Continuous machine learning on Euclidean graphs with unordered vertices

Abstract

Molecular graphs can change their chemical properties under non-rigid deformations in Euclidean space. Hence it is vitally important to distinguish rigid classes of molecular graphs under compositions of translations and rotations. Also, robust outputs of machine learning on molecular graphs embedded in Euclidean space should continuously change under perturbations, motivated by atomic vibrations and experimental noise. We developed a complete invariant that can be inverted back to an embedded graph, uniquely under rigid motion, and has a Lipschitz continuous distance satisfying all metric axioms. For a fixed dimension, the invariant and metric can be computed in polynomial time of the number m of unordered vertices and hence avoiding exponentially many permutations. The new invariants distinguish all chemically different graphs in the world's largest databases of 3D molecules in a few hours on a modest desktop.

1. Motivations for complete and continuous invariant inputs in application-driven ML

This paper formalizes necessary conditions for ML on real data with ambiguous representations and develops complete and Lipschitz continuous invariants satisfying these conditions on any Euclidean graphs and justifying a rigorous concept of a molecular structure. Many real structures from star constellations to molecules are represented by graphs embedded in a Euclidean space (Bonchev, 1991). A Euclidean graph $G \subset \mathbb{R}^n$ is a finite set of m unordered (unlabeled) vertices located at distinct points of \mathbb{R}^n and connected by straight-line edges. Forgetting all edges of $G \subset \mathbb{R}^n$ gives us the vertex set $V(G) \subset \mathbb{R}^n$ of m unordered points. A Euclidean graph can be disconnected and can have vertices v of any *degree* that is the number of edges whose endpoint is v. Loops and multiple edges (with the same endpoints) do not appear in Euclidean graphs because all edges are straight line segments and can also intersect in theory.

Graphs can be considered under any *equivalence* relation that should satisfy the axioms: 1) *reflexivity*: $G \sim G$, 2) *symmetry*: if $G \sim F$ then $F \sim G$, 3) *transitivity*: if $G \sim F$ and $F \sim H$ then $G \sim H$. In chemistry, the simplest equivalence is by chemical composition, which is insufficient in practice, e.g. *stereoisomers* in Fig. 1 (right) have the same chemical compositions and non-equivalent rigid shapes with different chemical properties (Rieder et al., 2023).



Figure 1. Top: graphs $T_1, T_2, T_3, T_4 \subset \mathbb{R}^3$ on the same vertices with solid edges are not isomorphic to each other. Bottom: stereoisomers are isomorphic combinatorially, not geometrically.

For molecules, the strongest equivalence (distinguishing as many graphs as practically possible) is a *geometric isomorphism* $G \cong F$, i.e. an orientation-preserving transformation of \mathbb{R}^n that bijectively maps the vertices and edges: $G \to F$. Geometric isomorphisms are also called *rigid motions* (compositions of translations and rotations), which form the special Euclidean group SE(n). The slightly weaker equivalence (not distinguishing mirror images) is an *isometry*, which is any distance-preserving transformation including reflections. Any geometrically isomorphic molecules have the same chemical properties. If a flexible molecule changes its rigid shape, its functional properties can change, so it is important to distinguish rigid shapes (Wilson et al., 1991).

To reliably distinguish at least some Euclidean graphs $G \subset \mathbb{R}^n$, we need an *invariant* I defined as a numerical descriptor preserved by any rigid motion in \mathbb{R}^n . Alternatively, if $I(G) \neq I(F)$, then $G \ncong F$, so any invariant has *no false negatives* that are pairs of different representatives of *rigidly equivalent graphs* (denoted by $G \cong F$) having equal values of a (non-invariant) descriptor. The number of vertices (or edges) of G is an integer-valued weak invariant that cannot separate any graphs in Fig. 1. The strongest invariant I separating all non-equivalent graphs is called *complete* meaning that if I(G) = I(F) then $G \cong F$. Alternatively, a

[.] Correspondence to: Anonymous Author anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

complete invariant *I* has *no false positives* that are pairs of non-equivalent graphs $G \not\cong F$ with I(G) = I(F).

057 Since all real data (such as inter-point distances) are noisy, 058 a more practically important answer is not binary ('same 059 or different') but should be continuously quantified by a 060 distance metric between isometry classes. The atomic 061 vibrations (Feynman, 1971) imply that rigid classes of 062 molecules graphs on m unordered atoms form a contin-063 uous Graph Isometry Space $GIS(\mathbb{R}^3; m)$. Only for tri-064 angular graphs with m = 3, their space was previously 065 known due to the side-side theorem saying that any 066 triangles are isometric if and only if they have the same 067 triple of sides (inter-point distances) a, b, c considered up 068 to 6 permutations. Hence the space of triangular graphs is 069 $\{0 < a \le b \le c \le a+b\} \subset \mathbb{R}^3$, where $c \le a+b$ guarantees 070 that distances a, b, c are realizable by a real triangle. 071

072**Problem 1.1** (complete invariant of Euclidean graphs with
a polynomial-time continuous metric). Find a function I :
 $\{Euclidean graphs with of unordered vertices in \mathbb{R}^n\} \rightarrow a$
space X with a distance d satisfying the conditions below:
076

(a) completeness of the invariant: any graphs G, F are related by rigid motion in \mathbb{R}^n if and only if I(G) = I(F);

1079 (b) Lipschitz continuity: there is a constant λ and a metric 1080 d satisfying the axioms 1) $d(\alpha, \beta) = 0$ if and only if $\alpha = \beta$, 1081 2) $d(\alpha, \beta) = d(\beta, \alpha)$, 3) $d(\alpha, \beta) + d(\beta, \gamma) \ge d(\alpha, \gamma)$ for 1082 all $\alpha, \beta, \gamma \in X$, such that if F is obtained by perturbing 1083 every vertex of G up to $\varepsilon > 0$, then $d(I(G), I(F)) \le \lambda \varepsilon$;

(c) invertibility: any Euclidean graph G can be reconstructed (uniquely up to rigid motion in \mathbb{R}^n) from I(G);

 $\begin{array}{ll} \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ I, d, and a reconstruction of <math>G \subset \mathbb{R}^n$ from I(G) can be obtained in polynomial time of the number of vertices. \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant \\ \textbf{(d) computability: } for a fixed dimension n, the invariant

091 Condition 1.1(a) means that a complete invariant I has the 092 strongest expressivity (Zhang et al., 2024) by uniquely iden-093 tifying any Euclidean graph under geometric isomorphism. 094 To be useful for noisy inputs, a complete invariant should 095 continuously change under perturbations in a suitable metric. 096 The axioms in 1.1(b) are the foundations of metric geometry 097 (Melter & Tomescu, 1984) and accepted in chemistry (Wein-098 hold, 1975). If the triangle axiom fails with any additive 099 error, the classical k-means and DBSCAN clustering are 100 open to adversarial attacks in (Rass et al., 2024). If the first axiom is ignored, $d \equiv 0$ satisfies all other axioms. The first axiom implies the completeness of I in 1.1(a) but the continuity is much stronger. Indeed, for any complete invariant 104 I, one can define the discrete metric d(I(G), I(F)) = 1 for 105 $G \ncong F$, which unhelpfully treats all non-equivalent graphs 106 (even near-duplicates) as equally distant. The Lipschitz continuity in 1.1(b) is necessary for smoothness, which is 108 implicitly assumed by any gradient-based optimization. 109

Condition 1.1(c) requires I to be not only complete and continuous but also efficient to explicitly reconstruct G, even better than a DNA code that is does not explain how to grow a living organism. Computability 1.1(d) prevents brute-force attempts, e.g. defining I(G) as the infinite set of images of G under all rigid motions or taking m! distance matrices over all permutations of m unordered vertices.

The main contribution is the new invariant Nested Centered Distribution, which solves Problem 1.1, including the new Lipschitz continuity, for all Euclidean graphs in \mathbb{R}^n .

2. Past work on distances for Euclidean graphs

Ordered clouds. The vertex set V(G) of a Euclidean graph $G \subset \mathbb{R}^n$ is called a *point cloud* C. If all points p_1, \ldots, p_m of C are ordered (not under the action of all m! permutations), a complete invariant of C under isometry (compositions of translations, rotations, reflections) is the classical $m \times m$ matrix (Li et al., 2023) of pairwise distances $|p_i - p_j|$ due to Theorem 9 in (Grinberg & Olver, 2019) or, after shifting the center of mass to the origin, the Gram matrix of scalar products $p_i \dot{p}_j$ by Theorem 1 in (Dekster & Wilker, 1987). This multidimensional scaling (Schoenberg, 1935) can also provide an embedding $C \subset \mathbb{R}^k$ preserving all distances of C for a dimension $k \leq m$. This embedding $C \subset \mathbb{R}^k$ uses eigenvectors whose ambiguity up to signs gives an exponential time that can be close to $O(2^m)$, not polynomial in the number m of ordered points as in 1.1(d).

Unordered clouds. Computational geometry developed many algorithms for detecting geometric isomorphism (or isometry, also called congruence) between point sets without edges (Huttenlocher et al., 1993; Chew & Kedem, 1992; Chew et al., 1999; Goodrich et al., 1999). For a set $A \subset \mathbb{Q}^n$ of m points, Theorem 3 in (Arvind & Rattan, 2016) computed in time $n^{O(n)} poly(mM)$ a canonizing function f(A), which can be considered a complete isometry invariant of A, where M upper bounds the binary encodings of the rational coordinates in the input. For point clouds under rigid motion (also distinguishing mirror images), Theorem 4.7 in (Widdowson & Kurlin, 2023) described a metric computable in time $O(n(m^{n-1}/n!)^3 \log m)$. (Hordan et al., 2024; Delle Rose et al., 2024; Nigam et al., 2024; Amir et al., 2024; Maennel et al., 2024) also achieved the completeness for point clouds but without a Lipschitz continuous metric as in 1.1(b). Energy potentials written as infinite series of spherical harmonics, are often considered complete representations of atomic environments, which holds in the limit but not for a finite size(Pozdnyakov et al., 2020). For a fixed set of m vertices in general position, one can choose any of m(m-1)/2 edges and produce $2^{m(m-1)/2}$ non-isometric graphs. Problem 1.1 for arbitrary graphs is computationally much harder than for point clouds due to exponentially many different graphs on the same vertex set. 110 The graph isomorphism problem (Grohe & Schweitzer, 111 2020) for abstract (non-Euclidean) graphs is another ver-112 sion of Problem 1.1 without continuous metrics. The lat-113 est advances (Babai, 2019; Helfgott et al., 2017) achieved 114 only quasipolynomial time. While many partial cases were 115 solved, e.g. for planar graphs (embedded in \mathbb{R}^2 without inter-116 secting edges), see (Kiefer et al., 2019), the k-dimensional 117 Weisfeiler-Leman test (Leman & Weisfeiler, 1968) fails for 118 3-regular graphs of size O(k). The key limitation of WL 119 tests is their local nature when invariants are gradually ex-120 panded from a vertex or a k-tuple. Then covering a graph on 121 m vertices needs O(m) expansions leading to exponential 122 sizes in m. Section 3.9 in (Dym & Gortler, 2024) discussed that a complete invariant (under all permutations of m ver-124 tices) that has a polynomial time in the dimension n would 125 also solve the graph isomorphism problem in polynomial 126 time. Condition 1.1(d) is easier for a fixed dimension n, 127 e.g. n = 2, 3 are practical cases. The number m of vertices 128 can be dozens or hundreds, e.g. for molecular graphs in 129 \mathbb{R}^3 , where vertices are centers of atoms and edges are inter-130 atomic bonds that keep atoms together in a stable molecule.

131 Geometric Deep Learning in (Bronstein et al., 2021) pio-132 neered an axiomatic approach to geometric classifications 133 beyond Euclidean space \mathbb{R}^n in (Bronstein et al., 2017). 134 Some neural networks were proved to be universal (Maron 135 et al., 2019; Zhou, 2020; Abbe & Sandon, 2020) in the sense 136 of approximating any continuous function on given data 137 with sufficiently many layers. This universality property has 138 been strengthened in Problem 1.1 to the full completeness 139 of an explicit invariant that should be computable in poly-140 nomial time and invertible to an original graph up to rigid 141 motion. The key challenge was to compute an exact (not 142 approximate) metric that is also Lipschitz continuous. 143

144 Equivariants (Kondor & Trivedi, 2018; Cohen et al., 2019; 145 Fuchs et al., 2020; Deng et al., 2021) are defined as de-146 scriptors E satisfying $E(f(G)) = T_f(E(G))$ for any rigid motion f and all graphs $G \subset \mathbb{R}^n$, where T_f can be any 147 148 map, not only the identity as for invariants. Any linear com-149 bination of points, e.g. the center of mass, is equivariant 150 but cannot distinguish graphs under translation. Equivari-151 ants (Gao et al., 2020; Qi & Luo, 2020; Tu et al., 2022; 152 Batzner et al., 2022) help predict forces acting on atoms to 153 move them to a more optimal configuration. These time-154 dependent graphs G_t can be studied directly by invariant 155 values $I(G_t)$ without computing intermediate atomic forces.

Many neural networks optimize millions of parameters, e.g.
see Table 4 (Goyal et al., 2021), to achieve great accuracies (Dong et al., 2018; Akhtar & Mian, 2018; Laidlaw & Feizi, 2019; Guo et al., 2019; Colbrook et al., 2022) but require re-training on any new data. All known descriptors of molecular graphs (Duvenaud et al., 2015; Choo et al., 2023) have no proofs of all conditions 1.1(a,b,c,d).

164

Gromov-Wasserstein metrics (Mémoli, 2011) are defined for any metric-measure spaces (Brécheteau, 2019) by minimizing over infinitely many correspondences between points, but cannot be approximated with a factor less than 3 in polynomial time unless P=NP by Corollary 3.8 in (Schmiedl, 2017) and Theorem 3.3 in (Agarwal et al., 2018), see fast algorithms for important cases in (Mémoli et al., 2021; Lim et al., 2023; Majhi et al., 2024). (Nikolentzos et al., 2017; Majhi & Wenk, 2022; Buchin et al., 2023) made significant advances in the related problems of matching and finding distances between fixed Euclidean graphs without considering isometry. Computing a metric between rigid classes is only a small part of Problem 1.1. Indeed, to efficiently navigate on Earth, in addition to distances between cities, we need a map of the planet and hence an invertible continuous invariant I similar to geographic coordinates.

3. Graph invariants: from fastest to complete

Let |p-q| denote the Euclidean distance between any points $p, q \in \mathbb{R}^n$. We always translate any graph $G \subset \mathbb{R}^n$ so that the *center of mass* $O(G) = \frac{1}{m} \sum_{p \in V(G)} p$ of the *vertex set* V(G) is at the origin $0 \in \mathbb{R}^n$. Then Problem 1.1 reduces to the SO(n)-invariance under orthogonal transformations.

