision level, we can effectively represent high resolution spherical images while still benefiting from the low-distortion icosahedral spherical approximation. We show that training standard convolutional neural networks on tangent images compares favorably to the many specialized spherical convolutional kernels that have been developed, while also scaling efficiently to handle significantly higher spherical resolutions. Furthermore, because our approach does not require specialized kernels, we show that we can transfer networks trained on perspective images to spherical data without fine-tuning and with limited performance drop-off. Finally, we demonstrate that tangent images can be used to improve the quality of sparse feature detection on spherical images, illustrating its usefulness for traditional computer vision tasks like structure-from-motion and SLAM.

Authors Note: This version of this paper has been updated with new network transfer results not present in the version published at CVPR 2020. These important new results demonstrate that network transfer to spherical images using our representation provides **equivalent performance** to perspective image networks after only a single epoch of fine-tuning and actually improves performance after 10 epochs of fine-tuning. The accompanying code has also been updated to coincide with this version of the paper. A full log of the changes is provided in Section 6.

1. Introduction

A number of methods have been proposed to address convolutions on spherical images. These techniques vary in design, encompassing learnable transformations [27, 28],

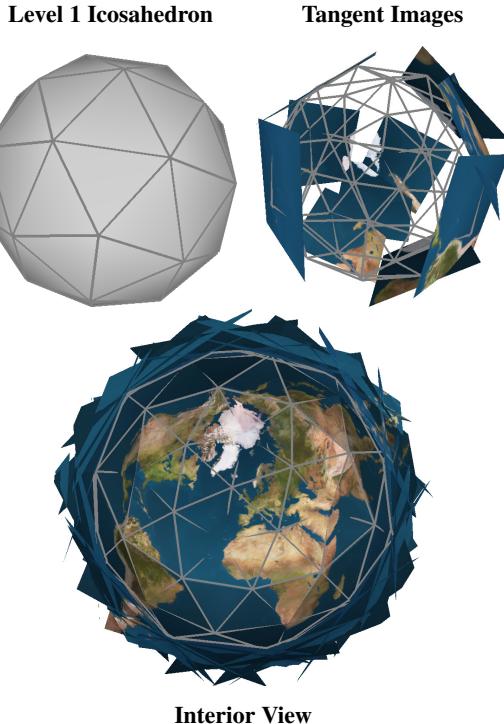


Figure 1: Using tangent images to represent a 4k Earth image [14]. **TL:** A base level 1 icosahedron. **TR:** Selection of tangent images rendered from the Earth image. **B:** Interior view of the tangent image spherical approximation.

generalizations and modifications of the convolution operation [8, 9, 11, 29], and specialized kernels for spherical representations [7, 17, 31]. In general, these spherical convolutions fall into two classes: those that operate on equirectangular projections and those that operate on a subdivided icosahedral representation of the sphere. The latter has been shown to significantly mitigate spherical distortion, which leads to significant improvements for dense prediction tasks [10, 11, 20]. It also has the useful property that icosahedron’s faces and vertices scale roughly by a factor of 4 at each subdivision, permitting a simple analogy to 2× up-sampling and down-sampling operations in standard convolutional neural networks (CNNs). Because of the perfor-

mance improvements provided by the subdivided icosahedron representation, we focus expressly on it in this paper.

Despite a growing body of work on these icosahedral convolutions, there are two significant impediments to further development: (1) the transferability of standard CNNs to spherical data on the icosahedron, and (2) the difficulty in scaling the proposed spherical convolution operations to high resolution spherical images. Prior work has implied [7, 11] or demonstrated [9, 29, 31] the transferability of networks trained on perspective images to different spherical representations. However, those who report results see a noticeable decrease in accuracy compared to CNN performance on perspective images and specialized networks that are trained natively on spherical data, leaving this important and desired behavior an unresolved question. Additionally, the proposed specialized convolutional kernels either require subsequent network tuning [7, 31] or are incompatible with the standard convolution [17].

Nearly all prior work on icosahedral convolutions has been built on the analogy between pixels and faces [7, 20] or pixels and vertices [11, 17, 31]. While elegant on the surface, this parallel has led to difficulties in scaling to higher resolution spherical images. Figure 2 depicts spherical image resolutions evaluated in the prior work. Notice that the highest resolution obtained so far is a level 8 subdivision, which is comparable to a 512×1024 equirectangular image. Superficially, this pixel resolution seems reasonably high, but the angular resolution per pixel is still quite low. A 512×1024 equirectangular image has an angular resolution of 0.352° . For comparison, a VGA resolution (480×640) perspective image with $45^\circ \times 60^\circ$ field of view (FOV) has an angular resolution of 0.094° . This is most similar to a 2048×4096 equirectangular image, which has an angular resolution of 0.088° and corresponds to a level 10 subdivided icosahedron. As this is a significantly higher resolution than prior work has been capable of demonstrating, this is the resolution on which we test our proposed approach.

In this work, we aim to address both transferability and scalability while leveraging efficient implementations of existing network architectures and operations. To this end, we propose a solution that decouples resolution from subdivision level using oriented, distortion-mitigated images that can be filtered with the standard grid convolution operation. Using these *tangent images*, standard CNN performance is competitive with specialized networks, yet they efficiently scale to high resolution spherical data and open the door to performance-preserving network transfer between perspective and spherical data. Furthermore, use of the standard convolution operation allows us to leverage highly-optimized convolution implementations, such as those from the cuDNN library [5], to train our networks. Additionally, the benefits of tangent images are not restricted to deep learning, as they address distortion through the data repre-

sentation rather than the data processing tools. This means that our approach can be used for traditional vision applications like structure-from-motion and SLAM as well.

We summarize our contributions as follows:

- We propose the tangent image spherical representation: a set of oriented, low-distortion images rendered tangent to faces of the icosahedron.
- We show that standard CNNs trained on tangent images perform competitively with specialized spherical convolutional kernels while also scaling effectively to high resolution spherical images.
- We demonstrate that tangent images facilitate network transfer between perspective and spherical images with no fine tuning and minimal performance drop-off.
- We illustrate the utility of tangent images for traditional computer vision tasks by using them to improve sparse keypoint matching on spherical images.

2. Related Work

Recently, there have been a number of efforts to close the gap between CNN performance on perspective images and spherical images. These efforts can be naturally divided based on the spherical image representation used.

2.1. Equirectangular images

Equirectangular images are a popular spherical image representation thanks to their simple relation between rectangular and spherical coordinates. However, they demonstrate severe image distortion as a result. A number of methods have been proposed to address this issue. Su and Grauman [27] develop a learnable, adaptive kernel to train a CNN to transfer models trained on perspective images to the equirectangular domain. Su *et al.* [28] extend this idea by developing a kernel that learns to transform a feature map according to local distortion properties. Cohen *et al.* [8, 6] develop spherical convolutions, which provides the rotational equivariance necessary for convolutions on the sphere. This method requires a specialized kernel, however, making it difficult to transfer the insights developed from years of research into traditional CNNs. Works from Coors *et al.* [9] and Tateno *et al.* [29] address equirectangular image distortion by warping the planar convolution kernel in a location-dependent manner. Because the equirectangular representation is so highly distorted, most recent work on this topic, has looked to leverage the distorted-reducing properties of the icosahedral spherical approximation.

2.2. Icosahedral representations

Representing the spherical image as a subdivided icosahedron mitigates spherical distortion, thus improving CNN accuracy compared to techniques that operate on equirectangular images. Eder and Frahm [10] motivate this representation using analysis from the field of cartography. Fur-

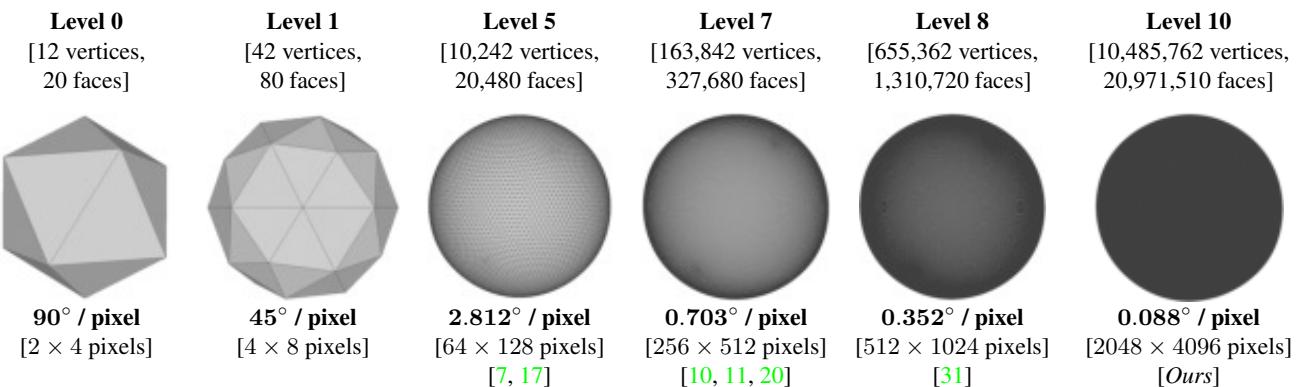


Figure 2: Demonstrating the number of elements, corresponding equirectangular image dimensions, and angular pixel resolution at various icosahedral subdivision levels. The citations beneath each denote the maximum resolution examined in those respective papers. Except for ours, they have all been limited by the pixel-to-face or pixel-to-vertex analogy.

ther research on this representation has primarily focused on the development of novel kernel designs to handle discretization and orientation challenges on the icosahedral manifold. Lee *et al.* [20] convolve on this representation by defining new, orientation-dependent, kernels to sample from triangular faces of the icosahedron. Jiang *et al.* [17] reparameterize the convolutional kernel as a linear combination of differential operators on the surface of an icosahedral mesh. Zhang *et al.* [31] present a method that applies a special hexagonal convolution on the icosahedral net. Cohen *et al.* [7] precompute an atlas of charts at different orientations that cover the icosahedral grid and use masked kernels along with an feature-orienting transform to convolve on these planar representations. Eder *et al.* [11] define the “mapped convolution” that allows the custom specification of convolution sampling patterns through a type of graph convolution. In this way, they specify the filters’ orientation and sample from the icosahedral surface. Our tangent image representation addresses data orientation by ensuring all tangent images are consistently oriented when rendering and circumvents the discretization issue by rendering to image pixel grids.

3. Mitigating Spherical Distortion

Image distortion is the reason that we cannot simply apply many state-of-the-art CNNs to spherical data. Distortion changes the representation of the image, resulting in local content deformation that violates translational equivariance, the key property of a signal required for convolution functionality. The graph in Figure 3 shows just how little distortion is required to produce a significant drop-off in CNN performance. Distortion in the most popular spherical image representations, equirectangular images and cube maps, is quite significant [10], and hence results in even worse performance. Although we can typically remove most lens distortion in perspective images using tools like the Brown-Conrady distortion model [2], spherical dis-

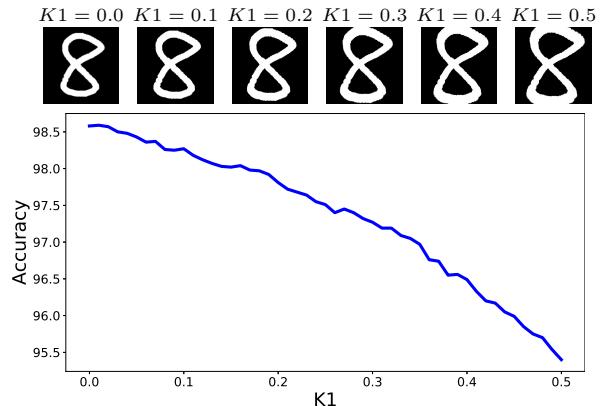


Figure 3: MNIST classification accuracy decreases as pincushion distortion is added to test images by varying the $K1$ parameter of the Brown-Conrady radial distortion model [2]. An example digit is shown at different distortion levels.

tortion is inescapable and is actually a function of the planar representation. This follows from Gauss’s Theorema Egregium, a consequence of which is that a spherical surface is not isometric to a plane. As such, any effort to represent a spherical image as a planar one will result in some degree of distortion. Thus, our objective, and one shared by cartographers for thousands of years, is limited to finding the optimal planar representation of the sphere for our use case.

