# EARTHMOVER-BASED MANIFOLD LEARNING FOR ANALYZING MOLECULAR CONFORMATION SPACES

*Nathan Zelesko*♠     *Amit Moscovich*◇     *Joe Kileel*◇     *Amit Singer*◇,♣

♠ Department of Mathematics, Brown University
◇ Program in Applied and Computational Mathematics, Princeton University
♣ Department of Mathematics, Princeton University

## ABSTRACT

In this paper, we propose a novel approach for manifold learning that combines the Earthmover's distance (EMD) with the diffusion maps method for dimensionality reduction. We demonstrate the potential benefits of this approach for learning shape spaces of proteins and other flexible macromolecules using a simulated dataset of 3-D density maps that mimic the non-uniform rotary motion of ATP synthase. Our results show that EMD-based diffusion maps require far fewer samples to recover the intrinsic geometry than the standard diffusion maps algorithm that is based on the Euclidean distance. To reduce the computational burden of calculating the EMD for all volume pairs, we employ a wavelet-based approximation to the EMD which reduces the computation of the pairwise EMDs to a computation of pairwise weighted-$\ell_1$ distances between wavelet coefficient vectors.

***Index Terms***— Shape space, dimensionality reduction, Wasserstein metric, diffusion maps, Laplacian eigenmaps, cryo-electron microscopy

## 1. INTRODUCTION

Proteins and other macromolecules are elastic structures that may deform in various ways. Since the spatial conformation of an organic molecule is known to play a key role in its biological function, the complete description of a molecule must include more than just a single static structure (as is traditionally produced by X-ray crystallography). Ideally, we would like to map the entire space of molecular conformations. However, understanding the topology and geometry of these conformation spaces remains one of the grand challenges in the field of structural biology [1].

One promising approach is to employ cryo-electron microscopy (cryo-EM) as a tool for structure determination in the presence of conformational heterogeneity [2]. In cryo-EM, multiple images of a particular macromolecule are taken by a transmission electron microscope and then processed using specialized algorithms. Traditionally, these algorithms construct an estimate of the mean molecular volume, in the
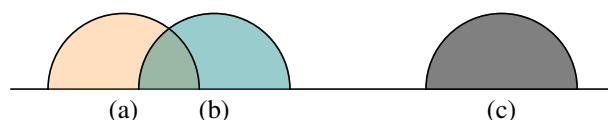


**Fig. 1**. *EMD vs. Euclidean distance for translational motion.* Euclidean (or $\ell_2$) distance is only meaningful for measuring small displacements. e.g., the $\ell_2$ distance between half-disks (c) and (a) is the same as between (c) and (b). By contrast, for any translational motion, the EMD is its magnitude.

form of a 3-D electrostatic density map. In particular, this process averages out any variability in the spatial conformations of the molecules in the sample. Recent works have applied techniques from the field of manifold learning to cryo-EM data sets, obtaining a low-dimensional representation of the molecular conformation space [3, 4]. Specifically, these works build affinity graphs based on the Euclidean distances between molecular volumes (or projection images) and then compute diffusion map embeddings [5, 6].

However, the Euclidean distance is suboptimal for capturing the distance between geometric conformations. Consider, for example, two conformations of a molecule that has only a single moving part. If the two conformations are distant, the support of the moving part in the two volumes may not intersect, rendering the Euclidean distance independent of the conformational distance. See Fig. 1. In such cases, in order to apply manifold learning based on a Euclidean metric, one need a dense cover of the conformation space by the molecules in the sample. Since the number of points in such a cover scales exponentially in the dimension, it may be infeasible to apply these methods, even using the largest existing experimental datasets, which consist of about $n \approx 10^6$ samples.

In this paper, we propose to use the *Earthmover's distance* (EMD), also known as the *Wasserstein metric*, instead of the commonly used Euclidean distance as input to manifold embedding algorithms. EMD has an intuitive geometric meaning: it measures the minimal amount of "work" needed to transform one pile of mass into another pile of equal mass,

where "work" is defined as the amount of mass moved times the distance by which it is moved. In particular, EMD provides a distance metric that is meaningful even between spatial conformations that are far from each other. Following the discussion above, this property should reduce the number of samples needed to learn the intrinsic manifold.

