Polynomial Width is Sufficient for Set Representation with High-dimensional Features

Anonymous Author(s) Affiliation Address email

Abstract

Set representation has become ubiquitous in deep learning for modeling the induc-1 tive bias of neural networks that are insensitive to the input order. DeepSets is the 2 most widely used neural network architecture for set representation. It involves 3 embedding each set element into a latent space with dimension L, followed by a 4 sum pooling to obtain a whole-set embedding, and finally mapping the whole-set 5 embedding to the output. In this work, we investigate the impact of the dimension 6 L on the expressive power of DeepSets. Previous analyses either oversimplified 7 high-dimensional features to be one-dimensional features or were limited to ana-8 lytic activations, thereby diverging from practical use or resulting in L that grows 9 exponentially with the set size N and feature dimension D. To investigate the 10 minimal value of L that achieves sufficient expressive power, we present two set-11 12 element embedding layers: (a) linear + power activation (LP) and (b) logarithm + linear + exponential activations (LLE). We demonstrate that L being poly(N, D)13 is sufficient for set representation using both embedding layers. We also provide a 14 lower bound of L for the LP embedding layer. Furthermore, we extend our results 15 to permutation-equivariant set functions and the complex field. 16

17 **1 Introduction**

Enforcing invariance into neural network architectures has become a widely-used principle to design deep learning models [1–7]. In particular, when a task is to learn a function with a set as the input, the architecture enforces permutation invariance that asks the output to be invariant to the permutation of the input set elements [8, 9]. Neural networks to learn a set function have found a variety of applications in particle physics [10, 11], computer vision [12, 13] and population statistics [14–16], and have recently become a fundamental module (the aggregation operation of neighbors' features in a graph [17–19]) in graph neural networks (GNNs) [20, 21] that show even broader applications.

Previous works have studied the expressive power of neural network architectures to represent set functions [8,9,22–26]. Formally, a set with N elements can be represented as $S = \{x^{(1)}, \dots, x^{(N)}\}$ where $x^{(i)}$ is in a feature space \mathcal{X} , typically $\mathcal{X} = \mathbb{R}^D$. To represent a set function that takes S and outputs a real value, the most widely used architecture DeepSets [9] follows Eq. (1).

$$f(\mathcal{S}) = \rho\left(\sum_{i=1}^{N} \phi(\boldsymbol{x}^{(i)})\right), \text{ where } \phi: \mathcal{X} \to \mathbb{R}^{L} \text{ and } \rho: \mathbb{R}^{L} \to \mathbb{R} \text{ are continuous functions.}$$
(1)

²⁹ DeepSets encodes each set element individually via ϕ , and then maps the encoded vectors after sum ³⁰ pooling to the output via ρ . The continuity of ϕ and ρ ensure that they can be well approximated ³¹ by fully-connected neural networks [27, 28], which has practical implications. DeepSets enforces

³² permutation invariance because of the sum pooling, as shuffling the order of $x^{(i)}$ does not change

Table 1: A comprehensive comparison among all prior works on expressiveness analysis with L. Our results achieve the tightest bound on L while being able to analyze high-dimensional set features and extend to the equivariance case.

Prior Arts	L	D > 1	Exact Rep.	Equivariance
DeepSets [9]	D+1	X	✓	✓
Wagstaff et al. [23]	D	×	1	1
Segol et al. [25]	$\binom{N+D}{N} - 1$	1	×	1
Zweig & Bruna [26]	$\exp(\min\{\sqrt{N}, D\})$	1	×	X
Our results	$\operatorname{poly}(N,D)$	\checkmark	\checkmark	\checkmark

the output. However, the sum pooling compresses the whole set into an *L*-dimension vector, which places an information bottleneck in the middle of the architecture. Therefore, a core question on using DeepSets for set function representation is that given the input feature dimension *D* and the set size *N*, what the minimal *L* is needed so that the architecture Eq. (1) can represent/universally approximate any continuous set functions. The question has attracted attention in many previous works [9, 23–26] and is the focus of the present work.

³⁹ An extensive understanding has been achieved for the case with one-dimensional features (D = 1).

⁴⁰ Zaheer et al. [9] proved that this architecture with bottleneck dimension L = N suffices to *accurately*

represent any continuous set functions when D = 1. Later, Wagstaff et al. proved that accurate

⁴² representations cannot be achieved when L < N [23] and further strengthened the statement to a

failure in approximation to arbitrary precision in the infinity norm when L < N [24].

However, for the case with high-dimensional features (D > 1), the characterization of the minimal 44 possible L is still missing. Most of previous works [9, 25, 29] proposed to generate multi-symmetric 45 polynomials to approximate permutation invariant functions [30]. As the algebraic basis of multi-46 symmetric polynomials is of size $L^* = {N+D \choose N} - 1$ [31] (exponential in min $\{D, N\}$), these works by default claim that if $L \ge L^*$, f in Eq. 1 can approximate any continuous set functions, while they do not check the possibility of using a smaller L. Zweig and Bruna [26] constructed a set function that f47 48 49 requires bottleneck dimension $L > N^{-2} \exp(O(\min\{D, \sqrt{N}\}))$ (still exponential in $\min\{D, \sqrt{N}\}$) 50 to approximate while it relies on the condition that ϕ , ρ only adopt analytic activations. This condition 51 is overly strict, as most of the practical neural networks allow using non-analytic activations, such as 52 ReLU. Zweig and Bruna thus left an open question whether the exponential dependence on N or D 53 of L is still necessary if ϕ , ρ allow using non-analytic activations. 54 Present work. The main contribution of this work is to confirm a negative response to the above 55 question. Specifically, we present the first theoretical justification that L being *polynomial* in N and 56 D is sufficient for DeepSets (Eq. (1)) like architecture to represent any *continuous* set functions 57

with *high-dimensional* features (D > 1). To mitigate the gap to the practical use, we consider two architectures to implement feature embedding ϕ (in Eq. 1) and specify the bounds on L accordingly:

• ϕ adopts a linear layer with power mapping: The minimal L holds a lower bound and an upper bound, which is $N(D+1) \leq L < N^5 D^2$.

• Constrained on the entry-wise positive input space $\mathbb{R}_{>0}^{N \times D}$, ϕ adopts *two layers with logarithmic* and exponential activations respectively: The minimal L holds a tighter upper bound $L \leq 2N^2D^2$.

We prove that if the function ρ could be any continuous function, the above two architectures 64 65 reproduce the precise construction of any set functions for high-dimensional features D > 1, akin to the result in [9] for D = 1. This result contrasts with [25, 26] which only present approximating 66 representations. If ρ adopts a fully-connected neural network that allows approximation of any 67 continuous functions on a bounded input space [27, 28], then the DeepSets architecture $f(\cdot)$ can 68 approximate any set functions universally on that bounded input space. Moreover, our theory can be 69 easily extended to permutation-equivariant functions and complex set functions, where the minimal 70 L shares the same bounds up to some multiplicative constants. 71

Another comment on our contributions is that Zweig and Bruna [26] use difference in the needed
dimension *L* to illustrate the gap between DeepSets [9] and Relational Network [32] in their expressive
powers, where the latter encodes set elements in a pairwise manner rather than in a separate manner.
The gap well explains the empirical observation that Relational Network achieves better expressive
power with smaller *L* [23,33]. Our theory does not violate such an observation while it shows that the

gap can be reduced from an exponential order in N and D to a polynomial order. Moreover, many real-77 world applications have computation constraints where only DeepSets instead of Relational Network 78 can be used, e.g., the neighbor aggregation operation in GNN being applied to large networks [21], 79 and the hypergraph neural diffusion operation in hypergraph neural networks [7]. Our theory points 80 out that in this case, it is sufficient to use polynomial L dimension to embed each element, while one 81 needs to adopt a function ρ with non-analytic activitions. 82

2 **Preliminaries** 83

2.1 Notations and Problem Setup 84

We are interested in the approximation and representation of functions defined over sets ¹. In 85 convention, an *N*-sized set $S = \{ \boldsymbol{x}^{(1)}, \cdots, \boldsymbol{x}^{(N)} \}$, where $\boldsymbol{x}^{(i)} \in \mathbb{R}^D, \forall i \in [N] (\triangleq \{1, 2, ..., N\})$, can be denoted by a data matrix $\boldsymbol{X} = \begin{bmatrix} \boldsymbol{x}^{(1)} & \cdots & \boldsymbol{x}^{(N)} \end{bmatrix}^\top \in \mathbb{R}^{N \times D}$. Note that we use the 86 87 superscript (i) to denote the *i*-th set element and the subscript *i* to denote the *i*-th column/feature 88 channel of \boldsymbol{X} , i.e., $\boldsymbol{x}_i = \begin{bmatrix} x_i^{(1)} & \cdots & x_i^{(N)} \end{bmatrix}^\top$. Let $\Pi(N)$ denote the set of all N-by-N permutation 89 matrices. To characterize the unorderedness of a set, we define an equivalence class over $\mathbb{R}^{N \times D}$:

90

Definition 2.1 (Equivalence Class). If matrices $X, X' \in \mathbb{R}^{N \times D}$ represent the same set \mathcal{X} , then they 91 are called equivalent up a row permutation, denoted as $X \sim X'$. Or equivalently, $X \sim X'$ if and 92 only if there exists a matrix $P \in \Pi(N)$ such that X = PX'. 93

Set functions can be in general considered as permutation-invariant or permutation-equivariant 94 functions, which process the input matrices regardless of the order by which rows are organized. The 95

formal definitions of permutation-invariant/equivariant functions are presented as below: 96

Definition 2.2. (Permutation Invariance) A function $f : \mathbb{R}^{N \times D} \to \mathbb{R}^{D'}$ is called permutation-97 invariant if $f(\mathbf{PX}) = f(\mathbf{X})$ for any $\mathbf{P} \in \Pi(N)$. 98

Definition 2.3. (Permutation Equivariance) A function $f : \mathbb{R}^{N \times D} \to \mathbb{R}^{N \times D'}$ is called permutation-99 equivariant if $f(\mathbf{P}\mathbf{X}) = \mathbf{P}f(\mathbf{X})$ for any $\mathbf{P} \in \Pi(N)$. 100

In this paper, we investigate the approach to design a neural network architecture with permutation in-101 variance/equivariance. Below we will first focus on permutation-invariant functions $f: \mathbb{R}^{N \times D} \to \mathbb{R}$. 102 Then, in Sec. 5, we show that we can easily extend the established results to permutation-equivariant 103 functions through the results provided in [7, 34] and to the complex field. The obtained results for 104 D' = 1 can also be easily extended to D' > 1 as otherwise f can be written as $[f_1 \cdots f_{D'}]^{\top}$ and 105 each f_i has single output feature channel. 106

2.2 DeepSets and The Difficulty in the High-Dimensional Case D > 1107

The seminal work [9] establishes the following result which induces a neural network architecture for 108 permutation-invariant functions. 109

Theorem 2.4 (DeepSets [9], D = 1). A continuous function $f : \mathbb{R}^N \to \mathbb{R}$ is permutation-invariant (*i.e.*, a set function) if and only if there exists continuous functions $\phi : \mathbb{R} \to \mathbb{R}^L$ and $\rho : \mathbb{R}^L \to \mathbb{R}$ 110 111

such that $f(\mathbf{X}) = \rho\left(\sum_{i=1}^{N} \phi(x^{(i)})\right)$, where L can be as small as N. Note that, here $x^{(i)} \in \mathbb{R}$. 112

Remark 2.5. The original result presented in [9] states the latent dimension should be as large as 113 N + 1. [23] tighten this dimension to exactly N. 114

Theorem 2.4 implies that as long as the latent space dimension L > N, any permutation-invariant 115

functions can be implemented by a unified manner as DeepSets (Eq.(1)). Furthermore, DeepSets 116

suggests a useful architecture for ϕ at the analysis convenience and empirical utility, which is formally 117

defined below ($\phi = \psi_L$): 118

Definition 2.6 (Power mapping). A power mapping of degree K is a function $\psi_K : \mathbb{R} \to \mathbb{R}^K$ which 119 transforms a scalar to a power series: $\psi_K(z) = \begin{bmatrix} z & z^2 & \cdots & z^K \end{bmatrix}^\top$. 120

¹In fact, we allow repeating elements in S, therefore, S should be more precisely called multiset. With a slight abuse of terminology, we interchangeably use terms multiset and set throughout the whole paper.



Figure 1: Illustration of the proposed linear + power mapping embedding layer (LP) and logarithm activation + linear + exponential activation embedding layer (LLE).

However, DeepSets [9] focuses on the case that the feature dimension of each set element is one (i.e., D = 1). To demonstrate the difficulty extending Theorem 2.4 to high-dimensional features, we reproduce the proof next, which simultaneously reveals its significance and limitation. Some intermediate results and mathematical tools will be recalled along the way later in our proof.

We begin by defining sum-of-power mapping (of degree K) $\Psi_K(\mathbf{X}) = \sum_{i=1}^N \psi_K(x_i)$, where ψ_K is the power mapping following Definition 2.6. Afterwards, we reveal that sum-of-power mapping $\Psi_K(\mathbf{X})$ has a continuous inverse. Before stating the formal argument, we formally define the injectivity of permutation-invariant mappings:

Definition 2.7 (Injectivity). A set function $h : \mathbb{R}^{N \times D} \to \mathbb{R}^{L}$ is injective if there exists a function $g : \mathbb{R}^{L} \to \mathbb{R}^{N \times D}$ such that for any $X \in \mathbb{R}^{N \times D}$, we have $g \circ f(X) \sim X$. Then g is an inverse of f.

And we summarize the existence of continuous inverse of $\Psi_K(x)$ into the following lemma shown by [9] and improved by [23]. This result comes from homeomorphism between roots and coefficients of monic polynomials [35].

Lemma 2.8 (Existence of Continuous Inverse of Sum-of-Power [9,23]). $\Psi_N : \mathbb{R}^N \to \mathbb{R}^N$ is injective, thus the inverse $\Psi_N^{-1} : \mathbb{R}^N \to \mathbb{R}^N$ exists. Moreover, Ψ_N^{-1} is continuous.

Now we are ready to prove necessity in Theorem 2.4 as sufficiency is easy to check. By choosing $\phi = \psi_N : \mathbb{R} \to \mathbb{R}^N$ to be the power mapping (cf. Definition 2.6), and $\rho = f \circ \Psi_N^{-1}$. For any scalar-

138 valued set
$$\boldsymbol{X} = \begin{bmatrix} x^{(1)} & \cdots & x^{(N)} \end{bmatrix}^{\top}$$
, $\rho \left(\sum_{i=1}^{N} \phi(x^{(i)}) \right) = f \circ \Psi_{N}^{-1} \circ \Psi_{N}(\boldsymbol{x}) = f(\boldsymbol{P}\boldsymbol{X}) = f(\boldsymbol{X})$

for some $P \in \Pi(N)$. The existence and continuity of Ψ_N^{-1} are due to Lemma 2.8.

Theorem 2.4 gives the *exact decomposable form* [23] for permutation-invariant functions, which is stricter than approximation error based expressiveness analysis. In summary, the key idea is to establish a mapping ϕ whose element-wise sum-pooling has a continuous inverse.