3.1. The icosahedral sphere

Consider the classical *method of exhaustion* of approximating a circle with inscribed regular polygons. It follows that, in three dimensions, we can approximate a sphere in the same way. Thus, the choice of planar spherical approximation ought to be the convex Platonic solid with the most faces: the icosahedron. The icosahedron has been used by cartographers to represent Earth at least as early as Buckminster Fuller’s Dymaxion map [3], which projects the globe onto the icosahedral net. Recent work in computer

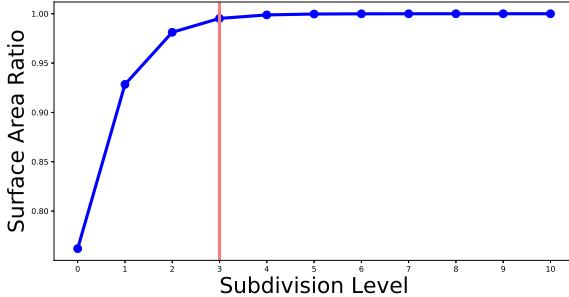


Figure 4: Ratio of the surface area of the subdivided icosahedron to the surface area of a sphere of the same radius at each subdivision level. This global metric demonstrates how closely the subdivision surface approximates a sphere and is drawn from established cartographic metrics [18]. Note the leveling off after the third subdivision level.

vision [7, 10, 11, 20, 17, 31] has demonstrated the shape’s utility for resolving the distortion problem for CNNs on spherical images as well.

While an improvement over single-plane image projections and its Platonic solid cousin, the cube, the 20-face icosahedron on its own is still limited in its distortion-mitigating properties. It can be improved by repeatedly applying Loop subdivision [23] to subdivide the faces and interpolate the vertices, producing increasingly close spherical approximations with decreasing amounts of local distortion on each face. Figure 4 demonstrates how distortion decreases at each subdivision level. Not all prior work takes advantage of this extra distortion reduction, though. There has largely been a trade-off between efficiency and representation. The charts used by Cohen *et al.* [7] and the net used by Zhang *et al.* [31] are efficient thanks to their planar image representations, but they are limited to the distortion properties of a level 0 icosahedron. On the other hand, the mapped convolution proposed by Eder *et al.* [11] operates on the mesh itself and thus can benefit from higher level subdivision, but it does not scale well to higher level meshes due to cache coherence problems when computing intermediate features on the mesh. Jiang *et al.* [17] provide efficient performance on the mesh, but do so by approximating convolution with a differential operator, which means existing networks can not be transferred. It is also interesting to note that the current top-performing method for many deep learning tasks, [31], uses the net of the level 0 icosahedron. This suggests that extensive subdivisions may not be necessary for all use cases.

Practical methods for processing spherical images must address the efficient scalability problem, but also should permit the transfer of well-researched, high-performance methods designed for perspective images. They should also provide the opportunity to modulate the level of acceptable distortion depending on the application. To address these constraints, we propose to break the coupling of subdivi-

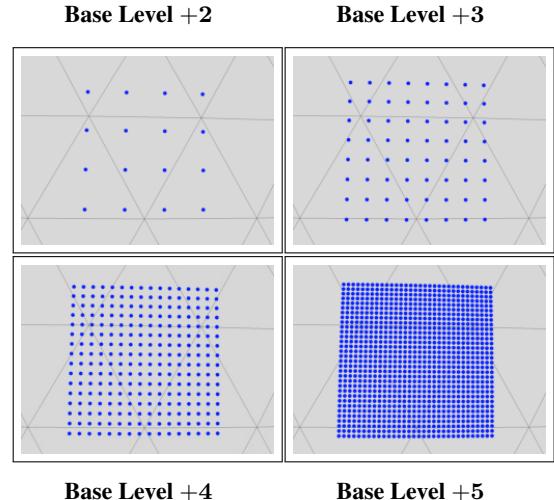


Figure 5: Illustrating how the tangent image resolution increases without changing the underlying subdivision level. The field-of-view of the tangent pixel grid remains unchanged, but its resolution increases by a factor of 2 in each dimension, demonstrated by the blue dots representing pixel samples on the sphere. This scaling maintains the angular pixel resolution of higher level icosahedrons without the need for additional subdivisions.

sion level and spherical image resolution by representing a spherical image as a collection of images with tunable resolution and distortion characteristics.

3.2. Tangent images

Subdividing the icosahedron provides diminishing returns rather quickly from a distortion-reduction perspective, as indicated by the red vertical line in Figure 4. Nonetheless, existing methods must continue to subdivide in order to match the spherical image resolution to the number of mesh elements. We untether these considerations by fixing a *base level* of subdivision, b , to define an acceptable degree of distortion, and then rendering the spherical image to square, oriented, planar pixel grids tangent to each face at that base level. The resolution of these *tangent images* is subsequently determined by the resolution of the spherical input. Given a subdivision level, s , corresponding to the spherical input resolution, the dimension of the tangent image, d , is given by the relation:

$$d = 2^{s-b} \quad (1)$$

This design preserves the same resolution scaling that would occur through further subdivisions by instead increasing the resolution of the tangent image. This relationship is illustrated in Figure 5.

Our tangent images are motivated by existing techniques in related fields. The approximation of sections of the sphere by low-distortion planar regions is similar to the Universal Transverse Mercator (UTM) geodetic coordinate

system, which divides the Earth into a number of nearly-Euclidean zones. Additionally, as tangent images can be thought of as rendering a spherical mesh to a set of quad textures, the high resolution benefits are similar to Ptex [4], a computer graphics technique that enables efficient high-resolution texturing by providing every quad of a 3D mesh with its own texture map. A visualization of the tangent image concept is provided in Figure 1.

Computing tangent images Tangent images are the **gnomonic projection** of the spherical data onto oriented, square planes centered at each face of a level b subdivided icosahedron. The number of tangent images, N , is determined by the faces of the base level icosahedron: $N = 20(4^b)$, while their spatial extent is a function of the vertex resolution, $R_v(b-1)$, of the level $b-1$ icosahedron and the resolution of the image grid, given by Equation (1). Let (ϕ_f, λ_f) be the barycenter of a triangular face of the icosahedron in spherical coordinates. We then compute the bounds of the plane in spherical coordinates as the inverse gnomonic projection at central latitude and longitude (ϕ_f, λ_f) of the points:

$$\left\{ \phi_f \pm \frac{d-1}{2d} R_v(b-1) \right\} \times \left\{ \lambda_f \pm \frac{d-1}{2d} R_v(b-1) \right\} \quad (2)$$

The vertex resolution, R_v , of a level b icosahedron, $\mathcal{S}(b)$, is computed as the mean angle between all vertices, v , and their neighbors, $\text{adj}(v)$:

$$R_v(b) = \frac{1}{|\mathcal{S}(b)|} \sum_{v \in \mathcal{S}(b)} \sum_{w \in \text{adj}(v)} \frac{\angle(v, w)}{|\text{adj}(v)|} \quad (3)$$

Using $R_v(b-1)$ ensures that the tangent images completely cover their associated triangular faces. Because vertex resolution roughly halves at each subsequent subdivision level, we define $R_v(-1) = 2R_v(0)$.

Using tangent images Tangent images require rendering from and to the sphere only once each. First, we create the tangent image set by rendering to the planes defined by Equation (2). Then, we apply the desired perspective image algorithm (e.g. a CNN or keypoint detector). Finally, we compute the regions on each plane visible to a spherical camera at the center of the icosahedron and render the algorithm output back to the sphere.

We have released our tangent image rendering code and associated experiments as a PyTorch extension¹.

4. Experiments

Prior research has established a common suite of experiments that have become the test bed for new research on spherical convolutions. This set typically includes some combination of spherical MNIST classification

Method	Filter	Acc.
Cohen <i>et al.</i> [8]	Spherical Correlation	85.0%
Esteves <i>et al.</i> [12]	Spectral Parameterization	88.9%
Jiang <i>et al.</i> [17]	MeshConv	90.5%
Ours	2D Convolution	89.1%

Table 1: Classification results on the ModelNet40 dataset [30]. Without any specialized convolution operations, our approach is competitive with the state of the art spherical convolution methods.

[8, 7, 17, 20, 31], shape classification [8, 12, 17], climate pattern segmentation [7, 17, 31], and semantic segmentation [7, 17, 20, 29, 31]. In order to benchmark against these prior works, we evaluate our method on the shape classification and semantic segmentation tasks. Additionally, we demonstrate our method’s fairly seamless transfer of CNNs trained on perspective images to spherical data. Finally, to show the versatility of the tangent image representation, we introduce a new benchmark, sparse keypoint detection on spherical images, and compare our representation to an equirectangular image baseline.

4.1. Classification

We first evaluate our proposed method on the shape classification task. As with prior work, we use the ModelNet40 dataset [30] rendered using the method described by Cohen *et al.* [8]. Because the data densely encompasses the entire sphere, unlike spherical MNIST, which is sparse and projected only on one hemisphere, we believe this task is more indicative of general classification performance.

Experimental setup We use the network architecture from Jiang *et al.* [17], but we replace the specialized kernels with simple 3×3 2D convolutions. A forward pass involves running the convolutional blocks on each patch separately and subsequently aggregating the patch features with average pooling. We train and test on level 5 resolution data as with the prior work.

Results and analysis Results of our experiments are shown in Table 1. Without any specialized convolutional kernels, we outperform most of the prior work on this task. The best performing method from Jiang *et al.* [17] leverages a specialized convolution approximation on the mesh, which inhibits the ability to fine-tune existing CNN models for the task. Our method can be thought of as using a traditional CNN in a multi-view approach to spherical images. This means that, for global inference tasks like classification, we could select our favorite pre-trained network and transfer it to spherical data. In this case, it is likely that some fine-tuning may be necessary to address the final patch aggregation step in our network design.

4.2. Semantic segmentation

We next consider the task of semantic segmentation in order to demonstrate dense prediction capabilities. To compare to prior work, we perform a baseline evaluation of our

¹<https://github.com/meder411/Tangent-Images>

Stanford2D3DS Dataset					
s	Method	Input	b	mAcc	mIOU
5	Cohen <i>et al.</i> [7]	RGB-D	0	55.9	39.4
	Jiang <i>et al.</i> [17]	RGB-D	5	54.7	38.3
	Zhang <i>et al.</i> [31]	RGB-D	0	58.6	43.3
	Ours	RGB-D	0	50.9	38.3
7	Tateno <i>et al.</i> [29]	RGB	ERP	-	34.6
	Lee <i>et al.</i> [20]	RGB	7	51.4	-
	Ours	RGB-D	0	59.1	44.9
10	Ours	RGB	0	61.0	44.3
	Ours	RGB	1	65.2	45.6
	Ours	RGB	2	61.5	42.7
	Ours	RGB-D	1	70.1	52.5

Table 2: Semantic segmentation results. s is the input resolution in terms of equivalent icosahedron level, b is the base subdivision level (ERP denotes equirectangular inputs), mIoU is the mean intersection-over-union metric, and mAcc is the weighted per-class mean prediction accuracy.

method at low icosahedron resolutions (5 and 7), but we also evaluate the performance of our method at a level 10 input resolution in order to demonstrate the usefulness of the tangent image representation for processing high resolution spherical data. No prior work has operated at this resolution. We hope that our work can serve as a benchmark for further research on high resolution spherical images.