Definition 3.1 (signed distance d(p,q) and invariants SRD, SPD, PDD). Let $G \subset \mathbb{R}^n$ be any Euclidean graph on *m* arbitrarily ordered vertices $p_1 \ldots, p_m$. If any $p_i, p_j \in V(G)$ are connected by an edge of *G*, define the signed distance as d(p,q) = |p-q|, else set d(p,q) = -|p-q|.

(a) The vector SRD(G) of sorted radial distances consists of m distances |p| for all $p \in V(G)$ in decreasing order.

(b) The vector SPD(G) of sorted pairwise distances consists of all distances $d(p_i, p_j)$ in decreasing order.

(c) Let D(G) be the $m \times (m-1)$ -matrix whose the *i*-th row consists of $d(p_i, p_j)$, $j \in \{1, ..., m\} \setminus \{i\}$, in increasing order. The Pointwise Distance Distribution PDD(G) consists of these unordered rows with equal weights 1/m.

If any k > 1 rows of D(G) are equal, they can be collapsed in PDD(G) to a single row with the *weight* k/m. The PDD was defined for clouds as a local distribution of distances in Definition 5.5 of (Mémoli, 2011) and for periodic sets in (Widdowson & Kurlin, 2022) but not for Euclidean graphs.

Table 1. Acronyms of all main invariants and metrics in the paper.

SRD	SORTED RADIAL VECTOR	Def 3.1
SPD	SORTED DISTANCE VECTOR	Def 3.1
PDD	POINTWISE DISTANCE DISTRIBUTION	Def 3.1
CR	CENTERED REPRESENTATION	Def 3.3
NCD	NESTED CENTERED DISTRIBUTION	Def 3.5
NBM	NESTED BOTTLENECK METRIC	Def 4.5

- 165 The PDD(G) includes every signed distance twice, once as 166 d(p,q) in the row of a vertex p, and as d(q,p) in the row of
- 167 a vertex q. Hence SPD(G) can be obtained from PDD(G)
- by (1) combining all distances into one vector, (2) sorting
- 169 them in decreasing order, and (3) keeping only one copy of
- 170 every two repeated distances. Example 3.2 shows that the

171 invariant PDD(G) is strictly stronger than SPD(G).

172 **Example 3.2** (invariants SRD, SPD, PDD for tetrahedral 173 graphs in Fig. 1). (a) Since the vertex sets of $T_i \,\subset\, \mathbb{R}^3$ 174 are regular tetrahedra with all pairwise distances 1, these 175 graphs have identical SRD(T_i) of 4 equal circumradii of 176 the same vertex set $V(T_i)$ independent of i = 1, ..., 4.

The first graph T_1 has two edges contributing +1 and four non-edges (dashed lines) contributing -1 to the Sorted Distance Vector $SPD(T_1) = (+1, +1, -1, -1, -1, -1)$. The graph T_2 also has two edges, so $SPD(T_2) = SPD(T_1)$ doesn't distinguish $T_1 \ncong T_2$ up to rigid motion. Similarly, the graphs $T_3 \ncong T_4$ are not distinguished by the invariants $SPD(T_3) = (+1, +1, +1, -1, -1, -1) = SPD(T_4)$.

185 (b) In T_1 , every vertex has exactly one edge and two 186 non-edges (dashed lines), hence its signed distances are 187 +1, -1, -1. The matrix $PDD(T_1) = (100\% \mid -1, -1, +1)$ 188 consists of a single row, where the weight 100% indicates 189 that all vertices of T_1 have the same row in PDD. The 190 graph T_2 has one vertex (25%) with no edges, two vertices 191 (50%) with one edge, and one vertex (25%) with two edges, $(25\% \mid -1 \quad -1 \quad -1)$ 193

¹⁹³
¹⁹⁴ so
$$PDD(T_2) = \begin{pmatrix} 50\% & -1 & -1 & +1 \\ 25\% & -1 & +1 & +1 \end{pmatrix} \neq PDD(T_1),$$

195 so PDD distinguishes the rigidly non-equivalent graphs 196 $T_1 \ncong T_2$ with $SPD(T_1) = SPD(T_2)$. The graph T_3 has one 197 vertex (25%) with no edges and three vertices (75%) with 198 two edges, so $PDD(T_3) = \begin{pmatrix} 25\% & -1 & -1 & -1 \\ 75\% & -1 & +1 & +1 \end{pmatrix}$. 199 200 The graph T_4 has two vertices (50%) with one edge and 201 two vertices (50%) with two edges. Then $PDD(T_4) =$ 202 $\begin{pmatrix} 50\% & -1 & -1 & +1 \\ 50\% & -1 & +1 & +1 \end{pmatrix}$, so PDD distinguishes the 203 204 graphs $T_3 \ncong T_4$ with equal $SPD(T_3) = SPD(T_4)$.

For a graph G with m unordered vertices, PDD(G) has m - 1 columns. The reduced version PDD(G; k) includes only the first k columns for $1 \le k < m - 1$. Though PDDs have unordered rows, they can be continuously compared by Earth Mover's Distance (Rubner et al., 2000).

211 Fig. S4 in (Pozdnyakov et al., 2020) described infinitely 212 many non-isometric pairs of clouds $C, C' \subset \mathbb{R}^3$ with 213 PDD(C) = PDD(C'). These counter-examples inspired 214 the stronger invariants for graphs below. For simplicity, we 215 will introduce all invariants and metrics in dimension n = 2. 216 All higher dimensions n > 2 are covered in appendices. 217 While PDD(G) includes signed distances to a single (arbi-218 trary) vertex $p_i \in V(G)$, a stronger invariant below include 219

triples of signed distances to three base points, one of which is the center of mass of V(G) because any point in \mathbb{R}^2 is uniquely determined by its distances to three fixed points.

Definition 3.3 (Centered Representation CR(G; A) of a graph with $A \subset V(G)$). Let $G \subset \mathbb{R}^2$ be a graph on m unordered points with the center of mass $p_0 = O(G) = 0$.

(a) For any vertex $p_1 \in V(G)$, the matrix $R(G; p_1)$ has m-1 unordered columns, one for each vertex $q \in V(G) \setminus \{p_1\}$, consisting of the signed distances $d(q, p_0)$ and $d(q, p_1)$. Here $p_0 = 0$ is not considered as a vertex of G, so $d(q, p_0) = -|q|$. The Centered Representation $CR(G; p_1)$ is the pair $[d(p_0, p_1), R(G; p_1)]$, where $d(p_0, p_1) = -|p_1|$.

(b) Fix a base pair A of ordered vertices $p_1, p_2 \in V(G)$. Let sign(A) be the sign of the 2×2 determinant on the vectors p_1, p_2 . Let D(A) be the matrix of signed distances between p_0, p_1, p_2 . The matrix R(G; A) has m - 2 unordered columns, one for each vertex $q \in V(G) \setminus A$, consisting of signed distances $d(q, p_0), d(q, p_1), d(q, p_2)$. The Centered Representation CR(G; A) is the triple [sign(A), D(A), R(G; A)].

After fixing $p_0 = 0$, the matrix D(A) and $\operatorname{sign}(A)$ help reconstruct base vertices $p_1, p_2 \in \mathbb{R}^2$, uniquely under rotation around 0. Any other $q \in V(G) \setminus A$ is fixed relative to p_0, p_1, p_2 by its column in R(G; A). A positive sign of $d(p_i, p_j)$ indicates an edge between vertices p_i, p_j . This argument will later be formalized in Theorem 4.6(b).

Example 3.4 (CRs for 2-vertex bases in \mathbb{R}^2). Let $G \subset \mathbb{R}^2$ be the triangular cycle on $p_1 = (2,0)$, $p_2 = (-1,1)$, $p_3 = (-1,-1)$, so O(G) = 0 and all signed distances are positive, see Fig. 2 (top left). For $A = (p_1, p_2)$, $\operatorname{sign}(A) = \operatorname{sign} \begin{vmatrix} 2 & -1 \\ 0 & 1 \end{vmatrix} = 1$. The distance matrix on $\begin{pmatrix} 0 & -2 & -\sqrt{2} \end{pmatrix}$

$$0, p_1, p_2 \text{ is } D(p_1, p_2) = \begin{pmatrix} 0 & -2 & -\sqrt{2} \\ -2 & 0 & \sqrt{10} \\ -\sqrt{2} & \sqrt{10} & 0 \end{pmatrix}. \text{ Then}$$

$$R(G; p_1, p_2) = \begin{pmatrix} -|p_3| \\ |p_3 - p_1| \\ |p_3 - p_2| \end{pmatrix} = \begin{pmatrix} -\sqrt{2} \\ \sqrt{10} \\ 2 \end{pmatrix}.$$
 Then

$$CR(G; p_1, p_2) = [+1, D(p_1, p_2), R(G; p_1, p_2)]. Replac-ing p_2 with p_3, we find sign(p_1, p_3) = sign \begin{vmatrix} 2 & -1 \\ 0 & -1 \end{vmatrix} = \begin{pmatrix} 0 & -2 & -\sqrt{2} \\ 0 & -2 & -\sqrt{2} \end{pmatrix}$$

$$-1, \quad D(p_1, p_3) = \begin{pmatrix} -2 & 0 & \sqrt{10} \\ -\sqrt{2} & \sqrt{10} & 0 \end{pmatrix}, \quad and$$

$$R(G; p_1, p_3) = \begin{pmatrix} -|p_2| \\ |p_2 - p_1| \\ |p_2 - p_3| \end{pmatrix} = \begin{pmatrix} -\sqrt{2} \\ \sqrt{10} \\ 2 \end{pmatrix}. \text{ The final}$$

triple is $CR(G; p_1, p_3) = [-1, D(p_1, p_3), R(G; p_1, p_3)].$

Though a Centered Representation $CR(G; p_1, p_2)$ will suffice to reconstruct $G \subset \mathbb{R}^2$ uniquely under rigid motion

in Theorem 4.6(b), $CR(G; p_1, p_2)$ for all vertices $p_1, p_2 \in V(G)$ should be considered in a joint unordered collection below to guarantee the independence of points p_1, p_2 .

223 224 **Definition 3.5** (Nested Centered Distribution NCD(G; h)). 225 Let $G \subset \mathbb{R}^2$ be any Euclidean graph with m unordered 226 vertices and the center of mass at the origin $0 \in \mathbb{R}^n$.

227 (a) The Nested Centered Distribution NCD(G; 1) of order 1 228 is the unordered set of Centered Representations $CR(G; p_1)$ 229 from Definition 3.3 for all vertices $p_1 \in V(G)$.

(b) For any vertex $p_1 \in V(G)$, the Centered Distribution $CD_1(G; p_1)$ is the unordered set of $CR(G; p_1, p_2)$ for all $p_2 \in V(G) \setminus \{p_1\}$. The Nested Centered Distribution NCD(G; 2) of order 2 is the unordered set of $CD_1(G; p_1)$ for all vertices $p_1 \in V(G)$, see Fig. 2 (top). The mirror image $\overline{NCD}(G; 2)$ is obtained from NCD(G; 2) by reversing $sign(p_1, p_2)$ of 2×2 determinants in all $CR(G; p_1, p_2)$.

238

247

269

270

271

272

273

274

The nested structure of NCD(G; 2) helps identify edges 239 between all vertices from G. After any vertex $q \in V(G) \setminus$ 240 $\{p_1, p_2\}$ is uniquely located by using one $CR(G; p_1, p_2)$, 241 we can use unsigned distances to associate any such q with 242 its unique $CR(G; p_1, q)$ in the collection $\{CR(G; p_1, p_i)\}$ 243 for all $p_i \in V(G) \setminus \{p_1\}$. The resulting $CR(G; p_1, q)$ con-244 tains signed distances and hence detects edges from q to all 245 other vertices, see details in the proof of Theorem 4.6(b). 246



Figure 2. **Top**: building the Nested Centered Distribution NCD in Definition 3.5 from Centered Representations in Definition 3.3 with metrics in section 4. **Bottom**: hierarchy of graph invariants.

SRD(G) can be considered NCD(G; 0) of order 0, containing signed distances from the center of mass p_0 to all vertices of G, additionally written in increasing order.

4. Continuous metrics on graph invariants

When points $0 \cup A = (p_0, p_1, p_2) \subset \mathbb{R}^2$ pass through a degenerate configuration in a straight line, i.e. p_1, p_2 become collinear, sign(A) discontinuously changes. To guarantee the Lipschitz continuity, we multiply such a sign by the strength σ below, which smooths the sign change, while the area of the triangle on p_0, p_1, p_2 is not Lipschitz continuous.

Definition 4.1 (strength $\sigma(C)$). Any triple $C = \{p_0, p_1, p_2\} \subset \mathbb{R}^2$ defines a triangle with inter-point distances a, b, c, and half-perimeter $p = \frac{1}{2}(a + b + c)$. The strength is $\sigma(C) = \frac{(p-a)(p-b)(p-c)}{p^2}$.

Lemma 4.2 (Theorem 4.4 in (Widdowson & Kurlin, 2023)). Let B be obtained from a set $C \subset \mathbb{R}^2$ of 3 points by perturbing every point within its ε -neighborhood. Then $|\sigma(B) - \sigma(C)| \leq 2\varepsilon\lambda_2$ for $\lambda_2 = 2\sqrt{3}$.

The strength $\sigma(A)$ will be normalized by λ_2 below to guarantee the final Lipschitz constant 2 for a metric in Theorem 4.6(c). For any $k \times k$ matrices M, N of real numbers, the metric L_{∞} is $\max_{i,j=1,...,k} |M_{ij} - N_{ij}|$. The *bottleneck* distance between any clouds A, B of (the same number of) m unordered points in a metric space with a distance d is $W_{\infty}(A, B) = \min_{\text{bijections } g:A \to B} \max_{p \in A} d(g(p), p).$

Definition 4.3 (max metric M_{∞} on CRs). Let Euclidean graphs $G, F \subset \mathbb{R}^n$ have m unordered vertices.

(a) For order h = 1, take any base vertices $p \in V(G)$ and $q \in V(F)$. Define the max metric $M_{\infty}(\operatorname{CR}(G;p), \operatorname{CR}(F;q))$ as the maximum of ||p| - |q|| and the bottleneck distance W_{∞} between the fixed clouds of unordered points $\{(-|p'|, d(p', p)) \mid p' \in V(G) - \{p\}\}$ and $\{(-|q'|, d(q', q)) \mid q' \in V(F) - \{q\}\}$ in \mathbb{R}^2 .