Curse of High-dimensional Features. We argue that the proof of Theorem 2.4 is not applicable 143 to high-dimensional set features $(D \ge 2)$. The main reason is that power mapping defined in 144 Definition 2.6 only receives scalar input. It remains elusive how to extend it to a multivariate version 145 that admits injectivity and a continuous inverse. A plausible idea seems to be applying power mapping 146 for each channel x_i independently, and due to the injectivity of sum-of-power mapping Ψ_N , each 147 channel can be uniquely recovered individually via the inverse Ψ_N^{-1} . However, we point out that each recovered feature channel $x'_i \sim x_i$, $\forall i \in [D]$, does not imply $[x'_1 \cdots x'_D] \sim X$, where 148 149 the alignment of features across channels gets lost. Hence, channel-wise power encoding no more 150 composes an injective mapping. Zaheer et al. [9] proposed to adopt multivariate polynomials as ϕ for 151 high-dimensional case, which leverages the fact that multivariate symmetric polynomials are dense in 152 the space of permutation invariant functions (akin to Stone-Wasserstein theorem) [30]. This idea later 153 got formalized in [25] by setting $\phi(\boldsymbol{x}^{(i)}) = \begin{bmatrix} \cdots & \prod_{j \in [D]} (x_j^{(i)})^{\alpha_j} & \cdots \end{bmatrix}$ where $\boldsymbol{\alpha} \in \mathbb{N}^D$ traverses all $\sum_{j \in [D]} \alpha_j \leq n$ and extended to permutation equivariant functions. Nevertheless, the dimension 154 155 $L = {\binom{N+D}{D}}$, i.e., exponential in min $\{N, D\}$ in this case, and unlike DeepSets [9] which exactly recovers f for D = 1, the architecture in [9, 25] can only approximate the desired function. 156 157

158 **3** Main Results

¹⁵⁹ In this section, we present our main result which extends Theorem 2.4 to high-dimensional features. ¹⁶⁰ Our conclusion is that to universally represent a set function on sets of length *N* and feature dimension

- D with the DeepSets architecture [9] (Eq. (1)), a dimension L at most polynomial in N and D is 161 needed for expressing the intermediate embedding space. 162
- Formally, we summarize our main result in the following theorem. 163

Theorem 3.1 (The main result). Suppose $D \ge 2$. For any continuous permutation-invariant function 164 $f: \mathcal{K}^{N \times D} \to \mathbb{R}, \ \mathcal{K} \subseteq \mathbb{R}, \ there \ exists \ two \ continuous \ mappings \ \phi: \mathbb{R}^D \to \mathbb{R}^L \ and \ \rho: \mathbb{R}^L \to \mathbb{R}$ such that for every $\mathbf{X} \in \mathcal{K}^{N \times D}, \ f(\mathbf{X}) = \rho\left(\sum_{i=1}^N \phi(\mathbf{x}^{(i)})\right)$ where 165

166

• For some $L \in [N(D+1), N^5D^2]$ when ϕ admits **linear layer + power mapping (LP)** architecture: 167

$$\phi(\boldsymbol{x}) = \begin{bmatrix} \psi_N(\boldsymbol{w}_1 \boldsymbol{x})^\top & \cdots & \psi_N(\boldsymbol{w}_K \boldsymbol{x})^\top \end{bmatrix}$$
(2)

for some $w_1, \dots, w_K \in \mathbb{R}^D$, and K = L/N. 168

• For some $L \in [ND, 2N^2D^2]$ when ϕ admits logarithm activations + linear layer + exponential 169 activations (LLE) architecture: 170

$$\phi(\boldsymbol{x}) = [\exp(\boldsymbol{w}_1 \log(\boldsymbol{x})) \quad \cdots \quad \exp(\boldsymbol{w}_L \log(\boldsymbol{x}))]$$
(3)

for some
$$w_1, \cdots, w_L \in \mathbb{R}^D$$
 and $\mathcal{K} \subseteq \mathbb{R}_{>0}$.

The bounds of L depend on the choice of the architecture of ϕ , which are illustrated in Fig. 1. In 172 the LP setting, we adopt a linear layer that maps each set element into K dimension. Then we apply 173 a channel-wise power mapping that separately transforms each value in the feature vector into an 174 N-order power series, and concatenates all the activations together, resulting in a KN dimension 175 feature. The LP architecture is closer to DeepSets [9] as they share the power mapping as the main 176 component. Theorem 3.1 guarantees the existence of ρ and ϕ (in the form of Eq. (2)) which satisfy 177 Eq. (1) without the need to set K larger than N^4D^2 while $K \ge D + 1$ is necessary. Therefore, the 178 total embedding size L = KN is bounded by N^5D^2 above and N(D+1) below. Note that this 179 lower bound is not trivial as ND is the degree of freedom of the input X. No matter how $w_1, ..., w_K$ 180 are adopted, one cannot achieve an injective mapping by just using ND dimension. 181

In the LLE architecture, we investigate the utilization of logarithmic and exponential activations in set 182 representation, which are also valid activations to build deep neural networks [36, 37]. Each set entry 183 will be squashed by a element-wise logarithm first, then linearly embedded into an L-dimensional 184 space via a group of weights, and finally transformed by an element-wise exponential activation. 185 Essentially, each $\exp(w_i \log(x)), i \in [L]$ gives a monomial of x. The LLE architecture requires the 186 feature space constrained on the positive orthant to ensure logarithmic operations are feasible. But 187 the advantage is that the upper bound of L is improved to be $2N^2D^2$. The lower bound ND for 188 the LLE architecture is a trivial bound due to the degree of freedom of the input X. Note that the 189 constraint on the positive orthant $\mathbb{R}_{>0}$ is not essential. If we are able to use monomial activations to 190 process a vector x as used in [25, 26], then, the constraint on the positive orthant can be removed. 191

Remark 3.2. The bounds in Theorem 3.1 are non-asymptotic. This implies the latent dimensions 192 specified by the corresponding architectures are precisely sufficient for expressing the input. 193

Remark 3.3. Unlike ϕ , the form of ρ cannot be explicitly specified, as it depends on the desired 194

function f. The complexity of ρ remains unexplored in this paper, which may be high in practice. 195

Importance of Continuity. We argue that the requirements of continuity on ρ and ϕ are essential 196 for our discussion. First, practical neural networks can only provably approximate continuous 197 functions [27,28]. Moreover, set representation without such requirements can be straightforward 198 (but likely meaningless in practice). This is due to the following lemma. 199

Lemma 3.4 ([38]). There exists a discontinuous bijective mapping between \mathbb{R}^D and \mathbb{R} if $D \geq 2$. 200

By Lemma 3.4, we can define a bijective mapping $r: \mathbb{R}^D \to \mathbb{R}$ which maps the high-dimensional 201 features to scalars, and its inverse exists. Then, the same proof of Theorem 2.4 goes through by letting $\phi = \psi_N \circ r$ and $\rho = f \circ r^{-1} \circ \Psi_N^{-1}$. However, we note both ρ and ϕ lose continuity. 202 203

Comparison with Prior Arts. Below we highlight the significance of Theorem 3.1 in contrast 204 to the existing literature. A quick overview is listed in Tab. 1 for illustration. The lower bound 205 in Theorem 3.1 corrects a natural misconception that the degree of freedom (i.e., L = ND for 206

multi-channel cases) is not enough for representing the embedding space. Fortunately, the upper bound in Theorem 3.1 shows the complexity of representing vector-valued sets is still manageable as it merely scales polynomially in N and D. Compared with Zweig and Bruna's finding [26], our result significantly improves this bound on L from exponential to polynomial by allowing non-analytic functions to amortize the expressiveness. Besides, Zweig and Bruna's work [26] is hard to be applied to the real domain, while ours are extensible to complex numbers and equivariant functions.

213 4 Proof Sketch

In this section, we introduce the proof techniques of Theorem 3.1, while deferring a full version and all missing proofs to the supplementary materials.

²¹⁶ The proof of Theorem 3.1 mainly consists of two steps below, which is completely constructive:

1. For the LP architecture, we construct a group of K linear weights $w_1 \cdots, w_K$ with $K \leq N^4 D^2$

such that the summation over the associated LP embedding (Eq. (2)): $\Psi(\mathbf{X}) = \sum_{i=1}^{N} \phi(\mathbf{x}^{(i)})$ is injective and has a continuous inverse. Moreover, if $K \leq D$, such weights do not exist, which

injective and has a continuous inverse. Moreover, if $K \le D$, such weights do not exist, which induces the lower bound.

221 2. Similarly, for the LLE architecture, we construct a group of L linear weights $w_1 \cdots, w_L$ with 222 $L \leq 2N^2D^2$ such that the summation over the associated LLE embedding (Eq. (3)) is injective 223 and has a continuous inverse. Trivially, if L < ND, such weights do not exist, which induces the 224 lower bound.

225 3. Then the proof of upper bounds can be concluded for both settings by letting $\rho = f \circ \Psi^{-1}$ since 226 $\rho\left(\sum_{i=1}^{N} \phi(\boldsymbol{x}^{(i)})\right) = f \circ \Psi^{-1} \circ \Psi(\boldsymbol{X}) = f(\boldsymbol{P}\boldsymbol{X}) = f(\boldsymbol{X})$ for some $\boldsymbol{P} \in \Pi(N)$.

227 Next, we elaborate on the construction idea which yields injectivity for both embedding layers in Sec.

4.1 and 4.2, respectively. To show injectivity, it is equivalent to establish the following statement for both Eq. (2) and Eq. (3), respectively:

$$\forall \boldsymbol{X}, \boldsymbol{X'} \in \mathbb{R}^{N \times D}, \sum_{i=1}^{N} \phi(\boldsymbol{x}^{(i)}) = \sum_{i=1}^{N} \phi(\boldsymbol{x'}^{(i)}) \Rightarrow \boldsymbol{X} \sim \boldsymbol{X'}$$
(4)

In Sec. 4.3, we prove the continuity of the inverse map for LP and LLE via arguments similar to [35].

231 4.1 Injectivity of LP

In this section, we consider ϕ follows the definition in Eq. (2), which amounts to first linearly transforming each set element and then applying channel-wise power mapping. This is, we seek a group of linear transformations w_1, \dots, w_K such that $X \sim X'$ can be induced from $Xw_i \sim$ $X'w_i, \forall i \in [K]$ for some K larger than N while being polynomial in N and D. The intuition is that linear mixing among each channel can encode relative positional information. Only if $X \sim X'$, the mixing information can be reproduced.

²³⁸ Formally, the first step accords to the property of power mapping (cf. Lemma 2.8), and we can obtain:

$$\sum_{i=1}^{N} \phi(\boldsymbol{x}^{(i)}) = \sum_{i=1}^{N} \phi(\boldsymbol{x'}^{(i)}) \Rightarrow \boldsymbol{X}\boldsymbol{w}_{i} \sim \boldsymbol{X'}\boldsymbol{w}_{i}, \forall i \in [K].$$
(5)

To induce $X \sim X'$ from $Xw_i \sim X'w_i$, $\forall i \in [K]$, our construction divides the weights $\{w_i, i \in [K]\}$ into three groups: $\{w_i^{(1)} : i \in [D]\}, \{w_j^{(2)} : j \in [K_1]\}$, and $\{w_{i,j,k}^{(3)} : i \in [D], j \in [K_1], k \in [K_2]\}$. Each block is outlined as below:

1. Let the first group of weights $w_1^{(1)} = e_1, \cdots, w_D^{(1)} = e_D$ to buffer the original features, where e_i is the *i*-th canonical basis.

244 2. Design the second group of linear weights, $w_1^{(2)}, \dots, w_{K_1}^{(2)}$ for K_1 as large as N(N-1)(D-1)

245 1)/2 + 1, which, by Lemma 4.4 latter, guarantees at least one of $Xw_j^{(2)}, j \in [K_1]$ forms an 246 anchor defined below:

Definition 4.1 (Anchor). Consider the data matrix $X \in \mathbb{R}^{N \times D}$, then $a \in \mathbb{R}^N$ is called an anchor of X if $a_i \neq a_j$ for any $i, j \in [N]$ such that $x^{(i)} \neq x^{(j)}$. 247 248

And suppose $a = X w_{j^*}^{(2)}$ is an anchor of X for some $j^* \in [K_1]$ and $a' = X' w_{j^*}^{(2)}$, then we show the following statement is true by Lemma 4.3 latter: 249 250

$$\begin{bmatrix} \boldsymbol{a} & \boldsymbol{x}_i \end{bmatrix} \sim \begin{bmatrix} \boldsymbol{a'} & \boldsymbol{x'}_i \end{bmatrix}, \forall i \in [D] \Rightarrow \boldsymbol{X} \sim \boldsymbol{X'}.$$
(6)

3. Design a group of weights $\boldsymbol{w}_{i,j,k}^{(3)}$ for $i \in [D], j \in [K_1], k \in [K_2]$ with $K_2 = N(N-1) + 1$ that mixes each original channel \boldsymbol{x}_i with each $\boldsymbol{X}\boldsymbol{w}_j^{(2)}, j \in [K_1]$ by $\boldsymbol{w}_{i,j,k}^{(3)} = \boldsymbol{e}_i - \gamma_k \boldsymbol{w}_j^{(2)}$. Then we 251

252 show in Lemma 4.5 that: 253

$$\boldsymbol{X}\boldsymbol{w}_{i} \sim \boldsymbol{X}'\boldsymbol{w}_{i}, \forall i \in [K] \Rightarrow \begin{bmatrix} \boldsymbol{X}\boldsymbol{w}_{j}^{(2)} & \boldsymbol{x}_{i} \end{bmatrix} \sim \begin{bmatrix} \boldsymbol{X}'\boldsymbol{w}_{j}^{(2)} & \boldsymbol{x}'_{i} \end{bmatrix}, \forall i \in [D], j \in [K_{1}] \quad (7)$$

With such configuration, injectivity can be concluded by the entailment along Eq. (5), (7), (6): Eq. (5) 254 guarantees the RHS of Eq. (7); The existence of the anchor in Lemma 4.4 paired with Eq. (6) guarantees $X \sim X'$. The total required number of weights $K = D + K_1 + DK_1K_2 \le N^4D^2$. 255 256

Below we provides a series of lemmas that demonstrate the desirable properties of anchors and 257 elaborate on the construction complexity. Detailed proofs are left in Appendix. In plain language, by 258 Definition 4.1, two entries in the anchor must be distinctive if the set elements at the corresponding 259 indices are not equal. As a consequence, we derive the following property of anchors: 260

Lemma 4.2. Consider the data matrix $X \in \mathbb{R}^{N \times D}$ and $a \in \mathbb{R}^N$ an anchor of X. Then if there exists $P \in \Pi(N)$ such that Pa = a then $Px_i = x_i$ for every $i \in [D]$. 261 262

With the above property, anchors defined in Definition 4.1 indeed have the entailment in Eq. (6): 263

Lemma 4.3 (Union Alignment based on Anchor Alignment). Consider the data matrix $\mathbf{X}, \mathbf{X}' \in \mathbb{R}^{N \times D}$, $\mathbf{a} \in \mathbb{R}^N$ is an anchor of \mathbf{X} and $\mathbf{a}' \in \mathbb{R}^N$ is an arbitrary vector. If $[\mathbf{a} \quad \mathbf{x}_i] \sim [\mathbf{a}' \quad \mathbf{x}'_i]$ for every $i \in [D]$, then $\mathbf{X} \sim \mathbf{X}'$. 264 265 266