Experimental setup We train and test our method on the Stanford 2D3DS dataset [1], as with prior work [8, 7, 17, 31]. We evaluate RGB-D inputs at levels 5, 7, and 10, the maximum resolution provided by the dataset. At level 10 we also evaluate using only RGB inputs to demonstrate the benefit of high resolution capabilities. For the level 5 and 7 experiments, we use the residual UNet-style architecture as in [17, 31], but we again replace the specialized kernels with 3×3 convolutions. The higher resolution of the level 10 inputs requires the larger receptive field of a deeper network, so we use a FCN-ResNet 101 [15, 22] model pre-trained on COCO [21] for those experiments. For level 5 data, we train on the entire set of tangent images, while for the higher resolution experiments, we randomly sample a subset of tangent images from each spherical input to expedite training. We found this sampling method to be useful without loss of accuracy. We liken it to training on multiple perspective views of a scene. Forward passes are run on all sampled tangent images before each backward pass. In this way, the computed gradients at every iteration come from the entire span of the spherical image’s field of view.

Results and analysis We report the results of our experiments in Table 2. Results on the Stanford2D3DS dataset are averaged over the 3 folds. Individual class results can be found in the supplementary material (Section 9). As expected, our method does not perform as well as prior work at the level 5 resolution. Recall that a level 5 resolution spherical image is equivalent to a 16×16 perspective image with 45° FOV. Our method takes that already low angular

resolution image and separates it into a set of low pixel resolution images. Although it had limited impact on classification, these dual low resolutions are problematic for dense prediction tasks. We expound on the low-resolution limitation further in the supplementary material (Section 7).

Where our tangent image representation excels is when scaling to high resolution images. What we sacrifice in low-resolution performance, we make up for by efficiently scaling to high resolution inputs. By scaling to the full resolution of the dataset, we are able to report the highest performing results ever on this spherical dataset by a wide margin using only RGB inputs. Adding the extra depth channel, we are able to increase the performance further (+4.9 mAcc, +6.9 mIOU). At input level 10, we find that base level 1 delivers the best trade-off between the lower FOV at higher base levels and the increased distortion present in lower ones. We elaborate on this trade-off in the supplementary material (Section 7).

4.3. Network transfer

Our contribution aims to address equivalent network performance regardless of the input data format. That is, for a given network, we strive to achieve *equal performance on both perspective and spherical data*. This objective is motivated by the limited number of spherical image datasets and the difficulty of collecting large scale spherical training data. If we can achieve high transferability of perspective image networks, we reduce the need for large amounts of spherical training data. Because generating tangent images inherently converts a spherical image into a collection of perspective ones, this representation facilitates the desired network transferability without requiring fine-tuning on the spherical data and with limited performance drop-off.

Experimental setup We evaluate the transferability of the tangent image representation in three experiments.

In the first experiment, we evaluate semantic segmentation performance on a *spherical* image test set using a network trained on the corresponding *perspective image* training set. We fine-tune the pre-trained, FCN-ResNet101 model [15, 22] provided by the PyTorch model zoo on the Stanford2D3DS dataset’s [1] perspective image training set. We then evaluate semantic segmentation performance on the spherical image test set at a level 8 resolution. As with the other semantic segmentation experiments, this experiment uses RGB-D inputs. During the dataset fine-tuning, we make sure to consider the desired angular resolution of the spherical test images. A network trained on perspective images with an angular resolution of 1° has learned filters accordingly. Should we apply those filters to an image captured at the identical position, at the same image resolution, but with a narrower FOV, the difference in angular resolution is effectively scale distortion. To match the angular resolution of our spherical evaluation set, we normalize the

camera matrices for all perspective images during training such they have the same angular resolution as the test images. Because this is effectively a center-crop of the data, we also randomly shift our new camera center in order capture all parts of the image. Details of this pre-processing are given in the supplementary material (Section 8). We evaluate performance without fine-tuning on spherical data, after 1 epoch of spherical fine-tuning, and again after 10 epochs. To control for the extra training, we also evaluate the perspective network after an additional 10 epochs of training.

The second experiment compares the transferability provided by tangent images to prior work that addresses this topic [31]. Using the network architecture from Zhang *et al.* [31], we train a model on the perspective images from the SYNTHIA dataset [26] that correspond to the OmniSYNTHIA dataset’s [31] training set. We again utilize the camera normalization procedure mentioned above. We evaluate performance on the OmniSYNTHIA test set at base level 1.

Finally, the third experiment studies the impact of matching angular resolution between training and testing. For this, we apply our FCN-ResNet 101 semantic segmentation model from the first experiment to the spherical test set at various resolutions.

Results and analysis Results for the first two experiments are given in Tables 3 and 4, respectively.

In the first experiment, note that both results are attained using a network trained only on perspective data. Without fine-tuning, we preserve about 97% of perspective network accuracy and 93% of the mean IOU. With only a single epoch of fine-tuning on spherical data, we see effective parity in performance, if not slightly improved performance on the spherical inputs. Finally, after 10 epochs of fine-tuning on the spherical format, the transferred network actually noticeably outperforms the original perspective network performance, even when compared to applying the same fine-tuning to the original network on perspective images. We surmise that this improvement comes from the greater FOV provided by the 360° image. These results suggest that tangent images sufficiently mitigate distortion, and, as a result, a network can begin to benefit from the ultra-wide field of view of 360° images. Although the tangent image representation breaks up the spherical FOV, the gradients are still computed from the full 360° image in our training routine, so the network still benefits from this extra information.

Additionally, we provide the results broken down by individual folds of the dataset because Fold 2 highlights the benefit of even 1 epoch of fine-tuning. Fold 2 is the hardest of the three at baseline, which might explain why tangent images provide less of a benefit than for the other two folds. However, after 10 epochs of fine-tuning on spherical data, we see Fold 2 performance increase significantly. We hypothesize that the ability to engage the wider FOV helps address particularly difficult scenes in the test set.

	Format	Res.	mAcc		mIOU	
Fold 1	P	$d = 128$	61.7	-	47.4	-
	S	L=8	60.2	-2.5%	45.3	-4.5%
	P-FT-1	$d = 128$	60.9	-	48.4	-
	S-FT-1	L=8	62.5	+2.6%	48.4	+0.0%
	P-FT-10	$d = 128$	60.6	-	48.9	-
Fold 2	S-FT-10	L=8	66.0	+8.9%	50.6	+3.5%
	P	$d = 128$	57.8	-	38.6	-
	S	L=8	47.8	-17.2%	34.4	-10.6%
	P-FT-1	$d = 128$	56.4	-	39.6	-
	S-FT-1	L=8	52.6	-6.7%	40.6	+2.5%
Fold 3	P-FT-10	$d = 128$	56.7	-	40.8	-
	S-FT-10	L=8	55.6	-1.9%	43.6	+7.1%
	P	$d = 128$	65.9	-	51.1	-
	S	L=8	64.2	-2.6%	47.3	-7.5%
	P-FT-1	$d = 128$	66.0	-	51.5	-
ALL FOLDS	S-FT-1	L=8	68.4	+3.6%	51.9	+0.6%
	P-FT-10	$d = 128$	65.6	-	51.7	-
	S-FT-10	L=8	70.3	+7.2%	54.6	+5.6%
	P	$d = 128$	65.9	-	51.1	-
	S	L=8	64.2	-2.6%	47.3	-7.5%
P-FT-1	P-FT-1	$d = 128$	66.0	-	51.5	-
	S-FT-1	L=8	68.4	+3.6%	51.9	+0.6%
	P-FT-10	$d = 128$	65.6	-	51.7	-
	S-FT-10	L=8	70.3	+7.2%	54.6	+5.6%

Table 3: Transfer learning using RGB-D data from the Stanford2D3DS dataset. “P” means the original network trained and evaluated on perspective images only, while “S” is that network evaluated on spherical data using tangent images, without any fine-tuning. “FT-#” denotes epochs of fine-tuning on a given format. The percentage next to the spherical results denotes how much of the original perspective network performance is maintained across input formats. Notice that, with fine-tuning, tangent images can actually lead to better performance than the corresponding central-perspective network.

The results of the second experiment demonstrate that the tangent image approach significantly outperforms the prior state-of-the-art without any specialized kernels or subsequent fine-tuning. Note that Zhang *et al.* [31] only report results after 10 epochs of fine-tuning on spherical images. Using tangent images actually provides noticeably better results without any such fine-tuning, although when we add 10 epochs of fine-tuning, we see an extra performance boost. It is also worth observing that our transfer results actually outperform the Zhang *et al.* [31] results trained natively on spherical data. Our experiments have been limited by the maximum resolution of available spherical image datasets, but this outcome suggests that network transfer with tangent images may permit even higher resolution spherical image inference.

Finally, the results of the third experiment are plotted in Figure 6. Recall that this model was trained on perspective images normalized to have a per-pixel angular resolution most similar to that of a level 8 icosahedron. This chart

s	Method	mAcc	mIOU
6	Ours (no FT)	55.2	43.2
	Ours (FT-10)	59.4	46.1
	Zhang <i>et al.</i> [31] (FT-10)	44.8	36.7
	Zhang <i>et al.</i> [31] (<i>native</i>)	52.2	43.6
7	Ours (no FT)	60.2	44.9
	Ours (FT-10)	65.3	50.2
	Zhang <i>et al.</i> [31] (FT-10)	47.2	38.0
	Zhang <i>et al.</i> [31] (<i>native</i>)	57.1	48.3
8	Ours (no FT)	70.8	54.9
	Ours (FT-10)	73.2	55.7
	Zhang <i>et al.</i> [31] (FT-10)	52.8	45.3
	Zhang <i>et al.</i> [31] (<i>native</i>)	55.1	47.1

Table 4: Comparing our transfer learning results to the prior work from Zhang *et al.* [31] on the OmniSYNTHIA dataset at different input resolutions, s. “no-FT” denotes no fine-tuning on spherical data, “FT-10” mean after 10 epochs of fine-tuning, and “native” means both trained and evaluated on spherical data. Even without fine-tuning tangent images significantly improve over the previous state-of-the-art. Fine-tuning provides a small, but noticeable, additional further improvement.