Methods for computing the EMD, based on off-the-shelf linear programming solvers, are expensive when the number of voxels is large. Therefore we used a fast approximation to the EMD, based on a wavelet representation [7].

To test our proposal, we compared the standard $\ell_2$-based diffusion maps to EMD-based diffusion maps on a synthetic dataset mimicking the motion of ATP synthase (Fig. 2). This dataset samples the underlying manifold in a non-uniform manner since ATP synthase has three dominant conformations that are 120° apart [8]. The approximate EMD-based approach yields a marked improvement in the number of samples required for learning the conformational manifold, while still offering a computationally feasible algorithm.

## 2. METHODS

In this section, we review the basic techniques that underlie Earthmover-based manifold learning. Our current focus is on learning shape spaces of 3-D volumes, but the same techniques may also be applied to analyze other types of datasets, such as 2-D image sets, 1-D histograms, etc. To start, let $X = \{\mathbf{x}_1, \ldots, \mathbf{x}_n\}$ be a set of 3-D voxel arrays in $\mathbb{R}^{L^3}$. We assume that $X$ obeys the *manifold hypothesis* [2, 9, 10], i.e., $\mathbf{x}_1, \ldots, \mathbf{x}_n$ form a (noisy) sample of a low-dimensional manifold $\mathcal{M} \subset \mathbb{R}^{L^3}$. Our task is to reorganize the data to better reflect the intrinsic geometry of $\mathcal{M}$.

For Riemannian manifolds, eigenfunctions of the *Laplace-Beltrami operator* provide an intrinsic coordinate system [11, 12]. Accordingly, several popular methods for dimensionality reduction and data representation methods are based on mapping input points using empirical estimates of Laplacian eigenfunctions [6, 5]. Under the manifold hypothesis, these estimates converge to eigenfunctions of the Laplace-Beltrami operator, or more generally to eigenfunctions of a weighted Laplacian, depending on the construction [13].

We now describe the diffusion maps method [6]. Let $w : \mathbb{R}^{L^3} \times \mathbb{R}^{L^3} \to \mathbb{R}$ denote a symmetric non-negative function that gives an affinity score for each pair of volumes. One common way of constructing affinities is to take a distance metric $d : X \times X \to \mathbb{R}$ and apply a Gaussian kernel with a suitably chosen width $\sigma$ to form the *affinity matrix* $W \in \mathbb{R}^{n \times n}$

$$W_{ij} = w(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-d(x_i, x_j)^2/(2\sigma^2)\right). \quad (1)$$

The *degree matrix* $D \in \mathbb{R}^{n \times n}$ is defined to be the diagonal matrix that satisfies $D_{ii} = \sum_{j=1}^{n} W_{ij}$. We use the *Coifman-Lafon normalized graph Laplacian* [14], which converges to the Laplace-Beltrami operator, regardless of the sampling
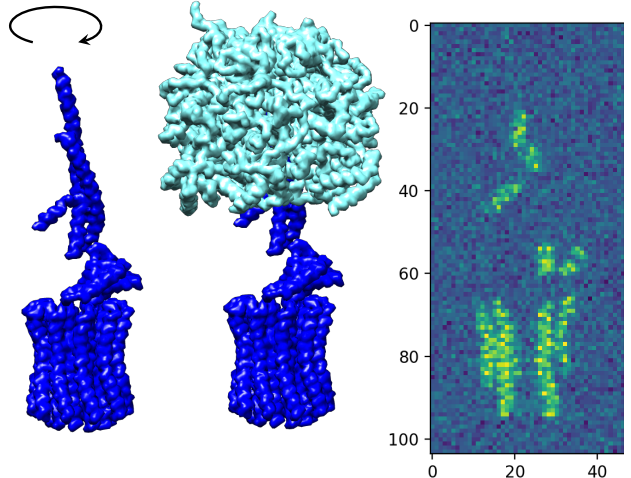