(b) For order h = 2, take any base sequences $A \subset V(G)$ and $B \subset V(F)$ of two vertices. Consider the m - 2 columns of R(G; A) from Definition 3.3 as a cloud of m - 2 unordered points in \mathbb{R}^2 , also for R(F; B). The max metric $M_{\infty}(\operatorname{CR}(G; A), \operatorname{CR}(F; B))$ is the maximum of $\frac{2}{\lambda_2} |\operatorname{sign}(A)\sigma(0 \cup A) - \operatorname{sign}(B)\sigma(0 \cup B)|$, $L_{\infty}(D(A), D(B))$, and $W_{\infty}(R(G; A), R(F; B))$.

The maximum of several distances in Definition 4.3 is needed to guarantee the first metric axiom, i.e. $M_{\infty}(\operatorname{CR}(G; A), \operatorname{CR}(F; B)) = 0$ should imply that $0 \cup A$ should be exactly matched by rotation with $0 \cup B$ and then $\operatorname{CR}(G; A) = \operatorname{CR}(F; B)$ up to a permutation of columns will imply that G coincides with F, see Lemma D.7.

To get a metric on Nested Centered Distributions, we will use the distance on bipartite graphs whose edge weights are the max metrics M_{∞} on Centered Representations.

Definition 4.4 (Bottleneck Matching Distance $BMD(\Gamma)$).

275 Let Γ be a complete bipartite graph with m white vertices 276 and m black vertices so that every white vertex is connected 277 to every black vertex by a single edge e of a weight $w(e) \ge$ 278 0. A vertex matching of the graph Γ is a collection E of 279 m disjoint edges with 2m distinct vertices. The weight 280 $W(E) = \max_{e \in E} w(e)$ is the largest weight of an edge in E.

²⁸¹ 282 The Bottleneck Matching Distance $BMD(\Gamma) = \min_{E} W(E)$

is the minimum weight of a vertex matching E of Γ .

Since a graph Γ is complete bipartite, any edge from a vertex matching E in Γ joins a white vertex with a black vertex. Then BMD(Γ) is minimized for all bijections E between all white vertices and all black vertices of Γ .

289 **Definition 4.5** (Nested Bottleneck Metric NBM on NCDs). 290 Let $G, F \subset \mathbb{R}^2$ be any graphs on m unordered vertices.

291 (a) For order h = 1, the Nested Bottleneck Metric NBM(NCD(G; 1), NCD(F; 1)) is the max metric M_∞(CR(G; p), CR(F; $\beta(p)$)) minimized for all bijections $\beta: V(G) \rightarrow V(F)$ between vertices of G and F.

(b) For order h = 2, any base vertices $p_1 \in V(G)$ and $q_1 \in$ 296 V(F), let the complete bipartite graph $\Gamma(G; p_1; F; q_1)$ have 297 m-1 white vertices and m-1 black vertices representing 298 $\operatorname{CR}(G; p_1, p_2)$ and $\operatorname{CR}(F; q_1, q_2)$ for all $p_2 \in V(G) - \{p_1\}$ 299 and $q_2 \in V(F) - \{q_1\}$, respectively. Set the weight w(e) of 300 an edge e joining the vertices represented by $CR(G; p_1, p_2)$ 301 and $CR(F; q_1, q_2)$ as the max metric M_{∞} between these 302 distributions, see Definition 4.3. Then Definition 4.4 gives 303 the bottleneck matching distance $BMD(\Gamma(G; p_1; F; q_1))$. 304

³⁰⁵ Let the complete bipartite graph $\Gamma(G, F)$ have weight ³⁰⁶ BMD($\Gamma(G; p_1; F; q_1)$) on each edge connecting vertices ³⁰⁷ representing $p_1 \in V(G)$ and $q_1 \in V(F)$. The Nested ³⁰⁸ Bottleneck Metric NBM(NCD(G; 2), NCD(F; 2)) is the ³⁰⁹ Bottleneck Matching Distance BMD($\Gamma(G, F)$).

311SRD(G) coincides with NCD(G; 0) after sorting, so NBM312can be defined as $L_{\infty}(SRD(G), SRD(F))$ for order h = 0.313The metrics W_{∞}, M_{∞} , NBM compare objects of the same314size. To compare graphs with different numbers of vertices,315 M_{∞} in Definition 4.5 can be replaced with Earth Mover's316Distance EMD in Definition C.2. All metric axioms and317main Theorem 4.6 below are proved in appendices C and D318for any dimension $n \ge 2$ and orders $1 \le h \le n$.

Theorem 4.6 (NCD solves Problem 1.1). (a) The Nested Centered Distribution NCD(G; h) in Definition 3.5 is invariant under any rigid motion for all Euclidean graph G on m unordered vertices and can be computed in time $O(n^2m^{h+1})$ with space $O(n^3 + hm^{h+1})$ for $h \le n = 2$.

325 **(b)** NCD(G; 2) is a complete invariant of all graphs $G \subset \mathbb{R}^2$ under rigid motion from the group SE(2).

(c) Perturbing each vertex of a graph $G \subset \mathbb{R}^2$ within its ε neighborhood changes NCD(G; h) up to 2ε in both metrics NBM and EMD for any order h = 1, 2.

(d) For any graphs $G, F \subset \mathbb{R}^2$ on m unordered vertices, the metrics NBM and EMD between the invariants NCD(G; h) and NCD(F; h) is computed in time $O(m^{2h+1.5} \log^{h+1} m)$ with space $O(m^{2h+1} \log^{h-1} m)$ for $h \leq n = 2$.

Theorem 4.6(b) implies that any graphs $G, F \subset \mathbb{R}^2$ are related by rigid motion *if and only if* NCD(G; 2) =NCD(F; 2). This equality is interpreted as a bijection $\text{NCD}(G; n) \to \text{NCD}(F; n)$ matching all CRs, which is equivalent to NBM = 0 by the first metric axiom. Since every CR can be stored in a vector form, the complete invariant NCD(G; 2) for n = 2 can be considered vectorial.

Table 2 emphasizes that most graphs should be first compared (or represented for machine learning) by simpler and faster invariants, so the complete NCD(G; n) is used only in rare cases but is still needed to distinguish all graphs.

Table 2. Invariants and metrics on graphs $G \subset \mathbb{R}^2$ with m unordered vertices: from the fastest (linear-time) to complete.

$\begin{array}{ccc} \text{SRD}(G) & O(m \log m) & L_{\infty} & O(m) \\ \text{SRD}(G) & O(m^2) & L & O(m^2) \end{array}$	
$\begin{array}{ccc} \operatorname{SFD}(G) & O(m^{2}) \\ \operatorname{PDD}(G) & O(m^{2}\log m) \\ \operatorname{NCD}(G;1) & O(m^{2}) \\ \operatorname{NCD}(G;2) & O(m^{3}) \\ \end{array} \\ \begin{array}{ccc} \operatorname{EMD} & O(m^{3}) \\ \operatorname{NBM} & O(m^{3.5}\log^{3}) \\ \operatorname{NBM} & O(m^{5.5}\log^{3}) \\ \end{array}$	m)

Example 4.7 (version of Theorem 4.6(b) for n = 1). For a graph $G \subset \mathbb{R}$ with the center of mass O(G) = 0, take any base vertex $p \in G$. Then $\operatorname{sign}(p)$ is the usual sign of $p \in \mathbb{R}$, D(p) is the signed distance -|p|, R(G;p) is the $2 \times (m - 1)$ matrix whose column for any vertex $q \in$ $V(G) - \{p\}$ consists of the signed distances d(q, 0) = -|q|and $d(q, p) = \pm |q - p|$, where the plus sign + indicates an edge between q, p, while the minus sign - means no edge.

For order h = 1, the Centered Representation is the pair CR(G; p) = [sign(p), -|p|, R(G; p)]. The base vertex pis fixed in the line \mathbb{R} by sign(p) and |p|. Any other vertex $q \in V(G) - \{p\}$ is uniquely determined in \mathbb{R} by its Euclidean distances |q|, |q - p| to the origin and the already fixed p. The location of any point $q \in \mathbb{R}$ is characterized by sign(q) and |q|, which helps unambiguously identify its Centered Representation CR(G; q) in the unordered collection NCD(G; 1) of all these CRs. The signs of d(q, q') in each R(G; q) determine the presence or absence of an edge of $G \subset \mathbb{R}$ between any vertices $q, q' \in V(G)$.

5. Experiments on largest molecular databases

The world's largest databases of 3D molecular geometry are QM9 (130K+ entries) (Ramakrishnan et al., 2014) and GD (GEOM_drugs of 31M+ entries) (Axelrod & Gomez-

Bombarelli, 2022), which have hundreds of 3D conformers
of *unordered* atoms for each of 621 and 61607 chemical
compositions, respectively. The Protein Data Bank has
backbones of *ordered* atoms classified by simpler invariants
(Anosova et al., 2025). All experiments took a few hours on
Ryzen 9 3950X 3.5 GHz, 64 MB of L3 cache, RAM 82GB.

The ICML guide for reviewing application-driven ML says that "novel ideas that are simple to apply may be especially valuable". To demonstrate the chemical importance of the linear-time invariant SRD, we extracted clouds of k = 10 neighbors around every atom, see their counts in Table 3.

338

339

340

341

342

343

353

354

355

357

358 359

360 361

362

363

367 368

Table 3. Counts of atoms by chemical elements in QM9 (2,407,753 atoms), GD0 (GEOM_drugs 0th conformers, 12,917,980 atoms).

2.4.5					
345 346	QM9: H	QM9: C	QM9: N	QM9: O	QM9: F
347	1,230,122	040,557	139,704	187,990	5,514
348	GD0: H	GD0: C	GD0: N	GD0: O	GD0: F
349	5,660,986	5,267,096	842,562	854,400	64,299
350	GD0: P	GD0: S	GD0: Cl	GD0: Br	GD0: I
351	1,350	159,648	53,404	14,010	225
352					

Though the data was skewed towards more popular elements H (hydrogen) and C (carbon), a default network in TensorFlow with 80/20 split for train/test achieved over 98% accuracy in predictions of the chemical element of a central atom by distances to only k = 3 nearest neighbours, see Table 4. Appendix A has all implementation details.

Table 4. Accuracies in percentages for predicting the chemical element of a central atom by a 4-layer network using *only the k shortest distances* to atomic neighbors within a molecular graph.

data	k = 2	k = 3	k = 4	k = 5	k = 6
QM9	94.63	98.64	98.24	98.54	98.77

In chemistry, both ML and non-ML predictions of elements
achieved only 86% on similar size data, see Table 7 summarized in (Vasylenko et al., 2025), because the underlying
descriptors were not invariant, e.g. under permutations of
atoms, which creates exponentially many representations of
the same molecule, incomplete, or their similarities failed
the triangle axiom, e.g. see (Steck et al., 2024).

376 High accuracies in Table 4 are rigorously explained by the 377 cascade comparisons on all atomic clouds (environments) 378 from QM9. Split all clouds from by the 1st distance (to the 379 nearest neighbor of a central atom p) rounded to 3 decimal 380 places in Å. This is a typical experimental precision, where 381 $1\text{\AA} = 10^{-10}m$ is approximately the smallest interatomic 382 distance. Second, split each subset with equal 1st distances 383 by 2nd distances, and so on up to k = 5 distances. All 2.4M+ 384

atomic clouds of different elements in QM9 were separated by the shorest distances to only 4 atomic neighbors.

The hierarchy of invariants in Fig. 2 and Table 2 transparently explained the reconstruction of chemical elements from distances to k nearest neghbors and inspired the harder task to reconstruct a chemical composition from a molecule-level (not atomwise) invariant of only atomic centers.

For molecular graphs from QM9, we computed the pseudometric L_{∞} (max absolute difference of corresponding coordinates) on all 873,527,974 pairs of SRDs, then 8,735,279 distances L_{∞} on the stronger SPDs for the 1% closest pairs, then 87,352 EMDs on PDDs for the 1% closest pairs, distances NBM on NCD(G; 1) and NCD(G; 2) for the top 10K closest pairs, and 64 NBMs on complete NCD(G, 3).

The invariants in Table 5 distinguish all chemically different molecules with NBM on complete invariants giving the largest separation. All chemical compositions in QM9 and GD were distinguished by the vector SRD of Euclidean distances (rounded to 3 decimal places in Å) from the molecular center of mass to 5 and 7 farthest atoms, respectively.

This transparent reconstruction of the full chemistry from precise enough atomic geometry gives hope to rigorously infer other molecular properties from geometric invariants.

Table 5. Chemically different molecules (given by QM9 ids) are geometrically distinguished by invariant metrics, see Fig. 3 (right).

smallest distances in Å, molecule A \neq molecule B				
SRD, $L_{\infty} = 0.021$, $H_3C_4N_3O_2(131923) \neq H_4C_5N_2O(5365)$				
SPD, $L_{\infty} = 0.055$, $H_3C_4N_5(123533) \neq H_3C_5N_3O(24547)$				
PDD, EMD = 0.051 , $H_3C_4N_5(123533) \neq H_3C_5N_3O(24521)$				
NCD, NBM = 0.071, $H_4C_5N_4(123532) \neq H_4C_6N_2O(24513)$				



Figure 3. Left: the smallest NBM ≈ 0.07 Å on NCD(G; 3) for chemically different molecules 123533 and 24521. **Right**: nearduplicate (almost flat) molecules 123532 and 24513 have the same composition and tiny EMD $\approx 2.37 \times 10^{-7}$ Å (not distinguishing mirror images) but a 100× higher NBM $\approx 2.95 \times 10^{-5}$ Å.

For QM9 molecul graphs, Fig. 4 and 9 NBM distances for different NCD invariants of orders h = 1, 2, 3.



Figure 4. Each dot is a comparison of molecular graphs from QM9: x = NBM on NCD(G; 1) vs y = NBM on NCD(G; 2).

406 407 **6. Discussion: conclusions and limitations**

403

404

405

408The comparisons of molecular graphs from QM9 and GD409imply that all chemically different molecules are rigidly410different, see the smallest distance NBM ≈ 0.07 Å on com-411plete invariants in Table 5. So the map {molecules} \rightarrow 412{graphs on atomic centers (without chemical elements)} is413injective on rigid classes and can be inverted on its image.

Hence the most important property (chemical composition)
is reconstructable from precise enough geometry. Using
only a few radial distances (5 at the atomic level and 7
at the molecular level, rounded to 3 decimal places) for
uniquely identifying all chemical elements in QM9 and GD
demonstrates the transparency of application-driven ML.