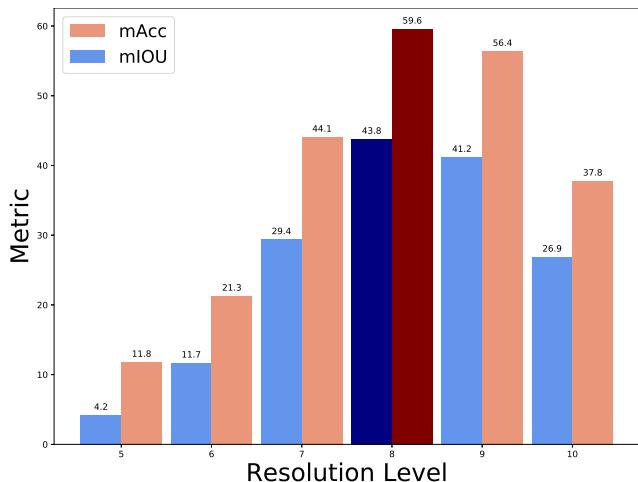


Figure 6: Results are shown for spherical semantic segmentation using a network trained on perspective images that are normalized to have an angular resolution equivalent to a level 8 spherical input. Performance drops off considerably as the angular resolution of the spherical inputs becomes more dissimilar to the training data. Level 8 results are darkened.

highlights the importance of camera normalization when training on perspective images with the purpose of transferring the network. Observe how performance deteriorates as the angular resolution of the spherical input moves further from the angular resolution of the training data.

4.4. Sparse keypoint correspondences

Recent research on spherical images has focused on deep learning tasks, primarily because many of those works have



Left Image **Right Image**
 Figure 7: FOV overlap visualized between an image pair from our keypoints benchmark derived from the Stanford 2D3DS dataset [1]. The red regions in the left image represent areas visible to the right camera, and the green regions in the right image represent areas visible to the left camera.

focused on the convolution operation. As our contribution relates to the representation of spherical data, not specifically convolution, we aim to show that our approach has applications beyond deep learning. To this end, we evaluate the use of tangent images for sparse keypoint detection, a critical step of structure-from-motion, SLAM, and a variety of other traditional computer vision applications.

Data As there is no existing benchmark for this task, we create a dataset using a subset of the spherical images provided by the Stanford2D3DS dataset [1]. To create this dataset, we first cluster the dataset’s Area 1 images according to the provided room information. Then, for each location, we compute SIFT features [24] in the equirectangular images and identify which image pairs have FOV overlap using the spherical structure-from-motion pipeline provided by the OpenMVG library [25]. Next, we compute the average volumetric FOV overlap for each overlapping image pair. Because we are dealing with 360° images, there are no image bounds to constrain “visible” regions. Instead, we use the ground truth depth maps and pose information to back-project each image pair into a canonical pose. We then compute the percentage of right image points visible to the left camera using the left image depth map to remove occluded points, and vice versa. We average the two values to provide an FOV overlap score for the image pair. This overlap is visualized in Figure 7. We define our keypoints dataset as the top 60 image pairs according to this overlap metric. Finally, we split the resulting dataset into an “Easy” set and “Hard” set, again based on FOV overlap. The resulting dataset statistics are shown in Table 5. All images are evaluated at their full, level 10 resolution. We provide the dataset details in the supplementary material (Section 10) to enable further research.

Experimental setup To evaluate our proposed representation, we detect and describe keypoints on the tangent image grids and then render those keypoints back to the spherical image. This rendering step ensures only keypoints visible to a spherical camera at the center of the icosahedron are rendered, as the tangent images have overlapping content. We then use OpenMVG [25] to compute putative correspondences and geometrically-consistent inlier matches.

Results and analysis We evaluate the quality of cor-

Split	# Pairs	Mean FOV Overlap	# Corr.
Hard	30	83.35%	298
Easy	30	89.35%	515

Table 5: Statistics of our keypoints benchmark. # Corr. is the number of inlier matches detected on the equirectangular images in that split. Statistics are averaged over the splits.

respondence matching at 3 different base levels using the equirectangular image format as a baseline. We compute the *putative matching ratio* (PMR), *matching score* (MS), and *precision* (P) metrics defined by Heinly *et al.* [16]. For an image set \mathcal{S} of image pairs, (L, R) , with p putative correspondences, f inlier matches, and $n_{\{L,R\}}$ detected keypoints visible to both images, the metrics over the image pairs as defined as follows:

$$\begin{aligned} \text{PMR} &= \frac{1}{2|\mathcal{S}|} \sum_{(L,R) \in \mathcal{S}} \left(\frac{p}{n_L} + \frac{p}{n_R} \right) \\ \text{MS} &= \frac{1}{2|\mathcal{S}|} \sum_{(L,R) \in \mathcal{S}} \left(\frac{f}{n_L} + \frac{f}{n_R} \right) \\ \text{P} &= \frac{1}{|\mathcal{S}|} \sum_{(L,R) \in \mathcal{S}} \frac{f}{p} \end{aligned} \quad (4)$$

In the same way that we compute the FOV overlap, we use the ground truth pose and depth information provided by the dataset to determine which keypoints in the left image should be visible to the right image (n_L) and vice versa (n_R), accounting for occlusion.

Results are given in Table 6. Our use of tangent images has a strong impact on the resulting correspondences, particularly on the hard split. Recall that this split has a lower FOV overlap and fewer inlier matches at the baseline equirectangular representation. Improved performance in this case is thus especially useful. We observe a significant improvement in PMR in both splits. We attribute this improvement to the computation of the SIFT feature vector on our less distorted representation. Like the convolution operation, SIFT descriptors also require translational equivariance in the detection domain. Tangent images restore this property with their low-distortion representation, which enables repeatable descriptors. The better localization of the keypoints affects the inlier matches as well, resulting in a better MS score. We attribute the leveling off in performance beyond level 1 to the reduced FOV of higher level subdivisions, which impedes the detector’s ability to find keypoints at larger scales.

5. Conclusion

We have presented tangent images, a spherical image representation that renders the image onto a oriented pixel grids tangent to a subdivided icosahedron. We have shown that these tangent images do not require specialized con-

Hard				
Metric	Equirect.	L0	L1	L2
PMR	22.2%	28.4%	30.1%	27.4%
MS	8.2%	11.1%	11.7%	10.9%
P	36.9%	39.5%	39.6%	40.2%
Easy				
Metric	Equirect.	L0	L1	L2
PMR	26.3%	32.4%	34.6%	31.9%
MS	13.6%	16.6%	17.7%	16.1%
P	46.0%	46.4%	47.5%	46.5%

Table 6: Keypoint evaluation metrics. We report the each metric’s average over all image pairs per split. L{0,1,2} are the subdivision levels at which we compute the keypoints.

volutional kernels for training CNNs and efficiently scale to represent high resolution data. We have also shown that they facilitate the transfer of networks trained on perspective images to spherical data with limited performance loss. These results further suggest that network transfer using tangent images can open the door to processing even higher resolution spherical images. Lastly, we have demonstrated the utility of tangent images for traditional computer vision tasks in addition to deep learning. Our results indicate that tangent images can be a very useful spherical representation for a wide variety of computer vision applications.

6. Differences from CVPR 2020 Publication

In order to present the improved performance of our experiments, ensure that our results match our publicly released code, and correct some minor errors in the CVPR 2020 version of this work, we have updated this version.

- An additional note about training with tangent images is added to Section 4.2.
- Table 2 in the CVPR version misreports results from Lee *et al.* [20]. This version corrects that mistake, but we note that it has no bearing on the analysis.
- We clarify that transfer learning experiments are performed on RGB-D images in Section 4.3. The CVPR version has mixed notes on this, saying in one place that they are performed on RGB images only and in another that they are performed on RGB-D inputs.
- We provide further exploration of network transfer by fine-tuning the network on spherical data. This additional work was done to be able to better compare to prior work [31] who report results only after fine-tuning. This extra fine-tuning is done for both network transfer experiments, and demonstrates the ability of tangent images to enable improved performance on spherical images.

Acknowledgements We would like to thank David Luebke, Pierre Moulon, Li Guan, and Jared Heinly for their consultation in support of this work. This research was funded in part by Zillow.

References

- [1] Iro Armeni, Sasha Sax, Amir R Zamir, and Silvio Savarese. Joint 2d-3d-semantic data for indoor scene understanding. *arXiv preprint arXiv:1702.01105*, 2017. 6, 8, 13, 14, 15, 17
- [2] Duane C Brown. Decentering distortion of lenses. *Photogrammetric Engineering and Remote Sensing*, 1966. 3
- [3] Fuller Richard Buckminster. Cartography, Jan. 29 1946. US Patent 2,393,676. 3
- [4] Brent Burley and Dylan Lacewell. Ptex: Per-face texture mapping for production rendering. In *Computer Graphics Forum*, volume 27, pages 1155–1164. Wiley Online Library, 2008. 5
- [5] Sharan Chetlur, Cliff Woolley, Philippe Vandermersch, Jonathan Cohen, John Tran, Bryan Catanzaro, and Evan Shelhamer. cudnn: Efficient primitives for deep learning. *arXiv preprint arXiv:1410.0759*, 2014. 2
- [6] Taco Cohen, Mario Geiger, Jonas Köhler, and Max Welling. Convolutional networks for spherical signals. *arXiv preprint arXiv:1709.04893*, 2017. 2
- [7] Taco Cohen, Maurice Weiler, Berkay Kicanaoglu, and Max Welling. Gauge equivariant convolutional networks and the icosahedral cnn. In *International Conference on Machine Learning*, pages 1321–1330, 2019. 1, 2, 3, 4, 5, 6
- [8] Taco S. Cohen, Mario Geiger, Jonas Khler, and Max Welling. Spherical CNNs. In *International Conference on Learning Representations*, 2018. 1, 2, 5, 6
- [9] Benjamin Coors, Alexandru Paul Condurache, and Andreas Geiger. Spherenet: Learning spherical representations for detection and classification in omnidirectional images. In *European Conference on Computer Vision*, pages 525–541. Springer, 2018. 1, 2
- [10] Marc Eder and Jan-Michael Frahm. Convolutions on spherical images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 1–5, 2019. 1, 2, 3, 4
- [11] Marc Eder, True Price, Thanh Vu, Akash Bapat, and Jan-Michael Frahm. Mapped convolutions. *arXiv preprint arXiv:1906.11096*, 2019. 1, 2, 3, 4
- [12] Carlos Esteves, Christine Allen-Blanchette, Ameesh Makadia, and Kostas Daniilidis. Learning so (3) equivariant representations with spherical cnns. In *European Conference on Computer Vision*, pages 52–68, 2018. 5
- [13] 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*, 24(6):381–395, 1981. 13
- [14] James Hastings-Trew Planet texture maps. Accessed: 2019-04-16. 1
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 6, 13, 14
- [16] Jared Heinly, Enrique Dunn, and Jan-Michael Frahm. Comparative evaluation of binary features. In *European Conference on Computer Vision*, pages 759–773, 2012. 9
- [17] Chiyou Max Jiang, Jingwei Huang, Karthik Kashinath, Prabhat, Philip Marcus, and Matthias Niessner. Spherical CNNs on unstructured grids. In *International Conference on Learning Representations*, 2019. 1, 2, 3, 4, 5, 6, 13, 14, 15
- [18] Jon A Kimerling, Kevin Sahr, Denis White, and Lian Song. Comparing geometrical properties of global grids. *Cartography and Geographic Information Science*, 26(4):271–288, 1999. 4
- [19] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015. 13, 14
- [20] Yeonkun Lee, Jaeseok Jeong, Jongseob Yun, Wonjune Cho, and Kuk-Jin Yoon. Spheredph: Applying cnns on a spherical polyhedron representation of 360deg images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9181–9189, 2019. 1, 2, 3, 4, 5, 6, 9
- [21] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755. Springer, 2014. 6, 13, 14
- [22] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015. 6
- [23] Charles Loop. Smooth subdivision surfaces based on triangles. *Master's thesis, University of Utah, Department of Mathematics*, 1987. 4
- [24] David G Lowe et al. Object recognition from local scale-invariant features. In *International Conference on Computer Vision*, volume 99, pages 1150–1157, 1999. 8
- [25] Pierre Moulon, Pascal Monasse, Romuald Perrot, and Renaud Marlet. Openmvg: Open multiple view geometry. In *International Workshop on Reproducible Research in Pattern Recognition*, pages 60–74. Springer, 2016. 8
- [26] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio M. Lopez. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. 7
- [27] Yu-Chuan Su and Kristen Grauman. Learning spherical convolution for fast features from 360 imagery. In *Advances in Neural Information Processing Systems*, pages 529–539, 2017. 1, 2
- [28] Yu-Chuan Su and Kristen Grauman. Kernel transformer networks for compact spherical convolution. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 1, 2
- [29] Keisuke Tateno, Nassir Navab, and Federico Tombari. Distortion-aware convolutional filters for dense prediction in panoramic images. In *European Conference on Computer Vision*, pages 732–750. Springer, 2018. 1, 2, 5, 6
- [30] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Ligang Zhang, Xiaou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1912–1920, 2015. 5, 13
- [31] Chao Zhang, Stephan Liwicki, William Smith, and Roberto Cipolla. Orientation-aware semantic segmentation on icosahedron spheres. *Proceedings of the IEEE International Conference on Computer Vision*, 2019. 1, 2, 3, 4, 5, 6, 7, 8, 9, 13, 15