However, the anchor a is required to be generated from X via a point-wise linear transformation. 267 The strategy to generate an anchor is to enumerate as many linear weights as needs, so that for any X, at least one j such that $Xw_j^{(2)}$ becomes an anchor. We show that at most N(N-1)(D-1)/2 + 1 linear weights are enough to guarantee the existence of an anchor for any X: 268 269 270

Lemma 4.4 (Anchor Construction). There exists a set of weights w_1, \dots, w_{K_1} where $K_1 = N(N-1)(D-1)/2 + 1$ such that for every data matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$, there exists $j \in [K_1]$, $\mathbf{X}w_j$ 271 272 is an anchor of X. 273

We wrap off the proof by presenting the following lemma which is applied to prove Eq. (7) by fixing 274 any $i \in [D], j \in [K_1]$ in Eq. (7) while checking the condition for all $k \in [K_2]$: 275

276

Lemma 4.5 (Anchor Matching). There exists a group of coefficients $\gamma_1, \dots, \gamma_{K_2}$ where $K_2 = N(N-1) + 1$ such that the following statement holds: Given any $\mathbf{x}, \mathbf{x}', \mathbf{y}, \mathbf{y}' \in \mathbb{R}^N$ such that $\mathbf{x} \sim \mathbf{x}'$ and $\mathbf{y} \sim \mathbf{y}'$, if $(\mathbf{x} - \gamma_k \mathbf{y}) \sim (\mathbf{x}' - \gamma_k \mathbf{y}')$ for every $k \in [K_2]$, then $[\mathbf{x} \ \mathbf{y}] \sim [\mathbf{x}' \ \mathbf{y}']$. 277

278

For completeness, we add the following lemma which implies LP-induced sum-pooling cannot be 279 injective if $K \leq ND$, when $D \geq 2$. 280

Theorem 4.6 (Lower Bound). Consider data matrices $\mathbf{X} \in \mathbb{R}^{N \times D}$ where $D \ge 2$. If $K \le D$, then for every $\mathbf{w}_1, \dots, \mathbf{w}_K$, there exists $\mathbf{X'} \in \mathbb{R}^{N \times D}$ such that $\mathbf{X} \not\sim \mathbf{X'}$ but $\mathbf{X}\mathbf{w}_i \sim \mathbf{X'}\mathbf{w}_i$ for every 281 282 $i \in [K].$ 283

Remark 4.7. Theorem 4.6 is significant in that with high-dimensional features, the injectivity is 284 provably not satisfied when the embedding space has dimension equal to the degree of freedom. 285

4.2 Injectivity of LLE 286

In this section, we consider ϕ follows the definition in Eq. (3). First of all, we note that each term in 287 the RHS of Eq. (3) can be rewritten as a monomial as shown in Eq. (8). Suppose we are able to use 288 monomial activations to process a vector $x^{(i)}$. Then, the constraint on the positive orthant $\mathbb{R}_{>0}$ in our 289 main result Theorem 3.1 can be even removed. 290

$$\phi(\boldsymbol{x}) = \begin{bmatrix} \cdots & \exp(\boldsymbol{w}_i \log(\boldsymbol{x})) & \cdots \end{bmatrix} = \begin{bmatrix} \cdots & \prod_{j=1}^{D} \boldsymbol{x}_j^{\boldsymbol{w}_{i,j}} & \cdots \end{bmatrix}$$
(8)

- Then, the assignment of w_1, \dots, w_L amounts to specifying the exponents for D power functions within the product. Next, we prepare our construction with the following two lemmas:
- **Lemma 4.8.** For any pair of vectors $\mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^N, \mathbf{y}_1, \mathbf{y}_2 \in \mathbb{R}^N$, if $\sum_{i \in [N]} \mathbf{x}_{1,i}^{l-k} \mathbf{x}_{2,i}^k = \sum_{i \in [N]} \mathbf{y}_{1,i}^{l-k} \mathbf{y}_{2,i}^k$ for every $l, k \in [N]$ such that $0 \le k \le l$, then $[\mathbf{x}_1 \ \mathbf{x}_2] \sim [\mathbf{y}_1 \ \mathbf{y}_2]$.

The above lemma is to show that we may use summations of monic bivariate monomials to align every two feature columns. The next lemma shows that such pairwise alignment yields union alignment.

Lemma 4.9 (Union Alignment based on Pairwise Alignment). Consider data matrices $X, X' \in \mathbb{R}^{N \times D}$. If $[x_i \quad x_j] \sim [x'_i \quad x'_j]$ for every $i, j \in [D]$, then $X \sim X'$.

Then the construction idea of w_1, \dots, w_L can be drawn from Lemma 4.8 and 4.9:

1. Lemma 4.8 indicates if the weights in Eq. (8) enumerate all the monic bivariate monomials in each pair of channels with degrees less or equal to N, i.e., $\boldsymbol{x}_{i}^{p}\boldsymbol{x}_{j}^{q}$ for all $i, j \in [D]$ and $p + q \leq N$, then we can yield:

$$\sum_{i=1}^{N} \phi(\boldsymbol{x}^{(i)}) = \sum_{i=1}^{N} \phi(\boldsymbol{x'}^{(i)}) \Rightarrow [\boldsymbol{x}_{i} \quad \boldsymbol{x}_{j}] \sim [\boldsymbol{x'}_{i} \quad \boldsymbol{x'}_{j}], \forall i, j \in [D].$$
(9)

The next step is to invoke Lemma 4.9 which implies if every pair of feature channels is aligned,
 then we can conclude all the channels are aligned with each other as well.

$$\begin{bmatrix} \boldsymbol{x}_i & \boldsymbol{x}_j \end{bmatrix} \sim \begin{bmatrix} \boldsymbol{x'}_i & \boldsymbol{x'}_j \end{bmatrix}, \forall i, j \in [D] \Rightarrow \boldsymbol{X} \sim \boldsymbol{X'}.$$
(10)

Based on these motivations, we assign the weights that induce all bivariate monic monomials with the degree no more than N. First of all, we reindex $\{w_i, i \in [L]\}$ as $\{w_{i,j,p,q}, i \in [D], j \in [D], p \in [N], q \in [p+1]\}$. Then weights can be explicitly specified as $w_{i,j,p,q} = (q-1)e_i + (p-q+1)e_j$, where e_i is the *i*-th canonical basis. With such weights, injectivity can be concluded by entailment along Eq. (9) and (10). Moreover, the total number of linear weights is $L = D^2(N+3)N/2 \le 2N^2D^2$, as desired.

311 4.3 Continuous Lemma

In this section, we show that the LP and LLE induced sum-pooling are both homeomorphic. We note that it is intractable to obtain the closed form of their inverse maps. Notably, the following remarkable result can get rid of inversing a functions explicitly by merely examining the topological relationship between the domain and image space.

Lemma 4.10. (*Theorem 1.2 [35]*) Let $(\mathcal{X}, d_{\mathcal{X}})$ and $(\mathcal{Y}, d_{\mathcal{Y}})$ be two metric spaces and $f : \mathcal{X} \to \mathcal{Y}$ is a bijection such that (**a**) each bounded and closed subset of \mathcal{X} is compact, (**b**) f is continuous, (**c**) f^{-1} maps each bounded set in \mathcal{Y} into a bounded set in \mathcal{X} . Then f^{-1} is continuous.

Subsequently, we show the continuity in an informal but more intuitive way while deferring a rigorous version to the supplementary materials. Denote $\Psi(\mathbf{X}) = \sum_{i \in [N]} \phi(\mathbf{x}^{(i)})$. To begin with, we set $\mathcal{X} = \mathbb{R}^{N \times D} / \sim$ with metric $d_{\mathcal{X}}(\mathbf{X}, \mathbf{X}') = \min_{\mathbf{P} \in \Pi(N)} ||\mathbf{X} - \mathbf{P}\mathbf{X}'||_1$ and $\mathcal{Y} = \{\Psi(\mathbf{X}) | \mathbf{X} \in \mathcal{X}\} \subseteq \mathbb{R}^L$ with metric $d_{\mathcal{Y}}(\mathbf{y}, \mathbf{y}') = ||\mathbf{y} - \mathbf{y}'||_{\infty}$. It is easy to show that \mathcal{X} satisfies the conditions (a) and $\Psi(\mathbf{X})$ satisfies (b) for both LP and LLE embedding layers. Then it remains to conclude the proof by verifying the condition (c) for the mapping $\mathcal{Y} \to \mathcal{X}$, i.e., the inverse of $\Psi(\mathbf{X})$. We visualize this mapping following our arguments on injectivity:

$$\begin{array}{cccc} (LP) & \Psi(\mathbf{X}) & \xrightarrow{\text{Eq. (5)}} & [\cdots & \mathbf{P}_i \mathbf{X} \mathbf{w}_i & \cdots], i \in [K] & \xrightarrow{\text{Eqs. (6)} + (I)} & \mathbf{P} \mathbf{X} \\ (LLE) & \underbrace{\Psi(\mathbf{X})}_{\mathcal{Y}} & \xrightarrow{\text{Eq. (9)}} & \underbrace{[\cdots & \mathbf{Q}_{i,j} \mathbf{x}_i & \mathbf{Q}_{i,j} \mathbf{x}_j & \cdots], i, j \in [D]}_{\mathcal{Z}} & \xrightarrow{\text{Eqs. (6)} + (I)} & \underbrace{\mathbf{P} \mathbf{X}}_{\mathcal{X}} & , \end{array}$$

for some X dependent P, Q. Here, $P_i, i \in [K]$ and $Q_{i,j}, i, j \in [D] \in \Pi(N)$. According to homeomorphism between polynomial coefficients and roots (Corollary 3.2 in [35]), any bounded set in \mathcal{Y} will induce a bound set in \mathcal{Z} . Moreover, since elements in \mathcal{Z} contains all the columns of \mathcal{X} (up to some changes of the entry orders), a bounded set in \mathcal{Z} also corresponds to a bounded set in \mathcal{X} . Through this line of arguments, we conclude the proof.

331 5 Extensions

In this section, we discuss two extensions to Theorem 3.1, which strengthen our main result.

Permutation Equivariance. Permutation-equivariant functions (cf. Definition 2.3) are considered as a more general family of set functions. Our main result does not lose generality to this class of functions. By Lemma 2 of [7], Theorem 3.1 can be directly extended to permutation-equivariant functions with *the same lower and upper bounds*, stated as follows:

Theorem 5.1 (Extension to Equivariance). For any permutation-equivariant function $f : \mathcal{K}^{N \times D} \rightarrow \mathbb{R}^N$, $\mathcal{K} \subseteq \mathbb{R}$, there exists continuous functions $\phi : \mathbb{R}^D \to \mathbb{R}^L$ and $\rho : \mathbb{R}^D \times \mathbb{R}^L \to \mathbb{R}$ such that $f(\mathbf{X})_j = \rho\left(\mathbf{x}^{(j)}, \sum_{i \in [N]} \phi(\mathbf{x}^{(i)})\right)$ for every $j \in [N]$, where $L \in [N(D+1), N^5D^2]$ when ϕ admits LP architecture, and $L \in [ND, 2N^2D^2]$ when ϕ admits LLE architecture ($\mathcal{K} \in \mathbb{R}_{>0}$).

Complex Domain. The upper bounds in Theorem 3.1 is also true to complex features up to a constant scale (i.e., $\mathcal{K} \subseteq \mathbb{C}$). When features are defined over $\mathbb{C}^{N \times D}$, our primary idea is to divide each channel into two real feature vectors, and recall Theorem 3.1 to conclude the arguments on an $\mathbb{R}^{N \times 2D}$ input. All of our proof strategies are still applied. This result directly contrasts to Zweig and Bruna's work [26] whose main arguments were established on complex numbers. We show that even moving to the complex domain, polynomial length of *L* is still sufficient for the DeepSets architecture [9]. We state a formal version of the theorem in the supplementary material.

348 6 Related Work

Works on neural networks to represent set functions have been discussed extensively in the Sec. 1. Here, we review other related works on the expressive power analysis of neural networks.

Early works studied the expressive power of feed-forward neural networks with different activations [27, 28]. Recent works focused on characterizing the benefits of the expressive power of deep architectures to explain their empirical success [39–43]. Modern neural networks often enforce some invariance properties into their architectures such as CNNs that capture spatial translation invariance. The expressive power of invariant neural networks has been analyzed recently [22, 44, 45].

The architectures studied in the above works allow universal approximation of continuous func-356 tions defined on their inputs. However, the family of practically useful architectures that enforce 357 permutation invariance often fail in achieving universal approximation. Graph Neural Networks 358 (GNNs) enforce permutation invariance and can be viewed as an extension of set neural networks 359 to encode a set of pair-wise relations instead of a set of individual elements [20, 21, 46, 47]. GNNs 360 suffer from limited expressive power [5, 17, 18] unless they adopt exponential-order tensors [48]. 361 Hence, previous studies often characterized GNNs' expressive power based on their capability of 362 distinguishing non-isomorphic graphs. Only a few works have ever discussed the function approxima-363 tion property of GNNs [49-51] while these works still miss characterizing such dependence on the 364 depth and width of the architectures [52]. As practical GNNs commonly adopt the architectures that 365 combine feed-forward neural networks with set operations (neighborhood aggregation), we believe 366 the characterization of the needed size for set function approximation studied in [26] and this work 367 may provide useful tools to study finer-grained characterizations of the expressive power of GNNs. 368

369 7 Conclusion

This work investigates how many neurons are needed to model the embedding space for set representation learning with the DeepSets architecture [9]. Our paper provides an affirmative answer that polynomial many neurons in the set size and feature dimension are sufficient. Compared with prior arts, our theory takes high-dimensional features into consideration while significantly advancing the state-of-the-art results from exponential to polynomial.

Limitations. The tightness of our bounds is not examined in this paper, and the complexity of ρ is uninvestigated and left for future exploration. Besides, deriving an embedding layer agnostic lower bound for the embedding space remains another widely open question.

378 A Formal Definitions

In this section, we begin by providing rigorous definitions to specify the topology of the input space of permutation-invariant functions.

Definition A.1. Equipped $\mathcal{K}^{N \times D}$ with the equivalence relation \sim (cf. Definition 2.1), we define metric space $(\mathcal{K}^{N \times D} / \sim, d_F)$, where $d_F : (\mathcal{K}^{N \times D} / \sim) \times (\mathcal{K}^{N \times D} / \sim) \rightarrow \mathbb{R}_{\geq 0}$ is the optimal transport distance:

$$d_F(\boldsymbol{X}, \boldsymbol{X'}) = \min_{\boldsymbol{P} \in \Pi(N)} \|\boldsymbol{P}\boldsymbol{X} - \boldsymbol{X'}\|_{\infty, \infty}, \qquad (11)$$

and \mathcal{K} can be either \mathbb{R} or \mathbb{C} .

Remark A.2. The $\|\cdot\|_{\infty,\infty}$ norm takes the absolute value of the maximal entry: $\max_{i \in [N], j \in [D]} |X_{i,j}|$. Other topologically equivalent matrix norms also apply.

Lemma A.3. The function $d_F : (\mathcal{K}^{N \times D} / \sim) \times (\mathcal{K}^{N \times D} / \sim) \rightarrow \mathbb{R}_{\geq 0}$ is a distance metric on $\mathcal{K}^{N \times D} / \sim$.