**Fig. 2**. *ATP synthase.* (left) $F_0$ and axle subunits. They rotate together in the presence of hydrogen ions, forming a tiny electric motor; (middle) the $F_1$ subunit (in cyan) envelops the axle. As the axle rotates, this subunit assembles ATP; (right) sample slice of the rotated $F_0$ and axle subunits with the additive Gaussian noise.

density. To compute this, one first performs a two-sided normalization of the affinity matrix, $\widetilde{W} = D^{-1} W D^{-1}$ and then computes the random-walk Laplacian, $\mathcal{L} = \widetilde{D}^{-1} \widetilde{W}$, where $\widetilde{D}$ is the degree matrix for $\widetilde{W}$. The random-walk Laplacian is similar to a positive semi-definite symmetric matrix and hence its eigenvectors are real and its eigenvalues are non-negative. The all-ones vector is an eigenvector of $\mathcal{L}$ with eigenvalue zero [15]. Let $\phi_0, \phi_1, \ldots, \phi_{n-1} \in \mathbb{R}^n$ be eigenvectors of $\mathcal{L}$ with corresponding eigenvalues $0 = \lambda_0 \leq \lambda_1 \leq \ldots \leq \lambda_{n-1}$. We think of the eigenvectors $\phi_\ell$ as real-valued functions on $X$, by identifying $\phi_\ell(\mathbf{x}_i) = (\phi_\ell)_i$. The *k-dimensional diffusion map* $\Psi_t^{(k)} : X \to \mathbb{R}^k$ is defined by:

$$\mathbf{x}_i \mapsto \left(\lambda_1^t \phi_1(\mathbf{x}_i), \ldots, \lambda_k^t \phi_k(\mathbf{x}_i)\right).$$

The mapping $\Psi_t^{(k)}$ gives a system of $k$ coordinates on $X$, which captures the intrinsic geometry of $\mathcal{M}$. In our simulations, we used $t = 0$, in which case diffusion maps coincide with *Laplacian eigenmaps* [5].

The diffusion map depends on the choice of affinity. The typical choice is a Gaussian kernel as defined in Eq. (1) that is based on a Euclidean (or $\ell_2$) distance function,

$$d_{\ell_2}(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\|_2.$$

We propose instead to base the Gaussian kernel of Eq. (1) on the *Earthmover's distance (EMD)*, also known as the *Wasserstein metric* [16]. EMD is popular in various applications, e.g., image retrieval [17], however, to the best of our knowledge, it has never been used to define affinities for manifold

learning algorithms. To define this distance, consider two 3-D density maps $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^{L^3}$ that are non-negative and normalized to unit mass. These densities define probability measures on the set of voxels, $[L]^3$, where $[L] = \{1, \ldots, L\}$. We set:

$$d_{\text{EMD}}(\mathbf{x}_i, \mathbf{x}_j) = \min_{\pi \in \Pi(\mathbf{x}_i, \mathbf{x}_j)} \sum_{\mathbf{u} \in [L]^3} \sum_{\mathbf{v} \in [L]^3} \pi(\mathbf{u}, \mathbf{v}) \|\mathbf{u} - \mathbf{v}\|_2,$$

where $\Pi(\mathbf{x}_i, \mathbf{x}_j)$ is the set of joint probability measures on $[L]^3 \times [L]^3$ with marginals $\mathbf{x}_i$ and $\mathbf{x}_j$, respectively.

Mathematically, EMD amounts to a linear program in $\mathcal{O}(L^6)$ variables subject to $\mathcal{O}(L^3)$ constraints, i.e., a significant computation. However, in the wavelet domain [18], EMD enjoys a fast (weighted-$\ell_1$) wavelet approximation [7], which we refer to as WEMD:

$$d_{\text{WEMD}}(\mathbf{x}_i, \mathbf{x}_j) = \sum_\lambda 2^{-5s/2} |\mathcal{W}\mathbf{x}_i(\lambda) - \mathcal{W}\mathbf{x}_j(\lambda)|. \quad (2)$$