Supplementary Material: Tangent Images for Mitigating Spherical Distortion

In this supplementary material we provide the following additional information:

- Expanded discussion of some of the current limitations of tangent images (Section 7)
- The details of the camera normalization process and class-level results of our transfer learning experiments (Section 8)
- Class-level and qualitative results for the semantic segmentation experiments at different input resolutions (Section 9)
- Details of our 2D3DS Keypoints dataset along with individual image pair results and a qualitative comparison of select image pairs (Section 10)
- Training and architecture details for all CNN experiments detailed in this paper (Section 11)
- An example of a spherical image represented as tangent images (Figure 8)

7. Limitations

We have demonstrated the usefulness of our proposed tangent images, but we have also exposed some limitations of the formulation and opportunities for further research.

b	0	1	2	3	4
FOV	73.1°	51.6°	31.5°	16.7°	8.4°

Table 7: Tangent image field of view at different base levels (b) for a level 10 input. There is a slight variation depending on the input resolution due to faces near the 12 icosahedral singular points, but the values stay mostly consistent.

Resolution When using tangent images, low angular resolution spherical data is processed on potentially low pixel resolution images. This can severely limit the receptive field of the networks, to which we attribute our poor performance on the level 5 semantic segmentation task, for example. However, this limitation is only notable in the context of the existing literature, because prior work has been restricted to low resolution spherical data, as shown in Figure 2 in the main paper. One viable solution is to incorporate the rendering into the convolution operation. In this way, we could regenerate the tangent images at every operation and effectively resolve the receptive field issue. However, as this is an issue for low resolution images, and our work is focused on addressing high resolution spherical data, we leave this modification for future study.

FOV The base subdivision level provides a constraint on the FOV of the tangent images. Table 7 shows the FOV of the tangent images at different base subdivision levels. As the FOV decreases, algorithms that rely on some sense of

context or neighborhood break down. We observe this effect for both the CNN and keypoint experiments. While this is certainly a trade-off with tangent images, we have demonstrated that base levels and resolutions can be modulated to fit the required application. Another important point to observe regarding tangent image FOV is that the relationship between FOV and subdivision level does not hold perfectly at lower subdivision levels due the outsize influence of faces near the 12 singular points on the icosahedron. This effect largely disappears after base level 2, but when normalizing camera matrices to match spherical angular resolution at base levels 0 and 1, it is necessary to choose the right base level for the data. We use a base level of 1 in our transfer learning experiments on the OmniSYNTHIA dataset for this reason.

8. Network Transfer

In this section, we detail the camera normalization process used when training the network for transferring to spherical data. We also provide class-level results for our experiments.

8.1. Camera normalization

In order to ensure angular resolution conformity between perspective training data and spherical test data, we normalize the intrinsic camera matrices of our training images to match the angular resolution of the spherical inputs, α_s . To do this, we resample all perspective image inputs to a common camera with the desired angular resolution. The angular resolution in radians-per-pixel of an image, α_x and α_y , is defined as:

$$\alpha_x = \frac{\Omega_x}{W} \quad \alpha_y = \frac{\Omega_y}{H} \quad (5)$$

where Ω_x and Ω_y are the fields of view of the image as a function of the image dimensions, W and H , and the focal lengths, f_x and f_y :

$$\begin{aligned} \Omega_x &= 2 \arctan \left(\frac{W}{2f_x} \right) \\ \Omega_y &= 2 \arctan \left(\frac{H}{2f_y} \right) \end{aligned} \quad (6)$$

Because spherical inputs have uniform angular resolution in every direction, we resample our perspective inputs likewise: $\alpha_x = \alpha_y = \alpha_s$.

Choosing camera properties For our camera-normalized perspective images, we want to choose fields of view, Ω'_x and Ω'_y , and image dimensions, W' and

H' that satisfy:

$$\frac{\Omega'_x}{W'} = \frac{\Omega'_y}{H'} = \alpha_s \quad (7)$$

While there are a variety of options that we could use, we choose to set $\Omega'_x = \Omega'_y = \frac{\pi}{4}$ because $\frac{\pi}{4}$ radians (45°) is a reasonable field of view for a perspective image. We select W' and H' accordingly. For a level 8 input, this results in $W' = H' = 128$.

Normalizing intrinsics Recall the definition of the intrinsic matrix, K , given focal lengths f_x and f_y and principal point (c_x, c_y) :

$$K = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \quad (8)$$

Given our choices for fields of view and image dimensions explained above, we compute a new, common intrinsic matrix. The new focal lengths, f'_x and f'_y , are computed as:

$$f'_x = \frac{W'}{2 \tan\left(\frac{\Omega'_x}{2}\right)} \quad f'_y = \frac{H'}{2 \tan\left(\frac{\Omega'_y}{2}\right)} \quad (9)$$

and, for simplicity, the new principal point is chosen to be:

$$c'_x = \frac{W'}{2} \quad c'_y = \frac{H'}{2} \quad (10)$$

Defining:

$$K' = \begin{bmatrix} f'_x & 0 & c'_x \\ 0 & f'_y & c'_y \\ 0 & 0 & 1 \end{bmatrix} \quad (11)$$

the camera intrinsics can be normalized using the relation:

$$x' = K' K^{-1} x \quad (12)$$

where x and x' are homogeneous pixel coordinates in the original and resampled images, respectively, and K^{-1} is the inverse of the intrinsic matrix associated with the original image.

Random shifts If we were to simply resample the image using Equation (12), we would end up with center crops of the original perspective images. In order to ensure that we do not discard useful information, we randomly shift the principle point of the original camera by some (δ_x, δ_y) before normalizing. This produces variations in where we crop the original image. Including this shift, we arrive at the formula we use for resampling the perspective training data:

$$x' = K' (K + \Delta)^{-1} x \quad (13)$$

where:

$$\Delta = \begin{bmatrix} 0 & 0 & \delta_x \\ 0 & 0 & \delta_y \\ 0 & 0 & 0 \end{bmatrix} \quad (14)$$

To ensure our crops stay within the bounds of the original image, we want:

$$\begin{aligned} \delta_x + P(0, 0)_x &\geq 0 \\ \delta_y + P(0, 0)_y &\geq 0 \\ \delta_x + P(W', H')_x &\leq W \\ \delta_y + P(W', H')_y &\leq H \end{aligned} \quad (15)$$

where $P(x', y')_{\{x, y\}}$ denotes the x - and y -dimensions of the new camera's coordinates projected into the original camera's coordinate system:

$$P(x', y') = K K'^{-1} \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} \quad (16)$$

Using this constraint, we sample crops over the entire image by randomly choosing δ_x and δ_y from the ranges:

$$\begin{aligned} \frac{f_x}{f'_x} c'_x - c_x &\leq \delta_x \leq W - c_x - \frac{f_x}{f'_x} (W' - c'_x) \\ \frac{f_y}{f'_y} c'_y - c_y &\leq \delta_y \leq H - c_y - \frac{f_y}{f'_y} (H' - c'_y) \end{aligned} \quad (17)$$

8.2. Per-class results

Table 8 gives the per-class results for our semantic segmentation transfer experiment. While perspective image performance should be considered an upper bound on spherical performance, note that in some classes, we appear to actually perform better on the spherical image. This is because the spherical evaluation is done on equirectangular images in order to be commensurate across representations. This means that certain labels are duplicated and others are reduced due to distortion, which can skew the per-class results.

9. Semantic Segmentation

We provide the per-class quantitative results for our semantic segmentation experiments from Section 4.2 in the main paper. Additionally, we qualitatively analyze the benefits of training higher resolution networks, made possible by the tangent image representation.

9.1. Class results

Per-class results are given for semantic segmentation in Table 9. Nearly every class benefits from the high-resolution inputs facilitated by tangent images. This is especially noticeable for classes with fine-grained detail and smaller objects, like *chair*, *window*, and *bookcase*. The *table* class is an interesting example of the benefit of our method. While prior work has higher accuracy, our high resolution classification has significantly better IOU. In other words, our high resolution inputs may not result in

correct classifications of every *table* pixel, but the classifications that are correct are much more precise. This increased precision is reflected almost across the board by mean IOU performance.

9.2. Qualitative results

Figure 9 gives 3 examples of semantic segmentation results at each resolution. The most obvious benefits of the higher resolution results are visible in the granularity of the output segmentation. Notice the fine detail preserved in the chairs in the level 10 output in the bottom row and even the doorway and whiteboard in the middle row. However, recall that our level 10 network uses a base level of 1. The effects of the smaller FOV of the tangent images are visible in the misclassifications of wall on the right of the level 10 output in the middle row. The level 5 network has no such problems classifying that surface, having been trained at a lower input resolution and using base level 0. Nevertheless, it is worth noting that large, homogeneous regions are going to be problematic for high resolution images, regardless of method, due to receptive field limitations of the network. If the region in question is larger than the receptive field of the network, there is no way for the network to infer context. As such, we are less concerned by errors on these large regions.

10. Stanford 2D3DS Keypoints Dataset

10.1. Details

Tables 10 and 11 give the details of the image pairs in our keypoints dataset. Tables 12 and 13 provide the individual metrics computed for each image pair.

10.2. Qualitative examples

We provide a qualitative comparison of keypoint detections in Figure 10. These images illustrate two interesting effects, in particular. First, in highly distorted regions of the image that have repeatable texture, like the floor in both images, detecting on the equirectangular format produces a number of redundant keypoints distributed over the distorted region. With tangent images, we see fewer redundant points as the base level increases and the ones that are detected are more accurate and robust, as indicated by the higher MS score. Additionally, the equirectangular representation results in more keypoint detections at larger scales. These outsize scales are an effect of distortion. Rotating the camera so that the corresponding keypoints are detected at different pixel locations with different distortion characteristics will produce a different scale, and consequently a difference descriptor. This demonstrates the need for translational equivariance in keypoint detection, which requires the lower distortion provided by our tangent images. This is reflected quantitatively by the higher PMR

scores.

