

Table 2: Hyperparameters for our 16M models.

Hyperparameter	PCQ	COCO	CODE	FLOW	MST	BRIDGES
Learning rate	1e-4	4e-4/3.5e-4	1e-4/1.5e-4	1e-4/2e-4	3e-4	2e-4/1e-4
Batch size	256	32	32	256	256	256
Optimizer	AdamW	AdamW	AdamW	AdamW	AdamW	AdamW
Grad. clip norm	1	1	1	1	1	1
Num. layers	16	16	16	16	16	16
Hidden dim.	384	384	384	384	384	384
Num. heads	16	16	16	16	16	16
Activation	ReLU	ReLU	ReLU	ReLU	ReLU	ReLU
RWSE/RRWP Steps	32	32	32	16	16	16
Num eigvals/eigvecs	32	32	32	16	16	16
Hidden dim. RWSE/RRWP	768	768	768	768	768	768
Hidden dim. LPE/SPE	384	384	384	384	384	384
Num. layers ϕ	2	2	2	2	2	2
Num. layers ρ	2	2	2	2	2	2
GNN type ρ (SPE)	GIN	GIN	GIN	GIN	GIN	GIN
Edge encoder	MLP	MLP	MLP	MLP	MLP	MLP
Weight decay	0.1	0.1	0.1	0.1	0.1	0.1
Dropout	0.1	0.1	0.1	0.1	0.1	0.1
Attention dropout	0.1	0.1	0.1	0.1	0.1	0.1
#Steps	2M	1M	200k	11718	11718	11718
#Warmup steps	20k	10k	2k	118	118	118

A Experimental details

Here we present hyperparameter choices, architecture design, and dataset selections for the empirical evaluation of our GT architecture.

A.1 Data sources and licenses

PCQM4MV2 is available at <https://ogb.stanford.edu/docs/lsc/pcqm4mv2/> under a CC BY 4.0 license. OGB-Code2 is available at <https://ogb.stanford.edu/docs/graphprop/#ogbg-code2> under a MIT license. The COCO-SP and PASCAL-VOC-SP datasets as part of the LRGB benchmark [Dwivedi et al., 2022] are available at <https://github.com/vijaydwivedi75/lrgb> under a CC BY 4.0 license. Statistics for all datasets, including the algorithmic reasoning datasets, are available in Table 4

A.2 Hyperparameters

Table 2 and Table 3 give an overview of the hyperparameters used for models highlighted in our work. Considering the large number of hyperparameters and scale of tasks, we did not perform a grid search or any other type of large-scale hyperparameter optimization. Nonetheless, we swept the learning rate for each task and model size. Across the experiments, we select the hyperparameters based on the best validation score and then evaluate on the test set. We search for suitable learning rates on the 16M models to determine the models we select for scaling. Due to the increased computational demand, we then reduce the learning rate for the 100M models.

For PCQ, we set the learning rate to 1e-4 after sweeping the learning rate over the set $\{7e-5, 1e-4, 3e-4\}$. Further, we set the learning rate for COCO to 4e-4 for eigen-information-based embeddings and 3.5e-4 for RWSE after sweeping over $\{7e-5, 1e-4, 3e-4, 4e-4\}$. For CODE, we reduce the learning rate to 1.5e-4 (RWSE) with the same sweep as with PCQ, and across all MST runs, we keep the learning rate at 3e-4 for the 16M models, reducing the learning rate to 1e-4 (LPE) and 7e-5 (RWSE) for the 90M model and 7e-5 for the 160M model at MST with the same initial sweep as for COCO. Furthermore, for FLOW we set the learning rate to 1e-4 (LPE,SPE,NoPE) and 2e-4 (RWSE, RRWP) respectively. We obtain similar results for BRIDGES with 1e-4 (RWSE, LPE)

Table 3: Hyperparameters for our 90M and 160M models.

Hyperparameter	PCQ(90M)	MST(90M)	MST(160M)
Learning rate	1e-4	1e-4/7e-5	7e-5
Batch size	256	256	256
Optimizer	AdamW	AdamW	AdamW
Grad. clip norm	1	1	1
Num. layers	24	24	24
Hidden dim.	768	768	1024
Num. heads	16	16	16
Activation	ReLU	ReLU	ReLU
RWSE/RRWP steps	32	32	32
Num. eigvals/eigvecs	32	32	32
Hidden dim. RWSE/RRWP	768	768	768
Hidden dim. LPE/SPE	384	384	384
Num. layers ϕ	2	2	2
Num. layers ρ	2	2	2
GNN type ρ (SPE)	GIN	GIN	GIN
Edge encoder	MLP	MLP	MLP
Weight decay	0.1	0.1	0.1
Dropout	0.1	0.1	0.1
Attention dropout	0.1	0.1	0.1
#Steps	2M	11718	11718
#Warmup steps	20k	118	118

Table 4: Dataset statistics.

Statistic	PCQ	COCO	CODE	FLOW	MST	BRIDGES
# Graphs	3,746M	123,286	452,741	1M	1M	1M
# Avg. Nodes	14.13	476.88	125.2	16/64	16/64	16/64
# Avg. Edges	14.56	3,815.08	124.2	48.11/213.586	31.66/209.34	48.46/ 395.02
Prediction level	graph	node	graph	graph	edge	edge
Metric	MAE	F1	F1	MAE	F1	F1

and $2e-4$ (RRWP,SPE,NoPE) as the learning rate using the initial sweep from PCQ. In addition, we evaluated each PE with $\{8,32\}$ random walk steps or eigenvectors and $\{4,16\}$ for algorithmic reasoning tasks. Regarding PE encoder design, we selected a simple MLP architecture with two layers where applicable. However, we use the same number of layers, heads, and embedding dimension for each dataset for our transformer architecture, thereby not changing the architecture. Otherwise, we follow previous literature for initial hyperparameter choices, namely the GraphGPS [Rampásek et al., 2022], GRIT [Ma et al., 2023], and Graphormer [Ying et al., 2021a] papers. We used an AdamW optimizer for each experiment with $\beta_1 = 0.9$ and $\beta_2 = 0.999$. Further, the learning rate scheduler is given by a cosine annealing learning rate scheduler with the warm-up steps set as one percent of the total number of steps. Additionally, we use an L1 loss for regression targets and a cross-entropy loss for classification targets, except CODE, where we use the proposed loss function. In Table 5 and Table 6, we further report runtimes and memory usage of all models evaluated in our work.

A.3 Architecture

In the following, we showcase the implementation of our GT architecture and the injection of PE information into the attention mechanism. We consider an Encoder, Processor, Decoder architecture with additional preprocessing for the PEs and graph-specific features.

Preprocessing First, we preprocess each dataset to include the respective eigenvalues and eigenvectors of the graph Laplacian and the powers of the random walk matrices. These are then applied to the respective PE encoder, consisting of an MLP with two layers, casting the PE features to the embedding dimension. We consider transformed graphs for edge-level tasks, such as BRIDGES and MST, as a special case. In this case, the graph is converted to constitute the edge-level graph corresponding to the original graph. Then, node features and PE features are computed on this transformed graph. For the GDT, we provide a maximum context size depending on the dataset. Tokens exceeding the context size are then removed.

Encoder The node and edge features of graphs in each dataset are then applied to a linear layer, mapping them to the embedding dimension. These feature embeddings are specific to each dataset and embed graph-specific features. For CODE we consider additional preprocessing steps, as described by [Hu et al., 2021] to derive the respective graph structure. Further, we add a [cls] token as it is a standard practice to read out graph-level representations [Ying et al., 2021b].

Processor Following our description of the GDT architecture as shown in Section 2, a single layer in the GDT architecture computes the expression shown in Definition 8 using GELU as a nonlinearity. The absolute PEs are added to the node embeddings before the initial layer, in the case of relative PEs to the attention bias B . The GT layer then computes full multi-head scaled-dot-product attention over node-level tokens, adding B to the unnormalized attention matrix before applying softmax. We refer to Appendix B.1 for a detailed discussion. From this representation, including node and edge features, relative and absolute PEs, and the embedded graph structure, the processor computes a representation of node and graph-level features. We then stack multiple GT layers together: 12 for the 16M model and 24 for the 100M model.

Decoder After the last layer, an MLP decoder with two layers is applied to provide the prediction head of the model. Since each dataset has its prediction target, we provide a decoder for each dataset mapping the last layer output to the prediction target, where $\mathbf{W}_1 \in \mathbb{R}^{d \times d}$ and $\mathbf{W}_2 \in \mathbb{R}^{d \times o}$ are learnable weight matrices and o is the respective output dimension for each task, i.e.,

$$\mathbf{W}_2 \text{LayerNorm}(\text{GELU}(\mathbf{W}_1 x)).$$

For clarity, we omit bias terms throughout this section. Each result is then passed to the respective loss function to compute the gradient step.

A.4 Algorithmic reasoning data

For our synthetic experiments, we evaluate on three out-of-distribution algorithmic reasoning tasks derived from the CLRS benchmark [Velickovic et al., 2022] for a total of 100M tokens. These tasks assess size generalization in a controlled synthetic setting with randomly generated graphs. Unlike CLRS, we do not train models with intermediate algorithmic steps. Here, we describe graph generation and each algorithmic reasoning task in detail.

Graph generation We develop a heuristic graph generation method that leads to graphs with desirable problem-specific properties, such as a reasonable distribution of shortest-path lengths or number of bridges. Concretely, we begin by sampling an Erdos-Renyi graph G with n nodes and edge probability p and denote the connected components of G with C_1, \dots, C_m . For each $i \in [m]$, we randomly choose a component C_j with $j \neq i$. Then, we select random nodes $v \in C_i$ and $w \in C_j$ and augment G with the edge (i, j) . We repeat this process several K times. We select parameters p and K for each task based on problem-specific characteristics. We detail these choices in the task descriptions, which we provide next.

Maximum flow In FLOW, the task is to predict the maximum flow value in an edge-weighted directed graph. The task uses discrete node features indicating whether a node is either the source of the flow, the sink of the flow, or neither. The task uses the flow capacity between two nodes as continuous scalar-valued edge features. FLOW is a graph-level regression task.

Minimum spanning tree In MST, the task is to predict the set of edges that forms the minimum spanning tree (MST) in an edge-weighted graph with mutually distinct edge weights, ensuring the

Table 5: 16M/90M/160M models runtime results of a single step, averaged across 1 000 steps. Each value is given in seconds/step.

PE	#Param.	PCQ	COCO	CODE	FLOW	MST	BRIDGES
NoPE	16M	0.079	0.105	0.0921	0.071	0.096	0.072
LPE	16M	0.084	0.109	0.104	0.071	0.106	0.088
SPE	16M	0.091	0.123	0.105	0.059	0.106	0.085
RWSE	16M	0.072	0.101	0.102	0.071	0.091	0.088
RRWP	16M	0.137	-	-	0.066	0.1452	0.101
LPE	90M	0.167	-	-	-	0.184	-
RWSE	90M	0.159	-	-	-	0.186	-
LPE	160M	0.219	-	-	-	0.246	-
RWSE	160M	0.219	-	-	-	0.237	-

uniqueness of the MST. The task uses the weight of each edge as continuous scalar-valued edge features. MST is a binary edge classification task where the class label indicates whether an edge is contained within the MST.

Bridges In BRIDGES, the task is to predict the set of edges that are bridges in an undirected graph. The task does not use any node or edge features. BRIDGES is a binary edge classification task where the class label indicates whether an edge is a bridge in the graph.

A.5 Runtime and memory

Here we provide additional information on the runtime and memory requirements of our GDT. We sample the runtime of each experiment by running multiple steps and averaging their runtime. For memory consumption, we consider the complete forward steps of the model and estimate the allocated memory using PyTorch functionality. All computations were made using bfloat16 precision during computation. We run the experiments on a single node consisting of one L40 GPU with 40GB VRAM, 12 CPU cores, and 120GB RAM for all runtime and memory computations. We further note that the presented runtimes are the final runtimes obtained from the selected experiments, and significantly more runtime was used to obtain the chosen hyperparameter choices. We note that the automatic compilation is performed automatically by torch.compile, improves the runtime and memory scaling significantly across all tasks.

Table 5 shows the runtime for a single step, averaged across 1 000 training steps obtained for each model. Timings were obtained using torch functionality. Further Table 6 shows the memory requirement for 1 000 steps of each model. We further note the runtime speed improvements during inference experiments from Section 4 while using FlashAttention [Dao et al., 2022].

Hardware optimizations Efficient compilation of neural networks is already available via CUDA implementations in PyTorch and programming languages such as Triton. We use torch.compile throughout all our experiments. In addition, we want to highlight FlashAttention [Dao et al., 2022], available for the standard transformer, and used in the GDT as an example of architecture-specific hardware optimizations that can reduce runtime and memory requirements.

A.6 Comparison with state-of-the-art

While our study focuses exclusively on the GDT, we provide SOTA performance numbers for our real-world tasks to understand whether the GDT performance is competitive with the best models in the literature. Concretely, for PCQ without 3D positions, the best models typically achieve between 0.0809 and 0.0859 MAE [Chen et al., 2023, Müller et al., 2024, Ma et al., 2023, Rampásek et al., 2022]. For COCO and PASCAL, we find models are generally evaluated on a 500K parameter budget and achieve up to 43.98 F1 and 49.12 F1, respectively [Chen et al., 2025]. Note that we do not adhere to this budget when training on COCO as we find it overly restrictive given the considerable size of this dataset. Consequently, we also use the pre-trained 15M model when performing few-shot

Table 6: 16M/90M/160M models memory requirements in MB for 1 000 steps of each model during training.

PE	#Param.	PCQ	COCO	CODE	FLOW	MST	BRIDGES
NoPE	16M	3120.63	5117.57	9749.29	1702.69	3255.37	2456.55
LPE	16M	3221.13	5239.89	9852.71	1763.87	3401.60	2499.54
SPE	16M	3161.32	5157.65	9763.58	1730.19	3290.09	2490.0
RWSE	16M	3131.57	5147.34	9766.30	1713.27	3276.84	2474.63
RRWP	16M	5419.56	-	-	2253.85	5522.04	3723.31
LPE	90M	9844.48	-	-	-	10197.77	-
RWSE	90M	9659.54	-	-	-	9947.07	-
LPE	160M	13513.76	-	-	-	13986.39	-
RWSE	160M	13266.44	-	-	-	13647.66	-

transfer from COCO to PASCAL. Finally, on CODE, the best models score somewhere between 19.37 [Chen et al., 2022] and 22.22 F1 [Geisler et al., 2023].

A.7 Scaling Results

Here we provide additional results for scaling the GDT to 90M and 160M parameters. These results correspond to the results seen in Figure 4.

Table 7: 90M and 160M parameter results for different PEs over 2 random seeds. PCQ MAE is in micro electron volt (meV) for clarity of presentation.

PE	PCQ (90M) MAE ↓	MST (90M) F1 ↑	MST (160M) F1 ↑
LPE	89.7 ± 0.4	92.86 ± 00.17	93.11 ± 01.01
RWSE	88.9 ± 0.7	94.29 ± 00.68	95.80 ± 00.18
RRWP	86.5 ± 0.3	-	-

B Background

Here, we provide background material on various concepts and definitions used in our work.

