

Convolutions on Spherical Images

Marc Eder Jan-Michael Frahm
University of North Carolina at Chapel Hill
{meder, jmf}@cs.unc.edu

Abstract

Applying convolutional neural networks to spherical images requires particular considerations. We look to the millennia of work on cartographic map projections to provide the tools to define an optimal representation of spherical images for the convolution operation. We propose a representation for deep spherical image inference based on the icosahedral Snyder equal-area (ISEA) projection, a projection onto a geodesic grid, and show that it vastly exceeds the state-of-the-art for convolution on spherical images, improving semantic segmentation results by 12.6%.

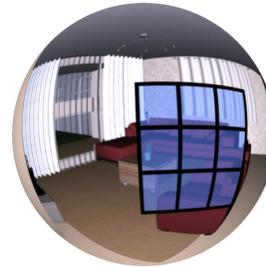


Figure 1: We show that convolving on a planar geodesic approximation to a sphere drastically improves results for the dense prediction task of semantic segmentation.

1. Introduction and Related Work

Omnidirectional imaging is becoming increasingly popular thanks to the proliferation of consumer-grade 360° cameras and the benefit of the wide field-of-view to a number of applications. At the same time, deep learning using convolutional neural networks (CNNs) has never been a more widely-used tool for computer vision tasks. Yet, applying CNNs to spherical images is not quite as trivial as simply training a vanilla CNN on spherical projections onto an image plane. Any planar representation of a sphere necessarily contains some degree of content deformation. This violates convolution’s translation equivariance, which requires the convolution kernel and the signal to be uniformly discretized. To address the unique problems of convolution on spherical images, several methods have been proposed.

Su and Grauman [13] train a CNN to transfer center-perspective-projection-trained models to the equirectangular domain using a learnable adaptive convolutional kernel to match the outputs. Observation that certain regions of cube-map faces contain less distortion than others, Xiong and Grauman [16] develop a content-aware rotation of the spherical data that “snaps-to” a view where the relevant information has minimal distortion. Zioulis *et al.* [17] explore the use of rectangular filter-banks, rather than square filters, to horizontally expand the network’s receptive field to account for distortion in equirectangular images. The spherical convolution derived by Cohen *et al.* [3, 4] uses a gen-

eralized formulation of convolution to filter rotations of the feature maps rather than translations. Finally, work from Coors *et al.* [5] and Tateno *et al.* [14] develop a location-dependent transformation of the planar convolutional kernel using a gnomonic (rectilinear) projection to adapt to the spatial distortion in spherical images.

We take a slightly different approach than these previous works. Rather than changing the tools of convolution, we look to find a better representation of the data. We first analyze the convolution operation to identify the important properties that must be considered when applying convolutions to spherical images and use this to determine a more optimal way to represent the data. We resolve that a planar geodesic approximation to the sphere is a more optimal representation and show that its use yields state-of-the-art performance for dense semantic segmentation.

2. Cartographic Map Projections

First, we briefly review cartographic map projection properties and explain their relevance to this problem.

2.1. Distortion in map projections

Gauss proved with his Theorema Egregium that it is impossible to represent a sphere on a plane without some degree of distortion. The nature of this distortion is dependent on the projection model used to map the sphere to the plane. There are numerous cartographic projections that satisfy