421 The solution to Problem 1.1 settled the long-standing chal-422 lenge of properly defining a *molecular structure*. A tradi-423 tional approach is to describe such a structure as "a set of unlabeled configurations that are relatively similar to each 424 other", quoted from the paragraph to the left of the caption 425 of Fig. 1 in (Lang et al., 2024). If this 'similarity' is treated 426 as an equivalence allowing perturbations of atoms up to a 427 positive threshold, sufficiently many perturbations can make 428 all molecules (of the same number of atoms) equivalent 429 430 by the transitivity axiom. A justified way to resolve this paradox is to embrace uncertainty and continuously quan-431 tify this similarity not by ignoring any perturbations up to a 432 threshold but by computing an exact distance satisfying all 433 434 metric axioms and Lipschitz continuity in Problem 1.1.

The question of whether to put close neighbors like nearduplicates in Fig. 3 (left) into one cluster of the "same"
molecules is rather administrative similar to assigning close
houses to one village (cluster) instead of different ones.

Studying molecules by fixing a composition is similar to drawing artificial boundaries between countries on Earth. Because some molecules of different compositions have close shapes as in Fig. 5, they should have similar properties. Now any properties of molecules should be possible to predict only from the complete invariant NCD(G; 3) even without chemistry in the same way as any precise geographic location uniquely determines all physical properties of this place such as the average annual temperature. Chemical compositions can be still helpful similar to the location's altitude, which easier predicts (say) the average temperature than theoretically sufficient geographic coordinates.

Any vertex p and edge of G can have an *attribute* and a *weight* respected by any isometry that maps one graph to another. These vertex attributes and edge weights can be incorporated as extra columns and rows in CRs from Definition 3.3, and then incorporated into NCD and NBM. We can compare graphs of different numbers of vertices because EMD works for both PDD and NCD as weighted distributions of any finite size. This comparison splits the vertices from V(G) into parts (subvertices) that are optimally 'transported' to a splitting of another vertex set V(F).

The main contribution is Theorem 4.6 and its extension in Theorem D.1 to all dimensions $n \ge 2$ fully solving Problem 1.1. The limitation is the time $O(n^2m^{n+1})$ of the complete invariant NCD(G; n) of any graphs $G \subset \mathbb{R}^n$. For a fixed dimension n, this polynomial complexity resolves two exponential-size challenges: m! permutations of m unordered vertices and up to $2^{m(m-1)/2}$ non-isometric graphs with up to m(m-1)/2 edges on m fixed vertices in \mathbb{R}^n .

In practice, all comparisons and property predictions can start from much faster (linear-time) invariants SRD and only in cases of close distances (potential confusions) progress to stronger invariants SPD, PDD, NCD. This hierarchical (cascade) computation can better address the curse of dimensionality instead of the one-size-fits-all approach.

A map f : objects \rightarrow descriptors \rightarrow properties is invertible only if objects are faithfully represented by complete invariants. Any non-invariant maps a single object to (usually infinitely) many values or representations. Any incomplete invariant can fail to differentiate between objects with different properties. Hence a generative approach (inverting f above) can succeed only after the *discriminative* problem is solved. The space $GRS(\mathbb{R}^3; m)$ of rigid classes of all graphs on m vertices in \mathbb{R}^3 contains all possible shapes of molecules (all already known and also all not yet discovered ones). The complete invariant NCD(G) of $G \subset \mathbb{R}^3$ defines geographic-style coordinates on a continuous map of $GRS(\mathbb{R}^3; m)$ containing QM9 and GD. Since the space $GRS(\mathbb{R}^3; m)$ is high-dimensional, we really need complete invariants to separate all known molecules and look for unexplored gaps containing new future molecules.

440 Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

441

442

443

444

445 446

447

448

449

450

451

452

453

454

455

456

- Abbe, E. and Sandon, C. On the universality of deep learning. *Advances in Neural Information Processing Systems*, 33:20061–20072, 2020.
- Agarwal, P. K., Fox, K., Nath, A., Sidiropoulos, A., and Wang, Y. Computing the gromov-hausdorff distance for metric trees. *ACM Transactions on Algorithms*, 14(2): 1–20, 2018.
- 457 Akhtar, N. and Mian, A. Threat of adversarial attacks on
 458 deep learning in computer vision: A survey. *IEEE Access*,
 459 6:14410–14430, 2018.
- 461 Amir, T., Gortler, S., Avni, I., Ravina, R., and Dym, N.
 462 Neural injective functions for multisets, measures and
 463 graphs via a finite witness theorem. *Advances in Neural*464 *Information Processing Systems*, 36, 2024.
- Anosova, O., Gorelov, A., Jeffcott, W., Jiang, Z., and Kurlin,
 V. A complete and bi-continuous invariant of protein backbones under rigid motion. *MATCH Communications in Mathematical and in Computer Chemistry (to appear), arxiv:2410.08203*, 2025.
- Antunes, L. M., Grau-Crespo, R., and Butler, K. T. Distributed representations of atoms and materials for machine learning. *npj Computational Materials*, 8(1):44, 2022.
- 476
 477 Arvind, V. and Rattan, G. The parameterized complexity
 477 of geometric graph isomorphism. *Algorithmica*, 75:258–
 478 276, 2016.
- 480 Axelrod, S. and Gomez-Bombarelli, R. Geom, energy481 annotated molecular conformations for property predic482 tion and molecular generation. *Scientific Data*, 9(1):185,
 483 2022.
- Babai, L. Canonical form for graphs in quasipolynomial time: preliminary report. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, pp. 1237–1246, 2019.
- Batzner, S., Musaelian, A., Sun, L., Geiger, M., Mailoa,
 J. P., Kornbluth, M., Molinari, N., Smidt, T. E., and
 Kozinsky, B. E(3)-equivariant graph neural networks for
 data-efficient and accurate interatomic potentials. *Nature communications*, 13(1):2453, 2022.

- Bonchev, D. Chemical graph theory: introduction and fundamentals, volume 1. CRC Press, 1991.
- Brécheteau, C. A statistical test of isomorphism between metric-measure spaces using the distance-to-a-measure signature. *Electronic J Statistics*, 13:795–849, 2019.
- Bronstein, M. M., Bruna, J., LeCun, Y., Szlam, A., and Vandergheynst, P. Geometric deep learning: going beyond Euclidean data. *IEEE Signal Processing Magazine*, 34 (4):18–42, 2017.
- Bronstein, M. M., Bruna, J., Cohen, T., and Veličković, P. Geometric deep learning: grids, groups, graphs, geodesics, and gauges. arXiv:2104.13478, 2021.
- Buchin, M., Chambers, E., Fang, P., Fasy, B. T., Gasparovic, E., Munch, E., and Wenk, C. Distances between immersed graphs: Metric properties. *La Matematica*, pp. 1–26, 2023.
- Bunch, J. R. and Hopcroft, J. E. Triangular factorization and inversion by fast matrix multiplication. *Mathematics* of Computation, 28(125):231–236, 1974.
- Chen, C., Ye, W., Zuo, Y., Zheng, C., and Ong, S. P. Graph networks as a universal machine learning framework for molecules and crystals. *Chemistry of Materials*, 31(9): 3564–3572, 2019.
- Chew, P. and Kedem, K. Improvements on geometric pattern matching problems. In *Scandinavian Workshop on Algorithm Theory*, pp. 318–325, 1992.
- Chew, P., Dor, D., Efrat, A., and Kedem, K. Geometric pattern matching in d-dimensional space. *Discrete & Computational Geometry*, 21(2):257–274, 1999.
- Choo, H. Y., Wee, J., Shen, C., and Xia, K. Fingerprintenhanced graph attention network (fingat) model for antibiotic discovery. *Journal of Chemical Information and Modeling*, 2023.
- Cohen, T. S., Geiger, M., and Weiler, M. A general theory of equivariant cnns on homogeneous spaces. *Advances in neural information processing systems*, 32, 2019.
- Colbrook, M. J., Antun, V., and Hansen, A. C. The difficulty of computing stable and accurate neural networks: On the barriers of deep learning and Smale's 18th problem. *Proc. National Academy of Sciences*, 119(12):e2107151119, 2022.
- Dekster, B. V. and Wilker, J. B. Edge lengths guaranteed to form a simplex. *Archiv der Mathematik*, 49(4):351–366, 1987.

- 495 Delle Rose, V., Kozachinskiy, A., Rojas, C., Petrache, M.,
 496 and Barceló, P. Three iterations of (d- 1)-wl test dis497 tinguish non isometric clouds of d-dimensional points.
 498 Advances in Neural Information Processing Systems, 36,
 499 2024.
- Deng, C., Litany, O., Duan, Y., Poulenard, A., Tagliasacchi,
 A., and Guibas, L. J. Vector neurons: A general framework for so(3)-equivariant networks. In *Proceedings of the International Conference on Computer Vision*, pp. 12200–12209, 2021.
- 506
 507
 508
 508
 Deza, E. and Deza, M. M. Encyclopedia of distances.
 Springer, 2009.
- 509 Dong, Y., Liao, F., Pang, T., Su, H., Zhu, J., Hu, X., and
 510 Li, J. Boosting adversarial attacks with momentum. In
 511 *Computer vision and pattern recognition*, pp. 9185–9193,
 512 2018.
- Duvenaud, D. K., Maclaurin, D., Iparraguirre, J., Bombarell, R., Hirzel, T., Aspuru-Guzik, A., and Adams, R. P. Convolutional networks on graphs for learning molecular fingerprints. *Advances in neural information processing systems*, 28, 2015.
- 519 Dym, N. and Gortler, S. J. Low-dimensional invariant em520 beddings for universal geometric learning. *Foundations*521 *of Computational Mathematics*, pp. 1–41, 2024.
- Efrat, A., Itai, A., and Katz, M. J. Geometry helps in
 bottleneck matching and related problems. *Algorithmica*,
 31(1):1–28, 2001.
 - Feynman, R. *The Feynman lectures on physics. Chapter 1: atoms in motion*, volume 1. 1971.

527

528

529

530

531

532

533

534

- Fisikopoulos, V. and Penaranda, L. Faster geometric algorithms via dynamic determinant computation. *Computational Geometry*, 54:1–16, 2016.
- Fredman, M. L. and Tarjan, R. E. Fibonacci heaps and their uses in improved network optimization algorithms. *Journal of the ACM*, 34(3):596–615, 1987.
- Fuchs, F., Worrall, D., Fischer, V., and Welling, M. Se(3)transformers: 3d roto-translation equivariant attention
 networks. *Advances in neural information processing systems*, 33:1970–1981, 2020.
- Gao, X., Hu, W., and Qi, G.-J. Graphter: Unsupervised
 learning of graph transformation equivariant representations via auto-encoding node-wise transformations. In *Proceedings of Computer Vision and Pattern Recognition*,
 pp. 7163–7172, 2020.
- Goldberg, A. and Tarjan, R. Solving minimum-cost flow problems by successive approximation. In *Proceedings of STOC*, pp. 7–18, 1987.

- Goodrich, M. T., Mitchell, J. S., and Orletsky, M. W. Approximate geometric pattern matching under rigid motions. *Transactions on Pattern Analysis and Machine Intelligence*, 21(4):371–379, 1999.
- Goyal, A., Law, H., Liu, B., Newell, A., and Deng, J. Revisiting point cloud shape classification with a simple and effective baseline. In *International Conference on Machine Learning*, pp. 3809–3820, 2021.
- Grinberg, D. and Olver, P. J. The n body matrix and its determinant. *SIAM Journal on Applied Algebra and Geometry*, 3(1):67–86, 2019.
- Grohe, M. and Schweitzer, P. The graph isomorphism problem. *Communications of the ACM*, 63(11):128–134, 2020.
- Guo, C., Gardner, J., You, Y., Wilson, A. G., and Weinberger, K. Simple black-box adversarial attacks. In *International Conference on Machine Learning*, pp. 2484–2493, 2019.
- Helfgott, H. A., Bajpai, J., and Dona, D. Graph isomorphisms in quasi-polynomial time. *arXiv:1710.04574*, 2017.
- Hopcroft, J. E. and Karp, R. M. An n⁵/2 algorithm for maximum matchings in bipartite graphs. *SIAM Journal* on Computing, 2(4):225–231, 1973.
- Hordan, S., Amir, T., Gortler, S. J., and Dym, N. Complete neural networks for euclidean graphs. In AAAI Conference on Artificial Intelligence, volume 38 (11), pp. 12482–12490, 2024.
- Horn, R. A. and Johnson, C. R. *Matrix analysis*. Cambridge University Press, 2012.
- Huttenlocher, D. P., Klanderman, G. A., and Rucklidge, W. J. Comparing images using the Hausdorff distance. *Transactions on pattern analysis and machine intelligence*, 15 (9):850–863, 1993.
- Kiefer, S., Ponomarenko, I., and Schweitzer, P. The weisfeiler–leman dimension of planar graphs is at most 3. *Journal of the ACM*, 66(6):1–31, 2019.
- Kondor, R. and Trivedi, S. On the generalization of equivariance and convolution in neural networks to the action of compact groups. In *International Conference on Machine Learning*, pp. 2747–2755, 2018.
- Laidlaw, C. and Feizi, S. Functional adversarial attacks. *Adv. Neural Information Proc. Systems*, 32, 2019.
- Lang, L., Cezar, H. M., Adamowicz, L., and Pedersen, T. B. Quantum definition of molecular structure. *Journal of the American Chemical Society*, 146(3):1760–1764, 2024.

- Leman, A. and Weisfeiler, B. A reduction of a graph to a
 canonical form and an algebra arising during this reduction. *Nauchno-Technicheskaya Informatsiya*, 2(9):12–16,
 1968.
- Li, Z., Wang, X., Huang, Y., and Zhang, M. Is distance matrix enough for geometric deep learning? *arXiv:2302.05743*, 2023.
- Liberti, L. and Lavor, C. *Euclidean distance geometry*.
 Springer, 2017.
- Lim, S., Mémoli, F., and Smith, Z. The gromov-hausdorff distance between spheres. *Geometry & Topology*, 27(9): 3733–3800, 2023.
- Maennel, H., Unke, O. T., and Müller, K.-R. Complete
 and efficient covariants for 3d point configurations with
 application to learning molecular quantum properties. *arXiv:2409.02730*, 2024.
- Majhi, S. and Wenk, C. Distance measures for geometric
 graphs. *arXiv:2209.12869*, 2022.