Proof. Identity, positivity, and symmetry trivially hold for d_F . It remains to show the triangle inequality as below: for arbitrary $X, X', X'' \in (\mathcal{K}^{N \times D} / \sim, d_F)$,

$$d_F(\boldsymbol{X}, \boldsymbol{X''}) = \min_{\boldsymbol{P} \in \Pi(N)} \|\boldsymbol{P}\boldsymbol{X} - \boldsymbol{X''}\|_{\infty,\infty} \leq \min_{\boldsymbol{P} \in \Pi(N)} \left(\|\boldsymbol{P}\boldsymbol{X} - \boldsymbol{Q}^*\boldsymbol{X'}\|_{\infty,\infty} + \|\boldsymbol{Q}^*\boldsymbol{X'} - \boldsymbol{X''}\|_{\infty,\infty} \right)$$

$$= \min_{\boldsymbol{P} \in \Pi(N)} \|\boldsymbol{P}\boldsymbol{X} - \boldsymbol{Q}^*\boldsymbol{X'}\|_{\infty,\infty} + \|\boldsymbol{Q}^*\boldsymbol{X'} - \boldsymbol{X''}\|_{\infty,\infty}$$

$$= d_F(\boldsymbol{X}, \boldsymbol{X'}) + d_F(\boldsymbol{X}, \boldsymbol{X''}),$$

391 where $Q^* = \operatorname{argmin}_{Q \in \Pi(N)} \|QX' - X''\|_{\infty,\infty}$.

Also we reveal a topological property for $(\mathcal{K}^{N \times D} / \sim, d_F)$ which is essential to show continuity later. **Lemma A.4.** *Each bounded and closed subset of* $(\mathcal{K}^{N \times D} / \sim, d_F)$ *is compact.*

Proof. Without loss of generality, the proof is done by extending Theorem 2.4 in [35] to highdimensional set elements. \Box

Then we can rephrase the definition of permutation invariant function as a proper function mapping between the two metric spaces: $f : (\mathcal{K}^{N \times D} / \sim, d_F) \to (\mathcal{K}^{D'}, d_{\infty})$, where $d_{\infty} : \mathcal{K}^{D'} \times \mathcal{K}^{D'} \to \mathbb{R}_{\geq 0}$ is the ℓ_{∞} -norm induced distance metric.

³⁹⁹ We also recall the definition of injectivity for permutation-invariant functions:

- 400 **Definition A.5** (Injectivity). The following statements are equivalent:
- 401 1. A permutation-invariant function $f : (\mathcal{K}^{N \times D} / \sim, d_F) \to (\mathcal{K}^{D'}, d_{\infty})$ is injective.
- 402 2. There exists a function $g: (\mathcal{K}^{D'}, d_{\infty}) \to (\mathcal{K}^{N \times D} / \sim, d_F)$ such that for every $\mathbf{X} \in \mathcal{K}^{N \times D}$, 403 $g \circ f(\mathbf{X}) \sim \mathbf{X}$.

404 3. For every
$$X, X' \in \mathcal{K}^{N \times D}$$
 such that $f(X) = f(X')$, then $X \sim X'$.

We give an intuitive definition of continuity for permutation-invariant functions via the epsilon-delta statement:

Definition A.6 (Continuity). A permutation-invariant function $f : (\mathcal{K}^{N \times D} / \sim, d_F) \to (\mathcal{K}, d_{\infty})$ is continuous if for arbitrary $\mathbf{X} \in \mathcal{K}^{N \times D}$ and $\epsilon > 0$, there exists $\delta > 0$ such that for every $\mathbf{X'} \in \mathcal{K}^{N \times D}$, $d_{F}(\mathbf{X}, \mathbf{X'}) < \delta$ then $d_{\infty}(f(\mathbf{X}) - f(\mathbf{X'})) < \epsilon$.

410 *Remark* A.7. Since d_F is a distance metric, other equivalent definitions of continuity still applies.

411 **B** Sum-of-Power Mapping

In this section, we extend Definition 2.6 and Lemma 2.8 to the complex version, which provides the mathematical tools for our later proof.

Definition B.1. Define (complex) power mapping: $\psi_N : \mathbb{C} \to \mathbb{C}^N$, $\psi_N(z) = \begin{bmatrix} z & z^2 & \cdots & z^N \end{bmatrix}^\top$ and (complex) sum-of-power mapping $\psi_N : (\mathbb{C}^N / \sim, d_F) \to (\mathbb{C}^N, d_\infty)$, $\psi_N(z) = \sum_{i=1}^N \psi_N(z_i)$.

Lemma B.2 (Existence of Continuous Inverse of Complex Sum-of-Power [35]). ψ_N is injective, thus the inverse $\psi_N^{-1} : (\mathbb{C}^N, d_\infty) \to (\mathbb{C}^N / \sim, d_F)$ exists. Moreover, ψ_N^{-1} is continuous.

Lemma B.3 (Corollary 3.2 [35]). Consider a function $\zeta : (\mathbb{C}^N, d_{\infty}) \to (\mathbb{C}^N / \sim, d_F)$ that maps the coefficients of a polynomial to its root multi-set. Then for any bounded subset $\mathcal{U} \subset (\mathbb{C}^N, d_{\infty})$, the image $\zeta(\mathcal{U}) = \{\zeta(\boldsymbol{z}) : \boldsymbol{z} \in \mathcal{U}\}$ is also bounded.

⁴²¹ *Remark* B.4. Lemma B.3 is also true for real numbers when we constrain the domain of ζ to be real ⁴²² coefficients such that the corresponding polynomial can fully split over the real domain.

Lemma B.5. Consider the N-degree sum-of-power mapping: $\psi_N : (\mathbb{C}^{\mathbb{C},N} / \sim, d_F) \to (\mathbb{C}^N, d_\infty)$, where $\psi_N(\boldsymbol{x}) = \sum_{i=1}^N \psi_N(x_i)$. Denote the range of ψ_N as $\mathcal{Z} \subseteq \mathbb{C}^N$ and its inverse mapping $\psi_N^{-1} : (\mathcal{Z}, d_\infty) \to (\mathbb{C}^N / \sim, d_F)$ (existence guaranteed by Lemma B.2). Then for every bounded set $\mathcal{U} \subset (\mathcal{Z}, d_\infty)$, the image $\Psi_N^{-1}(\mathcal{U}) = \{\Psi_N^{-1}(\boldsymbol{z}) : \boldsymbol{z} \in \mathcal{U}\}$ is also bounded.

Proof. We borrow the proof technique from [9] to reveal a polynomial mapping between $(\mathcal{Z}, d_{\infty})$ and $(\mathbb{C}^N / \sim, d_F)$. For every $\boldsymbol{\xi} \in (\mathbb{C}^N / \sim, d_F)$, let $\boldsymbol{z} = \psi_N(\boldsymbol{\xi})$ and construct a polynomial:

$$P_{\boldsymbol{\xi}}(x) = \prod_{i=1}^{N} (x - \xi_i) = x^N - a_1 x^{N-1} + \dots + (-1)^{N-1} a_{N-1} x + (-1)^N a_N, \qquad (12)$$

where $\boldsymbol{\xi}$ are the roots of $P_{\boldsymbol{\xi}}(x)$ and the coefficients can be written as elementary symmetric polynomials, i.e.,

$$a_n = \sum_{1 \le j_1 \le j_2 \le \dots \le j_n \le N} \xi_{j_1} \xi_{j_2} \cdots \xi_{j_n}, \forall n \in [N].$$
(13)

On the other hand, the elementary symmetric polynomials can be uniquely expressed as a function of z by Newton-Girard formula:

$$a_{n} = \frac{1}{n} \det \begin{bmatrix} z_{1} & 1 & 0 & 0 & \cdots & 0 \\ z_{2} & z_{1} & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ z_{n-1} & z_{n-2} & z_{n-3} & z_{n-4} & \cdots & 1 \\ z_{n} & z_{n-1} & z_{n-2} & z_{n-3} & \cdots & 1 \end{bmatrix} := Q(\boldsymbol{z}), \forall n \in [N]$$
(14)

where the determinant Q(z) is also a polynomial in z. Then the proof proceeds by observing that for any bounded subset $\mathcal{U} \subseteq (\mathcal{Z}, d_{\infty})$, the resulting $\mathcal{A} = Q(\mathcal{U})$ is also bounded in $(\mathbb{C}^N, d_{\infty})$. Therefore, by Lemma B.3, any bounded coefficient set \mathcal{A} will produce a bounded root multi-set in $(\mathbb{C}^N/\sim, d_F)$.

Remark B.6. Lemma B.3 is also true for real numbers. By Remark B.4, we can constrain the ambient space of A in Lemma B.3 to be real coefficients whose corresponding polynomial can split over real numbers, and the same proof proceeds.

440 C Proofs of LP Embedding Layer

In this section, we complete the proofs for the LP embedding layer (Eq. (2)). First we constructively show an upper bound that sufficiently achieves injectivity following our discussion in Sec. 4.1, and then prove Theorem 4.6 to reveal a lower bound that is necessary for injectivity. Finally, we show prove the continuity of the inverse of our constructed LP embedding layer with the techniques introduced in Sec. 4.3.

446 C.1 Upper Bound for Injectivity

To prove the upper bound, we construct an LP embedding layer with $L \le N^5 D^2$ output neurons such that its induced summation is injective. The main ingredient of our construction is anchor defined in Definition 4.1. Two key properties of anchors are stated in Lemma 4.2 and 4.3 (restated as follows) and proved below:

Lemma C.1. Consider the data matrix $X \in \mathbb{R}^{N \times D}$ and $a \in \mathbb{R}^N$ an anchor of X. Then if there exists $P \in \Pi(N)$ such that Pa = a then $Px_i = x_i$ for every $i \in [D]$.

453 Proof of Lemma 4.2. Prove by contradiction. Suppose $Px_i \neq x_i$ for some $i \in [D]$, then there exist 454 some $p, q \in [N]$ such that $x_i^{(p)} \neq x_i^{(q)}$ while $a_p = a_q$. However, this contradicts the definition of an 455 anchor (cf. Definition 4.1).

Lemma C.2 (Union Alignment based on Anchor Alignment). Consider the data matrix $X, X' \in \mathbb{R}^{N \times D}$, $a \in \mathbb{R}^N$ is an anchor of X and $a' \in \mathbb{R}^N$ is an arbitrary vector. If $[a \ x_i] \sim [a' \ x'_i]$ for every $i \in [D]$, then $X \sim X'$.

Proof of Lemma 4.3. According to definition of equivalence, there exists $Q_i \in \Pi(N)$ for every $i \in [D]$ such that $[a \ x_i] = [Q_i a' \ Q_i x'_i]$. Moreover, since $[a \ x_i] \sim [a' \ x'_i]$, it must hold that $a \sim a'$, i.e., there exists $P \in \Pi_N$ such that Pa = a'. Combined together, we have that $Q_i Pa = a$.

Next, we choose $Q'_i = P^{\top}Q_i^{\top}$ so $Q'_i a = Q'_iQ_iPa = a$. Due to the property of anchors (Lemma 464 4.2), we have $Q'_i x_i = x_i$. Notice that $x_i = Q'_i x_i = P^{\top}Q_i^{\top}Q_i x'_i = Px'_i$. Therefore, we can 465 conclude the proof as we have found a permutation matrix P that simultaneously aligns x_i and x'_i 466 for every $i \in [D]$, i.e., $X = [x_1 \cdots x_D] = [Px_1 \cdots Px_D] = PX'$.

Next, we need to examine how many weights are needed to construct an anchor via linear combining
 all the existing channels. We restate Lemma 4.4 with more specifications as well as a mathematical
 device to prove it as below:

Lemma C.3. Consider D linearly independent weight vectors $\{\boldsymbol{w}_1, \cdots, \boldsymbol{w}_D \in \mathbb{R}^D\}$. Then for every $p, q \in [N]$ such that $\boldsymbol{x}^{(p)} \neq \boldsymbol{x}^{(q)}$, there exists $\boldsymbol{w}_j, j \in [D]$, such that $\boldsymbol{x}^{(p)\top} \boldsymbol{w}_j \neq \boldsymbol{x}^{(q)\top} \boldsymbol{w}_j$.

Proof. This is the simple fact of full-rank linear systems. Prove by contradiction. Suppose for $\forall j \in [D]$ we have $\mathbf{x}^{(p)\top} \mathbf{w}_j = \mathbf{x}^{(q)\top} \mathbf{w}_j$. Then we form a linear system: $\mathbf{x}^{(p)\top} [\mathbf{w}_1 \cdots \mathbf{w}_D] =$ $\mathbf{x}^{(q)\top} [\mathbf{w}_1 \cdots \mathbf{w}_D]$. Since $\mathbf{w}_1, \cdots, \mathbf{w}_D$ are linearly independent, it yields $\mathbf{x}^{(p)} = \mathbf{x}^{(q)}$, which 475 reaches the contradiction.

Lemma C.4 (Anchor Construction). Consider a set of weight vectors $\{\boldsymbol{w}_1, \cdots, \boldsymbol{w}_{K_1} \in \mathbb{R}^D\}$ with $K_1 = N(N-1)(D-1)/2 + 1$, of which every D-length subset, i.e., $\{\boldsymbol{w}_j : \forall j \in \mathcal{J}\}, \forall \mathcal{J} \subseteq \mathcal{J}\}$ $[K_1], |\mathcal{J}| = D$, is linearly independent, then for every data matrix $\boldsymbol{X} \in \mathbb{R}^{N \times D}$, there exists $j^* \in [K_1], \boldsymbol{X} \boldsymbol{w}_{j^*}$ is an anchor of \boldsymbol{X} .

480 *Proof.* Define a set of pairs which an anchor needs to distinguish: $\mathcal{I} = \{(p,q) : \mathbf{x}^{(p)} \neq \mathbf{x}^{(q)}\} \subseteq [N]^2$ 481 Consider a *D*-length subset $\mathcal{J} \subseteq [K]$ with $|\mathcal{J}| = D$. Since $\{\mathbf{w}_j : \forall j \in \mathcal{J}\}$ are linear independent, 482 we assert by Lemma C.3 that for every pair $(p,q) \in \mathcal{I}$, there exists $j \in \mathcal{J}, \mathbf{x}^{(p)\top}\mathbf{w}_j \neq \mathbf{x}^{(q)\top}\mathbf{w}_j$. 483 It is equivalent to claim: for every pair $(p,q) \in \mathcal{I}$, at most D-1 many $\mathbf{w}_j, j \in [K_1]$ satisfy 484 $\mathbf{x}^{(p)\top}\mathbf{w}_j = \mathbf{x}^{(q)\top}\mathbf{w}_j$. Based on pigeon-hold principle, as long as $K_1 \ge N(N-1)(D-1)/2 + 1 =$ 485 $(D-1)\binom{N}{2} + 1 \ge (D-1)|\mathcal{I}| + 1$, there must exist $j^* \in [K_1]$ such that $\mathbf{x}^{(p)\top}\mathbf{w}_{j^*} \neq \mathbf{x}^{(q)\top}\mathbf{w}_{j^*}$ for 486 $\forall (p,q) \in \mathcal{I}$. By Definition 4.1, $\mathbf{X}\mathbf{w}_{j^*}$ generates an anchor.

Proposition C.5. The linear independence condition in Lemma C.4 can be satisfied with probability one by drawing i.i.d. Gaussian random vectors $w_1, \dots, w_{K_1} \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I})$.