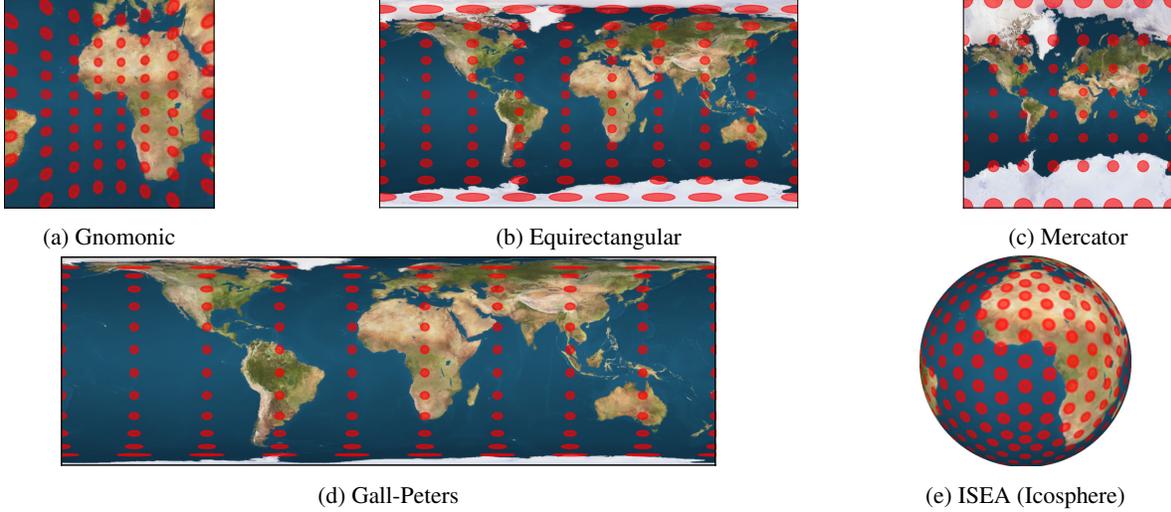


Figure 2: Different map projections with superimposed Tissot’s indicatrices to illustrate the distortion. Ellipse eccentricity shows conformality, size shows areal equality, and spacing shows equidistance. Note that the Gnomonic projection displayed only depicts a $90^\circ \times 90^\circ$ segment of the Earth, as would a face on a cube-map. Earth image from [1].

useful properties for specific applications, like navigation, but generally, there are three major projection properties to consider:

- **Equal area:** Preserves relative scale of objects
- **Conformal:** Preserves local angles
- **Equidistant:** Preserves distance between locations on the sphere

Resulting from Gauss’s theorem, it is impossible for a projection to be both conformal and equal area, and similarly a map cannot be equidistant everywhere [10]. However, certain projections may provide one or another of these properties. A typical method for visualizing the distortion properties of different map projections is to use Tissot’s indicatrices [10], which illustrate the local distortions at points on the map. A Tissot’s indicatrix represents the projection onto a map of a circle with infinitesimal radius on the surface of the sphere. The result is an ellipse where the major and minor axes relate the scale and angular distortions. By placing the circles with regular spacing on the sphere, the indicatrices also depict the distortion of distances in the projection.

Figures 2a-2d show Earth represented by the gnomonic (rectilinear) projection used for faces of cube-maps, the equirectangular (plate carée) projection, the common Mercator projection, and the Gall-Peters projection, respectively, with Tissot’s indicatrices overlaid. The gnomonic projection is neither conformal, equal areal, nor equidistant, demonstrated by the varied eccentricity and shape of the ellipses and the spacing between them. The equirectangular projection is an equidistant projection, preserving the distance between parallels of latitude, but is neither equal-area nor conformal, so the indicatrices are neither circular

nor equal in size. The Mercator projection is conformal but not equal-area, and thus relative scale varies throughout the map. Note how Antarctica appears significantly larger than Africa, despite actually having less than 50% of Africa’s area. All indicatrices remain circular as it is conformal, but they vary in size depending on location. Finally, the Gall-Peters projection is an equal-area projection, which addresses this scale imbalance by preserving the relative size of objects in the map. This method is not conformal, though, and therefore local angles are not preserved, which can be observed by the varying eccentricities of the ellipses.

2.2. Relevance to convolution

The 1D discrete convolution of a filter f of size K and a signal g is given in Eq. (1):

$$(f * g)[n] = \sum_{m=-\lfloor \frac{K}{2} \rfloor}^{\lfloor \frac{K}{2} \rfloor} f[m]g[n - m]. \quad (1)$$

The operation has two components: a sampling and a weighted summation, or more explicitly:

$$(f * g)[n] = \sum_{m=-\lfloor \frac{K}{2} \rfloor}^{\lfloor \frac{K}{2} \rfloor} f[m] \left(\sum_{l=-\infty}^{\infty} g[l] \delta[l - n + m] \right) \quad (2)$$

where $\delta[\cdot]$ is the Dirac delta function. The Dirac delta function can be expressed as the limit of a zero-centered, isotropic Gaussian distribution as the variance goes to zero, and thus in 2D can be represented by an infinitely small circle. Hence when convolving with a spherical image, each sample is an infinitesimal circle on the surface of a sphere.