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

- Majhi, S., Vitter, J., and Wenk, C. Approximating gromovhausdorff distance in euclidean space. *Computational Geometry*, 116:102034, 2024.
- Maron, H., Ben-Hamu, H., Serviansky, H., and Lipman, Y. Provably powerful graph networks. *Advances in neural information processing systems*, 32, 2019.
- Melter, R. A. and Tomescu, I. Metric bases in digital geometry. *Computer vision, graphics, and image Processing*, 25(1):113–121, 1984.
- Mémoli, F. Gromov–Wasserstein distances and the metric approach to object matching. *Foundations of computational mathematics*, 11:417–487, 2011.
- 588 Mémoli, F., Smith, Z., and Wan, Z. The Gromov-Hausdorff
 589 distance between ultrametric spaces: its structure and
 590 computation. *arXiv:2110.03136*, 2021.
- Merchant, A., Batzner, S., Schoenholz, S. S., Aykol, M.,
 Cheon, G., and Cubuk, E. D. Scaling deep learning for materials discovery. *Nature*, 624(7990):80–85, 2023.
- Nemec, L. Principal component analysis (pca):
 A physically intuitive mathematical introduction. https://towardsdatascience.com/
 principal-component-analysis-pca-8133b0
 2022.
- Nigam, J., Pozdnyakov, S. N., Huguenin-Dumittan, K. K., and Ceriotti, M. Completeness of atomic structure representations. *APL Machine Learning*, 2(1), 2024.

- Nikolentzos, G., Meladianos, P., and Vazirgiannis, M. Matching node embeddings for graph similarity. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.
- Oliynyk, A. O., Antono, E., Sparks, T. D., Ghadbeigi, L., Gaultois, M. W., Meredig, B., and Mar, A. Highthroughput machine-learning-driven synthesis of fullheusler compounds. *Chemistry of Materials*, 28(20): 7324–7331, 2016.
- Pozdnyakov, S. N., Willatt, M. J., Bartók, A. P., Ortner, C., Csányi, G., and Ceriotti, M. Incompleteness of atomic structure representations. *Phys. Rev. Lett.*, 125:166001, 2020. URL arXiv:2001.11696.
- Qi, G.-J. and Luo, J. Small data challenges in big data era: A survey of recent progress on unsupervised and semisupervised methods. *Transactions on Pattern Analysis and Machine Intelligence*, 44(4):2168–2187, 2020.
- Ramakrishnan, R., Dral, P. O., Rupp, M., and Von Lilienfeld,
 O. A. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific data*, 1(1):1–7, 2014.
- Rass, S., König, S., Ahmad, S., and Goman, M. Metricizing the euclidean space towards desired distance relations in point clouds. *IEEE Transactions on Information Forensics and Security*, 2024.
- Rieder, S. R., Oliveira, M. P., Riniker, S., and Hünenberger, P. H. Development of an open-source software for isomer enumeration. *Journal of Cheminformatics*, 15(1):10, 2023.
- Rubner, Y., Tomasi, C., and Guibas, L. The Earth Mover's Distance as a metric for image retrieval. *International Journal of Computer Vision*, 40(2):99–121, 2000.
- Sato, R., Cuturi, M., Yamada, M., and Kashima, H. Fast and robust comparison of probability measures in heterogeneous spaces. arXiv:2002.01615, 2020.
- Schmiedl, F. Computational aspects of the Gromov– Hausdorff distance and its application in non-rigid shape matching. *Discrete Comp. Geometry*, 57:854–880, 2017.
- Schoenberg, I. Remarks to Maurice Frechet's article "Sur la definition axiomatique d'une classe d'espace distances vectoriellement applicable sur l'espace de Hilbert. *Annals of Mathematics*, pp. 724–732, 1935.
- tion. https://towardsdatascience.com/ Shirdhonkar, S. and Jacobs, D. Approximate earth mover's principal-component-analysis-pca-8133b02f1diktance in linear time. In *Conference on Computer Vision and Pattern Recognition*, pp. 1–8, 2008.
 - Sippl, M. J. and Scheraga, H. A. Cayley-menger coordinates. *Proceedings of the National Academy of Sciences*, 83(8): 2283–2287, 1986.

- Steck, H., Ekanadham, C., and Kallus, N. Is cosinesimilarity of embeddings really about similarity? In *Companion Proceedings of the ACM on Web Conference* 2024, pp. 887–890, 2024.
- Tshitoyan, V., Dagdelen, J., Weston, L., Dunn, A., Rong,
 Z., Kononova, O., Persson, K. A., Ceder, G., and Jain, A.
 Unsupervised word embeddings capture latent knowledge
 from materials science literature. *Nature*, 571(7763):95–
 98, 2019.
- Tu, E., Wang, Z., Yang, J., and Kasabov, N. Deep semisupervised learning via dynamic anchor graph embedding in latent space. *Neural Networks*, 146:350–360, 2022.
- 619 Vasylenko, A., Antypov, D., Schewe, S., Daniels, L. M.,
 620 Claridge, J. B., Dyer, M. S., and Rosseinsky, M. J. Dig621 ital features of chemical elements extracted from local
 622 geometries in crystal structures. *Digital Discovery*, 2025.
- Ward, L., Agrawal, A., Choudhary, A., and Wolverton, C. A
 general-purpose machine learning framework for predicting properties of inorganic materials. *npj Computational Materials*, 2(1):1–7, 2016.
 - Weinhold, F. Metric geometry of equilibrium thermodynamics. *The Journal of Chemical Physics*, 63(6):2479–2483, 1975.

629

630

631

651

- Weston, L., Tshitoyan, V., Dagdelen, J., Kononova, O.,
 Trewartha, A., Persson, K. A., Ceder, G., and Jain, A.
 Named entity recognition and normalization applied to
 large-scale information extraction from the materials science literature. *Journal of chemical information and modeling*, 59(9):3692–3702, 2019.
- 638
 639
 640
 640 for periodic crystals. Advances in Neural Information Processing Systems, 35:24625–24638, 2022.
- Widdowson, D. and Kurlin, V. Recognizing rigid patterns of unlabeled point clouds by complete and continuous isometry invariants with no false negatives and no false positives. In *Proceedings of CVPR*, pp. 1275–1284, 2023.
- Wilson, S. R., Cui, W., Moskowitz, J. W., and Schmidt,
 K. E. Applications of simulated annealing to the conformational analysis of flexible molecules. *Journal of computational chemistry*, 12(3):342–349, 1991.
- Zhang, B., Fan, C., Liu, S., Huang, K., Zhao, X., Huang,
 J., and Liu, Z. The expressive power of graph neural networks: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 37:1455–1474, 2024.
- Zhou, D.-X. Universality of deep convolutional neural networks. *Applied and computational harmonic analysis*, 48 (2):787–794, 2020.

Zhou, Q., Tang, P., Liu, S., Pan, J., Yan, Q., and Zhang, S.-C. Learning atoms for materials discovery. *Proceedings* of the National Academy of Sciences, 115(28):E6411– E6417, 2018.

A. Extra details of experiments on the world's largest 3D molecular databases QM9 and GD

The default 4-layer network from TensorFlow used a "sequential" mode, 3 epochs, and the settings in Table 6.

Table 6. Parameters of the TensorFlow network for predictions in Table 4.

Layer (type)	OUTPUT SHAPE	NUMBER OF PARAMETERS
DENSE (DENSE)	(None, 32)	352
BATCH_NORMALIZATION	(None, 32)	128
RE_LU (RELU)	(None, 32)	0
DENSE_1 (DENSE)	(None, 5)	165

Past maps of QM9 in Fig. 5 based on eigenvalues are too dense without clear separation. Even if we zoom in, these two
 or three incomplete invariants will not provide any extra separation. The complete invariants NDP contain much more
 geometric information.



Figure 5. Left: each dot represents one QM9 molecule whose atomic cloud has two largest roots $l_1 \ge l_2$ of eigenvalues (moments of inertia (Nemec, 2022) or elongations in two principal directions) in Angstroms ($1\mathring{A} = 10^{-10}m \approx$ smallest interatomic distance). The color represents the free energy *G* characterizing molecular stability. **Right**: each dot represents one QM9 molecule whose atomic cloud has coordinates x, y expressed via the roots $l_1 \ge l_2 \ge l_3 \ge 0$ of three eigenvalues.

Fig. 7 shows the simplest geographic-style map of QM9 as a finite sample within $\bigcup_{m=3}^{29} \text{GRS}(\mathbb{R}^3; m)$ projected to the invariants $\text{SRD}_1 \ge \text{SRD}_2$. All molecules on the horizontal axis $y = \text{SRD}_1 - \text{SRD}_2 = 0$ have $\text{SRD}_1 = \text{SRD}_2$ (due to two equidistant atoms from the center of mass) and can be projected (like any subset of QM9) to other coordinates as in Fig. 8. Molecular properties can be visualized on these geographic maps as 'mountainous' landscapes.

Table 7. Past ML and non-MI	predictions of chemical elements have lower acc	curacies than by distance invariants in Table 4.
-----------------------------	---	--

Method	DESCRIPTION	ACCURACY	Reference
LEAF	LOCAL COORDINATION GEOMETRY	86%	(VASYLENKO ET AL., 2025)
MATSCHOLAR	ML-DERIVED FROM LITERATURE	81%	(WESTON ET AL., 2019)
MAT2VEC	ML-DERIVED FROM LITERATURE	80%	(TSHITOYAN ET AL., 2019)
ATOM2VEC	ML-DERIVED FROM COMPOSITIONAL CONTENT	79%	(ZHOU ET AL., 2018)
GNoME	FREQUENCY OF ELEMENTS AT THE SAME ATOMIC SITES	79%	(MERCHANT ET AL., 2023)
MAGPIE	ELEMENTAL PHYSICAL CHARACTERISTICS	78%	(WARD ET AL., 2016)
Oliynyk	ELEMENTAL PHYSICAL CHARACTERISTICS	75%	(OLIYNYK ET AL., 2016)
MEGNET	ML-DERIVED FROM ATOM, BOND AND GRAPH ATTRIBUTES	73%	(CHEN ET AL., 2019)
SkipAtom	ML-DERIVED FROM ATOM CONNECTIVITY GRAPHS	68%	(ANTUNES ET AL., 2022)

Continuous machine learning on Euclidean graphs



Figure 6. **Left**: the heatmap of all molecular graphs from QM9 in the simplest continuous invariants. **Right**: 18336 graphs with 19 atoms. The color indicates the number of molecules at every pixel.



Figure 7. Every dot represents a molecular graph with the invariant coordinates $x = \text{SRD}_1$, $y = \text{SRD}_1 - \text{SRD}_2$, all in Angstroms, where $1\text{\AA} = 10^{-10}m \approx$ the smallest interatomic distance.

Continuous machine learning on Euclidean graphs



B. Invariants and metrics on Euclidean graphs in any dimension $n \ge 2$

805

This section extends all new concepts and results from sections 3 and 4 to any dimension $n \ge 2$. Any *n* vectors $p_1, \ldots, p_n \in \mathbb{R}^n$ can be written as columns in the $n \times n$ matrix whose determinant has $sign(p_1, \ldots, p_n)$, which is ± 1 or 0 (if p_1, \ldots, p_n are linearly dependent).

Definition B.1 (Centered Representation CR(G; A) of a graph with a sequence $A \subset V(G)$). Let $G \subset \mathbb{R}^n$ be a graph on m unordered points with the center of mass O(G) = 0. For any $1 \le h \le n$, fix a base sequence A of ordered vertices $p_1, \ldots, p_h \in V(G)$. If h = n, let sign(A) be the sign of the $n \times n$ determinant on the vectors p_1, \ldots, p_n , else sign(A) = 0. Let D(A) be the matrix of signed distances between the ordered points $0 = p_0, p_1, \ldots, p_h$. The matrix R(G; A) has m - hunordered columns, one for each vertex $q \in V(G) - A$, consisting of h + 1 distances $d(q, p_i)$ for $i = 0, \ldots, h$, where $p_0 = 0$. The Centered Representation CR(G; A) is the triple [sign(A), D(A), R(G; A)].

Definition B.2 (Nested Centered Distribution NCD(G; h) of order h). Let $G \subset \mathbb{R}^n$ be any Euclidean graph on m unordered vertices and the center of mass at the origin $0 \in \mathbb{R}^n$. Fix an order $1 \le h \le n$.

(a) For any h - 1 distinct ordered vertices $p_1, \ldots, p_{h-1} \in V(G)$, the Centered Distribution $CD_{h-1}^{(h)}(G; p_1, \ldots, p_{h-1})$ of index h - 1 is the unordered set of Centered Representation $CR(G; p_1, \ldots, p_h)$ from Definition B.1 for all $p_h \in V(G) - \{p_1, \ldots, p_{h-1}\}$.

(b) Now we will iteratively decrement an integer k from h - 1 down to 1 and define $CD_{h-2}^{(h)}$ of index h - 2, and so on until $CD_1^{(h)}$ of index k = 1. For the initial k = h - 1, we use $CD_k^{(h)} = CD_{h-1}^{(h)}$ defined in part (a) above. For any k - 1 distinct



Figure 9. Each dot is a comparison of molecular graphs from QM9 by the distances on the progressively stronger invariants: NCD(G; 2) vs NCD(G; 3).

ordered vertices $p_1, \ldots, p_{k-1} \in V(G)$, the Centered Distribution $CD_{k-1}^{(h)}(G; p_1, \ldots, p_{k-1})$ of index k - 1 is the unordered collection of $CD_k^{(h)}(C; p_1, \ldots, p_k)$ of index k for all vertices $p_k \in V(G) - \{p_1, \ldots, p_{k-1}\}$.

(c) The Nested Centered Distribution NCD(G; h) of order h is the unordered collection of $CD_1^{(h)}(G; p_1)$ of index 1 for all vertices $p_1 \in V(G)$. For the order h = n, the mirror image $\overline{NCD}(G; n)$ is obtained from NCD(G; n) by reversing $sign(p_1, \ldots, p_n)$ of $n \times n$ determinants in all $CR; p_1, \ldots, p_n$).

If a sequence $0 \cup A = (p_0, p_1, \dots, p_n) \subset \mathbb{R}^n$ degenerates to a lower dimensional subspace, i.e. the vectors p_1, \dots, p_n become linearly dependent, then sign(A) of discontinuously changes. To guarantee the Lipschitz continuity, we multiply these signs by the strength σ below, while the volume $vol(0 \cup A)$ of the simplex on $0 \cup A$ is not Lipschitz continuous.

Definition B.3 (strength $\sigma(C)$). For any sequence C of n + 1 ordered points $p_0, \ldots, p_n \in \mathbb{R}^n$, the half-perimeter $p(C) = \frac{1}{2} \sum_{1 \le i < j \le n} |p_i - p_j|$ is the half-sum of pairwise distances between points of C. Let vol(C) denote the volume of the

n-dimensional simplex on C. The strength of the simplex C is $\sigma(C) = \frac{\operatorname{vol}^2(C)}{p^{2n-1}(C)}$.