Proof. We first note that generating a $D \times K_1$ Gaussian random matrix $(D \leq K_1)$ is equivalent to drawing a matrix with respect to a probability measure defined over $\mathcal{M} = \{ \mathbf{X} \in \mathbb{R}^{D \times K} : \operatorname{rank}(\mathbf{X}) \leq D \}$. Since rank-D matrices are dense in \mathcal{M} [53, 54], we can conclude that for $\forall \mathcal{J} \subseteq [K_1]$, $|\mathcal{J}| = D$, $\mathbb{P}(\{w_j : j \in \mathcal{J}\}\)$ are linearly independent) = 1. By union bound, $\mathbb{P}(\{w_j : j \in \mathcal{J}\}\)$ for all $\mathcal{J} \in [K]$, $|\mathcal{J}| = D$ are linearly independent) = 1.

We also restate Lemma 4.5 in the following lemma to demonstrate the weight construction for anchor matching:

Lemma C.6 (Anchor Matching). Consider a group of coefficients $\Gamma = \{\gamma_1, \dots, \gamma_{K_2} \in \mathbb{R}\}$ with $\gamma_i \neq 0, \forall i \in [K_2], \gamma_i \neq \gamma_j, \forall i, j \in [K_2], and K_2 = N(N-1) + 1$ such that for all 4-tuples $(\gamma_i, \gamma_j, \gamma_k, \gamma_l) \subset \Gamma$, if $\gamma_i \neq \gamma_j, \gamma_i \neq \gamma_k$ then $\gamma_i/\gamma_j \neq \gamma_k/\gamma_l$. It must hold that: Given any $x, x', y, y' \in \mathbb{R}^N$ such that $x \sim x'$ and $y \sim y'$, if $(x - \gamma_k y) \sim (x' - \gamma_k y')$ for every $k \in [K_2]$, then $[x \ y] \sim [x' \ y']$.

Proof. We note that $\boldsymbol{x} \sim \boldsymbol{x'}$ and $\boldsymbol{y} \sim \boldsymbol{y'}$ imply that there exist $\boldsymbol{P}_x, \boldsymbol{P}_y \in \Pi(N)$ such that $\boldsymbol{P}_x \boldsymbol{x} = \boldsymbol{x'}$ and $\boldsymbol{P}_y \boldsymbol{y} = \boldsymbol{y'}$. Also $(\boldsymbol{x} - \gamma_k \boldsymbol{y}) \sim (\boldsymbol{x'} - \gamma_k \boldsymbol{y'}), \forall k \in [K_2]$ implies there exists $\boldsymbol{Q}_k \in \Pi(N), \forall k \in [K_2]$ such that $\boldsymbol{Q}_k(\boldsymbol{x} - \gamma_k \boldsymbol{y}) = \boldsymbol{x'} - \gamma_k \boldsymbol{y'}$. Substituting the former to the latter, we can obtain:

$$\left(\boldsymbol{I} - \boldsymbol{Q}_{k}^{\top} \boldsymbol{P}_{x}\right) \boldsymbol{x} = \gamma_{k} \left(\boldsymbol{I} - \boldsymbol{Q}_{k}^{\top} \boldsymbol{P}_{x}\right) \boldsymbol{y}, \quad \forall k \in [K_{2}],$$
(15)

where for each $k \in [K_2]$, Eq. (15) corresponds to N equalities as follows. Here, we let $(Z)_i$ denote the *i*th column of the matrix Z.

$$(\boldsymbol{I} - \boldsymbol{Q}_{k}^{\top} \boldsymbol{P}_{x})_{1}^{\top} \boldsymbol{x} = \gamma_{k} (\boldsymbol{I} - \boldsymbol{Q}_{k}^{\top} \boldsymbol{P}_{x})_{1}^{\top} \boldsymbol{y}$$

$$\vdots$$

$$(\boldsymbol{I} - \boldsymbol{Q}_{k}^{\top} \boldsymbol{P}_{x})_{N}^{\top} \boldsymbol{x} = \gamma_{k} (\boldsymbol{I} - \boldsymbol{Q}_{k}^{\top} \boldsymbol{P}_{x})_{N}^{\top} \boldsymbol{y}$$
(16)

We compare and entries in $\boldsymbol{x} = [\cdots x_p \cdots]^\top$ and for each entry index $p \in [N]$, we define a set of nonzero pairwise differences between x_p and other entries in $\boldsymbol{x} : \mathcal{D}_x^{(p)} = \{x_p - x_q : q \in [N], x_p \neq x_q\}$. Similarly, for \boldsymbol{y} , we define $\mathcal{D}_y^{(p)} = \{y_p - y_q : q \in [N], y_p \neq y_q\}$. We note that for every $n \in [N]$, either $(\boldsymbol{I} - \boldsymbol{Q}_k^\top \boldsymbol{P}_x)_n^\top \boldsymbol{x} = 0$ or $(\boldsymbol{I} - \boldsymbol{Q}_k^\top \boldsymbol{P}_x)_n^\top \boldsymbol{x} \in \mathcal{D}_x^{(p)}$ for some $p \in [N]$ as $(\boldsymbol{Q}_k^\top \boldsymbol{P}_x)_n^\top \boldsymbol{x}$ is some x_q . Then it is sufficient to show there must exist $k \in [K_2]$ such that none of equations in Eq. (16) can be induced by some elements in $\mathcal{D}_x^{(p)}$, i.e.,

$$\exists k^* \in [K_2], \forall p, n \in [N] \text{ such that } (\boldsymbol{I} - \boldsymbol{Q}_{k^*}^\top \boldsymbol{P}_x)_n^\top \boldsymbol{x} \notin \mathcal{D}_x^{(p)}.$$
(17)

⁵¹² This is because Eq. (17) implies:

$$(\boldsymbol{I} - \boldsymbol{Q}_{k^*} \boldsymbol{P}_x)^\top \boldsymbol{x} = \boldsymbol{0} \quad \Rightarrow \quad \boldsymbol{x} = \boldsymbol{Q}_{k^*}^\top \boldsymbol{P}_x \boldsymbol{x} = \boldsymbol{Q}_{k^*}^\top \boldsymbol{x}',$$

(Since $\gamma_k \neq 0, \forall k \in [K_2]$) $(\boldsymbol{I} - \boldsymbol{Q}_{k^*} \boldsymbol{P}_y)^\top \boldsymbol{y} = \boldsymbol{0} \quad \Rightarrow \quad \boldsymbol{y} = \boldsymbol{Q}_{k^*}^\top \boldsymbol{P}_y \boldsymbol{y} = \boldsymbol{Q}_{k^*}^\top \boldsymbol{y}',$

513 which is $[\boldsymbol{x} \quad \boldsymbol{y}] = \boldsymbol{Q}_{k^*}^\top [\boldsymbol{x'} \quad \boldsymbol{y'}].$

To show Eq. (17), we construct N bipartite graphs $\mathcal{G}^{(p)} = (\mathcal{D}_x^{(p)}, \mathcal{D}_y^{(p)}, \mathcal{E}^{(p)})$ for $p \in [N]$ where each $\alpha \in \mathcal{D}_x^{(p)}$ or each $\beta \in \mathcal{D}_y^{(p)}$ is viewed as a node and $(\alpha, \beta) \in \mathcal{E}^{(p)}$ gives an edge if $\alpha = \gamma_k \beta$ for some $k \in [K_2]$. Then we prove the existence of k^* via seeing a contradiction that does counting the number of connected pairs (α, β) from two perspectives.

Perspective of $\mathcal{D}_x^{(p)}$. We argue that for $\forall p \in [N]$ and arbitrary $\alpha_1, \alpha_2 \in \mathcal{D}_x^{(p)}, \alpha_1 \neq \alpha_2$, there exists at most one $\beta \in \mathcal{D}_y^{(p)}$ such that $(\alpha_1, \beta) \in \mathcal{E}^{(p)}$ and $(\alpha_2, \beta) \in \mathcal{E}^{(p)}$. Otherwise, suppose there exists $\beta' \in \mathcal{D}_y^{(p)}, \beta' \neq \beta$ such that $(\alpha_1, \beta') \in \mathcal{E}^{(p)}$ and $(\alpha_2, \beta') \in \mathcal{E}^{(p)}$. Then we have $\alpha_1 = \gamma_i \beta$, $\alpha_2 = \gamma_j \beta, \alpha_1 = \gamma_k \beta'$, and $\alpha_2 = \gamma_l \beta'$ for some $\gamma_i, \gamma_j, \gamma_k, \gamma_l \in \Gamma$, which is $\gamma_i / \gamma_k = \gamma_k / \gamma_l$. As $\alpha_1 \neq \alpha_2$, it is obvious that $\gamma_i \neq \gamma_j$. Similarly, we have $\gamma_i \neq \gamma_k$. Altogether, it contradicts our assumption on Γ . Therefore, $|\mathcal{E}^{(p)}| \leq 2 \max\{|\mathcal{D}_x^{(p)}|, |\mathcal{D}_y^{(p)}|\} \leq 2(N-1)$. And the total edge number of all bipartite graphs should be less than 2N(N-1).

Perspective of Γ . We note that if for some $k \in [K_2]$ that makes $(\boldsymbol{I} - \boldsymbol{Q}_k^\top \boldsymbol{P}_x)_n^\top \boldsymbol{x} \in \mathcal{D}_x^{(p)}$ for some $p, n \in [N]$, i.e., $(\boldsymbol{I} - \boldsymbol{Q}_k^\top \boldsymbol{P}_x)_n^\top \boldsymbol{x} = \gamma_k (\boldsymbol{I} - \boldsymbol{Q}_k^\top \boldsymbol{P}_y)_n^\top \boldsymbol{y} \neq 0$, this γ_k contributes at least two edges in the entire bipartite graph, i.e., there being another $n' \in [N]$, $(\boldsymbol{I} - \boldsymbol{Q}_k^\top \boldsymbol{P}_x)_{n'}^\top \boldsymbol{x} = \gamma_k (\boldsymbol{I} - \boldsymbol{Q}_k^\top \boldsymbol{P}_y)_n^\top \boldsymbol{y} \neq 0$. 525 526 527 0. Otherwise, there exists a unique $n^* \in [N]$ such that $(\boldsymbol{I} - \boldsymbol{Q}_k^\top \boldsymbol{P}_x)_{n^*}^\top \boldsymbol{x} \in \mathcal{D}_x^{(p)}(\neq 0)$ and $(\boldsymbol{I} - \boldsymbol{Q}_k^\top \boldsymbol{P}_x)_n^\top \boldsymbol{x} = 0$ for all $n \neq n^*$. This cannot be true because $\mathbf{1}^\top (\boldsymbol{I} - \boldsymbol{Q}_k^\top \boldsymbol{P}_x) \boldsymbol{x} = 0$. By which, 528 529 if $\forall k \in [K_2], \exists p, n \in [N]$ such that $(\mathbf{I} - \mathbf{Q}_k^\top \mathbf{P}_x)_n^\top \mathbf{x} \in \mathcal{D}_x^{(p)}$ (i.e., Eq. (17) is always false), then the total number of edges is at least $2K_2 = 2N(N-1) + 2$. 530 531 Hereby, we conclude the proof by the contradiction, in which the minimal count of edges $2K_2$ by 532 **Perspective of** Γ already surpasses the maximal number 2N(N-1) by **Perspective of** $\mathcal{D}_x^{(p)}$. 533

- *Remark* C.7. A handy choice of Γ in Lemma C.6 are prime numbers, which are provably positive, 534 infinitely many, and not divisible by each other. 535
- We wrap off this section by formally stating and proving the injectivity statement of the LP embedding 536 537 layer.

Theorem C.8. Suppose $\phi : \mathbb{R}^D \to \mathbb{R}^L$ admits the form of Eq. (2) where $L = KN \leq N^5D^2$, $K = D + K_1 + DK_1K_2$ and $\mathbf{W} = \begin{bmatrix} \mathbf{w}_1^{(1)} & \cdots & \mathbf{w}_D^{(1)} & \mathbf{w}_1^{(2)} & \cdots & \mathbf{w}_{K_1}^{(2)} & \cdots & \mathbf{w}_{i,j,k}^{(3)} & \cdots \end{bmatrix}$ 538 539 is constructed as follows: 540

- 1. Let the first group of weights $w_1^{(1)} = e_1, \cdots, w_D^{(1)} = e_D$ to buffer the original features, where e_i 541 is the *i*-th canonical basis. 542
- 2. Choose the second group of linear weights, $w_1^{(2)}, \dots, w_{K_1}^{(2)}$ for K_1 as large as N(N-1)(D-1)/2 + 1, such that the conditions in Lemma C.4 are satisfied. 543 544

545 3. Design the third group of weights $\boldsymbol{w}_{i,j,k}^{(3)}$ for $i \in [D], j \in [K_1], k \in [K_2]$ where $\boldsymbol{w}_{i,j,k}^{(3)}$ $e_i - \gamma_k w_j^{(2)}$, $K_2 = N(N-1) + 1$, and $\gamma_k, k \in [K_2]$ are chosen such that conditions in Lemma C.6 are satisfied. 546 547

Then $\sum_{i=1}^{N} \phi(\boldsymbol{x}^{(i)})$ is injective (cf. Definition A.5). 548

Proof. Suppose $\sum_{i=1}^{N} \phi(\boldsymbol{x}^{(i)}) = \sum_{i=1}^{N} \phi(\boldsymbol{x'}^{(i)})$ for some $\boldsymbol{X}, \boldsymbol{X'} \in \mathbb{R}^{N \times D}$. Due to the property of power mapping (cf. Lemma 2.8): 549 550

$$\sum_{i=1}^{N} \phi(\boldsymbol{x}^{(i)}) = \sum_{i=1}^{N} \phi(\boldsymbol{x'}^{(i)}) \Rightarrow \boldsymbol{X} \boldsymbol{w}_{i}^{(1)} \sim \boldsymbol{X'} \boldsymbol{w}_{i}^{(1)}, \forall i \in [D], \boldsymbol{X} \boldsymbol{w}_{i}^{(2)} \sim \boldsymbol{X'} \boldsymbol{w}_{i}^{(2)}, \forall i \in [K_{1}], \quad (18)$$
$$\boldsymbol{X} \boldsymbol{w}_{i,j,k}^{(3)} \sim \boldsymbol{X'} \boldsymbol{w}_{i,j,k}^{(3)}, \forall i \in [D], j \in [K_{1}], k \in [K_{2}].$$

By Lemma C.4, it is guaranteed that there exists $j^* \in [K_1]$ such that $X w_{j^*}^{(2)}$ is an anchor, and according to Eq. (18), we have $X w_{i^*}^{(2)} \sim X' w_{i^*}^{(2)}$. By Lemma C.6, we induce: 552

$$\begin{aligned} \mathbf{X} \boldsymbol{w}_{i}^{(1)} \sim \mathbf{X'} \boldsymbol{w}_{i}^{(1)}, \forall i \in [D], \mathbf{X} \boldsymbol{w}_{j^{*}}^{(2)} \sim \mathbf{X'} \boldsymbol{w}_{j^{*}}^{(2)}, \mathbf{X} \boldsymbol{w}_{i,j,k}^{(3)} \sim \mathbf{X'} \boldsymbol{w}_{i,j,k}^{(3)}, \forall i \in [D], j \in [K_{1}], k \in [K_{2}] \\ \Rightarrow \begin{bmatrix} \mathbf{X} \boldsymbol{w}_{j^{*}}^{(2)} & \mathbf{x}_{i} \end{bmatrix} \sim \begin{bmatrix} \mathbf{X'} \boldsymbol{w}_{j^{*}}^{(2)} & \mathbf{x'}_{i} \end{bmatrix}, \forall i \in [D]. \end{aligned}$$
(19)