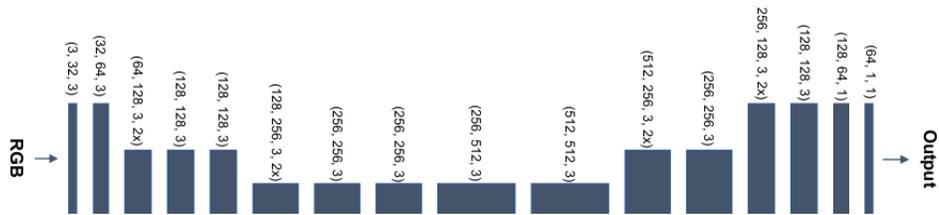


Figure 3: Simple encoder-decoder network architecture used for our experiments. Each layer but the last is followed by an exponential linear unit activation function [2] without batch normalization. The layers are described as $(input_channel, output_channel, kernel_size, up/downsampling)$ and ‘same’ padding is used for all layers.

As discussed previously, we can examine how this sampling is affected by distortion using Tissot’s indicatrix.

From the definition of convolution, the kernel must sample the same area at each location in the data. Thus, there is an implicit assumption that the data is undistorted (i.e. the indicatrices are all perfect circles of equal area). As illustrated by Figure 2b, this is not the case with equirectangular projections. Neither does it hold true for cube-maps, which use the gnomonic projection shown in Figure 2a. These two common spherical image formats violate this distortion-free assumption, which explains why the performance of traditional CNNs on these images falls short when compared to center perspective images that typically have little to no distortion. It follows that there are two viable solutions for resolving this issue: replace the Direct delta function in Eq. (2) with a bivariate, anisotropic Gaussian with location-dependent variance, or resample the image to a different projection with less distortion. The former approach would slow down convolution on spherical images, as a large variance at a sample location could require the accumulation of large regions of pixels, instead of just sampling at one. The latter approach is thus preferable. We still are blocked by the theoretical limitations of map projections, but we can select a compromise projection with the least distortion.

2.3. Icosahedral grids

To re-state, we desire a map projection that is: 1) nearly conformal, preserving the local shape of the data; 2) nearly equal-area, maintaining the local size of the data; and 3) nearly equidistant, allowing fixed-size convolutional filters to retain translational equivariance over the entire image. This suggests that the best projection should yield a very close planar approximation to the sphere itself. Stemming from R. Buckminster Fuller’s work on discretizing the sphere [6], researchers have developed various geodesic approximations for cartographic map projections. One such projection, the icosahedral Snyder equal area (ISEA) projection [11], is particularly well-suited for our task of representing spherical images. It projects the data on the sphere’s surface onto the planar faces of a recursively subdivided icosahedron using a triangle partitioning. This compromise projection is among the least distorted among geodesic grids [7], which satisfies our aforementioned criteria, and it

is easily represented by a triangular 3D mesh. Therefore, we propose that this ISEA projection is an ideal spherical image representation for inference with CNNs.

To perform the ISEA projection, we first create a 7th order ‘‘icosphere,’’ a regular icosahedron subdivided 7 times using Loop subdivision [9] with the vertices subsequently normalized to the unit sphere. We choose a 7th order icosphere because the number of vertices (163,842) is most similar to the number of pixels (131,072) in the 256×512 equirectangular images we use for comparison. Next, we resample the image to the vertices of the icosphere using barycentric interpolation on the faces. The result of this projection is shown in Figure 2e, along with the Tissot indicatrices of deformation on the geodesic’s surface. Notice how the indicatrices are all nearly equally-sized circles. This illustrates that the content deformation is very limited, and the convolution’s sampling, as described in Section 2.2, will be nearly unaffected. We use this representation for all spherical image inputs in our experiments.