In dimension n = 1, for any pair $C = \{p_0, p_1\} \subset \mathbb{R}$, the volume vol(C) is the length $|p_0 - p_1|$, the half-perimeter distance

⁸⁸⁰₈₈₁ p(C) is the half-distance $\frac{1}{2}|p_0 - p_1|$, so the strength is $\sigma(C) = \frac{\operatorname{vol}^2(C)}{p(C)} = 2|p_0 - p_1|$.

⁸⁸² **Lemma B.4** (Theorem 4.4 in (Widdowson & Kurlin, 2023)). Let *B* be obtained from a sequence $A \subset \mathbb{R}^n$ of *n* points by perturbing every point within its ε -neighborhood. Then $|\sigma(A) - \sigma(B)| \le 2\varepsilon\lambda_n$ for a constant λ_n , where $\lambda_1 = 2$, $\lambda_2 = 2\sqrt{3}$, $\lambda_3 \approx 0.43$.

Definition B.5 (max metric M_{∞} on CRs). Let Euclidean graphs $G, F \subset \mathbb{R}^n$ on m unordered vertices have base sequences $A, B \text{ of } h \leq n \text{ vertices. Consider the } m - h \text{ columns of } R(G; A) \text{ as a cloud of } m - h \text{ unordered points in } \mathbb{R}^h, \text{ also}$ for R(F; B). The max metric $M_{\infty}(\operatorname{CR}(G; A), \operatorname{CR}(F; B))$ is the maximum of $\frac{2}{\lambda_n}|\operatorname{sign}(A)\sigma(0 \cup A) - \operatorname{sign}(B)\sigma(0 \cup B)|,$ $L_{\infty}(D(A), D(B)), \text{ and the bottleneck distance } W_{\infty}(R(G; A), R(F; B)), \text{ where all signs are zeros for } h < n.$

⁸⁹¹ In Definition B.5, λ_n is the Lipschitz constant of σ from Lemma 4.2.

Definition B.6 (Nested Bottleneck Metric NBM on NCDs). Let $G, F \,\subset\, \mathbb{R}^n$ be any Euclidean graphs on m unordered vertices. For any ordered vertices $p_1 \ldots, p_{h-1} \in V(G)$ and $q_1 \ldots, q_{h-1} \in V(F)$, the complete bipartite graph $\Gamma(G; p_1, \ldots, p_{h-1}; F; q_1, \ldots, q_{h-1})$ has m - h + 1 white vertices and m - h + 1 black vertices representing $\operatorname{CR}(G; p_1, \ldots, p_h)$ and $\operatorname{CR}(F; q_1, \ldots, q_h)$ for all m - h + 1 vertices $p_h \in V(G) - \{p_1, \ldots, p_{h-1}\}$ and $q_h \in V(F) - \{q_1, \ldots, q_{h-1}\}$, respectively. Set the weight w(e) of an edge e joining the vertices represented by $\operatorname{CR}(G; p_1, \ldots, p_h)$, $\operatorname{CR}(F; q_1, \ldots, q_h)$ as the max metric M_{∞} between these distributions, see Definition B.5. Then Definition 4.4 gives the bottleneck matching distance $\operatorname{BMD}(\Gamma(G; p_1, \ldots, p_{h-1}; F; q_1, \ldots, q_{h-1}))$.

900 For any integer $1 \leq i < h$ and ordered vertices $p_1 \dots, p_{i-1} \in V(G)$ and $q_1 \dots, q_{i-1} \in V(F)$, the com-901 plete bipartite graph $\Gamma(G; p_1, \ldots, p_{i-1}; F; q_1, \ldots, q_{i-1})$ has m - i + 1 white vertices and m - i + 1 black vertices representing $\operatorname{CD}_{i}^{(h)}(G; p_{1}, \ldots, p_{i})$ and $\operatorname{CD}_{i}^{(h)}(F; q_{1}, \ldots, q_{i})$ for all m - i + 1 variable vertices $p_{i} \in V(G) - \{p_{1}, \ldots, p_{i-1}\}$ and $q_{i} \in V(F) - \{q_{1}, \ldots, q_{i-1}\}$, respectively. Set the weight w(e) of an edge e join-902 903 904 ing the vertices represented by $CD_i^{(h)}(G; p_1, \ldots, p_i)$ and $CD_i^{(h)}(F; q_1, \ldots, q_i)$ as the previously computed distance 905 BMD($\Gamma(G; p_1, \ldots, p_i; F; q_1, \ldots, q_i)$) for a smaller number *i* of fixed vertices. Then Definition 4.4 gives the bottle-906 neck matching distance BMD($\Gamma(G; p_1, \ldots, p_{i-1}; F; q_1, \ldots, q_{i-1})$). For i = 1, the graph $\Gamma(G, F)$ has m + m ver-907 tices representing $CD_1(G; p_1)$, $CD_1^{(h)}(F; q_1)$ for all $p_1 \in V(G)$ and $q_1 \in V(F)$. The Nested Bottleneck Metric 908 $\operatorname{NBM}(\operatorname{NCD}(G; h), \operatorname{NCD}(F; h))$ is the Bottleneck Matching Distance $\operatorname{BMD}(\Gamma(G, F))$. 909

911 C. Metrics on graphs and their continuity under perturbations

912913 This appendix verifies the axioms and Lipschitz continuity for all auxiliary metrics in section 4.

914 **Lemma C.1** (metric axioms for the bottleneck matching distance BMD). Let S, Q be any unordered distributions of the 915 same number of objects with a base metric d. Define the complete bipartite graph $\Gamma(S, Q)$ whose every edge e joining 916 objects $R_S \in S$ and $R_Q \in Q$ has the weight $w(e) = d(R_S, R_Q)$. Then the bottleneck matching distance BMD($\Gamma(S, Q)$) 917 from Definition 4.4 satisfies all metric axioms on such unordered distributions. 918

Proof of Lemma C.1. The coincidence axiom means that NBM(S, Q) = 0 if and only if the weighted distributions S, Q are equal in the sense that there is a bijection $g: S \to Q$ so that d(g(R), R) = 0 for any $R \in S$.

Indeed, if the weighted distributions S, Q can be matched by a bijection, we get a vertex matching E of $\Gamma(S, Q)$ whose all edges have weights w(e) = 0. Definition 4.4 implies that $BMD(\Gamma(S, Q)) = 0$ as required.

Conversely, if $BMD(\Gamma(S,Q)) = 0$, there is a vertex matching E in $\Gamma(S,Q)$ with all w(e) = 0. This matching E defines a required bijection $S \to Q$. The symmetry $BMD(\Gamma(S,Q)) = BMD(\Gamma(Q,S))$ follows from Definition 4.4 and the symmetry of the base metric d.

To prove the triangle inequality

910

919

920

921

924

925

926

927 928

929 930

933

934

 $BMD(\Gamma(S,Q)) + BMD(\Gamma(Q,T)) \ge BMD(\Gamma(S,T)),$

931 932 let E_{SQ}, E_{QT} be optimal vertex matchings in the graphs $\Gamma(S, Q), \Gamma(Q, T)$, respectively, such that

$$BMD(\Gamma(S,Q)) = W(E_{SQ}), BMD(\Gamma(Q,T)) = W(E_{QT}),$$

see Definition 4.4. The composition $E_{SQ} \circ E_{QT}$ is a vertex matching in $\Gamma(S, T)$, so $W(E_{SQ} \circ E_{QT}) \ge BMD(\Gamma(S, T))$. It suffices to prove that

$$W(E_{SQ}) + W(E_{QT}) \ge W(E_{SQ} \circ E_{QT}).$$

Let e_{ST} be an edge with a largest weight from $E_{SQ} \circ E_{QT}$, so $W(E_{SQ} \circ E_{QT}) = w(e_{ST})$. The edge e_{ST} can be considered the union of edges $e_{SQ} \in E_{SQ}$, $e_{QT} \in E_{QT}$.

By the triangle inequality for the base metric d,

$$w(e_{SQ}) + w(e_{QT}) \ge w(e_{ST}) = W(E_{SQ} \circ E_{QT})$$

implies that

$$W(E_{SQ}) + W(E_{QT}) \ge W(E_{SQ} \circ E_{QT})$$

because both terms on the left-hand side are maximized for all edges (not only e_{SQ}, e_{QT}) from E_{SQ}, E_{QT} .

Definition C.2 below makes sense for any distributions $\{[R_1, w_1], \ldots, [R_m, w_m]\}$, where R_1, \ldots, R_m are objects with a base metric d and weights $w_1, \ldots, w_m \in [0, 1]$. Each R_i can be CBR or CBD of any depth with a base metric M_{∞} or BMD from Definitions B.5, B.6.

Definition C.2 (EMD). Let $S = \{[R_i(S), w_i(S)]\}_{i=1}^{m(S)}$ and $Q = \{[R_j(Q), w_j(Q)]\}_{j=1}^{m(Q)}$ be weighted distributions of objects $R_i(S), R_j(Q)$, which live in a space with a metric d. A flow from S to Q is an $m(S) \times m(Q)$ matrix whose element $f_{ij} \in [0, 1]$ represents a partial flow from $R_i(S)$ to $R_j(Q)$. The Earth Mover's Distance is the minimum cost $\text{EMD}(S, Q) = \sum_{i=1}^{m(S)} \sum_{j=1}^{m(Q)} f_{ij} d(R_i(S), R_j(Q))$ for variable 'flows' $f_{ij} \in [0, 1]$ subject to the conditions $\sum_{j=1}^{m(Q)} f_{ij} \le w_i(S)$ for $i = 1, \ldots, m(S), \sum_{i=1}^{m(S)} f_{ij} \le w_j(Q)$ for $j = 1, \ldots, m(Q)$, and $\sum_{i=1}^{m(S)} \sum_{j=1}^{m(Q)} f_{ij} = 1$.

The first condition $\sum_{j=1}^{m(Q)} f_{ij} \leq w_i(S)$ means that not more than the weight $w_i(S)$ of $R_i(S)$ 'flows' into all $R_j(Q)$ via 'flows'

 $f_{ij}, j = 1, \dots, m(Q)$. The second condition $\sum_{i=1}^{m(S)} f_{ij} = w_j(Q)$ means that all 'flows' f_{ij} from $R_i(S)$ for $i = 1, \dots, m(S)$

'flow' into $R_j(Q)$ up to the maximum weight $w_j(Q)$. The last condition $\sum_{i=1}^{m(S)} \sum_{j=1}^{m(Q)} f_{ij} = 1$ forces to 'flow' all rows $R_i(S)$ to all rows $R_i(Q)$

to all rows $R_j(Q)$.

The EMD satisfies all metric axioms, see the appendix in (Rubner et al., 2000), needs $O(m^3 \log m)$ time for distributions of a maximum size m and is approximated in O(m) time, see (Shirdhonkar & Jacobs, 2008; Sato et al., 2020).

Definition C.2 can be adapted for the EMD between NDDs by (1) replacing the bottleneck distance W_{∞} in Definition B.5 with EMD between clouds of equally weighted points, and (2) replacing BMD(Γ) for a bipartite graph Γ with EMD(Γ) between the unordered sets (of potentially different sizes) of BDDs with weights on all white vertices and BDDs on all black vertices.

The Lipschitz continuity of NDD and EMD in Theorem D.1(c) needs Lemmas C.3, C.4, D.9.

If a metric graph G lives in an ambient metric space X, a natural perturbation of G is a shift of every vertex of G up to ε in the metric of X. Then the distance d(p,q) between any vertices p,q of G changes by at most 2ε .

We will prove the continuity in more general settings by only assuming that d(p,q) changes by at most 2ε for any $p, q \in V(G)$ without requiring an ambient space X.

Lemma C.3 (Lipschitz continuity of BMD). Let Γ be a complete bipartite graph with a vertex matching E such that any $e \in E$ has a weight $w(e) \leq \varepsilon$. Then BMD(Γ) $\leq \varepsilon$.

Proof of Lemma C.3. By Definition 4.4, the given matching E has the weight $W(E) = \max_{e \in E} w(e) \le \varepsilon$. Since $BMD(\Gamma) = \min_{E} W(E)$ is minimized for all vertex matchings, we get $BMD(\Gamma) \le \varepsilon$.

990 **Lemma C.4** (Lipschitz continuity of EMD). In Definition C.2, let distributions S, Q have a bijection $R_i(S) \leftrightarrow R_i(Q)$ 991 between equally weighted objects such that $d(R_i(S), R_i(Q)) \leq \varepsilon$ for all i = 1, ..., m, where m = m(S) = m(Q). Then 992 EMD $(S, Q) \leq \varepsilon$.

994 995 996 997 Proof of Lemma C.4. In Definition C.2, choose partial flows $f_{ij} = \frac{1}{m}$ for i = j, otherwise $f_{ij} = 0$. Then $\text{EMD}(S, Q) \leq \frac{1}{2} \sum_{i=1}^{m} \sum_{j=1}^{m} f_{ij} d(R_i(S), R_j(Q)) = \sum_{i=1}^{m} \frac{1}{m} d(R_i(S), R_i(Q)) \leq \frac{1}{m} \sum_{i=1}^{m} \varepsilon = \varepsilon.$

D. Proofs for Euclidean graphs from section 3

998

1035 1036

1040 1041

1001 This appendix rigorously proves all parts of Theorem D.1.

Theorem D.1 (NCD solves Problem 1.1). (a) The Nested Centered Distribution NCD(G; h) in Definition B.2 is invariant under any rigid motion for all Euclidean graph G on m unordered vertices and, for a fixed dimension n, can be computed in time $O(n^2m^{h+1})$ with space $O(n^2m^{h+1})$ for any order $1 \le h \le n$.

(b) NCD(G; 2) is a complete invariant of all graphs $G \subset \mathbb{R}^2$ under rigid motion from the group SE(n) in any dimension $n \geq 1$.

1008 (c) Perturbing each vertex of a graph $G \subset \mathbb{R}^n$ within its ε -neighborhood changes NCD(G;h) up to 2ε in both metrics 1009 NBM and EMD for any order $1 \le h \le n$.

(d) For any graphs $G, F \subset \mathbb{R}^n$ on m unordered vertices, the metrics NBM and EMD between the invariants NCD(G; h)and NCD(F; h) from Definition B.6 can be computed in time $O(m^{2h+1.5} \log^{h+1} m)$ with space $O(n^2m^{2h+1} \log^{h-1} m)$ for any order $1 \le h \le n$.

The *affine dimension* $0 \le \operatorname{aff}(A) \le n$ of a cloud $A = \{p_1, \ldots, p_m\} \subset \mathbb{R}^n$ is the maximum dimension of the vector space generated by all inter-point vectors $p_i - p_j$, $i, j \in \{1, \ldots, m\}$. Then $\operatorname{aff}(A)$ is an isometry invariant and is independent of an order of points of A. Any cloud A of 2 distinct points has $\operatorname{aff}(A) = 1$. Any cloud A of 3 points that are not in the same straight line has $\operatorname{aff}(A) = 2$.