Here, $\mathcal{W}\mathbf{x}$ denotes the *3-D wavelet transform* of $\mathbf{x}$, and the index $\lambda$ contains the shifts $(m_1, m_2, m_3) \in \mathbb{Z}^3$ and the scale $s \in \mathbb{Z}_{\geq 0}$. More explicitly, $\mathcal{W}$ decomposes $\mathbf{x} = \mathbf{x}[u_1, u_2, u_3]$ with respect to an orthonormal basis of functions,

$$2^{3s/2} f(2^s u_1 - m_1) \, g(2^s u_2 - m_2) \, h(2^s u_3 - m_3),$$

for varying $s \in \mathbb{Z}_{\geq 0}$, varying $(m_1, m_2, m_3) \in \mathbb{Z}^3$, and $(f, g, h)$ ranging over $\{\psi, \omega\}^3 \setminus \{(\omega, \omega, \omega)\}$ where $\psi, \omega$ are certain 1-D functions called the *mother and father wavelet* [18]. Formula (2) approximates EMD in the sense that $d_{\text{EMD}}$ and $d_{\text{WEMD}}$ are strongly equivalent metrics, i.e., there exist constants $C \geq c > 0$ such that for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{L^3}$, we have:

$$c \cdot d_{\text{WEMD}}(\mathbf{x}, \mathbf{y}) \leq d_{\text{EMD}}(\mathbf{x}, \mathbf{y}) \leq C \cdot d_{\text{WEMD}}(\mathbf{x}, \mathbf{y}).$$

Moreover, there are known bounds on the ratio $C/c$, depending on the type of wavelet used. We have chosen the Coiflets 3 wavelet since it gives a small ratio [7]. Wavelet transforms are computed in linear time, thus the same holds for the EMD approximation. We implemented the approximation (2), using the PyWavelets package [19]. We computed the wavelet transform up to scale $s = 5$ for accurate truncation of Eq. (2). This overparameterizes the volumes by a factor of $\approx 3$.

## 3. RESULTS

To test our methods, we generated two synthetic datasets of 3-D density maps that are simplified models of the conformation space of ATP synthase [8]. This enzyme is a molecular stepper motor with a central asymmetric axle that rotates in steps of $120°$ relative to the $F_1$ subunit, with short transient motion in-between the three dominant conformations. Here, the intrinsic geometry is a circle, with a sampling density concentrated around three equispaced angles. We simulated this motion by generating 3-D density maps in which the $F_1$ subunit is held in place while the $F_0$ and axle subunits are rotated
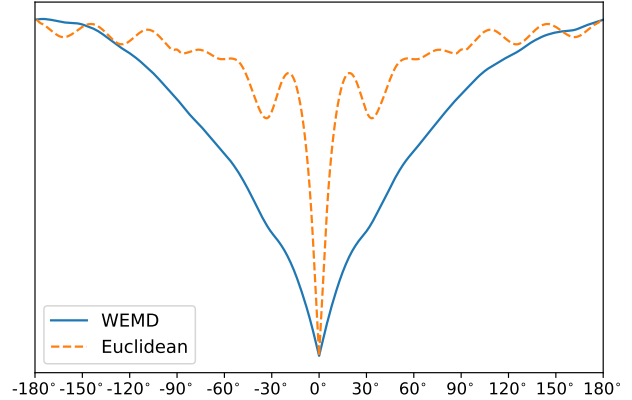


**Fig. 3**. *Euclidean distance vs. WEMD* as functions of the angle between two angles of the ATP synthase rotor (see Fig. 2). The $\ell_2$ distances are scaled to be comparable to the WEMD. The WEMD is monotone in the magnitude of the angular difference for almost the entire range whereas the Euclidean distance exhibits this behavior only up to about $\pm 19°$.

| $n$ | Wavelet transform | WEMD distances | $\ell_2$ distances |
|---|---|---|---|
| 25 | 60 | 0.9 | 0.8 |
| 50 | 121 | 3.4 | 1.2 |
| 100 | 243 | 13 | 2.4 |
| 200 | 481 | 50 | 4.8 |
| 400 | 965 | 221 | 9.4 |
| 800 | 1932 | 844 | 20.6 |

**Fig. 4**. *Running times [sec]* for computing the fast wavelet transform, all pairwise wavelet Earthmover approximations and all pairwise $\ell_2$ distances.