Basic notations Let $\mathbb{N} := \{1, 2, \dots\}$ and $\mathbb{N}_0 := \mathbb{N} \cup \{0\}$. The set \mathbb{R}^+ denotes the set of non-negative real numbers. For a set X , $A \subset X$ denotes the *strict* subset and $A \subseteq X$ denotes the subset. For $n \in \mathbb{N}$, let $[n] := \{1, \dots, n\} \subset \mathbb{N}$. We use $\{\{\dots\}\}$ to denote multisets, i.e., the generalization of sets allowing for multiple, finitely many instances for each of its elements. For two non-empty sets X and Y , let Y^X denote the set of functions from X to Y . Given a set X and a subset $A \subset X$, we define the indicator function $1_A: X \rightarrow \{0, 1\}$ such that $1_A(x) = 1$ if $x \in A$, and $1_A(x) = 0$ otherwise. Let M be an $n \times m$ matrix, $n > 0$ and $m > 0$, over \mathbb{R} , then $M_{i,\cdot}$, $M_{\cdot,j}$, $i \in [n]$, $j \in [m]$, are the i th row and j th column, respectively, of the matrix M . We denote with $\text{set}(M)$ the set of rows of M . Let N be an $n \times n$ matrix, $n > 0$, then the *trace* $\text{Tr}(N) := \sum_{i \in [n]} N_{ii}$. In what follows, $\mathbf{0}$ denotes an all-zero vector with an appropriate number of components.

Graphs An (*undirected*) graph G is a pair $(V(G), E(G))$ with *finite* sets of *vertices* $V(G)$ and *edges* $E(G) \subseteq \{\{u, v\} \subseteq V(G) \mid u \neq v\}$. *vertices* or *nodes* $V(G)$ and *edges* $E(G) \subseteq \{\{u, v\} \subseteq V(G) \mid u \neq v\}$. The *order* of a graph G is its number $|V(G)|$ of vertices. If not stated otherwise, we set $n := |V(G)|$ and call G an n -order graph. We denote the set of all n -order (undirected) graphs by \mathcal{G}_n and the set of all (undirected) graphs up to n vertices by $\mathcal{G}_{\leq n}$. In a *directed graph*, we define $E(G) \subseteq V(G)^2$, where each edge (u, v) has a direction from u to v . Given a directed graph G and vertices $u, v \in V(G)$, we say that v is a *child* of u if $(u, v) \in E(G)$. A (directed) graph G is called *connected* if, for any $u, v \in V(G)$, there exist $r \in \mathbb{N}$ and $\{u_1, \dots, u_r\} \subseteq V(G)$, such that $(u, u_1), (u_1, u_2), \dots, (u_r, v) \in E(G)$, and analogously for undirected graphs by replacing directed edges with undirected ones. We say that a graph G is *disconnected* if it is not connected.

For a graph G and an edge $e \in E(G)$, we denote by $G \setminus e$ the graph induced by removing edge e from G . For an n -order graph $G \in \mathcal{G}_n$, assuming $V(G) = [n]$, we denote its adjacency matrix by $A(G) \in \{0, 1\}^{n \times n}$, where $A(G)_{vw} = 1$ if and only if $\{v, w\} \in E(G)$. The neighborhood of a vertex $v \in V(G)$ is denoted by $N_G(v) := \{u \in V(G) \mid \{v, u\} \in E(G)\}$, where we usually omit the subscript for ease of notation, and the degree of a vertex v is $|N_G(v)|$. A graph G is a tree if connected, but $G \setminus e$ is disconnected for any $e \in E(G)$. A tree or a disjoint collection of trees is known as a forest.

A rooted tree (G, r) is a tree where a specific vertex r is marked as the root. For a rooted (undirected) tree, we can define an implicit direction on all edges as pointing away from the root; thus, when we refer to the children of a vertex u in a rooted tree, we implicitly consider this directed structure. For $S \subseteq V(G)$, the graph $G[S] := (S, E_S)$ is the subgraph induced by S , where $E_S := \{(u, v) \in E(G) \mid u, v \in S\}$. A (vertex-)labeled graph is a pair (G, ℓ_G) with a graph $G = (V(G), E(G))$ and a (vertex-)label function $\ell_G: V(G) \rightarrow \Sigma$, where Σ is an arbitrary countable label set. For a vertex $v \in V(G)$, $\ell_G(v)$ denotes its label. A Boolean (vertex-)d-labeled graph is a pair (G, ℓ_G) with a graph $G = (V(G), E(G))$ and a label function $\ell_G: V(G) \rightarrow \{0, 1\}^d$. We denote the set of all n -order Boolean d -labeled graphs as $\mathcal{G}_{n,d}^{\mathbb{B}}$. An attributed graph is a pair (G, a_G) with a graph $G = (V(G), E(G))$ and an (vertex-)attribute function $a_G: V(G) \rightarrow \mathbb{R}^{1 \times d}$, for $d > 0$. That is, contrary to labeled graphs, vertex annotations may be from an uncountable set. The attribute or feature of $v \in V(G)$ is $a_G(v)$. We denote the class of all n -order graphs with d -dimensional, real-valued vertex features by $\mathcal{G}_{n,d}^{\mathbb{R}}$.

Two graphs G and H are isomorphic if there exists a bijection $\varphi: V(G) \rightarrow V(H)$ that preserves adjacency, i.e., $(u, v) \in E(G)$ if and only if $(\varphi(u), \varphi(v)) \in E(H)$. In the case of labeled graphs, we additionally require that $\ell_G(v) = \ell_H(\varphi(v))$ for $v \in V(G)$. Moreover, we call the equivalence classes induced by \simeq isomorphism types and denote the isomorphism type of G by $\tau(G)$. A graph class is a set of graphs closed under isomorphism. Given two graphs G and H with disjoint vertex sets, we denote their disjoint union by $G \cup H$.

B.1 Transformers

Here, we will introduce attention with an additive attention bias and the transformer architecture.

Definition 5 (Attention (with bias)). Let $Q, K, V \in \mathbb{R}^{n \times d}$ and $B \in \mathbb{R}^{n \times n}$, with $n, d \in \mathbb{N}^+$. We define biased attention as

$$\text{Attention}(Q, K, V, B) := \text{softmax}(d^{-\frac{1}{2}} \cdot QK^T + B)V,$$

where softmax is applied row-wise and defined, for a vector $x \in \mathbb{R}^{1 \times n}$, as

$$\text{softmax}(x) := \left[\frac{\exp(x_1)}{\sum_{i \in [n]} \exp(x_i)} \quad \cdots \quad \frac{\exp(x_n)}{\sum_{i \in [n]} \exp(x_i)} \right].$$

Definition 6 (Multi-head attention (with bias)). Let $X \in \mathbb{R}^{n \times d}$, $B \in \mathbb{R}^{n \times n \times h}$, and let $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$, $W_O \in \mathbb{R}^{d \times d}$ be learnable parameters, with $n, d \in \mathbb{N}^+$. Let $h \in \mathbb{N}^+$ be the number of heads, such that a $d_h \in \mathbb{N}^+$ for which $d = h \cdot d_h$. We call d_h the head dimension and define h -head attention over X as

$$\text{MHA}(X, B) := [\tilde{X}_1 \quad \cdots \quad \tilde{X}_h] W_O,$$

where, for all $i \in [h]$,

$$\tilde{X}_i := \text{Attention}(XW_Q^{(i)}, XW_K^{(i)}, XW_V^{(i)}, B_i),$$

with $B_i \in \mathbb{Q}^{n \times n}$ denoting the attention bias for the i -head, indexed along the third dimension of B , $W_Q^{(i)}, W_K^{(i)}, W_V^{(i)} \in \mathbb{R}^{d \times d_h}$, and

$$\begin{aligned} W_Q &:= [W_Q^{(1)}, \dots, W_Q^{(h)}] \\ W_K &:= [W_K^{(1)}, \dots, W_K^{(h)}] \\ W_V &:= [W_V^{(1)}, \dots, W_V^{(h)}]. \end{aligned}$$

Table 8: Overview of learnable parameters in the GDT, excluding embedding parameters. Here, $d \in \mathbb{N}^+$ is the embedding dimension, $d_f \in \mathbb{N}^+$ is the hidden dimension, and $T \in \mathbb{N}^+$ is the number of layers. The suffix $\times T$ indicates that the parameters occur in each of the T layers.

Params.	Dims.	Module	Description
$\ell_V([\text{cls}])$	$1 \times d$	Token embeddings	Learnable embedding for the $[\text{cls}]$ token
$\ell_E([\text{cls}], \cdot)$	$1 \times d$		Learnable embedding for the out-going edges from the $[\text{cls}]$ token
$\ell_E(\cdot, [\text{cls}])$	$1 \times d$		Learnable embedding for the in-coming edges to the $[\text{cls}]$ token
\mathbf{W}_P	$d \times d$		Weight matrix for the node-level PEs
ρ	$(d \times d_f, d_f \times h)$	Attention bias	MLP applied to the edge embeddings
\mathbf{W}_U	$1 \times h$		Weight matrix for the relative PEs in the attention bias
$\mathbf{W}_Q \times T$	$d \times d$	MHA <small>Definition 6</small>	Query weight matrix in multi-head attention
$\mathbf{W}_K \times T$	$d \times d$		Key weight matrix in multi-head attention
$\mathbf{W}_V \times T$	$d \times d$		Value weight matrix in multi-head attention
$\mathbf{W}_1 \times T$	$d \times d_f$	MLP <small>Definition 7</small>	Input projection of the MLP
$\mathbf{W}_2 \times T$	$d_f \times d$		Output projection of the MLP

Definition 7 (Two-layer MLP). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ with $n, d, d_f \in \mathbb{N}^+$, where d_f is the hidden dimension. We define a two-layer MLP as

$$\text{MLP}(\mathbf{x}) := \sigma(\mathbf{x}\mathbf{W}_1)\mathbf{W}_2,$$

where MLP is applied independently to each row $\mathbf{x} \in \mathbb{R}^{1 \times d}$ in \mathbf{X} . Here, $\mathbf{W}_1 \in \mathbb{R}^{d \times d_f}$ is the in-projection matrix, $\mathbf{W}_2 \in \mathbb{R}^{d_f \times d}$ is the out-projection matrix, and $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is an element-wise activation function such as GELU [Hendrycks and Gimpel, 2016].

Definition 8 (Transformer architecture). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ be a token matrix and $\mathbf{B} \in \mathbb{R}^{n \times n}$ be an attention bias, with $n, d \in \mathbb{N}^+$. The $t + 1$ -th transformer layer updates token representations $\mathbf{X}^t \in \mathbb{R}^{n \times d}$ as

$$\begin{aligned} \mathbf{X}' &\leftarrow \mathbf{X}^t + \text{MHA}(\text{LayerNorm}(\mathbf{X}^t), \mathbf{B}), \\ \mathbf{X}^{t+1} &\leftarrow \mathbf{X}' + \text{MLP}(\text{LayerNorm}(\mathbf{X}')), \end{aligned}$$

where MLP is defined in Definition 7.

B.2 Extended notation for the theoretical analysis of the GDT

Here, we introduce some notation for the GDT that we will use in our theoretical analysis.

Learnable parameters We give an overview of all learnable parameters of the GDT in Table 8. In practice, node and edge features are typically present as integers or continuous feature vectors, and we embed them using learnable MLPs. We refer to such parameters as *embedding parameters*. Note that in Table 8, we exclude embedding parameters, as for simplicity, we assume in our framework that node and edge features are already embedded.

Let $d, d_f, T, h \in \mathbb{N}^+$ denote the number of embedding dimensions, the number of hidden dimensions, the number of layers, and the number of attention heads, respectively. Then, the number of learnable parameters (excluding embedding parameters) is given by

$$\begin{aligned} \# \text{ params} &= 3d + d^2 + dd_f + d_f h + h + 3Td^2 + 2Tdd_f \\ &= (3T + 1)d^2 + (2T + 1)dd_f + 3d + (d_f + 1)h. \end{aligned}$$

We denote the complete set of learnable parameters in Table 8 with $\Theta(d, d_f, T, h)$.

Graph transformer representations Here, we introduce some short-hand notation for graph transformer representations. Given token matrix \mathbf{X} and attention bias \mathbf{B} , we write $\hat{\mathbf{X}}^t(v)$ to denote the representation of node $v \in V(G)$ after t transformer layers with \mathbf{X} and \mathbf{B} as input. Note that, since we fix an arbitrary order of the nodes, if v is the i -th node in this order, $\hat{\mathbf{X}}^t(v) = \hat{\mathbf{X}}_i^t$.

B.3 Extended notation for the theoretical analysis of PEs

Here we introduce notations of Zhang et al. [2024] used in their paper on PE expressiveness. We adapt this notation to fit LPE [Müller and Morris, 2024], SAN [Kreuzer et al., 2021], and SignNet

[Lim et al., 2023] and use it in our proofs in Appendix D. We introduce the color refinement algorithm notation and propose a color refinement algorithm for each PE. The respective algorithms for BasisNet [Lim et al., 2023] and SPE [Huang et al., 2024] are given by Zhang et al. [2024].

Definition 9. [Zhang et al., 2024] We call any graph invariant a k -dim color mapping. The family of k -dim color mappings is denoted by M_k . Each color mapping defines an equivalence relation \sim_χ between rooted graphs G_u, H_v marking k vertices and $G^u \sim_\chi H^v$ iff $\chi_G(u) = \chi_H(v)$. Further, we denote the family of k -dim spectral color mappings by M_k^Λ . Similar to M_k the family of spectral color mappings is obtained from the color mappings acting on $\{(G_u, \lambda) : G_u \in G, \lambda \in \Lambda^M(G)\}$ where $\Lambda^M(G)$ denotes the eigenvalues of a matrix M .

Definition 10. [Zhang et al., 2024] A function T mapping from M_{k_1} to M_{k_2} is called a color transform. We assume all color transforms are order-preserving in terms of color mappings. Given $T(\chi) \preceq \chi$, a color transform is also called color refinement, and T^t denotes the t times composition of T . In addition T^∞ is the stable color refinement obtained from $T^{t'}$ with t' the smallest integer where further iterations do not induce a different partition of the underlying nodes in a graph resulting in $T \circ T^\infty \equiv T^\infty$. Following Zhang et al. [2024] T^∞ is well defined.

A coloring algorithm is then formed by concatenating a stable color transform $T^\infty : M_k \rightarrow M_k$ and a pooling function $U : M_k \rightarrow M_0$.

Definition 11. We say that color mappings χ_1, χ_2 are equivalent given that $G^u \sim_{\chi_1} H^v$ iff $G^u \sim_{\chi_2} H^v$. Furthermore, we say that a color mapping χ_1 is finer/more expressive than χ_2 if $G^u \sim_{\chi_1} H^v \Rightarrow G^u \sim_{\chi_2} H^v$, noted by \preceq .

Lemma 12. [Zhang et al., 2024] Let $T_1, T_2 : M_{k_1} \rightarrow M_{k_2}$ and $U_1, U_2 : M_{k_2} \rightarrow M_{k_3}$ be color refinements. If $T_1 \preceq T_2$ and $U_1 \preceq U_2$ then $U_1 \circ T_1 \preceq U_2 \circ T_2$.

Lemma 13. [Zhang et al., 2024] Let $T_1 : M_{k_1} \rightarrow M_{k_1}$ and $T_2 : M_{k_2} \rightarrow M_{k_2}$ be color refinements and $T^\infty : M_k \rightarrow M_k$ be the stable refinement of M_k . Further let $U_1 : M_{k_0} \rightarrow M_{k_1}$ and $U_2 : M_{k_1} \rightarrow M_{k_2}$ be color refinements. Then it follows. If $T_2 \circ U_2 \circ T_1^\infty \circ U_1 \equiv U_2 \circ T_1^\infty \circ U_1$ then $U_2 \circ T_1^\infty \circ U_1 \preceq T_2^\infty \circ U_2 \circ U_1$.

The two lemmas above provide a straightforward approach to determining whether architecture A_1 is more expressive than architecture A_2 [Zhang et al., 2024]. To prove that A_1 is more expressive than A_2 , we show that $T_2 \circ T_1^\infty \equiv T_1^\infty$ holds, with T_i being the color refinement of A_i respectively.

Definition 14. [Zhang et al., 2024] We define the following color refinements corresponding to the induced refinements of each algorithm. We define global pooling as providing an injective coloring using a hash function of a multiset. In this case, we consider the multiset over nodes in a graph. Other pooling operations are defined below.