¹⁰¹⁹ Lemma D.2 provides a simple criterion for a matrix to be realizable by squared distances of a point cloud in \mathbb{R}^n .

1021 **Lemma D.2** (realization of distances). (a) A symmetric $m \times m$ matrix of $s_{ij} \ge 0$ with $s_{ii} = 0$ is realizable as a 1022 matrix of squared distances between points $p_0 = 0, p_1, \ldots, p_{m-1} \in \mathbb{R}^n$ if and only if the $(m-1) \times (m-1)$ matrix 1023 $g_{ij} = \frac{s_{0i} + s_{0j} - s_{ij}}{2}$ has only non-negative eigenvalues.

1025 **(b)** If the condition in (a) holds, aff $(0, p_1, ..., p_{m-1})$ equals the number $k \le m-1 \le n$ of positive eigenvalues. Also in 1026 this case, $g_{ij} = p_i \cdot p_j$ define the Gram matrix GM of the vectors $p_1, ..., p_{m-1} \in \mathbb{R}^n$, which are uniquely determined in 1027 time $O(m^3)$ up to an orthogonal map in \mathbb{R}^n .

1029 **Proof of Lemma** D.2. (a) We extend Theorem 1 from (Dekster & Wilker, 1987) to the case m < n + 1 and justify the reconstruction of p_1, \ldots, p_{m-1} in time $O(m^3)$ uniquely in \mathbb{R}^n up to an orthogonal map from O(n).

1032 The part *only if* \Rightarrow . Let a symmetric matrix *S* consist of squared distances between points $p_0 = 0, p_1, \dots, p_{m-1} \in \mathbb{R}^n$. For 1033 $i, j = 1, \dots, m-1$, the matrix with the elements

$$g_{ij} = \frac{s_{0i} + s_{0j} - s_{ij}}{2} = \frac{p_i^2 + p_j^2 - |p_i - p_j|^2}{2} = p_i \cdot p_j$$

is the Gram matrix, which can be written as $GM = P^T P$, where the columns of the $n \times (m-1)$ matrix P are the vectors p_1, \ldots, p_{m-1} . For any vector $v \in \mathbb{R}^{m-1}$, we have

$$0 \le |Pv|^2 = (Pv)^T (Pv) = v^T (P^T P) v = v^T \mathbf{GM} v$$

Since the quadratic form $v^T GMv \ge 0$ for any $v \in \mathbb{R}^{m-1}$, the matrix GM is positive semi-definite meaning that GM has only non-negative eigenvalues, see Theorem 7.2.7 in (Horn & Johnson, 2012).

1045 The part $if \leftarrow$. For any positive semi-definite matrix GM, there is an orthogonal matrix Q such that $Q^T GMQ = D$ is the 1046 diagonal matrix, whose m - 1 diagonal elements are non-negative eigenvalues of GM. The diagonal matrix \sqrt{D} consists of 1047 the square roots of eigenvalues of GM.

(b) The number of positive eigenvalues of GM equals the dimension $k = \operatorname{aff}(\{0, p_1, \dots, p_{m-1}\})$ of the subspace in \mathbb{R}^n linearly spanned by p_1, \dots, p_{m-1} . We may assume that all $k \leq n$ positive eigenvalues of GM correspond to the first kcoordinates of \mathbb{R}^n . Since $Q^T = Q^{-1}$, the given matrix $\operatorname{GM} = QDQ^T = (Q\sqrt{D})(Q\sqrt{D})^T$ becomes the Gram matrix of the columns of $Q\sqrt{D}$. These columns become the reconstructed vectors $p_1, \dots, p_{m-1} \in \mathbb{R}^n$.

If there is another diagonalization $\tilde{Q}^T GM\tilde{Q} = \tilde{D}$ for $\tilde{Q} \in O(n)$, then \tilde{D} differs from D by a permutation of eigenvalues, which is realized by an orthogonal map, so we set $\tilde{D} = D$. Then $GM = \tilde{Q}D\tilde{Q}^T = (\tilde{Q}\sqrt{D})(\tilde{Q}\sqrt{D})^T$ is the Gram matrix of the columns of $\tilde{Q}\sqrt{D}$.

1057 The new columns differ from the previously reconstructed vectors $p_1, \ldots, p_{m-1} \in \mathbb{R}^n$ by the orthogonal map $Q\tilde{Q}^T$. Hence 1058 the reconstruction is unique up to O(n)-transformations. Computing eigenvectors p_1, \ldots, p_{m-1} requires a diagonalization 1059 of GM in time $O(m^3)$ (?)section 11.5]press2007numerical.

Though Lemma D.2 gives a two-sided criterion for realizability of distances by points $p_1, \ldots, p_m \in \mathbb{R}^n$, the space of distance matrices is highly singular and cannot be easily sampled. Even m = 4 points in \mathbb{R}^2 have 6 distances that should satisfy a polynomial equation saying that the tetrahedron with these 6 edge lengths has volume 0. So a randomly sampled matrix of potential distances for m > n + 1 is unlikely to be realizable by a cloud of m ordered points in \mathbb{R}^n .

1066 Chapter 3 in (Liberti & Lavor, 2017) discusses realizations of a complete graph given by a distance matrix in \mathbb{R}^n . 1067 Lemma D.3(a) and later results hold for all clouds including degenerate ones, e.g. for 3 points in a straight line.

Any points $p_1, \ldots, p_{n-1} \in A$ have $\operatorname{aff}(p_1, \ldots, p_{n-1}) \leq n-2$. For example, any two distinct points in $A \subset \mathbb{R}^3$ generate a straight line. In \mathbb{R}^2 , any point $p_1 \neq O(A)$ forms a suitable $\{p_1\}$. In \mathbb{R}^3 , one can choose any distinct points $p_1, p_2 \in A$ so that the infinite straight line via p_1, p_2 avoids O(A).

If there are no such p_1, p_2 , then $A \subset \mathbb{R}^3$ is contained in a straight line L, so aff(A) = 1. In this degenerate case, the stronger condition aff $(O(A) \cup \{p_1, \dots, p_{n-1}\}) = aff(A)$ will help reconstruct $A \subset L$ by using any point $p_1 \neq O(A)$. The first step is to reconstruct any ordered sequence from its distance matrix in Lemma D.3(a).

1076 Lemma D.3(a) holds for all degenerate clouds, e.g. for three points are in a straight line.

1077 **Lemma D.3** (reconstruction of ordered points). (a) Any sequence of ordered points $A = (p_1, ..., p_m)$ in \mathbb{R}^n can be 1078 reconstructed (uniquely up to isometry) from the matrix of the Euclidean distances $|p_i - p_j|$ in time $O(m^3)$. If all distances 1079 are divided by $R = \max_{i=1,...,m} |p_i|$, the reconstruction of $A \subset \mathbb{R}^n$ is unique up to isometry and uniform scaling.

(b) If $m \le n$, the uniqueness of reconstructions in part (a) holds if we replace isometry with rigid motion. Hence any n-1ordered points p_1, \ldots, p_{n-1} can be uniquely reconstructed from all pairwise distances between $0, p_1, \ldots, p_{n-1}$ up to SO(n) rotation around the origin $0 \in \mathbb{R}^n$.

Proof of Lemma D.3. (a) By translation, we can put p_1 at the origin $0 \in \mathbb{R}^n$. Let GM be the $(m-1) \times (m-1)$ matrix $g_{ij} = \frac{p_i^2 + p_j^2 - |p_i - p_j|^2}{2} = p_i \cdot p_j$ constructed from squared distances between $p_1 = 0, \ldots, p_m$ for $i, j = 2, \ldots, m$. By Lemma D.2(b) if GM has $k \le n$ positive eigenvalues, then $p_1 = 0, \ldots, p_m$ can be uniquely determined up to isometry in $\mathbb{R}^k \subset \mathbb{R}^n$ in time $O(m^3)$. If all distances are divided by the same radius R, the above construction guarantees uniqueness up to isometry and uniform scaling.

1092 (b) If $m \le n$, any mirror image of $A \subset \mathbb{R}^n$ after a suitable rigid motion in \mathbb{R}^n can be assumed to belong to an 1093 (n-1)-dimensional hyperspace $H \subset \mathbb{R}^n$, where they are matched by a mirror reflection $H \to H$ with respect to an 1094 (n-2)-dimensional subspace $S \subset H$. This reflection is realized by the SO(n) rotation through 180° around S. 1095

1096 Lemma D.3(b) for m = n = 3 implies that any triangle is determined by its sides up to rigid motion in \mathbb{R}^3 . For example, 1097 the sides 3, 4, 5 define a right-angled triangle whose mirror images are not related by rigid motion inside a plane $H \subset \mathbb{R}^3$, 1098 but are matched by composing a suitable rigid motion in H and a 180° rotation of \mathbb{R}^3 around a line in H. 1100 **Lemma D.4** (time of determinant). Any $n \times n$ determinant can be computed in time $O(n^3)$ with space $O(n^3)$.

1101

1127

1141

Proof of Lemma D.4. Any $n \times n$ determinant can be computed by Gaussian elimination in time $O(n^3)$ with space $O(n^3)$, see (Bunch & Hopcroft, 1974). The more recent theoretical estimate is $O(n^{2.373})$ by (Fisikopoulos & Penaranda, 2016).

Proof of Theorem D.1(a). Any rigid motion of \mathbb{R}^n mapping a Euclidean graph $G \subset \mathbb{R}^n$ to another graph F is a bijection preserving distances and signs of determinants, and hence induces a bijection $\text{CBR}(G; p_1, \ldots, p_i) \to \text{CBR}(F; q_1, \ldots, q_i)$ for all $p_1, \ldots, p_i \in V(G)$ and corresponding vertices $q_1, \ldots, q_i \in V(F)$ for any $i = 1, \ldots, h$, which implies a bijection NCD(G; h) \to NCD(F; h). By Definition 3.5, if G has m unordered vertices, the NCD(G) consists of $m(m-1) \ldots (m-h+1) = O(m^h)$ Centered Base Representations CBR(G; A) for all base sequences $A \subset V(G)$ of h ordered vertices.

1111 Every CBR(G; A) consists of the three components sign(A), CD(A), CR(G; A). For h = n, sign(A) is the $n \times n$ 1112 determinant computable in time $O(n^3)$ with space $O(n^3)$ by Lemma D.4. The distance matrix CD(A) needs $O(h^2)$ time 1113 and $O(h^2)$ space. The $(h + 1) \times (m - h)$ matrix CR(G; A) has O(hm) distances, each computable in time O(n). So 1114 CBR(G; A) can be computed in time $O(n^2m)$ with space $O(n^3 + hm)$, where $n \le m$. Multiplying these complexities by 1115 the number $O(m^h)$ of base sequences gives the final time $O(n^2m^{h+1})$ and space $O(n^3 + hm^{h+1})$ for NCD(G).

The proof of Theorem D.1(b) will use the fact that any point in \mathbb{R}^n is uniquely determined by n + 1 distances to n + 1 ordered points that affinely span \mathbb{R}^n , and also Lemma D.5.

1120 **Lemma D.5** (equal CBRs). Let a Euclidean graph $G \subset \mathbb{R}^n$ have the vertex set V(G) with the center of mass at 1121 $p_0 = 0 \in \mathbb{R}^n$. Let n - 1 ordered vertices p_1, \ldots, p_{n-1} linearly span an (n - 1)-dimensional subspace $S \subset \mathbb{R}^n$. Let 1122 $G(p_1, \ldots, p_{n-1})$ be the subgraph of G on the vertex set V(G) and all edges of G at p_1, \ldots, p_{n-1} . For any other vertex p, let 1123 $\operatorname{CBR}'(G; p_1, \ldots, p_{n-1}, p)$ be obtained from the Centered Base Representation $\operatorname{BR}(G; p_1, \ldots, p_{n-1}, p)$ by removing signs 1124 of distances from all vertices $q \in V(G) \setminus \{p_1, \ldots, p_{n-1}, p\}$ to p. If $\operatorname{BR}'(G; p_1, \ldots, p_{n-1}, p) = \operatorname{BR}'(G; p_1, \ldots, p_{n-1}, p')$ 1125 for some vertices $p, p' \in V(G) \setminus \{p_1, \ldots, p_{n-1}\}$, the mirror reflection with respect to S maps $G(p_1, \ldots, p_{n-1})$ to itself 1126 and p to p'.

1128 **Proof of Lemma** D.5. Under the reflection f_S of \mathbb{R}^n with respect to the subspace $S \subset \mathbb{R}^n$, the vertices p, p' should be swapped because they have equal (signed) distances to the ordered points $p_0, \ldots, p_{n-1} \in S$. The equality of given 1129 CBR's means that $V' = V(G) \setminus \{p_1, \dots, p_{n-1}, p, p'\}$ bijectively maps to itself via $q \mapsto q'$ so that any matched q, q'1130 1131 have the same distances to the n + 1 ordered points p_0, \ldots, p_{n-1}, p as to p_0, \ldots, p_{n-1}, p' , respectively. Any point 1132 in \mathbb{R}^n is determined by its distances to the *n* affinely independent points p_0, \ldots, p_{n-1} up to the mirror reflection f_S . 1133 Since f_S fixes p_0, \ldots, p_{n-1} , the reflection f_S should swap q, q' in such pairs and all their edges, so we conclude that 1134 $f_S(G(p_1,\ldots,p_{n-1})) = G(p_1,\ldots,p_{n-1})$ and $f_S(p) = p'$. 1135

Proof of Theorem D.1(b). The completeness is proved by reconstructing any Euclidean graph $G \subset \mathbb{R}^n$ from NCD(G; n) uniquely up to rigid motion.

We prove that any Euclidean graph $G \subset \mathbb{R}^n$ can be reconstructed from its Nested Distance Distribution NCD(G; n) by induction on the dimension n.

The inductive base n = 1 is Example 4.7. Assume that any graph G on m unordered vertices can be reconstructed in \mathbb{R}^k in time $O(k^3m)$ for any k < n. Below we prove the inductive step for the dimension n > 1. Start from any CBR(G; A) = [sign(A), CD(A), CR(G; A)] from Definition 3.5, where A is a sequence of some n ordered (not yet geometrically fixed) vertices $p_0, \ldots, p_n \in V(G)$. The first point p_0 is fixed at the origin $0 \in \mathbb{R}^n$ as usual by translation.