together by a random angle. The angles were drawn i.i.d. according to the following mixture model:

$$\frac{2}{5} U[0, 360] + \frac{1}{5} \mathcal{N}(0, 1) + \frac{1}{5} \mathcal{N}(120, 1) + \frac{1}{5} \mathcal{N}(240, 1),$$

where $U$ and $\mathcal{N}$ denote uniform and Gaussian distributions, respectively. To form our datasets, we downloaded entry 1QO1 [20] from the Protein Data Bank [21], produced 3-D density maps at a 6Å resolution with array dimensions $47 \times 47 \times 107$ using the `molmap` command in UCSF Chimera [22], and then took random rotations of the $F_0$ and axle subunits. From this, we generated a clean dataset and a noisy dataset. For the latter, i.i.d. Gaussian noise was added with mean zero and standard deviation equal to one-tenth of the maximum voxel value.

We first tested the plausibility of our proposal by comparing the EMD approximation to the Euclidean distance for a range of angular differences using the noiseless dataset (Fig. 3). We then performed 2-dimensional diffusion maps for various sample sizes, using both the Euclidean distance
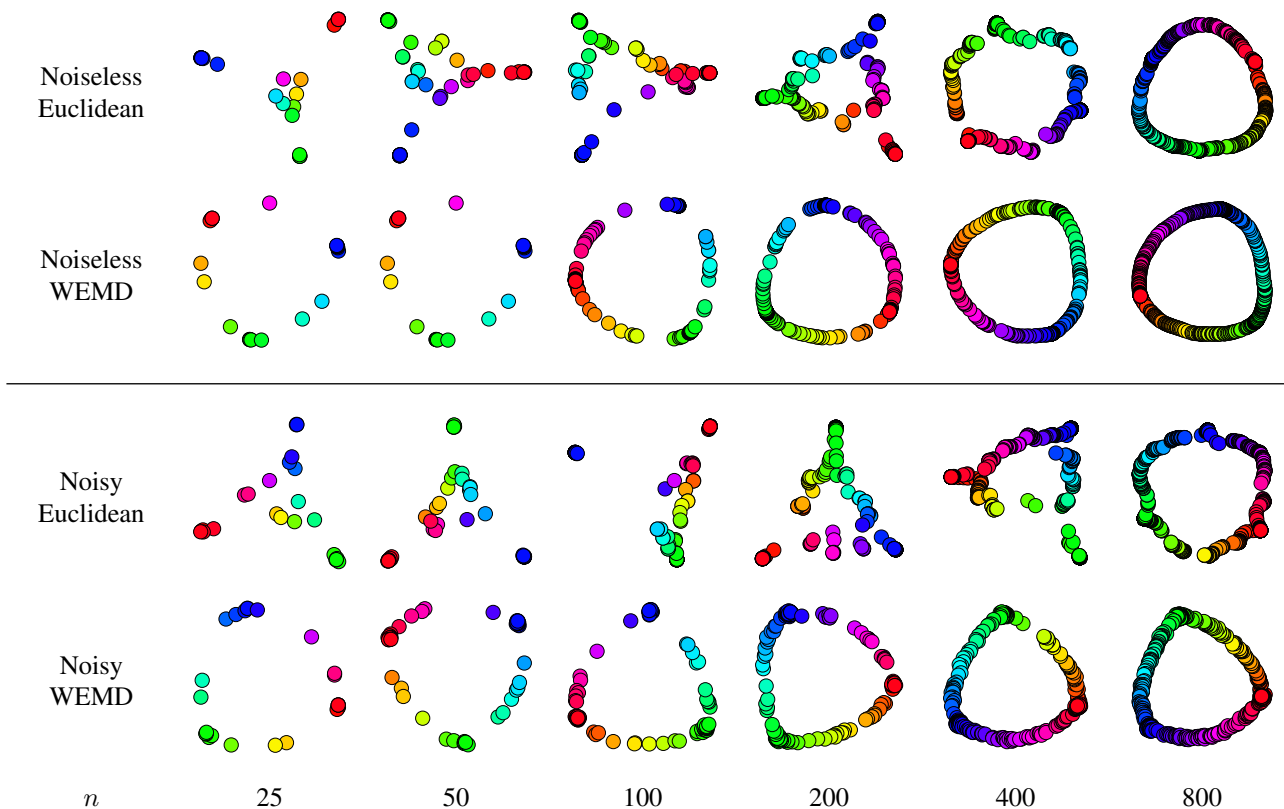
**Fig. 5**. *Main results.* Euclidean vs. EMD-based diffusion mappings on the clean and noisy ATP synthase datasets for sample sizes $n = 25, 50, 100, 200, 400, 800$. The Euclidean diffusion maps need more than $400$ samples to capture the intrinsic geometry whereas WEMD manages to do so with merely $n = 25$ samples. The colors encode the (ground truth) angle.

