
Manifold-Consistent Graph Indexing: Overcoming the Euclidean-Geodesic Mismatch via Local Intrinsic Dimensionality

Dongfang Zhao¹

Abstract

Retrieval-augmented generation (RAG) and approximate nearest neighbor (ANN) search have been critical components of modern large language model (LLM) serving services as they enable efficient and effective retrieval of relevant information to reduce LLM’s hallucination. However, state-of-the-art methods are mostly based on graph indexing techniques that are agnostic to the intrinsic geometry of the data, and thus often perform poorly in high-dimensional spaces due to a Euclidean-Geodesic mismatch. To that end, we propose a new graph indexing method called Manifold-Consistent Graph Indexing (MCGI). The key idea of MCGI is to leverage the local intrinsic dimensionality (LID) of the data to construct a graph that is consistent with the underlying manifold structure, thereby reducing the mismatch and improving performance. Our theoretical analysis shows that MCGI achieves improved approximation guarantees comparing to existing methods, such as HNSW and DiskANN. We also report experimental results demonstrating that MCGI outperforms existing methods in various benchmarks and real-world applications.

1. Introduction

2. Related Work

3. Methodology

3.1. Definitions

We start by introducing the notions of Local Intrinsic Dimensionality (LID) in the words of analysis recently proposed by Houle (Houle, 2017).

Definition 3.1 (Local Intrinsic Dimensionality). Let \mathcal{X} be a

¹University of Washington, Tacoma School of Engineering & Technology and Paul G. Allen School of Computer Science & Engineering. Correspondence to: Dongfang Zhao <dzhao@uw.edu>.

domain equipped with a distance measure $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}^+$. For a reference point $x \in \mathcal{X}$, let $F_x(r) = \mathbb{P}(d(x, Y) \leq r)$ denote the cumulative distribution function (CDF) of the distance between x and a random variable Y drawn from the underlying data distribution. The Local Intrinsic Dimensionality (LID) of x , denoted as $ID(x)$, is defined as the intrinsic growth rate of the probability measure within the neighborhood of x :

$$ID(x) \triangleq \lim_{r \rightarrow 0} \frac{r \cdot F'_x(r)}{F_x(r)} = \lim_{r \rightarrow 0} \frac{d \ln F_x(r)}{d \ln r}, \quad (1)$$

provided the limit exists and $F_x(r)$ is continuously differentiable for $r > 0$.

Remark 3.2 (Institution of LID). The definition of LID can be understood as a measure of the multiplicative growth rate of the volume of a ball centered at x with radius r as r approaches 0. Let D denote the dimensionality of the ambient space. If the data lies on a local D -dimensional manifold, then the CDF around an infinitely small neighborhood of x satisfies:

$$F_x(r) \approx C \cdot r^D, \quad (2)$$

where C is a constant. Thus, the following holds:

$$F'_x(r) \approx C \cdot D \cdot r^{D-1}. \quad (3)$$

Combining equations (2) and (3), we get:

$$D \approx \frac{F'_x(r)}{F_x(r)} \cdot r, \quad (4)$$

thus Eq. (1).

While Eq. 3.1 provides an intuitive closed-form formula for intrinsic dimensionality, in practice we usually do not have access to the true CDF $F_x(r)$. Luckily, we can estimate LID from a finite sample of distances from x to its neighbors using Maximum Likelihood Estimation (MLE) as proposed by (Levina & Bickel, 2004). According to (Amsaleg et al., 2015), LID can be estimated as follows.

Definition 3.3 (LID Maximum Likelihood Estimator). Given a reference point x and its k -nearest neighbors determined by the distance measure d , let $r_i = d(x, v_i)$ denote the distance to the i -th nearest neighbor, sorted such that

$r_1 \leq \dots \leq r_k$. Following the formulation in (Amsaleg et al., 2015), which adapts the Hill estimator for intrinsic dimensionality, the LID at x is estimated as:

$$\widehat{\text{LID}}(x) = -\left(\frac{1}{k} \sum_{i=1}^k \ln \frac{r_i}{r_k}\right)^{-1}. \quad (5)$$

3.2. Mapping Function

The primary goal of Manifold-Consistent Graph Indexing is that the graph topology should adapt to the local geometric complexity. In regions where the Local Intrinsic Dimensionality (LID) is low, the data manifold approximates a flat Euclidean subspace. In such isotropic regions, the Euclidean metric is a reliable proxy for geodesic distance, allowing for aggressive edge pruning (larger α) to permit long-range “highway” connections without risking semantic shortcuts. Conversely, regions with high LID typically exhibit significant curvature, noise, or singularity. Here, the Euclidean distance often violates the manifold geodesic structure. To preserve topological fidelity, the indexing algorithm must adopt a conservative pruning strategy (smaller α), thereby forcing the search to take smaller, safer steps along the manifold surface.