Since $Xw_{i^*}^{(2)}$ is an anchor, by union alignment (Lemma 4.3), we have: 553

$$\begin{bmatrix} \boldsymbol{X}\boldsymbol{w}_{j^*}^{(2)} & \boldsymbol{x}_i \end{bmatrix} \sim \begin{bmatrix} \boldsymbol{X}'\boldsymbol{w}_{j^*}^{(2)} & \boldsymbol{x}'_i \end{bmatrix}, \forall i \in [D] \Rightarrow \boldsymbol{X} \sim \boldsymbol{X}'.$$
(20)

Here $K = D + K_1 + DK_1K_2 \le N^4D^2$, thus $L = KN \le N^5D^2$, which concludes the proof. \Box 554

C.2 Continuity 555

- Next, we show that under the construction of Theorem C.8, the inverse of $\sum_{i=1}^{N} \phi(x^{(i)})$ is continuous. The main idea is to check the three conditions provided in Lemma 4.10: 556
- 557

Corollary C.9. Consider channel-wise high-dimensional sum-of-power $\widehat{\Psi_N}(\mathbf{X})$: $(\mathbb{R}^{N \times K} / \sim d_F) \rightarrow (\mathbb{R}^{NK}, d_{\infty})$ defined as below:

$$\widehat{\Psi_N}(\boldsymbol{X}) = \begin{bmatrix} \Psi_N(\boldsymbol{x}_1)^\top & \cdots & \Psi_N(\boldsymbol{x}_K)^\top \end{bmatrix}^\top \in (\mathbb{R}^{NK}, d_\infty),$$
(21)

and an associated mapping $\widehat{\Phi_N}$: $(\mathbb{R}^{NK}, d_\infty) \to (\mathbb{R}^N / \sim, d_F)^K$:

$$\widehat{\Phi_N}(\boldsymbol{Z}) = \begin{bmatrix} \Psi_N^{-1}(\boldsymbol{z}_1) & \cdots & \Psi_N^{-1}(\boldsymbol{z}_K) \end{bmatrix},$$
(22)

where $\mathbf{Z} = \begin{bmatrix} \mathbf{z}_1^\top & \cdots & \mathbf{z}_K^\top \end{bmatrix}^\top$, $\mathbf{z}_i \in \mathbb{R}^N, \forall i \in [K]$. We denote the induced product metric over ($\mathbb{R}^N / \sim, d_F$)^K as d_F^K : $(\mathbb{R}^N / \sim)^K \times (\mathbb{R}^N / \sim)^K \to \mathbb{R}_{\geq 0}$:

$$d_F^K(\boldsymbol{Z}, \boldsymbol{Z'}) = \max_{i \in [K]} d_F(\boldsymbol{z}_i, \boldsymbol{z'}_i).$$
(23)

Then the mapping $\widehat{\Phi_N}$ maps any bounded set in $(\mathbb{R}^{NK}, d_\infty)$ to a bounded set in $(\mathbb{R}^N / \sim, d_F)^K$.

Proof. Proved by noting that if $d_{\infty}(\boldsymbol{z}_i, \boldsymbol{z'}_i) \leq C_1$ for some $\boldsymbol{z}_i, \boldsymbol{z'}_i \in (\mathbb{R}^N, d_{\infty}), \forall i \in [K]$ and a constant $C_1 \geq 0$, then $d_F(\Psi_N^{-1}(\boldsymbol{z}_i), \Psi_N^{-1}(\boldsymbol{z'}_i)) \leq C_2, \forall i \in [K]$ for some constant $C_2 \geq 0$ by Lemma B.5 and Remark B.6. Finally, we have:

$$d_F^K(\widehat{\Phi_N}(\boldsymbol{Z}), \widehat{\Phi_N}(\boldsymbol{Z'}) = \max_{i \in [K]} d_F(\Psi_N^{-1}(\boldsymbol{z}_i), \Psi_N^{-1}(\boldsymbol{z'}_i)) \le C_2,$$

⁵⁶⁷ which is also bounded above.

Now we are ready to present and prove the continuity of the LP embedding layer.

Theorem C.10. Suppose ϕ admits the form of Eq. (2) and follows the construction in Theorem C.8, then the inverse of LP-induced sum-pooling $\sum_{i=1}^{N} \phi(\mathbf{x}^{(i)})$ is continuous.

Proof. The proof is done by invoking Lemma 4.10. First of all, the inverse of $\Psi(\mathbf{X}) = \sum_{i=1}^{N} \phi(\mathbf{x}^{(i)})$, denoted as $\Psi^{-1} : (\mathbb{R}^{NK}, d_{\infty}) \to (\mathbb{R}^{N \times D} / \sim, d_F)$, exists due to Theorem C.8. By Lemma A.4, any closed and bounded subset of $(\mathbb{R}^{N \times D} / \sim, d_F)$ is compact. Trivially, $\Psi(\mathbf{X})$ is continuous. Then it remains to show the condition (**c**) in Lemma 4.10. We decompose Ψ^{-1} into two mappings following the similar idea of proving its existence:

$$(\mathbb{R}^{NK}, d_{\infty}) \xrightarrow{\Phi_N} (\mathbb{R}^N / \sim, d_F)^K \xrightarrow{\pi} (\mathbb{R}^{N \times D} / \sim, d_F) :$$

where $\widehat{\Phi_N}$ is defined in Eq. (22) and π exists due to Eqs. (19) and (20) in Theorem C.8. Also according to our construction in Theorem C.8, the first identity block induces that: for any $\mathbf{Z} \in (\mathbb{R}^N / \sim, d_F)^K$, $z_i \sim \pi(\mathbf{Z})_i$ for every $i \in [D]$. Therefore, $\forall \mathbf{Z}, \mathbf{Z'} \in (\mathbb{R}^N / \sim, d_F)^K$ such that $d_F^K(\mathbf{Z}, \mathbf{Z'}) \leq C$ for some constant C > 0, we have:

$$d_F(\pi(\boldsymbol{Z}), \pi(\boldsymbol{Z'})) \le \max_{i \in [D]} d_F(\boldsymbol{z}_i, \boldsymbol{z'}_i) \le d_F^K(\boldsymbol{Z}, \boldsymbol{Z'}) \le C,$$
(24)

which implies π maps every bounded set in $(\mathbb{R}^N / \sim, d_F)^K$ to a bounded set in $(\mathbb{R}^{N \times D} / \sim, d_F)$. Now we conclude the proof by the following chain of argument:

$$\mathcal{Z} \subseteq (\mathbb{R}^{NK}, d_{\infty})$$
 is bounded $\xrightarrow{\text{Corollary C.9}} \widehat{\Phi_N}(\mathcal{Z})$ is bounded $\xrightarrow{\text{Eq. (24)}} \pi \circ \widehat{\Phi_N}(\mathcal{Z})$ is bounded \Box

582

583 C.3 Lower Bound for Injectivity

In this section, we prove Theorem 4.6 which shows that $K \ge D + 1$ is necessary for injectivity of LP-induced sum-pooling when $D \ge 2$. Our argument mainly generalizes Lemma 2 of [55] to our equivalence class. To proceed our argument, we define the linear subspace V by vectorizing $[Xw_1 \cdots Xw_K]$ as below:

$$\mathcal{V} := \left\{ \begin{bmatrix} \boldsymbol{X} \boldsymbol{w}_1 \\ \vdots \\ \boldsymbol{X} \boldsymbol{w}_K \end{bmatrix} : \boldsymbol{X} \in \mathbb{R}^{N \times D} \right\} = \mathcal{R} \left(\begin{bmatrix} (\boldsymbol{w}_1 \otimes \boldsymbol{I}_N) \\ \vdots \\ (\boldsymbol{w}_K \otimes \boldsymbol{I}_N) \end{bmatrix} \right), \quad (25)$$

where $\mathcal{R}(\mathbf{Z})$ denotes the column space of \mathbf{Z} and \otimes is the Kronecker product. \mathcal{V} is a linear subspace of \mathbb{R}^{NK} with dimension at most \mathbb{R}^{ND} , characterized by $\mathbf{w}_1, \cdots, \mathbf{w}_K \in \mathbb{R}^D$. For the sake of notation simplicity, we denote $\Pi(N)^{\otimes K} = \{ \operatorname{diag}(\mathbf{Q}_1, \cdots, \mathbf{Q}_K) : \forall \mathbf{Q}_1, \cdots, \mathbf{Q}_K \in \Pi(N) \}$, and $\mathbf{I}_K \otimes \Pi(N) = \{ \mathbf{I}_K \otimes \mathbf{Q} : \forall \mathbf{Q} \in \Pi(N) \}$. Next, we define the notion of unique recoverability [55] 588 589 590 591 as below: 592

Definition C.11 (Unique Recoverability). The subspace \mathcal{V} is called uniquely recoverable under 593 $Q \in \Pi(N)^{\otimes K}$ if whenever $x, x' \in \mathcal{V}$ satisfy Qx = x', there exists $P \in I_K \otimes \Pi(N), Px = x'$. 594

Subsequently, we derive a necessary condition for the unque recoverability: 595

Lemma C.12. A linear subspace $\mathcal{V} \subseteq \mathbb{R}^{NK}$ is uniquely recoverable under $Q \in \Pi(N)^{\otimes K}$ only 596

if there exists $P \in I_K \otimes \Pi(N)$, $Q(\overline{\mathcal{V})} \cap \mathcal{V} \subset \mathcal{E}_{QP^{\top},\lambda=1}$, where $\mathcal{E}_{QP^{\top},\lambda}$ denotes the eigenspace 597

corresponding to the eigenvalue λ . 598

Proof. We first show that $Q(\mathcal{V}) \cap \mathcal{V} \subseteq \bigcup_{P \in I_K \otimes \Pi(K)} \mathcal{E}_{QP^{\top},\lambda=1}$. Since the LHS is a subspace, while the RHS is a union of subspaces, there exists $P \in I_K \otimes \Pi(K)$ such that $Q(\mathcal{V}) \cap \mathcal{V} \subseteq \mathcal{E}_{QP^{\top},\lambda=1}$. 599 600 Then it remains to show $Q(\mathcal{V}) \cap \mathcal{V} \subseteq \bigcup_{P \in I_K \otimes \Pi(K)} \mathcal{E}_{QP^{\top}, \lambda=1}$. 601

Proved by contradiction. Suppose there exists $x \in Q(\mathcal{V}) \cap \mathcal{V}$ but $x \notin \bigcup_{P \in I_K \otimes \Pi(K)} \mathcal{E}_{QP^{\top}, \lambda=1}$. Or 602 equivalently, there exists $x' \in \mathcal{V}$ and x = Qx', while for $\forall P \in I_K \otimes \Pi(N), x \neq QP^{\top}x$. This 603 implies $Q^{\top} x = x' \neq P x$ for $\forall P \in I_K \otimes \Pi(N)$. However, this contradicts the fact that $\mathcal{V} \subseteq \mathbb{R}^{NK}$ 604 is uniquely recoverable (cf. Definition C.11). 605

We also introduce a useful Lemma C.13 that gets rid of the discussion on Q in the inclusion: 606

Lemma C.13. Suppose $\mathcal{V} \subseteq \mathbb{R}^N$ is a linear subspace, and A is a linear mapping. $A(\mathcal{V}) \cap \mathcal{V} \cap \mathcal{E}_{A,\lambda} =$ 607 **0** if and only if $\mathcal{V} \cap \mathcal{E}_{\mathbf{A},\lambda} = \mathbf{0}$. 608

Proof. The sufficiency is straightforward. The necessity is shown by contradiction: Suppose $\mathcal{V} \cap$ 609 $\mathcal{E}_{A,\lambda} \neq \mathbf{0}$, then there exists $\mathbf{x} \in \mathcal{V} \cap \mathcal{E}_{A,\lambda}$ such that $\mathbf{x} \neq \mathbf{0}$. Then $A\mathbf{x} = \lambda \mathbf{x}$ implies $\mathbf{x} \in A(\mathcal{V})$. Hence, $\mathbf{x} \in A(\mathcal{V}) \cap \mathcal{V} \cap \mathcal{E}_{A,\lambda}$ which reaches the contradiction. 610

611

Now we are ready to present the proof of Theorem 4.6: 612

Proof of Theorem 4.6. Proved by contrapositive. First notice that, $\forall X, X' \in \mathbb{R}^{N \times D}, Xw_i \sim X'w_i, \forall i \in [K] \Rightarrow X \sim X'$ holds if and only if dim $\mathcal{V} = ND$ and \mathcal{V} is uniquely recoverable under all possible $Q \in \Pi(N)^{\otimes K}$. By Lemma C.12, for every $Q \in \Pi(N)^{\otimes K}$, there exists $P \in I_K \otimes \Pi(N)$ such that $Q(\mathcal{V}) \cap \mathcal{V} \subset \mathcal{E}_{QP^{\top},\lambda=1}$. This is $Q(\mathcal{V}) \cap \mathcal{V} \cap \mathcal{E}_{QP^{\top},\lambda} = 0$ for all $\lambda \neq 1$. By Lemma C.13, we have $\mathcal{V} \cap \mathcal{E}_{QP^{\top},\lambda} = 0$ for all $\lambda \neq 1$. Then proof is concluded by discussing the dimension of ambient space \mathbb{R}^{NK} such that an ND-dimensional subspace \mathcal{V} can reside. To ensure $\mathcal{V} \cap \mathcal{E}_{QP^{\top},\lambda} = 0$ 613 614 615 616 617 618 for all $\lambda \neq 1$, it is necessary that $\dim \mathcal{V} \leq \min_{\lambda \neq 1} \operatorname{codim} \mathcal{E}_{QP^{\top},\lambda}$ for every $Q \in \Pi(N)^{\otimes K}$ and its 619 associated $P \in I_K \otimes \Pi(N)$. Relaxing the dependence between Q and P, we derive the inequality: 620

$$ND = \dim \mathcal{V} \le \min_{\mathbf{Q} \in \Pi(N) \otimes K} \max_{\mathbf{P} \in \mathbf{I}_K \otimes \Pi(N)} \min_{\lambda \neq 1} \operatorname{codim} \mathcal{E}_{\mathbf{Q}\mathbf{P}^{\top},\lambda} \le NK - 1,$$
(26)

where the last inequality considers the scenario where every non-one eigenspace is one-dimensional, 621 which is achievable when $K \ge 2$. Hence, we can bound $K \ge D + 1/K$, i.e., $K \ge D + 1$. 622

Proofs for LLE Embedding Layer D 623

In this section, we present the complete proof for the LLE embedding layer (Eq. (3)). Similar to 624 the LP embedding layer, we construct an LLE whose induced sum-pooling is injective following 625 arguments in Sec. 4.2 and has continuous inverse with the techniques introduced in Sec. 4.3. 626

627 D.1 Upper Bound for Injectivity

To prove the upper bound, we construct an LLE embedding layer with $L \le 2N^2D^2$ output neurons such that its induced sum-pooling is injective. The main proof technique is to bind every pair of channel with complex numbers and invoke the injectivity of sum-of-power mapping over the complex domain.