Global pooling: Define $T_{GP} : M_1 \rightarrow M_0$ and $\chi \in M_1$ such that for a graph G and a color mapping $\chi \in M_1$ it holds:

$$[T_{GP}(\chi)](G) = \text{hash}(\{\chi_G(u) : u \in V(G)\}).$$

The 1-WL refinement gives us the 1-WL coloring update generalized to all nodes in the graph instead of neighboring graphs.

1-WL refinement: Given $T_{WL} : M_1 \rightarrow M_1$ such that for any choice of $\chi \in M_1$:

$$[T_{WL}(\chi)]_G(u) = \text{hash}(\chi_G(u), \{(\chi_G(v), \text{atp}_G(u, v)) : v \in V(G)\}).$$

We then define one—and two-dimensional spectral pooling, which allows for pooling over distinct eigenvalues similar to the global pooling refinement.

Spectral Pooling: Define $T_{SP2} : M_2^\Lambda \rightarrow M_1$ and $T_{SP1} : M_1^\Lambda \rightarrow M_1$ such that for $\chi \in M_2^\Lambda$ and $\chi' \in M_1^\Lambda$:

$$[T_{SP1}(\chi')]_G(u) = \text{hash}(\{\chi'_G(\lambda, u) : \lambda \in \Lambda^M(G)\}).$$

$$[T_{SP2}(\chi)]_G(u, v) = \text{hash}(\{\chi_G(\lambda, u, v) : \lambda \in \Lambda^M(G)\}).$$

A pooling variant without the spectrum is denoted by T_{P2} and T_{P1} .

To allow for an examination of BasisNet and SPE, we consider the 2-IGN refinement. This refinement is obtained from evaluating the expressiveness of a 2-IGN and its basis functions as defined by Maron et al. [2019a].

820 **2-IGN refinement:** With $T_{\text{IGN}}: M_2 \rightarrow M_2$, $\chi \in M_2$ as any color mapping and $\delta_{uv}(c) = c$ given
 821 $u = v$, otherwise 0:

$$\begin{aligned} [T_{\text{IGN}}(\chi)]_G(u, v) = & \text{hash}(\chi_G(u, v), \chi_G(u, u), \chi_G(v, v), \chi_G(v, u), \delta_{uv}(\chi_G(u, u)), \\ & \{\!\!\{ \chi_G(u, w) : w \in V(G) \}\!\!\}, \{\!\!\{ \chi_G(w, u) : w \in V(G) \}\!\!\}, \\ & \{\!\!\{ \chi_G(v, w) : w \in V(G) \}\!\!\}, \{\!\!\{ \chi_G(w, v) : w \in V(G) \}\!\!\}, \\ & \{\!\!\{ \chi_G(w, w) : w \in V(G) \}\!\!\}, \{\!\!\{ \chi_G(w, x) : w, x \in V(G) \}\!\!\}, \\ & \delta_{uv}(\{\!\!\{ \chi_G(u, w) : w \in V(G) \}\!\!\}), \delta_{uv}(\{\!\!\{ \chi_G(w, u) : w \in V(G) \}\!\!\}), \\ & \delta_{uv}(\{\!\!\{ \chi_G(w, w) : w \in V(G) \}\!\!\}), \delta_{uv}(\{\!\!\{ \chi_G(w, x) : w, x \in V(G) \}\!\!\})). \end{aligned}$$

822 Further, we use the BasisNet pooling refinement and Siamese IGN refinement to describe the process
 823 of the BasisNet computation.

824 **BasisNet Pooling:** Given $T_{\text{BP}}: M_2^\Lambda \rightarrow M_1^\Lambda$ and $\chi \in M_2^\Lambda$:

$$\begin{aligned} [T_{\text{BP}}(\chi)]_G(\lambda, u) = & \text{hash}(\chi_G(\lambda, u, v), \{\!\!\{ \chi_G(\lambda, u, v) : v \in V(G) \}\!\!\}, \\ & \{\!\!\{ \chi_G(\lambda, v, u) : v \in V(G) \}\!\!\}, \{\!\!\{ \chi_G(\lambda, v, v) : v \in V(G) \}\!\!\}, \{\!\!\{ \chi_G(\lambda, v, w) : v, w \in V(G) \}\!\!\}). \end{aligned}$$

825 **Siamese IGN refinement:** Given $T_{\text{SIAM}}: M_2^\Lambda \rightarrow M_2^\Lambda$ and $\chi \in M_2^\Lambda$:

$$[T_{\text{SIAM}}(\chi)]_G(\lambda, u, v) = [T_{\text{IGN}}(\chi(\lambda, \cdot, \cdot))]_G(u, v).$$

826 We further provide additional color refinement algorithms based on the encodings introduced through-
 827 out this work: Given the initial color refinement of SAN as $\chi_{\text{SAN}}(\lambda, u) = (\lambda_1, \dots, \lambda_m, v_{1:m}^u)$ where
 828 λ_i denote the eigenvalues and v^u the eigenvector of the graph Laplacian associated to the node u .
 829 Then we define the SAN color refinement alongside the existing refinements as follows:

$$[T_{\text{SAN}}(\chi)]_G = T_{\text{GP}} \circ T_{\text{ENC}} \circ T_L(\chi_{\text{SAN}}).$$

830 with T_{ENC} denoting the transformer encoder layer. The complete BasisNet refinement is given by the
 831 concatenation of refinements given in Definition Definition 14:

$$[T_{\text{BasisNet}}]_G = T_{\text{GP}} \circ T_{\text{WL}} \circ T_{\text{SP1}} \circ T_{\text{BP}} \circ T_{\text{SIAM}}(\chi_{\text{Basis}}).$$

832 Following the definition of BasisNet as a color refinement in Definition 14 and assuming a message
 833 passing GNN for ρ , the color refinement of SignNet using the initial SignNet color refinement
 834 $\chi_{\text{Sign}}(\lambda, u, w) = (\lambda, \mathbf{V}^u, \mathbf{V}^w)$ and $T_\phi: M_2^\Lambda \rightarrow M_2^\Lambda$ is given by:

$$[T_{\text{Sign}}(\chi)]_G = T_{\text{GP}} \circ T_{\text{WL}}^\infty \circ T_{\text{SP2}} \circ T_\phi(\chi_{\text{Sign}}),$$

835 where \mathbf{V}^u denotes the eigenvector associated to the eigenvalue λ and the node u and T_ϕ be the
 836 refinement depending on the choice of ϕ . However, we note $T_{\text{IGN}} \preceq T_\phi$, by definition of SignNet and
 837 BasisNet.

838 The refinement for the LPE encoding is given similarly to the BasisNet and SPE refinement by
 839 replacing ρ with a color refinement which is 1-WL expressive. Furthermore, we assume ϕ to be a
 840 spectral color refinement with expressiveness up to a 2-IGN. Common choices for ϕ include MLPs or
 841 GNNs known to be less expressive than a 2-IGN. With initial colorings $\chi_{\text{LPE}}(\lambda, u, w) = (\lambda, \mathbf{V}^u, \mathbf{V}^w)$
 842 and $T_\phi^{\text{LPE}}: M_2^\Lambda \rightarrow M_2^\Lambda$ the refinement follows:

$$[T_{\text{LPE}}(\chi)]_G = T_{\text{GP}} \circ T_{\text{WL}}^\infty \circ T_{\text{SP2}} \circ T_\phi^{\text{LPE}}(\chi_{\text{LPE}}),$$

843 where \mathbf{V}^u denotes the eigenvector associated to node u .

844 B.4 Positional encodings

845 In the following, we define positional encodings.

846 RWSE

847 **Definition 15.** Let $\mathbf{R} := \mathbf{D}^{-1} \mathbf{A}$ be the random walk matrix, with \mathbf{D} denoting the degree matrix and
 848 \mathbf{A} the adjacency matrix of a graph G . The random walk structural encodings (RWSE) are given by:

$$\mathbf{P}_{i,i} = [\mathbf{I}, \mathbf{R}, \mathbf{R}^2, \dots, \mathbf{R}^{k-1}]_{i,i}.$$

849 **Definition 16.** Let $\mathbf{P}_{i,i}$ be the RWSE encoding vector and $\mathbf{F}: \mathbb{R}^k \rightarrow \mathbb{R}^d$ a MLP with two layers, and
 850 d denoting the encoding dimension. Then the RWSE encoding is computed by $\mathbf{F}(\mathbf{P}_{i,i})$ and denoted
 851 by $\mathbf{P}_k^{\text{RW}}(G)$ for a graph G with random walk length k .

852 RRWP

853 **Definition 17.** Let $\mathbf{R} := \mathbf{D}^{-1}\mathbf{A}$ with \mathbf{D} as the diagonal degree matrix and \mathbf{A} as the adjacency
854 matrix, be the random walk operator, and k the maximum length of the random walk. Then the
855 relative random walk probabilities (RRWP) are defined as:

$$\mathbf{P}_{i,j} = [\mathbf{I}, \mathbf{R}, \mathbf{R}^2, \dots, \mathbf{R}^{k-1}]_{i,j}.$$

856 The initial node encoding p_0 is then defined as \mathbf{R}_{ii} for each node i in the graph.

857 **Definition 18** (RRWP Encoding Computation). Let \mathbf{P} be the RRWP encoding tensor and MLP :
858 $\mathbb{R}^k \rightarrow \mathbb{R}^d$, where d denotes the encoding dimension, be a multi-layer neural network. Then the
859 encoding MLP($\mathbf{P}_{i,j,:}$) is computed element-wise by the multi-layer neural network. RRWP is then
860 denoted by $\mathbf{P}_k^{\text{RR}}(G)$ for a graph G with random walk length k

861 **Spectral Attention Networks (SAN)** Kreuzer et al. [2021] propose incorporating eigenvalues
862 and eigenvectors in a positional encoding neural network. SAN encoding can be computed using
863 row-wise applied neural networks by selecting the k lowest eigenvalues and associated eigenvectors.

864 **Definition 19.** Let $\phi: \mathbb{R} \rightarrow \mathbb{R}$ be a linear layer and $\rho: \mathbb{R} \rightarrow \mathbb{R}$ be a transformer encoder layer with
865 sum aggregation. Further \mathbf{V}_i denotes the i -th column of the eigenvector matrix \mathbf{V} . Then the SAN
866 encoding is defined as follows:

$$\text{SAN}(\mathbf{V}, \lambda) = \rho([\phi(\mathbf{V}_1, \lambda) \dots \phi(\mathbf{V}_k, \lambda)]).$$

867 A generalization of the SAN encoding is given by the LPE encoding of Müller and Morris [2024] in
868 Definition 23.

869 **SignNet** Since the computation of eigenvectors using the eigenvector decomposition is not sign
870 invariant, and both \mathbf{V}_i and $-\mathbf{V}_i$ are valid eigenvectors of the graph Laplacian Lim et al. [2023]
871 propose the construction of a sign invariant encoding using eigenvector information. Considering
872 the k smallest eigenvalues and associated eigenvectors from the eigenvalue decomposition, SignNet
873 computes the corresponding encoding using a neural network architecture.

874 **Definition 20.** Let $\phi_1, \dots, \phi_k: \mathbb{R} \rightarrow \mathbb{R}$ and $\rho: \mathbb{R} \rightarrow \mathbb{R}$ be permutation equivariant neural networks
875 from vectors to vectors. Then the SignNet encodings are computed using:

$$\text{SignNet}(\mathbf{V}) = \rho([\phi_1(\mathbf{V}_1) + \phi_1(-\mathbf{V}_1) \dots \phi_k(\mathbf{V}_k) + \phi_k(-\mathbf{V}_k)]).$$

876 Commonly, ϕ_1, \dots, ϕ_k are selected as element-wise MLPs or DeepSets [Lim et al., 2023] and ρ as a
877 GIN with sum aggregation and the adjacency matrix of the original graph.

878 **BasisNet** Proposed as an extension of SignNet by Lim et al. [2023], BasisNet encodings provide
879 an encoding invariant to the basis of eigenspaces obtained from the graph Laplacian. Since the
880 orthogonal group $O(1)$ denotes sign invariance, BasisNet also incorporates sign invariance.

881 **Definition 21.** Let \mathbf{V}_i denote the orthonormal basis of an d_i -dimensional eigenspace of the graph
882 Laplacian. Further, l denotes the number of eigenspaces. Given unrestricted neural networks
883 $\phi_{d_1}, \dots, \phi_{d_l}: \mathbb{R} \rightarrow \mathbb{R}$, shared across the subspaces with the same dimension d_i , and a permutation
884 equivariant neural network $\rho: \mathbb{R} \rightarrow \mathbb{R}$ BasisNet encodings are computed the following:

$$\text{BasisNet}(\mathbf{V}) = \rho([\phi_{d_1}(\mathbf{V}_1 \mathbf{V}_1^T) \dots \phi_{d_l}(\mathbf{V}_l \mathbf{V}_l^T)]).$$

885 Implementation wise Lim et al. [2023] propose 2-IGNs [Maron et al., 2019a] for ϕ_{d_i} and a FFN with
886 sum aggregation for ρ . They note that all neural networks could be replaced with k -IGNs, however
887 they deemed it infeasible for efficient computation. This reduces the computation to the following
888 with $\rho: \mathbb{R} \rightarrow \mathbb{R}$ and $\text{IGN}_{d_i}: \mathbb{R}^{n^2} \rightarrow \mathbb{R}^n$ denoting an IGN from matrices to vectors:

$$\text{BasisNet}(\mathbf{V}) = \text{FFN}([\text{IGN}_{d_1}(\mathbf{V}_1 \mathbf{V}_1^T) \dots \text{IGN}_{d_l}(\mathbf{V}_l \mathbf{V}_l^T)]).$$

889 **SPE** Following notation from Huang et al. [2024], the SPE encoding is computed using the k
890 smallest eigenvalues and associated eigenvectors obtained from an eigenvalue decomposition. With
891 sufficient conditions for neural networks ϕ_1, \dots, ϕ_k and ρ , SPE is stable with respect to the graph
892 Laplacian.

Definition 22. Given $\phi_1 \dots \phi_k : \mathbb{R} \rightarrow \mathbb{R}$ as Lipschitz continuous, equivariant FFNs and $\rho : \mathbb{R} \rightarrow \mathbb{R}$ a Lipschitz continuous, permutation equivariant neural network. Then the SPE encoding is computed by:

$$\text{SPE}(\mathbf{V}, \lambda) = \rho([\mathbf{V} \text{diag}(\phi_1(\lambda)) \mathbf{V}^T \dots \mathbf{V} \text{diag}(\phi_k(\lambda)) \mathbf{V}^T]),$$

with ϕ_1, \dots, ϕ_k and ρ applied row wise. Further, we denote the SPE embedding on a graph G by $\mathbf{P}_k^{\text{SPE}}$.

Commonly, ϕ_i is considered an element-wise MLP, and ρ is a GIN using the adjacency matrix of the original graph. Huang et al. [2024] propose to split tensor

$$\mathbf{Q} = [\mathbf{V} \text{diag}(\phi_1(\lambda)) \mathbf{V}^T \dots \mathbf{V} \text{diag}(\phi_k(\lambda)) \mathbf{V}^T] \in \mathbb{R}^{n \times n \times l}$$

into n matrices of shape $n \times l$ which are then passed into the GIN ρ and aggregated using sum aggregation into a single $n \times d$ matrix.

LPE Initially introduced by Kreuzer et al. [2021] and generalized by Müller and Morris [2024], the LPE encodings are computed similarly to the previously introduced SPE encodings. Instead of using the eigenvector matrix $\mathbf{V} \in \mathbb{R}^{n \times l}$, each i -th column consisting of one eigenvector denoted by $\mathbf{V}_i \in \mathbb{R}^l$ is used.