Let $u \in V$ be a node in the graph, and $\widehat{\text{LID}}(u)$ be its estimated local intrinsic dimensionality. We define the pruning parameter $\alpha(u)$ as:

$$\alpha(u) \triangleq \Phi(\widehat{\text{LID}}(u)). \quad (6)$$

The function $\Phi : \mathbb{R}^+ \rightarrow [\alpha_{\min}, \alpha_{\max}]$ is designed to satisfy the following geometric intuition: in regions with high LID, the graph should enforce a stricter connectivity constraint (smaller α) to avoid short-circuiting the manifold; conversely, in low-LID regions, the constraint can be relaxed (larger α).

To ensure the mapping is robust across datasets with varying complexity scales, we employ Z-score normalization based on the empirical distribution of the LID estimates. We first compute the normalized score $z(u)$:

$$z(u) = \frac{\widehat{\text{LID}}(u) - \mu_{\widehat{\text{LID}}}}{\sigma_{\widehat{\text{LID}}}}, \quad (7)$$

where $\mu_{\widehat{\text{LID}}}$ and $\sigma_{\widehat{\text{LID}}}$ denote the mean and standard deviation of the set of estimated LID values $\{\widehat{\text{LID}}(v) \mid v \in V\}$ computed across the entire graph.

We then formulate Φ using a logistic function to smoothly map the Z-score to the operational range $[\alpha_{\min}, \alpha_{\max}]$:

$$\Phi(\widehat{\text{LID}}(u)) = \alpha_{\min} + \frac{\alpha_{\max} - \alpha_{\min}}{1 + \exp(z(u))}. \quad (8)$$

We set $\alpha_{\min} = 1.0$ and $\alpha_{\max} = 1.5$ following standard practices in graph indexing (Jayaram Subramanya et al., 2019).

This formulation ensures that nodes with average complexity ($z(u) \approx 0$) are assigned $\alpha \approx 1.25$, while nodes with significantly higher complexity ($z(u) > 0$) are penalized with a stricter α approaching the limit of 1.0.

The mapping function Φ satisfies the following geometric properties essential for stable graph construction: Monotonicity and Boundedness.

Proposition 3.4 (Monotonicity). *The mapping function Φ is strictly decreasing with respect to the estimated local intrinsic dimensionality. Formally, given that the standard deviation of the LID estimates $\sigma_{\widehat{\text{LID}}} > 0$ and the pruning range $\alpha_{\max} > \alpha_{\min}$, the derivative satisfies:*

$$\frac{d\Phi}{d\widehat{\text{LID}}(u)} < 0. \quad (9)$$

Proof. Let $L = \widehat{\text{LID}}(u)$ be the independent variable. We define the normalized Z-score z as a function of L :

$$z(L) = \frac{L - \mu_{\widehat{\text{LID}}}}{\sigma_{\widehat{\text{LID}}}}. \quad (10)$$

The mapping function is defined as:

$$\Phi(L) = \alpha_{\min} + \frac{C}{1 + \exp(z(L))}, \quad (11)$$

where $C = \alpha_{\max} - \alpha_{\min}$. Since we strictly set $\alpha_{\max} = 1.5$ and $\alpha_{\min} = 1.0$, it follows that $C > 0$. To determine the sign of the gradient, we apply the chain rule:

$$\frac{d\Phi}{dL} = \frac{d\Phi}{dz} \cdot \frac{dz}{dL}. \quad (12)$$

First, we differentiate the Z-score term with respect to L :

$$\frac{dz}{dL} = \frac{1}{\sigma_{\widehat{\text{LID}}}}. \quad (13)$$

Next, we differentiate the logistic component Φ with respect to z :

$$\frac{d\Phi}{dz} = \frac{d}{dz} (\alpha_{\min} + C(1 + e^z)^{-1}) \quad (14)$$

$$= C \cdot (-1) \cdot (1 + e^z)^{-2} \cdot \frac{d}{dz}(1 + e^z) \quad (15)$$

$$= -C \cdot \frac{e^z}{(1 + e^z)^2}. \quad (16)$$

Combining these terms yields the full derivative:

$$\frac{d\Phi}{dL} = -\frac{C}{\sigma_{\widehat{\text{LID}}}} \cdot \frac{e^z}{(1 + e^z)^2}. \quad (17)$$

We analyze the sign of each component:

- The operational range constant $C > 0$.

- The standard deviation $\sigma_{\widehat{\text{LID}}} > 0$, assuming the dataset exhibits non-zero geometric variance.
- The exponential function $e^z > 0$ for all $z \in \mathbb{R}$.
- The denominator $(1 + e^z)^2 > 0$.