and the wavelet-based approximation to the EMD, as described in the previous section. The resulting embeddings are shown in Fig. 5. The value of the width parameter $\sigma$ in the Gaussian kernel (1) was hand-picked to yield the best results. We note that for the Euclidean diffusion maps, careful tuning of $\sigma$ was required. However, this was not necessary for the EMD approximation, where a wide range of $\sigma$ values gave excellent results. Running times (on an Intel Core i7) for the computation of EMD and Euclidean-based diffusion maps are listed in Fig. 4.

## 4. CONCLUSION

In this paper, we proposed to use Earthmover-based affinities in the diffusion maps framework to analyze molecular conformation spaces. We showed that this results in a marked decrease in the number of samples needed to capture the intrinsic conformation space of ATP synthase. The method is computationally tractable, thanks to a fast wavelet approximation, and robust to noise. Our results show promise, particularly for the analysis of cryo-EM datasets with continuous heterogeneity. More broadly, EMD-based manifold learning could be applied to analyze the variability of other collec-

tions of 3-D shapes [23], 2-D images [17], videos and other signals, e.g., to better model animal motion [24]. Our work also raises several interesting theoretical questions: in which cases can one prove that EMD-based manifold learning has a lower sample complexity than manifold learning based on the Euclidean distance? More ambitiously, are there reasonable generative models for variability where EMD is the optimal distance metric?

## 5. REPRODUCIBILITY

Code for reproducing the results in this paper is available at http://github.com/nathanzelesko/earthmover

## 6. ACKNOWLEDGMENTS

# 7. REFERENCES

[1] J. Frank, "New opportunities created by single-particle cryo-EM: the mapping of conformational space," *Biochemistry*, vol. 57, no. 6, pp. 888, 2018. doi:10.1021/acs.biochem.8b00064

[2] C.O.S. Sorzano et al., "Survey of the analysis of continuous conformational variability of biological macromolecules by electron microscopy," *Acta Crystallogr. Sect. F Struct. Biol. Commun.*, vol. 75, no. 1, pp. 19–32, 2019. doi:10.1107/S2053230X18015108

[3] P. Schwander, R. Fung, and A. Ourmazd, "Conformations of macromolecules and their complexes from heterogeneous datasets," *Philos. Trans. R. Soc. B Biol. Sci.*, vol. 369, no. 1647, pp. 1–8, 2014. doi:10.1098/rstb.2013.0567

[4] A. Moscovich, A. Halevi, J. Andén, and A. Singer, "Cryo-EM reconstruction of continuous heterogeneity by Laplacian spectral volumes," *Inverse Probl.,* submitted, 2019.

[5] M. Belkin and P. Niyogi, "Laplacian eigenmaps for dimensionality reduction and data representation," *Neural Comput.*, vol. 15, no. 6, pp. 1373–1396, 2003. doi:10.1162/089976603321780317

[6] R.R. Coifman et al., "Geometric diffusions as a tool for harmonic analysis and structure definition of data: diffusion maps," *PNAS*, vol. 102, no. 21, pp. 7426–7431, 2005. doi:10.1073/pnas.0500334102

[7] S. Shirdhonkar and D.W. Jacobs, "Approximate Earthmover's distance in linear time," in *CVPR 2008*. doi:10.1109/CVPR.2008.4587662