Figure 11 shows an example of inlier correspondences computed on the equirectangular images and at different base levels for an image pair from the hard split. Even though we detect fewer keypoints using tangent planes, we still have the same quality or better inlier correspondences. Distortion in the equirectangular format results in keypoint over-detection, which can potentially strain the subsequent inlier model fitting. Using tangent images, we detect fewer, but higher quality, samples. This results in more efficient and reliable RANSAC [13] model fitting. This is why tangent images perform noticeably better on the hard set, where there are fewer correspondences to be found.

11. Network Training Details

We detail the training parameters and network architectures used in our experiments to encourage reproducible research. All deep learning experiments were performed using the PyTorch library.

11.1. Shape classification

For the shape classification experiment, we use the network architecture from [17], replacing their MeshConv layers with 3×3 2D convolutions with unit padding. For downsampling operations, we bilinearly interpolate the tangent images. We first render the ModelNet40 [30] shapes to equirectangular images to be compatible with our tangent image generation implementation. The equirectangular image dimensions are 64×128 , which is equivalent to the level 5 icosahedron. We train the network with a batch size of 16 and learning rate of 5×10^{-3} using the Adam optimizer [19].

11.2. Semantic segmentation

We use the residual U-Net-style architecture from [17, 31] for our semantic segmentation results at levels 5 and 7. Instead of MeshConv [17] and HexConv [31], we swap in 3×3 2D convolutions with unit padding. For the level 10 results, we use the fully-convolutional ResNet101 [15] pre-trained on COCO [21] provided in the PyTorch model zoo. We train and test on each of the standard folds of the Stanford 2D3DS dataset [1]. For all spherical data, evaluation metrics are computed on the re-rendered spherical image, not on the tangent images. As mentioned in the main paper, when training with tangent images, forward passes are run on all tangent images in a batch before a backward pass is run. This ensures that the whole 360° image is incorporated into the gradient computation.

Level 5, 7 parameters For the level 5 and 7 experiments, our tangent images were base level 0 RGB-D images, and we use all 20 tangent images from each image in a batch of 8 spherical images, resulting in an effective batch size of 160 tangent images. We use Adam optimization [19] with

an initial learning rate of 10^{-2} and decay by 0.9 every 20 epochs, as in [17].

Level 10 parameters For the level 10 experiments, we use RGB-D images at base level 1 and randomly sample 4 tangent images from each image in a batch of 4 spherical inputs, resulting in an effective batch size of 16 tangent images. Because the pre-trained network does not have a depth channel, we initialize the depth channel filter with zero weights. This has the effect of slowly adding the depth information to the model. Similarly, the last layer is randomly initialized, as Stanford 2D3DS has a different number of classes than COCO. We use Adam optimization [19] with a learning rate of 10^{-4} .

11.3. Transfer learning

We again use the fully-convolutional ResNet101 [15] architecture pre-trained on COCO [21]. We fine tune for 10 epochs on the perspective images of the Standford2D3DS dataset [1]. We use a batch size of 16 and a learning rate of 10^{-4} . When fine-tuning, we again use a batch size of 16, but we reduce the learning rate to 10^{-5} .

mAcc						
Class	P	S	Perf. %	P-FT-10	S-FT-10	Perf. %
<i>beam</i>	28.0	22.9	-18.1%	15.3	28.1	+83.0%
<i>board</i>	71.6	57.1	-20.2%	68.9	65.3	-5.2%
<i>bookcase</i>	66.5	53.5	-19.6%	65.8	63.0	-4.3%
<i>ceiling</i>	88.0	88.8	+0.9%	89.3	87.8	-1.7%
<i>chair</i>	66.7	66.6	-0.3%	67.7	73.8	+9.1%
<i>clutter</i>	57.4	45.9	-20.2%	61.1	56.2	-8.1%
<i>column</i>	21.1	24.5	+16.4%	14.7	25.3	+72.7%
<i>door</i>	60.2	67.1	+11.4%	58.3	75.9	+30.2%
<i>floor</i>	90.0	94.8	+5.4%	91.3	94.1	+3.1%
<i>sofa</i>	36.3	36.2	-0.1%	41.1	50.1	+21.9%
<i>table</i>	75.7	54.9	-27.4%	73.7	63.3	-14.1%
<i>wall</i>	77.9	72.4	-7.1%	82.7	81.9	-0.9%
<i>window</i>	64.1	61.7	-3.7%	62.7	67.0	+6.9%

mIOU						
Class	P	S	Perf. %	P-FT-10	S-FT-10	Perf. %
<i>beam</i>	8.6	8.4	-2.3%	6.7	10.8	+62.2%
<i>board</i>	50.5	48.1	-4.7%	56.2	54.3	-3.4%
<i>bookcase</i>	45.5	42.5	-6.6%	47.0	49.2	+4.6%
<i>ceiling</i>	72.6	73.1	+0.6%	73.6	83.9	+13.9%
<i>chair</i>	50.1	50.3	+0.4%	51.2	57.9	+13.2%
<i>clutter</i>	36.8	34.3	-6.9%	38.4	39.4	+2.7%
<i>column</i>	11.5	9.9	-14.2%	10.0	10.2	+2.8%
<i>door</i>	49.0	43.2	-11.7%	48.9	51.3	+4.9%
<i>floor</i>	82.0	71.3	-13.0%	84.0	90.5	+7.7%
<i>sofa</i>	22.2	23.3	+4.8%	24.9	30.5	+22.5%
<i>table</i>	50.5	46.2	-8.5%	53.7	53.8	+0.3%
<i>wall</i>	64.6	60.0	-7.0%	66.5	68.1	+2.5%
<i>window</i>	50.1	39.7	-20.8%	51.6	45.0	-12.7%

Table 8: Per-class results for the semantic segmentation transfer learning experiment on the Stanford 2D3DS dataset [1]. “Perf. %” denotes how much better (+) or worse (-) the transferred network performance is as a percentage of the corresponding perspective image network. “P” and “S” refer to evaluation on perspective images and spherical images, respectively. “FT-10” refers to the network fine-tuned on the associated format for 10 epochs. These metrics are averaged over all folds.

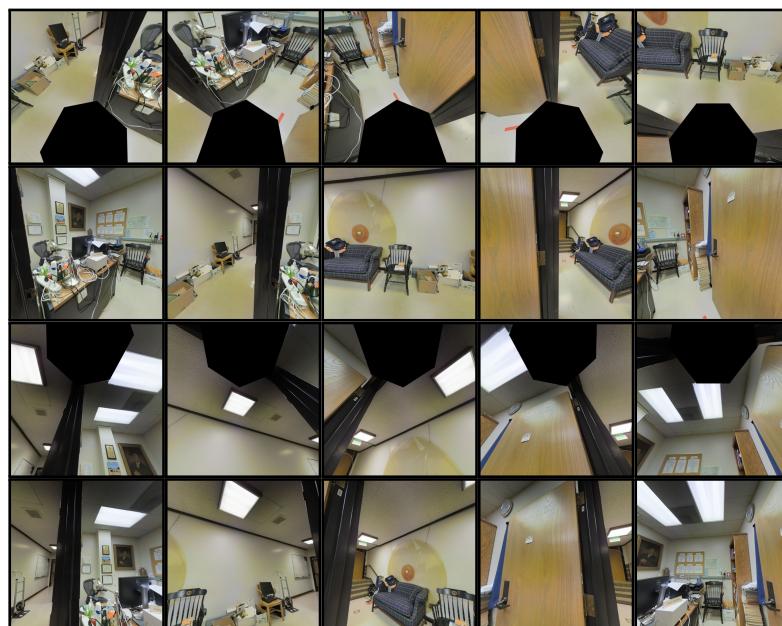
mAcc													
Method	<i>beam</i>	<i>board</i>	<i>bookcase</i>	<i>ceiling</i>	<i>chair</i>	<i>clutter</i>	<i>column</i>	<i>door</i>	<i>floor</i>	<i>sofa</i>	<i>table</i>	<i>wall</i>	<i>window</i>
Jiang <i>et al.</i> [17]	19.6	48.6	49.6	93.6	63.8	43.1	28.0	63.2	96.4	21.0	70.0	74.6	39.0
Zhang <i>et al.</i> [31]	23.2	56.5	62.1	94.6	66.7	41.5	18.3	64.5	96.2	41.1	79.7	77.2	41.1
Ours L5	28.5	39.3	48.0	85.8	51.3	54.2	11.9	63.5	92.7	16.6	51.1	77.5	41.6
Ours L7	30.4	57.2	56.8	88.6	72.7	56.6	18.4	68.9	95.2	27.8	57.9	78.6	58.6
Ours L10	26.6	76.6	63.0	90.2	85.0	62.5	37.4	72.0	97.3	67.5	73.6	78.6	80.5

mIOU													
Method	<i>beam</i>	<i>board</i>	<i>bookcase</i>	<i>ceiling</i>	<i>chair</i>	<i>clutter</i>	<i>column</i>	<i>door</i>	<i>floor</i>	<i>sofa</i>	<i>table</i>	<i>wall</i>	<i>window</i>
Jiang <i>et al.</i> [17]	8.7	32.7	33.4	82.2	42.0	25.6	10.1	41.6	87.0	7.6	41.7	61.7	23.5
Zhang <i>et al.</i> [31]	10.9	39.7	37.2	84.8	50.5	29.2	11.5	45.3	92.9	19.1	49.1	63.8	29.4
Ours L5	9.0	27.5	34.7	81.4	38.4	30.2	5.2	42.6	89.2	10.2	42.0	58.3	29.5
Ours L7	10.3	41.3	40.5	84.8	51.0	40.0	7.3	47.2	92.6	16.0	49.7	66.0	37.2
Ours L10	6.0	49.3	49.7	85.4	71.7	44.4	16.0	52.2	94.5	33.1	62.1	70.0	48.5

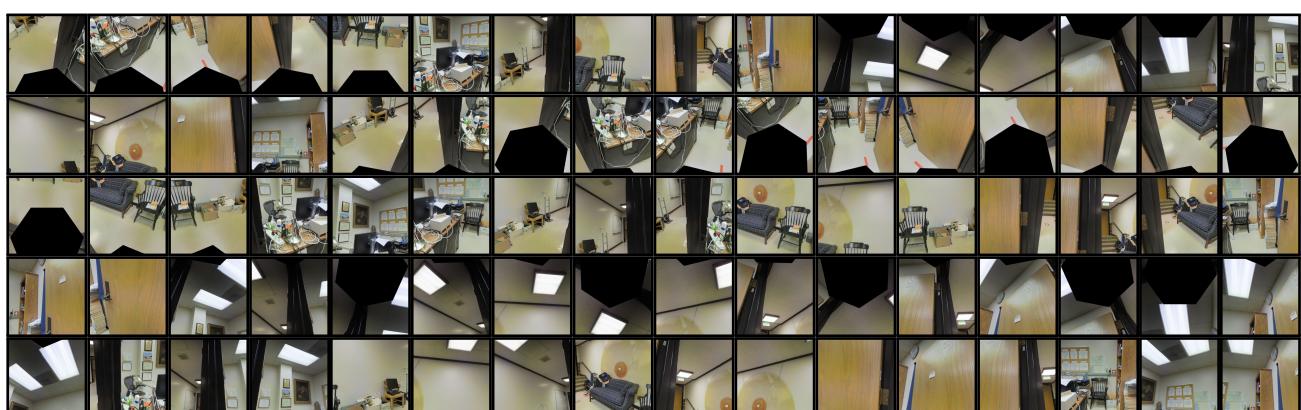
Table 9: Per-class results for RGB-D inputs on the Stanford 2D3DS dataset [1]



(a) Spherical image (equirectangular format)



(b) Base level 0



(c) Base level 1

Figure 8: Example of tangent images at base levels 0 and 1.