Lemma D.2(b) for the matrix CD(A) gives the number $k \le n$ of positive eigenvalues of the Gram matrix of the *n* vectors p_1, \ldots, p_n in time $O(n^3)$. If aff(A) = k < n, we use the nested structure of NCD(G; n) to take another CBR($G; p_1, \ldots, p_k, q, \ldots, p_n$) for a new vertex $q \in V(G) - A$. Check if $aff(p_1, \ldots, p_k, q) = k+1$ again by Lemma D.2(b) using the matrix $D(p_1, \ldots, p_k, q)$. If the affine dimension has not increased, we take another CBR with the same points p_1, \ldots, p_k and a new (k + 1)-st point from $V(G) - \{A \cup q\}$ and so on.

This search through Centered Base Representations involving the remaining vertices of G requires a maximum of m - n - 1steps with $O(n^3)$ time for every computation of the affine dimension. Hence in time $O(n^3m)$, we can find a Centered Base 1155 Representation CBR(G; A) whose base sequence A affinely generates the subspace of dimension k = aff(V(G)) in \mathbb{R}^n . If 1156 k < n, the proof follows from the inductive hypothesis for the smaller dimension k. 1157 If aff(V(G)) = n, use the same notations for the fixed vertices $0 = p_0, \ldots, p_n$ that linearly generate \mathbb{R}^n . Lemma D.3(a) 1158 for m = n + 1 and the distance matrix D(A) allow us to reconstruct n + 1 ordered points $0 = p_0, \ldots, p_n$ up to isometry 1159 in \mathbb{R}^n in time $O(n^3)$. By Definition 3.5 every column of CR(G; A) contains Euclidean distances from the vertices 1160 $0 = p_0, \ldots, p_n \in \mathbb{R}^n$, which affinely generate \mathbb{R}^n , to another vertex $q \in V(G) - A$. 1161 1162 These n + 1 distances uniquely determine the position of q in \mathbb{R}^n whose coordinates can be found as follows. Each scalar product $q \cdot p_i$ can be computed as $|q| \cdot |p_i| \cos \angle (q, 0, p_i) = \frac{|q|^2 + |p_i|^2 - |q - p_i|^2}{2}$ for i = 1, ..., n. On another hand, 1163 1164 product $q \cdot p_i$ can be computed as $|q| \cdot |p_i| \cos 2(q, 0, p_i) = \frac{2}{2}$ for i = 1, ..., n. On another hand, $q \cdot p_i$ is a linear combination of unknown coordinates of q with coefficients equal to the coordinates of p_i . One can find all 1165 1166 coordinates of q in time $O(n^3)$ by solving the system of linear equations, where the $n \times n$ determinant on the linear basis 1167 p_1, \ldots, p_n is not zero. The total time is $O(n^3m)$. 1168 1169 Since all vertices $q \in V(G) - A$ are geometrically unique, they can be (arbitrarily) ordered, say p_{n+1}, \ldots, p_m , following p_0, \ldots, p_n . The signs of distances in the matrix CR(G; A) also tell us about (present or absent) edges from p_0, \ldots, p_n to 1170 all other vertices $q \in V(G) - A$. 1171 1172 The nested structure of NCD(G; n) allows us to consider m - n unordered Base Representations $CBR(G; p_1, \ldots, p_{n-1}, p_i)$ 1173 for all vertices p_j with $j = n, \dots, m$. Every vertex $p_j \in V(G)$ is uniquely determined in \mathbb{R}^n by the column of its signed 1174 distances to p_0, \ldots, p_n in the $(n+1) \times (m-n-1)$ matrix $R(G; p_0, \ldots, p_n)$ for $j = n+1, \ldots, m$. 1175 1176 By Lemma D.5, this distance list of p_j (without edges between p_j, p_k for j, k > n) suffices to identify one or maximum two 1177 Base Representations among all m-n unordered CBRs with the fixed n points p_0, \ldots, p_{n-1} and variable n-th vertices. If 1178 there is a choice of two CBRs, we can take any of them for p_j . Indeed, choosing another vertex p_k , which should be mirror 1179 symmetric to p_i , will produce a mirror image of the reconstructed subgraph $G(p_1, \ldots, p_n, p_i)$ by Lemma D.5. 1180 1181 The matrix $CR(G; p_1, \ldots, p_{n-1}, p_i)$ from the found CBRs contains signs that determine the (present or absent) edges from 1182 p_i to all other vertices p_k for $k = n + 1, \ldots, m$. 1183 To guarantee the uniqueness of $G \subset \mathbb{R}^n$ under rigid motion and not only under isometry, we additionally use sign (p_1, \ldots, p_n) 1184 from CBR to fix an orientation of the simplex on p_0, \ldots, p_n . 1185 1186 The strength $\sigma(A)$ depends only on the distance matrix D(A), we write $\sigma(A)$ for brevity. When the simplex on A 1187 degenerates, the strength $\sigma(A)$ vanishes and is Lipschitz continuous by Lemma 4.2, while the volume of the simplex on B 1188 is not Lipschitz continuous as shown below. 1189 1190 In \mathbb{R}^2 , consider the triangle with two vertices fixed at $(\pm l, 0)$ and one moving vertex $(0, t\varepsilon)$ for $t \in [-1, 1]$. The signed area 1191 of the triangle changes from $-l\varepsilon$ (unbounded because l can be large for any fixed small ε) to 0 (when t = 0 and the triangle 1192 degenerates), then to $l\varepsilon$ (when t = 1). The area changes by $2l\varepsilon$ while only one vertex moves by 2ε , so the ratio of the area 1193 change over a point perturbation can be as large as a half-distance between given points. 1194 1195 **Lemma D.6** (time of strength). For any base sequence A of n ordered points $p_1, \ldots, p_n \in \mathbb{R}^n$, the strength $\sigma(A)$ can be 1196 computed in time $O(n^3)$. 1197 1198 **Proof of Lemma** D.6. The half-perimeter p(A) is computable via all pairwise distances in time $O(n^2)$. The squared volume 1199 $vol^2(A)$ can be expressed by the Cayley-Menger $(n+2) \times (n+2)$ determinant from (Sippl & Scheraga, 1986) in inter-point 1200 distances, which can be computed in time $O(n^3)$ by Lemma D.4. \square 1201 **Lemma D.7** (axioms and time of M_{∞} on CBRs). Let $G, F \subset \mathbb{R}^n$ be Euclidean graphs with m unordered vertices and base sequences $A \subset V(G)$ and $B \subset V(F)$ of $h \leq n$ ordered vertices. The metric $M_{\infty}(\text{CBR}(G; A), \text{CBR}(F; B))$ from Definition B.5 satisfies all metric axioms and is computable in time $O(h^2 + m^{1.5} \log^{h+1} m)$ with space $O(h^2 + m \log^{h-1} m)$. 1203 1204 1205 **Proof of Lemma** D.7. The metric axioms for M_{∞} follow from the same axioms for the metrics L_{∞} and W_{∞} because 1206 the maximum of metrics is still a metric, see metric transforms in section 4.1 of (Deza & Deza, 2009). The first metric 1207 $\frac{2}{\lambda_n}|\operatorname{sign}(A)\sigma(A) - \operatorname{sign}(B)\sigma(B)| \text{ can be computed in time } O(n^3) \text{ by Lemma D.6. The metric } L_{\infty}(\operatorname{CD}(A), \operatorname{CD}(B))$ 1208 1209 22

1210 requires time $O(h^2)$ and space $O(h^2)$. The bottleneck distance $W_{\infty}(\operatorname{CR}(G; A)), \operatorname{CR}(F; B))$ between $(h + 1) \times (m - h)$ 1211 matrices $\operatorname{CR}(G; A), \operatorname{CR}(F; B)$ with unordered columns (considered as clouds of m - h unordered points in \mathbb{R}^{h+1}) needs 1212 time $O(m^{1.5} \log^{h+1} m)$ and space $O(m \log^{h-1} m)$ by Theorem 6.5 in (Efrat et al., 2001).

Lemma D.8 (metric axioms for NBM on NCDs). The Nested Bottleneck Metric NBM from Definition B.6 satisfies all
 metric axioms on Nested Distance Distributions.

Proof of Lemma D.8. Induction on the depth i = n, ..., 1. The inductive base i = n follows from the metric axioms in Lemma D.7 for M_{∞} in Definition B.5.

The inductive step from a depth i (between 1, n) to the smaller value i - 1 follows from Lemma C.1 and the metric axioms in the inductive hypothesis for the depth i.

1223 **Lemma D.9** (Lipschitz continuity of M_{∞}). Let A be a base sequence of $1 \le h \le n$ ordered vertices in a Euclidean graph 1224 $G \subset \mathbb{R}^n$. Let B, F be obtained from A, G, respectively, by perturbing every vertex of G within its ε -neighborhood in \mathbb{R}^n . 1225 Then CBR(G; A) changes in M_{∞} from Definition B.5 by at most 2ε , so $M_{\infty}(\text{CBR}(G; A), \text{CBR}(F; B)) \le 2\varepsilon$.

1226

1233

1238

1264

1227 **Proof of Lemma D.9.** Order all vertices of the graphs G, F so that every vertex $p_i \in V(G)$ has the same index as its 1228 perturbation $q_i \in V(F)$. The bijection $p_i \leftrightarrow q_i$ induces the bijections between the corresponding elements of the 1229 matrices $CD(A) \leftrightarrow CD(B)$ and $CR(G; A) \leftrightarrow CR(F; B)$, which all differ by at most 2ε . Lemma 4.2 implies that 1230 $\frac{2}{\lambda_n} |\operatorname{sign}(A)\sigma(A) - \operatorname{sign}(B)\sigma(B)| \le 2\varepsilon$ Since all three components of the max metric M_{∞} in Definition B.5 have the 1231 upper bound 2ε , conclude that $M_{\infty} \le 2\varepsilon$.

1234 Definition C.2 can be adapted for the EMD between NCDs by (1) replacing the bottleneck distance W_{∞} in Definition B.5 1235 with EMD between clouds of equally weighted points, and (2) replacing BMD(Γ) for a bipartite graph Γ with EMD(Γ) 1236 between the unordered sets (of potentially different sizes) of CBDs with weights on all white vertices and CBDs on all 1237 black vertices.

1239 **Proof of Theorem** D.1(c). We first prove the Lipschitz continuity of the metric NBM on NCDs. Order all vertices of the 1240 graphs G, F so that every $p_i \in V(G)$ has the same index as its ε -perturbation $q_i \in V(F)$. In Definition B.6, for any base 1241 sequence A of $p_1, \ldots, p_h \in V(G)$, there is a base sequence B of vertices $q_1, \ldots, q_h \in V(F)$, which are ε -perturbations of 1242 p_1, \ldots, p_h , respectively, such that $M_{\infty}(\text{CBR}(G; A), \text{CBR}(F; B)) \leq 2\varepsilon$ by Lemma D.9.

These distances M_{∞} are weights of edges in the index-preserving vertex matching E of the complete bipartite graph $\Gamma(G; p_1, \ldots, p_{h-1}; F; q_1, \ldots, q_{h-1})$ for any p_1, \ldots, p_{h-1} and their ε -perturbations q_1, \ldots, q_{h-1} . Then BMD($\Gamma(G; p_1, \ldots, p_{h-1}; F; q_1, \ldots, q_{h-1})$) $\leq 2\varepsilon$ by Lemma C.3. Since this conclusion holds for all (choices of) $p_1, \ldots, p_{h-1} \in V(G)$, we iteratively apply this argument for the bipartite graphs $\Gamma(G; p_1, \ldots, p_{i-1}; F; q_1, \ldots, q_{i-1})$ for $1 \leq i < n$ and finally conclude that NBM(NCD(G; h), NCD(F; h)) $\leq 2\varepsilon$. The proof that EMD(NCD(G; h), NCD(F; h)) $\leq 2\varepsilon$ is similar by using Lemma C.4 instead of C.3.

1250 1251 **Proof of Theorem D.1**(d). In Definition B.6, for any fixed $1 \le i \le h$ and ordered vertices $p_1 \ldots, p_{i-1} \in V(G)$ and 1252 $q_1 \ldots, q_{i-1} \in V(F)$, the complete bipartite graph $\Gamma(G; p_1, \ldots, p_{i-1}; F; q_1, \ldots, q_{i-1})$ has V = 2(m - i + 1) = O(m)1253 vertices and $E = (m - i + 1)^2 = O(m^2)$ edges.

For i = h, the weight w(e) of each edge e equals M_{∞} , which needs time $O(m^{1.5} \log^{h+1} m)$ and space $O(m \log^{h-1} m)$ by Lemma D.7 for any $h \le n \le m$. For all $O(m^2)$ edges of $\Gamma(G; p_1, \ldots, p_{h-1}; F; q_1, \ldots, q_{h-1})$, the time is $O(m^{3.5} \log^{h+1} m)$, the space is $O(m^3 \log^{h-1} m)$. The bottleneck matching distance BMD for such a graph is computed by (Hopcroft & Karp, 1973) in time $O(E\sqrt{V}) = O(m^{2.5})$, which is dominated by the time $O(m^{3.5} \log^{h+1} m)$ preparing the weighted graph.

For all $O(m^{2(h-1)})$ choices of ordered vertices $p_1, \ldots, p_{h-1} \in V(G)$ and $q_1, \ldots, q_{h-1} \in V(F)$, the Bottleneck Matching Distance for all graphs $\Gamma(G; p_1, \ldots, p_{h-1}; F; q_1, \ldots, q_{h-1})$ are found in time

 $O(m^{2(h-1)})O(m^{3.5}\log^{h+1}m) = O(m^{2h+1.5}\log^{h+1}m)$

with space $O(m^{2h+1}\log^{h-1}m)$. For every next iteration i = h - 2, ..., 1, the parameter i goes down by 1 every time. We can compute all distances $BMD(\Gamma(G; p_1, ..., p_{i-1}; F; q_1, ..., q_{i-1}))$ in time

$$O(m^{2(i-1)})O(m^{3.5}\log^{h+1}m) = O(m^{2i+1.5}\log^{h+1}m).$$

1270 The sum of all these times for i = 1, ..., h - 1 is still $O(m^{2h+1.5} \log^{h+1} m)$ from the first step.

All CBDs in Definition 3.5 have sizes at most m, which is the maximum number of points in the given clouds. The EMD between weighted distributions of a maximum size m can be computed in near-cubic time $O(m^3 \log m)$, see (Fredman & Tarjan, 1987; Goldberg & Tarjan, 1987). Since this complexity is dominated by the time $O(m^{3.5} \log^{h+1} m)$ for computing $O(m^2)$ weights M_{∞} , each in time $O(m^{1.5} \log^{h+1} m)$ by Lemma D.7, the total time for the EMD is the same as for the NBM, similarly for space complexities