Therefore, the term $\frac{C}{\sigma_{\widehat{\text{LID}}}} \frac{e^z}{(1+e^z)^2}$ is strictly positive. The leading negative sign guarantees that $\frac{d\Phi}{dL} < 0$. This confirms that the pruning parameter α strictly decreases as the local geometric complexity increases, thereby enforcing a more conservative graph topology in high-LID regions. \square

Proposition 3.5 (Boundedness). *The pruning parameter $\alpha(u)$ derived from the mapping function is strictly bounded within the prescribed operational interval. For any node u with a finite LID estimate:*

$$\alpha_{\min} < \alpha(u) < \alpha_{\max}. \quad (18)$$

Proof. Let $S(u)$ denote the logistic component of the mapping function:

$$S(u) = \frac{1}{1 + \exp(z(u))}. \quad (19)$$

For any finite input $\widehat{\text{LID}}(u)$, the Z-score $z(u)$ is finite. The exponential function maps the real line to the positive real line, i.e., $\exp(z(u)) \in (0, \infty)$. Consequently, the denominator lies in the interval $(1, \infty)$. Taking the reciprocal yields the bounds for the logistic component:

$$0 < S(u) < 1. \quad (20)$$

Substituting $S(u)$ back into the definition of Φ :

$$\alpha(u) = \alpha_{\min} + (\alpha_{\max} - \alpha_{\min}) \cdot S(u). \quad (21)$$

Since $(\alpha_{\max} - \alpha_{\min}) > 0$, we can apply the inequality boundaries:

$$\alpha(u) > \alpha_{\min} + (\alpha_{\max} - \alpha_{\min}) \cdot 0 = \alpha_{\min}, \quad (22)$$

$$\alpha(u) < \alpha_{\min} + (\alpha_{\max} - \alpha_{\min}) \cdot 1 = \alpha_{\max}. \quad (23)$$

This proves that the topology is strictly confined. The pruning behavior never exceeds the relaxation upper limit (α_{\max}) and never becomes stricter than the lower limit (α_{\min}), ensuring graph connectivity and preventing degree explosion. \square

3.3. Algorithm

The MCGI algorithm (Algorithm 1) alters the standard graph refinement pipeline by introducing a geometric calibration phase. Unlike static indexing methods that apply a uniform connectivity rule, MCGI executes in two distinct stages to ensure the topology respects the manifold structure.

Algorithm 1 Manifold-Consistent Graph Indexing (MCGI)

Input: Dataset X , Max Degree R , Beam Width L

Output: Optimized Graph G

{Phase 1: Geometric Calibration}

$\mathcal{L} \leftarrow \text{ParallelEstimateLID}(X)$

$\mu \leftarrow \text{Mean}(\mathcal{L})$

$\sigma \leftarrow \text{StdDev}(\mathcal{L})$

for each node $u \in V$ in parallel do

$z_u \leftarrow (\mathcal{L}[u] - \mu)/\sigma$

$\alpha_u \leftarrow \alpha_{\min} + (\alpha_{\max} - \alpha_{\min})/(1 + \exp(z_u))$

end for

{Phase 2: Topology Refinement}

$G \leftarrow \text{RandomGraph}(X, R)$

for $\text{iter} \leftarrow 1$ to MaxIter do

for each node $u \in G$ in parallel do

$\mathcal{C} \leftarrow \text{GreedySearch}(u, G, L)$

$\mathcal{N}_{\text{new}} \leftarrow \emptyset$

for $v \in \text{SortByDistance}(\mathcal{C} \cup \mathcal{N}(u))$ do

$\text{pruned} \leftarrow \text{False}$

for $n \in \mathcal{N}_{\text{new}}$ do

if $\alpha_u \cdot d(n, v) \leq d(u, v)$ then

$\text{pruned} \leftarrow \text{True}; \text{break}$

end if

end for

if not pruned and $|\mathcal{N}_{\text{new}}| < R$ then

$\mathcal{N}_{\text{new}}.\text{add}(v)$

end if

end for

$\mathcal{N}(u) \leftarrow \mathcal{N}_{\text{new}}$

end for

end for

Phase 1: Geometric Calibration. Before modifying the graph topology, the system first performs a global analysis of the dataset geometry. We estimate the LID for every point and aggregate the population statistics (μ, σ) defined in Section 3.2. This phase “freezes” the geometric profile of the dataset. By pre-computing these statistics, we decouple the complexity estimation from the graph update loop, ensuring that the mapping function Φ remains stable and computationally efficient during the intensive edge-selection process.

Phase 2: Manifold-Consistent Refinement. The index construction follows an iterative refinement strategy. Let $\mathcal{N}(u)$ denote the set of neighbors for node u in the graph G . In each iteration, the algorithm dynamically updates $\mathcal{N}(u)$ by:

1. Queries the pre-computed geometric profile to determine the node-specific constraint $\alpha(u)$.
2. Explores the graph to identify a candidate pool \mathcal{C} .