Definition 23. Let $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^k$ be a row-wise applied FFN and $\rho : \mathbb{R} \rightarrow \mathbb{R}$ a permutation equivariant network. Furthermore, $\epsilon \in \mathbb{R}$ denotes a learnable parameter. Then the LPE are given by:

$$\text{LPE}(\mathbf{V}, \lambda) = \rho([\phi(\mathbf{V}_1^T, \lambda + \epsilon), \dots, \phi(\mathbf{V}_k^T, \lambda + \epsilon)]).$$

Setting $\epsilon = 0$ reduces the LPE to the encoding provided by Kreuzer et al. [2021]. As proposed by Müller and Morris [2024] ρ sums the input tensor with regard to its first dimension and applies an FFN. In contrast to the SAN encoding, no transformer encoder is used to compute the encoding. Similar to previous embeddings, we denote LPE embeddings for a graph G by $\mathbf{P}_k^{\text{LPE}}$ with k as the number of eigenvalues used.

C Proving that the GDT can simulate the GD-WL

Here, we prove Theorem 1. Concretely, we formally state and prove both statements in Theorem 1 in Appendix C.1 and Appendix C.2, respectively.

C.1 Lower-bound on the expressivity of the GDT

Here, we prove that we can compute the GD-WL [Zhang et al., 2023] with the GDT, our GT defined in Section 2. Concretely, given a graph $G := (V(G), E(G))$, recall from Section 2.3 the GD-WL as updating the color $\chi_G^t(v)$ of node $v \in V(G)$, as

$$\chi_G^{t+1}(v) := \text{hash}(\{(d_G(v, w), \chi_G^t(w)) : w \in V(G)\}),$$

where d_G is a distance between nodes in G and hash is an injective map. In the transformer, we will represent node colors as one-hot vectors of some arbitrary but fixed dimension d . Furthermore, we will incorporate pairwise distances through the attention bias. We then show that a single GT layer can compute the color update in Equation (1). We show this result by leveraging specific characteristics of softmax-attention. For notational convenience, we will denote with $\mathbf{X}^t \in \{0, 1\}^{L \times d}$ the one-hot color matrix of the GD-WL after t iterations.

We begin by stating our main result. Afterwards, we develop our proof techniques and prove the result.

Stating the main result We will formally state our theorem, showing that our GT can simulate the GD-WL. Afterwards, we give an overview of the proof, including the key challenges and ideas.

Theorem 24. Let $G := (V(G), E(G))$ be a graph with $n \in \mathbb{N}^+$ nodes and node distance function $d_G : V(G)^2 \rightarrow \mathbb{Q}$. Let $d, d_f, T, h \in \mathbb{N}^+$ denote the number of embedding dimensions, the number of hidden dimensions, the number of layers, and the number of attention heads, respectively. Let $L := n + 1$ and let $\hat{\mathbf{X}}^0 \in \mathbb{R}^{L \times d}$ and $\mathbf{B} \in \mathbb{R}^{L \times L \times h}$ be initial token embeddings and attention bias constructed according to Section 2 using node distance d_G . Then, there exists weights for the parameters in $\Theta(d, d_f, T, h)$ such that $\hat{\mathbf{X}}^t = \mathbf{X}^t$, for all $t \geq 0$, an arbitrary but fixed hash and using d_G as the distance function.

Proof overview The central problem we face when proving the above theorem is how to injectively encode the multisets in Equation (1) with softmax-attention. This is because softmax-attention computes a weighted mean, whereas existing results for encoding multisets use sums [Xu et al., 2019, Morris et al., 2019, Zhang et al., 2023]. Because these multisets are at the core of our proof, we introduce them here formally.

Definition 25 (Distance-paired multisets). Given a graph $G := (V(G), E(G))$ with $n \in \mathbb{N}^+$ nodes and let $L := n + 1$, for each token $v \in V(G) \cup \{\text{[cls]}\}$, we construct a vector $\mathbf{v} \in \mathbb{Q}^{1 \times L}$ from the distances of v to tokens in $V(G) \cup \{\text{[cls]}\}$, such that

$$\mathbf{v}_i := d_G(v, w_i),$$

where $w_i \in V(G)$, for $i \in [n]$, is the i -th token in an arbitrary but fixed ordering of nodes in $V(G)$ and w_L is the [cls] token. We fix the distance of [cls] to all tokens as $\max_{v, w \in V(G)} d_G(v, w) + 1$. We represent node colors as one-hot vectors and stack them into a matrix $\mathbf{X} \in \{0, 1\}^{L \times d}$ with $d \in \mathbb{N}^+$ and where \mathbf{X}_L , representing the color of the [cls] token, receives a special color, not used by any node. We then write the distance-paired multiset in Equation (1) as

$$[\mathbf{v}]_{\mathbf{X}} := \{(\mathbf{v}_i, \mathbf{X}_i)\}_{i \in [L]}.$$

We can then restate the update of token $v \in V(G) \cup \{\text{[cls]}\}$ by the GD-WL as

$$\chi_G^{t+1}(v) := \text{hash}([\mathbf{v}]_{\mathbf{X}^t}), \quad (3)$$

For notational convenience, for every $\mathbf{x} \in \text{set}(\mathbf{X})$, we define $A(\mathbf{x}) := \{i \in [L] \mid \mathbf{X}_i = \mathbf{x}\}$ as the set of token indices with token representation \mathbf{x} . Further, we write

$$[\mathbf{v}]_{\mathbf{x}} := \{\mathbf{v}_i \mid i \in A(\mathbf{x})\}$$

and

$$[\mathbf{v}]_{\mathbf{x}, \mathbf{w}} := \{\mathbf{v} + \mathbf{w}_j \mid \mathbf{v} \in [\mathbf{v}]_{\mathbf{x}}, j \in [L]\},$$

again for notational convenience, where $\mathbf{w} \in \mathbb{Q}^{1 \times L}$ is the distance vector corresponding to another node $w \in V(G) \cup \{\text{[cls]}\}$.

Recall that we introduced the distance function d_G into the attention via the attention bias \mathbf{B} . Now, to injectively encode distance-paired multisets, we want to prove that there exists weights $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V$ such that for two tokens $v, w \in V(G) \cup \{\text{[cls]}\}$ with corresponding distance vectors \mathbf{v}, \mathbf{w} ,

$$\text{softmax}(\mathbf{X}(v)\mathbf{W}_Q(\mathbf{X}\mathbf{W}_K)^T + \mathbf{v})\mathbf{X}\mathbf{W}_V = \text{softmax}(\mathbf{X}(w)\mathbf{W}_Q(\mathbf{X}\mathbf{W}_K)^T + \mathbf{w})\mathbf{X}\mathbf{W}_V,$$

if and only if $[\mathbf{v}]_{\mathbf{X}} = [\mathbf{w}]_{\mathbf{X}}$. Note that for simplicity, we omit the scaling factor in the attention and that we wrote \mathbf{v} and \mathbf{w} to indicate the corresponding row of \mathbf{B} for tokens v and w , respectively. We will now simplify, by setting $\mathbf{W}_Q = \mathbf{W}_K = \mathbf{0}$ and $\mathbf{W}_V = \mathbf{I}$ and arrive at the condition

$$\text{softmax}(\mathbf{v})\mathbf{X} = \text{softmax}(\mathbf{w})\mathbf{X} \iff [\mathbf{v}]_{\mathbf{X}} = [\mathbf{w}]_{\mathbf{X}}.$$

Here, we prove the above holds under mild conditions in the following lemma. Note that we split up the forward and backward directions of the lemma, as we will use different proof strategies for each direction.

Lemma 26. Let $\mathbf{v}, \mathbf{w} \in \mathbb{Q}^{1 \times L}$ with $\max_i \mathbf{v}_i = \max_i \mathbf{w}_i$ and let $\mathbf{X} \in \{0, 1\}^{L \times d}$ be a matrix whose rows are one-hot vectors, for some $L, d \in \mathbb{N}^+$. Further, we require \mathbf{X} to have at least two distinct rows. Then,

$$\text{softmax}(\mathbf{v})\mathbf{X} = \text{softmax}(\mathbf{w})\mathbf{X} \implies [\mathbf{v}]_{\mathbf{X}} = [\mathbf{w}]_{\mathbf{X}} \quad (4)$$

and

$$\text{softmax}(\mathbf{v})\mathbf{X} = \text{softmax}(\mathbf{w})\mathbf{X} \longleftarrow [\mathbf{v}]_{\mathbf{X}} = [\mathbf{w}]_{\mathbf{X}}. \quad (5)$$

As mentioned above, we will treat the forward and backward directions differently. The backward direction is fairly straightforward, seeing that $\text{softmax}(\mathbf{v})\mathbf{X}$ is a function over $[\mathbf{v}]_{\mathbf{X}}$. For the forward direction, the idea is first to notice that the condition $[\mathbf{v}]_{\mathbf{X}} = [\mathbf{w}]_{\mathbf{X}}$, on the right side of Equation (4), is equivalent to comparing the multiset of distances paired with each distinct one-hot vector in \mathbf{X} independently, as distinct one-hot vectors do not have common non-zero channels; see the following lemma for a precise statement of this property and see Appendix E for the proof.

976 **Lemma 27.** Let $\mathbf{v}, \mathbf{w} \in \mathbb{Q}^{1 \times L}$ and let $\mathbf{X} \in \{0, 1\}^{L \times d}$ be a matrix whose rows are one-hot vectors,
 977 for some $L, d \in \mathbb{N}^+$. Then, $[\mathbf{v}]_{\mathbf{X}} = [\mathbf{w}]_{\mathbf{X}}$, if and only if, for every $\mathbf{x} \in \text{set}(\mathbf{X})$, $[\mathbf{v}]_{\mathbf{x}} = [\mathbf{w}]_{\mathbf{x}}$.

978 To understand the implication of this result in the context of proving Lemma 26, let us first rearrange
 979 the left side of Equation (4) as follows.

980 **Lemma 28.** Let $\mathbf{v}, \mathbf{w} \in \mathbb{Q}^{1 \times L}$ and let $\mathbf{X} \in \{0, 1\}^{L \times d}$ be a matrix whose rows are one-hot vectors,
 981 for some $L, d \in \mathbb{N}^+$. Then, $\text{softmax}(\mathbf{v})\mathbf{X} = \text{softmax}(\mathbf{w})\mathbf{X}$, if and only if, for every $\mathbf{x} \in \text{set}(\mathbf{X})$,

$$\sum_{i \in A(\mathbf{x})} (\alpha_i - \beta_i) = 0,$$

982 where $\alpha_i := \text{softmax}(\mathbf{v})_i$ and $\beta_i := \text{softmax}(\mathbf{w})_i$.

983 Lemma 28 and Lemma 27 can be seen as complementary decompositions of the left and right side of
 984 Equation (4) for each unique one-hot vector in \mathbf{X} . As a result, we can restate Lemma 26 as follows.

985 **Lemma 29** (Decomposed Lemma 26). Let $\mathbf{v}, \mathbf{w} \in \mathbb{Q}^{1 \times L}$ with $\max_i \mathbf{v}_i = \max_i \mathbf{w}_i$ and let $\mathbf{X} \in$
 986 $\{0, 1\}^{L \times d}$ be a matrix whose rows are one-hot vectors, for some $L, d \in \mathbb{N}^+$. Further, we require \mathbf{X}
 987 to have at least two distinct rows. Then,

$$\sum_{i \in A(\mathbf{x})} (\alpha_i - \beta_i) = 0 \implies [\mathbf{v}]_{\mathbf{x}} = [\mathbf{w}]_{\mathbf{x}}, \quad (6)$$

988 for all $\mathbf{x} \in \text{set}(\mathbf{X})$, where $\alpha_i := \text{softmax}(\mathbf{v})_i$ and $\beta_i := \text{softmax}(\mathbf{w})_i$, and

$$\text{softmax}(\mathbf{v})\mathbf{X} = \text{softmax}(\mathbf{w})\mathbf{X} \iff [\mathbf{v}]_{\mathbf{X}} = [\mathbf{w}]_{\mathbf{X}}. \quad (7)$$

989 To prove Lemma 29, we leverage a known result about exponential numbers (as used within softmax)
 990 from transcendental number theory, namely that a set of exponential numbers with distinct rational
 991 coefficients is linearly independent, also known as the *Lindemann-Weierstrass theorem* [Baker, 1990].
 992 To understand intuitively how this theorem is used, let us assume for simplicity that the softmax is
 993 unnormalized, meaning that we can write the left side of Equation (6) as

$$\sum_{i \in A(\mathbf{x})} (\exp(\mathbf{v}_i) - \exp(\mathbf{w}_i)) = 0.$$

994 With the help of the Lindemann-Weierstrass theorem, we obtain the following claim, which we prove
 995 in Appendix E.

996 **Claim 30.** Let $A, B \subset \mathbb{Q}$ be finite multisets with $|A| = |B|$. Then, the sum

$$\sum_{a \in A} \exp(a) - \sum_{b \in B} \exp(b) = 0,$$

997 if, and only if, $A = B$.

998 Hence, with the unnormalized softmax, the left side of Equation (6) holds if and only if $[\mathbf{v}]_{\mathbf{x}} = [\mathbf{w}]_{\mathbf{x}}$.
 999 However, the full softmax also introduces normalization, which we denote with $Z_{\alpha} := \sum_{k=1}^n \exp(\mathbf{v}_k)$
 1000 and $Z_{\beta} := \sum_{k=1}^n \exp(\mathbf{w}_k)$, respectively. As a result, we have the condition

$$\sum_{i \in A(\mathbf{x})} (\alpha_i - \beta_i) = 0 \quad (8)$$

$$\Leftrightarrow \sum_{i \in A(\mathbf{x})} \frac{1}{Z_{\alpha}} \exp(\mathbf{v}_i) - \frac{1}{Z_{\beta}} \exp(\mathbf{w}_i) = 0 \quad (9)$$

$$\Leftrightarrow \sum_{i \in A(\mathbf{x})} \frac{\exp(\mathbf{v}_i) \cdot Z_{\beta} - \exp(\mathbf{w}_i) \cdot Z_{\alpha}}{Z_{\alpha} Z_{\beta}} = 0 \quad (10)$$

$$\Leftrightarrow \sum_{i \in A(\mathbf{x})} \exp(\mathbf{v}_i) \cdot Z_{\beta} - \exp(\mathbf{w}_i) \cdot Z_{\alpha} = 0 \quad (11)$$

$$\Leftrightarrow \sum_{i \in A(\mathbf{x})} \sum_{k=1}^L \exp(\mathbf{v}_i + \mathbf{w}_k) - \exp(\mathbf{w}_i + \mathbf{v}_k) = 0. \quad (12)$$

$$\Leftrightarrow \sum_{i \in A(\mathbf{x})} \sum_{k=1}^L \exp(\mathbf{v}_i + \mathbf{w}_k) - \sum_{j \in A(\mathbf{x})} \sum_{l=1}^L \exp(\mathbf{w}_j + \mathbf{v}_l) = 0. \quad (13)$$

Note that the multiset of exponents in the positive exponentials is $[v]_{x,w}$ and the set of exponents in the negative exponentials is $[w]_{x,v}$. Using Claim 30, we can now restate Lemma 29 once more as follows.

Lemma 31 (Multi-set only version of Lemma 29). *Let $v, w \in \mathbb{Q}^{1 \times L}$ with $\max_i v_i = \max_i w_i$ and let $X \in \{0, 1\}^{L \times d}$ be a matrix whose rows are one-hot vectors, for some $L, d \in \mathbb{N}^+$. Further, we require X to have at least two distinct rows. Then,*

$$[v]_{x,w} = [w]_{x,v} \implies [v]_x = [w]_x, \quad (14)$$

for all $x \in \text{set}(X)$ and

$$\text{softmax}(v)X = \text{softmax}(w)X \iff [v]_X = [w]_X. \quad (15)$$

We will give the full proof with all details in this section.

First, we will review some number theory background, formally state the Lindemann-Weierstrass theorem and its implications and then give the proof of Lemma 31. Using Lemma 31 and in particular, the equivalent Lemma 26, we finally prove Theorem 24.