- 632 With Lemma B.2, we can prove Lemma 4.8 as below:
- ⁶³³ Proof of Lemma 4.8. If for any pair of vectors $\boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathbb{R}^N, \boldsymbol{y}_1, \boldsymbol{y}_2 \in \mathbb{R}^N$ such that ⁶³⁴ $\sum_{i \in [N]} \boldsymbol{x}_{1,i}^{l-k} \boldsymbol{x}_{2,i}^k = \sum_{i \in [N]} \boldsymbol{y}_{1,i}^{l-k} \boldsymbol{y}_{2,i}^k$ for every $l, k \in [N], 0 \le k \le l$, then for $\forall l \in [N]$,

$$\sum_{i=1}^{N} (\boldsymbol{x}_{1,i} + \boldsymbol{x}_{2,i}\sqrt{-1})^{l} = \sum_{i=1}^{N} \sum_{k=0}^{l} (\sqrt{-1})^{k} \boldsymbol{x}_{1,i}^{l-k} \boldsymbol{x}_{2,i}^{k} = \sum_{k=0}^{l} (\sqrt{-1})^{k} \left(\sum_{i=1}^{N} \boldsymbol{x}_{1,i}^{l-k} \boldsymbol{x}_{2,i}^{k}\right)$$

$$= \sum_{k=0}^{l} (\sqrt{-1})^{k} \left(\sum_{i=1}^{N} \boldsymbol{y}_{1,i}^{l-k} \boldsymbol{y}_{2,i}^{k}\right)$$

$$= \sum_{i=1}^{N} (\boldsymbol{y}_{1,i} + \boldsymbol{y}_{2,i}\sqrt{-1})^{l}$$

$$= \psi_{N}(\boldsymbol{y}_{1} + \boldsymbol{y}_{2}\sqrt{-1})$$
(27)

- ⁶³⁵ Then by Lemma B.2, we have $(x_1 + x_2\sqrt{-1}) \sim (y_1 + y_2\sqrt{-1})$, which is essentially $[x_1 \ x_2] \sim$ ⁶³⁶ $[y_1 \ y_2]$.
- Now we are ready to prove the injectivity of the LLE layer.
- **Theorem D.1.** Suppose $\phi : \mathbb{R}^D \to \mathbb{R}^L$ admits the form of Eq. (3) and $\mathbf{W} = [\cdots \quad \mathbf{w}_{i,j,p,q} \quad \cdots] \in \mathbb{R}^{D \times L}$, $i \in [D], j \in [D], p \in [N], q \in [p+1]$ is constructed as follows:

$$\boldsymbol{w}_{i,j,p,q} = (q-1)\boldsymbol{e}_i + (p-q+1)\boldsymbol{e}_j,$$
(28)

where e_i is the *i*-th canonical basis. Then $\sum_{i=1}^{N} \phi(x^{(i)})$ is injective (cf. Definition A.5).

Proof. First of all, notice that $L = D^2 \sum_{p=1}^{N} (p+1) = D^2 (N+3)N/2 \le 2N^2 D^2$. According to Eq. 8, we can rewrite for $\forall i \in [D], j \in [D], p \in [N], q \in [p+1]$

$$\phi(\boldsymbol{x})_{i,j,p,q} = \exp(\boldsymbol{w}_{i,j,p,q}^{\top} \log(\boldsymbol{x})) = \prod_{k=1}^{D} \boldsymbol{x}_{k}^{\boldsymbol{w}_{i,j,p,q,k}} = \boldsymbol{x}_{i}^{q-1} \boldsymbol{x}_{j}^{p-q+1},$$
(29)

Then for $\boldsymbol{X}, \boldsymbol{X'} \in \mathbb{R}^{N \times D}$, $\sum_{i \in [N]} \phi(\boldsymbol{x}^{(i)} = \sum_{i \in [N]} \phi(\boldsymbol{x'}^{(i)})$ implies $\sum_{i \in [N]} \boldsymbol{x}_i^{q-1} \boldsymbol{x}_j^{p-q+1} = \sum_{i \in [N]} \boldsymbol{x'}_i^{q-1} \boldsymbol{x'}_j^{p-q+1}$ for $\forall i \in [D], j \in [D], p \in [N], q \in [p+1]$. By Lemma 4.8, we have for $[\boldsymbol{x}_i \ \boldsymbol{x}_j] \sim [\boldsymbol{x'}_i \ \boldsymbol{x'}_j]$ for $\forall i \in [D], j \in [D]$. Finally, Lemma 4.9 directly yields $\boldsymbol{X} \sim \boldsymbol{X'}$. \Box

646 D.2 Continuity

⁶⁴⁷ The proof idea of continuity for LLE layer shares the same outline with the LP layer.

Corollary D.2. Consider the mapping $\widetilde{\Phi_N} = \psi_N^{-1} \circ \tau : (\mathbb{R}^{N(N+3)/2}, d_\infty) \to (\mathbb{C}^N / \sim, d_F)$, where $\tau : (\mathbb{R}^{N(N+3)/2}, d_\infty) \to (\mathbb{C}^N, d_\infty)$ is a linear mapping that combines sum of monomials to polynomials following Eq. (27). Construct the mapping $\widehat{\Phi_N} : (\mathbb{R}^L, d_\infty) \to (\mathbb{C}^N / \sim, d_F)^{D^2}$ based on $\widetilde{\Phi_N}$:

$$\widehat{\Phi_N}(\mathbf{Z}) = \begin{bmatrix} \overline{\Phi_N}(\mathbf{z}_1) & \cdots & \overline{\Phi_N}(\mathbf{z}_{D^2}) \end{bmatrix},$$
(30)

where $\mathbf{Z} = \begin{bmatrix} \mathbf{z}_1^\top & \cdots & \mathbf{z}_{D^2}^\top \end{bmatrix}^\top \in \mathbb{R}^L$, $\mathbf{z}_i \in \mathbb{R}^{(N(N+3)/2)}, \forall i \in [D^2]$. We denote the induced product metric over $(\mathbb{C}^N/\sim, d_F)^{D^2}$ as $d_F^{D^2} : (\mathbb{C}^N/\sim)^{D^2} \times (\mathbb{C}^N/\sim)^{D^2} \to \mathbb{R}_{\geq 0}$:

$$d_F^{D^2}(\boldsymbol{Z}, \boldsymbol{Z'}) = \max_{i \in [D^2]} d_F(\boldsymbol{z}_i, \boldsymbol{z'}_i).$$
(31)

Then the mapping $\widehat{\Phi_N}$ maps any bounded set in (\mathbb{R}^L, d_∞) to a bounded set in $(\mathbb{C}^N / \sim, d_F)^{D^2}$.

Proof. Since τ is a linear mapping, any $Z, Z' \in \mathbb{R}^L$ such that $d_{\infty}(z_i, z'_i) \leq C_1, \forall i \in [D^2]$ for some constant $C_1 \geq 0$, then $d_F(\tau(z_i), \tau(z'_i)) \leq C_2$ for some constant $C_2 \geq 0$. By Lemma B.3, ψ_N^{-1} maps any bounded set in $(\mathbb{C}^N, d_{\infty})$ to a bounded set in $(\mathbb{C}^N / \sim, d_F)$. This is $d_F(\widetilde{\Phi}_N(z_i), \widetilde{\Phi}_N(z'_i)) \leq C_2, \forall i \in [D^2]$ for some constant $C_2 \geq 0$. Finally, we have:

$$d_F^{D^2}(\widehat{\Phi_N}(\boldsymbol{Z}), \widehat{\Phi_N}(\boldsymbol{Z'}) = \max_{i \in [D^2]} d_F(\widetilde{\Phi_N}(\boldsymbol{z}_{i,j}), \widetilde{\Phi_N}(\boldsymbol{z'}_{i,j})) \le C_2,$$

659 which is also bounded above.

Theorem D.3. Suppose ϕ admits the form of Eq. (3) and follows the construction in Theorem D.1, then the inverse of LLE-induced sum-pooling $\sum_{i=1}^{N} \phi(\mathbf{x}^{(i)})$ is continuous.

Proof. It is sufficient to verify three conditions in Lemma 4.10. First of all, we denote the inverse of $\Psi(\mathbf{X}) = \sum_{i=1}^{N} \phi(\mathbf{x}^{(i)})$, denoted as $\Psi^{-1} : (\mathbb{R}^{L}, d_{\infty}) \to (\mathbb{R}^{N \times D} / \sim, d_{F})$, which exists thanks to Theorem D.1. By Lemma A.4, any closed and bounded subset of $(\mathbb{R}^{N \times D} / \sim, d_{F})$ is compact. Obviously, $\Psi(\mathbf{X})$ is continuous. Then it remains to show the condition (**c**) in Lemma 4.10. Similar to Theorem C.10, we decompose Ψ^{-1} into two mappings following the clue of proving its existence:

$$(\mathbb{R}^{NK}, d_{\infty}) \xrightarrow{\Phi_N} (\mathbb{C}^N / \sim, d_F)^{D^2} \xrightarrow{\pi} (\mathbb{R}^{N \times D} / \sim, d_F) ,$$

where $\widehat{\Phi_N}$ is defined in Eq. (30) and π exists due to Theorem D.1. Also according to our construction in Theorem D.1, for any $\mathbf{Z} \subset (\mathbb{C}^N/\sim, d_F)^{D^2}$ and $\forall i, j \in [D]$, there exists $k \in [D^2]$ such that $(\pi(\mathbf{z})_i + \pi(\mathbf{z})_j \sqrt{-1}) \sim \mathbf{z}_k$. Therefore, $\forall \mathbf{Z}, \mathbf{Z'} \in (\mathbb{C}^N/\sim, d_F)^{D^2}$ such that $d_F^{D^2}(\mathbf{Z}, \mathbf{Z'}) \leq C$ for some constant C > 0, we have:

$$d_F(\pi(\boldsymbol{Z}), \pi(\boldsymbol{Z'})) \le \max_{i \in [D^2]} d_F(\boldsymbol{z}_i, \boldsymbol{z'}_i) \le d_F^{D^2}(\boldsymbol{Z}, \boldsymbol{Z'}) \le C,$$
(32)

which implies π maps every bounded set in $(\mathbb{C}^N/\sim, d_F)^K$ to a bounded set in $(\mathbb{R}^{N\times D}/\sim, d_F)$. Now we conclude the proof by the following chain of argument:

$$\mathcal{Z} \subseteq (\mathbb{R}^{NK}, d_{\infty})$$
 is bounded $\xrightarrow{\text{Corollary D.2}} \widehat{\Phi_N}(\mathcal{Z})$ is bounded $\xrightarrow{\text{Eq. (32)}} \pi \circ \widehat{\Phi_N}(\mathcal{Z})$ is bounded

673

674 E Extension to Permutation Equivariance

In this section, we prove Theorem 5.1, the extension of Theorem 3.1 to equivariant functions, following the similar arguments with [7]:

Lemma E.1 ([7,34]). $f : \mathbb{R}^{N \times D} \to \mathbb{R}^N$ is a permutation-equivariant function if and only if there is a function $\rho : \mathbb{R}^{N \times D} \to \mathbb{R}$ that is permutation invariant to the last N - 1 entries, such that f(\mathbf{Z})_i = $\rho(\mathbf{z}^{(i)}, \underbrace{\mathbf{z}^{(i+1)}, \cdots, \mathbf{z}^{(N)}, \cdots, \mathbf{z}^{(i-1)}}_{N-1})$ for any $i \in [N]$.

Proof. (Sufficiency) Define $\pi : [N] \to [N]$ be an index mapping associated with the permutation matrix $\boldsymbol{P} \in \Pi(N)$ such that $\boldsymbol{P}\boldsymbol{Z} = [\boldsymbol{z}^{(\pi(1))}, \cdots, \boldsymbol{z}^{(\pi(N))}]^{\top}$. Then we have:

$$f\left(\boldsymbol{z}^{(\pi(1))}, \cdots, \boldsymbol{z}^{(\pi(N))}\right)_{i} = \rho\left(\boldsymbol{z}^{(\pi(i))}, \boldsymbol{z}^{(\pi(i+1))}, \cdots, \boldsymbol{z}^{(\pi(N))}, \cdots, \boldsymbol{z}^{(\pi(i-1))}\right).$$

Since $\rho(\cdot)$ is invariant to the last N-1 entries, it can shown that:

$$f(\boldsymbol{P}\boldsymbol{Z})_i =
ho\left(\boldsymbol{z}^{(\pi(i))}, \boldsymbol{z}^{(\pi(i+1))}, \cdots, \boldsymbol{z}^{(\pi(N))}, \cdots, \boldsymbol{z}^{(\pi(i-1))}\right) = f(\boldsymbol{Z})_{\pi(i)}.$$

(Necessity) Given a permutation-equivariant function $f : \mathbb{R}^{N \times D} \to \mathbb{R}^N$, we first expand it to the following form: $f(\mathbf{Z})_i = \rho_i(\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(N)})$. Permutation-equivariance means $\rho_{\pi(i)}(\boldsymbol{z}^{(1)}, \dots, \boldsymbol{z}^{(N)}) = \rho_i(\boldsymbol{z}^{\pi(1)}, \dots, \boldsymbol{z}^{\pi(N)})$ for any permutation mapping π . Suppose given an index $i \in [N]$, consider any permutation $\pi : [N] \to [N]$ such that $\pi(i) = i$. Then, we have:

$$\rho_i\left(\boldsymbol{z}^{(1)},\cdots,\boldsymbol{z}^{(i)},\cdots,\boldsymbol{z}^{(N)}\right) = \rho_{\pi(i)}\left(\boldsymbol{z}^{(1)},\cdots,\boldsymbol{z}^{(i)},\cdots,\boldsymbol{z}^{(N)}\right) = \rho_i\left(\boldsymbol{z}^{(\pi(1))},\cdots,\boldsymbol{z}_i,\cdots,\boldsymbol{z}^{(\pi(N))}\right)$$

which implies $\rho_i : \mathbb{R}^{N \times D} \to \mathbb{R}$ must be invariant to the N - 1 elements other than the *i*-th element. Now, consider a permutation π where $\pi(1) = i$. Then we have:

$$\rho_i\left(\boldsymbol{z}^{(1)}, \boldsymbol{z}^{(2)}, \cdots, \boldsymbol{z}^{(N)}\right) = \rho_{\pi(1)}\left(\boldsymbol{z}^{(1)}, \boldsymbol{z}^{(2)}, \cdots, \boldsymbol{z}^{(N)}\right) = \rho_1\left(\boldsymbol{z}^{(\pi(1))}, \boldsymbol{z}^{(\pi(2))}, \cdots, \boldsymbol{z}^{(\pi(N))}\right)$$
$$= \rho_1\left(\boldsymbol{z}^{(i)}, \boldsymbol{z}^{(i+1)}, \cdots, \boldsymbol{z}^{(N)}, \cdots, \boldsymbol{z}^{(i-1)}\right),$$

where the last equality is due to the invariance to N-1 elements, stated beforehand. This implies two results. First, for all i, $\rho_i(\boldsymbol{z}^{(1)}, \boldsymbol{z}^{(2)}, \cdots, \boldsymbol{z}^{(i)}, \cdots, \boldsymbol{z}^{(N)}), \forall i \in [N]$ should be written in terms of $\rho_1(\boldsymbol{z}^{(i)}, \boldsymbol{z}^{(i+1)}, \cdots, \boldsymbol{z}^{(N)}, \cdots, \boldsymbol{z}^{(i-1)})$. Moreover, ρ_1 is permutation invariant to its last N-1entries. Therefore, we just need to set $\rho = \rho_1$ and broadcast it accordingly to all entries. We conclude the proof.