3. Filters connections using the dynamic occlusion criterion.

3.4. Theoretical Analysis

We analyze the proposed method regarding its computational overhead and topological guarantees.

Computational Complexity. The MCGI construction adds a pre-processing step (Phase 1) but maintains the same asymptotic complexity as standard graph construction algorithms (Phase 2).

- *Calibration Overhead:* The LID estimation requires a k -nearest neighbor search for a subset of points. Assuming a sample size of N and a fixed reference set size, this operation is bounded by $O(N \log N)$. The subsequent Z-score calculation and α mapping are linear scan operations, $O(N)$.
- *Construction Cost:* The core refinement loop (Algorithm 1, Phase 2) has a complexity of $O(T \cdot N \cdot R \cdot \log L)$, where T is the number of iterations and R is the degree bound. The node-specific parameter $\alpha(u)$ is accessed via a constant-time lookup $O(1)$ during the edge pruning check.

Since the calibration is a non-iterative, one-time operation, it does not alter the asymptotic complexity class of the indexing pipeline. The total time complexity remains dominated by the iterative distance computations required for edge selection, ensuring that the method scales linearly with N similar to baseline iterative algorithms.

Connectivity Preservation. A primary concern with dynamic pruning is whether strict constraints fragment the graph into disconnected components. To address this, we observe that the strictest constraint applied by our algorithm (where $\alpha(u) \rightarrow 1.0$) corresponds to the definition of the Relative Neighborhood Graph (RNG). It has been theoretically established that the RNG is a supergraph of the Euclidean Minimum Spanning Tree (EMST) for any set of points in general position (Toussaint, 1980). Since the EMST is connected, the RNG is necessarily connected. In our method, the pruning parameter satisfies $\alpha(u) \geq 1.0$ for all u . Let E_{MCGI} be the edge set produced by our algorithm and E_{RNG} be the edge set of the exact RNG. The condition implies:

$$E_{RNG} \subseteq E_{MCGI}. \quad (24)$$

Thus, provided that the greedy search sufficiently explores the candidate space, the resulting topology retains the connectivity guarantees of the RNG, ensuring that every node remains reachable from the entry point.

4. Evaluation

5. Conclusion

References

- Amsaleg, L., Chelly, O., Furon, T., Girard, S., Houle, M. E., Kawarabayashi, K.-i., and Nett, M. Estimating local intrinsic dimensionality. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’15, pp. 29–38, New York, NY, USA, 2015. Association for Computing Machinery. ISBN 9781450336642. doi: 10.1145/2783258.2783405. URL <https://doi.org/10.1145/2783258.2783405>.
- Fu, C., Cai, C., Zhou, D., Liu, W., and Wang, C. Fast approximate nearest neighbor search with the navigating spreading-out graph. *Proceedings of the VLDB Endowment*, 12(5):461–474, 2019.
- Houle, M. E. Local intrinsic dimensionality I: an extreme-value-theoretic foundation for similarity applications. In Beecks, C., Borutta, F., Kröger, P., and Seidl, T. (eds.), *Similarity Search and Applications - 10th International Conference, SISAP 2017, Munich, Germany, October 4-6, 2017, Proceedings*, volume 10609 of *Lecture Notes in Computer Science*, pp. 64–79. Springer, 2017. doi: 10.1007/978-3-319-68474-1\5. URL https://doi.org/10.1007/978-3-319-68474-1_5.
- Jayaram Subramanya, S., Devvrit, F., Simhadri, H. V., Krishnawamy, R., and Kadekodi, R. Diskann: Fast accurate billion-point nearest neighbor search on a single node. *Advances in Neural Information Processing Systems*, 32, 2019.
- Levina, E. and Bickel, P. Maximum likelihood estimation of intrinsic dimension. In Saul, L., Weiss, Y., and Bottou, L. (eds.), *Advances in Neural Information Processing Systems*, volume 17. MIT Press, 2004. URL https://proceedings.neurips.cc/paper_files/paper/2004/file/74934548253bcab8490ebd74afed7031-Paper.pdf.
- Malkov, Y. A. and Yashunin, D. A. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. In *IEEE transactions on pattern analysis and machine intelligence*, volume 42, pp. 824–836. IEEE, 2018.
- Toussaint, G. T. The relative neighbourhood graph of a finite planar set. *Pattern Recognition*, 12(4):261–268, 1980. ISSN 0031-3203. doi: [https://doi.org/10.1016/0031-3203\(80\)90066-7](https://doi.org/10.1016/0031-3203(80)90066-7). URL <https://www.sciencedirect.com/science/article/pii/0031320380900667>.