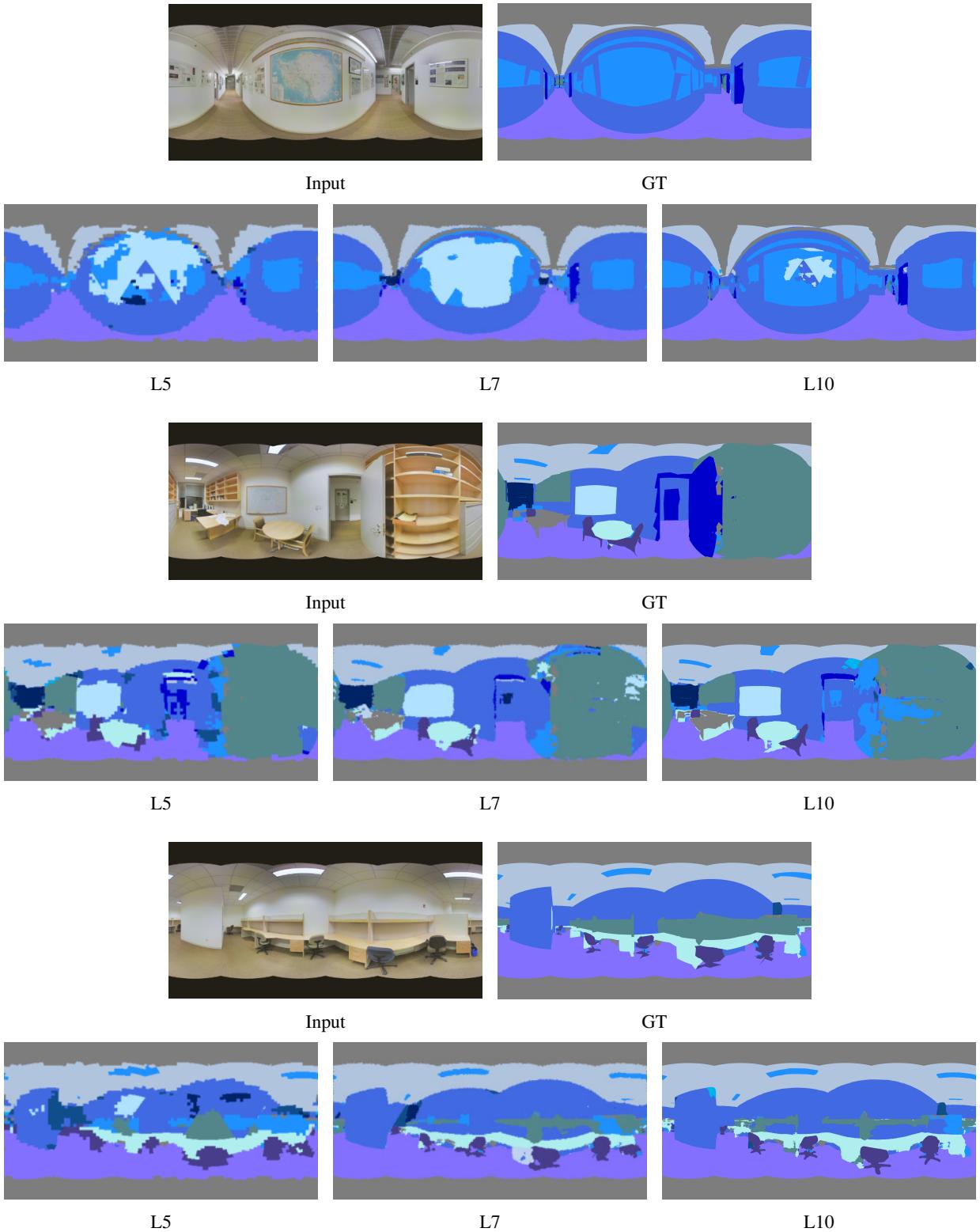


Figure 9: Qualitative results of semantic segmentation on the Stanford 2D3DS dataset [1] at different input resolutions. These results illustrate the performance gains we access by being able to scale to high resolution spherical inputs.

Easy			
Pair ID	Left Image	Right Image	FOV Overlap
1	0f65e09_hallway_7	1ebdfef_hallway_7	0.87
2	08a99a5_hallway_7	251a331_hallway_7	0.97
3	08a99a5_hallway_7	f7c6c2a_hallway_7	0.89
4	251a331_hallway_7	f7c6c2a_hallway_7	0.87
5	251a331_hallway_7	b261c3b_hallway_7	0.97
10	f7c6c2a_hallway_7	b261c3b_hallway_7	0.89
20	9178f6a_hallway_6	29abbc1_hallway_6	0.89
23	bc12865_hallway_6	7d58331_hallway_6	0.88
24	ee20957_hallway_6	bed890d_hallway_6	0.86
25	ee20957_hallway_6	eaba8c8_hallway_6	1.00
28	bed890d_hallway_6	eaba8c8_hallway_6	0.86
29	077f181_hallway_6	83baa70_hallway_6	0.86
30	97ab30c_hallway_6	eaba8c8_hallway_6	0.86
31	fc19236_office_18	e7d9e58_office_18	0.88
34	09ad38a_office_26	04a59ce_office_26	0.96
35	04a59ce_office_26	c16a90f_office_26	0.96
37	c40ca55_office_31	7b74e08_office_31	0.87
39	4a41b27_office_31	7b74e08_office_31	0.87
43	5512025_office_23	7f04c9b_office_23	0.92
44	5512025_office_23	5a18aa0_office_23	0.86
45	7f04c9b_office_23	5a18aa0_office_23	0.87
47	433548f_hallway_3	dcab252_hallway_3	0.91
49	d31d981_office_8	54e6de3_office_8	0.89
50	f85a909_office_3	c9feabc_office_3	0.88
51	f85a909_office_3	97be01e_office_3	0.89
52	c9feabc_office_3	97be01e_office_3	0.90
54	8fd8146_office_10	ab03f88_office_10	0.88
55	7c870c2_hallway_8	4de69cf_hallway_8	0.87
56	33e598f_office_15	8910cb1_office_15	0.88
58	46b4538_office_1	db2e53f_office_1	0.92

Table 10: Easy split of Stanford2D3DS keypoints dataset image pairs.

Hard			
<i>Pair ID</i>	<i>Left Image</i>	<i>Right Image</i>	<i>FOV Overlap</i>
0	c14611b_hallway_7	37a4f42_hallway_7	0.83
6	1b253d2_hallway_7	6e945c8_hallway_7	0.84
7	5d3a59a_hallway_7	ec0b9b3_hallway_7	0.81
8	ac01e35_hallway_7	649838b_hallway_7	0.85
9	f6c6ce3_hallway_7	5221e31_hallway_7	0.85
11	438c5fb_hallway_7	ec0b9b3_hallway_7	0.82
12	ec0b9b3_hallway_7	531efee_hallway_7	0.85
13	724bbea_hallway_7	c8c806b_hallway_7	0.85
14	724bbea_hallway_7	55db392_hallway_7	0.82
15	32d9e73_hallway_7	55db392_hallway_7	0.85
16	fcd2380_office_22	2d842ce_office_22	0.85
17	2d842ce_office_22	ffd2cca_office_22	0.86
18	89d9828_hallway_6	87e6e35_hallway_6	0.81
19	89d9828_hallway_6	7d58331_hallway_6	0.84
21	75acaa8_hallway_6	87e6e35_hallway_6	0.84
22	92b146f_hallway_6	8c78856_hallway_6	0.86
26	b640b47_hallway_6	87e6e35_hallway_6	0.80
27	bed890d_hallway_6	97ab30c_hallway_6	0.85
32	af50002_WC_1	36dd48f_WC_1	0.84
33	1edba7e_WC_1	e0c041d_WC_1	0.84
36	c40ca55_office_31	a77fba5_office_31	0.85
38	4a41b27_office_31	da4629d_office_31	0.82
40	da4629d_office_31	9084f21_office_31	0.84
41	75361af_office_31	ecf7fb4_office_31	0.82
42	2100dd9_office_4	26c24c7_office_4	0.83
46	84cdc9a_conferenceRoom_1	0d600f9_conferenceRoom_1	0.83
48	dcab252_hallway_3	a9cda4d_hallway_3	0.82
53	6549526_office_21	08aa476_office_21	0.83
57	dbcdb33_office_20	f02c98c_office_20	0.83
59	24f42d6_hallway_5	684b940_hallway_5	0.84

Table 11: Hard split of Stanford2D3DS keypoints dataset image pairs.

Pair ID	Easy											
	Equirect.			L0			L1			L2		
	PMR	MS	P	PMR	MS	P	PMR	MS	P	PMR	MS	P
1	0.27	0.07	0.27	0.35	0.09	0.26	0.40	0.13	0.33	0.34	0.10	0.28
2	0.41	0.17	0.42	0.45	0.27	0.60	0.50	0.28	0.57	0.48	0.28	0.58
3	0.34	0.19	0.55	0.40	0.23	0.56	0.46	0.27	0.60	0.44	0.25	0.58
4	0.21	0.08	0.36	0.27	0.11	0.41	0.30	0.14	0.48	0.25	0.12	0.48
5	0.59	0.48	0.80	0.67	0.47	0.70	0.70	0.45	0.64	0.66	0.35	0.54
10	0.29	0.10	0.35	0.35	0.12	0.35	0.36	0.14	0.38	0.34	0.11	0.33
20	0.30	0.09	0.30	0.46	0.14	0.32	0.50	0.19	0.38	0.43	0.14	0.33
23	0.23	0.07	0.32	0.33	0.13	0.39	0.34	0.14	0.42	0.34	0.13	0.38
24	0.16	0.07	0.43	0.22	0.09	0.42	0.24	0.10	0.43	0.23	0.10	0.44
25	0.83	0.48	0.58	0.99	0.71	0.72	1.01	0.58	0.58	0.95	0.58	0.61
28	0.18	0.06	0.34	0.25	0.08	0.33	0.24	0.10	0.40	0.24	0.10	0.40
29	0.29	0.13	0.46	0.35	0.16	0.46	0.38	0.21	0.55	0.34	0.17	0.51
30	0.17	0.07	0.42	0.25	0.08	0.33	0.24	0.08	0.32	0.21	0.08	0.37
31	0.14	0.06	0.47	0.15	0.07	0.48	0.16	0.08	0.51	0.14	0.07	0.51
34	0.46	0.39	0.86	0.48	0.38	0.78	0.51	0.42	0.82	0.52	0.43	0.84
35	0.42	0.34	0.82	0.43	0.33	0.77	0.46	0.38	0.82	0.47	0.38	0.82
37	0.23	0.09	0.40	0.24	0.10	0.42	0.25	0.12	0.47	0.25	0.11	0.43
39	0.22	0.08	0.39	0.23	0.08	0.34	0.24	0.09	0.36	0.23	0.07	0.32
43	0.27	0.19	0.70	0.38	0.24	0.63	0.39	0.25	0.63	0.37	0.24	0.65
44	0.13	0.04	0.30	0.18	0.07	0.39	0.23	0.08	0.36	0.16	0.06	0.36
45	0.14	0.05	0.36	0.17	0.07	0.40	0.21	0.09	0.44	0.17	0.06	0.39
47	0.31	0.21	0.67	0.40	0.25	0.64	0.42	0.22	0.52	0.36	0.21	0.57
49	0.10	0.04	0.41	0.15	0.05	0.36	0.15	0.06	0.37	0.15	0.06	0.42
50	0.18	0.08	0.46	0.21	0.11	0.50	0.23	0.11	0.47	0.22	0.11	0.52
51	0.15	0.04	0.31	0.19	0.06	0.32	0.22	0.07	0.31	0.19	0.07	0.38
52	0.15	0.05	0.32	0.18	0.06	0.35	0.19	0.08	0.39	0.18	0.06	0.34
54	0.17	0.05	0.31	0.24	0.09	0.35	0.25	0.08	0.33	0.22	0.08	0.38
55	0.18	0.10	0.53	0.24	0.11	0.45	0.26	0.13	0.49	0.22	0.07	0.32
56	0.22	0.11	0.50	0.32	0.16	0.50	0.33	0.16	0.48	0.29	0.14	0.49
58	0.16	0.06	0.37	0.19	0.07	0.37	0.22	0.09	0.40	0.20	0.08	0.38