3. Semantic Segmentation Evaluation

We evaluate the icosphere representation using semantic segmentation and demonstrate that limiting distortion in the representation spawns significant accuracy improvements.

3.1. Dataset

For our experiments, we use the SUMO [15] dataset, a derivative of SunCG [12] containing 58,631 RGB-D cube-maps representing synthetic spherical image captures. We divide the SUMO dataset into 80% training, 10% validation, and 10% testing splits and pre-process the data by resampling the images onto an 7th order icosphere. An example of this resampling is displayed with a super-imposed convolutional kernel in Figure 1. As the dataset contains 132 semantic classes, some with high inter-class similarity (e.g. ‘bunk bed’ and ‘baby bed’), we aggregated these similar classes to balance the data for these experiments. There is still a very high variance in class frequency even after aggregation, so we only evaluate performance on the 15 most common classes. For exact comparison with existing work, we render the predictions on the icosphere back into an equirectangular image to assess per-pixel accuracy.

Representation	Floor	Ceiling	Wall	Door	Cabinet	Rug	Window	Curtain
Equirect. (gnom. kernel) [5, 14]	0.9315	0.9710	0.8597	0.6466	0.6376	0.7284	0.7012	0.4703
Equirect. (equirect. kernel)	0.9327	0.9654	0.8566	0.6692	0.6638	0.7003	0.7180	0.4248
Icosphere (ours)	0.9352	0.9703	0.8797	0.6890	0.7037	0.6970	0.7562	0.5744
	Sofa	Partition	Bed	Chair	Table	Shelving	Chandelier	All Classes
Equirect. (gnom. kernel) [5, 14]	0.7114	0.4172	0.7133	0.4219	0.4587	0.3278	0.4491	0.5904
Equirect. (equirect. kernel)	0.7337	0.4448	0.7228	0.4415	0.5363	0.2843	0.4566	0.5969
Icosphere (ours)	0.7374	0.4683	0.7776	0.4375	0.5018	0.3733	0.4472	0.6639

Table 1: Semantic segmentation mIOU for the different representations. Classes ordered by descending frequency.

3.2. CNN Training

We follow the method of Coors *et al.* [5] and Tateno *et al.* [14] to sample the spherical data using a gnomonic projection of the planar filter. However, as our data is no longer in a 2D image, we project the planar filter onto the icosphere’s surface and sample from the faces using barycentric interpolation. We set the filter resolution to the mean distance between adjacent vertices on the icosphere, to maintain analogy with the notion of resolution of the planar filter being set to the distance between adjacent pixels. During convolution, we apply this filter at each vertex of the icosphere and sample from the faces using barycentric interpolation of the data at the 3 vertices. For up- and down-sampling, we leverage different subdivision orders, which, similar to pixels in 2D images, vary roughly by a factor of 4 in the number of vertices. We compare convolving with our icosphere representation to the method of [5, 14], which convolves with equirectangular images. We use the encoder-decoder architecture given in Figure 3 for both methods. We train for 10 epochs using Adam [8] with an initial learning rate of 10^{-4} , reduced by half every 3 epochs. Our training criterion is the cross-entropy loss weighted by inverse class frequency.

3.3. Results and Discussion

Results of our semantic segmentation experiments are given in Table 1 as the per-class mean intersection-over-union (mIOU). While certain classes perform marginally better in the equirectangular format, we see an 12.6% improvement in the overall mIOU compared to [5, 14] simply by resampling the image. There are two reasons for this jump. The first is that the gnomonic projection used for the kernel in [5, 14] is not optimal for modeling the distortion in an equirectangular image. In fact, an equirectangular projection of the kernel would be a more appropriate, as it would correctly adapt the sampling locations to equirectangular distortion. The difference in the sampling patterns is slight, though, as shown in Figure 4, and the quantitative result is only slightly better (second row in Table 1). It is worth mentioning that the gnomonic kernel does make sense for sampling from the icosphere, as it represents the projection of a tangent plane onto a sphere. In our case, the planar filter is tangent to our icosahedral spherical approx-