Proof of Theorem 5.1 [7]. Sufficiency can be shown by verifying the equivariance. We conclude the proof by showing the necessity with Lemma E.1. First we rewrite any permutation equivariant function $f(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}) : \mathbb{R}^{N \times D} \to \mathbb{R}^N$ as:

$$f\left(\boldsymbol{x}^{(1)},\cdots,\boldsymbol{x}^{(N)}\right)_{i}=\tau\left(\boldsymbol{x}^{(i)},\boldsymbol{x}^{(i+1)},\cdots,\boldsymbol{x}^{(N)},\cdots,\boldsymbol{x}^{(i-1)}\right),$$
(33)

where π is invariant to the lask N - 1 elements, according to Lemma E.1. Given ϕ with either LP or LLE architectures, $\Psi(\mathbf{X}) = \sum_{i=1}^{N} \phi(\mathbf{x}^{(i)})$ is injective and has continuous inverse if:

• $L \in [N(D+1), N^5D^2]$ when ϕ admits LP architecture. (By Theorem C.8 and C.10).

• $L \in [ND, 2N^2D^2]$ when ϕ admits LLE architecture. (By Theorem D.1 and D.3).

The proof proceeds by letting $\rho : \mathbb{R}^D \times \mathbb{R}^L \to \mathbb{R}$ take the form $\rho(\boldsymbol{x}, \boldsymbol{z}) = \tau(\boldsymbol{x}, \Phi^{-1}(\boldsymbol{z} - \phi(\boldsymbol{x})))$, and observe that:

$$\begin{aligned} \tau\left(\boldsymbol{x}^{(i)}, \boldsymbol{x}^{(i+1)}, \cdots, \boldsymbol{x}^{(N)}, \cdots, \boldsymbol{x}^{(i-1)}\right) &= \tau\left(\boldsymbol{x}^{(i)}, \Phi^{-1} \circ \Phi(\boldsymbol{x}^{(i+1)}, \cdots, \boldsymbol{x}^{(N)}, \cdots, \boldsymbol{x}^{(i-1)})\right) \\ &= \tau\left(\boldsymbol{x}^{(i)}, \Phi^{-1}\left(\Phi\left(\boldsymbol{x}^{(1)}, \cdots, \boldsymbol{x}^{(N)}\right) - \phi(\boldsymbol{x}^{(i)})\right)\right) \\ &= \rho\left(\boldsymbol{x}^{(i)}, \sum_{i=1}^{N} \phi(\boldsymbol{x}^{(i)})\right) \end{aligned}$$

703

704 **F** Extension to Complex Numbers

In this section, we formally introduce the nature extension of our Theorem 3.1 to the complex numbers:

Corollary F.1 (Extension to Complex Domain). For any permutation-invariant function f: $\mathcal{K}^{N \times D} \to \mathbb{R}, \ \mathcal{K} \subseteq \mathbb{C}, \ there \ exists \ continuous \ functions \ \phi : \mathbb{C}^D \to \mathbb{R}^L \ and \ \rho : \mathbb{R}^L \to \mathbb{C} \ such$ that $f(\mathbf{X}) = \rho\left(\sum_{i \in [N]} \phi(\mathbf{x}^{(i)})\right)$ for every $j \in [N]$, where $L \in [2N(D+1), 4N^5D^2]$ when ϕ admits LP architecture, and $L \in [2ND, 8N^2D^2]$ when ϕ admits LLE architecture ($\mathcal{K} \in \mathbb{C}_{>0}$).

Proof. We let ϕ first map complex features $\mathbf{x}^{(i)} \in \mathbb{C}^D, \forall i \in [N]$ to real features $\tilde{\mathbf{x}}^{(i)} = [\Re(\mathbf{x}^{(i)})^\top \quad \Im(\mathbf{x}^{(i)})^\top] \in \mathbb{R}^{2D}, \forall i \in [N]$ by divide the real and imaginary parts into separate channels, then utilize either LP or LLE embedding layer to map $\tilde{\mathbf{x}}^{(i)}$ to the latent space. The upper bounds of desired latent space dimension are scaled by 4 for both architectures due to the quadratic dependence on D. Then the same proof of Theorems C.8, C.10, D.1, and D.3 applies.

716 **G** Connection to Unlabled Sensing

Unlabeled sensing [56], also known as linear regression without correspondence [55, 57–59], solves 717 the linear system y = PAx with a given measurement matrix $A \in \mathbb{R}^{M \times N}$ and an unknown 718 permutation $P \in \Pi(M)$. [55, 56] show that as long as A is over-determinant $(M \ge 2N)$, such 719 problem is well-posed (i.e., has a unique solution) for almost all cases. Unlabeled sensing shares 720 the similar structure with our LP embedding layer in which a linear layer lifts the feature space 721 to a higher-dimensional ambient space, ensuring the solvability of alignment across each channel. 722 However, our invertibility is defined between the set and embedding spaces, which differs from 723 724 exact recovery of unknown variables desired in unlabeled sensing [55]. In fact, the well-posedness 725 of unlabeled PCA [54], studying matrix completion with shuffle perturbations, shares the identical definition with our injectivity. But it is noteworthy that the results in [54] are only drawn over a726 dense subset of the input space, while ours are stronger in considering all possible inputs. Hence, our 727 theory could potentially bring new insights into the field of unlabeled sensing, which may be of an 728 independent interest. 729

730 **References**

- [1] Yann LeCun, Yoshua Bengio, et al. Convolutional networks for images, speech, and time series.
 The handbook of brain theory and neural networks, 3361(10):1995, 1995.
- [2] Taco Cohen and Max Welling. Group equivariant convolutional networks. In *International conference on machine learning*, pages 2990–2999. PMLR, 2016.
- [3] Michael M Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst.
 Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42, 2017.
- [4] Risi Kondor and Shubhendu Trivedi. On the generalization of equivariance and convolution
 in neural networks to the action of compact groups. In *International Conference on Machine Learning*, pages 2747–2755, 2018.
- [5] Haggai Maron, Heli Ben-Hamu, Nadav Shamir, and Yaron Lipman. Invariant and equivariant
 graph networks. In *International Conference on Learning Representations (ICLR)*, 2018.
- [6] Alexander Bogatskiy, Brandon Anderson, Jan Offermann, Marwah Roussi, David Miller, and
 Risi Kondor. Lorentz group equivariant neural network for particle physics. In *International Conference on Machine Learning*, pages 992–1002. PMLR, 2020.
- [7] Peihao Wang, Shenghao Yang, Yunyu Liu, Zhangyang Wang, and Pan Li. Equivariant hyper graph diffusion neural operators. In *International Conference on Learning Representations* (*ICLR*), 2023.
- [8] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017.
- [9] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov,
 and Alexander J Smola. Deep sets. In *Advances in Neural Information Processing Systems* (*NeurIPS*), 2017.
- [10] Vinicius Mikuni and Florencia Canelli. Point cloud transformers applied to collider physics.
 Machine Learning: Science and Technology, 2(3):035027, 2021.
- [11] Huilin Qu and Loukas Gouskos. Jet tagging via particle clouds. *Physical Review D*, 101(5):056019, 2020.
- [12] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In
 Proceedings of the IEEE/CVF international conference on computer vision, pages 16259–16268, 2021.
- [13] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosiorek, Seungjin Choi, and Yee Whye Teh.
 Set transformer: A framework for attention-based permutation-invariant neural networks. In
 International conference on machine learning, pages 3744–3753. PMLR, 2019.
- [14] Yan Zhang, Jonathon Hare, and Adam Prugel-Bennett. Deep set prediction networks. *Advances in Neural Information Processing Systems*, 32, 2019.

- [15] Yan Zhang, Jonathon Hare, and Adam Prügel-Bennett. Fspool: Learning set representations
 with featurewise sort pooling. In *International Conference on Learning Representations*, 2020.
- [16] Aditya Grover, Eric Wang, Aaron Zweig, and Stefano Ermon. Stochastic optimization of sorting
 networks via continuous relaxations. In *International Conference on Learning Representations*,
 2020.
- [17] Christopher Morris, Martin Ritzert, Matthias Fey, William L Hamilton, Jan Eric Lenssen,
 Gaurav Rattan, and Martin Grohe. Weisfeiler and leman go neural: Higher-order graph neural
- networks. In *the AAAI Conference on Artificial Intelligence*, volume 33, pages 4602–4609, 2019.
- [18] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019.
- [19] Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Veličković. Principal
 neighbourhood aggregation for graph nets. *Advances in Neural Information Processing Systems*,
 33:13260–13271, 2020.
- [20] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini.
 The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2008.
- [21] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large
 graphs. In *Advances in Neural Information Processing Systems*, 2017.
- [22] Haggai Maron, Ethan Fetaya, Nimrod Segol, and Yaron Lipman. On the universality of invariant networks. In *International conference on machine learning*, pages 4363–4371. PMLR, 2019.
- [23] Edward Wagstaff, Fabian Fuchs, Martin Engelcke, Ingmar Posner, and Michael A Osborne.
 On the limitations of representing functions on sets. In *International Conference on Machine Learning*, pages 6487–6494. PMLR, 2019.
- [24] Edward Wagstaff, Fabian B Fuchs, Martin Engelcke, Michael A Osborne, and Ingmar Posner.
 Universal approximation of functions on sets. *Journal of Machine Learning Research*, 23(151):1–
 56, 2022.
- [25] Nimrod Segol and Yaron Lipman. On universal equivariant set networks. In *International Conference on Learning Representations (ICLR)*, 2020.
- [26] Aaron Zweig and Joan Bruna. Exponential separations in symmetric neural networks. *arXiv preprint arXiv:2206.01266*, 2022.
- [27] George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.
- [28] Kurt Hornik, Maxwell Stinchcombe, Halbert White, et al. Multilayer feedforward networks are
 universal approximators. *Neural Networks*, 2(5):359–366, 1989.
- [29] Shupeng Gui, Xiangliang Zhang, Pan Zhong, Shuang Qiu, Mingrui Wu, Jieping Ye, Zhengdao
 Wang, and Ji Liu. Pine: Universal deep embedding for graph nodes via partial permutation
 invariant set functions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,
 44(2):770–782, 2021.
- [30] Nicolas Bourbaki. *Éléments d'histoire des mathématiques*, volume 4. Springer Science &
 Business Media, 2007.
- [31] David Rydh. A minimal set of generators for the ring of multisymmetric functions. In *Annales de l'institut Fourier*, volume 57, pages 1741–1769, 2007.
- [32] Adam Santoro, David Raposo, David G Barrett, Mateusz Malinowski, Razvan Pascanu, Peter
 Battaglia, and Timothy Lillicrap. A simple neural network module for relational reasoning.
 Advances in neural information processing systems, 30, 2017.
- [33] R. L. Murphy, B. Srinivasan, V. Rao, and B. Ribeiro. Janossy pooling: Learning deep
 permutation-invariant functions for variable-size inputs. In *International Conference on Learn- ing Representations (ICLR)*, 2018.
- [34] Akiyoshi Sannai, Yuuki Takai, and Matthieu Cordonnier. Universal approximations of permutation invariant/equivariant functions by deep neural networks. *arXiv preprint arXiv:1903.01939*, 2019.
- [35] Branko Ćurgus and Vania Mascioni. Roots and polynomials as homeomorphic spaces. *Expositiones Mathematicae*, 24(1):81–95, 2006.

- [36] Andrew R Barron. Universal approximation bounds for superpositions of a sigmoidal function.
 IEEE Transactions on Information theory, 39(3):930–945, 1993.
- [37] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network
 learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 2015.
- [38] Fernando Q Gouvêa. Was cantor surprised? *The American Mathematical Monthly*, 118(3):198–209, 2011.
- [39] Dmitry Yarotsky. Error bounds for approximations with deep relu networks. *Neural Networks*,
 94:103–114, 2017.
- [40] Shiyu Liang and R Srikant. Why deep neural networks for function approximation? In International Conference on Learning Representations, 2017.
- [41] Joe Kileel, Matthew Trager, and Joan Bruna. On the expressive power of deep polynomial neural networks. *Advances in neural information processing systems*, 32, 2019.
- [42] Nadav Cohen, Or Sharir, and Amnon Shashua. On the expressive power of deep learning: A
 tensor analysis. In *Conference on learning theory*, pages 698–728. PMLR, 2016.
- [43] Maithra Raghu, Ben Poole, Jon Kleinberg, Surya Ganguli, and Jascha Sohl-Dickstein. On the
 expressive power of deep neural networks. In *international conference on machine learning*,
 pages 2847–2854. PMLR, 2017.
- [44] Dmitry Yarotsky. Universal approximations of invariant maps by neural networks. *Constructive Approximation*, 55(1):407–474, 2022.
- [45] Ding-Xuan Zhou. Universality of deep convolutional neural networks. *Applied and computational harmonic analysis*, 48(2):787–794, 2020.
- [46] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural
 message passing for quantum chemistry. In *International Conference on Machine Learning* (*ICML*), 2017.
- [47] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*, 2017.
- [48] Nicolas Keriven and Gabriel Peyré. Universal invariant and equivariant graph neural networks.
 In Advances in Neural Information Processing Systems, pages 7090–7099, 2019.
- [49] Zhengdao Chen, Soledad Villar, Lei Chen, and Joan Bruna. On the equivalence between graph
 isomorphism testing and function approximation with gnns. In *Advances in Neural Information Processing Systems*, pages 15868–15876, 2019.
- [50] Zhengdao Chen, Lei Chen, Soledad Villar, and Joan Bruna. Can graph neural networks count
 substructures? volume 33, 2020.
- [51] Waïss Azizian and Marc Lelarge. Expressive power of invariant and equivariant graph neural
 networks. In *ICLR 2021-International Conference on Learning Representations*, 2021.
- [52] Andreas Loukas. What graph neural networks cannot learn: depth vs width. In *International Conference on Learning Representations*, 2020.
- [53] Manolis C Tsakiris. Low-rank matrix completion theory via plücker coordinates. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [54] Yunzhen Yao, Liangzu Peng, and Manolis Tsakiris. Unlabeled principal component analysis.
 Advances in Neural Information Processing Systems, 34:30452–30464, 2021.
- [55] Manolis Tsakiris and Liangzu Peng. Homomorphic sensing. In *International Conference on Machine Learning*, pages 6335–6344. PMLR, 2019.
- Isomorphic and Second Se
- [57] Daniel J Hsu, Kevin Shi, and Xiaorui Sun. Linear regression without correspondence. *Advances in Neural Information Processing Systems*, 30, 2017.
- [58] Manolis C Tsakiris, Liangzu Peng, Aldo Conca, Laurent Kneip, Yuanming Shi, and Hayoung
 Choi. An algebraic-geometric approach for linear regression without correspondences. *IEEE Transactions on Information Theory*, 66(8):5130–5144, 2020.

- [59] Liangzu Peng and Manolis C Tsakiris. Linear regression without correspondences via concave minimization. *IEEE Signal Processing Letters*, 27:1580–1584, 2020.