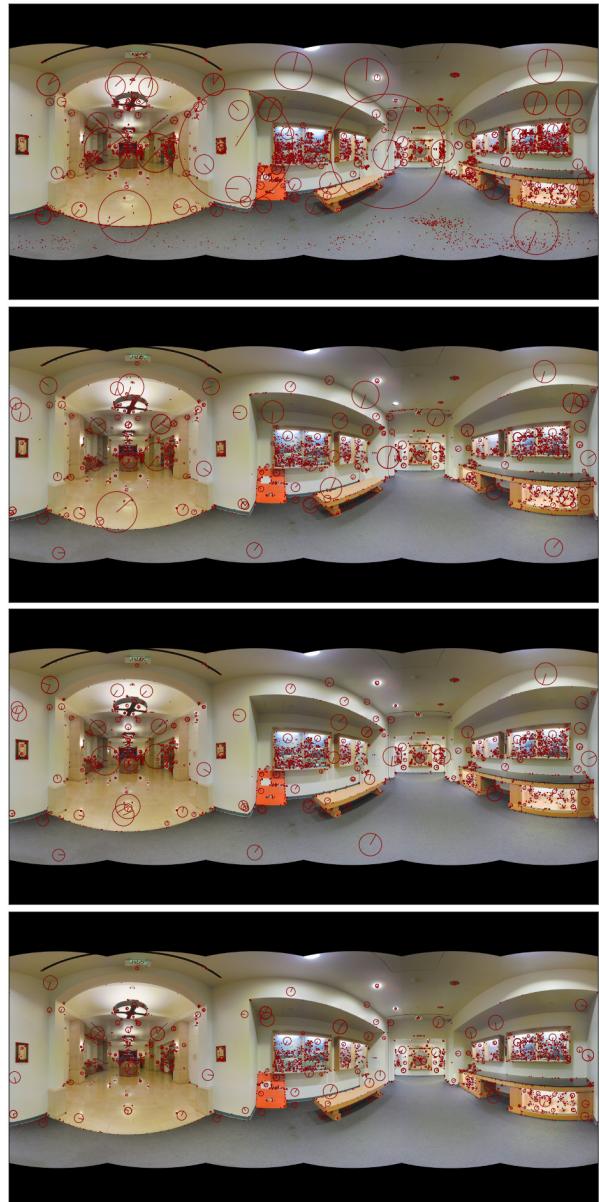
Table 12: Keypoint matching results on individual image pairs in the easy split.

Pair ID	Hard											
	Equirect.			L0			L1			L2		
	PMR	MS	P	PMR	MS	P	PMR	MS	P	PMR	MS	P
0	0.25	0.12	0.51	0.28	0.14	0.50	0.30	0.15	0.49	0.30	0.15	0.49
6	0.26	0.09	0.33	0.35	0.15	0.43	0.38	0.16	0.43	0.34	0.14	0.40
7	0.22	0.07	0.32	0.28	0.09	0.34	0.29	0.09	0.32	0.29	0.10	0.34
8	0.23	0.08	0.38	0.30	0.14	0.46	0.30	0.11	0.37	0.29	0.11	0.39
9	0.31	0.08	0.26	0.42	0.12	0.28	0.42	0.13	0.32	0.39	0.11	0.27
11	0.23	0.07	0.32	0.28	0.11	0.41	0.32	0.08	0.25	0.27	0.08	0.29
12	0.20	0.05	0.24	0.34	0.09	0.26	0.33	0.07	0.21	0.29	0.09	0.32
13	0.24	0.07	0.30	0.35	0.10	0.29	0.37	0.11	0.30	0.31	0.10	0.31
14	0.26	0.08	0.32	0.43	0.10	0.24	0.50	0.17	0.33	0.40	0.12	0.30
15	0.30	0.12	0.40	0.40	0.19	0.47	0.42	0.21	0.51	0.36	0.17	0.49
16	0.16	0.05	0.34	0.19	0.06	0.35	0.19	0.07	0.37	0.17	0.06	0.36
17	0.19	0.09	0.47	0.21	0.11	0.51	0.24	0.12	0.49	0.21	0.11	0.51
18	0.24	0.06	0.26	0.36	0.12	0.33	0.36	0.10	0.28	0.31	0.12	0.38
19	0.20	0.06	0.28	0.31	0.09	0.29	0.34	0.12	0.34	0.28	0.12	0.42
21	0.22	0.08	0.35	0.30	0.11	0.37	0.31	0.12	0.38	0.29	0.10	0.34
22	0.25	0.07	0.29	0.35	0.12	0.35	0.36	0.13	0.37	0.33	0.13	0.41
26	0.21	0.06	0.31	0.31	0.10	0.31	0.33	0.10	0.30	0.29	0.08	0.29
27	0.16	0.06	0.37	0.24	0.11	0.46	0.25	0.12	0.48	0.22	0.10	0.46
32	0.25	0.09	0.37	0.30	0.12	0.39	0.34	0.15	0.43	0.30	0.12	0.39
33	0.19	0.09	0.49	0.24	0.12	0.50	0.25	0.13	0.51	0.26	0.14	0.53
36	0.23	0.10	0.42	0.25	0.11	0.44	0.26	0.11	0.44	0.25	0.10	0.42
38	0.22	0.09	0.39	0.23	0.08	0.37	0.23	0.09	0.37	0.23	0.09	0.41
40	0.20	0.10	0.50	0.22	0.11	0.51	0.23	0.11	0.50	0.22	0.11	0.52
41	0.23	0.12	0.52	0.25	0.14	0.54	0.26	0.14	0.54	0.27	0.15	0.56
42	0.17	0.05	0.30	0.21	0.08	0.37	0.21	0.09	0.41	0.20	0.08	0.42
46	0.21	0.10	0.50	0.25	0.10	0.40	0.25	0.10	0.39	0.24	0.08	0.35
48	0.27	0.08	0.29	0.32	0.12	0.38	0.32	0.13	0.41	0.29	0.13	0.44
53	0.15	0.05	0.33	0.15	0.05	0.32	0.17	0.06	0.32	0.15	0.04	0.26
57	0.17	0.07	0.43	0.20	0.10	0.47	0.20	0.10	0.51	0.20	0.10	0.49
59	0.26	0.13	0.50	0.28	0.15	0.53	0.27	0.15	0.53	0.29	0.16	0.54

Table 13: Keypoint matching results on individual image pairs in the hard split.



(a) From image pair 58



(b) From image pair 33

Figure 10: Comparison of SIFT keypoint detections. Each column, top to bottom: equirectangular, L0, L1, L2

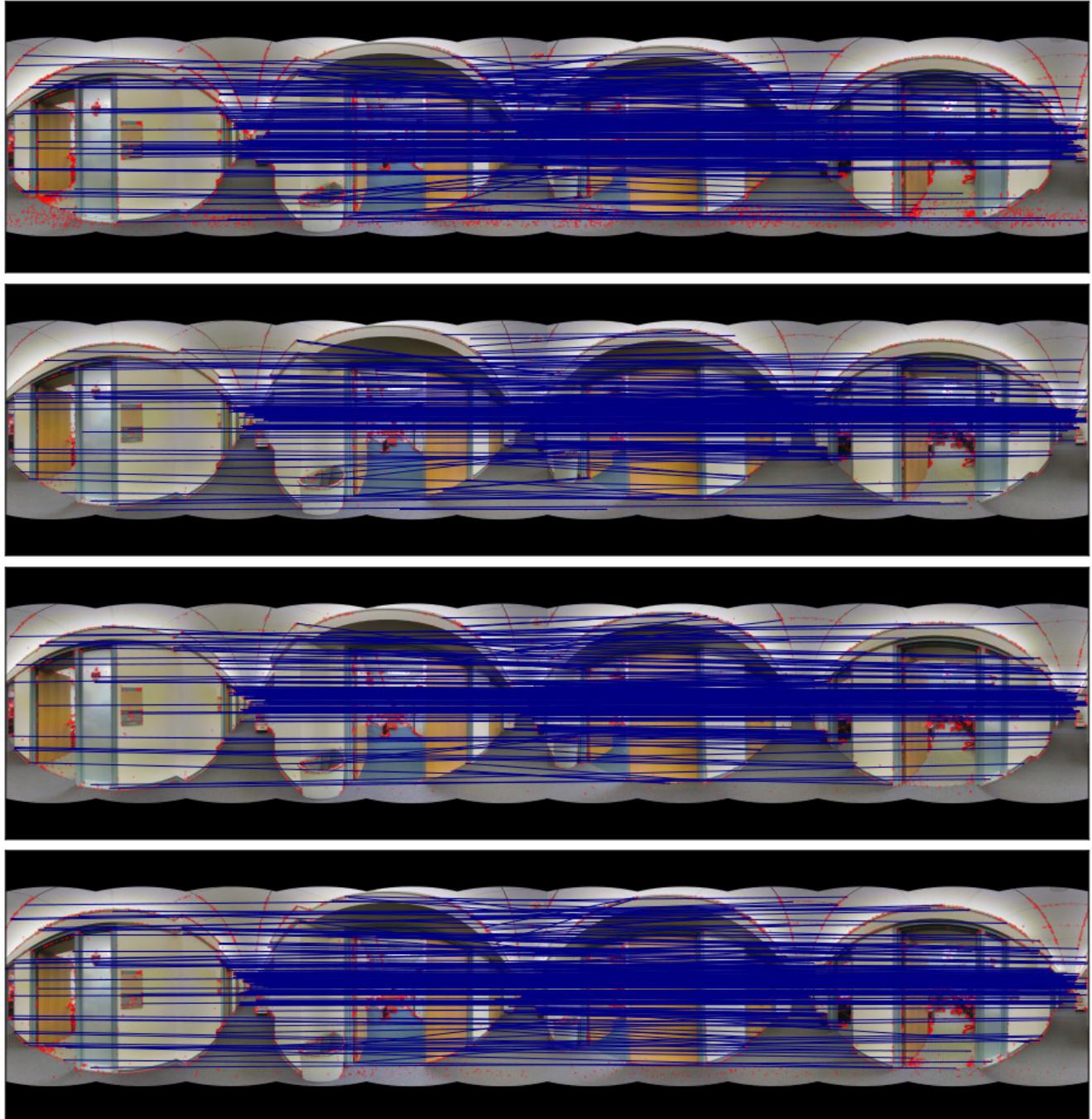


Figure 11: Comparison of SIFT matches on image pair 15. From top to bottom: equirectangular, L0, L1, L2.