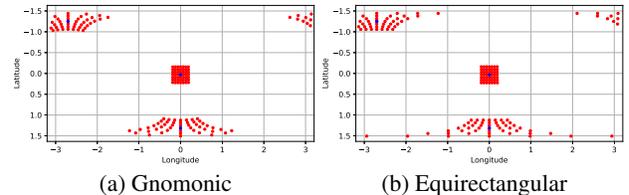


Figure 4: Comparing convolutional kernel sampling patterns using gnomonic and equirectangular projections.

imation. Nevertheless, the more important consideration is that neither the gnomonic nor equirectangular filter projections correct for the information imbalance caused by the distortion. The filter projection method of [5, 14] dynamically adjusts where the kernel samples, but it only accounts for the conformality and equidistance properties of distortion. It does not address the areal inequality. Unequal area projections disrupt the distribution of information throughout the image. For example, all of the pixels along the top row in an equirectangular images are essentially redundant, having all been sampled from the pole. That whole row carries the same information capacity of a single pixel at the equator. For convolutions, this means that the amount of information aggregated at each pixel varies throughout the image. The ISEA projection, being nearly equal-area, corrects for this, ensuring that predictions at each vertex of the icosphere are made by accumulating equivalent quantities of information. This translates to better inference. It is also worth noting that the objects for which the equirectangular representation perform marginally better are typically found at the top (ceiling, chandelier) or bottom (rug) of an image. The equirectangular format’s pixel redundancy increases the number of training samples for these objects, which would explain the better performance at test time.

4. Conclusion

Borrowing from the field of cartography, we have identified three key principles of local distortion that impact our ability to convolve on spherical data. Applying these principles to our choice of data representation, we have proposed the ISEA projection on an icosphere as an ideal spherical image representation. We demonstrated that by significantly limiting distortion, this representation provides a 12.6% boost in semantic segmentation performance.

References

- [1] Planet texture maps. <http://planetpixelemposium.com/earth8081.html>. Accessed: 2018-04-16. 2
- [2] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 2015. 3
- [3] Taco Cohen, Mario Geiger, Jonas Köhler, and Max Welling. Convolutional networks for spherical signals. *arXiv preprint arXiv:1709.04893*, 2017. 1
- [4] Taco S Cohen, Mario Geiger, Jonas Köhler, and Max Welling. Spherical cnns. *arXiv preprint arXiv:1801.10130*, 2018. 1
- [5] 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, 4
- [6] R Buckminster Fuller. *Synergetics: explorations in the geometry of thinking*. Estate of R. Buckminster Fuller, 1982. 3
- [7] 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. 3
- [8] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 4
- [9] Charles Loop. Smooth subdivision surfaces based on triangles. *Master’s thesis, University of Utah, Department of Mathematics*, 1987. 3
- [10] John Parr Snyder. *Map projections—A working manual*, volume 1395. US Government Printing Office, 1987. 2
- [11] John P Snyder. An equal-area map projection for polyhedral globes. *Cartographica: The International Journal for Geographic Information and Geovisualization*, 29(1):10–21, 1992. 3
- [12] Shuran Song, Fisher Yu, Andy Zeng, Angel X Chang, Manolis Savva, and Thomas Funkhouser. Semantic scene completion from a single depth image. *IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 3
- [13] 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
- [14] 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, 4
- [15] Lyne Tchapmi and Daniel Huber. The sumo challenge. 3
- [16] Bo Xiong and Kristen Grauman. Snap angle prediction for 360 panoramas. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–18, 2018. 1
- [17] Nikolaos Zioulis, Antonis Karakottas, Dimitrios Zarpalas, and Petros Daras. Omnidepth: Dense depth estimation for indoors spherical panoramas. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 448–465, 2018. 1