[8] M. Yoshida, E. Muneyuki, and T. Hisabori, "ATP synthase – a marvellous rotary engine of the cell," *Nat. Rev. Mol. Cell Biol.*, vol. 2, no. 9, pp. 669–677, 2001. doi:10.1038/35089509

[9] A. Moscovich, A. Jaffe, and B. Nadler, "Minimax-optimal semi-supervised regression on unknown manifolds," in *AISTATS 2017*. http://proceedings.mlr.press/v54/moscovich17a.html

[10] A.B. Lee, K.S. Pedersen, and D. Mumford, "The nonlinear statistics of high-contrast patches in natural images," *Int. J. Comput. Vis.*, vol. 54, no. 1-3, pp. 83–103, 2003. doi:10.1023/A:1023705401078

[11] P. Bérard, G. Besson, and S. Gallot, "Embedding Riemannian manifolds by their heat kernel," *Geom. Funct. Anal.*, vol. 4, no. 4, pp. 373–398, 1994. doi:10.1007/BF01896401

[12] P.W. Jones, M. Maggioni, and R. Schul, "Manifold parametrizations by eigenfunctions of the Laplacian and heat kernels," *PNAS*, vol. 105, no. 6, pp. 1803–1808, 2008. doi:10.1073/pnas.0710175104

[13] D. Ting, L. Huang, and M. Jordan, "An analysis of the convergence of graph Laplacians," in *ICML 2010*. https://icml.cc/Conferences/2010/papers/554.pdf

[14] R.R. Coifman and S. Lafon, "Diffusion maps," *Appl. Comput. Harmon. Anal.*, vol. 21, no. 1, pp. 5–30, 2006. doi:10.1016/j.acha.2006.04.006

[15] U. von Luxburg, "A tutorial on spectral clustering," *Stat. Comput.*, vol. 17, no. 4, pp. 395–416, 2007. doi:10.1007/s11222-007-9033-z

[16] C. Villani, *Optimal Transport: Old and New*, Springer Berlin Heidelberg, 2009. doi:10.1007/978-3-540-71050-9

[17] Y. Rubner, C. Tomasi, and L.J. Guibas, "The Earthmover's distance as a metric for image retrieval," *Int. J. Comput. Vis.*, vol. 40, no. 2, pp. 99–121, 2000. doi:10.1023/A:1026543900054

[18] S. Mallat, *A Wavelet Tour of Signal Processing*, Elsevier, 3rd edition, 2009. doi:10.1016/B978-0-12-374370-1.X0001-8

[19] G. Lee, R. Gommers, F. Waselewski, K. Wohlfahrt, and A. O'Leary, "PyWavelets: a Python package for wavelet analysis," *J. Open Source Softw.*, vol. 4, no. 36, pp. 1237, 2019. doi:10.21105/joss.01237

[20] D. Stock, A.G. Leslie, and J.E. Walker, "Molecular architecture of the rotary motor in ATP synthase," *Science*, vol. 286, no. 5445, pp. 1700–1705, 1999. doi:10.1126/science.286.5445.1700

[21] P.W. Rose et al., "The RCSB protein data bank: integrative view of protein, gene and 3D structural information," *Nucleic Acids Res.*, vol. 45, no. D1, pp. D271–D281, 2017. doi:10.1093/nar/gkw1000

[22] E.F. Pettersen at al., "UCSF Chimera – a visualization system for exploratory research and analysis," *J. Comput. Chem.*, vol. 25, no. 13, pp. 1605–1612, 2004. doi:10.1002/jcc.20084

[23] M. Ovsjanikov, W. Li, L.J. Guibas, and N.J. Mitra, "Exploration of continuous variability in collections of 3D shapes," *ACM Trans. Graph.*, vol. 30, no. 4, pp. 1–10, 2011. doi:10.1145/2010324.1964928

[24] D.L. Hu, J. Nirody, T. Scott, and M.J. Shelley, "The mechanics of slithering locomotion," *PNAS*, vol. 106, no. 25, pp. 10081–10085, 2009. doi:10.1073/pnas.0812533106