Number theory We will formally introduce the necessary background on number theory and the Lindemann-Weierstrass theorem. A number is **algebraic** if it is the root of a non-zero single-variable polynomial with finite degree and rational coefficients. For example, all rational numbers $\frac{a}{b}$ with $a, b \in \mathbb{N}^+$ are algebraic, as they are the root of the polynomial $ax - b$ with integer coefficients. On the other hand, a number is **transcendental** if and only if it is not algebraic. For example, it is known that $\exp(a)$ is transcendental if a is algebraic and non-zero. This last fact follows from the Lindemann-Weierstrass theorem, which we state next [Baker, 1990].

Theorem 32 (Baker [1990], Theorem 1.4). *Let a_1, \dots, a_n be distinct **algebraic** numbers. Then, $\exp(a_1), \dots, \exp(a_n)$ are linearly independent with algebraic rational coefficients.*

Here, we will use the fact that attention uses the \exp function in the softmax and use Theorem 32 to compute injective representations of the GD-WL multisets by expressing them as sums of exponential numbers.

Proving Lemma 26 We now prove Lemma 31, equivalent to Lemma 26.

Lemma 33 (Proof of Lemma 31). *Let $v, w \in \mathbb{Q}^{1 \times L}$ with $\max_i v_i = \max_i w_i$ and let $X \in \{0, 1\}^{L \times d}$ be a matrix whose rows are one-hot vectors, for some $L, d \in \mathbb{N}^+$. Further, we require X to have at least two distinct rows. Then,*

$$[v]_{x,w} = [w]_{x,v} \implies [v]_x = [w]_x,$$

for all $x \in \text{set}(X)$ and

$$\text{softmax}(v)X = \text{softmax}(w)X \iff [v]_X = [w]_X.$$

Proof. Note that by assumption X has at least two distinct rows and hence, $A(x) \subset [L]$. As a result, the forward implication follows from the following claim.

Claim 34. For all $x \in \text{set}(X)$, if $\max_i v_i = \max_i w_i$ and $A(x) \subset [n]$, then, $[v]_{x,w} = [w]_{x,v} \implies [v]_x = [w]_x$.

Proof. Let $K := \max_i v_i = \max_i w_i$. We begin by sorting the entries in $[v]_x$ and $[w]_x$ in descending order, obtaining sorted vectors v^* and w^* . By assumption, we have that $v_1^* = w_1^* = K$. Now, let $i \in [|v]_x|$ be the smallest number for which $v_i^* \neq w_i^*$. If no such i exists, then $[v]_x = [w]_x$. Otherwise, without loss of generality, we assume that $v_i^* > w_i^*$. We now show that then, the sum $v_i^* + K$ appears at least once more in $[v]_{x,w}$ than in $[w]_{x,v}$.

First, note that there cannot exist some $j > i$ for which $w_j^* + K = v_i^* + K$. Second, for each $j < i$, $v_j^* = w_j^*$, meaning that for each such j where $v_j^* + K = v_i^* + K$ appears in $[v]_{x,w}$, $w_j^* + K = v_i^* + K$ appears in $[w]_{x,v}$.

Hence, $v_i^* + K$ appears at least once more in $[v]_{x,w}$ than in $[w]_{x,v}$, implying $[v]_{x,w} \neq [w]_{x,v}$. As a result, we have that $[v]_x = [w]_x \vee [v]_{x,w} \neq [w]_{x,v}$ which is logically equivalent to $[v]_{x,w} = [w]_{x,v} \implies [v]_x = [w]_x$. This shows the statement. \square

1044 To see why in Claim 34 it is important that $A(x)$ is a strict subset of $[n]$, we note that $A(x) = [L]$
 1045 implies $[v]_{x,w} = [w]_{x,v}$, irrespective of whether $[v]_x = [w]_x$. Notably, the proof holds if there
 1046 exists at least one $i \in [L] \setminus A(x)$, irrespective of whether $v_i = w_i$.

1047 The backward direction follows directly from the fact that $\text{softmax}(v)X$ and $\text{softmax}(w)X$ are
 1048 functions over $[v]_X$ and $[w]_X$, respectively.

1049 Together with Claim 34, this shows the statement. \square

1050 **Proving the GD-WL simulation result** Now that Lemma 26 has been proven, we will prove the
 1051 main result, Theorem 24, next.

1052 **Theorem 35** (Proof of Theorem 24). *Let $G := (V(G), E(G))$ be a graph with $n \in \mathbb{N}^+$ nodes and*
 1053 *node distance function $d_G : V(G)^2 \rightarrow \mathbb{Q}$. Let $d, d_f, T, h \in \mathbb{N}^+$ denote the number of embedding*
 1054 *dimensions, the number of hidden dimensions, the number of layers, and the number of attention*
 1055 *heads, respectively. Let $L := n + 1$ and let $\hat{X}^0 \in \mathbb{R}^{L \times d}$ and $B \in \mathbb{R}^{L \times L \times h}$ be initial token*
 1056 *embeddings and attention bias constructed according to Section 2 using node distance d_G . Then,*
 1057 *there exists weights for the parameters in $\Theta(d, d_f, T, h)$ such that $\hat{X}^t = X^t$, for all $t \geq 0$, an*
 1058 *arbitrary but fixed hash and using d_G as the distance function.*

1059 *Proof.* Note that the GD-WL produces a finite number of colors at each iteration and B is constructed
 1060 from d_G whose co-domain is compact for graphs with finite size. Hence, and since each transformer
 1061 layer is a composition of continuous functions, the domain of each transformer layer is compact.
 1062 Recall that we want to show that the t -th transformer layer can simulate

$$\chi_G^{t+1}(v) := \text{hash}([v]_{X^t}),$$

1063 for all $v \in V(G) \cup \{\text{[cls]}\}$, where $X^t \in \{0, 1\}^{L \times d}$ is a one-hot color matrix of the GD-WL colors
 1064 at iteration t . Let v be arbitrary but fixed. We say that v is the i -th node in an arbitrary but fixed
 1065 ordering of $V(G)$.

1066 We restate the definition of the transformer layer in a simplified form, omitting multiple heads, the
 1067 residual streams and LayerNorm. In particular, we state the layer only to update the i -th row of the
 1068 token matrix.

$$\hat{X}(v)^{t+1} = \text{MLP}(\text{softmax}(\hat{X}(v)^t W_Q (\hat{X}^t W_K)^T + v) \hat{X}^t W_V),$$

1069 where we recall that the i -th row of B is v . We now set $W_Q = W_K$ to all-zeros and W_V to the
 1070 identity matrix and obtain

$$\hat{X}(v)^{t+1} = \text{MLP}(\text{softmax}(v) \hat{X}^t).$$

1071 We prove the statement by induction over t . For the base case at $t = 0$, the token matrix X contains
 1072 the one-hot colors of the nodes in $V(G)$ as well as the special one-hot color of the [cls] token.
 1073 Setting $\hat{X}^0 := X$, we have that $\hat{X}^0 \in \{0, 1\}^{L \times d}$ and $\hat{X}^0 = X^0$. Further, due to the [cls] token,
 1074 we know that \hat{X}^0 has at least two distinct rows.

1075 Finally, note that, by construction, for each pair of distance vectors $\max_i v_i = \max_i w_i =$
 1076 $\max_{v,w \in V(G)} d_G(v, w) + 1$ and that every distance vector $v \in \mathbb{Q}^{1 \times L}$. These two conditions hold
 1077 throughout the induction and we will use them in the induction step to apply Lemma 26.

1078 In the induction step for $t > 0$, we assume that

- 1079 1. $\hat{X}^t \in \{0, 1\}^{L \times d}$
- 1080 2. $\hat{X}^t = X^t$
- 1081 3. \hat{X}^t has at least two distinct rows

1082 We want to prove that the same holds for $t + 1$. Let $v, w \in V(G) \cup \{\text{[cls]}\}$ be arbitrary but fixed.
 1083 Note that $\chi^{t+1}(v) = \chi^{t+1}(w)$ if and only if $[v]_{X^t} = [w]_{X^t}$. By the induction hypothesis, we have
 1084 that $[v]_{X^t} = [w]_{X^t}$ if and only if $[v]_{\hat{X}^t} = [w]_{\hat{X}^t}$. Further, by Lemma 26, we have that

$$\text{softmax}(v) \hat{X}^t = \text{softmax}(w) \hat{X}^t \iff [v]_{\hat{X}^t} = [w]_{\hat{X}^t},$$

1085 and as a consequence,

$$\text{softmax}(\mathbf{v})\hat{\mathbf{X}}^t = \text{softmax}(\mathbf{w})\hat{\mathbf{X}}^t \iff \chi^{t+1}(v) = \chi^{t+1}(w).$$

1086 Hence, there exists an injective function f that maps, for each token $v \in V(G) \cup \{\text{[cls]}\}$ with
 1087 distance vector \mathbf{v} , the vector $\text{softmax}(\mathbf{v})\hat{\mathbf{X}}^t$ to a one-hot vector of $\chi^{t+1}(v)$ with d dimensions. Since
 1088 the domain of the t -th transformer layer is compact, f is continuous. Hence, by universal function
 1089 approximation, there exists weights of the MLP such that, for each $v \in V(G) \cup \{\text{[cls]}\}$, $\hat{\mathbf{X}}^{t+1}(v)$
 1090 is a one-hot vector of $\chi^{t+1}(v)$. As a result,

- 1091 1. $\hat{\mathbf{X}}^{t+1} \in \{0, 1\}^{L \times d}$
- 1092 2. $\hat{\mathbf{X}}^{t+1} = \mathbf{X}^{t+1}$
- 1093 3. $\hat{\mathbf{X}}^{t+1}$ has at least two distinct rows.

1094 This completes the induction and proves the statement. \square

1095 C.2 Upper-bound on expressivity of the GDT

1096 Moreover, we can prove an expressivity upper-bound for the GDT using a technique adapted from
 1097 Müller et al. [2024]. We begin by showing the following result.

1098 **Lemma 36.** *Let $G := (V(G), E(G), \ell_V)$ be a graph with n nodes and without edge embeddings, let
 1099 $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d \times d}$ be arbitrary but fixed weight matrices with $d \in \mathbb{N}^+$, and let $\mathbf{B} \in \mathbb{Q}^{n \times n}$ be
 1100 an attention bias. Let*

$$\alpha(\mathbf{X}, \mathbf{U}) := \text{Attention}(\mathbf{X}\mathbf{W}_Q, \mathbf{X}\mathbf{W}_K, \mathbf{X}\mathbf{W}_V, \mathbf{B}).$$

1101 *There exists a distance function d_G over $V(G) \cup \{\text{[cls]}\}$ and functions f, h with*

$$f(\mathbf{X}_i) := h(\{(d_G(i, j), \mathbf{X}_j)\}),$$

1102 *such that for all \mathbf{X} and all i , $\alpha(\mathbf{X}, \mathbf{U})_i = f(\mathbf{X}_i)$.*

1103 *Proof.* We define d_G with have co-domain \mathbb{Q}^2 such that

$$d_G(i, j) = [\mathbf{B}_{ij}, I(i = j)],$$

1104 for all $i, j \in V(G) \cup \{\text{[cls]}\}$, where $[\cdot]$ is the concatenation operation and $I(i = j)$ is the indicator
 1105 function. We denote with $d_G(i, j)_k$ the k -th element in $d_G(i, j)$, for $k \in \{1, 2\}$. Let

$$g(\mathbf{X}_i, \mathbf{X}_j) := \exp(\mathbf{X}_i \mathbf{W}_Q (\mathbf{X}_j \mathbf{W}_K)^T).$$

1106 We choose h as follows. We note that by definition, $1 = d_G(i, i)_2 > d_G(i, j)_2$ for all $i \neq j$. Hence,
 1107 h can decompose its input into three arguments:

- 1108 1. \mathbf{X}_i , identified from the tuple $(d_G(i, j), \mathbf{X}_j)$ where $d_G(i, j)_2 = 1$, i.e., $i = j$,
- 1109 2. the multiset of distances $\{(d_G(i, j)_1)\}$,
- 1110 3. the multiset of vectors $\{\mathbf{X}_j\}$.

1111 Then, h computes

$$w_{ij} := \exp(g(\mathbf{X}_i, \mathbf{X}_j) + d_G(i, j)_1)$$

1112 and

$$\tilde{w}_{ij} := \frac{w_{ij}}{\sum_k w_{ik}},$$

1113 for all i, j . Finally, h computes

$$\sum_j \tilde{w}_{ij} \mathbf{X}_j \mathbf{W}_V,$$

1114 for all i , to obtain $\alpha(\mathbf{X}, \mathbf{U})_i$. \square

1115 Intuitively, the above lemma shows that biased attention can be written as a function over the multiset
 1116 in the GD-WL if the distance function is a metric. We use this result to show that the GD-WL is at
 1117 least as expressive as a GDT with relative PEs.

Proposition 37. Let $G := (V(G), E(G), \ell_V)$ be a graph without edge embeddings and let $\mathbf{B} \in \mathbb{Q}^{n \times n}$ be an attention bias. Let $d, d_f, T, h \in \mathbb{N}^+$ denote the number of embedding dimensions, the number of hidden dimensions, the number of layers, and the number of attention heads, respectively. For any choice of parameters $\Theta(d, d_f, T, h)$ for the GDT, there exists a distance function d_G over $V(G) \cup \{\text{cls}\}$ and a hash function hash for the GD-WL such that for all $t \geq 0$ and all pairs of nodes $i, j \in V(G)$, $\chi^t(i) = \chi^t(j)$ if and only if $\mathbf{X}_i^t = \mathbf{X}_j^t$.

Proof. We prove the statement by induction over t . For $t = 0$, the statement holds by definition, as the initial token embeddings \mathbf{X}^0 without any absolute PE are simply the node embeddings ℓ_V and the initial colors of the GD-WL are chosen to be consistent with the node embeddings ℓ_V . For $t > 0$, we assume by the induction hypothesis that for all pairs of nodes $i, j \in V(G)$, $\chi^{t-1}(i) = \chi^{t-1}(j)$ if and only if $\mathbf{X}_i^{t-1} = \mathbf{X}_j^{t-1}$. By definition, \mathbf{X}^t is computed via Equation (2), namely

$$\mathbf{X}^t := \text{MLP}(\text{Attention}(\mathbf{X}^{t-1}\mathbf{W}_Q, \mathbf{X}^{t-1}\mathbf{W}_K, \mathbf{X}^{t-1}\mathbf{W}_V, \mathbf{B})),$$

where the MLP is applied row-wise. Let

1. f denote the function in Lemma 36 consistent with projections $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V$ and attention bias \mathbf{B} ,
2. $\text{onehot}: [n] \rightarrow \{0, 1\}^n$ denote the function that maps numbers $1, \dots, n$ to their corresponding n -dimensional one-hot vector.

Then, for all i, j , we have that

$$\text{MLP} \circ f \circ \text{onehot}(\chi^{t-1}(i)) = \text{MLP} \circ f \circ \text{onehot}(\chi^{t-1}(j))$$

if and only if

$$\mathbf{X}_i^t = \mathbf{X}_j^t.$$

Finally, there at most n distinct rows in \mathbf{X}^t . Let h be a function that injectively maps each unique row in \mathbf{X}^t to a color in $[n]$. We choose $\text{hash} := h \circ \text{MLP} \circ f \circ \text{onehot}$, and have that, for all pairs of nodes $i, j \in V(G)$, $\chi^t(i) = \chi^t(j)$ if and only if $\mathbf{X}_i^t = \mathbf{X}_j^t$. This completes the induction and hence concludes the proof. \square

C.3 Expressivity of the GDT

Together, Theorem 24 and Proposition 37 correspond to the first and second statements in Theorem 1, respectively. A consequence of Theorem 24 is the fact that the GDT with NoPE is equivalent to the 1-WL. In particular, let $G := (V(G), E(G), \ell_V, \ell_E)$ be a graph with n nodes. Let $\ell_E(v, w) := \mathbf{1}$, for all $v, w \in V(G)$, where $\mathbf{1}$ is the vector containing 1 in every element. Further, let $\ell_E(v, v) := \mathbf{2}$, for all $v \in V(G)$, where $\mathbf{2}$ is the vector containing 2 in every element. Then, the GDT with NoPE is equivalent, according to Theorem 24, to the following update of the GD-WL:

$$\chi_G^{t+1}(v) := \text{hash}(\{\!\!\{d_G(v, w), \chi_G^t(w) : w \in V(G)\}\!\!\}),$$

where $d_G(v, v) = 2$, $d_G(v, w) = 1$ if $(v, w) \in E(G)$, and $d_G(v, w) = 0$, else, for all $v, w \in V(G)$. This can be equivalently written as

$$\chi_G^{t+1}(v) := \text{hash}((\chi_G^t(v), \{\!\!\{\chi_G^t(w) : w \in N_G(v)\}\!\!\})),$$

giving the 1-WL update rule.

D Proofs of Section 3

To guide our proofs of Section 3, we introduce CSL-graphs, obtaining the fact that RWSE cannot distinguish all of these graphs. These results are then expanded to provide an introduction to Theorem 2. We give the result and proof of Theorem 2 first with a simplified version for a minimum number of random walk steps, leveraging results obtained by Tönshoff et al. [2023]. In addition, the proofs of Proposition 3 and Proposition 47 are then given with additional details on the expressiveness hierarchy obtained from the definition of each PE. We provide further incremental results for our selection of PEs, complementing the results from Section 3.

Warm-up: CSL graphs We begin by introducing a class of simple and intuitive, yet not 1-WL distinguishable graphs, so-called CSL graphs. These graphs consist of a n node cycle with skip connections of length k, l originating from each node. CSL graphs are a canonical example of a graph class requiring a distance measure, motivating additional PEs [Rampásek et al., 2022, Müller et al., 2024]. Here, we show that they cannot be fully distinguished by RWSE and provide guidance on how to find pairs of CSL graphs indistinguishable by RWSE. We first introduce CSL graphs $G_{(n,k)}$ and their properties to prove the following results.

Definition 38. Let n, k be natural numbers and $k < n - 1$. $G_{(n,k)}$ defines an undirected graph which is 4-regular. The set of nodes is given by $V(G_{(n,k)}) = \{0, \dots, n - 1\}$. The edges are given by a two-step process. First, to construct a cycle in the CSL graph, every edge $\{i, i + 1\} \in E(G_{(n,k)})$ for $j \in \{0, \dots, n - 2\}$. Additionally $\{n - 1, 0\} \in E(G_{(n,k)})$ holds. Furthermore, the skip links are introduced by defining the sequence $s_1 = 0$ and $s_{i+1} = (s_i + k) \bmod n$ and deriving the edges with $\{s_i, s_{i+1}\} \in E(G_{(n,k)})$.

In addition, we introduce the notation used throughout the following proofs. Considering the skip links introduced in the CSL graphs we denote such a skip link by the mapping $s^k : V(G_{(n,k)}) \rightarrow V(G_{(n,k)})$ with $s(v_i) := v_{(i+k) \bmod n}$ for nodes $\{v_1, \dots, v_n\}$. A traversal to the next node v_{i+1} or v_{i-1} from node v_i is denoted by s^1 and s^{-1} respectively. We further provide specific random walks using a tuple of the visited nodes in a graph.

Proposition 39. Two CSL graphs $G_{(n_1,k_1)}$ and $H_{(n_2,k_2)}$ with $n_1 = n_2$ are non isomorphic if k_1, k_2 are co-prime natural numbers.

Given the definition of CSL graphs, it is possible to derive isomorphism results concerning such graphs. Furthermore, we note that CSL graphs are 1-WL indistinguishable, but can be distinguished by various WL variants such as the GD-WL [Zhang et al., 2023].

Proposition 40. There exists at least one pair of CSL graphs that RWSE cannot distinguish for any choice of random walk length.

Nonetheless, we note that the expressive power of RWSE is sufficient to distinguish many 1-WL indistinguishable graphs; for example, most CSL graphs can already be distinguished by RWSE. Furthermore, a minimum step number is given depending on the skip length of each CSL graph.

Proposition 41. Let $n, k \in \mathbb{N}^+$ with $n > k(k + 1) + 1$. RWSE can distinguish any pair of CSL graphs with n nodes and skip length $k, k + 1$, with a random walk length of $k + 1$.

With Definition 38 we derive the proofs for each lemma individually.

Proof of Proposition 41 For Proposition 41 we consider a subclass of CSL graphs. The minimum node number is specifically chosen to prevent a random walk of k steps to complete a cycle in the graph, even with using s^{k+1} exclusively for each step. We follow a two-step process for the proof: First, we gather paths existing in one graph but not the other. Then, in a second step, we show that all different paths of length at most k exist in both graphs. Further, we highlight that for random walks of length less than k , the paths are equal in both graphs.

Proof. Let $G_{(n,k)}$ and $H_{(n,k+1)}$ be two CSL graphs with skip link mappings s^k and s^{k+1} . Further let $n > k(k + 1) + 1$. To distinguish them, we denote the same node in both graphs with v_0 and w_0 . Then, for k random walk steps, we first examine whether there exist paths in $G_{(n,k)}$ which are not present in $H_{(n,k+1)}$. These include $(v_0, \dots, v_{k-1}, v_0)$ and $(v_0, \dots, v_{n-k+1}, v_0)$ as valid paths in $G_{(n,k)}$, which provide two k step walks with one skip link each. However, we note that such paths are not possible in $H_{(n,k+1)}$ as the skip link s^{k+1} has a length $k + 1$ and therefore no corresponding walk exists. In addition, we show that there exists no other walks in $H_{(n,k+1)}$, which are not present in $G_{(n,k)}$. For this, we must consider the cases of k even or odd.

Case 1: k is even: In this case, we must consider all combinations of skip links and s^1, s^{-1} functions. Since we assume an even number of steps, we know that all return walks with an even number of skip links or no skip links exist in both graphs. However, walks with an uneven number of skip links cannot return to v_0 or w_0 , as with the skip link size of k or $k + 1$, no return walks with at most k steps exist. As can be seen, the skip length does not influence the existence of any of the proposed walks; therefore, they exist both in $G_{(n,k)}$ and $H_{(n,k+1)}$.

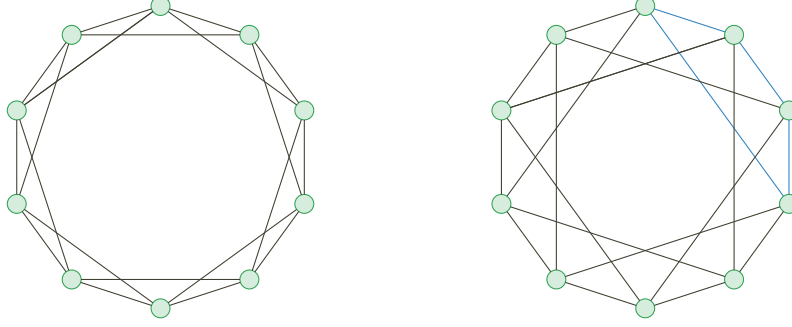


Figure 5: A pair of CSL graphs $G_{10,2}, H_{10,3}$. We note that the path marked in blue does not exist in graph G or has no replacement path.

Case 2: k is uneven: Given an uneven k , we conclude that no return walks with a length of k exist since steps can return neither a skip link nor an uneven number of any combination of skip links and steps can return to the origin node.

For step numbers lower than k the same cases apply in permuted order, depending on the step number. Also, due to the minimal n chosen, no instances exist where a circumference of the graph circle occurs in any combination of steps. As for both cases there exist no paths which are not present in $G_{(n,k)}$, it follows that there exist additional walks for $G_{(n,k)}$ which do not occur in $H_{(n,k+1)}$. Due to the structure of CSL graphs, the number of random walks, including both returning and non-returning random walks, is the same across all nodes in both graphs, resulting in the statement of Proposition 41.

Since the RWPE encoding is different if a single step returns a different return walk probability, RWPE can distinguish $G_{(n,k)}, H_{(n,k+1)}$ with the given random walk steps. \square

For the following pair of CSL graphs, we consider the computation of the number of returning walks. This enables us to directly compute the random walk probabilities for a single node. Due to the design of CSL graphs, each node in a graph has the same random walk return probability, thereby allowing us to derive the proof of Proposition 40.

Proof of Proposition 40 We consider two CSL graphs $G_{(11,3)}$ and $H_{(11,4)}$. With the computation of the number of returning walks of length r , W^r given by $W^r = \sum_{i=1}^n \lambda_i^r$, where λ_i denotes the eigenvalues of the adjacency matrix, we can compute the number of returning walks for each node in both graphs, since due to the graph structure the number of returning walks is equal for each node. Furthermore, we know that the number of total walks is equal in both graphs, given the graph structure as seen in Figure 5. From this, we compute the fraction of returning walks of varying length r for each node, resulting in the computation of the RWSE embedding. Since both number of returning walks and total number of walks are equal for each node in both graphs, as seen in the proof of Proposition 41, we receive the same RWSE embedding for $G_{(11,3)}$ and $H_{(11,4)}$.

In addition to the indistinguishability results obtained for RWSE, we want to propose a result that initially differentiates RWSE and RRWP. We then refine this result in Proposition 42, showcasing RRWP to be strictly more expressive than RWSE.

Proposition 42. *RRWP can distinguish all pairs of non-isomorphic CSL graphs.*

Proof of Proposition 42 In contrast to RWSE, the RRWP embedding uses the complete random walk matrix and information concerning the random walks between any two nodes. Because of this, RRWP can capture information that is not visible to RWSE due to the restriction to the diagonal of the random walk matrix.

Proof. We consider two arbitrary CSL graphs $G_{(n,i)}$ and $H_{(n,j)}$ with i, j co-prime and $i \neq j$. Then we can compute the RRWP encoding for each node by computing the random walk matrix. We saw from previous CSL graphs that the random walk matrix diagonal can be equal for both graphs, depending on the choice of i, j , and the number of random walk steps. However, for RRWP, we also

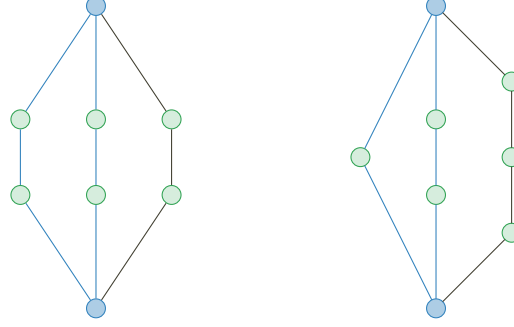


Figure 6: A pair of graphs from the construction method provided by Tönshoff et al. [2023]. Note that these graphs can only be distinguished by a returning random walk of length $\mathcal{O}(n)$ while being 1-WL distinguishable.

consider non-diagonal elements. Due to $i \neq j$, the RRWP tensor elements differ for each random walk of length 1, since different nodes are connected. Given any injective MLP layer, different RRWP tensor elements result in different RRWP embeddings for the nodes of both graphs, allowing the CSL graphs to be distinguished by RRWP. \square

Using the above results for CSL graphs, we derive a first bound associated with RWSE, depending on the graph structure of the CSL graphs. However, CSL graphs are not distinguishable by 1-WL, leaving a comparison between RWSE and 1-WL open. In the following, we want further to improve our understanding of RWSE and its expressiveness. We first provide an introduction and intermediate result given by Lemma 43 to limit the expressiveness to random walks of sufficient length. Afterwards, we provide the proof of Theorem 2, concluding our examination of RWSE.

Introduction to Theorem 2 To prove Theorem 2, we first provide an intermediate result from leveraging the graphs introduced by Tönshoff et al. [2023] for their work. This allows us to derive graphs that need a certain number of random walk steps, depending on their graph structure. We then expand on this concept in our proof of Theorem 2 by giving an example of graphs not distinguishable by RWSE.

Lemma 43. *There exists at least one pair of non-isomorphic graphs with order $3n - 1$ that can be distinguished by RWSE only with a random walk length of at least $\mathcal{O}(n)$*

Since RWSE requires returning random walks to construct the respective embedding, graphs exist that are only distinguishable by random walks of specific length. However, we want to provide a class of graphs requiring at least $\mathcal{O}(n)$ steps. Furthermore, we want these graphs to be 1-WL distinguishable, providing a first step to Theorem 2, where we show the indistinguishability of RWSE and 1-WL. Following Tönshoff et al. [2023] and their evaluation of another random walk-based GNN architecture (CraWL), we adapt their counterexample to our evaluation of RWSE.

Proof. We provide a proof by giving a constructed graph only distinguishable by random walks of length greater than $n - 1$, fulfilling the necessary condition of $\mathcal{O}(n)$ for the walk length. Following Tönshoff et al. [2023] with their counterexample for the CraWL algorithm, we adapt the corresponding graphs for the RWSE embedding. Since we only consider returning random walk probabilities in the RWSE embedding, we disregard any information from other random walks.

Given the graphs in Figure 6 with n nodes, we consider the blue marked nodes. Since the returning walks are the same for both nodes for any walk of length $r \leq n - 1$, the graphs cannot be distinguished for any RWSE embedding with a walk length of r . Due to the graph construction, the corresponding walks for all other nodes are the same in both graphs. However, due to the cycle colored in blue, the return walk probabilities differ in both graphs for walks with a length $r' \geq n$. This results in a pair of graphs only distinguishable by random walks of length at least n . Results from Tönshoff et al. [2023] allow for constructing further examples with n nodes and order $3n - 1$. By construction, these graphs are distinguishable by 1-WL [Tönshoff et al., 2023], resulting in the stated lemma. \square

1282 With the results from Lemma 43, we can now expand the set of graphs not distinguishable by RWSE,
 1283 while 1-WL distinguishable. Combining both results, we can determine a set of graphs limiting the
 1284 expressiveness of RWSE and further investigate the expressive power of random walks.

1285 We provide an example pair of trees not distinguishable by the RWSE encoding. In contrast, all trees
 1286 are known to be distinguishable by the 1-WL algorithm [Cai et al., 1992]. Combining the findings
 1287 from Theorem 2 with the CSL graph results in Proposition 41, it follows that RWSE embeddings are
 1288 incomparable to the 1-WL color refinement algorithm.

1289 **Proof of Theorem 2** For this proof, we first consider a pair of trees shown by Cvetković [1988].
 1290 Originally introduced as an example of trees with differing eigenvalues and graph angles, thereby
 1291 being distinguishable by the EA-invariant, these graphs are not distinguishable by the RWSE embed-
 1292 ding for an arbitrary number of random walk steps. The proof is split into multiple steps, containing
 1293 parts of both trees that need to be considered.

1294 *Proof.* We separate the graph as shown in figure 7 into backbones I, J and subtrees X, Y_1, Y_2 and
 1295 the original tree G . X stays the same throughout both graphs, whereas Y_1, Y_2 change their edges
 1296 connecting them to the graph. In the first step, we consider the return probability to one of the
 1297 backbone nodes. Going from either of the backbone nodes to the original tree G has a probability of
 1298 $\frac{1}{5}$. Furthermore, going to X from either backbone node has the probability $\frac{2}{5}$. Once in either part of
 1299 X , the return probability to the backbone node is given by p_X , which is equal in both graphs. Finally,
 1300 the probability of going to Y_1, Y_2 is denoted by p_{Y_1}, p_{Y_2} respectively, again equal for both graphs.
 1301 Further, the return probabilities Y_1 and Y_2 are equal, as can be easily seen from the combination of
 1302 Y_1 and Y_2 with the respective backbone node. This allows for a complete evaluation of the backbone
 1303 nodes, assigning them probabilities p_I and p_J .

1304 Secondly, we consider random walks only on the nodes of the subtree X , not returning to either
 1305 of the backbone nodes. Since X is equal in both graphs, the return probabilities are also the same,
 1306 denoted by x . However, they differ from the return probabilities in Y_1 and Y_2 , given by y_1 and y_2
 1307 respectively. Because of that RWSE can distinguish X, Y_1, Y_2 without connections to the backbone
 1308 nodes.

1309 Combining our previous knowledge, we now consider return probabilities for nodes in Y_1 and Y_2
 1310 without restricting ourselves to the subtrees. We know that the probability of walking towards the
 1311 backbone node is given by the respective position of each node in Y_1, Y_2 , equal in both trees. Once
 1312 the random walk arrives at either backbone node, the probability to return to said backbone is given
 1313 by either p_I or p_J and a probability of $\frac{2}{5}$ to return to the originating subgraph. As p_I and p_J are
 1314 equal in both graphs and p_{Y_1}, p_{Y_2} are equal, it follows that the return probability for nodes $y_1 \in Y_1$
 1315 are equal across both graphs, denoted by p_{y_1T} . The same holds for Y_2 with the return probabilities
 1316 denoted by p_{y_2T} .

1317 Finally, for both graphs, the nodes are assigned the RWSE probability vectors p_I, p_J once, p_X 17
 1318 times, and p_{y_1T}, p_{y_2T} 8 times each. Furthermore, the two backbone nodes of both graphs cannot be
 1319 distinguished using the RWSE embedding. Therefore, the graphs are equal under RWSE, however as
 1320 shown by [Cai et al., 1992] every pair of trees can be distinguished using the 1-WL test, resulting in
 1321 the incomparability of RWSE and the 1-WL test. \square

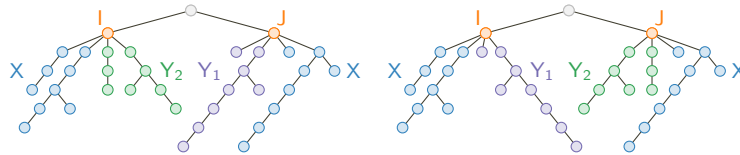


Figure 7: A pair of trees given by Cvetković [1988]. The grey node denotes an arbitrary tree concatenated to the existing tree. Each part of the tree is colored according to its respective appearance. We further label both backbone nodes in orange, denoting the left backbone node with I and the right backbone node with J .

Combining the results of Proposition 40, Lemma 43, and Theorem 2, we obtain fine-grained observations of graph structures not distinguishable by RWSE. While random walks are sufficiently robust to distinguish many common graph structures and graphs not distinguishable by 1-WL, RWSE still fails to distinguish specific trees and graphs originally proposed by Tönshoff et al. [2023].

In the following, we provide proofs supplementing the theoretical expressiveness hierarchy introduced by Zhang et al. [2024] and Black et al. [2024]. We first give an intermediate result relating RWSE and RRWP, thereby providing a lower bound for RRWP and an upper bound for RWSE. Further, we propose adapted results for LPE and SPE, comparing them to other PEs and to each other. The significant results are shown in Proposition 3 and Proposition 4.

Proof of Proposition 3 The proof of Proposition 3 is given by simply evaluating the corresponding embeddings for both RWSE and RRWP. Since both embeddings use the same MLP encoder layer, we restrict ourselves to directly evaluating random walk matrices. We first state an extended version of Proposition 3 and provide a proof.

Lemma 44 (Proposition 3 in the main paper). *Let G, H be two non-isomorphic graphs and $\mathbf{P}_k^{RW}(G), \mathbf{P}_k^{RW}(H)$ the generated RWSE encodings for both graphs with a random walk length k . Then for the generated RRWP encodings $\mathbf{P}_k^{RR}(G), \mathbf{P}_k^{RR}(H)$ it follows:*

$$\mathbf{P}_k^{RR}(G) = \mathbf{P}_k^{RR}(H) \Rightarrow \mathbf{P}_k^{RW}(G) = \mathbf{P}_k^{RW}(H).$$

In addition, at least one pair of graphs exists, distinguishable by RRWP, which is not distinguishable by RWSE.

Proof. Given random walk matrices $\mathbf{R}_G := \mathbf{D}_G^{-1} \mathbf{A}_G, \mathbf{R}_H := \mathbf{D}_H^{-1} \mathbf{A}_H$, its power matrices up to the power of k and the corresponding RRWP embeddings $\mathbf{P}_k^{RR}(G), \mathbf{P}_k^{RR}(H)$ for two non-isomorphic graphs we can directly deduce that for each tensor in the RRWP embeddings, corresponding to the RRWP embedding for a single node in each graph, the diagonal elements of the random walk matrix are the same for each power up to k . Therefore, it directly follows that $\mathbf{P}_k^{RR}(G) = \mathbf{P}_k^{RR}(H)$ results in $\mathbf{P}_k^{RW}(G) = \mathbf{P}_k^{RW}(H)$, with $\mathbf{P}_k^{RW}(G), \mathbf{P}_k^{RW}(H)$ denoting the RWSE encodings obtained from the same random walk matrices. We provide a simple example of a pair of graphs not distinguishable by the RWSE embedding, whereas RRWP can distinguish between the two graphs. This example follows directly from Proposition 40 since it is proven there that RWSE cannot distinguish this pair of graphs. However, it directly follows from the definition of the random walk matrix. It assumes that the MLP encoder preserves identity and that both graphs can be distinguished due to their differences in skip lengths as determined by the definition of CSL graphs. \square

We now want to consider the expressive power of LPE concerning RWSE. For this, we use results obtained by Black et al. [2024] and Zhang et al. [2024] and combine them with similar results obtained for SignNet and BasisNet [Lim et al., 2023]. Our additional observations are highlighted in Proposition 4.

Proof of Proposition 4 We provide an expanded version of Proposition 4 split into three parts (Lemma 45, Lemma 46, and Proposition 47), considering the expressive power of SAN [Kreuzer et al., 2021] and RWSE first, expanding the proof to LPE and SPE in Lemma 46, and finally upper bounding RRWP by using SPE. With these three parts, we are then able to derive the proposition.

Lemma 45 (RWSE and SAN). *Let the number k of eigenvalues $\lambda \in \mathbb{R}^n$ and eigenvectors $\mathbf{V} \in \mathbb{R}^{n \times n}$ used in the SAN encoding be equal to the number of nodes in non-isomorphic graphs G, H . Then given the encodings $\mathbf{P}_k^{SAN}(G), \mathbf{P}_k^{SAN}(H)$ with it follows:*

$$\mathbf{P}_k^{SAN}(G) = \mathbf{P}_k^{SAN}(H) \Rightarrow \mathbf{P}^{RW}(G) = \mathbf{P}^{RW}(H),$$

for a pair of RWSE encodings $\mathbf{P}^{RW}(G), \mathbf{P}^{RW}(H)$ and an arbitrary number of random walk steps.

Lemma 46. (RWSE and LPE) *Given Lemma 45 and the LPE embeddings $\mathbf{P}_k^{LPE}(G), \mathbf{P}_k^{LPE}(H)$ for two non-isomorphic graphs G, H with k nodes it follows:*

$$\mathbf{P}_k^{LPE}(G) = \mathbf{P}_k^{LPE}(H) \Rightarrow \mathbf{P}^{RW}(G) = \mathbf{P}^{RW}(H),$$

for a pair of RWSE embeddings $\mathbf{P}^{RW}(G), \mathbf{P}^{RW}(H)$ and an arbitrary number of random walk steps. The same result follows by replacing LPE with SignNet, BasisNet, or SPE as the eigenvector-based embedding.

Proof of Lemma 45 and Lemma 46 In the following, we provide the proofs for both lemmas. Since Lemma 46 proposes an extension of the previous lemma, we first show the specialized case for the SAN embedding and expand it to the more general case of eigenvector-based encodings. With proofs provided for both lemmas, we can directly derive Proposition 4 by combining them with Proposition 47.

Proof. The proof follows the proof for the comparison between SignNet and RWSE made by Lim et al. [2023]. Since the RWSE embedding is determined by the random walk matrix and its powers, we first determine a corresponding relation between the random walk matrix ($\mathbf{D}^{-1}\mathbf{A}$) and eigenvalues and eigenvectors of the normalized graph Laplacian. Due to the definition of the random walk matrix, the eigenvectors of said matrix are determined by $v_i^R = D^{-1/2}v_i$, where v_i denotes the corresponding eigenvector of the normalized graph Laplacian. This results in the following equation relating the random walk matrix diagonal to the eigenvectors of the graph Laplacian [Lim et al., 2023]:

$$(\text{diag}(\mathbf{D}^{-1}\mathbf{A}))^k = \text{diag}\left(\sum_{i=1}^k (1 - \lambda_i)^k v_i v_i^T\right). \quad (16)$$

Following Lim et al. [2023], the linear layer can approximate $\sum_{i=0}^k$ and the transformer encoder to approximate $(1 - \lambda_i)^k$ as both are permutation equivariant functions from vectors to vectors. Since eigenvalues and eigenvectors are directly given to the SAN embedding and the linear and transformer encoder layer being able to approximate $(1 - \lambda_i)^k$ for each λ_i the approximation directly follows. This approximation assumes using all k eigenvalues and the complete eigenvectors obtained from the decomposition.

For Lemma 46 we consider Equation (16). However, for each embedding, we must consider whether eigenvalues and eigenvectors can be recovered and $(1 - \lambda_i)^k$ can be approximated. We split the following proof for each encoding and assume we use all eigenvalues and eigenvectors.

For LPE, we can directly recover eigenvalues and eigenvectors from the input to each embedding by using the eigenvectors \mathbf{V}_i passed to LPE. Using a sufficiently expressive ϕ and ρ LPE is able to approximate $(1 - \lambda_i)^k$. This follows directly from the assumption of ϕ and ρ to be permutation equivariant MLPs or more expressive neural network architectures, thereby able to approximate the given functions as in the SAN case [Lim et al., 2023].

For SPE, a slightly different case has to be considered. Since SPE uses the projection matrices obtained from $\mathbf{V}\mathbf{V}^T$, the eigenvectors must be recovered.

Instead of directly recovering eigenvectors, we use the properties of the underlying projection matrices. From this, we can directly recover the eigenvectors needed from $\mathbf{V}\text{diag}(\phi_i(\lambda))\mathbf{V}^T$ for a suitable ϕ , which can be reverted by ρ to retain the eigenvectors. For the eigenvalues, we consider ϕ_i to be eigenvalue-preserving functions, allowing us to recover the eigenvalues from the diagonalized representation. The remaining proof follows from the observations made by Lim et al. [2023] for SignNet and BasisNet. In addition, SignNet and BasisNet Lim et al. [2023] prove that both embeddings can approximate the RWSE embedding given suitable ϕ and ρ . \square

Proposition 47 (RRWP and SPE). *Given the SPE embeddings $\mathbf{P}_k^{\text{SPE}}(G)$, $\mathbf{P}_k^{\text{LPE}}(H)$ for two non-isomorphic graphs G, H with k nodes it follows:*

$$\mathbf{P}_k^{\text{SPE}}(G) = \mathbf{P}_k^{\text{SPE}}(H) \Rightarrow \mathbf{P}^{\text{RR}}(G) = \mathbf{P}^{\text{RR}}(H),$$

for a pair of RRWP embeddings $\mathbf{P}^{\text{RR}}(G)$, $\mathbf{P}^{\text{RR}}(H)$ and an arbitrary number of random walk steps.

With the partial hierarchy obtained for LPE and SPE, we want to look further at random walk-based PEs. Since RWSE is upper bounded by LPE and incomparable to the 1-WL, it remains to propose an upper bound of RRWP, known to be more expressive than RWSE from Proposition 3.

Proof of Proposition 47 Following Zhang et al. [2024] with their proof of a representation of the page rank distance using projection matrices, we show that RRWP can be represented using the page rank distance and that such distance can be approximated using information recovered from the SPE embedding. We note that a proof of SPE being more expressive than GRIT is provided by Zhang et al. [2024]. Nonetheless, we reduce the proof only to involve the RRWP embedding to align with our theory framework.

1416 First, we consider the representation of the RRWP embedding using the generalized PageRank
 1417 distance. For this, we consider RRWP as a distance-based embedding of the form

$$\mathbf{P}_k^{\text{RR}}(u, v) = [\mathbf{D}^{-1}\mathbf{A}(u, v), (\mathbf{D}^{-1}\mathbf{A})^2(u, v), \dots, (\mathbf{D}^{-1}\mathbf{A})^k(u, v)],$$

1418 for nodes $u, v \in V(G)$. Thereby, the RRWP embedding for a selection of nodes can be represented
 1419 by the multi-dimensional page rank distance PR for a given weight sequence γ_i :

$$PR(u, v) = \sum_{i=0}^{\infty} \gamma_i (\mathbf{D}^{-1}\mathbf{A})^i(u, v).$$

1420 From Zhang et al. [2024] we obtain the following equality satisfying the needed relation between
 1421 page rank distance and projection matrices.

$$\sum_{k=0}^{\infty} \gamma_k (\mathbf{D}^{-1}\mathbf{A})^k = \sum_i \left(\sum_{k=0}^{\infty} \gamma_k (1 - \lambda_i)^k \right) \mathbf{P}_i(u, v) (\deg(u)^{-1/2}) (\deg(v)^{-1/2})$$

1422 with $\mathbf{P}_i(u, v)$ denoting the element at position (u, v) of the i -th projection matrix. However, we
 1423 still need to recover node degree information from the SPE encoding and prove that SPE retains the
 1424 projection matrix information.

1425 Given the property of projection matrices to recover the underlying matrix using eigenvalue decom-
 1426 position [Zhang et al., 2024], we recover node degree information using the diagonal of the graph
 1427 Laplacian. Since SPE uses the graph Laplacian to compute eigenvalues and eigenvectors, we can
 1428 directly recover relevant degree information via the following equation.

$$\begin{aligned} \mathbf{L} &= \sum_{i=1}^n \lambda_i \mathbf{P}_i \\ \text{diag}(\mathbf{L}) &= \text{diag}\left(\sum_{i=1}^n \lambda_i \mathbf{P}_i\right) = \sum_{i=1}^n \lambda_i \text{diag}(\mathbf{P}_i) \\ \mathbf{L}_{uu} &= \sum_{i=1}^n \lambda_i (\mathbf{P}_i(u, u)) \end{aligned}$$

1429 Following the definition of the SPE encoding, eigenvalues and the elements of the projection matrix
 1430 can be recovered using suitably expressive ϕ and ρ . Given ϕ to be a 2-IGN and ρ to be a MLP or
 1431 1-WL expressive GNN, $\deg(u)^{-1/2}$ can be approximated by ρ , whereas $\sum_i \left(\sum_{k=0}^{\infty} \gamma_k (1 - \lambda_i)^k \right)$
 1432 can be approximated by a 2-IGN as shown by Maron et al. [2019a], Lim et al. [2023]. This allows
 1433 for approximating the RRWP embedding via the PageRank distance using SPE as an upper bound,
 1434 concluding our proof.

1435 Combining the results of Lemma 45, Proposition 47, and Proposition 3, we obtain Proposition 4
 1436 directly, supplementing the hierarchy of PEs in their theoretical expressiveness. These results provide
 1437 a comprehensive theoretical expressiveness hierarchy, showing that random walk-based embeddings
 1438 are expressive compared to the 1-WL test but are bounded by eigeninformation-based embeddings.
 1439 We note that all embeddings are bounded by the 3-WL test as shown by Zhang et al. [2024].

1440 **Additional results on theoretical expressiveness** Using notation established in section Ap-
 1441 pendix B.3 we provide additional proofs to complement the framework established by Black et al.
 1442 [2024] and Zhang et al. [2024] concerning theoretical expressiveness of PEs. At first, we consider
 1443 the proof of Lemma 48, highlighting the connection between SAN and SignNet. Then we consider
 1444 SignNet and BasisNet, expanding on the results of Lim et al. [2023] and adapting them to LPE. We
 1445 provide these additional results to improve the hierarchy of theoretical expressiveness in PEs and
 1446 additional results concerning LPE, relating it to other eigenvector-based embeddings.

1447 Throughout these proofs, we consider the respective ϕ and ρ to be selected as MLPs or GNNs. For
 1448 ϕ , we choose, based on previous analysis by Zhang et al. [2024], a function mapping at most as
 1449 expressive as a 2-IGN. Similarly, for ρ we select any 1-WL expressive GNN or MLP. Note that these
 1450 assumptions differ from the selections made in our empirical evaluation.

1451 **Lemma 48.** *Given a sufficiently expressive ϕ and ρ for SignNet, aligning with the implementation*
 1452 *of Lim et al. [2023] and the original implementation of the SAN embedding, SignNet is at least as*
 1453 *expressive as the SAN embedding.*

1454 *Proof.* Let SAN and SignNet be represented by the respective color refinement algorithms shown in
 1455 Definition 14. To show Lemma 48, we need to show that

$$T_{GP} \circ T_{WL}^\infty \circ T_{SP2} \circ T_\phi(\chi\text{Sign}) \preceq T_{GP} \circ T_{ENC} \circ T_L(\chi_{SAN}).$$

1456 Since we assume a standard transformer encoder to be at most 1-WL expressive and knowing that
 1457 T_{GP} is order preserving concerning Definition 10, we can reduce the above equation to the following
 1458 expression:

$$T_{WL}^\infty \circ T_{SP2} \circ T_\phi(\chi\text{Sign}) \preceq T_{WL}^\infty \circ T_L(\chi_{SAN}).$$

1459 This expression can now be evaluated. Given two non-isomorphic graphs G, H and arbitrary nodes
 1460 $u, v \in V(G)$ and $x, y \in V(H)$ the following holds true:

$$\begin{aligned} T_{WL}^\infty \circ T_{SP2} \circ T_\phi(\chi\text{Sign}(u, v)) &= T_{WL}^\infty \circ T_{SP2} \circ T_\phi(\chi\text{Sign}(x, y)) \\ T_{SP2} \circ T_\phi(\chi\text{Sign}(u, v)) &= T_{SP2} \circ T_\phi(\chi\text{Sign}(x, y)) \\ \{\!\{ T_\phi(\chi\text{Sign}(u, v)) \}\!\} &= \{\!\{ T_\phi(\chi\text{Sign}(x, y)) \}\!\} \\ \{\!\{ T_\phi(\lambda_G, \mathbf{V}^u, \mathbf{V}^v) \}\!\} &= \{\!\{ T_\phi(\lambda_H, \mathbf{V}^x, \mathbf{V}^y) \}\!\}. \end{aligned}$$

1461 From the equivalence of the multisets, it follows directly that $\chi_{SAN}(u, v) = \chi_{SAN}(x, y)$ holds for any
 1462 choice of nodes given an injective T_ϕ . \square

1463 With the conclusion of the proof for the SAN embedding, we further evaluate the connection between
 1464 BasisNet and SignNet and the LPE embedding. First, we show that BasisNet can approximate
 1465 SignNet, an observation highlighting the differences in expressiveness noted by Lim et al. [2023].
 1466 In the second part of the proof, we then conclude our comparison of eigenvector-based embeddings
 1467 with the comparison of LPE to BasisNet. Throughout the proof we again assume ϕ_{SN}, ϕ_{LPE} to be at
 1468 most 2-IGN expressive and ρ_{SN}, ρ_{LPE} to be 1-WL expressive.

1469 *Proof.* Let SignNet and BasisNet be represented by the color refinement algorithms from Defini-
 1470 tion 14. Then for two non-isomorphic graphs G, H with nodes $u, v \in V(G)$ and $x, y \in V(H)$ we
 1471 consider the color refinement algorithms for both encodings. With this it follows that we have to
 1472 show $T_{GP} \circ T_{WL} \circ T_{SP1} \circ T_{BP} \circ T_{SIAM}(\chi_{Basis}) \preceq T_{GP} \circ T_{WL} \circ T_{SP2} \circ T_\phi(\chi\text{Sign})$. Since $T_{GP} \circ T_{WL}$ is
 1473 order preserving we only consider the relation $T_{SP1} \circ T_{BP} \circ T_{SIAM}(\chi_{Basis}) \preceq T_{SP2} \circ T_\phi(\chi\text{Sign})$. First
 1474 of all, we show that $T_{SIAM}(\chi_{Basis}) \preceq T_\phi(\chi\text{Sign})$:

$$\begin{aligned} T_{SIAM}(\chi_{Basis})(\lambda_G, u, v) &= T_{SIAM}(\chi_{Basis})(\lambda_H, x, y) \\ \Rightarrow [T_{IGN}(\chi_{Basis}(\lambda_G, \cdot, \cdot))]_G(u, v) &= [T_{IGN}(\chi_{Basis}(\lambda_H, \cdot, \cdot))]_H(x, y). \end{aligned}$$

1475 Using the definition of the IGN color refinement, we can directly approximate the eigenvalues used
 1476 in the initial encoding of SignNet. Furthermore, a 2-IGN architecture is at least as expressive as
 1477 the architectures used for ϕ in SignNet. Since the multisets of the projection matrices allow us to
 1478 approximate the eigenvectors used by the SignNet encoding, the initial encoding of SignNet can be
 1479 approximated, allowing for the approximation of $T_\phi(\chi\text{Sign})$,

$$\begin{aligned} [T_{IGN}(\chi_{Basis}(\lambda_G, \cdot, \cdot))]_G(u, v) &= [T_{IGN}(\chi_{Basis}(\lambda_H, \cdot, \cdot))]_H(x, y) \\ \Rightarrow \chi\text{Sign}(\lambda_G, u, v) &= \chi\text{Sign}(\lambda_H, x, y) \Rightarrow T_\phi(\chi\text{Sign})(\lambda_G, u, v) = T_\phi(\chi\text{Sign})(\lambda_H, x, y). \end{aligned}$$

1480 Given that $\bar{\chi} = T_{SIAM}(\chi_{Basis})$, we now only have to show that $T_{BP}(\bar{\chi}) \preceq T_{SP2}(\bar{\chi})$. Using the same
 1481 nodes as above:

$$\begin{aligned} T_{BP}(\bar{\chi})(\lambda, u) &= T_{BP}(\bar{\chi})(\lambda, x) \\ \Rightarrow \bar{\chi}_G(\lambda, u, u) &= \bar{\chi}_H(\lambda, x, x) \wedge \{\!\{ \bar{\chi}_G(\lambda, u, v) : v \in V(G) \}\!\} = \{\!\{ \bar{\chi}_G(\lambda, x, v) : v \in V(H) \}\!\} \wedge \\ &\quad \{\!\{ \bar{\chi}_G(\lambda, v, u) : v \in V(G) \}\!\} = \{\!\{ \bar{\chi}_H(\lambda, v, x) : v \in V(H) \}\!\} \wedge \\ &\quad \{\!\{ \bar{\chi}_G(\lambda, v, v) : v \in V(G) \}\!\} = \{\!\{ \bar{\chi}_H(\lambda, v, v) : v \in V(H) \}\!\} \wedge \\ &\quad \{\!\{ \bar{\chi}_G(\lambda, v, w) : v, w \in V(G) \}\!\} = \{\!\{ \bar{\chi}_H(\lambda, v, w) : v, w \in V(H) \}\!\} \wedge \\ &\quad \Rightarrow T_{SP2}(\bar{\chi})(\lambda, u, v) = T_{SP2}(\bar{\chi})(\lambda, x, y). \end{aligned}$$

1482 Since both parts of the relation $T_{BP} \circ T_{SIAM}(\chi_{\text{Basis}}) \preceq T_{SP2} \circ T_\phi(\chi_{\text{Sign}})$ hold and all color refinements
 1483 are considered to be order preserving and expressiveness preserving the proof directly follows.

1484 In case of the LPE embedding, the proof follows the same structure with T_ϕ being replaced with
 1485 T_ϕ^{LPE} as given in Definition 14. Since we do not assume T_ϕ to be more expressive than T_ϕ^{LPE} and
 1486 both being bounded by a 2-IGN in expressiveness, we can replace T_ϕ , and therefore, we omit the
 1487 proof. \square

1488 E Additional technical proofs

1489 **Multiset operations** Let D be a finite set with an arbitrary but fixed order. We denote with D_i the i -
 1490 th element in the order over D . Let A be a finite multiset over D . We write $A := \{(a_i, D_i) \mid i \in [|D|]\}$
 1491 with $a_i \geq 0$, the multiplicity of element D_i in A .

1492 We define $|A| := \sum_i a_i$. Further, let $B := \{(b_i, D_i) \mid i \in [|D|]\}$ be another finite multiset over D .
 1493 We define

$$A \cap B := \{(\min\{a_i, b_i\}, D_i) \mid i \in [|D|]\}$$

1494 and

$$A \setminus B := \{(\max\{a_i - b_i, 0\}, D_i) \mid i \in [|D|]\}.$$

1495 We note that $A \cap B$ is symmetric while $B \setminus A$ is not symmetric. Nonetheless, we prove that if
 1496 $|A| = |B|$, then $|A \setminus B|$ is symmetric.

1497 **Claim 49.** Let A, B be two multisets over a finite domain. If $|A| = |B|$, then $|A \setminus B| = |B \setminus A|$.

1498 *Proof.* We have that

$$|A \setminus B| = \sum_i \max\{a_i - b_i, 0\} = \sum_i a_i - \min\{a_i, b_i\} = |A| - |A \cap B|.$$

1499 Hence, if $|A| = |B|$, then $|A \setminus B| = |A| - |A \cap B| = |B| - |A \cap B|$ and since $|A \cap B|$ is symmetric,
 1500 $|A \setminus B| = |B| - |B \cap A| = |B \setminus A|$. \square

1501 **Claim 50** (Proof of Claim 30). Let $A, B \subset \mathbb{Q}$ be finite multisets with $|A| = |B|$. Then, the sum

$$\sum_{a \in A} \exp(a) - \sum_{b \in B} \exp(b) = 0,$$

1502 if, and only if, $A = B$.

1503 *Proof.* Let $f(A, B) = \sum_{a \in A} \exp(a) - \sum_{b \in B} \exp(b)$. Note that for each element a in A that also
 1504 appears as b in B , we have that $\exp(a) - \exp(b) = 0$. Hence, we define $A^* := A \setminus B$ and $B^* := B \setminus A$
 1505 and have that $f(A, B) = f(A^*, B^*)$. Further, since according to the lemma statement $|A| = |B|$, we
 1506 have that $|A^*| = |B^*|$; see Claim 49.

1507 We first show that $f(A, B) = 0$ if and only if $A = B$. To this end, note that the sum is 0 if the
 1508 positive and the negative summands cancel out, that is, if $A^* = B^* = \emptyset$ and hence, $A = B$. If
 1509 $A \neq B$, then the above sum is a non-zero sum of exponentials with algebraic exponents, and thus, by
 1510 Theorem 32, non-zero. Hence, we have $A = B \Leftrightarrow f(A, B) = 0$. \square

1511 **Lemma 51** (Proof of Lemma 28). Let $\mathbf{v}, \mathbf{w} \in \mathbb{Q}^{1 \times L}$ and let $\mathbf{X} \in \{0, 1\}^{L \times d}$ be a matrix whose rows
 1512 are one-hot vectors, for some $L, d \in \mathbb{N}^+$. Then, $\text{softmax}(\mathbf{v})\mathbf{X} = \text{softmax}(\mathbf{w})\mathbf{X}$, if and only if, for
 1513 every $\mathbf{x} \in \text{set}(\mathbf{X})$,

$$\sum_{i \in A(\mathbf{x})} (\alpha_i - \beta_i) = 0,$$

1514 where $\alpha_i := \text{softmax}(\mathbf{v})_i$ and $\beta_i := \text{softmax}(\mathbf{w})_i$.

1515 *Proof.* We have

$$\text{softmax}(\mathbf{v})\mathbf{X} - \text{softmax}(\mathbf{w})\mathbf{X} = \sum_{i=1}^n (\alpha_i - \beta_i) \cdot \mathbf{X}_i = \sum_{\mathbf{x} \in \text{set}(\mathbf{X})} \sum_{i \in A(\mathbf{x})} (\alpha_i - \beta_i) \cdot \mathbf{x}.$$

1516 Since the rows of \mathbf{X} are one-hot vectors, $\text{set}(\mathbf{X})$ is linearly independent we have that

$$\sum_{\mathbf{x} \in \text{set}(\mathbf{X})} \sum_{i \in A(\mathbf{x})} (\alpha_i - \beta_i) \cdot \mathbf{x} = 0,$$

1517 if, and only if, $\sum_{i \in A(\mathbf{x})} (\alpha_i - \beta_i) = 0$, for all $\mathbf{x} \in \text{set}(\mathbf{X})$. \square

1518 **Lemma 52** (Proof of Lemma 27). *Let $\mathbf{v}, \mathbf{w} \in \mathbb{Q}^{1 \times L}$ and let $\mathbf{X} \in \{0, 1\}^{L \times d}$ be a matrix whose rows*
 1519 *are one-hot vectors, for some $L, d \in \mathbb{N}^+$. Then, $[\mathbf{v}]_{\mathbf{X}} = [\mathbf{w}]_{\mathbf{X}}$, if and only if for every $\mathbf{x} \in \text{set}(\mathbf{X})$,*

$$\{\{\mathbf{v}_i \mid i \in [n] \wedge \mathbf{X}_i = \mathbf{x}\}\} = \{\{\mathbf{w}_i \mid i \in [n] \wedge \mathbf{X}_i = \mathbf{x}\}\}.$$

1520 *Proof.* We define, for each $\mathbf{x} \in \text{set}(\mathbf{X})$,

$$\begin{aligned} V(\mathbf{x}) &:= \{\{\mathbf{v}_i \mid i \in [n] \wedge \mathbf{X}_i = \mathbf{x}\}\} \\ W(\mathbf{x}) &:= \{\{\mathbf{w}_i \mid i \in [n] \wedge \mathbf{X}_i = \mathbf{x}\}\}. \end{aligned}$$

1521 For the forward implication, assume towards a contradiction that $[\mathbf{v}]_{\mathbf{X}} = [\mathbf{w}]_{\mathbf{X}}$ but there exists an
 1522 $\mathbf{x} \in \text{set}(\mathbf{X})$ such that $V(\mathbf{x}) \neq W(\mathbf{x})$. However, then there also exists a number $v \in V(\mathbf{x})$ that
 1523 appears x times in $V(\mathbf{x})$ but y times in $W(\mathbf{x})$, with $x \neq y$. Without loss of generality, we assume
 1524 that $x < y$. Then, the tuple (v, \mathbf{x}) appears fewer times in $[\mathbf{v}]_{\mathbf{X}}$ than in $[\mathbf{w}]_{\mathbf{X}}$, implying $[\mathbf{v}]_{\mathbf{X}} \neq [\mathbf{w}]_{\mathbf{X}}$,
 1525 a contradiction.

1526 For the backward implication, assume towards a contradiction that for all $\mathbf{x} \in \text{set}(\mathbf{X})$, $V(\mathbf{x}) = W(\mathbf{x})$
 1527 but $[\mathbf{v}]_{\mathbf{X}} \neq [\mathbf{w}]_{\mathbf{X}}$. Then, there exists a tuple (v, \mathbf{x}) that appears x times in $[\mathbf{v}]_{\mathbf{X}}$ but y times in $[\mathbf{w}]_{\mathbf{X}}$,
 1528 with $x \neq y$. Without loss of generality, we assume that $x < y$. But then, for the vector \mathbf{x} , there
 1529 exists a number v that appears fewer times in $V(\mathbf{X})$ than in $W(\mathbf{X})$, implying $V(\mathbf{X}) \neq W(\mathbf{X})$, a
 1530 contradiction. This shows the statement. \square